

Statistical Relational Learning Notes

Ganesh Ramakrishnan and Ashwin Srinivasan

January 10, 2009

Contents

1	Sets, Relations and Logic	3
1.1	Sets and Relations	3
1.1.1	Sets	3
1.1.2	Relations	4
1.2	Logic	13
1.3	Propositional Logic	14
1.3.1	Syntax	15
1.3.2	Semantics	17
1.3.3	Inference	27
1.3.4	Resolution	28
1.3.5	Davis-Putnam Procedure	40
1.3.6	Local Search Methods	44
1.3.7	Default inference under closed world assumption	45
1.3.8	Lattice of Models	50
1.4	First-Order Logic	55
1.4.1	Syntax	56
1.4.2	Semantics	64
1.4.3	From Datalog to Prolog	71
1.4.4	Lattice of Herbrand Models	74
1.4.5	Inference	75
1.4.6	Subsumption Revisited	80
1.4.7	Subsumption Lattice over Atoms	81
1.4.8	Covers of Atoms	86
1.4.9	The Subsumption Theorem Again	89
1.4.10	Proof Strategies Once Again	90
1.4.11	Execution of Logic Programs	93
1.4.12	Answer Substitutions	97
1.4.13	Theorem Proving for Full First-Order Logic	99
1.5	Adequacy	99
2	Exploring a Structured Search Space	101
2.1	Subsumption Lattice over Clauses	101
2.1.1	Lattice Structure of Clauses	103
2.1.2	Covers in the Subsume Order	108

2.2	The Implication Order	109
2.3	Inverse Reduction	112
2.4	Generality order and ILP	112
2.5	Incorporating Background Knowledge	115
2.5.1	Plotkin's relative subsumption ($\succeq_{\mathcal{B}}$)	116
2.5.2	Relative implication ($\models_{\mathcal{B}}$)	122
2.5.3	Buntine's generalized subsumption ($\geq_{\mathcal{B}}$)	124
2.6	Using Generalization and Specialization	126
2.7	Using the Cover Structure: Refinement Operators	126
2.7.1	Example	127
2.7.2	Ideal Refinement Operators	128
2.7.3	Refinement Operators on Clauses for Subsumption	130
2.7.4	Refinement Operators on Theories	133
2.8	Inverse Resolution	133
2.8.1	V-operator	133
2.8.2	Predicate Invention: W-operator	135
2.9	Summary	139
3	Linear Algebra	143
3.1	Linear Equations	143
3.2	Vectors and Matrices	145
3.3	Solution of Linear Equations by Elimination	149
3.3.1	Elimination as Matrix Multiplication	151
3.4	Solution of Linear Equations by Matrix Inversion	156
3.4.1	Inverse Matrices	156
3.4.2	Gauss-Jordan Elimination	158
3.5	Solution of Linear Equations using Vector Spaces	160
3.5.1	Vector Spaces and Matrices	160
3.5.2	Null Space	162
3.6	Elimination for Computing the Null Space ($A\mathbf{x} = 0$)	163
3.6.1	Pivot Variables and Row Echelon Form	163
3.6.2	Reduced Row Echelon Form	166
3.6.3	Solving $A\mathbf{x} = \mathbf{b}$	168
3.7	Independence, Basis, and Dimension	173
3.7.1	Independence	174
3.7.2	Basis and Dimension	175
3.8	Matrix Spaces	178
3.9	Orthogonality and Projection	180
3.9.1	Projection Matrices	181
3.9.2	Least Squares	181
3.9.3	Orthonormal Vectors	183
3.9.4	Gram-Schmidt Orthonormalization	185
3.9.5	Fourier and Wavelet Basis	186
3.10	Determinants	187
3.10.1	Formula for determinant	192
3.10.2	Formula for Inverse	194

3.11	Eigenvalues and Eigenvectors	195
3.11.1	Solving for Eigenvalues	196
3.11.2	Some Properties of Eigenvalues and Eigenvectors	197
3.11.3	Matrix Factorization using Eigenvectors	202
3.12	Positive Definite Matrices	203
3.12.1	Equivalent Conditions	204
3.12.2	Some properties	205
3.13	Singular Value Decomposition	206
3.13.1	Pseudoinverse	208
4	Convex Optimization	211
4.1	Introduction	211
4.1.1	Mathematical Optimization	211
4.1.2	Some Topological Concepts in \mathbb{R}^n	212
4.1.3	Optimization Principles for Univariate Functions	214
4.1.4	Optimization Principles for Multivariate Functions	229
4.1.5	Absolute extrema and Convexity	250
4.2	Convex Optimization Problem	251
4.2.1	Why Convex Optimization?	251
4.2.2	History	252
4.2.3	Affine Set	253
4.2.4	Convex Set	253
4.2.5	Examples of Convex Sets	255
4.2.6	Convexity preserving operations	258
4.2.7	Convex Functions	263
4.2.8	Convexity and Global Minimum	264
4.2.9	Properties of Convex Functions	266
4.2.10	Convexity Preserving Operations on Functions	277
4.3	Convex Optimization Problem	278
4.4	Duality Theory	280
4.4.1	Lagrange Multipliers	280
4.4.2	The Dual Theory for Constrained Optimization	286
4.4.3	Geometry of the Dual	290
4.4.4	Complementary slackness and KKT Conditions	293
4.5	Algorithms for Unconstrained Minimization	297
4.5.1	Descent Methods	298
4.5.2	Newton's Method	303
4.5.3	Variants of Newton's Method	307
4.5.4	Gauss Newton Approximation	307
4.5.5	Levenberg-Marquardt	309
4.5.6	BFGS	309
4.5.7	Solving Systems Large Sparse Systems	311
4.5.8	Conjugate Gradient	319
4.6	Algorithms for Constrained Minimization	322
4.6.1	Equality Constrained Minimization	322
4.6.2	Inequality Constrained Minimization	325

4.7	Linear Programming	328
4.7.1	Simplex Method	331
4.7.2	Interior point barrier method	338
4.8	Least Squares	342
4.8.1	Linear Least Squares	343
4.8.2	Least Squares with Linear Constraints	348
4.8.3	Least Squares with Quadratic Constraints	351
4.8.4	Total Least Squares	352
5	Statistics	353
5.1	Variables	353
5.2	Descriptive Statistics	354
5.2.1	Histograms	354
5.2.2	Scatter Diagram	357
5.3	Summary statistics	357
5.4	Standard measures of location	357
5.4.1	Standard measures of spread and association	359
5.4.2	Correlation coefficient	360
5.4.3	Ecological Correlation	363
5.5	Normal approximation	364
5.5.1	Standard Normal Curve	364
5.5.2	General Normal Curve	365
5.5.3	Standard Units	365
5.6	Regression	366
5.6.1	The Regression Effect	371
5.6.2	SD Line	371
5.7	Probability and sampling	371
5.8	Binomial distribution	371
5.9	Interval estimation	371
5.10	Some standard significance tests	371
6	Searching on Graphs	373
6.1	Search	373
6.1.1	Depth First Search	375
6.1.2	Breadth First Search (BFS)	376
6.1.3	Hill Climbing	376
6.1.4	Beam Search	378
6.2	Optimal Search	379
6.2.1	Branch and Bound	380
6.2.2	A* Search	381
6.3	Constraint Satisfaction	381
6.3.1	Algorithms for map coloring	383
6.3.2	Resource Scheduling Problem	385

7	Graphical Models	387
7.1	Semantics of Graphical Models	388
7.1.1	Directed Graphical Models	388
7.1.2	Undirected Graphical Models	394
7.1.3	Comparison between directed and undirected graphical models	397
7.2	Inference	399
7.2.1	Elimination Algorithm	400
7.2.2	Sum-product Algorithm	402
7.2.3	Max Product Algorithm	406
7.2.4	Junction Tree Algorithm	411
7.2.5	Junction Tree Propagation	414
7.2.6	Constructing Junction Trees	415
7.2.7	Approximate Inference using Sampling	419
7.3	Factor Graphs	429
7.4	Exponential Family	431
7.5	A Look at Modeling Continuous Random Variables	433
7.6	Exponential Family and Graphical Models	437
7.6.1	Exponential Models and Maximum Entropy	439
7.7	Model fitting	443
7.7.1	Sufficient Statistics	444
7.7.2	Maximum Likelihood	446
7.7.3	What the Bayesians Do	448
7.7.4	Model Fitting for Exponential Family	450
7.7.5	Maximum Likelihood for Graphical Models	452
7.7.6	Iterative Algorithms	454
7.7.7	Maximum Entropy Revisted	457
7.8	Learning with Incomplete Observations	459
7.8.1	Parameter Estimation for Mixture Models	460
7.8.2	Expectation Maximization	462
7.9	Variational Methods for Inference	466
8	Statistical Relational Learning	473
8.1	Role of Logical Abstraction in SRL	474
8.2	Inductive Logic Programming	475
8.3	Probabilistic ILP Settings	477
8.3.1	Probabilistic setting for justification	479
8.4	Learning from Entailment	483
8.4.1	Constraining the ILP Search	486
8.4.2	Example: Aleph	487
8.5	Probabilistic Learning from Entailment	495
8.5.1	Structure and Parameter Learning	496
8.5.2	Example: Extending FOIL	497
8.5.3	Example: Stochastic Logic Programming	498
8.6	Learning from Interpretations	498
8.7	Learning from Probabilistic Interpretations	499

8.7.1	Example: Markov Logic Networks	500
8.7.2	Example: Bayesian Logic Programs	502
8.8	Learning from Proofs	506
8.9	Probabilistic Proofs	507
8.9.1	Example: Extending GOLEM	509
8.10	Summary	511

Chapter 1

Sets, Relations and Logic

‘Crime is common. Logic is rare. Therefore it is upon the logic rather than the crime that you should dwell.’ Sherlock Holmes in Conan Doyle’s *The Copper Breeches*.

1.1 Sets and Relations

1.1.1 Sets

A *set* is a fundamental concept in mathematics. Simply speaking, it consists of some objects, usually called its *elements*. Here are some basic notions about sets that you must already know about:

- A set S with elements a, b and c is usually written as $S = \{a, b, c\}$. The fact that a is an element of S is usually denoted by $a \in S$.
- A set with no elements is called the “empty set” and is denoted by \emptyset .
- Two sets S and T are equal ($S = T$) if and only if they contain precisely the same elements. Otherwise $S \neq T$.
- A set T is a subset of a set S ($T \subseteq S$) if and only if every element of T is also an element of S . If $T \subseteq S$ and $S \subseteq T$ then $S = T$. Sometimes $T \subseteq S$ may sometimes also be written as $S \supseteq T$. If $T \subseteq S$ and S has at least one element not in T , then $T \subset S$ (T is said to be “proper subset” of S). Again, $T \subset S$ may sometimes be written as $S \supset T$.

We now look at the meanings of the *union*, *intersection*, and *equivalence* of sets. The intersection, or product, of sets S and T , denoted by $S \cap T$ or ST or $S \cdot T$ consists of all elements common to both S and T . $ST \subset S$ and $ST \subseteq T$ for all sets S and T . Now, if S and T have no elements in common, then they are said to be *disjoint* and $ST = \emptyset$. It should be easy for you to see that $\emptyset \subseteq S$ for all S and $\emptyset \cdot S = \emptyset$ for all S . The union, or sum, of sets S and T , denoted

by $S \cup T$ or $S + T$, is the set consisting of elements that belong at least to S or T . Once again, it should be a straightforward matter to see $S \subseteq S + T$ and $T \subseteq S + T$ for all S and T . Also, $S + \emptyset = S$ for all S . Finally, if there is a one-to-one correspondence between the elements of set S and set T , then S and T are said to be equivalent ($S \sim T$). Equivalence and subsets form the basis of the definition of an infinite set: if $T \subset S$ and $S \sim T$ then S is said to be an infinite set. The set of natural numbers \mathcal{N} is an example of an infinite set (any set $S \sim \mathcal{N}$ is said to be *countable set*).

1.1.2 Relations

A finite sequence is simply a set of n elements with a 1 – 1 correspondence with the set $\{1, \dots, n\}$ arranged in order of succession (an *ordered pair*, for example, is just a finite sequence with 2 elements). Finite sequences allow us to formalise the concept of a relation. If A and B are sets, then the set $A \times B$ is called the *cartesian product* of A and B and is denoted by all ordered pairs (a, b) such that $a \in A$ and $b \in B$. Any subset of $A \times B$ is a binary relation, and is called a relation from A to B . If $(a, b) \in R$, then aRb means “ a is in relation R to b ” or, “relation R holds for the ordered pair (a, b) ” or “relation R holds between a and b .” A special case arises from binary relations within elements of a single set (that is, subsets of $A \times A$). Such a relation is called a “relation in A ” or a “relation over A ”. There are some important kinds of properties that may hold for a relation R in a set A :

Reflexive. The relation is said to be reflexive if the ordered pair $(a, a) \in R$ for every $a \in A$.

Symmetric. The relation is said to be symmetric if $(a, b) \in R$ iff $(b, a) \in R$ for $a, b \in A$.

Transitive. The relation is said to be transitive if $(a, b) \in R$ and $(b, c) \in R$, then $(a, c) \in R$ for $a, b, c \in A$.

Here are some examples:

- The relation \leq on the set of integers is reflexive and transitive, but not symmetric.
- The relation $R = \{(1, 1), (1, 2), (2, 1), (2, 2), (3, 3), (4, 4)\}$ on the set $A = \{1, 2, 3, 4\}$ is reflexive, symmetric and transitive.
- The relation \div on the set \mathcal{N} defined as the set $\{(x, y) : \exists z \in \mathcal{N} \text{ s.t. } xz = y\}$ is reflexive and transitive, but not symmetric.
- The relation \perp on the set of lines in a plane is symmetric but neither reflexive nor transitive.

It should be easy to see a relation like R above is just a set of ordered pairs. Functions are just a special kind of binary relation F which is such that if $(a, b) \in F$ and $(a, c) \in F$ then $b = c$. Our familiar notion of a function F from a set A to a set B is one which associates with each $a \in A$ exactly one element $b \in B$ such that $(a, b) \in F$. Now, a function from a set A to itself is usually called a *unary operation* in A . In a similar manner, a *binary operation* in A is a function from $A \times A$ to A (recall $A \times A$ is the Cartesian product of A with itself: it is sometimes written as A^2). For example, if $A = \mathcal{N}$, then addition $(+)$ is a binary operation in A . In general, an n -ary operation F in A is a function from A^n to A , and if it is defined for every element of A^n , then A is said to be *closed* with respect to the operation F . A set which is closed for one or more n -ary operations is called an *algebra*, and a sub-algebra is a subset of such a set that remains closed with respect to those operations. For example:

- \mathcal{N} is closed wrt the binary operations of $+$ and \times , and \mathcal{N} along with $+$, \times form an algebra.
- The set \mathcal{E} of even numbers is a subalgebra of algebra of \mathcal{N} with $+$, \times . The set \mathcal{O} of odd numbers is not a subalgebra.
- Let $S \subseteq U$ and $S' \subseteq U$ be the set with elements of U not in S (the unary operation of complementation). Let $U = \{a, b, c, d\}$. The subsets of U with the operations of complementation, intersection and union form an algebra. (How many subalgebras are there of this algebra?)

Equivalence Relations

Any relation R in a set A for which all three properties hold (that is, R is reflexive, symmetric, and transitive) is said to be an “equivalence relation”. Suppose, for example, we are looking at the relation R over the set of natural numbers \mathcal{N} , which consists of ordered pairs (a, b) such that $a + b$ is even¹ You should be able to verify that R is an equivalence relation over \mathcal{N} . In fact, R allows us to split \mathcal{N} into two disjoint subsets: the set of odd numbers \mathcal{O} and the set of even numbers \mathcal{E} such that $\mathcal{N} = \mathcal{O} \cup \mathcal{E}$ and R is an equivalence relation over each of \mathcal{O} and \mathcal{E} . This brings us to an important property of equivalence relations:

Theorem 1 *Any equivalence relation E over a set S partitions S into disjoint non-empty subsets S_1, \dots, S_k such that $S = S_1 \cup \dots \cup S_k$.*

Let us see how E can be used to partition S by constructing subsets of S in the following way. For every $a \in S$, if $(a, b) \in E$ then a and b are put in the same subset. Let there be k such subsets. Now, since $(a, a) \in E$ for every $a \in S$, every element of S is in some subset. So, $S = S_1 \cup \dots \cup S_k$. It also follows that the subsets are disjoint. Otherwise there must be some $c \in S_i, S_j$. Clearly, S_i

¹Equivalence is often denoted by \approx . Thus, for an equivalence relation E , if $(a, b) \in E$, then $a \approx b$.

and S_j are not singleton sets. Suppose S_i contains at least a and c . Further let there be a $b \notin S_i$ but $b \in S_j$. Since $a, c \in S_i$, $(a, c) \in E$ and since $c, b \in S_j$, $(c, b) \in E$. Thus, we have $(a, c) \in E$ and $(c, b) \in E$, which must mean that $(a, b) \in E$ (E is transitive). But in this case b must be in the same subset as a by construction of the subsets, which contradicts our assumption that $b \notin S_i$. The converse of this is also true:

Theorem 2 *Any partition of a set S partitions into disjoint non-empty subsets S_1, \dots, S_k such that $S = S_1 \cup \dots \cup S_k$ results in an equivalence relation over S .*

(Can you prove that this is the case? Start by constructing a relation E , with $(a, b) \in E$ if and only if a and b are in the same block, and prove that E is an equivalence relation.)

Each of the disjoint subsets S_1, S_2, \dots are called "equivalence classes", and we will denote the equivalence class of an element a in a set S by $[a]$. That is, for an equivalence relation E over a set S :

$$[a] = \{x : x \in S, (a, x) \in E\}$$

What we are saying above is that the collection of all equivalence classes of elements of S forms a partition of S ; and conversely, given a partition of the set S , there is an equivalence relation E on S such that the sets in the partition (sometimes also called its "blocks") are the equivalence classes of S .

Partial Orders

Given an equality relation $=$ over elements of a set S , a partial order \preceq over S is a relation over S that satisfies the following properties:

Reflexive. For every $a \in S$, $a \preceq a$

Anti-Symmetric. If $a \preceq b$ and $b \preceq a$ then $a = b$

Transitive. If $a \preceq b$ and $b \preceq c$ then $a \preceq c$

Here are some properties about partial orders that you should know (you will be able to understand them immediately if you take, as a special case, \preceq as meaning \leq and \prec as meaning $<$):

- If $a \preceq b$ and $a \neq b$ then $a \prec b$
- $b \succ a$ means $a \preceq b$, $b \succ a$ means $a \prec b$
- If $a \preceq b$ or $b \preceq a$ then a, b are comparable, otherwise they are not comparable.

A set S over which a relation of partial order is defined is called a *partially ordered set*. It is sometimes convenient to refer to a set S and a relation R defined over S together by the pair $\langle S, R \rangle$. So, here are some examples of partially ordered sets $\langle S, \preceq \rangle$:

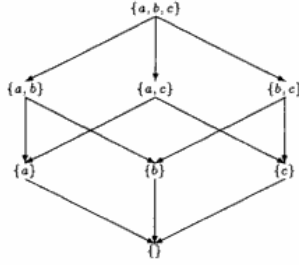


Figure 1.1: The lattice structure of $\langle S, \preceq \rangle$, where S is the power set of $\{a, b, c\}$.

- S is a set of sets, $S_1 \preceq S_2$ means $S_1 \subseteq S_2$
- $S = \mathcal{N}$, $n_1 \preceq n_2$ means $n_1 = n_2$ or there is a $n_3 \in \mathcal{N}$ such that $n_1 + n_3 = n_2$
- S is the set of equivalence relations E_1, \dots over some set T , $E_L \preceq E_M$ means for $u, v \in T$, $uE_L v$ means $uE_M v$ (that is, $(u, v) \in E_L$ means $(u, v) \in E_M$).

Given a set $S = \{a, b, \dots\}$ if $a \prec b$ and there is no $x \in S$ such that $a \prec x \prec b$ then we will say b *covers* a or that a is a *downward cover* of b . Now, suppose S_{down} be a set of downward covers of $b \in S$. If for all $x \in S$, $x \prec b$ implies there is an $a \in S_{down}$ s.t. $x \preceq a \prec b$, then S_{down} is said to be a *complete* set of downward covers of b . Partially ordered sets are usually shown as diagrams like in Figure 1.1.

The diagrams, as you can see, are graphs (sometimes called Hasse graphs or Hasse diagrams). In the graph, vertices represent elements of the partially ordered set. A vertex v_2 is at a higher level than vertex v_1 whenever $v_1 \prec v_2$, and there is an edge between the two vertices only if v_2 covers v_1 (that is, v_2 is an immediate predecessor). The graph is therefore really a directed one, in which there is a directed edge from a vertex v_2 to v_1 whenever v_2 covers v_1 . Also, since the relation is anti-symmetric, there can be no cycles. So, the graph is a directed acyclic graph, or DAG.

In the diagram in Figure ?? on the left, S is the set of non-empty subsets of $\{a, b, c\}$ and \preceq denotes the subset relationship (that is, $S_1 \preceq S_2$ if and only if $S_1 \subset S_2$). The diagram on the right is an example of a *chain*, or a *totally ordered* set.

You should be able to see that a finite chain of length n can be put in a one-to-one correspondence to a finite sequence of natural numbers $(1, \dots, n)$ (the correct way to say this is that a finite chain is isomorphic with a finite sequence of natural numbers). In general, a partially ordered set S is a chain if for every pair $a, b \in S$, $a \prec b$ or $b \prec a$. There is a close relationship between a partially ordered set and a chain. Suppose S is a partially ordered set. We

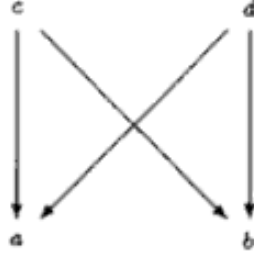
can always associate a function f from the elements of S to \mathcal{N} (the set of natural numbers), so that if $a \prec b$ for $a, b \in S$, then $f(a) < f(b)$. f is called a *consistent enumeration* of S , and is not unique and we can use it to define a chain consistent with S . (We will leave the proof of the existence a consistent enumeration for you. One way would be to use the method of induction on the number of elements in S : clearly there is such an enumeration for $|S| = 1$. Assume that an enumeration exists for $|S| = n - 1$ and prove it for $|S| = n$.)

Some elements of a partially ordered set have special properties. Let $\langle S, \preceq \rangle$ be a p.o. set and $T \subseteq S$. Then (in the following, you should read the symbol \exists as being shorthand for “there exists”, and \forall as “for all”):

<ul style="list-style-type: none"> – Least element of T $a \in T \text{ s.t. } \forall t \in T \ a \preceq t$ 	<ul style="list-style-type: none"> – Greatest element of T $a \in T \text{ s.t. } \forall t \in T \ a \succeq t$
<ul style="list-style-type: none"> – Least element, if it exists, is unique. If $T = S$ this is the “zero” element 	<ul style="list-style-type: none"> – Greatest element, if it exists is unique. If $T = S$ then this is the “unity” element
<ul style="list-style-type: none"> – Minimal element of T $a \in T \ \nexists t \in T \text{ s.t. } t \prec a$ 	<ul style="list-style-type: none"> – Maximal element of T $a \in T \ \nexists t \in T \text{ s.t. } t \succ a$
<ul style="list-style-type: none"> – Minimal element need not be unique 	<ul style="list-style-type: none"> – Maximal element need not be unique
<ul style="list-style-type: none"> – Lower bound of T $b \in S \text{ s.t. } b \preceq t \ \forall t \in T$ 	<ul style="list-style-type: none"> – Upper bound of T $b \in S \text{ s.t. } b \succeq t \ \forall t \in T$
<ul style="list-style-type: none"> – Glb g of T $b \preceq g \ \forall b, g: \text{ lbs of } T$ 	<ul style="list-style-type: none"> – Lub g of T $b \succeq g \ \forall b, g: \text{ ub of } T$
<ul style="list-style-type: none"> – If it exists, the glb is unique 	<ul style="list-style-type: none"> – If it exists the lub is unique

As you would have observed, there is a difference between a least element and a minimal element (and correspondingly, between greatest and maximal elements). The requirement of a minimal (maximal) upper bound is, in some sense, a weakening of the requirement of a least (greatest) upper bound. If x and y are both lub’s of some set $T \subseteq S$, then $y \preceq x$ and $x \preceq y$, so then $x \approx y$. This means that all lub’s of T are equivalent. Dually, if x and y are glb’s of some T , then also $x \approx y$. Thus, if a least element exists, then it is unique: this is not necessarily the case with a minimal element. Also, least and greatest elements must belong to the set T , but lower and upper bounds need not.

For this example, S has: (1) one upper bound b ; (2) no lower bound; (3) a greatest element b ; (4) no least element; (5) no greatest lower bound; (6) two minimal elements a and e ; and (7) one maximal element b . Can you identify what the corresponding statements are for T ?

Figure 1.2: $\{a, b\}$ has no lub here.

The glb and lub are sometimes also taken to be binary operations on a partially ordered set S , that assigns to an ordered pair in S^2 the corresponding glb or lub. The first operation is called the *product* or *meet* and is denoted by \cdot or \sqcap . The second operation is sometimes called the *sum* or *join* and is denoted by $+$ or \sqcup .

In a quasi-ordered set, a subset need not have a lub or glb. We will take an example to illustrate this. Let $S = \{a, b, c, d\}$, and let \preceq be defined as $a \preceq c$, $b \preceq c$, $a \preceq d$ and $d \preceq b$. Then since c and d are incomparable, the set a, b has no lub in this quasi-order. See Figure 1.2.

Similarly, a set need not have a maximal or a minimal, nor upward or downward covers. For instance, let S be the infinite set $\{y, x_1, x_2, x_3, \dots\}$, and let \preceq be a quasi-order on S , defined as $y \prec \dots \prec x_{n+1} \prec x_n \prec \dots \prec x_2 \prec x_1$. Then there is no upward cover of y : for every x_n , there always is an x_{n+1} such that $y \prec x_{n+1} \prec x_n$. In this case, y has no complete set of upward covers.

Note that a complete set of upward covers for y need not contain all upward covers of y . However, in order to be complete, it should contain at least one element from each equivalence class of upward covers. On the other hand, even the set of all upward covers of y need not be complete for y . For the example given above, the set of all upward covers of y is empty, but obviously not complete.

A notion of some relevance later is that of a function f defined on a partially ordered set $\langle S, \preceq \rangle$. Specifically, we would like to know if the function is: (a) monotonic; and (b) continuous. Monotonicity first:

A function f on $\langle S, \preceq \rangle$ is monotonic if and only if for all $u, v \in S$,
 $u \preceq v$ means $f(u) \preceq f(v)$

Now, suppose a subset S_1 of S have a least upper bound $\text{lub}(S_1)$ (with some abuse of notation: here $\text{lub}(X)$ is taken to be the lub of the elements in set X). Such subsets are called “directed” subsets of S . Then:

A function f on $\langle S, \preceq \rangle$ is continuous if and only if for all directed subsets S_i of S , $f(\text{lub}(S_i)) = \text{lub}(\{f(x) : x \in S_i\})$.

That is, if a directed set S_i has a least upper bound $\text{lub}(S_i)$, then the set obtained by applying a continuous function f to the elements of S_i has least upper bound $f(\text{lub}(S_i))$. Functions that are both monotonic and continuous on some partially ordered set $\langle S, \preceq \rangle$ are of interest to us because they can be used, for some kinds of orderings, to guarantee that for some $s \in S$, $f(s) = s$. That is, f is said to have a “fixpoint”.

Lattices

A lattice is just a partially ordered set $\langle S, \preceq \rangle$ in which every pair of elements $a, b \in S$ has a glb (represented by \sqcap) and a lub (represented by \sqcup). From the definitions of lower and upper bounds, we are able to show that in any such partially ordered set, the operations will have the following properties:

- $a \sqcap b = b \sqcap a$, and $a \sqcup b = b \sqcup a$ (that is, they are commutative).
- $a \sqcap (b \sqcap c) = (a \sqcap b) \sqcap c$, and $a \sqcup (b \sqcup c) = (a \sqcup b) \sqcup c$ (that is, they are associative).
- $a \sqcap (a \sqcup b) = a$, and $a \sqcup (a \sqcap b) = a$ (that is, they are “absorptive”).
- $a \sqcap b = a$ and $a \sqcup b = b$.

We will not go into all the proofs here, but show one for illustration. Since $a \sqcap b$ is the glb of a and b , $a \sqcap b \preceq a$. Clearly then $a \sqcup (a \sqcap b)$, which is the lub of a and $a \sqcap b$, is a . This is one of the absorptive properties above. You should also be able to see, from these properties, that a lattice can also be seen simply as an algebra with two binary operations \sqcap and \sqcup that are commutative, associative and absorptive.

Theorem 3 *A lattice is an algebra with the binary operations of \sqcup and \sqcap .*

Here is an example of a lattice: let S be all the subsets of $\{a, b, c\}$, and for $X, Y \in S$, $X \preceq Y$ means $X \subseteq Y$, $X \sqcap Y = X \cap Y$ and $X \sqcup Y = X \cup Y$. Then $\langle S, \subseteq \rangle$ is a lattice. The empty set \emptyset is the zero element, and S is the unity element of the lattice. More generally, a lattice that has a zero or least element (which we will denote \perp), and a unity or greatest element (which we will denote \top) is called a *bounded* lattice. In such lattices, the following necessarily hold: $a \sqcup \top = \top$; $a \sqcap \top = a$; $a \sqcup \perp = a$; and $a \sqcap \perp = \perp$. A little thought should convince you that a finite lattice will always be bounded: if the lattice is the set $S = \{a_1, \dots, a_n\}$ then $\top = a_1 \sqcup \dots \sqcup a_n$ and $\perp = a_1 \sqcap \dots \sqcap a_n$. (But, does the reverse hold: will a bounded lattice always be finite?)

Two properties of subsets of lattices are of interest to us. First, a subset M of a lattice L is called a *sublattice* of L if M is also closed under the same binary operations of \sqcup and \sqcap defined for L (that is, M is a lattice with the same operations as those of L). Second, if a lattice L has the property that every subset of L has a lub and a glb, then the L is said to be a *complete* lattice. Clearly, every finite lattice is complete. Further, since every subset of

L has a lub and a glb, this must certainly be true of L itself. So, L has a lub, which must necessarily be the greatest element of L . Similarly, L has a glb, which must necessarily be the least element of L . In fact, the elements of L are ordered in such a way that each element is on some path from \top to \perp in the Hasse diagram. An example of an ordered set that is always a complete lattice is the set of all subsets of a set S , ordered by \subseteq , with binary operations \cap and \cup for the glb and lub. This set, the “powerset” of S , is often denoted by 2^S . So, if $S = \{a, b, c\}$, 2^S is the set $\{\emptyset, \{a\}, \{b\}, \{c\}, \{a, b\}, \{a, c\}, \{b, c\}, \{a, b, c\}\}$. Clearly, every subset of s of 2^S has both a glb and a lub in S .

There are two important results concerning complete lattices and functions defined on them. The Knaster-Tarski Theorem tells us that every monotonic function on a complete lattice $\langle S, \preceq \rangle$ has a least fixpoint.

Theorem 4 *Let $\langle S, \preceq \rangle$ be a complete lattice and let $f : S \rightarrow S$ be a monotonic function. Then the set of fixed points of f in L is also a complete lattice $\langle P, \preceq \rangle$ (which obviously means that f has a greatest as well as a least fixpoint).*

Proof Sketch:² Let $D = \{x \mid x \preceq f(x)\}$. From the very definition of D , it follows that every fixpoint is in D . Consider some $x \in D$. Then because f is monotone we have $f(x) \preceq f(f(x))$. Thus,

$$\forall x \in D, f(x) \in D \quad (1.1)$$

Let $u = \text{lub}(D)$ (which should exist according to our assumption that $\langle S, \preceq \rangle$ is a complete lattice. Then $x \preceq u$ and $f(x) \preceq f(u)$, so $x \preceq f(x) \preceq f(u)$. Therefore $f(u)$ is an upper bound of D . However, u is the least upper bound, hence $u \preceq f(u)$, which in turn implies that, $u \in D$. From (1.1), it follows that $f(u) \in D$. From $u = \text{lub}(D)$, $f(u) \in D$ and $u \preceq f(u)$, it follows that $f(u) = u$. Because every fixpoint is in D we have that u is the greatest fixpoint of f . Similarly, it can be proved that if $E = \{x \mid f(x) \preceq x\}$, then $v = \text{glb}(E)$ is a fixed point and therefore the smallest fixpoint of f . \square

Kleene’s First Recursion Theorem tells us how to find the element $s \in S$ that is the least fixpoint, by incrementally constructing lubs starting from applying a continuous function to the least element of the lattice (\perp).

Theorem 5 *Let S be a complete partial order and let $f : S \rightarrow S$ be a continuous (and therefore monotone) function. Then the least fixed point of f is the supremum of the ascending Kleene chain of f :*

$$\perp \preceq f(\perp) \preceq f(f(\perp)) \preceq \dots \preceq f^n(\perp) \preceq \dots$$

In the special case that \preceq is \subseteq , the incremental procedure starts with the empty set \emptyset , and progressive lub’s are obtained by application of the set-union operation \cup . We will not give the proofs of this result here.

²Can you complete the proof?

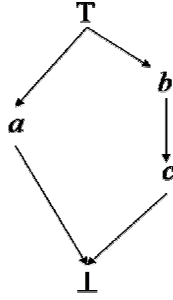


Figure 1.3: Example lattice for illustrating the concept of lattice length.

A final concept we will need is the concept of the length of a lattice. For a pair of elements a, b in a lattice L such that $a \preceq b$, the interval $[a, b]$ is the set $\{x : x \in L, a \preceq x \preceq b\}$. Now, consider a subset of $[a, b]$ that contains both a and b , and is such that any pair of elements in the subset are comparable. Then that subset is a chain from a to b : if the number of elements in the subset is n , then the length of the chain is $n - 1$. Maximal chains from a to b are those of the form $a = x_1 \prec x_2 \prec \cdots \prec x_n = b$ such that each x_i is covered by x_{i+1} . If all maximal chains from a to b are finite, then the longest of these defines the length of the interval $[a, b]$. For a bounded lattice, the length of the interval $[\perp, \top]$ defines the length of the lattice. So, in the lattice in Figure 1.3, there are two maximal chains between \perp and \top , of lengths 2 and 3 (what are these?). The length of lattice is thus equal to 3. Now, it should be evident that finite lattices will always have a finite length, but it is possible for lattices to have a finite length, but have infinitely many elements. For example, the lattice $L = \{\perp, \top, x_1, x_2, \dots\}$ such that $\perp \prec x_i \prec \top$ has a finite length (all maximal chains are of length 2). (Indeed, it is even possible to have an infinite set in which maximal chains are of finite, but increasing in lengths of $1, 2, \dots$)

Quasi-Orders

A quasi-order Q in a set S is a binary relation over S that satisfies the following properties:

Reflexive. For every $a \in S$, aQa

Transitive. If aQb and bQc then aQc

You can see that a quasi-order differs from an equivalence relation in that symmetry is not required. Further, it differs from a partial order because no equality is defined, and therefore the property of anti-symmetry property cannot be defined either. There are two important properties of quasi-orders, which we will present here without proof:

- If a quasi-order Q is defined on a set $S = \{a, b, \dots\}$, and we define a binary relation E as follows: aEb iff aQb and bQa . Then E is an equivalence relation.
- Let E partition S into subsets X, Y, \dots of equivalent elements. Let $T = \{X, Y, \dots\}$ and \preceq be a binary relation in T meaning $X \preceq Y$ in T if and only if xQy in S for some $x \in X, y \in Y$. Then T is partially ordered by \preceq .

What these two properties say is simply this:

A quasi-order Q over a set S results in a partial ordering over a set of equivalence classes of elements in S .

1.2 Logic

Logic, the study of arguments and ‘correct reasoning’, has been with us for at least the better part of two thousand years. In Greece, we associate its origins with Aristotle (384 B.C.–322 B.C.); in India with Gautama and the Nyaya school of philosophy (3rd Century B.C.?); and in China with Mo Ti (479 B.C.–381 B.C.) who started the Mohist school. Most of this dealt with the use and manipulation of *syllogisms*. It would only be a small injustice to say that little progress was made until Gottfried Wilhelm von Leibniz (1646–1716). He made a significant advance in the use of logic in mathematics by introducing symbols to represent statements and relations. Leibniz hoped to reduce all errors in human reasoning to minor calculational mistakes. Later, George Boole (1815–1864) established the connection between logic and sets, forming the basis of *Boolean algebra*. This link was developed further by John Venn (1834–1923) and Augustus de Morgan (1806–1872). It was around this time that Charles Dodgson (1832–1898), writing under the pseudonym Lewis Carroll, wrote a number of popular logic textbooks. Fundamental changes in logic were brought about by Friedrich Ludwig Gottlob Frege (1848–1925), who strongly rejected the idea that the laws of logic are synonymous with the laws of thought. For Frege, the former were laws of *truth*, having little to say on the processes by which human beings represent and reason with reality. Frege developed a logical framework that incorporated propositions with relations and the validity of arguments depended on the relations involved. Frege also introduced the device of *quantifiers* and *bound variables*, thus laying the basis for *predicate logic*, which forms the basis of all modern logical systems. All this and more is described by Bertrand Russell (1872–1970) and Alfred North Whitehead (1861–1947) in their monumental work, *Principia Mathematica*. And then in 1931, Kurt Gödel (1906–1978) showed much to the dismay of mathematicians everywhere that formal systems of arithmetic would remain incomplete.

Rational agents require knowledge of their world in order to make rational decisions. With the help of a declarative (knowledge representation) language, this knowledge (or a portion of it) is represented and stored in a knowledge

base. A knowledge-base is composed of sentences in a language with a truth theory (logic), so that someone external to the system can interpret sentences as statements about the world (semantics). Thus, to express knowledge, we need a precise, *declarative* language. By a *declarative language*, we mean that

1. The system believes a statement S *iff* it considers S to be true, since one cannot believe S without an idea of what it means for the world to fulfill S .
2. The knowledge-based must be precise enough so that we must know, (1) which symbols represent sentences, (2) what it means for a sentence to be true, and (3) when a sentence follows from other sentences.

Two declarative languages will be discussed in this chapter: (0 order or) propositional logic and first order logic.

1.3 Propositional Logic

Formal logic is concerned with statements of fact, as opposed to opinions, commands, questions, exclamations *etc.* Statements of fact are assertions that are either true or false, the simplest form of which are called *propositions*. Here are some examples of propositions:

The earth is flat.
Humans are monkeys.
 $1 + 1 = 2$

At this stage, we are not saying anything about whether these are true or false: just that they are sentences that are one or the other. Here are some examples of sentences that are not propositions:

Who goes there?
Eat your broccoli.
This statement is false.

It is normal to represent propositions by letters like P, Q, \dots . For example, P could represent the proposition ‘Humans are monkeys.’ Often, simple statements of fact are insufficient to express complex ideas. *Compound statements* can be combining two or more propositions with *logical connectives* (or simply, connectives). The connectives we will look at here will allow us to form sentences like the following:

It is not the case that P
 P and Q
 P or Q
 P if Q

The P 's and Q 's above are propositions, and the words underlined are the connectives. They have special symbols and names when written formally:

<u>Statement</u>	<u>Formally</u>	<u>Name</u>
It is <u>not</u> the case that P	$\neg P$	Negation
P <u>and</u> Q	$P \wedge Q$	Conjunction
P <u>or</u> Q	$P \vee Q$	Disjunction
P <u>if</u> Q	$P \leftarrow Q$	Conditional

There is, for example, a form of argument known to logicians as the *disjunctive syllogism*. Here is one due to the Stoic philosopher Chrysippus, about a dog chasing a rabbit. The dog arrives at a fork in the road, sniffs at one path and then dashes down the other. Chrysippus used formal logic to describe this:³

<u>Statement</u>	<u>Formally</u>
The rabbit either went down Path A or Path B.	$P \vee Q$
It did not go down Path A.	$\neg P$
Therefore it went down Path B.	$\therefore Q$

Here P represents the proposition 'The rabbit went down Path A' and Q the proposition 'The rabbit went down Path B.' To argue like Chrysippus requires us to know how to write correct logical sentences, ascribe truth or falsity to propositions, and use these to derive valid consequences. We will look at all these aspects in the sections that follow.

1.3.1 Syntax

Every language needs a *vocabulary*. For the language of propositional logic, we will restrict the vocabulary to the following:

Propositional symbols:	P, Q, \dots
Logical connectives⁴:	$\neg, \wedge, \vee, \leftarrow$
Brackets:	$(,)$

The next step is to specify the rules that decide how legal sentences are to be formed within the language. For propositional logic, legal sentences or *well-formed formulae* (wffs for short) are formed using the following rules:

1. Any propositional symbol is a wff;
2. If α is a wff then $\neg\alpha$ is a wff; and

³There is no suggestion that the principal agent in the anecdote employed similar means of reasoning.

3. If α and β are wffs then $(\alpha \wedge \beta)$, $(\alpha \vee \beta)$, and $(\alpha \leftarrow \beta)$ are wffs.

Wffs consisting simply of propositional symbols (Rule 1) are sometimes called *atomic* wffs and others *compound* wffs. Informally, it is acceptable to drop outermost brackets. Here are some examples of wffs and ‘non-wffs’:

<u>Formula</u>	<u>Comment</u>
$(\neg P)$	Not a wff. Parentheses are only allowed with the connectives in Rule 3
$\neg\neg P$	P is wff (Rule 1), $\neg P$ is wff (Rule 2), $\therefore \neg\neg P$ is wff (Rule 2)
$(P \leftarrow (Q \wedge R))$	P, Q, R are wffs (Rule 1), $\therefore (Q \wedge R)$ is a wff (Rule 3), $\therefore (P \leftarrow (Q \wedge R))$ is a wff (Rule 3)
$P \leftarrow (Q \wedge R)$	Not a wff, but acceptable informally
$((P) \wedge (Q))$	Not a wff. Parentheses are only allowed with the connectives in Rule 3
$(P \wedge Q \wedge R)$	Not a wff. Rule 3 only allows two symbols within a pair of brackets

One further kind of informal notation is widespread and quite readable. The conditional $(P \leftarrow ((Q_1 \wedge Q_2) \dots Q_n))$ is often written as $(P \leftarrow Q_1, Q_2, \dots Q_n)$ or even $P \leftarrow Q_1, Q_2, \dots Q_n$.

It is one thing to be able to write legal sentences, and quite another matter to be able to assess their truth or falsity. This latter requires a knowledge of semantics, which we shall look at shortly.

Normal Forms

Every formulae in propositional logic is equivalent to a formula that can be written as a conjunction of disjunctions. That is, something like $(A \vee B) \wedge (C \vee D) \wedge \dots$. When written in this way the formula is said to be in *conjunctive normal form* or CNF. There is another form, which consists of a disjunction of conjunctions, like $(A \wedge B) \vee (C \wedge D) \vee \dots$, called the *disjunctive normal form* or DNF. In general, a formula F in CNF can be written somewhat more cryptically as:

$$F = \bigwedge_{i=1}^n \left(\bigvee_{j=1}^m L_{i,j} \right)$$

and a formula G in DNF as:

$$G = \bigvee_{i=1}^n \left(\bigwedge_{j=1}^m L_{i,j} \right)$$

Here, $\bigvee F_i$ is short for $F_1 \vee F_2 \vee \dots$ and $\bigwedge F_i$ is short for $F_1 \wedge F_2 \wedge \dots$. In both CNF and DNF forms above, the $L_{i,j}$ are either propositions or negations of propositions (we shall shortly call these “literals”).

1.3.2 Semantics

There are three important concepts to be understood in the study of semantics of well-formed formulæ: *interpretations*, *models*, and *logical consequence*.

Interpretations

For propositional logic, an *interpretation* is simply an assignment of either *true* or *false* to all propositional symbols in the formula. For example, given the wff $(P \leftarrow (Q \wedge R))$ here are two different interpretations:

	P	Q	R
I_1 :	true	false	true
I_2 :	false	true	true

You can think of I_1 and I_2 as representing two different ‘worlds’ or ‘contexts’. After a moment’s thought, it should be evident that for a formula with N propositional symbols, there can never be more than 2^N possible interpretations.

Truth or falsity of a wff only makes sense given an interpretation (by the principle of bivalence, any interpretation can only result in a wff being either *true* or *false*). Clearly, if the wff simply consists of a single propositional symbol (recall that this was called an atomic wff), then the truth-value is simply that given by the interpretation. Thus, the wff P is *true* in interpretation I_1 and *false* in interpretation I_2 . To obtain the truth-value of compound wffs like $(P \leftarrow (Q \wedge R))$ requires a knowledge the semantics of the connectives. These are usually summarised in a tabular form known as *truth tables*. The truth tables for the connectives of interest to us are given below.

Negation. Let α be a wff⁵. Then the truth table for $\neg\alpha$ is as follows:

α	$\neg\alpha$
false	true
true	false

Conjunction. Let α and β be wffs. The truth table for $(\alpha \wedge \beta)$ is as follows:

⁵We will use Greek characters like α, β to stand generically for any wff.

α	β	$(\alpha \wedge \beta)$
false	false	false
false	true	false
true	false	false
true	true	true

Disjunction. Let α and β be wffs. The truth table for $(\alpha \vee \beta)$ is as follows:

α	β	$(\alpha \vee \beta)$
false	false	false
false	true	true
true	false	true
true	true	true

Conditional. Let α and β be wffs. The truth table for $(\alpha \leftarrow \beta)$ is as follows:

α	β	$(\alpha \leftarrow \beta)$
false	false	true
false	true	false
true	false	true
true	true	true

We are now in a position to obtain the truth-value of a compound wff. The procedure is straightforward: given an interpretation, we find the truth-values of the smallest ‘sub-wffs’ and then use the truth tables for the connectives to obtain truth-values for increasingly complex sub-wffs. For $(P \leftarrow (Q \wedge R))$ this means:

1. First, obtain the truth-values of P, Q, R using the interpretation;
2. Next, obtain the truth-value of $(Q \wedge R)$ using the truth table for ‘Conjunction’ and the truth-values of Q and R (Step 1); and
3. Finally, obtain the truth-value $(P \leftarrow (Q \wedge R))$ using the truth table for ‘Conditional’ and the truth-values of P (Step 1) and $(Q \wedge R)$ (Step 2).

For the interpretations I_1 and I_2 earlier these truth-values are as follows:

	P	Q	R	$(Q \wedge R)$	$(P \leftarrow (Q \wedge R))$
I_1 :	true	false	true	false	true
I_2 :	false	true	true	true	false

Thus, $(P \leftarrow (Q \wedge R))$ is *true* in interpretation I_1 and *false* in I_2 .

Models

Every interpretation (that is, an assignment of truth-values to propositional symbols) that makes a well-formed formula *true* is said to be a *model* for that formula. Take for example, the two interpretations I_1 and I_2 above. We have already seen that I_1 is a model for $(P \leftarrow (Q \wedge R))$; and that I_2 is not a model for the same formula. In fact, I_1 is also a model for several other wffs like: P , $(P \wedge R)$, $(Q \vee R)$, $(P \leftarrow Q)$, *etc.* Similarly, I_2 is a model for Q , $(Q \wedge R)$, $(P \vee Q)$, $(Q \leftarrow P)$, *etc.*

As another example, let $\{P, Q, R\}$ be the set of all atoms in the language, and α be the formula $((P \wedge Q) \leftrightarrow (R \rightarrow Q))$. Let I be the interpretation that makes P and R true, and Q false (so $I = \{P, R\}$). We determine whether α is true or false under I as follows:

1. P is true under I , and Q is false under I , so $(P \wedge Q)$ is false under I .
2. R is true under I , Q is false under I , so $(R \rightarrow Q)$ is false under I .
3. $(P \wedge Q)$ and $(R \rightarrow Q)$ are both false under I , so α is true under I .

Since α is true under I , I is a model of α . Let $I' = \{P\}$. Then $(P \wedge Q)$ is false, and $(R \rightarrow Q)$ is true under I' . Thus α is false under I' , and I' is not a model of α .

The definition of model can be extended to a set of formulæ; an interpretation I is said to be a model of a set of formulæ Σ if I is a model of all formulæ $\alpha \in \Sigma$. Σ is then said to have I as a model. We will offer an example to illustrate this extended definition. Let $\Sigma = \{P, (Q \vee I), (Q \rightarrow R)\}$, and let $I = \{P, R\}$, $I' = \{P, Q, R\}$, and $I'' = \{P, Q\}$ be interpretations. I and I' satisfy all formulas in Σ , so I and I' are models of Σ . On the other hand, I'' falsifies $(Q \rightarrow R)$, so I'' is not a model of Σ .

At this point, we can distinguish amongst two kinds of formulæ:

1. A wff may be such that *every* interpretation is a model. An example is $(P \vee \neg P)$. Since there is only one propositional symbol involved (P), there are at most $2^1 = 2$ interpretations possible. The truth table summarising the truth-values for this formula is:

	P	$\neg P$	$(P \vee \neg P)$
I_1 :	false	true	true
I_2 :	true	false	true

$(P \vee \neg P)$ is thus *true* in every possible ‘context’. Formulæ like these, for which every interpretation is a model are called *valid* or *tautologies*

2. A wff may be such that *none* of the interpretations is a model. An example is $(P \wedge \neg P)$. Again there is only one propositional symbol involved (P), and thus only two interpretations possible. The truth table summarising the truth-values for this formula is:

	P	$\neg P$	$(P \wedge \neg P)$
$I_1 :$	false	true	false
$I_2 :$	true	false	false

$(P \wedge \neg P)$ is thus *false* in every possible ‘context’. Formulae like these, for which none of the interpretations is a model are called *unsatisfiable* or *inconsistent*

Finally, any wff that has at least *one* interpretation as a model is said to be *satisfiable*.

Logical Consequence

We are often interested in establishing the truth-value of a formula given that of some others. Recall the Chrysippus argument:

<u>Statement</u>	<u>Formally</u>
The rabbit either went down Path A or Path B.	$P \vee Q$
It did not go down Path A.	$\neg P$
Therefore it went down Path B.	$\therefore Q$

Here, we want to establish that if the first two statements are true, then the third follows. The formal notion underlying all this is that of *logical consequence*. In particular, what we are trying to establish is that some well-formed formula α is the *logical consequence* of a conjunction of other well-formed formulae Σ (or, that Σ *logically implies* α). This relationship is usually written thus:

$$\Sigma \models \alpha$$

Σ being the conjunction of several wffs, it is itself a well-formed formula⁶. Logical consequence can therefore also be written as the following relationship between a pair of wffs:

$$((\beta_1 \wedge \beta_2) \dots \beta_n) \models \alpha$$

It is sometimes convenient to write Σ as the set $\{\beta_1, \beta_2, \dots, \beta_n\}$ which is understood to stand for the conjunctive formula above. But how do we determine if this relationship between Σ and α does indeed hold? What we want is the following: whenever the statements in Σ are true, α must also be true. In formal terms, this means: $\Sigma \models \alpha$ *if every model of Σ is also model of α* . Decoded:

⁶There is therefore nothing special needed to extend the concepts of validity and unsatisfiability to conjunctions of formulae like Σ . Thus, Σ is valid if and only if every interpretation is a model of the conjunctive wff (in other words, a model for each wff in the conjunction); and it is unsatisfiable if and only if none of the interpretations is a model of the conjunctive wff. It should be apparent after some reflection that if Σ is valid, then all logical consequences of it are also valid. On the other hand, if Σ is unsatisfiable, then any well-formed formula is a logical consequence.

- Recall that a model for a formula is an interpretation (assignment of truth-values to propositions) that makes that formula *true*;
- Therefore, a model for Σ is an interpretation that makes $((\beta_1 \wedge \beta_2) \dots \beta_n)$ *true*. Clearly, such an interpretation will make each of $\beta_1, \beta_2, \dots, \beta_n$ *true*;
- Let I_1, I_2, \dots, I_k be all the interpretations that satisfy the requirement above: that is, each is a model for Σ and there are no other models for Σ (recall that if there are N propositional symbols in Σ and α together, then there can be no more than 2^N such interpretations);
- Then to establish $\Sigma \models \alpha$, we have to check that each of I_1, I_2, \dots, I_k is also a model for α (that is, each of them make α *true*).

The definition of logical entailment can be extended to the entailment of sets of formulae. Let Σ and Γ be sets of formulas. Then Γ is said to be a logical consequence of Σ (written as $\Sigma \models \Gamma$), if $\Sigma \models \alpha$, for every formula $\alpha \in \Gamma$. We also say Σ (logically) implies Γ .

We are now in a position to see if Chrysippus was correct. We wish to see if $((P \vee Q) \wedge \neg P) \models Q$. From the truth tables on page 17, we can construct a truth table for $((P \vee Q) \wedge \neg P)$:

	P	Q	$(P \vee Q)$	$\neg P$	$((P \vee Q) \wedge \neg P)$
I_1 :	false	false	false	true	false
I_2 :	false	true	true	true	true
I_3 :	true	false	true	false	false
I_4 :	true	true	true	false	false

It is evident that of the four interpretations possible only one is a model for $((P \vee Q) \wedge \neg P)$, namely: I_2 . Clearly I_2 is also a model for Q . Therefore, every model for $((P \vee Q) \wedge \neg P)$ is also a model for Q ⁷. It is therefore indeed true that $((P \vee Q) \wedge \neg P) \models Q$. In fact, you will find you can ‘move’ formulae from left to right in a particular manner. Thus if:

$$((P \vee Q) \wedge \neg P) \models Q$$

then the following also hold:

$$(P \vee Q) \models (Q \leftarrow \neg P) \quad \text{and} \quad \neg P \models (Q \leftarrow (P \vee Q))$$

These are consequences of a more general result known as the *deduction theorem*, which we look at now. Using a set-based notation, let $\Sigma = \{\beta_1, \beta_2, \dots, \beta_i, \dots, \beta_n\}$. Then, the deduction theorem states:

⁷Although I_4 is also a model for Q , the test for logical consequence only requires us to examine those interpretations that are models of $((P \vee Q) \wedge \neg P)$. This precludes I_4 .

Theorem 6

$$\Sigma \models \alpha \text{ if and only if } \Sigma - \{\beta_i\} \models (\alpha \leftarrow \beta_i)$$

Proof: Consider first the case that $\Sigma \models \alpha$. That is, every model of Σ is a model of α . Now assume $\Sigma - \{\beta_i\} \not\models (\alpha \leftarrow \beta_i)$. That is, there is some model, say M , of $\Sigma - \{\beta_i\}$ that is not a model of $(\alpha \leftarrow \beta_i)$. That is, β_i is true and α is false in M . That is M is a model for $\Sigma - \{\beta_i\}$ and for β_i , but not a model for α . In other words, M is a model for Σ but not a model for α which is not possible. Therefore, if $\Sigma \models \alpha$ then $\Sigma - \{\beta_i\} \models (\alpha \leftarrow \beta_i)$. Now for the “only if” part. That is, let $\Sigma - \{\beta_i\} \models (\alpha \leftarrow \beta_i)$. We want to show that $\Sigma \models \alpha$. Once again, let us assume the contrary (that is, $\Sigma \not\models \alpha$. This means there must be a model M for Σ that is not a model for α . However, since $\Sigma = \{\beta_1, \beta_2, \dots, \beta_i, \dots, \beta_n\}$, M is both a model for $\Sigma - \{\beta_i\}$ and a model for each of the β_i . So, M cannot be a model for $(\alpha \leftarrow \beta_i)$. We are therefore in a position that there is a model M for $\Sigma - \{\beta_i\}$ that is not a model of $(\alpha \leftarrow \beta_i)$, which contradicts what was given. \square

The deduction theorem isn’t restricted to propositional logic, and holds for first-order logic as well. It can be invoked repeatedly. Here is an example of using it twice:

$$\Sigma \models \alpha \text{ if and only if } \Sigma - \{\beta_i, \beta_j\} \models (\alpha \leftarrow (\beta_i \wedge \beta_j))$$

With Chrysippus, applying the deduction theorem twice results in:

$$\{(P \vee Q), \neg P\} \models Q \text{ if and only if } \emptyset \models (Q \leftarrow ((P \vee Q) \wedge \neg P))$$

If $\emptyset \models (Q \leftarrow ((P \vee Q) \wedge \neg P))$ then every model for \emptyset must be a model for $(Q \leftarrow ((P \vee Q) \wedge \neg P))$. By convention, every interpretation is a model for \emptyset ⁸. It follows that every interpretation must be a model for $(Q \leftarrow ((P \vee Q) \wedge \neg P))$. Recall that this is just another way of stating that $(Q \leftarrow ((P \vee Q) \wedge \neg P))$ is valid (page 19)⁹.

What is the difference between the concepts of logical consequence denoted by \models and the connective \rightarrow in a statement such as $\Sigma \models \Gamma$? where, $\Sigma = \{(P \wedge Q), (P \rightarrow R)\}$ and $\Gamma = \{P, Q, R\}$? And how do these two notions of implication relate to the phrase ‘if...then’, often used in propositions or theorems? We delineate the differences below:

⁸That is, we take the empty set to denote a distinguished proposition *True* that is *true* in every interpretation. Correctly then, the formula considered is not $((\beta_1 \wedge \beta_2) \dots \beta_n)$ but $(True \wedge ((\beta_1 \wedge \beta_2) \dots \beta_n))$.

⁹To translate declarative knowledge into action (as in the case of the dog from Chrysippus’s anecdote), one of two possible strategies can be adopted. The first is called ‘Negative selection’ which involves *excluding any provably futile* actions. The second is called ‘Positive selection’ which involves *suggesting only actions that are provably safe*. There can be some actions that are neither provably safe nor provably futile.

1. The connective \rightarrow is a *syntactical symbol* called ‘if ... then’ or ‘implication’, which appears within formulæ. The truth value of the formula $(\alpha \rightarrow \xi)$ depends on the *particular* interpretation I we happen to be considering: according to the truth table, $(\alpha \rightarrow \xi)$ is true under I if α is false under I and/or ξ is true under I ; $(\alpha \rightarrow \xi)$ is false otherwise.
2. The concept of ‘logical consequence’ or ‘(logical) implication’, denoted by ‘ \models ’ describes a *semantical relation* between formulæ. It is defined in terms of *all* interpretations: ‘ $\alpha \models \xi$ ’ is true if every interpretation that is a model of α , is also a model of ξ .
3. The phrase ‘if. .. then’, which is used when stating, for example, propositions or theorems is also sometimes called ‘implication’. This describes a relation between assertions which are phrased in (more or less) natural language. It is used for instance in proofs of theorems, when we state that some assertion implies another assertion. Sometimes we use the symbols ‘ \vdash ’ or ‘ \vDash ’ for this. If assertion A implies assertion B, we say that B is a necessary condition for A (i.e., if A is true, B must necessarily be true), and A is a sufficient condition for B (i.e., the truth of B is sufficient to make A true). In other words, if A implies B, and B implies A, we write “A iff B”, where ‘iff’ abbreviates ‘if, and only if’.

Closely related to logical consequence is the notion of *logical equivalence*. A pair of wffs α and β are logically equivalent means:

$$\alpha \models \beta \quad \text{and} \quad \beta \models \alpha$$

This means the truth values for α and β are the same in all cases, and is usually written more concisely as:

$$\alpha \equiv \beta$$

Examples of logically equivalent formulæ are provided by De Morgan’s laws:

$$\neg(\alpha \vee \beta) \equiv \neg\alpha \wedge \neg\beta$$

$$\neg(\alpha \wedge \beta) \equiv \neg\alpha \vee \neg\beta$$

Also, if *True* denotes the formula that is *true* in every interpretation and *False* the formula that is *false* in every interpretation, then the following equivalences should be self-evident:

$$\alpha \equiv (\alpha \wedge \text{True})$$

$$\alpha \equiv (\alpha \vee \text{False})$$

More on the Conditional

We are mostly concerned with rules that utilise the logical connective \leftarrow . This makes this particular connective more interesting than the others, and it is worth noting some further details about it. Although we will present these here using examples from the propositional logic, the main points are just as applicable to formulæ in the predicate logic.

Recall the truth table for the conditional from page 18:

α	β	$(\alpha \leftarrow \beta)$
false	false	true
false	true	false
true	false	true
true	true	true

There is, therefore, only one interpretation that makes $(\alpha \leftarrow \beta)$ *false*. This may come as a surprise. Consider for example the statement:

The earth is flat \leftarrow Humans are monkeys

An interpretation that assigns *false* to both ‘The earth is flat’ and ‘Humans are monkeys’ makes this statement *true* (line 1 of the truth table). In fact, the only world in which the statement is false is one in which the earth is not flat, and humans are monkeys¹⁰. Consider now the truth table for $(\alpha \vee \neg\beta)$:

α	β	$\neg\beta$	$(\alpha \vee \neg\beta)$
false	false	true	true
false	true	false	false
true	false	true	true
true	true	false	true

It is evident from these truth tables that every model for $(\alpha \leftarrow \beta)$ is a model for $(\alpha \vee \neg\beta)$ and vice-versa. Thus:

$$(\alpha \leftarrow \beta) \equiv (\alpha \vee \neg\beta)$$

Thus, the conditional:

(Fred is human \leftarrow (Fred walks upright \wedge Fred has a large brain))

is equivalent to:

¹⁰The unusual nature of the conditional is due to the fact that it allows premises and conclusions to be completely unrelated. This is not what we would expect from conditional statements in normal day-to-day discourse. For this reason, the \leftarrow connective is sometimes referred to as the *material conditional* to distinguish it from a more intuitive notion.

(Fred is human $\vee \neg$ (Fred walks upright \wedge Fred has a large brain))

Or, using De Morgan's Law (page 23) and dropping some brackets for clarity:

Fred is human $\vee \neg$ Fred walks upright $\vee \neg$ Fred has a large brain

In this form, each of the premises on the right-hand side of the the original conditional (Fred walks upright, Fred has a large brain) appear negated in the final disjunction; and the conclusion (Fred is human) is unchanged. For a reason that will become apparent later we will use the term *clauses* to denote formulæ that contain propositions or negated propositions joined together by disjunction (\vee). We will also use the term *literals* to denote propositions or negated propositions. Clauses are thus disjunctions of literals.

It is common practice to represent a clause as a set of literals, with the disjunctions understood. Thus, the clause above can be written as:

{ Fred is human, \neg Fred walks upright, \neg Fred has a large brain }

The equivalence $\alpha \leftarrow \beta \equiv \alpha \vee \neg\beta$ also provides an alternative way of presenting the deduction theorem.

On page 22 the statement of this theorem was:

$$\Sigma \models \alpha \text{ if and only if } \Sigma - \{\beta_i\} \models (\alpha \leftarrow \beta_i)$$

This can now be restated as:

$$\Sigma \models \alpha \text{ if and only if } \Sigma - \{\beta_i\} \models (\alpha \vee \neg\beta_i)$$

The theorem thus operates as follows: when a formula moves from the left of \models to the right, it is negated and disjoined (using \vee) with whatever exists on the right. The theorem can also be used in the "other direction": when a formula moves from the right of \models to the left, it is negated and conjoined (using \wedge or \cup in the set notation) to whatever exists on the left. Thus:

$$\Sigma \models (\alpha \vee \neg\beta) \text{ if and only if } \Sigma \cup \{\neg\alpha\} \models \neg\beta$$

A special case of this arises from the use of the equivalence $\alpha \equiv (\alpha \vee False)$ (page 23):

$$\Sigma \models \alpha \text{ if and only if } \Sigma \models (\alpha \vee False) \text{ if and only if } \Sigma \cup \{\neg\alpha\} \models False$$

The formula *False* is commonly written as \square and the result above as:

$$\Sigma \models \alpha \text{ if and only if } \Sigma \cup \{\neg\alpha\} \models \square$$

The conditional ($\alpha \leftarrow \beta$) is sometimes mistaken to mean the same as ($\alpha \wedge \beta$). Comparison against the truth table for ($\alpha \wedge \beta$) shows that these two formulæ are not equivalent:

α	β	$(\alpha \wedge \beta)$
false	false	false
false	true	false
true	false	false
true	true	true

There are a number of ways in which $(\alpha \leftarrow \beta)$ can be translated in English. Some of the more popular ones are:

If β , then α α , if β β implies α
 β only if α β is sufficient for α α is necessary for β
 All β 's are α 's

Note the following related statements:

Conditional $(\alpha \leftarrow \beta)$
 Contrapositive $(\neg\beta \leftarrow \neg\alpha)$

It should be easy to verify the following equivalence:

Conditional \equiv Contrapositive $(\alpha \leftarrow \beta) \equiv (\neg\beta \leftarrow \neg\alpha)$

Errors of reasoning arise by assuming other equivalences. Consider for example the pair of statements:

S_1 : Fred is an ape \leftarrow Fred is human
 S_2 : Fred is human \leftarrow Fred is an ape

S_2 is the sometimes called the *converse* of S_1 . An interpretation that assigns *true* to 'Fred is an ape' and *false* to 'Fred is human' is a model for S_1 but not a model for S_2 . The two statements are thus not equivalent.

More on Normal Forms

We are now able to state two properties concerning normal forms:

1. If F is a formula in CNF and G is a formula in DNF, then $\neg F$ is a formula in DNF and $\neg G$ is a formula in CNF. This is a generalisation of De Morgan's laws and can be proved using the technique of mathematical induction (that is, show truth for a formula with a single literal; assume truth for a formula with n literals; and then show that it holds for a formula with $n + 1$ literals.)

2. Every formula F can be written as a formula F_1 in CNF and a formula F_2 in DNF. It is straightforward to see that any formula F can be written as a DNF formula by examining the rows of the truth table for F for which F is true. Suppose F consists of the propositions A_1, A_2, \dots, A_n . Then each such row is equivalent to some conjunction of literals L_1, L_2, \dots, L_n , where L_i is equal to A_i if A_i is true in that row and equal to $\neg A_i$ otherwise. Clearly, the disjunction of each row for which F is true gives the DNF formula for F . We can get the corresponding CNF formula G by negating the DNF formula (using the property above), or by examining the rows for which F is false in the truth table.

It should now be clear that a CNF expression is nothing more than a conjunction of a set of clauses (recall a clause is simply a disjunction of literals). It is therefore possible to convert any propositional formula F into CNF—either using the truth table as described, or using the following procedure:

1. Replace all conditional statements of the form $A \leftarrow B$ by the equivalent form using disjunction (that is, $A \vee \neg B$). Similarly replace all $A \leftrightarrow B$ with $(A \vee \neg B) \wedge (\neg A \vee B)$.
2. Eliminate double negations ($\neg\neg A$ replaced by A) and use De Morgan's laws wherever possible (that is, $\neg(A \wedge B)$ replaced by $(\neg A \vee \neg B)$ and $\neg(A \vee B)$ replaced by $(\neg A \wedge \neg B)$).
3. Distribute the disjunct \vee . For example, $(A \vee (B \wedge C))$ is replaced by $(A \vee B) \wedge (A \vee C)$.

An analogous process converts any formula to an equivalent formula in DNF. We should note that during conversion, formulae can expand exponentially. However, if only satisfiability should be preserved, conversion to CNF formula can be done polynomially.

1.3.3 Inference

Enumerating and comparing models is one way of determining whether one formula is a logical consequence of another. While the procedure is straightforward, it can be tedious, often requiring the construction of entire truth tables. A different approach makes no explicit reference to truth values at all. Instead, if α is a logical consequence of Σ , then we try to show that we can *infer* or *derive* α from Σ using a set of well-understood rules. Step-by-step application of these rules results in a *proof* that deduces that α follows from Σ . The rules, called *rules of inference*, thus form a system of performing calculations with propositions, called the *propositional calculus*¹¹. Logical implication can be mechanized by using a propositional calculus. We will first concentrate on a particular inference rule called *resolution*.

¹¹In general, a set of inference rules (potentially including so called logical axioms) is called a *calculus*.

1.3.4 Resolution

Before proceeding further, some basic terminology from proof theory may be helpful (this is not specially confined to the propositional calculus). Proof theory considers the *derivability* of formulæ, given a set of inference rules \mathcal{R} . Formulæ given initially are called *axioms* and those derived are *theorems*. That formula α is a theorem of a set of axioms Σ using inference rules \mathcal{R} is denoted by:

$$\Sigma \vdash_{\mathcal{R}} \alpha$$

When \mathcal{R} is obvious, this is simply written as $\Sigma \vdash \alpha$. The axioms can be valid (that is, all interpretations are models), or problem-specific (that is, only some interpretations may be models). The axioms together with the inference rules constitute what is called an *inference system*. The axioms together with all the theorems that are derivable from it is called a *theory*. A theory is said to be *inconsistent* if there is a formula α such that the theory contains both α and $\neg\alpha$.

We would like the theorems derived to be logical consequences of the axioms provided. For, if this were the case then by definition, the theorems will be *true* in all models of the axioms (recall that this is what logical consequence means). They will certainly be true, therefore, in any particular ‘intended’ interpretation of the axioms. Ensuring this property of the theorems depends entirely on the inference rules chosen: those that have this property are called *sound*. That is:

$$\text{if } \Sigma \vdash_{\mathcal{R}} \alpha \text{ then } \Sigma \models \alpha$$

Some well-known sound inference rules are:

$$\text{Modus ponens: } \{\beta, \alpha \leftarrow \beta\} \vdash \alpha$$

$$\text{Modus tollens: } \{\neg\alpha, \alpha \leftarrow \beta\} \vdash \neg\beta$$

Theorems derived by the use of sound inference rules can be added to the original set of axioms. That is, given a set of axioms Σ and a theorem α derived using a sound inference rule, $\Sigma \equiv \Sigma \cup \{\alpha\}$.

We would also like to derive *all* logical consequences of a set of axioms and rules with this property at said to be *complete*:

$$\text{if } \Sigma \models \alpha \text{ then } \Sigma \vdash_{\mathcal{R}} \alpha$$

Axioms and inference rules are not enough: we also need a strategy to select and apply the rules. An inference system (that is, axioms and inference rules) along with a strategy is called a *proof procedure*. We are especially interested here in a special inference rule called *resolution* and a strategy called SLD (the meaning of this is not important at this point: we will come to it later). The result is a proof procedure called *SLD-resolution*. Here we will simply illustrate the rule of resolution for manipulating propositional formulæ, and use an unconstrained proof strategy. A description of SLD will be left for a later section.

Suppose we are given as axioms the conditional formulæ (using the informal notation that replaces \wedge with commas):

β_1 : Fred is an ape \leftarrow Fred is human
 β_2 : Fred is human \leftarrow Fred walks upright, Fred has a large brain

Then the following is a theorem resulting from the use of resolution:

α : Fred is an ape \leftarrow Fred walks upright, Fred has a large brain

That α is indeed a logical consequence of $\beta_1 \wedge \beta_2$ can be checked by constructing truth tables for the formulæ: you will find that every interpretation that makes $\beta_1 \wedge \beta_2$ true will also make α true. More generally, here is the rule of resolution when applied to a pair of conditional statements:

$$\{(P \leftarrow Q_1, \dots, Q_i, \dots, Q_n), (Q_i \leftarrow R_1, \dots, R_m)\} \vdash \\ P \leftarrow Q_1, \dots, Q_{i-1}, R_1, \dots, R_m, Q_{i+1}, \dots, Q_n$$

The equivalence from page 24 ($\alpha \leftarrow \beta \equiv \alpha \vee \neg\beta$) allows resolution to be presented in a different manner (we have taken the liberty of dropping some brackets here):

$$\{(P \vee \neg Q_1 \vee \dots \vee \neg Q_i \vee \dots \vee \neg Q_n), (Q_i \vee \neg R_1 \vee \dots \vee \neg R_m)\} \vdash \\ P \vee \neg Q_1 \vee \dots \vee \neg Q_{i-1} \vee \neg R_1 \vee \dots \vee \neg R_m \vee \neg Q_{i+1} \vee \dots \vee \neg Q_n$$

On page 25, we introduced the terms clauses and literals. Thus, resolution applies to a pair of ‘parent’ clauses that contain *complementary* literals $\neg L$ and L . The result (the ‘resolvent’) is a clause containing all literals from each clause, except the complementary pair. Or, more abstractly, let C_1 and C_2 be a pair of clauses, and let $L \in C_1$ and $\neg L \in C_2$. Then, the resolvent of C_1 and C_2 is the clause:

$$R = (C_1 - \{L\}) \cup (C_2 - \{\neg L\})$$

Resolution of a pair of *unit clauses*—those that contain just single literals L and $\neg L$ —results in the the *empty clause*¹², or \square , which means that the parent clauses were inconsistent.

We can show that resolution is a sound inference rule.

Theorem 7 *Suppose R is the resolvent of clauses C_1 and C_2 . That is, $\{C_1, C_2\} \vdash R$. The resolution is sound, that is, $\{C_1, C_2\} \models R$.*

Proof: We want to show that if C_1 and C_2 are true and R is a resolvent of C_1 and C_2 then R is true. Let us assume C_1 and C_2 are true, and that R was obtained by resolving on some literal L in C_1 and C_2 . Further, let $C_1 = C \vee L$

¹²Note the difference between an empty clause \square and empty set of clauses $\{\}$. An interpretation I logically entails C iff there exists an $l \in C$ such that $I \models l$. I logically entails Σ if for all $C \in \Sigma$, $I \models C$. Thus, by definition, for all interpretations I , $I \not\models \square$ and $I \not\models \{\}$, whereas $I \models \{\}$.

and $C_2 = D \vee \neg L$, giving $R = C \vee D$. Now, L is either true or false. Suppose L is true. Then clearly C_1 is true, but since $\neg L$ is false and C_2 is true (by assumption), it must be that D , and hence R must be true. It is easy to see that we similarly arrive to the same conclusion about R even if L was false. \square

So, the theorems obtained by applying resolution to a set of axioms are all logical consequences of the axioms. In general, we will denote a clause C derived from a set of clauses Σ using resolution by $\Sigma \vdash_R C$. This means that there is a finite sequences of clauses $R_1, \dots, R_k = C$ such that each C_i (where C_i is a clause being resolved upon in the i^{th} resolution step) is either in Σ or is a resolvent of a pair of clauses already derived (that is, from $\{R_1, \dots, R_{i-1}\}$). Now, although it is the case that if $\Sigma \vdash_R \alpha$ then $\Sigma \models \alpha$, the reverse does not hold. For example, a moment's thought should convince you that:

$$\{\text{Fred is an ape, Fred is human}\} \models \text{Fred is an ape} \leftarrow \text{Fred is human}$$

However, using resolution, there is no way of deriving $\text{Fred is an ape} \leftarrow \text{Fred is human}$ from Fred is an ape and Fred is human . As an inference rule, resolution is thus incomplete¹³. However, it does have an extremely useful property known as *refutation completeness*. This is that if a formula Σ is inconsistent, then the empty clause \square will be eventually derivable by resolution. We will distinguish between the two completeness by referring to general completeness as affirmation completeness.

Thus, since $\text{Fred is an ape} \leftarrow \text{Fred is human}$ is a logical consequence of Fred is an ape and Fred is human , then the formula:

$$\Sigma : \{\text{Fred is an ape, Fred is human, } \neg(\text{Fred is an ape} \leftarrow \text{Fred is human})\}$$

must be inconsistent. This can be verified using resolution. First, the clausal form of $(\text{Fred is an ape} \leftarrow \text{Fred is human})$ is $(\text{Fred is an ape} \vee \neg \text{Fred is human})$. Using De Morgan's Law on this clausal form, we can see that $\neg(\text{Fred is an ape} \leftarrow \text{Fred is human})$ is equivalent to $\neg \text{Fred is an ape} \wedge \text{Fred is human}$. We can now rewrite Σ :

$$\Sigma' : \{\text{Fred is an ape, Fred is human, } \neg \text{Fred is an ape}\}$$

Resolution of the pair Fred is an ape , $\neg \text{Fred is an ape}$ would immediately result in the empty clause \square . The general steps in a refutation proof procedure using resolution are therefore:

- Let S be a set of clauses and α be a propositional formula. Let $C = S \cup \{\neg\alpha\}$.
- Repeatedly do the following:
 1. Select a pair of clauses C_1 and C_2 from C that can be resolved on some proposition P .

¹³It can be however proved that resolution is affirmation complete with respect to atomic conclusions.

2. Resolve C_1 and C_2 to give R .
3. If $R = \square$ then stop. Otherwise, if R contains both a proposition Q and its negation $\neg Q$ then discard R . Otherwise add R to C .

In general, we know that any formula F can be converted to a conjunction of clauses. We can distinguish between the following sets. $Res^0(F)$, which is simply the set of clauses in F . $Res^n(F)$, for $n > 0$, which is the clauses containing all clauses in $Res^{n-1}(F)$ and all clauses obtained by resolving a pair of clauses from $Res^{n-1}(F)$. Since there are only a finite number of propositional symbols in F and a finite number of clauses in its CNF, we can see that there will only be a finite number of clauses that can be obtained using resolution. That is, there is some m such that $Res^m(F) = Res^{m-1}(F)$. Let us call this final set consisting of all the original clauses and all possible resolvents $Res^*(F)$. Then, the property of refutation-completeness for resolution can be stated more formally as follows:

Theorem 8 *For some formula F , $\square \in Res^*(F)$ if and only if F is unsatisfiable.*

Though the resolution rule by itself is not (affirmation) complete for clauses in general, this property states that it is complete with respect to unsatisfiable sets of clauses. The complete proof of this will be provided on page 33. To get you started however, we show that if $\square \in Res^*(F)$ then F is unsatisfiable. We can assume that $\square \notin Res^0(F)$, since \square is not a disjunction of literals. Therefore there must be some k for which $\square \notin Res^k(F)$ and $\square \in Res^{k+1}(F)$. This can only mean that both L and $\neg L$ are in $Res^k(F)$. That is L and $\neg L$ are obtained from F by resolution. By the property of soundness of resolution, this means that $F \models (L \wedge \neg L)$. That is, F is unsatisfiable.

There are also other proof processes that are refutation-complete. Examples of such processes are the Davis-Putnam Procedure¹⁴, Tableaux Procedure, *etc.* In the worst case, the resolution search procedure can take exponential time. This, however, very probably holds for all other proof procedures. For CNF formulae in propositional logic, a type of resolution process called the Davis-Putnam Procedure (backtracking over all truth values) is probably (in practice) the fastest refutation-complete process.

The Subsumption Theorem

A property related to logical implication is that of subsumption. A propositional clause C *subsumes* a propositional clause D if $C \subseteq D$. What does this mean? It just means that every literal in C appears in D . Here are a pair of clauses C and D such that C subsumes D :

C : Fred is an ape

D : Fred is an ape \leftarrow Fred is human

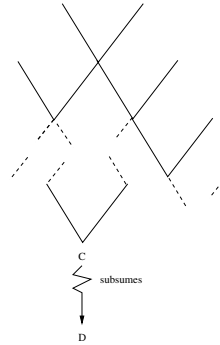
¹⁴It can be proved that the Davis-Putnam procedure is sound as well as complete.

In general, it should be easy to see that if C and D are clauses such that $C \subseteq D$, then $C \models D$. In fact, for propositional logic, it is also the case that if $C \models D$ then C subsumes D (we see why this is so shortly).

The notion of subsumption acts as the basis for an important result linking resolution and logical implication, called the *subsumption theorem*:

Theorem 9 *If Σ is a set of clauses and D is a clause. Then $\Sigma \models D$ if and only if D is a tautology or there is a clause C such that there is a derivation of C from Σ using resolution ($\Sigma \vdash_R C$) and C subsumes D .*

By “derivation of a clause C ” here, we mean the same as on page 30, that is, there is a sequence of clauses $R_1, \dots, R_k = C$ such that each R_i is either in Σ or is a resolvent of a pair of clauses in $\{R_1, \dots, R_{i-1}\}$. In effect, the Subsumption Theorem tells us that logical implication can be decomposed into a sequence of resolution steps, followed by a subsumption step:



Proof of Subsumption Theorem:

We will show that “only if” part of the theorem holds using the method of induction on the size of Σ :

1. Let $|\Sigma| = 1$. That is $\Sigma = \{\alpha_1\}$. Let $\Sigma \models D$. Since the only result of applying resolution to Σ is α_1 , we need to show that $\alpha_1 \subseteq D$. Suppose $\alpha_1 \not\subseteq D$. Let L be a literal in α_1 that is not in D . Let I be an interpretation that assigns L to true and all literals in D to false. Clearly I is a model for α_1 but not a model for D , which is not possible since $\Sigma \models D$. Therefore, $\alpha_1 \subseteq D$.
2. Let the theorem hold for $|\Sigma| = n$. We will see that it follows that it holds for $|\Sigma| = n + 1$. Let $\Sigma = \{\alpha_1, \dots, \alpha_{n+1}\}$ and $\Sigma \models D$. By the Deduction Theorem, we know that this means $\Sigma - \{\alpha_{n+1}\} \models (D \leftarrow \alpha_{n+1})$, or $\Sigma - \{\alpha_{n+1}\} \models (D \vee \neg\alpha_{n+1})$. Let us set $\Sigma' = \Sigma - \{\alpha_{n+1}\}$, and let L_1, \dots, L_k be the literals in α_{n+1} that do not appear in D . That is $\alpha_{n+1} = L_1 \vee \dots \vee L_k \vee D'$, where $D' \subseteq D$. $\Sigma' \models (D \vee \neg\alpha_{n+1})$, you should be able to see that $\Sigma' \models (D \vee \neg L_i)$ for $1 \leq i \leq k$. Since $|\Sigma'| = n$ and we

believe, by the induction hypothesis, that the subsumption theorem holds for $|\Sigma'| = n$, there must be some β_i for each L_i , such that $\Sigma' \vdash_R \beta_i$ and $\beta_i \subseteq (D \vee \neg L_i)$. Suppose $\neg L_i \notin \beta_i$. Then $\beta_i \subseteq D$, which means that β_i subsumes D . Since $\Sigma' \models \beta_i$ and $\Sigma \models \Sigma'$, the result follows. Now suppose $\neg L_i \in \beta_i$. That is, $\beta_i = \neg L_i \vee \beta'_i$, where $\beta'_i \subseteq D$. Clearly, we can resolve this with $\alpha_{n+1} = L_1 \vee \dots \vee L_i \vee \dots \vee L_k \vee D'$ to give $L_1 \vee \dots \vee L_{i-1} \vee \beta'_i \vee \dots \vee L_k \vee D'$. Progressively resolving against each of the β_i , we will be left with the clause $C = \beta'_1 \vee \beta'_2 \vee \dots \vee \beta'_k \vee D'$. Since C is the result of resolutions using a clause from Σ (that is α_{n+1}) and clauses derivable from $\Sigma' \subseteq \Sigma$, it is evident that $\Sigma \vdash_R C$. Also, since $\beta'_i \subseteq D$ and $D' \subseteq D$, $C \subseteq D$ and the result follows.

You should be able to see that the proof in the other direction (the “if” part) follows easily enough from the soundness of resolution. \square

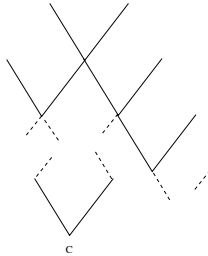
An immediate consequence of the Subsumption Theorem is that refutation-completeness of resolution follows.

Proof of Theorem 8

Recall what refutation completeness of resolution means: if Σ is a set of clauses that is unsatisfiable, then the empty clause \square is derivable using resolution. If Σ is unsatisfiable, then $\Sigma \models \square$. From the Subsumption Theorem, we know that if $\Sigma \models \square$, there must be a clause C such that $\Sigma \vdash_R C$ and C subsumes \square . But the only clause subsuming \square is \square itself. Hence $C = \square$, which means that if $\Sigma \models \square$ then $\Sigma \vdash_R \square$.

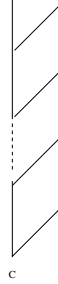
Proofs Using Resolution

So far, we have no strategy for directing the clauses obtained using resolution. Clauses are derived using any pair of clauses with complementary literals, and the process simply continues until we find the clause we want (for example, \square if we are interested in a proof of unsatisfiability). This procedure is clearly quite inefficient, since there is almost nothing constraining a proof, other than the presence of complementary literals. A “proof” for a clause C then ends up looking something like this:



Being creatures of limited patience and resources, we would like more directed approach. We can formalise this by changing our notion of a derivation. Recall

what we have been using so far: the derivation of a clause C from a set of clauses Σ means there is a sequence of clauses $R_1, \dots, R_k = C$ such that each C_i is either in Σ or is a resolvent of a pair of clauses in $\{R_1, \dots, R_{i-1}\}$. This results in the unconstrained form of a proof for C . We will say that there is a *linear* derivation for C from Σ if there is a sequence $R_0, \dots, R_k = C$ such that $R_0 \in \Sigma$ and each R_i ($1 \leq i \leq k$) is a resolvent of R_{i-1} and a clause $C_i \in \Sigma \cup \{R_0, \dots, R_{i-2}\}$. With a little thought, you should be able to convince yourself that this will result in a derivation with a “linear” look:



Here, you can see that each new resolvent forms one of the clauses for the next resolution step. The other clause—sometimes called the “side clause”—can be any one of the clauses in Σ or a previous resolvent. For reasons evident from the diagram above, the proof strategy is called linear resolution, and we will extend our notation to indicate both the inference rule and the proof strategy. Thus $\Sigma \vdash_{LR} C$ will mean that C is derived from Σ using linear resolution. We can restrict things even further, by requiring side clauses to be only from Σ . The resulting proof strategy, called *input resolution* is important as it is a generalised form of SLD resolution, first mentioned on page 28.

While the restrictions imposed by the proof strategies ensure that proofs are more directed (and hence efficient), it is important at this point to ask: at what cost? Of course, since we are still using resolution as an inference rule, the individual (and overall) inference steps remain sound. But what about completeness? By this we mean refutation-completeness, since this is the only kind of completeness we were able to show with unconstrained resolution. In fact, it is the case that linear resolution retains the property of refutation-completeness, but input resolution for arbitrary clauses does not. That input resolution is not refutation complete can be proved using a simple counter-example:

$$\begin{aligned} C_0 &: \text{Fred is an ape} \leftarrow \text{Fred is human} \\ C_1 &: \text{Fred is an ape} \leftarrow \neg \text{Fred is human} \\ C_2 &: \neg \text{Fred is an ape} \leftarrow \text{Fred is human} \\ C_3 &: \neg \text{Fred is an ape} \leftarrow \neg \text{Fred is human} \end{aligned}$$

Now, a little effort should convince you that this set of clauses is unsatisfiable. But input resolution will simply yield a sequence of resolvents: Fred is human, Fred is an ape, Fred is human, . . .

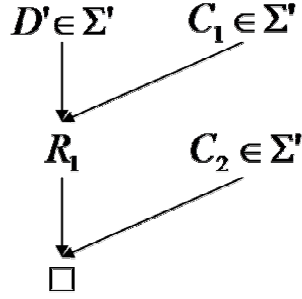


Figure 1.4: Part of case 1 of the proof for theorem 10.

Theorem 10 *If Σ is an unsatisfiable set of clauses, and $C \in \Sigma$ such that $\Sigma \setminus \{C\}$ is satisfiable, then there is a linear refutation of Σ with C as top clause.*

Proof:

We can assume Σ is finite. Let n be the number of distinct atoms occurring in literals in literals in Σ . We prove the lemma by induction on n

1. If $n = 0$, then $\Sigma = \{\square\}$. Since $\Sigma \setminus \{C\}$ is satisfiable, $C = \square$.
2. Suppose the lemma holds for $n \leq m$, and suppose $m + 1$ distinct atoms appear in Σ . We distinguish two cases.

- *Case 1:* Suppose $C = L$, where L is a literal. We first delete all clauses from Σ which contain the literal L (so we also delete C itself from Σ). Then we replace clauses which contain the literal $\neg L$ by clauses constructed by deleting these $\neg L$ (so for example, $L_1 \vee \neg L \vee L_2$ will be replaced by $L_1 \vee L_2$). Call the finite set obtained in this way Γ .

Note that neither the literal L , nor its negation, appears in clauses in Γ . If M were a Herbrand model of Γ , then $M \cup \{L\}$ (*i.e.*, the Herbrand interpretation which makes L true, and is the same as M for other literals) would be a Herbrand model of Σ . Thus since Σ is unsatisfiable, Γ must be unsatisfiable.

Now let Σ' be an unsatisfiable subset of Γ , such that every proper subset of Σ' is satisfiable. Σ' must contain a clause D' obtained from a member of Σ which contained $\neg L$, for otherwise the unsatisfiable set Σ' would be a subset of $\Sigma \setminus \{C\}$, contradicting the assumption that $\Sigma \setminus \{C\}$ is satisfiable. By construction of Σ' , we have that $\Sigma' \setminus \{D'\}$ is satisfiable. Furthermore, Σ' contains at most m distinct atoms, so by the induction hypothesis there exists a linear refutation of Σ' with top clause D' . See the Figure 1.4 for illustration.

Each side clause in this refutation that is not equal to a previous center clause, is either a member of Σ or is obtained from a member

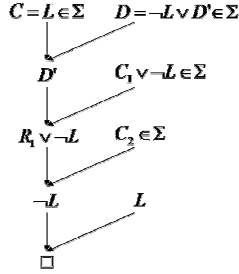


Figure 1.5: The complete picture of case 1 of the proof for theorem 10.

of Σ by means of the deletion of $\neg L$. In the latter kind of side clauses, put back the deleted $\neg L$ literals, and add these $\neg L$ to all later center clauses. Note that afterwards, these center clauses may contain multiple copies of $\neg L$. In particular, the last center clause changes from \square to $\neg L \vee \dots \vee \neg L$. Since D' is a resolvent of C and $D = \neg L \vee D' \in \Sigma$, we can add C and D as parent clauses on top of the previous top clause D' . That way, we get a linear derivation of $\neg L \vee \dots \vee \neg L$ from Σ , with top clause C . Finally, the literals in $\neg L \vee \dots \vee \neg L$ can be resolved away using the top clause $C = L$ as side clause. This yields a linear refutation of Σ with top clause C (see Figure 1.5).

- *Case 2:*

Suppose $C = L \vee C'$, where C' is a non-empty clause. C' cannot contain $\neg L$, for otherwise C would be a tautology, contradicting the assumption that Σ is unsatisfiable while $\Sigma \setminus \{C\}$ is satisfiable. Obtain Σ' from Σ by deleting clauses containing $\neg L$, and by removing the literal L from the remaining clauses. Note that $C' \in \Sigma'$. If M were a Herbrand model of Σ' , then $M \cup \{\neg L\}$ is a Herbrand model of Σ . Thus since Σ is unsatisfiable, Σ' is unsatisfiable.

Furthermore, because $\Sigma \setminus \{C\}$ is satisfiable, there is a Herbrand model M' of $\Sigma \setminus \{C\}$. Since Σ is unsatisfiable, M' is not a model of C . L is a literal in C , hence L must be false under M' . Every clause in $\Sigma' \setminus \{C'\}$ is obtained from a clause in $\Sigma \setminus C$ by deleting L from it. Since M' is a model of every clause in $\Sigma \setminus \{C\}$ and L is false under M' , every clause in $\Sigma' \setminus \{C'\}$ is true under M' . Therefore M' is a model of $\Sigma' \setminus \{C'\}$, which shows that $\Sigma' \setminus \{C'\}$ is satisfiable.

Then by the induction hypothesis, there exists a linear refutation of Σ' with top clause C' . Now similar to case 1, put back the previously deleted L literals to the top and side clauses, and to the appropriate center clauses. This gives a linear derivation of $L \vee \dots \vee L$ from Σ with top clause C .

Note that $\{L\} \cup (\Sigma \setminus \{C\})$ is unsatisfiable, because L is false in any Herbrand model of $\Sigma \setminus \{C\}$, as shown above. On the other hand, $\Sigma \setminus \{C\}$ is satisfiable. Thus by case 1 of this proof, there exists a linear refutation of $\{L\} \cup (\Sigma \setminus \{C\})$ with top clause L . Since L is a factor of $L \vee \dots \vee L$, we can put out linear derivation of $L \vee \dots \vee L$ “on top” of this linear refutation of $\{L\} \cup (\Sigma \setminus \{C\})$ with top clause L , thus obtaining a linear refutation of Σ with top clause C .

□

The incompleteness of input also means that the Subsumption Theorem will not hold for input resolution in general. What, then, can we say about SLD resolution? The short answer is that it too is incomplete. But, for a restricted form of clauses, input and SLD resolution *are* complete. The restriction is to Horn clauses: recall that these are clauses that have at most 1 positive literal. Indeed, it is this restriction that forms the basis of theorem-proving in the PROLOG language, which is restricted (at least in its pure form) to Horn clauses, albeit in first-order logic (but the result still holds in that case as well).

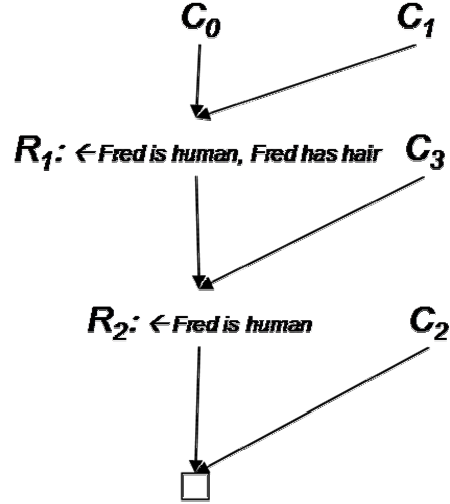
Proofs Using SLD Resolution

Before we get to SLD, we first make our description of input resolution a little more precise: the derivation of a clause C from a set of clauses Σ using input resolution means there is a sequence of clauses $R_0, \dots, R_k = C$ such that $R_0 \in \Sigma$ and each R_i ($1 \leq i \leq k$) is a resolvent of R_{i-1} and a clause $C_i \in \Sigma$. Now, we add further restrictions. Let Σ be a set of Horn clauses. Further, let R_i be a resolvent of a selected negative literal in R_{i-1} and the positive literal of a definite clause $C_i \in \Sigma$. The selection rule is called the “computation rule” and the resulting proof strategy is called SLD resolution (“SelecteLiner Definite” resolution). We illustrate this with an example. Let Σ be the set of clauses:

$C_0 : \neg \text{Fred is an ape}$
 $C_1 : \text{Fred is an ape} \leftarrow \text{Fred is human, Fred has hair}$
 $C_2 : \text{Fred is human}$
 $C_3 : \text{Fred has hair}$

A little thought should convince you that $\Sigma \models \square$. We want to see if $\Sigma \vdash_{SLD} \square$. It is evident that C_0 and C_1 resolve. The resolvent is R_1 :

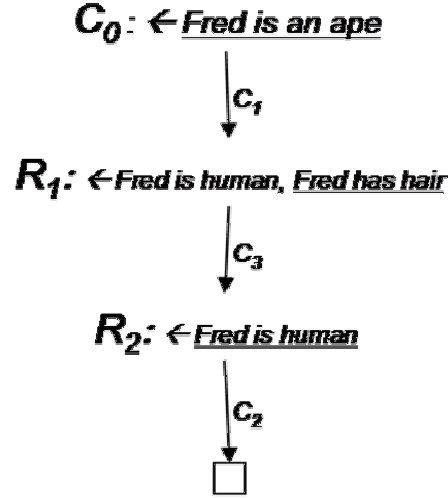
$C_0 : \neg \text{Fred is an ape}$
 $C_1 : \text{Fred is an ape} \leftarrow \text{Fred is human, Fred has hair}$
 $R_1 : \leftarrow \text{Fred is human, Fred has hair}$
 $C_2 : \text{Fred is human}$
 $C_3 : \text{Fred has hair}$

Figure 1.6: Example of SLD-deduction of \square from Σ .

Since we are using SLD, one of the resolvents for the next step has to be R_1 . The other resolvent has to be one of the C_i 's. Suppose our selection rule selects the “rightmost” literal first for resolution (that is, \neg Fred has hair in R_1). This resolves with C_3 , giving $R_2: \neg$ Fred is human, which in turn resolves with C_2 to give \square . The SLD (input) resolution diagram for this is presented in Figure 1.6.

It is more common, especially in the logic-programming literature, to present instead the search process confronting a SLD-resolution theorem prover in the form of a tree-diagram, called an *SLD-tree*. Such a tree effectively contains all possible derivations that can be obtained using a particular literal selection rule. Each node in the tree is a “goal” of the form $\leftarrow L_1, L_2, \dots, L_k$. That is, it is a clause of the form $(\neg L_1 \vee \neg L_2 \vee \dots \vee \neg L_k)$. Given a set of clauses Σ , the children of a node in the SLD-tree are the result of resolving with clauses in Σ (nodes representing the empty clause \square have no children). The SLD-tree for the example we just looked at is shown in Figure 1.7

We can now see what refutation-completeness for Horn clauses for SLD-resolution means in terms of SLD-trees. In effect, this means that if a set of clauses is unsatisfiable then there will be a leaf in the SLD-tree with the empty clause \square . Further, the computation rule will not alter this (informally, you can see that different computation rules will simply move the location of the \square around). We will have more to say on SLD-resolution with first-order logic in a later section. There, we will see that in addition to the computation rule, we will also need a “search” rule that determines how the SLD-tree is searched. Search

Figure 1.7: Example SLD-tree with C_0 at root.

trees there can have infinite branches, and although completeness is unaffected by the choice of the computation rule (that is, there will be a \square in the tree if the set of first-order clauses is unsatisfiable), we may not be able to reach it with a fixed search rule.

Theorem 11 *If Σ is an unsatisfiable set of horn clauses, then there is an SLD refutation of Σ .*

Proof:

We can assume Σ is finite. Let n be the number of facts (clauses consisting of a single positive literal) in Σ . We prove the lemma by induction on n

1. If $n = 0$, then $\square \in \Sigma$ for otherwise the empty set would be a Herbrand model of I .
2. Suppose the lemma holds for $0 \leq n \leq m$, and suppose $m + 1$ distinct facts appear in Σ . If $\square \in \Sigma$ the lemma is obvious, so suppose $\square \notin \Sigma$.

Let, A be a fact in Σ . We first delete all clauses from Σ ; which have A as head (so we also delete the fact A from Σ). Then we replace clauses which have A in their body by clauses constructed by deleting these atoms A from the body (so for example, $B \leftarrow A, B_1, \dots, B_k$ will be replaced by $B \leftarrow B_1, \dots, B_k$). Call the set obtained in this way Σ' . If M is a model of Σ' , then $M \cup \{A\}$ is a Herbrand model of Σ . Thus since Σ is unsatisfiable,

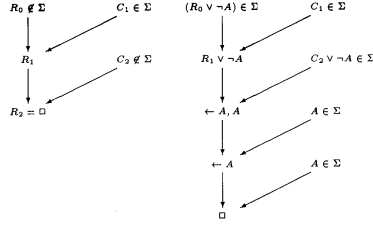


Figure 1.8: Illustration of the proof for theorem 11. The SLD refutation of Σ' is on the left and that for Σ is on the right.

Σ' must be unsatisfiable. Σ' only contains m facts, so by the induction hypothesis, there is an SLD-refutation of Σ' . If this refutation only uses clauses from Σ' which were also in Σ then this is also an SLD-refutation of Σ , so then we are done. Otherwise, if C is the top clause or an input clause in this refutation and $C \notin \Sigma$, then C was obtained from some $C' \in \Sigma$ by deleting all atoms A from the body of Σ . For all such C , do the following: restore the previously deleted copies of A to the body of C (which turns C into C' again), and add these atoms A to all later resolvents. This way, we can turn the SLD-refutation of Σ' into an SLD-derivation of $\leftarrow A, \dots, A$ from Σ . See Figure 1.8 for illustration, where we add previously deleted atoms A to the bodies of R_0 and C_2 . Since also $A \in \Sigma$, we can construct an SLD-refutation of Σ , using A a number of times as input clause to resolve away all members of the goal $\leftarrow A, \dots, A$.

□

1.3.5 Davis-Putnam Procedure

The inference problem addressed so far (particularly through the resolution procedure) is to determine if a proposition α logically follows from a given logical theory Σ . As we saw, this is achieved by reducing the problem to a coNP-complete¹⁵ unsatisfiability problem; based on the contradiction theorem, it amounts to negating the goal formula α , add it to the theory and test the conjunction for unsatisfiability.

However, often, one is faced with the requirement for a model M for a logical theory Σ . This can turn out to be easier problem than the usual problem of

¹⁵A decision problem C is Co-NP-complete if it is in Co-NP and if every problem in Co-NP is polynomial-time many-one reducible to it. The problem of determining whether a given boolean formula is tautology is a coNP-complete problem as well. A problem C is a member of co-NP if and only if its complement \overline{C} is in complexity class NP. For example, the satisfiability problem is an NP-complete problem. Therefore the unsatisfiability problem is a coNP-complete problem.

inference, since it is enough to find one model for the theory, as against trying all possible truth assignments as in the case of solving an unsatisfiability problem. For example, theory might describe constraints on the different parts of a car. And you are interested in a model that satisfies all the constraints. In terms of search, you search the space of assignment and stop when you find an assignment that satisfies the theory.

While the resolution procedure can be modified so that it gives you a model, the Davis-Putnam-Logemann-Loveland (DPLL) procedure is a more efficient procedure for solving SAT problems. Given a set of clauses Σ defined over a set of variables \mathcal{V} , the Davis-Putnam procedure $DPLL(\Sigma)$ returns ‘satisfiable’ if it is satisfiable. Otherwise return ‘unsatisfiable’.

The $DPLL(\Sigma)$ procedure consists of the following steps. The first two steps specify termination conditions. The last two rules actually work on the clauses in Σ .

1. If $\Sigma = \emptyset$, return ‘satisfiable’. This convention was introduced on pages 22 when we introduced logical entailment, as also in the footnote on page 29.
2. If $\square \in \Sigma$ return ‘unsatisfiable’. This convention was discussed in the footnote on page 29.
3. **Unit-propagation Rule:** If Σ contains any unit-clause $C = \{c\}$ (*c.f.* page 29 for definition), assign a truth-value to the variable in literal c that satisfies c , ‘simplify’ Σ to Σ' and (recursively) return $DPLL(\Sigma')$. The rationale here is that if Σ has any unit clause $C = \{c\}$, the only way to satisfy C is to make c true. The simplification of Σ to Σ' is achieved by:
 - (a) removing all clauses from Σ that contain the literal c (since all such clauses will now be true)
 - (b) removing the negation of literal c (*i.e.*, $\neg c$) from every clause in Σ that contains it
4. **Splitting Rule:** Select from \mathcal{V} , a variable v which has not been assigned a truth-value. Assign one truth value t to it, simplify Σ to Σ' and (recursively) call $DPLL(\Sigma')$.
 - (a) If the call returns ‘satisfiable’ (*i.e.*, we made a right choice for truth value of v), then return ‘satisfiable’.
 - (b) Otherwise (that is, if we made a wrong choice for truth value of v), assign the other truth-value to v in Σ , simplify to Σ'' and return $DPLL(\Sigma'')$.

The DPLL procedure can construct a model (if there exists one) by doing a book-keeping over all the assignments. This procedure is complete (that is, it constructs a model if there exists one), correct (the procedure always finds a truth assignment that is a model) and guaranteed to terminate (since the

space of possible assignments is finite and since DPLL explores that space systematically). In the worst case, DPLL requires exponential time, owing to the splitting rule. This is not surprising, given the NP-completeness of the SAT problem. Heuristics are needed to determine (i) which variable should be instantiated next and (ii) to what value the instantiated variable should be set. In all SAT competitions¹⁶ so far, Davis Putnam-based procedures have shown the best performance.

As an illustration, we will use the DPLL procedure to determine a model for $\Sigma = \{\{a, b, \neg c\}, \{\neg a, \neg b\}, \{c\}, \{a, \neg b\}\}$. Since Σ is neither an empty set nor contains the empty clause, we move on to step 3 of the procedure. Σ contains a single unit clause $\{c\}$, which we will set to true and simplify the theory to $\Sigma^1 = \{\{a, b\}, \{\neg a, \neg b\}, \{a, \neg b\}\}$. Next, we apply the splitting rule. Let us choose a and set it to ‘false’. This yields $\Sigma^2 = \{\{b\}, \{\neg b\}\}$. It can be verified that $\Sigma^2 \equiv \{\square\}$. We therefore backtrack and set a to ‘true’. This yields $\Sigma^3 = \{\{\neg b\}\}$. Thereafter, application of unit propagation yields $\Sigma^4 = \{\}$, which is satisfiable according to step 1 of the procedure. Thus, using the DPLL procedure, we obtain a model for Σ as $M = \{a, c\}$.

As another example, consider $\Sigma = \{\{a, \neg b, \neg c, \neg d\}, \{b, \neg d\}, \{c, \neg d\}, \{d\}\}$. As we will see, application to DPLL to this problem does not involve any backtracking, making the task relatively easier. Σ is neither an empty set nor contains the empty clause. Hence we apply the unit propagation rule 3 on $\{d\}$ to obtain $\Sigma^1 = \{\{a, \neg b, \neg c\}, \{b\}, \{c\}\}$. We can again apply unit propagation to Σ^1 on $\{b\}$ to obtain $\Sigma^2 = \{\{a, \neg c\}, \{c\}\}$. Finally, we can apply unit propagation to Σ^2 on $\{c\}$ and then to $\{a\}$ obtain $\Sigma^4 = \{\}$, which is satisfiable. This yields a model $M = \{d, b, c, a\}$ of Σ .

The DPLL procedure is similar to the traditional constraint propagation procedure. The splitting rule in DPLL is similar to the backtracking step in constraint propagation, while the unit propagation rule is similar to the unavoidable steps of consistency checking and forward propagation in traditional CP.

Davis Putnam on Horn Clauses

The Davis Putnam procedure is polynomial on horn clauses. Why this is so will become apparent through the following sequence of arguments:

1. Since the DPLL steps (1 through 4) could at most throw out a clause or throw out a literal within a clause, it is obvious that simplifications in DPLL on Horn clauses always generate Horn clauses.
2. A set of Horn clauses without unit clauses can be always satisfied by an assignment of ‘false’ to all variables (since all clauses will have at least one negative literal).
3. The first sequence of applications of the unit propagation rule in DPLL should either lead to the empty clause (which is satisfiable) or should lead

¹⁶<http://www.satcompetition.org/>

to a set of Horn clauses that has no unit clause (which is also satisfiable according to statement (2) above).

4. Thus, for Horn clauses, it suffices to apply steps 1 through 3 of the DPLL procedure. Since no backtracking is involved and since steps 1 through 3 can apply at most n times (n being the number of variables in the theory), for Horn clauses, DPLL is polynomial in the number of variables.
5. Even if step 4 of DPLL were applied in the case of Horn clauses, the time taken would still be polynomial in the number of variables. This is because, assigning a value ‘false’ to the variable chosen in step 4 cannot change the satisfiability, whereas, assigning a value ‘true’ can either lead to an immediate contradiction (after unit propagation) or will not affect satisfiability at all.

Phase Transitions

We saw that in the worst case, DPLL requires exponential time. Couldn’t we do better in the average case? For CNF-formulae in which the probability for a positive appearance, negative appearance and non-appearance in a clause is $1/3$, DP needs on average quadratic time [Gol79]. In retrospect, it was discovered that the formulae in [Gol79] have a very high probability of being satisfiable. Thus, these formulae are not representative of those encountered in practice. The idea of *phase transition* was conjectured by [CKT91] to identify hard to solve problem instances:

All NP-complete problems have at least one order parameter and the hard to solve problems are around a critical value of this order parameter. This critical value (a phase transition) separates one region from another, such as over-constrained and under-constrained regions of the problem space.

This conjecture was initially confirmed for the graph coloring and Hamilton path problems and later for other NP-complete problems, including SAT. In the case of SAT problem, an example of the order parameter is ratio of the number of variables to the number of clauses.

1. For higher settings of the parameter, the problem is over-constrained (the formulae are unsatisfiable). If the probability of a solution is close to 0, this fact can usually be determined early in the search.
2. For lower settings of the parameter, the problem is under-constrained (the formulae are easily satisfiable). When the probability of a solution is close to 1, there are many solutions, and the first search path of a backtracking search is usually successful.

When this parameter is varied, the problem moves from the over-constrained to the under-constrained region (or vice versa). At phase the transition points,

half of the problems are satisfiable and half are not. It is typically in this region that algorithms have difficulty in solving the problem¹⁷. Cook and Mitchell [CM97] empirically found that for a 3-SAT problem, the phase transition occurs at a clause:variables ratio of around 4.3 (in this experiment, clauses were generated by choosing variables for a clause and complementing each variable with probability 0.5). As an illustration, in the 2003 version of the SAT competition, the largest instances solved using greedy SAT solvers consisted of 100,000 variables and 1,000,000 clauses (clause:variable ratio of 10), whereas the smallest unsolved instances comprised 200 variables and 1,000 clauses (clause:variable ratio of 5). It was also reported in [CM97] that the runtime for the DPLL procedure peaks at the phase transition. In the phase transition region, the DPLL algorithm often near successes. Many benchmark problems are located in the phase transition region, though they have a special structure in addition.

1.3.6 Local Search Methods

Local search methods are standard search procedures for optimization problems. A local search method explores the neighborhood of the current solution and tries to enhance the solution till it cannot do any better. The hope is to produce better configurations through local modifications. The value of a configuration in a logical problem could be measured using the number of satisfied constraints/clauses. However, for logical problems, local maxima are inappropriate; it is required to satisfy all clauses in the theory and not just some. But through random restarts or by noise injection, local maxima can be escaped. In practice, local search performs quite well for finding satisfying assignments of CNF formulæ, especially for under-constrained or over-constrained SAT problems.

GSAT and WalkSat [SKC93] are two local search algorithms to solve boolean satisfiability problems in CNF. They start by assigning a random value to each variable. If the assignment satisfies all clauses, the algorithm terminates, returning the assignment. Otherwise, an unsatisfied clause is selected and the value of exactly one variable changed. Due to the conjunctive normal form, flipping one variable will result in that clause becoming satisfied. The above is then repeated until all the clauses are satisfied. WalkSAT and GSAT differ in the methods used to select the variable to flip. While GSAT makes the change which minimizes the number of unsatisfied clauses in the new assignment, WalkSAT selects the variable that, when flipped, results in no previously satisfied clauses becoming unsatisfied (some sort of downward compatibility requirement). MaxWalkSat is a variant of WalkSat designed to solve the weighted satisfiability problem, in which each clause is associated with a weight. The goal in MaxWalkSat is to find an assignment (which may or may not satisfy the entire formula) that maximizes the total weight of the clauses satisfied by that assignment. These algorithms

¹⁷Note that hard instances can also exist in regions of the more easily satisfiable or unsatisfiable instances.

perform very well on the randomly generated formulæ in the phase transition region. Monitoring the search procedure of these greedy solvers reveals that in the beginning, each procedure is very good at reducing the number of unsatisfied clauses. However, it takes a long time to satisfy the few remaining clauses (called plateaus). The GSAT algorithm is outlined in Figure 1.9.

<p>INPUT: A set of clauses Σ, MAX-FLIPS, and MAX-TRIES. OUTPUT: A satisfying truth assignment of Σ, if found. for $i = 1$ to MAX-TRIES do T = a randomly-generated truth assignment. for $j:=1$ to MAX-FLIPS do if T satisfies Σ then return T end if v = a propositional variable such that a change in its truth assignment gives the largest increase in the number of clauses of Σ that are satisfied by T. $T = T$ with the truth assignment of v reversed. end for end for return “Unsatisfiable”.</p>
--

Figure 1.9: Procedure GSAT.

1.3.7 Default inference under closed world assumption

Given any set of formulae Σ , the closed-world assumption is the assumption that Σ determines all the knowledge there is to be had about the formulae in the language. Thus, if we consider any proposition¹⁸ A then A is taken to be true exactly when Σ (logically) implies A , but is otherwise taken to be false.

The closed-world assumption underlies the mode of reasoning known as *default inference*. There are many situations, both in ordinary daily life and in specific technical computing matters (such as explaining the theory of finite failure and the relationship it bears to reasoning with negative information), when default inference is a necessary supplement to deductive inference. Consider this simple example of a single clause axiom $\Sigma = \{A \leftarrow B\}$ for which $B(\Sigma)$ is just $\{A, B\}$. Under the closed-world assumption, we may infer $\neg A$ for any ground atom $A \in B(\Sigma)$ that is not implied by Σ . This is a constrained¹⁹ application of the rule of default inference:

¹⁸In the first order logic programming context, the propositions in which we are primarily interested are the atoms of the Herbrand base.

¹⁹Constrained because (i) Σ is assumed to be inconsistent, else Σ would necessarily imply both A and $\neg A$ (ii) only the case where A is atomic is considered, otherwise if Σ implied, say, neither A nor $\neg A$, then the default rule would infer both $\neg A$ and A , which would again be inconsistent.

Infer $\neg A$ in default of Σ implying A

Motivated by the desire to draw sound conclusions about negative information, we will consider two constructions that provide consequence-oriented meaning for default inference under the closed-world assumption.

1. $CWA(\Sigma)$: The combination of Σ with the default conclusions inferred from it is denoted by $CWA(\Sigma)$ and is defined by

$$CWA(\Sigma) = \Sigma \cup \{\neg A \mid A \in B(\Sigma) \text{ and not } \Sigma \models A\}$$

For the above example, $CWA(\Sigma) = \{A \leftarrow B, \neg A, \neg B\}$.

Soundness for the default inference of negative conclusions under the closed-world assumption can be referred to the the logical construction $CWA(\Sigma)$ as follows:

$$\text{for all } A \in B(\Sigma), CWA(\Sigma) \models \neg A \text{ if } \Sigma \vdash_{CWA} \neg A$$

There some practical problems with CWA:

- (a) One problem is that, we have no immediate way of writing down $CWA(\Sigma)$, for we cannot directly discern which particular negative facts $\neg A$ to include in it; we would have to infer them all first.
 - (b) A more serious defect is that, whilst $CWA(\Sigma)$ is always consistent when Σ is definite²⁰, it is likely to be inconsistent otherwise. For example, if Σ comprises just the clause $\{A \vee B\}$, then $CWA(\Sigma) = \{A \vee B, \neg A, \neg B\}$ which is inconsistent.
2. $CWA(\Sigma)$: For default inference applied to indefinite programs, we refer the soundness criterion to *completion* $COMP(\Sigma)$. It is also known as the completed database of Σ . Unlike $CWA(\Sigma)$, it can be written down more-or-less directly and is consistent for all well-structured programs. In the case of propositional logic, the construction is particularly simple:
 - (a) Initialize $COMP(\Sigma) = \emptyset$.
 - (b) Assume that Σ is any set of clauses of the form $(A \leftarrow \text{body})$.
 - (c) For each A mentioned in Σ but not defined in Σ , construct $\neg A$ and add it to $COMP(\Sigma)$.
 - (d) For each A having a definition in Σ of the form

$$\begin{aligned} A &\leftarrow \text{body} - 1 \\ &\dots \\ &\dots \\ A &\leftarrow \text{body} - n \end{aligned}$$

²⁰Exercise: Prove.

construct the clause A iff $(body - 1 \vee \dots \vee body - n)$ and add it to $COMP(\Sigma)$,

$Comp(\Sigma)$ can also be viewed as the result of simplifying $\Sigma \cup \text{only-if}(\Sigma)$, where $\text{only-if}(\Sigma)$ comprises every completed definition of the form $\neg A$ for $A \notin \Sigma$ as well as $(q \rightarrow \alpha)$ for every completed definition $(q \text{ iff } \alpha) \in COMP(\Sigma)$.

Following are some characteristics of completions

- (a) In general it is also a more economical construction in the sense that it implies, but does not necessarily declare, various negative propositions which $CWA(\Sigma)$ would have to declare explicitly. As an example, if $\Sigma = \{A \leftarrow B\}$, $COMP(\Sigma) = \{\neg B, A \text{ iff } B\}$. While Σ neither implies A nor implies B , $Comp(P)$ implies both $\neg A$ and $\neg B$, which is exactly the same outcome obtained using $CWA(\Sigma)$ instead. Observe, however, that $CWA(\Sigma)$ declares $\neg A$ **explicitly** whereas $Comp(\Sigma)$ does not.
- (b) Generally, completion is more conservative than the closed-world assumption in the negative facts that it implies. For example, with $\Sigma_1 = \{A \leftarrow A\}$, $CWA(\Sigma_1) \models \neg A$, whereas $COMP(\Sigma_1) \not\models \neg A$. Similarly, with $\Sigma_2 = \{A \leftarrow \neg B\}$, $CWA(\Sigma_2) \models \neg A \wedge \neg B$ (is inconsistent), whereas $COMP(\Sigma_2) \not\models \neg A$ though $COMP(\Sigma_2) \models \neg B$.
- (c) On the other hand, $COMP$ can be sometimes less conservative; for $\Sigma_3 = \{A \leftarrow \neg A\}$, $COMP(\Sigma_3) \models \text{everything}$, whereas $CWA(\Sigma_3) \models A$ and $CWA(\Sigma_3) \not\models \neg A$.
- (d) Completed definitions capture the programmer's intentions more fully than do uncompleted definitions in the original program. For computational purposes however, the program alone happens to be sufficient for deducing all intended answers.
- (e) In cases where Σ is indefinite, the elementary consequences of $COMP(\Sigma)$ may arise from the joint contributions of Σ and $\text{only-if}(\Sigma)$ and not from either of them alone.
- (f) There is a syntactic bias built into the process of program completion. For example, if $\Sigma_1 = \{A \leftarrow \neg B\}$ and $\Sigma_2 = \{B \leftarrow \neg A\}$ their respective completions are $COMP(\Sigma_1) = \{A \text{ iff } \neg B, \neg B\}$ and $COMP(\Sigma_2) = \{B \text{ iff } \neg A, \neg A\}$. These completions are not logically equivalent despite the fact that the original programs are. The difference between their completions deliberately reflects the difference between the procedural intentions suggested by the programs' syntaxes: the first program anticipates queries of the form $?A$ whilst the second anticipates queries of the form $?B$.

Finite Failure Extension

We look at one further extension of SLD-resolution that is of special interest to us, namely, the idea of negation as "failure to prove". **We know that SLD**

alone is able to deal soundly and completely with all positive atomic queries A posed to this database. But with the finite failure facility we can now also ask directly, in our extended language, whether some atom is *not* in the database.

This extension enables one to express, via a call ‘*fail A*’, the condition that a call A shall finitely fail. In Prolog the fail operator is denoted by ‘*not*’ and is referred to as ‘negation by failure’. The operational meaning of ‘*fail*’ is summed up by the finite failure rule:

A call ‘*fail A*’ succeeds *iff* its subcall A finitely fails.

This rule can be incorporated directly into the execution strategy for virtually no cost in terms of implementation overheads - to evaluate a call ‘*fail A*’, the interpreter merely evaluates A in the standard way and then responds to the outcome as just prescribed. The phrase ‘ A finitely fails’ means that the execution tree generated by the evaluation of A must be a finitely failed tree - that is, it must have finite depth, finite breadth and all the computations contained within it must terminate with failure. By contrast, should one or more of those computations terminate with success then the call ‘*fail A*’ is itself deemed to have finitely failed. Finally, if the evaluation of A neither succeeds nor finitely fails, then the same holds for the evaluation of ‘*fail A*’

The rule of “negation as finite failure” states that if all branches of an SLD-tree finitely fail for $\Sigma \cup \{\leftarrow A\}$ for a set of definite clauses Σ and a ground literal A , we can derive the ground literal $\neg A$ as a result. This combination of SLD-resolution and negation-as-failure results in the proof strategy we called SLDNF-resolution. Once again, let us look at an example:

C_0 : \neg Fred is an ape
 C_1 : Fred is an ape \leftarrow *not* Fred is human, Fred is a primate
 C_3 : Fred is a primate

Here, *not* stands for “not provable” (which is not the same as \neg). C_0 and C_1 resolve to give R_1 : \leftarrow *not* Fred is human, Fred is a primate. With a rightmost literal computation rule as before, the next resolvent is R_2 : \leftarrow *not* Fred is human. It is evident that Fred is human is not provable and \square results.

SLDNF is sound and complete for propositional definite clauses. However, there are some important issues in extending SLD with finite failure to first order logic such as (i) incompleteness when applied to activated non-ground fail calls and (ii) unsoundness in certain cases. To solve the latter, more serious problem, a selection policy called *safe computation rule* is used in practice; this policy sacrifices completeness for the sake of soundness. These and other concepts such as *floundering* will be discussed in a later section.

The SLD finite failure set

Relative to a given Σ and a given inference system R , the SLD finite failure $FF(\Sigma, R)$ set comprises exactly propositional atoms for which the query $?q$

finitely fails. We shall restrict our attention throughout to the case where Σ is definite and R is SLD. Once we know what the finite failure set is relative to SLD inference, we shall also be able to say something about the execution by SLDNF of queries containing ‘fail calls’ - subject, of course, to the assumption that the chosen computation rule is safe. The more general situation where ‘fail’ calls may occur also in clause bodies, is significantly more complicated, and will not be addressed here. For queries of the form $?q$ using a definite Σ under some SLD computation rule R partitions $B(\Sigma)$ into three species of atoms

1. $SS(\Sigma, R)$, or those for which $?q$ succeeds. Since, if $?q$ succeeded in one SLD tree, it succeeded in all SLD trees, the capacity for $?q$ to succeed is independent of the computation rule. So this set can be simply denoted by $SS(\Sigma)$.
2. $FF(\Sigma, R)$, or those for which $?q$ finitely fails
3. $IF(\Sigma, R)$, or those for which $?q$ infinitely fails.

However, $?q$ may finitely fail under one computation rule yet infinitely fail under another. A simple example is where Σ comprises just the clause $q \leftarrow B \wedge q$. Under the standard leftmost call rule, the query $?q$ finitely fails, whereas it infinitely fails under a rule which always selects the rightmost call. Thus the boundary between $FF(\Sigma, R)$ and $IF(\Sigma, R)$ depends upon R . There is a simple theoretical way of eliminating this dependence upon R . We invoke the idea of a particular sort of computation rule which ensures that any call introduced into a computation is selected after some arbitrary but finite number of execution steps. Such a rule is called a *fair computation rule*. Considering the example again, we might allow a fair computation rule to select ‘ q ’ calls any finite number of times, but sooner or later the fairness requirement would demand that an ‘ A ’ call be selected, thus immediately forcing finite failure. With this new concept, we can now state the following facts: $?q$ finitely fails under some computation rule *iff* it finitely fails under all fair computation rules²¹ We can denote by $FF(\Sigma)$, the set of all atoms q for which $?q$ has some finitely failed fair SLD-tree, and by $IF(\Sigma)$, the set of all atoms q for which $?q$ has some infinitely failed fair SLD-tree. $FF(\Sigma)$ is called the (fair-) SLD finite failure set of Σ . Referring again to the example above, we shall have $FF(\Sigma) = \{q, A\}$ and $IF(\Sigma) = \emptyset$.

Suppose queries were also allowed to contain ‘*fail*’ calls besides atomic ones, but with Σ still restricted to be definite. The following statements then hold true:

$$\begin{aligned} \text{for all } q \in B(\Sigma), \quad q \in FF(\Sigma) & \text{ iff } ?fail \text{ } q \text{ succeeds under SLDNF} \\ q \in SS(\Sigma) & \text{ iff } ?q \text{ succeeds under SLDNF} \end{aligned}$$

We can interpret finite failure (i) in terms of the position of $FF(\Sigma)$ within the lattice of interpretations and (ii) in terms of the classical negation (\neg), which

²¹Although fair computation rules are not normally implemented, they certainly tidy up our mathematical account of finite failure.; equivalently, at least one of its SLD-trees is finitely failed *iff* all of its fair SLD-trees are finitely failed.

provides a semantics for ‘fail’, and some soundness and completeness results for SLDNF, based upon the logical consequences of $\text{COMP}(\Sigma)$.

Completion Semantics for SLDNF

For pure Horn-clause programs and queries, we had a particularly simple connection between logical meaning and operational meaning:

$$\text{for all } q \in B(\Sigma), \Sigma \models q \text{ iff } q \in SS(\Sigma)$$

The construction of a consequence-oriented semantics for the finite failure extension is more problematic. A program containing fail is not a construct of classical logic and so is not amenable to the notion of classical logical consequence. Nevertheless, a variety of analogous connections have been suggested. The best-known of these is based upon the so-called *completion semantics*, which relies upon two ideas:

1. Interpreting ‘fail’ as the classical negation connective \neg (this is why fail is commonly referred to as ‘negation by failure’).
2. Relating the success or finite failure of calls to logical consequences of $\text{COMP}(\Sigma)$ rather than of Σ alone.

On this basis we can then characterize, in logical terms, the *soundness* of execution of atomic queries under the (safe-)SLDNF implementation of fail:

$$\begin{array}{lll} \text{for all } q \in B(\Sigma), & \text{COMP}(\Sigma) \models q & \text{if } ?q \text{ succeeds under SLDNF} \\ & \text{COMP}(\Sigma) \models \neg q & \text{if } ? \text{fail } q \text{ succeeds under SLDNF} \end{array}$$

These results hold even if ‘fail;’ calls occur in Σ . The ‘only-if’ part in the two statements above, which characterizes the *completeness* for SLDNF holds only when Σ is definite (so Σ will obviously not contain ‘fail’ calls).

There are two caveats in the above statement(s): (i) $\text{COMP}(\Sigma)$ may itself be inconsistent and (ii) the SLDNF uses some safe computation rule. Also, the simplification of analysis of finite failure using *fair* computation rules is not possible when the query contains ‘fail’ calls.

1.3.8 Lattice of Models

For any Σ that is a set of definite clausal formula, it can be shown that the set $\mathcal{M}(\Sigma)$ of models is a complete lattice ordered by set-inclusion (that is, for models $M1, M2 \in \mathcal{M}(\Sigma)$, $M1 \preceq M2$ if and only if $M1 \subseteq M2$), and with the binary operations of \cap and \cup as the glb and lub respectively. Recall that a complete lattice has a unique least upper bound and a unique greatest lower bound. Since we are really talking about sets that are ordered by set-inclusion, this means that there is a unique *smallest* one (the size being measured by the number of elements). This is called the *minimal model* of the formula, and it can be shown that this must be the intersection of all models for the formula. We will illustrate with an example. Consider a Σ consisting of the following clauses:

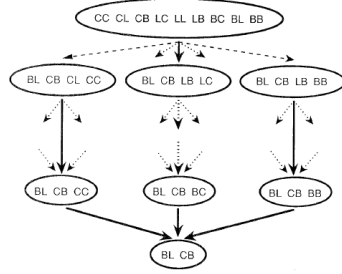


Figure 1.10: A complete lattice of models.

$$\begin{aligned}
 C_0 &: CC \leftarrow CL \\
 C_1 &: CB \leftarrow \neg BL \\
 C_2 &: CL \leftarrow LL \\
 C_3 &: BL
 \end{aligned}$$

A fail-safe way of getting a model is to choose the entire set of propositions $B(\Sigma) = \{CC, CL, CB, LC, LL, LB, BC, BL, BB\}$, which is called the *base* of Σ . For then every atom in Σ is assigned true and this in turn makes all of its clauses true. This model is the (unique) maximal model for Σ . However, such a choice is clearly excessive—many of the atoms in $B(\Sigma)$ do not even occur in $G(\Sigma)$ and so their truth values are irrelevant. Stripping out those irrelevant atoms leaves a somewhat smaller model $\{CC, CL, CB, BL, LL\}$. Which of these atoms must appear in any model? Clearly BL must, in order to make the program's fourth clause true. Then, since BL must, so also must CB in order to make the second clause true. The first and third clauses employ only the remaining three atoms CC, CL and LL , and both those clauses can be made true by making those atoms false. In conclusion, only BL and CB *must* be true – thus the (unique) minimal model for Σ is $\{BL, CB\}$.

Some indication of the complete lattice of models for the program is shown in Figure 1.10 where, edges stand for the *covers* relationship. Moreover, only a few of the models are shown—there are 64 models in the entire lattice. Note also that, in general, most subsets of $B(\Sigma)$ will be counter-models—for our example there are 448 of these, of which the smallest is \emptyset . It can be proved that if Σ is a set of definite clauses, it must always be satisfiable, and therefore, $\emptyset \notin \mathcal{M}(\Sigma)$.

The reason why the discussion above has focused solely upon definite clauses is that, in general, a set of indefinite clauses does not yield a complete lattice and may therefore have multiple minimal models. The lack of a unique minimal model makes it harder to assign an unambiguous meaning to such a set.

There is an important result relating a set of definite clausal formulae Σ , its minimal model $\mathcal{MM}(\Sigma)$ and the atoms that are logical consequences of Σ :

Theorem 12 If α is an atom then $\Sigma \models \alpha$ if and only if $\alpha \in \mathcal{MM}(\Sigma)$.

The proof of theorem 12 makes use of the so-called model intersection property:

Theorem 13 *If M_1, \dots, M_n are any models for a definite clause set Σ then their intersection is a model of Σ .*

Proof: It is easy to show this by induction. For any $k < n$, let I_k denote $M_1 \cap \dots \cap M_k$. Now consider any clause $C \in \Sigma$:

$$C : A \leftarrow B_1 \wedge \dots \wedge B_m$$

We shall prove that, for all $k \leq n$, I_k satisfies C .

1. *Base case ($k = 1$):* $I_1 = M_1$ and therefore satisfies C .
2. *Induction step ($1 \leq k \leq n$):* Assume that M_k satisfies C ;

if I_k satisfies C then	either	$B_i \notin I_k$ for some i
	or	$B_i \in I_k$ for all i , and $A \in I_k$
if M_{k+1} satisfies C then	either	$B_i \notin M_{k+1}$ for some i
	or	$B_i \in M_{k+1}$ for all i , and $A \in M_{k+1}$
it then follows that	either	$B_i \notin I_k \cap M_{k+1}$ for some i
	OR	$B_i \in I_k \cap M_{k+1}$ for all i , and $A \in I_k \cap M_{k+1}$
And hence		I_{k+1} satisfies C

Thus for all $k \leq n$, I_k satisfies C and-by a similar argument – every other clause in Σ .

□

We next prove theorem 12 making use of theorem 13.

Proof of theorem 12: EXERCISE

We will first prove that the intersection I^* of all the models of Σ is the minimal model $MM(\Sigma)$, using contradiction as follows:

Suppose I^* is not the minimal model
 then there must exist some model M_j such that $M_j \subset I^*$
 then there must exist some atom q such that $q \notin M_j$ and $q \in I^*$
 but $q \in I^*$ implies that $q \in M_i$, for all i , contradicting $q \in M_j$
 therefore the initial assumption is false.

Using this result it is now easy to prove the relationship between $MM(\Sigma)$ and any atomic consequence q of Σ :

If $\Sigma \models q$	then	q is true in every model of Σ
	then	$q \in I^*$
	then	$q \in MM(\Sigma)$
if $q \in MM(\Sigma)$	then	q is true in every model of Σ
	then	$\Sigma \models q$

□

Thus, the minimal model of a definite clausal formula is identical to the set of all ground atoms logically implied by that formula, which was defined as the success set $SS(\Sigma)$. Thus, the minimal model provides, in effect, the meaning (or semantics) of the formula. We shall see later that this is just one of several ways of giving significance to a program's minimal model.

We can envisage a procedure for enumerating the models of a formula. Consider the powerset of the base $B(\Sigma)$. Now, we know that this powerset ordered by \subseteq necessarily forms a complete lattice, with binary operations \cap and \cup . Some subset of this powerset is the set of all models, which we know is also a lattice ordered by \subseteq with the same binary operations. So, the model lattice is a sublattice of the lattice obtained from the powerset of the base. Suppose now we start at some point s in this sublattice, and we move to a new point that consists only of those atoms of the formula made true by the model s . Let us call these atoms s_1 . Then, a little thought should convince you that s_1 is also a member of the sublattice of models. Repeating the process with s_2 we can move to models s_2, s_3 and so on. Will this procedure converge eventually on the minimal model? Not necessarily, since we could end up moving back-and-forth between points of the sub-lattice. (When will this happen, and how can we ensure that we do converge on the minimal model?).

A slightly more general process can be formalised as the application of a function T_Σ that, for a set of clauses Σ , generates an interpretation (not necessarily a model) from another. That is:

$$I_{k+1} = T_\Sigma(I_k)$$

where

$$T_\Sigma(I) = \{a : a \leftarrow \text{body} \in \Sigma \text{ and } \text{body} \in I\}$$

It can be shown that T_Σ is both monotonic and continuous on the complete lattice obtained by ordering the powerset of the base by \subseteq . So, we know from the Knaster-Tarski Theorem mentioned on page 11, that there must be a least fixpoint for T_Σ in this lattice. We can prove that the procedure of obtaining I_{k+1} from application of T_Σ to I_k will yield that fixpoint, and further, that this fixpoint will be the minimal model.

Theorem 14 *$MM(\Sigma)$ is the least fixpoint of T_Σ .*

Proof: The definition of T_Σ requires that

for all $I \subseteq B(\Sigma)$ and for all $q, q \in T_\Sigma(I)$ iff $(\exists \text{body})[(q \leftarrow \text{body}) \in \Sigma \text{ and } \text{body} \in I]$

whereas, the definition of a model requires that

for all $I \subseteq B(\Sigma)$, I is a model for Σ iff for all $q, q \in I$ if $(\exists \text{body})[(q \leftarrow \text{body}) \in \Sigma \text{ and } \text{body} \in I]$.

These two sentences jointly imply

$$\begin{aligned}
& \text{for all } I \subseteq B(\Sigma) \\
I \text{ is a model for } \Sigma & \text{ iff } \text{for all } q, q \in I \text{ if } q \in T_\Sigma(I) \\
& \text{iff } T_\Sigma(I) \subseteq I \\
& \text{iff } I \text{ is a pre-fixpoint of } T_\Sigma
\end{aligned}$$

Where, for a function f on $\langle S, \preceq \rangle$, an element $u \in S$ is a pre-fixpoint of f if and only if $f(u) \preceq u$. Then, we can recall from the proof of the Knaster-Tarski theorem on page 11 that the least fix point is also the least pre-fixpoint. Hence, since the models are exactly the pre-fixpoints, the least model $MM(\Sigma)$, which is the glb of all the pre-fix points, must be the least pre-fixpoint. Finally, by the Knaster-Tarski theorem, $MM(\Sigma)$ must also be the least fixpoint. \square

So we now have several equivalent characterizations of the minimal model, including the success set $SS(\Sigma)$ and the new one posed in terms of the least fixpoint of T_Σ , denoted $LFP(T_\Sigma)$

Mathematical Characterization of Finite Failure

It so happens that there is a somewhat related (to the T_Σ function) method of constructing the finite failure set $FF(\Sigma)$. The set of atoms which occur as headings in Σ is simply $T_\Sigma(B(\Sigma))$, as is plain from the definition of T_Σ . The simplest way for $?q$ to fail finitely is for there to be no clause in Σ who heading unifies with q . In this case the failure is said to occur within depth $k = 1$ and the set of all such atoms $q \in B(\Sigma)$ is denoted by $FF(\Sigma, 1)$ and can be characterized very easily using the T_Σ function as

$$FF(\Sigma, 1) = B(\Sigma) - T_\Sigma(B(\Sigma))$$

Generalizing this principle, $FF(\Sigma)$ just contains each atom q for which $?q$ finitely fails within some depth $k \in \mathcal{N}$.

$$FF(\Sigma) = \cup_{k \in \mathcal{N}} FF(\Sigma, k) = B(\Sigma) - \cap_{k \in \mathcal{N}} T_\Sigma^k(B(\Sigma))$$

When we start at the top element $B(\Sigma)$ of our lattice of interpretations, repeated application of the T_Σ function generates a monotonically decreasing sequence $B(\Sigma) \supseteq T_\Sigma(B(\Sigma)) \supseteq T_\Sigma^2(B(\Sigma)) \dots$, whose limit is the greatest lower bound (glb) $T_\Sigma \downarrow$ of $\{T_\Sigma^k(B(\Sigma)) \mid k \in \mathcal{N}\}$. Thus, the mathematical characterization of $FF(\Sigma)$ (independent of the execution mechanism) is

$$FF(\Sigma) = B(\Sigma) - T_\Sigma \downarrow$$

Based on this equivalence, what is the value of $FF(\Sigma)$ for the last example that we discussed?

It can be proved that $T_\Sigma \uparrow = LFP(T_\Sigma) = MM(\Sigma) \subseteq T_\Sigma \downarrow$. There is an assymetry in the relationship between the limits and the extremal fixpoints of T_Σ . Whereas, $T_\Sigma \uparrow$ is always equal to the least fix point $LFP(T_\Sigma)$, $T_\Sigma \downarrow$ does not always equal the greatest fix point $GFP(T_\Sigma)$, though it does hold for most 'sensible' Σ (at the least for non-recursive clauses and in the case of

first order logic, for function-free programs). The region between $T_\Sigma \uparrow$ and $T_\Sigma \downarrow$ corresponds to $IF(\Sigma)$, which comprises exactly those atoms which fail infinitely. In practice, most sensible programs have $IF(\Sigma) = \emptyset$, leading to $T_\Sigma \uparrow = T_\Sigma \downarrow$ and therefore $B(\Sigma) = SS(\Sigma) \cup FF(\Sigma)$.

1.4 First-Order Logic

Suppose you wanted to express logically the statement: ‘All humans are apes.’ One of two ways can be used to formalise this in propositional logic. We can use a single proposition that stands for the entire statement, or with a well-formed formula consisting of a lot of conjunctions: **Human1 is an ape** \wedge **Human2 is an ape** . . . Using a single proposition does not give any indication of the structure inherent in the statement (that, for example, it is a statement about two sets of objects—humans and apes—one of which is entirely contained in the other). The conjunctive expression is clearly tedious in a world with a lot of humans. Things can get worse. Consider the following argument:

Some animals are humans.
All humans are apes.
Therefore some animals are apes.

That the argument is valid is evident: yet it is beyond the power of propositional logic to establish it. If, for example, we elected to represent each of the statements with single propositions then all we would end up with is:

<u>Statement</u>	<u>Formally</u>
Some animals are humans.	P
All humans are apes.	Q
Therefore some animals are apes.	$\therefore R$

But the formal argument is clearly invalid, as it is easy to think up arguments where P, Q are *true* and R is *false*. What is needed is in fact something along the following lines:

<u>Statement</u>	<u>Formally</u>
Some animals are humans.	Some P are Q
All humans are apes.	All Q are R
Therefore some animals are apes.	\therefore some P are R

Here, P, Q , and R do not stand for propositions, but for terms like *animals*, *humans* and *apes*. The use of terms like these related to each other by the expressions ‘some’ and ‘all’ will allow us to form sentences like the following:

All P are Q
 No P are Q
 Some P are Q
 Some P are not Q

The expressions ‘some’ and ‘all’ are called *quantifiers*, which when combined with the logical connectives introduced in connection with proposition logic ($\neg, \wedge, \vee, \leftrightarrow$), results in the powerful framework of first-order or *predicate logic*.

1.4.1 Syntax

The language of predicate logic introduces many new constructs that are not found in the simpler, propositional case. We will first introduce these informally.

Constants. It is conventional in predicate logic to use lowercase letters to denote proper names of objects. For example, in the sentence ‘Fred is human’, Fred could be represented as *fred*.

Variables. Consider the statements:

All humans are apes
 Some apes are not human

Using the letter x as a variable that can stand for individual objects, these can be expressed as:

For all x , if x is human then x is an ape
 For some x , x is an ape and x is not human

Quantifiers. The language of predicate logic introduces the symbol \forall , called the *universal quantifier*, to denote ‘for all.’ The symbol \exists , called the *existential quantifier*, is used to denote ‘for some’ or, more precisely, ‘for at least one.’ The sentences above can therefore be written as:

$\forall x$ (if x is human then x is an ape)
 $\exists x$ (x is an ape and x is not human)

Predicates. In their simplest case, these are symbols used to attribute properties to particular objects. It is conventional in logic (but ungrammatical in English) to write the subject after the predicate. Thus the sentence ‘Fred is human’ would be formalised as $Human(fred)$ ²². More generally, predicate symbols can be used to represent relations between two or more objects. Thus,

²²Logicians are a parsimonious lot: they would represent ‘Fred is human’ as Hf . The representation here is non-standard, but preferred for clarity.

‘Fred likes bananas’ can be represented as: $Likes(fred, bananas)$. The general form is therefore a predicate symbol, followed by one or more *arguments* separated by commas and enclosed by brackets. The number of arguments is sometimes called the *arity* of the predicate symbol, and the predicate symbol is often written along with its arity (for example, $Likes/2$). Formalising sentences like those above would result in quantified variables being arguments:

$$\begin{aligned} \forall x (\text{if } Human(x) \text{ then } Ape(x)) \\ \exists x (Ape(x) \text{ and not } Human(x)) \end{aligned}$$

Or, using the logical connectives that we have already come across:

$$\begin{aligned} \forall x (Ape(x) \leftarrow Human(x)) \\ \exists x (Ape(x) \wedge \neg Human(x)) \end{aligned}$$

Functions. Consider the statement: ‘The father of Fred is human.’ Although we have not named Fred’s father, it is evident that a unique individual is being referred to, and it possible to denote him by using a *function* symbol. One way to formalise the statement is: $Human(father(fred))$. Here, it is understood that $father(fred)$ denotes Fred’s father. A function symbol is one which, when attached to one or more terms denoting objects produces an expression that denotes a single object. It is important that that the result is unique: a function symbol could not be used to represent, for example, ‘parent of Fred.’ As with predicates, the number of arguments of the function is sometimes called its arity.

The following points would not be evident from this informal presentation:

1. Variables need not designate different objects. Thus, in $\forall x \forall y Likes(x, y)$, x and y could refer to the same object;
2. The choice of variable names is unimportant. Thus, $\forall x \forall y Likes(x, y)$ has the same meaning as $\forall y \forall z Likes(y, z)$;
3. The same variable name, if quantified differently, need not designate the same object. Thus, in $\forall x \forall y Likes(x, y) \wedge \forall x \forall y Hates(y, x)$ the x, y in $Likes(\dots)$ need not be same as the x, y in $Hates(\dots)$;
4. The order of quantifiers can matter when \forall and \exists are mixed. Thus, $\exists x \forall y Likes(x, y)$ has a different meaning to $\forall y \exists x Likes(x, y)$. However changing the order has no effect if the quantifiers are all of the same type. Thus, $\forall x \forall y Likes(x, y)$ has the same meaning as $\forall y \forall x Likes(x, y)$;
5. “Free” variables in a formula are those that are not quantified. For example, in the formula $\forall x Likes(x, y)$, y is a free variable. In contrast, quantified variables are called “bound” variables. It may not be immediately apparent that a variable can have both free and bound occurrences in a formula. For example in $\exists x (Likes(x, y) \wedge \exists y DisLikes(y, x))$, the variable y is free in the $Likes$ and bound in $Dislikes$. x on the other hand is

bound in both (by the outermost quantifier). It is normal to call a formula with no free variables a *sentence*, and it only really makes sense to ask about the truth of sentences;

6. Negation should be treated with caution. Thus, in ‘Some apes are not humans’, ‘not’ plays the role of *complementation*, by stating that the set of apes and the set of non-humans have at least one member in common. This can be formalised as $\exists x(Ape(x) \wedge \neg Human(x))$. On the other hand, ‘not’ plays the role of true negation in ‘It is not true that some apes are humans’ formalised as $\neg \exists x(Ape(x) \wedge Human(x))$;
7. It can be tricky to match English sentences to ones that use \forall and \exists . Thus, in ‘If something has a tail then it is not an ape’, the use of ‘something’ suggests that formalisation would involve \exists . The statement is, in fact, a general one about apes not having tails, and involves universal quantification: $\forall x(\neg Ape(x) \leftarrow Tail(x))$;
8. By denoting ‘at least one’, the existential quantifier \exists includes ‘exactly one’ and ‘all’. This does not coincide exactly with the usual English notion of ‘some’, which denotes more than one, but less than all.

We can now examine the formal rules for constructing well-formed formulae in predicate logic. For the language of predicate logic, we will restrict the vocabulary to the following:

Constant symbols:	A string of one or more lowercase letters (except those denoting variables)
Variable symbols:	A lowercase letter (except those denoting constants)
Predicate symbols:	Uppercase letter, followed by zero or more letters
Function symbols:	Lowercase letter, followed by zero or more letters (except those denoting constants or variables)
Quantifier symbols:	\forall, \exists
Logical connectives:	$\neg, \wedge, \vee, \leftarrow$
Brackets:	(,)

In addition, we will sometimes employ the device of using subscripts to denote unique symbols (for example, x_1, x_2, \dots for a string of variables).

With this vocabulary, a *term* is simply a constant, variable or a functional expression (that is, a function applied to a tuple of terms). The following are all examples of terms: x , $fred$, $father(fred)$, $father(father(fred))$, $father(x)$.

These, however, are not terms: $Likes(fred, bananas)$, $Likes(fred, father(fred))$, $father(Likes(fred, bananas))$. An *atomic formula*, sometimes simply called an *atom* is a predicate symbol applied to a tuple of terms. Thus, $Likes(fred, bananas)$, $Likes(fred, father(fred))$, $Likes(x, father(x))$ are all examples of atoms. Finally, a *ground atomic formula* or a *ground atom* is an atom without any variables. Well-formed formulæ (wffs) are then formed using the following rules:

1. Any ground atomic formula is a wff;
2. If α is a wff then $\neg\alpha$ is a wff;
3. If α and β are wffs then $(\alpha \wedge \beta)$, $(\alpha \vee \beta)$, and $(\alpha \leftarrow \beta)$ are wffs; and
4. If α is wff containing a constant c and $\alpha^{c/x}$ be the result of replacing one or more occurrences of c with a variable x that does not appear in α . Then $\forall x\alpha^{c/x}$ and $\exists x\alpha^{c/x}$ are wffs.

Rules 1–3 are like their propositional counterparts (page 15). Rule 4 is new, and requires further explanation. It is the only way variables are introduced into a formula. As an example, take the following statement: $(Human(fred) \wedge Likes(fred, bananas))$. That this a wff follows from an application of Rules 1 and 3. The following formulæ are all wffs, following a single application of Rule 4:

$$\begin{aligned} &\forall x(Human(x) \wedge Likes(fred, bananas)) \\ &\exists x(Human(x) \wedge Likes(fred, bananas)) \\ &\forall x(Human(fred) \wedge Likes(x, bananas)) \\ &\exists x(Human(fred) \wedge Likes(x, bananas)) \\ &\forall x(Human(x) \wedge Likes(x, bananas)) \\ &\exists x(Human(x) \wedge Likes(x, bananas)) \\ &\forall x(Human(fred) \wedge Likes(fred, x)) \\ &\exists x(Human(fred) \wedge Likes(fred, x)) \end{aligned}$$

A single application of Rule 4 therefore only introduces a single new variable. Subsequent applications will introduce more. For example, take the first formula above: $\forall x(Human(x) \wedge Likes(fred, bananas))$. The following wffs all result from applying Rule 4 to this statement:

$$\begin{aligned} &\forall y\forall x(Human(x) \wedge Likes(y, bananas)) \\ &\exists y\forall x(Human(x) \wedge Likes(y, bananas)) \\ &\forall y\forall x(Human(x) \wedge Likes(fred, y)) \\ &\exists y\forall x(Human(x) \wedge Likes(fred, y)) \end{aligned}$$

As with propositional logic, it is acceptable to drop outermost brackets:

$$(\forall x(Ape(x) \leftarrow Human(x)) \wedge \exists x(Ape(x) \wedge \neg Human(x)))$$

can be written as:

$$\forall x(Ape(x) \leftarrow Human(x)) \wedge \exists x(Ape(x) \wedge \neg Human(x))$$

Clausal Form

We are now in a position to expand on the notion of clauses and literals, first introduced on page 25. Consider the conditional statement:

$$\forall x(Ape(x) \leftarrow Human(x))$$

Recall the following from page 24:

$$(\alpha \leftarrow \beta) \equiv (\alpha \vee \neg\beta)$$

This means:

$$\forall x(Ape(x) \leftarrow (Human(x))) \equiv \forall x(Ape(x) \vee \neg Human(x))$$

The term in brackets on right-hand side is an example of a *clause* in first-order logic. In general, formulæ consisting of universally-quantified clauses all look alike:²³

$$\forall x_1 \forall x_2 \dots (\alpha_1 \wedge \alpha_2 \dots)$$

That is, they consist of a prefix that consists only of universally quantifiers, and each α_i , or clause, is a quantifier-free formula that looks like:

$$\alpha_i = (\beta_1 \vee \beta_2 \vee \dots \vee \beta_n)$$

where each β_j , or *literal*. Thus, as in propositional logic, a clause is a disjunction of literals. Each literal, however, is not a proposition, but is either an atomic formula (like $Ape(x)$, sometimes called a *positive* literal) or a negated atomic formula (like $\neg Human(x)$, sometimes called a *negative* literal). It is sometimes convenient to use $\forall \mathbf{x}$ to denote $\forall x_1 \forall x_2 \dots$ and to adopt a set-based notation to represent a clausal formula: $\{\forall \mathbf{x} \alpha_1, \forall \mathbf{x} \alpha_2, \dots\}$. Here it is understood that the formula stands for a conjunction of clauses. Often, the quantification is taken to be understood and left out. Further, individual clauses are themselves sometimes written as sets of literals:

$$\alpha_i = \{\beta_1, \beta_2, \dots, \beta_n\}$$

Clausal forms are of particular interest, computationally speaking. The language of logic programs (usually written in the Prolog language). is, at least in its ‘pure’ form, equivalent to clausal-form logic. Here is an example:

²³We will take a few liberties here by not including some brackets.

<u>Logic Program in Prolog</u>	<u>Clausal Form</u>
grandfather(X,Y):- father(X,Z), parent(Z,Y).	$\{\forall x\forall y\forall z(Grandfather(x,y) \vee$ $\neg Father(x,z) \vee$ $\neg Parent(z,y)),$
father(henry,jane).	$Father(henry, jane),$
parent(jane,john).	$Parent(jane, john)\}$

There are three clauses in this example. The first clause has three literals and the remainder have one literal each. Further, each clause has exactly one positive literal: such clauses are called *definite* clauses. More generally, clauses that contain *at most* one positive literal are called *Horn* clauses. From now on, we may sometimes write clauses in a lazy manner that is somewhere in between the syntax of Prolog and the true clausal form:

Logic Program in Prolog

grandfather(X,Y):- father(X,Z), parent(Z,Y).

Clausal Form

$\{\forall x\forall y\forall z(Grandfather(x,y) \vee \neg Father(x,z) \vee \neg Parent(z,y))\}$

Lazy Clausal Form

$Grandfather(x,y) \leftarrow Father(x,z), Parent(z,y)$

Skolem Functions and Constants

Skolemization refers to a process of replacing an existentially quantified variable in a formula by a new term; it is merely the process of providing a *name* for something that already exists. Whether the new term is a functional expression or a constant depends on where the existential quantifier appears in the formula. If it is preceded by one or more universal quantifiers, like:

$$\forall x_1 \dots \forall x_n \exists y \alpha$$

then a single skolemization step replaces all occurrences of y in α by the functional expression $\mathbf{f}(x_1, \dots, x_n)$. Here $\mathbf{f}(\dots)$ is a function symbol, called a *Skolem function*, that does not appear anywhere in the formula. Thus, skolemization of the formula $\forall x \exists y Likes(x, y)$ results in $\forall x Likes(x, \mathbf{f}(x))$. In general, if the existential quantifier appears in between some universal quantifiers:

$$\forall x_1 \dots \forall x_{i-1} \exists y \forall x_{i+1} \dots \forall x_n \alpha$$

then all occurrences of y in α can be replaced by the functional expression $\mathbf{f}(x_1, \dots, x_{i-1})$. A special case arises if the existential quantifier precedes zero or more universal quantifiers:

$$\exists y \forall x_1 \dots \forall x_n \alpha$$

In this case, a single skolemization step replaces all occurrences of y in α by a Skolem function of arity 0 (that is, a constant) \mathbf{c} . Here \mathbf{c} is a constant symbol, called a *Skolem constant*, that does not appear anywhere in the formula. Thus, skolemization of the formula $\exists y \forall x \text{Likes}(x, y)$ results in $\forall x \text{Likes}(x, \mathbf{c})$. You can see that in both cases, a single step of Skolemization reduces the number of existential quantifiers in a formula ϕ by 1. Let us denote this single step by $s(\phi)$. It should be easy to see that repeatedly performing Skolemization steps (that is, $s(s(\dots s(\phi)))$) will result in a formula with just universal quantifiers.

Normal Forms

Universally-quantified clausal forms are a special case of specific kind of normal form for first-order formula, which we are now able to present, having described the process of Skolemization. A formula is said to be in *prenex normal form* or PNF if all its quantifiers are in front. That is, the formula looks something like $Q_1 x_1 \dots Q_n x_n \phi$, where the Q_i is either a \forall or \exists . Further, if ϕ is a conjunction of disjunction of literals, then the formula is said to be in *conjunctive prenex normal form*. It can be shown that every first-order formula ϕ can be expressed by an equivalent one ϕ' in conjunctive PNF. Here is an example: suppose we want to find the conjunctive PNF for $\phi : (\exists x \forall y \text{DisLikes}(x, y) \wedge \forall x \exists y \text{Likes}(x, y))$. We first rename variables to give $\phi' : (\exists x \forall y \text{DisLikes}(x, y) \wedge \forall u \exists v \text{Likes}(u, v))$. We can then move the quantifiers to the left giving: $\phi'' : \exists x \forall y \forall u \exists v (\text{DisLikes}(x, y) \wedge \text{Likes}(u, v))$, which is in conjunctive PNF.

For reasons that will become apparent, we are further interested here only in formulae in conjunctive PNFs in which all the quantifiers are universal ones (that is, \forall). Such formulae are said to be in *Skolem normal form*, or SNF. A SNF can be obtained from the conjunctive PNF by applying the Skolemization process to eliminate \exists quantifiers. For example, suppose we want to “Skolemize” the formula ϕ above. We first find the conjunctive PNF (ϕ'' above). Using the Skolemization procedure described earlier, we replace x by a Skolem constant and v by a Skolem function. The SNF of ϕ is $\phi^S : \forall y \forall u (\text{Dislikes}(\mathbf{c}, y) \wedge \text{Likes}(u, \mathbf{f}(u)))$.

We can now return to clausal forms. Here are the steps for converting a formula into a set of clauses in clausal form:

1. Rename variables to ensure there are no variables with the same names in different quantifiers;
2. Eliminate \leftarrow 's and iff's;
 - $\alpha_1 \text{ iff } \alpha_2 \Rightarrow (\alpha_1 \leftarrow \alpha_2) \wedge (\alpha_2 \leftarrow \alpha_1)$
 - $\alpha_1 \leftarrow \alpha_2 \Rightarrow \alpha_1 \vee \neg \alpha_2$
3. Move \neg 's inwards;

- $\neg(\exists \mathbf{X})\alpha \Rightarrow (\forall \mathbf{X})\neg\alpha$
- $\neg(\forall \mathbf{X})\alpha \Rightarrow (\exists \mathbf{X})\neg\alpha$
- $\neg(\alpha_1 \vee \alpha_2) \Rightarrow \neg\alpha_1 \wedge \neg\alpha_2$
- $\neg(\alpha_1 \wedge \alpha_2) \Rightarrow \neg\alpha_1 \vee \neg\alpha_2$
- $\neg\neg\alpha \Rightarrow \alpha$

4. Distribute \vee 's over \wedge 's;

- $\alpha \vee (\alpha_1 \wedge \alpha_2) \Rightarrow (\alpha \vee \alpha_1) \wedge (\alpha \vee \alpha_2)$
- Assuming \mathbf{X} does not occur in α_1 : $\alpha_1 \vee (\forall \alpha_2) \Rightarrow (\forall \mathbf{X})(\alpha_1 \vee \alpha_2)$
- Assuming \mathbf{X} does not occur in α_1 : $\alpha_1 \vee (\exists \alpha_2) \Rightarrow (\exists \mathbf{X})(\alpha_1 \vee \alpha_2)$

5. Distribute \forall 's;

- $(\forall \mathbf{X})(\alpha_1 \wedge \alpha_2) \Rightarrow (\forall \mathbf{X})\alpha_1 \wedge (\forall \mathbf{X})\alpha_2$

At this stage, the PNF form is generated.

6. Skolemise existentially quantified variables;

- $(\forall X_1)(\forall X_2) \dots (\forall X_n)(\exists Y)\alpha(Y) \Rightarrow (\forall X_1)(\forall X_2) \dots (\forall X_n)\alpha(f(X_1, X_2, \dots, X_n))$

If further applications of step 6 are possible, then they should be carried out.

7. Rewrite as clauses by dropping universal quantifiers; and

8. Standardise variables apart.

Here are some statements from a book by Lewis Carroll, written in first-order logic:

$$S_1 : \forall x(\text{Scented}(x) \leftarrow \text{Coloured}(x)) \wedge$$

$$S_2 : \forall x(\text{DisLike}(x) \leftarrow \neg\text{GrownOpen}(x)) \wedge$$

$$S_3 : \neg(\exists x(\text{GrownOpen}(x) \wedge \neg\text{Coloured}(x))) \wedge$$

$$S_4 : \neg\forall x((\text{DisLike}(x) \leftarrow \neg\text{Scented}(x)))$$

You should find that these sentences, when converted by the steps above, give the following set of clauses:

$$C_1 : \{\neg\text{Coloured}(x), \text{Scented}(x)\}$$

$$C_2 : \{\text{GrownOpen}(y), \text{DisLike}(y)\}$$

$$C_3 : \{\neg\text{GrownOpen}(z), \text{Coloured}(z)\}$$

$$C_4 : \{\neg\text{Scented}(c)\}$$

$$C_5 : \{\neg \text{DisLike}(c)\}$$

Recall that representing a clause by a set $\{L_1, L_2, \dots, L_k\}$, is just short-form for the disjunction $L_1 \vee L_2 \vee \dots \vee L_k$; and the actual clausal form for the formula $F : S_1 \wedge S_2 \wedge S_3 \wedge S_4 \wedge S_5$ is $C : \forall x \forall y \forall z (C_1 \wedge C_2 \wedge C_3 \wedge C_4 \wedge C_5)$. Is F equivalent to C ? To answer this, we need to understand how meanings are assigned to first-order formulæ.²⁴

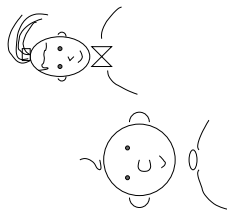
1.4.2 Semantics

As with propositional logic, the semantics of predicate logic is primarily concerned with interpretations, models, and logical consequence. Of these, it is only the notion of interpretation that requires a re-examination.

Interpretations

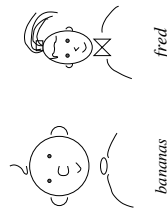
Recall that interpretations in propositional logic were simply assignments of *true* or *false* to propositional symbols. Matters are not as simple in predicate logic for two reasons. First, we have to deal with the additional complexity of expressions arising from a richer vocabulary. Thus, the truth-value (meaning) of $\text{Likes}(\text{fred}, \text{bananas})$ will depend on the meanings of each of the symbols Likes , fred and bananas . Second, we have to interpret sentences that contain quantifiers.

Informally, let us see what is required in constructing an interpretation that allows us to understand $\text{Likes}(\text{fred}, \text{bananas})$. Before we begin, remove any pre-conceptions of what the words in the statement mean in English: for us they are simply a predicate symbol (Likes) and two constant symbols (fred , bananas) in some ‘formal-world’. The first step then is to identify a domain of objects in the ‘real-world’.



Next, we associate constant symbols in our formal-world to objects in the real-world (just to avoid any pre-conceptions, we have scrambled things a little bit):

²⁴But the answer is “no”: the problem, as you might have guessed, comes about because of the Skolemization step.



and associate the predicate symbols to a relation in the real-world:



Likes

We can now see that $Likes(fred, bananas)$ is *false* as the objects corresponding to the ordered pair $\langle fred, bananas \rangle$ are not in the real-world relation represented by *Likes*.

Formally, an interpretation in predicate logic is a specification of:

1. A domain D ;
2. A mapping of constants to elements in D ;
3. A mapping of each n -argument predicate symbol to a relation on D^n , where $D^n = \{\langle d_1, \dots, d_n \rangle \mid d_i \in D\}$ is the n -fold Cartesian product; and
4. A mapping of each n -argument function symbol to a function from $D^n \rightarrow D$

You will also see sometimes that the term “structure” is used to describe what we have called an interpretation. In such cases, a distinction is made between the *vocabulary*, which consists of the constants, functions and predicate symbols; the and the *structure*, which consists of the domain D and the three mappings. We will continue to use “interpretation” to retain a similarity to propositional logic.

Given an interpretation, every atom—a predicate symbol with a tuple of terms as arguments—is assigned a truth-value according to whether the objects designated by the arguments are in the relation designated by the predicate symbol.

Well-formed formulæ in predicate logic consist of more than atoms. We also need rules for assigning truth-values to formulæ that contain logical connectives ($\neg, \wedge, \vee, \leftarrow$) and quantifiers (\forall, \exists). The semantics of the logical connectives in predicate logic are the same as those in propositional logic. Thus, as before, assigning meanings to formulæ with these connectives requires the use of the truth tables on page 17. For example, the formula $Human(fred) \wedge Likes(fred, bananas)$ is *true* only if the interpretation results in both the atoms $Human(fred)$ and $Likes(fred, bananas)$ being *true*. What though, of formulæ that contain quantified variables? The rules for these are:

1. Any wff $\forall x\alpha$ is *true* if and only if for every domain element that we can associate with x , α is *true*;
2. Any wff $\exists x\alpha$ is *true* if and only if for some domain element that we can associate with x , α is *true*.

Models and Logical Consequence

The meanings of these, and related concepts, are unchanged from propositional logic. Thus:

Models. Any interpretation that makes a wff *true* is called a model for that formula;

Validity and Unsatisfiability. A formula for which all interpretations are models is said to be valid. A formula for which none of the interpretations are a model is said to be unsatisfiable. A formula that has at least one model is said to be satisfiable;

Consequence. Given a conjunction of wffs Σ represented as the set $\{\beta_1, \dots, \beta_n\}$ and a wff α if $\Sigma \models \alpha$ then every model for Σ is a model for α ;

Deduction Theorem. Given a conjunction of wffs $\Sigma = \{\beta_1, \dots, \beta_n\}$ and a wff α , $\Sigma \models \alpha$ if and only if $\Sigma - \{\beta_i\} \models (\alpha \leftarrow \beta_i)$. The proof is the same as was for that in the propositional case (*c.f.* theorem 6);

Equivalence. Given a pair of wffs α and β , if $\alpha \models \beta$ and $\beta \models \alpha$ then α and β are equivalent ($\alpha \equiv \beta$).

The following relations hold in the predicate logic (as usual, α, β are wffs):

$\neg(\alpha \vee \beta)$	\equiv	$\neg\alpha \wedge \neg\beta$	De Morgan's law
$\neg(\alpha \wedge \beta)$	\equiv	$\neg\alpha \vee \neg\beta$	De Morgan's law
$(\alpha \leftarrow \beta)$	\equiv	$(\alpha \vee \neg\beta)$	
$(\alpha \leftarrow \beta)$	\equiv	$(\neg\beta \leftarrow \neg\alpha)$	Conditional \equiv Contrapositive
$\neg\forall x\alpha$	\equiv	$\exists x\neg\alpha$	
$\neg\exists x\alpha$	\equiv	$\forall x\neg\alpha$	
$\forall x\alpha$	\equiv	$\forall y\alpha^{x/y}$	Renaming of x by y
$\exists x\alpha$	\equiv	$\exists y\alpha^{x/y}$	Renaming of x by y
$\forall x\forall y\alpha$	\equiv	$\forall y\forall x\alpha$	
$\exists x\exists y\alpha$	\equiv	$\exists y\exists x\alpha$	
$\forall x(\alpha \wedge \beta)$	\equiv	$(\forall x\alpha \wedge \forall x\beta)$	Distributivity of \forall
$\exists x(\alpha \vee \beta)$	\equiv	$(\exists x\alpha \vee \exists x\beta)$	Distributivity of \exists
$\forall x\alpha$	\models	$\exists x\alpha$	
$(\forall x\alpha \vee \forall x\beta)$	\models	$\forall x(\alpha \vee \beta)$	
$\exists x(\alpha \wedge \beta)$	\models	$(\exists x\alpha \wedge \exists x\beta)$	
$\forall x\alpha^{y/\mathbf{f}(x)}$	\models	$\forall x\exists y\alpha$	Non-equivalence of Skolemized form
$\forall x\alpha^{y/\mathbf{c}}$	\models	$\exists y\forall x\alpha$	Non-equivalence of of Skolemized form

More on Normal Forms

We stated earlier that any first-order formula could be converted to a conjunctive prenex normal form, or conjunctive PNF. We further saw how a formula $\phi = Q_1x_1 \dots Q_nx_n\phi_0(x_1, \dots, x_n)$ could be “Skolemized” to give a formula ϕ^S in Skolem Normal Form, or SNF. The interest in SNFs lies in the following fact:

Theorem 15 *ϕ is satisfiable if and only if ϕ^S is satisfiable (you should be able to convince yourself that checking for satisfiability is equivalent to checking for logical consequence).*

Proof sketch: Recall that $\phi^S = s(s(\dots s(\phi)))$ where $s(\cdot)$ denotes a single step of Skolemization. Now, it is sufficient to show that ϕ is satisfiable if and only if $s(\phi)$ is satisfiable (the full proof will follow by induction). We can also assume that ϕ is in conjunctive PNF. Now $s(\phi)$ results in either replacing a variable by a Skolem constant or by a Skolem function. Since the former is just a special case of the latter, we will just consider the case when $s(\phi)$ results in replacing a variable by a Skolem function. Let Q_i be the existential quantifier removed by the Skolemization step, and let $\psi(x_1, \dots, x_i) = Q_{i+1} \dots Q_n\phi_0(x_1, \dots, x_n)$. Suppose $\forall x_1 \dots \forall x_{i-1} \exists x_i \psi(x_1, \dots, x_i)$ has some model M . Let $M_{\mathbf{f}}$ extend M by interpreting \mathbf{f} in such a way that $M_{\mathbf{f}}$ is a model for $\psi(c_1, \dots, c_{i-1}, \mathbf{f}(c_1, \dots, c_{i-1}))$ for all possible values $c_1, \dots, c_{i-1} \in M$ for the variables x_1, \dots, x_{i-1} . Then, clearly, $M_{\mathbf{f}}$ is a model for $\forall x_1 \dots \forall x_{i-1} \psi(x_1, \dots, x_{i-1}, \mathbf{f}(x_1, \dots, x_{i-1}))$. Now

consider the converse. Suppose $\forall x_1 \dots \forall x_{i-1} \psi(x_1, \dots, x_{i-1}, \mathbf{f}(x_1, \dots, x_{i-1}))$ has some model M then, it follows from the meaning of \exists that M is also a model for $\forall x_1 \dots \forall x_{i-1} \exists x_i \psi(x_1, \dots, x_i)$. It follows therefore that ϕ is satisfiable if and only if $s(\phi)$ is satisfiable. \square

We end this section on a note of caution: ϕ and $s(\phi)$ are not equivalent. That is, Skolemization does not preserve logical equivalence. A simple example should convince you of this. Let $\phi = \exists x \text{First}(x)$ and $s(\phi) = \text{First}(\mathbf{c})$. Clearly, we can find models for $\exists \text{First}(x)$ that are not models of $\text{First}(\mathbf{c})$.

Herbrand Interpretations

The 4-step specification above makes an interpretation in predicate logic much more elaborate than its counterpart in propositional logic (which was simply an assignment of *true* or *false* to propositions). The reference to ‘real-world’ objects, relations, and functions adds a further degree of complexity: how is all this to be conveyed to an automated procedure? In fact, many of these problems can be side-stepped by confining attention only to a domain that consists solely of formal symbols. Called the *Herbrand universe* (U_L), this is simply all the ground (or variable-free) terms that can be constructed using the constants and function symbols available in a first order language L . Consider as example a language that consists of:

Constant symbol:	<i>zero</i>
Predicate symbol:	<i>Nat/1</i>
Function symbols:	<i>pred, succ</i>

The Herbrand universe U_L in this instance consists of terms like *zero*, *pred(zero)*, *succ(zero)*, *pred(succ(zero))*, *succ(pred(zero))* and so on. The *Herbrand base* B_L is the set of all ground atoms that can be constructed using the predicate symbols and terms from the Herbrand universe U_L . Here, the Herbrand base B_L consists of atoms like *Nat(zero)*, *Nat(pred(zero))*, *Nat(succ(zero))*, *Nat(pred(succ(zero)))* and so on. A Herbrand interpretation I_L is—quite like the propositional case—simply an assignment of *true* to some subset of B_L and *false* to the rest. In fact, it is common practice to associate ‘Herbrand interpretation’ only with the subset assigned *true*: it being understood that all other atoms in the Herbrand base are assigned *false*. Thus, $\{\text{Nat}(\text{zero})\}$ is an I_L that assigns *true* to *Nat(zero)* and *false* to all other atoms in B_L .

Herbrand Models

Since a model is an interpretation that makes a well-formed formula *true*, a *Herbrand model* M_L is simply a Herbrand interpretation I_L that makes a well-formed formula *true*. Let us return to the example presented earlier:

Constant symbol:	<i>zero</i>
Predicate symbol:	<i>Nat/1</i>
Function symbols:	<i>pred, succ</i>

Recall from page 68, that:

$$\begin{aligned} \text{Herbrand universe } (U_L): & \quad \{zero, pred(zero), succ(zero), pred(succ(zero)), \dots\} \\ \text{Herbrand base } (B_L): & \quad \{Nat(zero), Nat(pred(zero)), Nat(succ(zero)), \dots\} \end{aligned}$$

Further, a Herbrand interpretation is simply a subset of the Herbrand base containing all atoms that are *true*. Thus, $I_1 = \{Nat(zero)\}$ is a Herbrand interpretation in which $Nat(zero)$ is *true* and all other atoms in the Herbrand base are *false*. We can now examine whether I_1 is a Herbrand model for the formula:

$$\Sigma_1 : Nat(zero) \wedge \forall x(Nat(succ(x)) \leftarrow Nat(x))$$

Being a conjunctive expression, we require I_1 to be a Herbrand model for both $Nat(zero)$ and $\forall x(Nat(succ(x)) \leftarrow Nat(x))$. I_1 is clearly a model for $Nat(zero)$ as this atom is assigned *true* in the interpretation. But what about the conditional? The rule for the universal quantifier (page 66) dictates that the conditional statement is *true* if it is *true* for every element that the variable x can be associated with. In other words, the interpretation is a model for every element that x can be associated with. In the Herbrand world, x can be associated with any element of the Herbrand universe ($zero, succ(zero), pred(zero)$ and so on). Suppose x was associated with $zero$. Then we would require I_1 to be a model for $Nat(succ(zero)) \leftarrow Nat(zero)$. Since I_1 assigns $Nat(succ(zero))$ to *false* and $Nat(zero)$ to *true*, I_1 is not a model for $Nat(succ(zero)) \leftarrow Nat(zero)$ (line 2 in the truth-table for the conditional on page 17). Thus, I_1 is not a Herbrand model for $\forall x(Nat(succ(x)) \leftarrow Nat(x))$ and in turn for Σ_1 . Consider, on the other hand, the formula:

$$\Sigma_2 : Nat(zero) \wedge \forall x(Nat(x) \leftarrow Nat(pred(x)))$$

As before, suppose x was associated with $zero$. The conditional then becomes $Nat(zero) \leftarrow Nat(pred(zero))$. With interpretation I_1 , $Nat(zero)$ is *true* and $Nat(pred(zero))$ is *false*. I_1 is therefore a model for this formula (line 3 in the truth table for the conditional). All other associations for x result in both sides of the conditional being *false* and I_1 being a model for each such formula (line 1 in the truth table). Thus, I_1 makes $\forall x(Nat(x) \leftarrow Nat(pred(x)))$ *true* for every element that x can be associated with, and is a model for it and in turn for Σ_2 .

Herbrand models are particularly relevant to the study of clausal forms (page 60). Recall that these are conjunctions of clauses, each of which contains only universally quantified variables and consists of a disjunction of literals. Both Σ_1 and Σ_2 above can be written in clausal form:

$$\begin{aligned} \Sigma_1' : & \quad Nat(zero) \wedge \forall x(Nat(succ(x)) \vee \neg Nat(x)) \\ \Sigma_2' : & \quad Nat(zero) \wedge \forall x(Nat(x) \vee \neg Nat(pred(x))) \end{aligned}$$

The *ground instantiation* of a clausal formula is the conjunction of ground (variable-free) clauses that result by replacing variables with terms from the Herbrand universe. For example, the ground instantiation of Σ_2' is:

$$\begin{aligned} \mathcal{G}(\Sigma_2') : \quad & \text{Nat}(\text{zero}) && \wedge \\ & (\text{Nat}(\text{zero}) \vee \neg \text{Nat}(\text{pred}(\text{zero}))) && \wedge \\ & (\text{Nat}(\text{pred}(\text{zero})) \vee \neg \text{Nat}(\text{pred}(\text{pred}(\text{zero})))) && \wedge \\ & (\text{Nat}(\text{succ}(\text{zero})) \vee \neg \text{Nat}(\text{pred}(\text{succ}(\text{zero})))) && \wedge \\ & \dots \end{aligned}$$

You can therefore think of the ground instantiation as making explicit the meaning of the universal quantifier \forall . Now, it should be clear that a Herbrand interpretation will determine the truth-value for all clauses in the ground instantiation of a clausal formula.

We present another example illustrating Herbrand models. Consider the following program P :

$$\begin{aligned} \text{likes}(\text{john}, X) &\leftarrow \text{likes}(X, \text{apples}) \\ \text{likes}(\text{mary}, \text{apples}) &\leftarrow \end{aligned}$$

Suppose the language \mathcal{L} contained no symbols other than those in P . Then, $\mathcal{B}(P)$ is the set $\{\text{likes}(\text{john}, \text{john}), \text{likes}(\text{john}, \text{apples}), \text{likes}(\text{apples}, \text{john}), \text{likes}(\text{john}, \text{mary}), \text{likes}(\text{mary}, \text{john}), \text{likes}(\text{mary}, \text{apples}), \text{likes}(\text{apples}, \text{mary}), \text{likes}(\text{mary}, \text{mary}), \text{likes}(\text{apples}, \text{apples})\}$. Now, $\{\text{likes}(\text{mary}, \text{apples}), \text{likes}(\text{john}, \text{mary})\}$ is a subset of $\mathcal{B}(P)$, and is a Herbrand interpretation. Moreover, it is also a Herbrand model for P . Similarly, $\{\text{likes}(\text{mary}, \text{apples}), \text{likes}(\text{john}, \text{mary}), \text{likes}(\text{mary}, \text{john})\}$ is also a model for P . The ground instantiation $\mathcal{G}(P)$ for this program is:

$$\begin{aligned} \text{likes}(\text{john}, \text{john}) &\leftarrow \text{likes}(\text{john}, \text{apples}) \\ \text{likes}(\text{john}, \text{mary}) &\leftarrow \text{likes}(\text{mary}, \text{apples}) \\ \text{likes}(\text{john}, \text{apples}) &\leftarrow \text{likes}(\text{apples}, \text{apples}) \\ \text{likes}(\text{mary}, \text{apples}) &\leftarrow \end{aligned}$$

It can be verified²⁵ that the interpretation $\{\text{likes}(\text{mary}, \text{apples}), \text{likes}(\text{john}, \text{mary})\}$ is a model for the $\mathcal{G}(P)$ above.

The importance of Herbrand models for clausal formulæ stems from the following property:

Theorem 16 *A clausal formula Σ has a model if and only if its ground instantiation $\mathcal{G}(\Sigma)$ has a Herbrand model.*

Proof: \Rightarrow : Suppose Σ has a model M . Then we define the following Herbrand interpretation I as follows. Let P be an n -ary predicate symbol occurring in Σ . Then we define the function I_P from U_L^n to $\{T, F\}$ as follows: $I_P(t_1, \dots, t_n) = T$ if $P(t_1, \dots, t_n)$ is true under M , and $I_P(t_1, \dots, t_n) = F$ otherwise. It can easily be shown that $I = \cup_{P \in \Sigma} I_P$ is a Herbrand model of Σ .

\Leftarrow : This is obvious (a Herbrand model is a model). \square

²⁵EXERCISE.

In other words, there must be *some* assignment of truth-values to atoms in the Herbrand base that makes all clauses in Σ *true*. In the example above, the Herbrand interpretation $I_1 = \{Nat(zero)\}$ that assigns $Nat(zero)$ to *true* and everything else to *false*, is clearly a model for $\mathcal{G}(\Sigma_2')$. Therefore, from the property stated here, we can say that Σ_2' has a model.

If we are dealing only with a *definite clausal formula*—a clausal formula in which all clauses have exactly one positive literal (Σ_1' and Σ_2' are both of this type)—then more is known about the Herbrand models of the formula. Recall that a Herbrand model is nothing more than a set of ground atoms, which when assigned *true*, make the formula *true*.

1.4.3 From Datalog to Prolog

The statement “*Any animal that has hair is a mammal*” can be written as a clause using monadic predicates (*i.e.* predicates with arity 1):

$$\forall X \text{ is_mammal}(X) \leftarrow \text{has_hair}(X)$$

Usually clauses are written without explicit mention of the quantifiers:

$$\text{is_mammal}(X) \leftarrow \text{has_hair}(X)$$

$$\text{is_mammal}(X) \leftarrow \text{has_milk}(X)$$

$$\text{is_bird}(X) \leftarrow \text{has_feathers}(X)$$

...

Datalog

Datalog is a subset of the language of first order language; it has all the components of first order logic (variables, constants and recursion), except functions. A Datalog “expert” system will encode these rules using monadic predicates as:

```
is_mammal(X) :- has_hair(X).
is_mammal(X) :- has_milk(X).
is_bird(X) :- has_feathers(X).
is_bird(X) :- can_fly(X), has_eggs(X).
is_carnivore(X) :- is_mammal(X), eats_meat(X).
is_carnivore(X) :- has_pointed_teeth(X), has_claws(X), has_pointy_eyes(X).
cheetah(X) :- is_carnivore(X), has_tawny_colour(X), has_dark_spots(X).
tiger(X) :- is_carnivore(X), has_tawny_colour(X), has_black_stripes(X).
penguin(X) :- is_bird(X), cannot_fly(X), can_swim(X).
```

Now here are some statements²⁶ particular to animals:

has_hair(peter).	fat(peter).
has_green_eyes(peter).	has_tawny_colour(peter).
eats_meat(peter).	has_black_stripes(peter).
has_milk(bob).	eats_meat(bob).
has_tawny_colour(bob).	has_dark_spots(bob).
can_fly(bob).	

²⁶EXERCISE: What are the logical consequences of all the clauses?

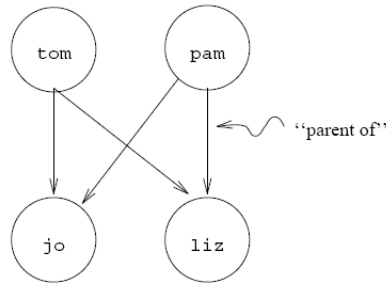


Figure 1.11: Graph representing ‘parent of’ relation.

However, monadic predicates: not expressive enough. While monadic predicates lets us make statements like “Every son has a parent”:

$$\forall X \exists Y \text{ parent}(Y) \leftarrow \text{son}(X)$$

for more complex relationships, we will need predicates of arity > 1 . Usually, relationships can be described pictorially by a directed acyclic graph (DAG) as in Figure 1.11 The parent-child relation could also be specified as a set of ordered pairs $\langle X, Y \rangle$, or, as a set of definite clauses

$$\begin{aligned} \text{parent}(\text{tom}, \text{jo}) &\leftarrow \\ \text{parent}(\text{pam}, \text{jo}) &\leftarrow \\ \text{parent}(\text{tom}, \text{liz}) &\leftarrow \\ \text{parent}(\text{pam}, \text{liz}) &\leftarrow \end{aligned}$$

Consider the *predecessor* relation, namely, all ordered tuples $\langle X, Y \rangle$ s.t. X is an ancestor of Y . This set will include Y ’s parents, Y ’s grandparents, Y ’s grandparents’ parents, etc.

$$\begin{aligned} \text{pred}(X, Y) &\leftarrow \text{parent}(X, Y) \\ \text{pred}(X, Z) &\leftarrow \text{parent}(X, Y), \text{parent}(Y, Z) \\ \text{pred}(X, Z) &\leftarrow \text{parent}(X, Y1), \text{parent}(Y1, Y2), \text{parent}(Y2, Z) \\ &\dots \end{aligned}$$

As can be seen through this example, variables and constants are not enough: we need *recursion*:

$$\begin{aligned} \forall X, Z \text{ X is a predecessor of Z if} \\ 1. \text{ X is a parent of Z; or} \\ 2. \text{ X is a parent of some Y, and Y is a predecessor of Z} \end{aligned}$$

The predecessor relation is thus

$$\begin{aligned} \text{pred}(X, Y) &\leftarrow \text{parent}(X, Y) \\ \text{pred}(X, Z) &\leftarrow \text{parent}(X, Y), \text{pred}(Y, Z) \end{aligned}$$

and can be pictorially depicted as in 1.12

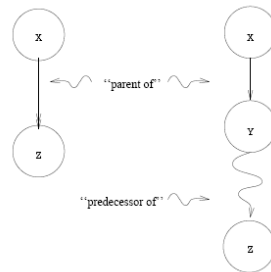


Figure 1.12: The predecessor relation.

Prolog = Predicates + Variables + Constants + Functions

Datalog (first order logic without functions) is however not expressive enough. To express arithmetic operations, lists of objects, etc. it is not enough to simply allow variables and constants as terms. We will also need *function* symbols as supported in Prolog.

Consider Peano's postulates for the set of natural numbers \mathcal{N} .

1. The constant 0 is in \mathcal{N}
2. if X is in \mathcal{N} then $s(X)$ is in \mathcal{N}
3. There are no other elements in \mathcal{N}
4. There is no X in \mathcal{N} s.t. $s(X) = 0$
5. There are no X, Y in \mathcal{N} s.t. $s(X) = s(Y)$ and $X \neq Y$

We can write a definite clause definition using 1 constant symbol and 1 unary function symbol for enumerating the elements of \mathcal{N} :

$$\begin{aligned} \text{natural}(0) &\leftarrow \\ \text{natural}(s(X)) &\leftarrow \text{natural}(X) \end{aligned}$$

The elements of \mathcal{N} can be now generated by asking:

natural(N)?

Prolog also supports lists. Lists are simply collections of objects. For e.g. 1, 2, 3... or 1, a, dog, ... Lists are defined as follows:

1. The constant *nil* is a list
2. If X is a term, and Y is a list then $.(X, Y)$ is a list

So the list 1, 2, 3 is represented as:

$$.(1,.(2,.(3,nil)))$$

Usually logic programming systems use a “[]” notation, in which the constant *nil* is represented as [] and the list 1, 2, 3 is [1, 2, 3]. In this notation, the symbol | is used to separate a list into a “head” (the elements to the left of the |) and a “tail” (the list to the right of the |). Thus:

List	Represented as	Values of variables
[1, 2, 3]	[X Y]	$X = 1, Y = [2, 3]$
[[1, 2], 3]	[X Y]	$X = [1, 2], Y = [3]$
[1]	[X Y]	$X = 1, Y = []$
[1 2]	[X Y]	$X = 1, Y = 2$
[1]	[X, Y]	
[1, 2, 3]	[X, Y Z]	$X = 1, Y = 2, Z = [3]$

1.4.4 Lattice of Herbrand Models

The discussion in this section is more or less similar to the discussion in Section 1.3.8 and the reader is referred to the proofs in that Section for proofs of most statements that will be made in this section. The only difference is that while Section 1.3.8, in this section, we will talk about Herbrand models.

For a definite clausal formula, it can be shown that the set H of Herbrand models is a complete lattice ordered by set-inclusion (that is, for Herbrand models $M1, M2 \in H$, $M1 \preceq M2$ if and only if $M1 \subseteq M2$), and with the binary operations of \cap and \cup as the glb and lub respectively. Recall that a complete lattice has a unique least upper bound and a unique greatest lower bound. Since we are really talking about sets that are ordered by set-inclusion, this means that there is a unique *smallest* one (the size being measured by the number of elements). This is called the *minimal model* of the formula, and it can be shown that this must be the intersection of all Herbrand models for the formula.

In the example above, Σ_2' has several Herbrand models ($\{Nat(Zero)\}$ and $\{Nat(Zero), Nat(pred(Zero))\}$ are two examples). Of these $\{Nat(Zero)\}$ is the smallest, and is the minimal model. There is an important result relating a definite clausal formula Σ , its minimal model $MM(\Sigma)$ and the ground atoms that are logical consequences of Σ :

Theorem 17 *If α is a ground atom then $\Sigma \models \alpha$ if and only if $\alpha \in MM(\Sigma)$.*

Here $MM(\cdot)$ denotes the minimal model. Thus, the minimal model of a definite clausal formula is identical to the set of all ground atoms logically implied by that formula. Thus, the minimal model provides, in effect, denotes the meaning (or semantics) of the formula. The proof of this theorem follows nearly from theorem 12 that was proved earlier.

We can envisage a procedure for enumerating the Herbrand models of a formula. Consider the powerset of the Herbrand base of the formula. Now,

we know that this powerset ordered by \subseteq necessarily forms a complete lattice, with binary operations \cap and \cup . Some subset of this powerset is the set of all Herbrand models, which we know is also a lattice ordered by \subseteq with the same binary operations. So, the model lattice is a sublattice of the lattice obtained from the powerset of the Herbrand base. Suppose now we start at some point s in this sublattice, and we move to a new point that consists only of those ground atoms of the formula made true by the model s . Let us call these atoms s_1 . Then, a little thought should convince you that s_1 is also a member of the sublattice of Herbrand models. Repeating the process with s_2 we can move to models s_2, s_3 and so on. Will this procedure converge eventually on the minimal model? Not necessarily, since we could end up moving back-and-forth between points of the sub-lattice. (When will this happen, and how can we ensure that we do converge on the minimal model?).

A slightly more general process can be formalised as the application of a function T_P that, for a clausal formula P , generates an interpretation (not necessarily a model) from another. That is:

$$I_{k+1} = T_P(I_k)$$

where

$$T_P(I) = \{a : a \leftarrow \text{body} \in \mathcal{G}(P) \text{ and } \text{body} \in I\}$$

where $\mathcal{G}(P)$ is the ground instantiation of P as before. It can be shown that T_P is both monotonic and continuous on the complete lattice obtained by ordering the powerset of the Herbrand base by \subseteq . So, we know from the Knaster-Tarski Theorem mentioned on page 11, that there must be a least fixpoint for T_P in this lattice. We can prove that the procedure of obtaining I_{k+1} from application of T_P to I_k will yield that fixpoint, and further, that this fixpoint will be the minimal model. As an inference procedure though, it is not really very practical: especially if all we needed to do is check if a particular atom was a logical consequence. It gets worse if the minimal model is not finite, in which case the procedure may not terminate in a finite number of steps. For all these reasons, we will need to do better.

1.4.5 Inference

Consider the following set of clauses S :

```
likes(john, flowers) ←
likes(mary, food) ←
likes(mary, wine) ←
likes(john, wine) ←
likes(john, mary) ←
likes(paul, mary) ←
```

If you entered these clauses into a program capable of executing logic programs (some implementation of Prolog), and asked:

likes(john, X)?

you will get a number of answers:

$$\begin{aligned} X &= \textit{flowers} \\ X &= \textit{wine} \\ X &= \textit{mary} \end{aligned}$$

On the other hand, if the query were

$$\textit{likes}(\textit{john}, X), \textit{likes}(\textit{mary}, X)?$$

the answer should be:

$$X = \textit{wine}$$

How this works will be examined in shortly. For now, consider $\textit{likes}(\textit{john}, X)?$. An intuitive procedure will be:

1. Start search from 1st clause
2. Search for any clause whose head has predicate $\textit{likes}/2$, and 1st argument is \textit{john}
3. If no clause is found *return* otherwise *goto* 4
4. X is associated (“instantiated”) with the 2nd argument of the head literal, the clause position marked, and the value associated with X is output
5. Start search from clause marked, and *goto* 2

As in the propositional case, we will only be concerned here with the rule of resolution. In a broad sense, this remains similar to its propositional counterpart (page 29: it applies to clauses with a pair of complementary literals, and the result (or resolvent) is a clause with the complementary pair removed. However the intricacies of predicate logic require a bit more care. Take the following pair of conditionals (and their clausal forms):

<u>Conditional</u>	<u>Clausal Form</u>
$\forall x(\textit{Ape}(x) \leftarrow \textit{Human}(x))$	$\forall x(\textit{Ape}(x) \vee \neg \textit{Human}(x))$
$\textit{Human}(\textit{fred}) \leftarrow$ $\textit{Human}(\textit{father}(\textit{fred}))$	$\textit{Human}(\textit{fred}) \vee \neg \textit{Human}(\textit{father}(\textit{fred}))$

For resolution to apply, we require the clausal forms to contain a pair of complementary literals. We nearly do have such a pair: $\neg \textit{Human}(x)$ in the first clause and $\textit{Human}(\textit{fred})$ in the second. It is apparent that if variable x in the first clause were to be restricted to the term \textit{fred} , then we would indeed have a complementary pair, and the resolvent is:

<u>Resolvent</u>	<u>Clausal Form</u>
$\textit{Ape}(\textit{fred}) \leftarrow$ $\textit{Human}(\textit{father}(\textit{fred}))$	$\textit{Ape}(\textit{fred}) \vee \neg \textit{Human}(\textit{father}(\textit{fred}))$

A single resolution step in predicate logic thus involves ‘substituting’ terms for variables so that a complementary pair of literals results. Here, such a pair would result if we could somehow ‘match’ the literals $Human(x)$ and $Human(fred)$. The resulting mapping of variables to terms is called the *unifier* of the two literals. Thus, mapping x to $fred$ is a unifier for the literals $Human(x)$ and $Human(fred)$.

Substitution

More generally, a *substitution* is a mapping from variables to terms that is usually denoted as $\theta = \{v_1/t_1, v_2/t_2, \dots, v_n/t_n\}$. Applying a substitution θ to a well-formed formula α results in a *substitution instance*, usually denoted by $\alpha\theta$. Thus, applying the substitution $\theta = \{x/fred\}$ to $\alpha : \forall x(Ape(x) \vee \neg Human(x))$ results in the substitution instance $\alpha\theta : (Ape(fred) \vee \neg Human(fred))$. We usually require substitutions to have the following properties:

1. They should be *functions*. That is, each variable to the left of the / should be distinct. Thus, $\{x/fred, x/bill\}$ is not a legal substitution; and
2. They should be *idempotent*. That is, each term to the right of the / should not contain a variable that appears to the left of the /. Thus, $\{x/father(x)\}$ is not a legal substitution. This test is sometimes called the “occurs-check”. The occur-check disallows self-referential bindings such as $X/f(X)$. However, the temptation to omit the occur-check in unification algorithms is very strong, owing to the high processing cost of including it; it is the only test in the comparison cycle which has to scrutinize the inner contents of terms, whereas all other tests examine only the terms’ principal (outermost) symbols.

A pair of substitutions can be *composed* (‘joined together’). For example, composing $\{x/father(y)\}$ with $\{y/fred\}$ results in $\{x/father(fred)\}$. In general, the result of composing substitutions

$$\theta_1 = \{u_1/s_1, \dots, u_m/s_m\}$$

$$\theta_2 = \{v_1/t_1, \dots, v_n/t_n\}$$

is (this may not be a legal substitution):

$$\theta_1 \circ \theta_2 = \{u_1/s_1\theta_2, \dots, u_m/s_m\theta_2\} \cup \{v_i/t_i \mid v_i \notin \{u_1, \dots, u_m\}\}$$

Theorem 18 *If α is a universally quantified expression that is not a term (i.e., a literal or a conjunction or disjunction of literals), and θ is a substitution, then the following holds: $\alpha \models \alpha\theta$. For example, $P(x) \vee \neg Q(y) \models P(a) \vee \neg Q(y)$, where we have used the substitution $\{x/a\}$.*

Proof sketch: The proof for this example is easy: suppose I is a model, with domain D , of $P(x) \vee \neg Q(y)$. Then for all $d_1 \in D$, and for all $d_2 \in D$, $I_P(d_1) = T$ or $I_Q(d_2) = F$. Suppose a is mapped to domain element d by I , then for all $d \in D$, $I_P(d) = T$ or $I_Q(d) = F$. Hence I is a model of $P(a) \vee \neg Q(y)$. It is clear that for different α or θ , a similar proof can always be given. Hence always $\alpha \models \alpha\theta$. \square

Unifiers

We are now in a position to state more formally the notion of unifiers. To say that a substitution θ is a unifier for formulæ α_1 and α_2 means $\alpha_1\theta = \alpha_2\theta$. However, there can be many unifiers. For example, the formulæ $\alpha_1 : \forall x \forall z \text{Parent}(\text{father}(x), z)$ and $\alpha_2 : \forall y \text{Parent}(y, \text{fred})$ have as unifiers $\theta_1 = \{x/\text{fred}, y/\text{father}(\text{fred}), z/\text{fred}\}$ and $\theta_2 = \{y/\text{father}(x), z/\text{fred}\}$. In the first case $\alpha_1\theta_1 = \alpha_2\theta_1 = \text{Parent}(\text{father}(\text{fred}), \text{fred})$; and in the second case $\alpha_1\theta_2 = \alpha_2\theta_2 = \forall x \text{Parent}(\text{father}(x), x)$. Notice that θ_2 is, in some sense, more ‘general’ than θ_1 as it imposes less severe constraints on the variables. There is, in fact, a *most general unifier* (or mgu) for a pair of formulæ. The substitution θ is a most general unifier for α_1 and α_2 if and only if:

1. $\alpha_1\theta = \alpha_2\theta$ (that is, θ is a unifier for α_1 and α_2); and
2. For any other unifier σ for α_1 and α_2 , there is a substitution μ such that $\sigma = \theta \circ \mu$ (that is, $\alpha_1\sigma$ is a substitution instance of $\alpha_1\theta$).

In the example just shown, θ_2 is the most general unifier.

Returning now to resolution, we can state the main steps involved for a pair of clauses C_1 and C_2 :

1. Rename all variables in clause C_2 so that they cannot be confused with those in C_1 (for the variables in C_2 are independent of those in C_1 and the renamed clause is equivalent to C_2). This is sometimes called “standardising the clauses apart”;
2. Identify complementary literals and see if an mgu exists;
3. Apply mgu and form the resolvent C .

Here is an example:

Formula

$C_1 : \forall x (\text{Ape}(x) \leftarrow \text{Human}(x))$

$C_2 : \forall x (\text{Human}(x) \leftarrow \text{Human}(\text{father}(x)))$

Clausal Form

$\forall x (\text{Ape}(x) \vee \neg \text{Human}(x))$

$\forall x (\text{Human}(x) \vee \neg \text{Human}(\text{father}(x)))$

The 3 steps above are:

1. Standardise apart. The two clauses are now:

$$C_1 : \forall x(Ape(x) \vee \neg Human(x))$$

$$C_2 : \forall y(Human(y) \vee \neg Human(father(y)))$$

2. Identify complementary literals and mgu. It is evident that $\neg Human(x)$ in C_1 and $Human(y)$ in C_2 are complementary. Their mgu is $\theta = \{x/y\}$;
3. Apply mgu and form resolvent. The resolvent C is as shown below:

$$C : \forall x(Ape(x) \vee \neg Human(father(x)))$$

As with propositional logic, the set-based notation used for clauses (page 60) allows us to present resolution in a compact (algebraic) form:

$$R = (C_1 - \{L\})\theta \cup (C_2 - \{M\})\theta$$

The difference to propositional logic is, of course, the appearance of θ , the mgu of literals L and $\neg M$. In fact, there is another problem that we have avoided. Suppose our clauses C_1 and C_2 are $C_1 : \forall x \forall y (Human(x) \vee Human(y))$ and $C_2 : \forall u \forall v (\neg Human(u) \vee \neg Human(v))$. Now it is clear that $\{C_1, C_2\}$ is unsatisfiable. But, unfortunately, we will not be able to get to the empty clause \square using resolution as we have just described it. Here is one possible resolvent: $R : \forall y \forall v (Human(y) \vee \neg Human(v))$. In fact, every possible resolvent of the two clauses will contain two literals, as will resolvents using those resolvents, and so on. What we really want to do is to eliminate redundant literals in any clause. For example, C_1 should really just be $\forall x Human(x)$ and C_2 should really just be $\forall u \neg Human(u)$. The procedure that removes redundant literals in this manner is called *factoring*.

Factoring

Formally, if C is a clause, L_1, \dots, L_n ($n \geq 1$) some unifiable literals from C , and θ an mgu for the set $\{L_1, \dots, L_n\}$, then the clause obtained by deleting $L_1\theta, \dots, L_n\theta$ from $C\theta$ is called a *factor* of C . For example, $Q(a) \vee P(f(a))$ is a factor of the clause $\neg Q(a) \vee P(f(a)) \vee P(y)$ using $\{y/f(a)\}$ as an mgu for $\{P(f(a)), P(y)\}$. Also, $Q(x) \vee P(x, a)$ is a factor of $Q(x) \vee Q(y) \vee Q(z) \vee P(z, a)$.

Operationally, it finds a substitution that unifies one or more literals in a clause, and retains only a single copy of the unified literals. Semantically speaking, a literal L is redundant in a clause C , if it is equivalent to a clause without that literal. That is $C - \{L\} \equiv C$. Note that every non-empty clause C is a factor of C itself, using the empty substitution \emptyset as mgu for one literal in C . It can easily be shown if C' is a factor of C , then $C \models C'$. We leave this to the reader to prove²⁷ From now on, we will assume that this elimination procedure has been executed on clauses, and we are only dealing with their “factors”.

²⁷Exercise.

Resolution

The rule of resolution remains sound for clauses in the predicate logic. That is, if C_1 and C_2 are clauses and R is a resolvent, then $\{C_1, C_2\} \models R$. The presence of variables and substitutions makes the proof of this a little more involved.

Theorem 19 *Suppose R is the result of resolving on literal L in C_1 and M in C_2 . Let θ be the most general unifier of L and $\neg M$ that is used to obtain R . Then, the soundness of a single step of resolution means $\{C_1, C_2\} \models (C_1 - \{L\})\theta \cup (C_2 - \{M\})\theta$.*

Proof: Let M be a model for C_1 and C_2 . Now, we know that either (a) $L\theta$ is true and $M\theta$ is false in M ; or (b) $L\theta$ is false and $M\theta$ is true in M . Suppose the former. Since M is a model for C_2 , it is a model for $C_2\theta$ (based on theorem 18). Therefore, at least one other literal $(C_2 - \{M\})\theta$ must be true in M . In other words, M is a model for $(C_1 - \{L\})\theta \cup (C_2 - \{M\})\theta$. Case (b) similarly results in M being a model for $(C_1 - \{L\})\theta$ and hence for R . So, a single resolution step is sound - the soundness of a proof consisting of several resolutions steps can be shown quite easily using the technique of induction. \square

Recall the second property of resolution from propositional logic, namely that of refutation-completeness. In other words, if a formula (or a set of formulae) is inconsistent, then the empty clause \square is derivable by the use of resolution. This property continues to hold for resolution in first-order logic. But before we look at that, we revisit an important result.

1.4.6 Subsumption Revisited

Recall that in propositional logic, a clause C subsumed a clause D if $C \subseteq D$. In first-order logic, this generalises as follows. A clause C *subsumes* a clause D if there is some substitution θ such that $C\theta \subseteq D$. What does this mean? It means that after applying the substitution θ to C , every literal in C appears in D . Here are a pair of clauses C and D such that C subsumes D :

$$C : \text{Primate}(x) \leftarrow \text{Ape}(x)$$

$$D : \text{Primate}(\text{Henry}) \leftarrow \text{Ape}(\text{Henry}), \text{Human}(\text{Henry})$$

Here, a substitution of $\theta = \{x/\text{Henry}\}$ applied to C makes $C\theta \subseteq D$. In general:

Theorem 20 *If C and D are clauses such that $C\theta \subseteq D$ for some substitution θ , then $C \models D$.*

Proof: Since C is a universally quantified formula, by theorem 18, we must have $C \models C\theta$. Also, since clauses are disjunctions of literals and $C\theta \subseteq D$, clearly, $C\theta \models D$ and the result follows. \square

However, unlike propositional logic, the reverse does not hold. That is, $C \models D$ does not necessarily mean that C subsumes D . Here is an example of this:

$$C : \text{Human}(x) \leftarrow \text{Human}(\text{father}(x))$$

$$D : \text{Human}(y) \leftarrow \text{Human}(\text{father}(\text{father}(y)))$$

With a little thought (let us not get too entangled in the species problem here), you should be able to convince yourself that $C \models D$. But you will find it impossible to find a substitution θ that will make $C\theta \subseteq D$. What makes the difference to the propositional case? The difference between implication and subsumption in first-order logic arises because of self-recursive clauses of the kind shown: a short, but influential paper by Georg Gottlob shows that it is indeed only the self-recursive case that results in the difference.

1.4.7 Subsumption Lattice over Atoms

The subsumption relation is an example of a quasi-order. Let us take the simple case of definite clauses with a single literal (that is, atoms). Consider the set \mathcal{A} of all atoms in some language, and $\mathcal{A}^+ = \mathcal{A} \cup \{\top, \perp\}$. Let the binary relation \succeq be such that:

- $\top \succeq \mathbf{l}$ for all $\mathbf{l} \in \mathcal{A}^+$
- $\mathbf{l} \succeq \perp$ for all $\mathbf{l} \in \mathcal{A}^+$
- $\mathbf{l} \succeq \mathbf{m}$ iff there is a substitution θ such that $\mathbf{l}\theta = \mathbf{m}$, for $\mathbf{l}, \mathbf{m} \in \mathcal{A}$

We will represent a list of elements e_1, \dots, e_n as the (as the language Prolog does) by $[e_1, \dots, e_n]$, and let $\mathbf{l} = \text{Mem}(x, [x, y])$ and $\mathbf{m} = \text{Mem}(1, [1, 2])$ then $\mathbf{l} \succeq \mathbf{m}$ with $\theta = \{x/1, y/2\}$. It is easy to see that \succeq is a quasi-order over \mathcal{A}^+ : clearly $\mathbf{l} \succeq \mathbf{l}$, with the empty substitution $\theta = \emptyset$ (that is, \succeq is reflexive). Now, let $\mathbf{l} \succeq \mathbf{m}$ and $\mathbf{m} \succeq \mathbf{l}$. That is, there are some substitutions θ_1 and θ_2 such that $\mathbf{l}\theta_1 = \mathbf{m}$ and $\mathbf{m}\theta_2 = \mathbf{l}$. That is, $(\mathbf{l}\theta_1) \circ \theta_2 = \mathbf{l}$. With $\theta = \theta_1 \circ \theta_2$ it follows that $\mathbf{l} \succeq \mathbf{l}$.

Since \succeq is a quasi-order, we know a partial ordering must result from the partition of \mathcal{A}^+ into a set of equivalence classes \mathcal{A}_E^+ . In fact, the partitions are $\{\{\top\}\}, \{\{\perp\}\}, X_1, \dots$ where $[l]$ denotes all atoms that are alphabetic variants²⁸ of \mathbf{l} . That is, if $\mathbf{l}, \mathbf{m} \in X_i$ then there are substitutions μ and σ s.t. $\mathbf{l}\mu = \mathbf{m}$ and $\mathbf{m}\sigma = \mathbf{l}$. That is, \succeq is a partial ordering over the set of equivalence classes of atoms (\mathcal{A}_E^+). ($\text{Mem}(x_1, [x_1, y_1]), \text{Mem}(x_2, [x_2, y_2]) \dots$ are examples of members of an equivalence class.)

Recall that the difference between subsumption and implication in first-order logic arose with the appearance of self-recursive clauses. Since there is no possibility of this with atoms in first-order logic, subsumption and implication are equivalent, and we can see that logical implication (*models*) over atoms is also a quasi-order over atoms.

²⁸Two atoms are subsume-equivalent iff they are variants. This is not true for clauses in general.

As soon as we have a quasi-order, we can effectively construct a partial-order over equivalence classes. So, the quasi-order of subsumption over atoms results in a partial order over equivalence classes of atoms. In fact, \mathcal{A}_E^+ is a lattice with the binary operations \sqcap and \sqcup defined on elements of \mathcal{A}_E^+ as follows (here, we have used $[\cdot]$ to represent an equivalence class):

- $[\perp] \sqcap [\mathbf{l}] = [\perp]$, and $[\top] \sqcap [\mathbf{l}] = [\mathbf{l}]$
- If $\mathbf{l}_1, \mathbf{l}_2 \in \mathcal{A}$ have a most general unifier (see page 78) θ then $[\mathbf{l}_1] \sqcap [\mathbf{l}_2] = [\mathbf{l}_1\theta] = [\mathbf{l}_2\theta]$.

This can be proved as follows. Let $[\mathbf{u}] \in \mathcal{A}_E^+$ such that $[\mathbf{l}_1] \succeq [\mathbf{u}]$ and $[\mathbf{l}_2] \succeq [\mathbf{u}]$, then we need to show that $[\mathbf{l}_1\theta] \succeq [\mathbf{u}]$. If $[\mathbf{u}] = [\perp]$, this is obvious. If $[\mathbf{u}]$ is conventional, then there are substitutions σ_1 and σ_2 such that $[\mathbf{l}_1\sigma_1] = [\mathbf{u}] = [\mathbf{l}_2\sigma_2]$. Here we can assume σ_1 only acts on variables in \mathbf{l}_1 , and σ_2 only acts on variables in \mathbf{l}_2 . Let $\sigma = \sigma_1 \cup \sigma_2$. Notice that σ is a unifier for $\{[\mathbf{l}_1], [\mathbf{l}_2]\}$. Since θ is an mgu for $\{[\mathbf{l}_1\sigma_1], [\mathbf{l}_2\sigma_2]\}$, there is a γ such that $\theta\gamma = \sigma$. Now $[\mathbf{l}_1\theta\gamma] = [\mathbf{l}_1\sigma] = [\mathbf{l}_1\sigma_1] = [\mathbf{u}]$, so $[\mathbf{l}_1\theta] \succeq [\mathbf{u}]$.

- If $\mathbf{l}_1, \mathbf{l}_2 \in \mathcal{A}$ do not have a most general unifier θ then $[\mathbf{l}_1] \sqcap [\mathbf{l}_2] = [\perp]$.
Since \mathbf{l}_1 and \mathbf{l}_2 are not unifiable, there is no conventional atom \mathbf{u} such that $[\mathbf{l}_1] \succeq [\mathbf{u}]$ and $[\mathbf{l}_2] \succeq [\mathbf{u}]$. Hence $[\mathbf{l}_1] \sqcap [\mathbf{l}_2] = [\perp]$.
- $[\perp] \sqcup [\mathbf{l}] = [\mathbf{l}]$, and $[\top] \sqcup [\mathbf{l}] = [\top]$
- If \mathbf{l}_1 and \mathbf{l}_2 have an “anti-unifier” \mathbf{m} then $[\mathbf{l}_1] \sqcup [\mathbf{l}_2] = [\mathbf{m}]$; otherwise $[\mathbf{l}_1] \sqcup [\mathbf{l}_2] = [\top]$

Henceforth, we will drop the square brackets $[\cdot]$ to denote equivalence classes and will instead implicitly assume their presence. The “anti-unifier” in the join operation is not something we have come across before, and needs some explanation. To get started, let us look at the atom $Mem(1, [1, 2])$. The list $[1, 2]$ written out in long-hand is really a term composed of the constants 1, 2 and the empty list, which we will denote by the constant nil . That is, $[1, 2]$ is really the term $list(1, list(2, list(nil)))$, where $list$ is a function, and $Mem(1, [1, 2])$ is really $Mem(1, list(1, list(2, nil)))$. Now, we can devise a “term-place” notation to identify the occurrence of each term in any atom. In $Mem(1, list(1, list(2, nil)))$, the 1 is a term that occurs in two places: in the first argument (or “place”) of Mem , and as the first argument of the second place of Mem . We can denote these two occurrences as $(1, \langle 1 \rangle)$ and $(1, \langle 2, 1 \rangle)$. Similarly, we can encode the occurrences of other terms: $(2, \langle 2, 2, 1 \rangle)$ and $(nil, \langle 2, 2, 2 \rangle)$.

You should convince yourself that the occurrence of every term t in an atom can indeed be represented by the pair (t, p) , where p is a sequence of places. We now have all we need to be able to describe the anti-unification algorithm for a pair of literals with the same predicate symbol (adapted from Plotkin, 1970):

Input: A pair of atoms \mathbf{l}_1 and \mathbf{l}_2 with the same predicate symbol

Output: $\mathbf{l}_1 \sqcup \mathbf{l}_2$

1. Let $\mathbf{l} = \mathbf{l}_1$ and $\mathbf{m} = \mathbf{l}_2$, $\theta = \emptyset$, $\sigma = \emptyset$
2. If $\mathbf{l} = \mathbf{m}$ return \mathbf{l} and stop.
3. Try to find terms t_1 and t_2 that have the same (leftmost) place in \mathbf{l} and \mathbf{m} respectively, such that $t_1 \neq t_2$ and either t_1 and t_2 begin with different function symbols, or at least one of them is a variable.
4. If there is no such t_1, t_2 , return \mathbf{l} and stop.
5. Choose a variable x that does not occur in either \mathbf{l} or \mathbf{m} and wherever t_1 and t_2 occur in the same place in \mathbf{l} and \mathbf{m} , replace each of them by x
6. Set θ to $\theta \cup \{x/t_1\}$ and σ to $\sigma \cup \{x/t_2\}$
7. Go to Step 3

The Table 1.1 shows the progressive construction of lubs starting with terms and culminating in literals.

Lub	Definition	Examples
lub of terms $lub(t_1, t_2)$	<ol style="list-style-type: none"> 1. $lub(t, t) = t$, 2. $lub(f(s_1, \dots, s_n), f(t_1, \dots, t_n)) = f(lub(s_1, t_1), \dots, lub(s_n, t_n))$, 3. $lub(f(s_1, \dots, s_m), g(t_1, \dots, t_n)) = V$, where $f \neq g$, and V is a variable which represents $lub(f(s_1, \dots, s_m), g(t_1, \dots, t_n))$, 4. $lub(s, t) = V$, where $s \neq t$ and at least one of s and t is a variable; in this case, V is a variable which represents $lub(s, t)$. 	<ul style="list-style-type: none"> • $lub([a, b, c], [a, c, d]) = [a, X, Y]$. • $lub(f(a, a), f(b, b)) = f(lub(a, b), lub(a, b)) = f(V, V)$ where V stands for $lub(a, b)$. • When computing lggs one must be careful to use the same variable for multiple occurrences of the lubs of subterms, i.e., $lub(a, b)$ in this example. This holds for lubs of terms, atoms and clauses alike.
lub of atoms $lub(\mathbf{a}_1, \mathbf{a}_2)$	<ol style="list-style-type: none"> 1. $lub(P(s_1, \dots, s_n), P(t_1, \dots, t_n)) = P(lub(s_1, t_1), \dots, lub(s_n, t_n))$, if atoms have the same predicate symbol P, 2. $lub(P(s_1, \dots, s_m), Q(t_1, \dots, t_n))$ is undefined if $P \neq Q$. 	
lub of literals $lub(\mathbf{l}_1, \mathbf{l}_2)$	<ol style="list-style-type: none"> 1. if \mathbf{l}_1 and \mathbf{l}_2 are atoms, then $lub(\mathbf{l}_1, \mathbf{l}_2)$ is computed as defined above, 2. if both \mathbf{l}_1 and \mathbf{l}_2 are negative literals, $\mathbf{l}_1 = \overline{\mathbf{a}_1}$, $\mathbf{l}_2 = \overline{\mathbf{a}_2}$, then $lub(\mathbf{l}_1, \mathbf{l}_2) = lub(\overline{\mathbf{a}_1}, \overline{\mathbf{a}_2}) = \overline{lub(\mathbf{a}_1, \mathbf{a}_2)}$, 3. if \mathbf{l}_1 is a positive and \mathbf{l}_2 is a negative literal, or vice versa, $lub(\mathbf{l}_1, \mathbf{l}_2)$ is undefined. 	<ul style="list-style-type: none"> • $lub(\text{Parent}(\text{ann}, \text{mary}), \text{Parent}(\text{ann}, \text{tom})) = \text{Parent}(\text{ann}, X)$. • $lub(\text{Parent}(\text{ann}, \text{mary}), \overline{\text{Parent}(\text{ann}, \text{tom})}) = \text{undefined}$. • $lub(\text{Parent}(\text{ann}, X), \text{Daughter}(\text{mary}, \text{ann})) = \text{undefined}$.

Table 1.1: Table showing progressive definitions of lubs, starting with terms and culminating in literals.

Let us look at an example of constructing the anti-unifier of $Mem(1, [1, 2])$ and $Mem(2, [2, 4])$. That is, $\mathbf{l}_1 = Mem(1, list(1, list(1, list(2, nil))))$ and $\mathbf{l}_2 = Mem(2, list(2, list(2, list(4, nil))))$. You should be able to work through the

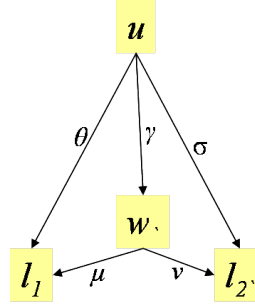


Figure 1.13: Illustration of the proof of theorem 22.

steps of the algorithm to find that it terminates with $\mathbf{l} = \mathbf{m} = Mem(x, list(x, list(y, nil)))$ with $\theta = \{x/1, y/2\}$ and $\sigma = \{x/2, y/4\}$.

But is the procedure correct? That is, does it really return a lub of a pair of atoms \mathbf{l}_1 and \mathbf{l}_2 ? Suppose the procedure returned an atom \mathbf{l} , and let $\theta = \{x_1/s_1, \dots, x_k/s_k\}$ and $\sigma = \{x_1/t_1, \dots, x_k/t_k\}$. That is $\mathbf{l}\theta = \mathbf{l}_1$ and $\mathbf{l}\sigma = \mathbf{l}_2$. Now suppose there is some other atom \mathbf{l}' such that $\mathbf{l}' \succeq \mathbf{l}_1$ and $\mathbf{l}' \succeq \mathbf{l}_2$. Then, we have to show that $\mathbf{l}' \succeq \mathbf{l}$ for any such \mathbf{l}' .

The proof of this is a bit laborious and will be dealt with subsequently. The truth of the next lemma is easy to see:

Theorem 21 *After each iteration of the Anti-Unification Algorithm, there are terms s_1, \dots, s_i and t_1, \dots, t_i such that:*

1. $\theta = \{z_1/s_1, \dots, z_i/s_i\}$ and $\sigma = \{z_1/t_1, \dots, z_i/t_i\}$.
2. $\mathbf{l}\theta = \mathbf{l}_1$ and $\mathbf{m}\sigma = \mathbf{l}_2$.
3. For every $1 \leq j \leq i$, s_j and t_j differ in their first symbol.
4. There are no $1 \leq j, k \leq i$ such that $j \neq k$, $s_j = s_k$ and $t_j = t_k$.

Theorem 22 *Let \mathbf{l}_1 and \mathbf{l}_2 be two atoms with the same predicate symbol. Then the Anti-Unification Algorithm with \mathbf{l}_1 and \mathbf{l}_2 as inputs returns $\mathbf{l}_1 \sqcup \mathbf{l}_2$.*

Proof: It is easy to see that the algorithm terminates after a finite number of steps, for any $\mathbf{l}_1, \mathbf{l}_2$. Let \mathbf{u} be the atom that the algorithm returns, and let $\theta = \{z_1/s_1, \dots, z_i/s_i\}$ and $\sigma = \{z_1/t_1, \dots, z_i/t_i\}$ be the final values of θ and σ in the computation of \mathbf{u} (so \mathbf{u} equals the final values of \mathbf{l} and \mathbf{m} in the execution of the algorithm). Then $\mathbf{u}\theta = \mathbf{l}_1$ and $\mathbf{u}\sigma = \mathbf{l}_2$ by Theorem 21, part 2. Suppose \mathbf{v} is an atom such that $\mathbf{v} \succeq \mathbf{l}_1$ and $\mathbf{v} \succeq \mathbf{l}_2$. In order to show that $\mathbf{u} = \mathbf{l}_1 \sqcup \mathbf{l}_2$, we have to prove $\mathbf{v} \succeq \mathbf{u}$.

Let $\mathbf{w} = \mathbf{u} \sqcap \mathbf{v}$, which exists by the proof on page 82. Then $\mathbf{u} \succeq \mathbf{w}$ and $\mathbf{v} \succeq \mathbf{w}$. Since $\mathbf{w} = \mathbf{u} \sqcap \mathbf{v}$, $\mathbf{u} \succeq \mathbf{l}_1$ and $\mathbf{v} \succeq \mathbf{l}_1$, we must have $\mathbf{w} \succeq \mathbf{l}_1$. Similarly

$\mathbf{w} \succeq \mathbf{l}_2$. Thus there are substitutions γ, μ, ν , such that $\mathbf{u}\gamma = \mathbf{w}$, $\mathbf{l}_1 = \mathbf{w}\mu = \mathbf{u}\gamma\mu$ and $\mathbf{l}_2 = \mathbf{w}\nu = \mathbf{u}\gamma\nu$. Then $\mathbf{u}\theta = \mathbf{l}_1 = \mathbf{u}\gamma\mu$ and $\mathbf{u}\sigma = \mathbf{l}_2 = \mathbf{u}\gamma\nu$ (see Figure 1.13 for illustration). Hence, if x is a variable occurring in \mathbf{u} , then $x\theta = x\gamma\mu$ and $x\sigma = x\gamma\nu$.

We will now show that \mathbf{u} and $\mathbf{w} = \mathbf{u}\gamma$ are variants, by showing that γ is a renaming substitution for \mathbf{u} . Suppose it is not. Then γ maps some variable x in \mathbf{u} to a term that is not a variable, or γ unifies two distinct variables x, y in \mathbf{u} .

Suppose x is a variable in \mathbf{u} , such that $x\gamma = t$, where t is a term that is not a variable. If x is not one of the z_j 's, then $x\gamma\mu = x\theta = x$, contradicting the assumption that $x\gamma = t$ is not a variable. But on the other hand, if x equals some z_j , then $t\mu = x\gamma\mu = x\theta = s_j$ and $t\nu = x\gamma\nu = x\sigma = t_j$. Then s_j and t_j would both start with the first symbol of t , contradicting theorem 21, part 3. So this case leads to a contradiction.

Suppose x and y are distinct variables in \mathbf{u} such that γ unifies x and y . Then,

1. If neither x nor y is one of the z_j 's, then $x\gamma\mu = x\theta = x \neq y = y\theta = y\gamma\mu$, contradicting $x\gamma = y\gamma$.
2. If x equals some z_j and y does not, then $x\gamma\mu = x\theta = s_j$ and $x\gamma\nu = x\sigma = t_j$, so $x\gamma\mu \neq x\gamma\nu$ by theorem 21, part 3. But $y\gamma\mu = y\theta = y = y\sigma = y\gamma\nu$, contradicting $x\gamma = y\gamma$.
3. Similarly for the case where y equals some z_j and x does not.
4. If $x = z_j$ and $y = z_k$, then $j \neq k$, since $x \neq y$. Furthermore, $s_j = x\theta = x\gamma\mu = y\gamma\mu = y\theta = s_k$ and $t_j = x\sigma = x\gamma\nu = y\gamma\nu = y\sigma = t_k$. But this contradicts theorem 21, part 4.

Thus, the assumption that γ unifies two variables in \mathbf{u} also leads to a contradiction. Thus γ is a renaming substitution for \mathbf{u} , and hence \mathbf{u} and \mathbf{w} are variants. Finally, since $\mathbf{v} \succeq \mathbf{w}$, we have $\mathbf{v} \succeq \mathbf{u}$. \square

It is however a more straightforward matter to see the following:

Theorem 23 *If \mathcal{A}_E^+ is a set of equivalence classes of atoms \mathcal{A} augmented by the two elements \top and \perp , then for any pair of elements $\mathbf{l}_1, \mathbf{l}_2 \in \mathcal{A}_E^+$, $\mathbf{l}_1 \sqcup \mathbf{l}_2$ always exists.*

Proof: The possibilities for each of \mathbf{l}_1 and \mathbf{l}_2 are that they are either: (1) some variant of \top ; (2) some variant of \perp ; or (3) an atom from S . It can be verified that $\mathbf{l}_1 \sqcap \mathbf{l}_2$ is defined in all 9 cases.

1. If $\mathbf{l}_1 = \top$ or $\mathbf{l}_2 = \top$, then $\mathbf{l}_1 \sqcup \mathbf{l}_2 = \top$. If $\mathbf{l}_1 = \perp$, then $\mathbf{l}_1 \sqcup \mathbf{l}_2 = \mathbf{l}_2$. If $\mathbf{l}_2 = \perp$, then $\mathbf{l}_1 \sqcup \mathbf{l}_2 = \mathbf{l}_1$.
2. If \mathbf{l}_1 and \mathbf{l}_2 are conventional atoms with the same predicate symbol, $\mathbf{l}_1 \sqcup \mathbf{l}_2$ is given by the Anti-Unification Algorithm.

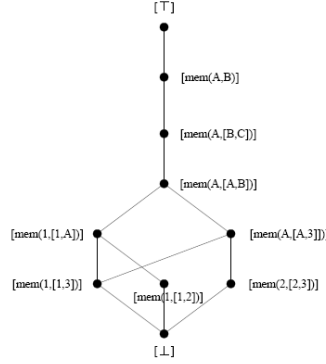


Figure 1.14: An example subsumption lattice over atoms.

3. If \mathbf{l}_1 and \mathbf{l}_2 are conventional atoms with different predicate symbols, then $\mathbf{l}_1 \sqcup \mathbf{l}_2 = \top$.

□

Now that we have established the existence of an lub and glb of any $\mathbf{l}_1, \mathbf{l}_2 \in \mathcal{A}_E^+$, we have shown that the set of atoms ordered by subsumption, is a lattice.

Theorem 24 *Let \mathcal{A} be the set of atoms. Then $\langle \mathcal{A}_E^+, \succeq \rangle$ is a lattice.*

Figure 1.14 shows an example subsumption lattice over atoms in $S^+ = \{ \top, \perp, mem(1, [1, 3]), mem(1, [1, 2]), mem(2, [2, 3]), mem(1, [1, A]), mem(A, [A, B]), mem(A, [A, 3]), mem(A, [B, C]), mem(A, [B|C]), mem(A, B) \}$. Note that

- $l = mem(A, [A, B]) \succeq mem(1, [1, 2]) = m$ since with $\theta = \{A/1, B/2\}$, $l\theta = m$
- $mem(A1, [A1, B1]), mem(A2, [A2, B2]) \dots$ are all members of the same equivalence class

Recall that, for atoms $l, m \in S$, subsumption is equivalent to implication. That is, if $l \models m$ then $l \succeq m$. Least-general-generalisation of atoms will be encountered in Lab Nos. 5, 6.

What about subsumption over clauses with more than just one literal? Is this still a quasi-order, with a lattice structure over equivalence classes of clauses? The short answer is “yes”, but more on this in Chapter 2.

1.4.8 Covers of Atoms

What about the covers relation in the subsumption lattice of atoms? Recall that covers are the smallest non-trivial steps between individual atoms that we can take in the lattice. Since \mathbf{l}_2 is a downward cover of \mathbf{l}_1 iff \mathbf{l}_1 is an upward cover of \mathbf{l}_2 , we will first restrict attention to downward covers.

Downward Covers

Theorem 25 *Let \mathbf{l}_1 be a conventional atom, f an n -ary function symbol (recall that f can be of zero arity and therefore a constant), z a variable in \mathbf{l}_1 , and x_1, \dots, x_n , distinct variables not appearing in \mathbf{l}_1 . Let*

1. $\theta = \{z/f(x_1, \dots, x_n)\}$ and
2. $\sigma = \{x_i/x_j\}$, $i \neq j$

Then $\mathbf{l}_2 = \mathbf{l}_1\theta$ and $\mathbf{l}_3 = \mathbf{l}_1\sigma$ are both downward covers of \mathbf{l}_1 . In fact, every downward cover of \mathbf{l}_1 must of one of these two forms (note that a special instance of the first case is when the function f has arity 0 and is therefore a constant). The substitutions θ and σ are termed elementary substitutions. In ILP, these 3 operations define a “downward refinement operator”

Proof:

Proof for (1): It is clear that \mathbf{l}_1 and \mathbf{l}_2 are not variants, so $\mathbf{l}_1 \succ \mathbf{l}_2$. Suppose there is a \mathbf{m} such that $\mathbf{l}_1 \succ \mathbf{m}\mathbf{l}_2$. Then there are γ, μ , such that $\mathbf{l}_1\gamma = \mathbf{m}$ and $\mathbf{m}\mu = \mathbf{l}_2$, hence $\mathbf{l}_1\gamma\mu = \mathbf{l}_2 = \mathbf{l}_1\theta$. Here γ only acts on variables in \mathbf{l}_1 , and μ only acts on variables in \mathbf{m} .

Let (x, p) be a term occurrence in \mathbf{l}_1 , where x is a variable. Suppose $x \neq z$, then $x\theta = x$, so (x, p) must also be a term occurrence in \mathbf{l}_2 . Hence $x\gamma$ must be a variable, for otherwise $(x\gamma\mu, p)$ in \mathbf{l}_2 would contain a constant or a function. Thus γ must map all variables other than z to variables. Furthermore, γ cannot unify two distinct variables in \mathbf{l}_1 , for then \mathbf{l}_2 would also have to unify these two variables, which is not the case.

If $z\gamma$ is also a variable, then γ would map all variables to variables, and since γ cannot unify distinct variables, it would map all distinct variables in \mathbf{l}_1 to distinct variables. But then γ would be a renaming substitution for \mathbf{l}_1 , contradicting $\mathbf{l}_1 \succ \mathbf{m}$. Hence γ must map z to some term containing a function symbol.

Now the only way we can have $\mathbf{l}_1\gamma\mu = \mathbf{l}_2$, is if $z\gamma = f(y_1, \dots, y_n)$ for distinct y_i not appearing in \mathbf{l}_1 , and no variable in \mathbf{l}_1 is mapped to some y_i by γ . But then $\mathbf{l}_1\gamma$ and \mathbf{l}_2 would be variants, contradicting $\mathbf{l}_1\gamma = \mathbf{m} \succ \mathbf{l}_2$. Therefore such a \mathbf{m} does not exist, and \mathbf{l}_2 is a downward cover of \mathbf{l}_1 .

Proof for (2): It is clear that $\mathbf{l}_1 \succ \mathbf{l}_3$. Suppose there is a \mathbf{m} such that $\mathbf{l}_1 \succ \mathbf{m} \succ \mathbf{l}_3$. Then there are γ, μ , such that $\mathbf{l}_1\gamma = \mathbf{m}$ and $\mathbf{m}\mu = \mathbf{l}_3$, hence $\mathbf{l}_1\gamma\mu = \mathbf{l}_3 = \mathbf{l}_1\sigma$. Here γ only acts on variables in \mathbf{l}_1 , and μ only on variables in \mathbf{m} . Note that γ and μ can only map variables to variables, since otherwise $\mathbf{l}_1\gamma\mu = \mathbf{l}_3$ would contain more occurrences of functions or constants than \mathbf{l}_1 , contradicting $\mathbf{l}_1\sigma = \mathbf{l}_3$, since σ does not add any occurrences of function symbols to \mathbf{l}_1 .

If γ does not unify any variables in \mathbf{l}_1 , then \mathbf{l}_1 and \mathbf{m} would be variants, contradicting $\mathbf{l}_1 \succ \mathbf{m}$. If γ unifies any other variables than z and x , then we could not have $\mathbf{l}_1\gamma\mu = \mathbf{l}_3$. Hence γ must unify z and x , and cannot unify any other variables. But then $\mathbf{l}_1\gamma$ and \mathbf{l}_3 would be variants, contradicting $\mathbf{l}_1\gamma = \mathbf{m} \succ \mathbf{l}_3$. Therefore such a \mathbf{m} does not exist, and \mathbf{l}_3 is a downward cover of \mathbf{l}_1 \square

Application of the elementary substitutions on \top results in its downward covers called *most general atoms*, that consist of all n-ary predicate symbols, each with n-distinct variables as arguments. Dually, \mathbf{l}_2 is an upward cover of \mathbf{l}_1 iff \mathbf{l}_1 is a downward cover of \mathbf{l}_2 . Thus the upward covers of some conventional atom \mathbf{l}_1 are also of two types, which can be constructed by inverting the two elementary substitutions discussed in theorem 25. Trivially, every ground atom²⁹ is an upward cover of \perp .

Further, it can be proved that given two atoms \mathbf{l}_1 and \mathbf{l}_2 such that $\mathbf{l}_1 \succ \mathbf{l}_2$ ($\mathbf{l}_2 \succ \mathbf{l}_1$), there is a finite sequence of downward (upward) covers from \mathbf{l}_1 (\mathbf{l}_2) to a variant of \mathbf{l}_2 (\mathbf{l}_1). This means that if we want to get from \mathbf{l}_1 to (a variant of) \mathbf{l}_2 , we only need to consider downward (upward) covers of \mathbf{l}_1 , downward (upward) covers of downward (upward) covers of \mathbf{l}_1 , etc. In fact, there is a *finite downward cover chain algorithm* for this purpose, which is outlined in Figure 1.15.

INPUT: Conventional atoms \mathbf{l}, \mathbf{m} , such that $\mathbf{l} \succ \mathbf{m}$.
OUTPUT: A finite chain $\mathbf{l} = \mathbf{l}_0 \succ \mathbf{l}_1 \succ \dots \succ \mathbf{l}_{n-1} \succ \mathbf{l}_n = \mathbf{m}$, where each \mathbf{l}_{i+1} is a downward cover of \mathbf{l}_i .
Set $\mathbf{l}_0 = \mathbf{l}$ and $i = 0$, let θ_0 be such that $\mathbf{l}\theta_0 = \mathbf{m}$; **(1)**
if No term in θ_i contains a function or a constant; **(2) then**
 Goto 3.
else if $x/f(t_1, \dots, t_n)$ is a binding in θ_i ($n \geq 0$) **then**
 Choose new distinct variables z_1, \dots, z_n ;
 Set $\mathbf{l}_{i+1} = \mathbf{l}_i\{x/f(z_1, \dots, z_n)\}$;
 Set $\theta_{i+1} = (\theta_i \setminus \{x/f(t_1, \dots, t_n)\}) \cup \{z_1/t_1, \dots, z_n/t_n\}$;
 Set i to $i + 1$ and goto 2;
end if
if There are distinct variables x, y in \mathbf{l}_i , such that $x\theta_i = y\theta_i$ **(3) then**
 Set $\mathbf{l}_{i+1} = \mathbf{l}_i\{x/y\}$;
 Set $\theta_{i+1} = \theta_i \setminus \{x/x\theta_i\}$;
 Set i to $i + 1$ and goto 3;
else if Such x, y do not exist **then**
 Set $n = i$ and stop;
end if

Figure 1.15: Finite Downward Cover Chain Algorithm.

The subsumption ordering on atoms can be summarized through the following example:

- $l = mem(A, [A, B]) \succeq mem(1, [1, 2]) = m$ since with $\theta = \{A/1, B/2\}$, $l\theta = m$
- $mem(A1, [A1, B1]), mem(A2, [A2, B2]) \dots$ are all members of the same equivalence class

²⁹Number of ground atoms can be infinite if the language consists of a function symbol of arity ≥ 1 .

Upward Covers

To construct a finite chain of upward covers to \mathbf{l}_1 , starting from \mathbf{l}_2 , where $\mathbf{l}_1 \succ \mathbf{l}_2$, the algorithm 1.15 needs to be reversed. While algorithm 1.15 first instantiates variables to functions and constants, and then unifies some variables, the reverse algorithm "undoes" unifications and instantiations using *inverse substitution* (which is, strictly speaking, not a function).

One asymmetry of the downward and upward cases concerns the upward covers of \perp . In case of a language without constants but with at least one function symbol of arity ≥ 1 , the bottom element \perp has no upward covers at all, let alone a finite complete set of upward covers. In case of a language with at least one constant and at least one function symbol of arity ≥ 1 , there are an infinite number of conventional ground atoms, each of which is an upward cover of \perp . Together these ground atoms comprise a complete set of upward covers of \perp , but again \perp has no finite complete set of upward covers in this case. However, each conventional atom does have a finite complete set of upward covers. The top element \top does not have any upward covers at all, but it has the empty set as a finite complete set of upward covers, since no element lies "above" \top .

1.4.9 The Subsumption Theorem Again

The Subsumption Theorem holds for first-order logic, just as it did for propositional logic:

If Σ is a set of first-order clauses and D is a first-order clause. Then $\Sigma \models D$ if and only if D is a tautology or there is a clause C such that there is a derivation of C from Σ using resolution ($\Sigma \vdash_R C$) and C subsumes D .

By "derivation of a clause C " here, we mean the same as in propositional logic (page 30), that is, there is a sequence of clauses $R_1, \dots, R_k = C$ such that each R_i is either in Σ or is a resolvent of a pair of clauses in $\{R_1, \dots, R_{i-1}\}$. While extending the proof of theorem 9, the proof of this is a bit involved and we do not present it here: we refer you to [NCdW97] for a complete proof that shows that the result does indeed hold.

An immediate consequence is that the refutation-completeness of resolution follows for first-order logic as well. It is not possible, therefore, to decide, using resolution, whether a set Σ of Horn clauses is, in fact, unsatisfiable (that is, $\Sigma \models \square$: in fact, we will see later that the roots of this undecidability is more fundamental than just to do with Horn clauses or resolution).³⁰ All that we are saying with refutation-completeness is that if $\Sigma \models \square$ then there is a resolution

³⁰This gives us another difference between implication and subsumption. Unlike implication, subsumption between a pair of clauses *is* decidable, although not necessarily efficiently in all cases. We can informally show that it is decidable whether a clause C subsumes a clause D . If $C \succeq D$, then there is a substitution θ which maps each $\mathbf{l}_i \in C$ to some $\mathbf{l}_j \in D$. If C contains n literals, and D contains m literals, then there are m^n ways in which the literals in C can be paired up with literals in D . Then we can decide $C \succeq D$ by checking whether for at least one of those m^n ways of pairing the n literals in C to some of the m literals in D , there is a

proof that ends in \square . We cannot, however use an algorithm with a resolution-based theorem prover to determine if some set Σ is, in fact, unsatisfiable. What can go wrong? Well, if a \square is found, then the algorithm can terminate with “success”. But, it may end up with cases where it will be impossible to tell if the resolution process will terminate. We will see an example of this shortly.

1.4.10 Proof Strategies Once Again

Recall what we have been using so far: the derivation of a clause C from a set of clauses Σ means there is a sequence of clauses $R_1, \dots, R_k = C$ such that each C_i is either in Σ or is a resolvent of a pair of clauses in $\{R_1, \dots, R_{i-1}\}$. This remains unchanged for first-order logic, and results in the unconstrained form of a proof for C . So, just as in the propositional case, we will say that there is a *linear* derivation for C from Σ if there is a sequence $R_0, \dots, R_k = C$ such that $R_0 \in \Sigma$ and each R_i ($1 \leq i \leq k$) is a resolvent of R_{i-1} and a clause $C_i \in \Sigma \cup \{R_0, \dots, R_{i-2}\}$. By requiring side clauses to be only from Σ , we obtain (as before) the strategy called input resolution. Also, as in propositional logic, linear resolution is refutation-complete but input (and SLD) resolution for arbitrary clauses is not, except when restricted to Horn clauses.

To recap refutation, consider $P \models q$. Based on the deduction theorem:

$$P \models q \equiv P \models (q \leftarrow) \text{ iff:}$$

$$P \models (FALSE \leftarrow \sim q) \text{ iff:}$$

$$P \cup \{\sim q\} \models FALSE$$

That is $P \models q$ iff $P \cup \{\sim q\}$ is unsatisfiable. Thus, logical consequence can be checked by refutation.

We will illustrate general resolution in first order logic with an example. First of all, resolving a pair of clauses requires a substitution that unifies a pair of complementary literals.

$$\text{Let } D = \{\text{logical}(A), \neg\text{android}(A)\} \text{ and } \theta = \{A/Y\}$$

$$D\theta =$$

$$\text{The resolvent of } C, D \text{ is } E = \{\text{likes}(X, Y), \neg\text{vulcan}(X), \neg\text{android}(Y)\}$$

$$\mathbf{E} = (C - \{l\})\theta \cup (D - \{m\})\theta = (C\theta - \{l\}\theta) \cup (D\theta - \{m\}\theta) \text{ where } l\theta = \neg m\theta$$

The typical resolution steps are:

Step 0. Given a pair of clauses:

$$\begin{aligned} C_1 &: \text{likes}(\text{steve}, X) \leftarrow \text{buys}(X, \text{ilp_book}) \\ C_2 &: \text{buys}(X, \text{ilp_book}) \leftarrow \text{sensible}(X), \text{rich}(X) \end{aligned}$$

Step 1. Rename all variables apart.

θ such that $\mathbf{l}_i\theta = \mathbf{m}_j$, for each $(\mathbf{l}_i, \mathbf{m}_j)$ in the pairing. If so, there is a θ such that $C\theta \subseteq D$, and hence $C \succeq D$. If not, then there is no such θ , and $C \not\succeq D$.

$$C_1 : \text{likes}(\text{steve}, A) \leftarrow \text{buys}(A, \text{ilp_book})$$

$$C_2 : \text{buys}(B, \text{ilp_book}) \leftarrow \text{sensible}(B), \text{rich}(B)$$

Step 2. Identify complementary literals and see if mgu exists.

$$\text{buys}(B, \text{ilp_book})\theta = \text{buys}(A, \text{ilp_book})\theta$$

$$\theta = \{A/B\}$$

Step 3. Apply θ and form resolvent C .

1. Let $C_1\theta = h_1 \vee \sim l_1 \vee \sim l_2 \dots \vee \sim l_j$
2. Let $C_2\theta = l_1 \vee \sim m_1 \vee \sim m_2 \dots \vee \sim m_k$
3. Then $C = h_1 \vee \sim m_1 \vee \dots \vee \sim m_k \vee \sim l_2 \dots \vee \sim l_j$

Two questions arise with SLD-resolution: how is the negative literal to be selected from the R_{i-1} ; and if more than one clause in Σ can resolve with the selected literal, which one should be selected? We illustrate this with an example. Let Σ be the set of clauses:

$$C_0 : \forall x_1 \forall x_2 \neg \text{Grandparent}(x_1, x_2)$$

$$C_1 : \forall x \forall y \forall z (\text{Grandparent}(x, y) \vee$$

$$\quad \neg \text{Parent}(x, z) \vee \neg \text{Parent}(z, y))$$

$$C_2 : \text{Parent}(\text{henry}, \text{jane})$$

$$C_3 : \text{Parent}(\text{jane}, \text{john})$$

A little thought should convince you that $\Sigma \models \square$. We want to see if $\Sigma \vdash_{SLD} \square$. It is evident that C_0 and C_1 resolve with mgu $\{x_1/x, x_2/y\}$. The resolvent is R_1 :

$$C_0 : \forall x_1 \forall x_2 \neg \text{Grandparent}(x_1, x_2)$$

$$C_1 : \forall x \forall y \forall z (\text{Grandparent}(x, y) \vee$$

$$\quad \neg \text{Parent}(x, z) \vee \neg \text{Parent}(z, y))$$

$$R_1 : \forall u \forall v \forall w (\neg \text{Parent}(u, w) \vee \neg \text{Parent}(w, v))$$

$$C_2 : \text{Parent}(\text{henry}, \text{jane})$$

$$C_3 : \text{Parent}(\text{jane}, \text{john})$$

Since we are using SLD, one of the resolvents for the next step has to be R_1 . The other resolvent has to be one of the C_i 's. Clearly, R_1 can resolve with either C_2 or C_3 . Which one should be selected? Having selected one of C_2 or C_3 , it is clear that either of the negative literals in R_1 ($\neg \text{Parent}(u, w)$ or $\neg \text{Parent}(w, v)$) could act as the complementary literal. Which one should be selected? These two issues—which clause and which literal—are decided by a *search rule* and a *computation rule* respectively. A simple search rule is to select clauses in the order they have been shown. A simple computation rule is to select the “leftmost” literal first. In fact, these are the rules used by most PROLOG systems. With this search and computation rule, you should be able to derive the empty clause \square . Here is the input resolution diagram for this (Exercise - fill this in):

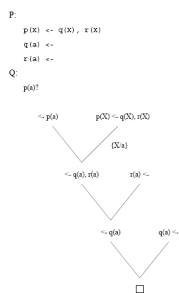


Figure 1.16: An example SLD-tree.

You should be able to draw the diagram for a different choice of search and computation rule: for example a search rule that select clauses in reverse order of appearance, and a “rightmost” literal first computation rule. It is more common, especially in the logic-programming literature, to present instead the search process confronting a SLD-resolution theorem prover in the form of a tree-diagram, called an *SLD-tree*. Such a tree effectively contains all possible derivations that can be obtained using a particular computation rule. Each node in the tree is a “goal” of the form $\leftarrow L_1, L_2, \dots, L_k$. That is, it is a clause of the form $\forall x_1 \forall x_2 \dots (\neg L_1 \vee \neg L_2 \vee \dots \vee \neg L_k)$. Given a set of clauses Σ , the children of a node in the SLD-tree are the result of resolving with clauses in Σ (nodes representing the empty clause \square have no children). SLD-trees have three kinds of branches: those that end in \square (“success” branches); those that end in goals (clauses) that cannot be resolved any further (“failure” branches); and those that continue indefinitely (infinite branches).³¹ An example SLD tree is produced in Figure 1.16

We can now visualise the effect of the proof strategy, search and computation rules. The SLD proof strategy effectively determines the nodes that can appear in the SLD-tree. Assuming some convention for drawing the tree, the effect of

³¹We do not like these.

changing the computation rules is then to alter the ordering of nodes in an SLD-tree. The search rule represents the procedural aspect of actually conducting the search (for example, in forcing a depth-first search of the tree). We can also clarify how completeness is affected by these choices. First, the refutation-completeness for Horn (and definite) clauses for SLD-resolution can be restated as meaning that if a set of clauses is unsatisfiable then there will be a leaf in the SLD-tree with the empty clause \square . It can be shown that the choice of computation rule will not alter this (informally, you can see that different computation rules will simply move the location of the \square around). If there is a \square in the tree, will a specific search rule always find it? The answer is “no”: an example is shown below, in which a search rule that does a depth-first search (enumerating clauses in the conventional leftmost first manner), will never recover from the infinite branch. This was illustrated on page 97.

1.4.11 Execution of Logic Programs

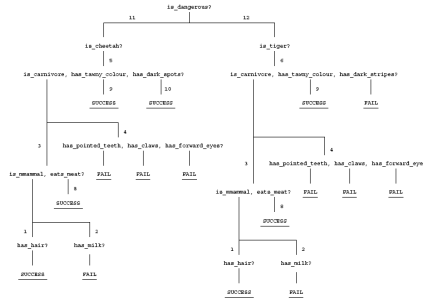
Executing definite-clause definitions can sometimes lead to *non-termination* (“infinite loops”) or even *unsound* behaviour (recall the idiosyncratic behaviour of *not/1* in the propositional case). How are logic programs executed?

1. Execution of propositional logic programs
2. Execution of programs without recursion or negation
3. Execution of programs with recursion but no negation
4. Execution of programs with recursion and negation

Searching for answers: Prolog

Consider the following program:

1. *is_mammal* \leftarrow *has_hair*
2. *is_mammal* \leftarrow *has_milk*
3. *is_carnivore* \leftarrow *is_mammal, eats_meat*
4. *is_carnivore* \leftarrow *has_pointed_teeth, has_claws, has_forward_eyes*
5. *is_cheetah* \leftarrow *is_carnivore, has_tawny_colour, has_dark_spots*
6. *is_tiger* \leftarrow *is_carnivore, has_tawny_colour, has_dark_stripes*
7. *has_hair* \leftarrow
8. *eats_meat* \leftarrow
9. *has_tawny_colour* \leftarrow
10. *has_dark_spots* \leftarrow
11. *is_dangerous* \leftarrow *is_cheetah*
12. *is_dangerous* \leftarrow *is_tiger*

Figure 1.17: The search space for query *is_dangerous?*.

Recall, based on the discussion of SLD resolution starting on page 37 that the search space for the query *is_dangerous?* will be as shown in Figure 1.17

Even with this simple Prolog program, a number of choices have to be when searching this space to see if *is_dangerous* is a logical consequence of clauses 1 – 12. As discussed earlier, research into *proof procedures* in logic programming has been concerned with searching such spaces *efficiently* keeping in mind the properties of *soundness* or *completeness*:

- Anything that is derived should be a logical consequence
- Any logical consequence should be derivable

Typically, executing a logic program involves solving queries of the form: $\mathbf{l}_1, \mathbf{l}_2, \dots, \mathbf{l}_n?$ where the \mathbf{l}_i are literals. Two problems confront us when solving this query:

1. Which literal of the \mathbf{l}_i should be solved first?
 - the rule governing this is called the *computation* rule
2. Which clause should be selected first, when more than one can be used to solve the literal selected?
 - the rule governing this is called the *search* rule

For a given program and query, the *computation* rule determines a tree of choices. The *search* rule determines the order in which this tree is searched (i.e. depth-first, breadth-first *etc.*). Computation proceeds as follows. Given a program P and a query: $\mathbf{l}_1, \mathbf{l}_2, \dots, \mathbf{l}_{j-1}, \mathbf{l}_j, \dots, \mathbf{l}_n?$

1. Use the computation rule to select \mathbf{l}_j
2. Use the search rule to select a clause: $\mathbf{l}_j \leftarrow \mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_k$ in P that can solve \mathbf{l}_j . If none found, *STOP*
3. Solve the query: $\mathbf{l}_1, \mathbf{l}_2, \dots, \mathbf{l}_{j-1}, \mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_k \dots \mathbf{l}_n?$

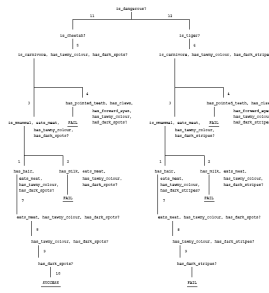


Figure 1.18: Search space for query *is_important?* using the leftmost literal first computation rule in Figure 1.17.

- As will be seen shortly, the head of the clause selected does not have to match *exactly* the literal selected. It will be enough if the two can *unify*, *i.e.* there is some substitution of variables for terms in the two literals that makes them the same.
- The step of replacing the literal selected with the literals comprising the body of the clause (Step 3) is an application of the rule of inference, *resolution*

In Figure 1.18 is the earlier query with a computation rule that selects the leftmost literal first in the query: The search rule will determine how this tree is searched for a leaf terminating in *SUCCESS* (for e.g. depth-first left-to-right, depth-first, right-to-left *etc.*)

- Different choices will affect efficiency, and sometimes even the ability to find the *SUCCESS* leaf
- Trees like this are known as *SLD-trees*: a reference to the trees obtained using a particular computation rule in conjunction with the inference rule of resolution for definite clause programs.

What about proofs for datalog programs without recursion? Consider again the example

1. $gparent(X, Z) \leftarrow parent(X, Y), parent(Y, Z)$
2. $parent(tom, jo) \leftarrow$
3. $parent(jo, bob) \leftarrow$
4. $parent(jo, jim) \leftarrow$

The rightmost literal first computation rule yields the proof tree in Figure 1.19 What about Datalog programs with recursion? Consider the following program:

1. $less_than(X, Y) \leftarrow succ(Y, X)$
2. $less_than(X, Y) \leftarrow succ(Z, Y), less_than(Z, Y)$

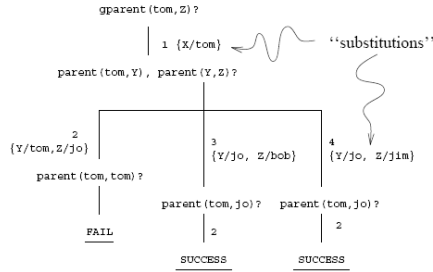


Figure 1.19: Search space for query $gparent(tom, Z)?$ using the rightmost literal first computation rule.

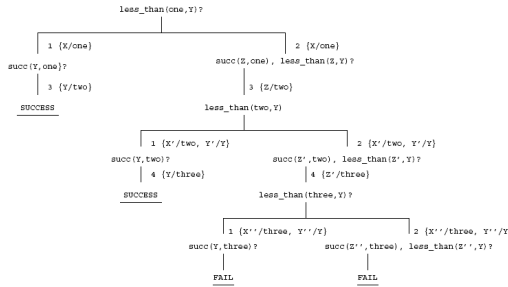


Figure 1.20: Search space for query $less_than(one, Y)?$ using the leftmost literal first computation rule.

- 3. $succ(two, one) \leftarrow$
- 4. $succ(three, two) \leftarrow$

Leftmost literal rule for the query $less_than(one, Y)$ yield the search tree in Figure 1.20

What about completeness with respect to the choice of computation and search rules? One way to search the trees obtained so far is depth-first, left-to-right. Since clauses that appear first (reading top to bottom) in the program have been drawn on the left, this search rule selects clauses in order of appearance in the program. In fact, most logic programs are executed using the following:

Computation rule. Leftmost literal first

Search rule. Depth first search for clauses in order of appearance

However, will a logic-programming system with an arbitrary computation rule, and a depth-first search of clauses in some fixed order always find a leaf terminating in *SUCCESS* (if one exists)? The answer to this is a **No**.

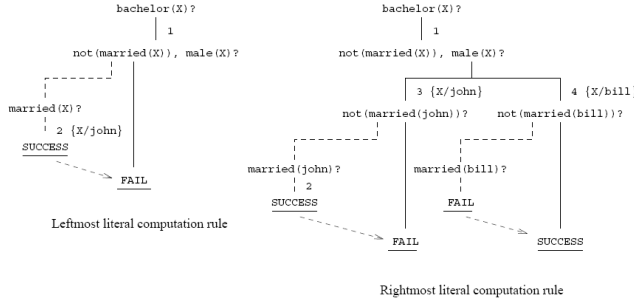


Figure 1.21: Two proof trees with different computation rules for the query $bachelor(X)?$.

Finite failure: programs with negation

Consider the following program:

1. $bachelor(X) \leftarrow not(married(X)), male(X)$
2. $married(john) \leftarrow$
3. $male(john) \leftarrow$
4. $male(bill) \leftarrow$

Two trees for the query $bachelor(X)?$ with different computation rules are outlined in Figure 1.21. where, the query $not(q)?$ succeeds iff $q?$ “finitely fails”. A finitely failed tree is such that it has finite depth, finite depth and all computations end in *FAIL*. As discussed for propositional logic in Section 1.22, the query $not(q)?$ succeeds iff $q?$ “finitely fails”. A finitely failed tree is a derivation tree of finite depth, finite depth and for which, all computations end in *FAIL*. Given a program P , a computation rule R and a search rule S , one could think of all the ground queries that can be asked of P . For *e.g.* with the program with $less_than/2$ and $succ/2$ clauses, these are of the form: $less_than(one,two)?$, $less_than(two,one)? \dots, succ(one,two)? \dots$. Such queries fall into 3 categories:

1. Those which are answered *yes* because a *SUCCESS* branch is found
2. Those which are answered *no* because the query finitely fails
3. Those that are neither in categories 1 or 2

The three categories of ground clauses are depicted in Figure 1.22.

1.4.12 Answer Substitutions

The example we have just seen is a simple form of the kind of refutation-based theorem proving using resolution adopted by systems based on logic programming. There, we are typically interested in “answer substitutions” for a query, given some set of definite clauses Σ (recall that these are a special kind of Horn

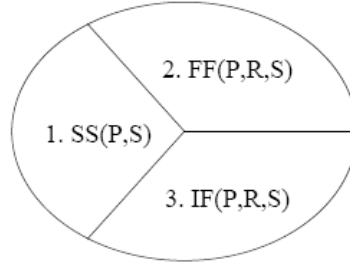


Figure 1.22: Three categories of ground clauses.

clause, in which there is exactly one positive literal). That is, we want to know if there are any substitutions for x_1, \dots, x_n such that $Q(x_1, \dots, x_n)$ is a logical consequence of Σ ? That is, we are seeking substitutions for x_1, \dots, x_n such that:

$$\Sigma \models \exists x_1 \exists x_2 \dots \exists x_n Q(x_1, x_2, \dots, x_n)$$

holds. Substitutions for the x_1, \dots, x_n are then said to be *correct* answers.

Using $\exists \mathbf{x}Q(\mathbf{x})$ to denote the formula on the right, we can apply the deduction theorem in manner shown on page 25:

$$\Sigma \models \exists \mathbf{x}Q(\mathbf{x}) \text{ if and only if } \Sigma \cup \{\neg \exists \mathbf{x}Q(\mathbf{x})\} \models \square$$

Since $\neg \exists x \alpha \equiv \forall x \neg \alpha$ (page 66):

$$\Sigma \models \exists \mathbf{x}Q(\mathbf{x}) \text{ if and only if } \Sigma \cup \{\forall \mathbf{x} \neg Q(\mathbf{x})\} \models \square$$

Now, Σ is a set definite clauses and $\forall \mathbf{x} \neg Q(\mathbf{x})$ is a Horn clause. Taking as an article of faith what we said earlier about the refutation completeness of SLD-resolution: $\Sigma \cup \{\forall \mathbf{x} \neg Q(\mathbf{x})\} \models \square$ then $\Sigma \cup \{\forall \mathbf{x} \neg Q(\mathbf{x})\} \vdash_{SLD} \square$. This provides one way of determining *computed* answers for x_1, \dots, x_n such that $Q(x_1, \dots, x_n)$ is a logical consequence of a set of Horn clauses Σ ? Answers are obtained by:

1. Adding $\forall \mathbf{x} \neg Q(\mathbf{x})$ to Σ . The result is a set of clauses;
2. Finding all substitutions for x_1, \dots, x_n that result in a derivation of \square using SLD-resolution.

Are computed answers always correct? If some care is taken to ensure there are no variables shared between queries and clauses in Σ before commencing an SLD derivation, the answer to this question is “yes”. (In fact, if we are only in checking for a proof for \square , then these additional precautions are not needed.) We know, of course, that resolution is refutation-complete. And we also know that the minimal model of Σ , or $MM(\Sigma)$ is identical to the ground atomic logical consequences of Σ . That is, $MM(\Sigma)$ consists of all ground atoms q such

that $\Sigma \models q$. The set of ground atoms q for which $\Sigma \cup \{-q\} \vdash_{SLD} \square$ is called the *success set* of Σ , or $SS(\Sigma)$. You should now be able to convince yourself that $SS(\Sigma) = MM(\Sigma)$.

The resolution procedure as described here has a limitation concerned with term evaluation

- Consider a function $sqr/1$ that accepts a natural number and returns its square
- The mgu algorithm cannot unify $p(sqr(2))$ and $p(2)$

Extensions are possible to overcome this

- Resolution with “paramodulation” performs term rewrites to achieve this
- But, logic programming systems use a special predicate that forces term evaluation
- Thus, $p(sqr(2))$ is usually written as $X \text{ is } 2 * 2, p(X)$. $p(X)$ unifies with $p(4)$ after forced evaluation the value of X by the $is/2$ predicate

Finally, the extension of SLD to negation-as-failure requires some care. SLDNF-resolution for first-order clauses is only sound if we are checking the proof of a ground literal.

1.4.13 Theorem Proving for Full First-Order Logic

1.5 Adequacy

We now turn to the questions of what can, and cannot be achieved by using a representation based on Horn clauses, and inference based on SLD-resolution. First of all, the representation problem. Based on the work of Stephen Kleene, we know that the class of *partial recursive functions* is effectively computable. By the latter, we mean that these functions can be calculated by a Turing machine. By the former, we mean a function (more on functions later) that takes a fixed number of natural numbers as arguments and returns a single natural number as a result (the “partial” part comes from the fact that the function may not be defined for all values of its arguments). In fact, Kleene actually proved that the class of partially recursive functions is identical to the class of functions computable by a Turing machine, which are the only problems solvable by a computer. Now, it has been shown that every partial recursive function can be represented by a definite clause program, and every evaluation of such a function can be encoded as a query in the form we have seen earlier. So, in principle at least, the language of Horn clauses is adequate for representing all computable problems.

But representability does not mean convenient representability. It is for the latter reasons it may be preferable to use different logics, or even other programming languages. We have had a glimpse earlier of how theorem-proving for first-order logic can be accomplished by converting these statements to Horn clauses. So we know that such conversions can be done, and attention turns to the soundness and completeness of the inference mechanisms which we have

already considered—specifically, SLD-based resolution—for Horn clause programs. To summarise, ignoring for the moment, the business of search and computation rules:

- Without the “occurs”-check (page 77), SLD resolution can be unsound. In practice logic programming systems like PROLOG omit this check, hoping that the situations where it was needed would not come up.
- The extension of SLD resolution with negation-as-failure can be unsound if attempting to prove atoms that are non-ground.
- SLD-resolution is refutation-complete for Horn clauses.
- SLD-resolution is not affirmation-complete for Horn clauses.
- For general clauses, resolution is refutation-complete if we perform factoring.
- For general clauses, resolution is not affirmation-complete.

When we add on the constraints imposed by the search rule, we lose even refutation-completeness with some kinds of “unfair” search rules (we saw an example of this earlier, when a depth-first search of the SLD-tree went down an infinite path). But, assuming we are using a fair rule, we can confirm the validity any sentence that *is* valid. On the other hand, we cannot deal with problems that are unsolvable, or reliably identify an implication that does not hold. That is, the property of validity is only *semi-decidable*. This is not especially because of the use of Horn clauses or SLD-resolution, but because it is an unavoidable aspect of first-order logic. It is known, for example, from a result due to Schmidt-Schauss in 1988 that implication between a pair of general clauses in first-order logic is undecidable (that is, we will never be able to construct a procedure that will terminate in all cases with the answer). In 1992, it was shown by Marcinkowski and Pacholski that even implication between a pair of Horn clauses is undecidable.

Chapter 2

Exploring a Structured Search Space

ILP is concerned with the automatic construction of “general” logical statements from “specific” ones. For example, given $mem(1, [1, 2]) \leftarrow construct\ mem(A, [A|B]) \leftarrow$. What do the words “general” and “specific” mean in a logical setting? Can statements of increasing (decreasing) generality be enumerated in an orderly manner? These are questions about the mathematics of “generality” ILP identifies “generality” with \models . That is, C_1 is “more general” than C_2 iff $C_1 \models C_2$. The relation \models results in a quasi-ordering over a set of clauses. ILP systems are programs that search such quasi-ordered sets.

The result (theorem 24) that $\langle \mathcal{A}_E^+, \succeq \rangle$ is a lattice shows that the set of atoms is well structured. The more structured a set is, the better it is suited to be searched for candidates to include in a theory. This search usually proceeds by small upward steps (generalization) or downward steps (specialization) in the lattice. In this chapter, we will first dwell upon the subsumption lattice over clauses. If we want to generalize or specialize a set of clauses to a single clause, we can use a least generalization or greatest specialization of this set. On the other hand, we may also want to generalize or specialize an individual clause to another individual clause using the idea of covers - the smallest non-trivial steps between individual clauses that we can take in the lattice.

Subsumption is the generality order that is used most often in ILP. It is used much more than logical implication. The reasons for this are mainly practical: subsumption is more tractable and more efficiently implementable than implication. For instance, subsumption between clauses is decidable whereas implication is not (Section 1.4.9).

2.1 Subsumption Lattice over Clauses

Let C and D be clauses. We say $C \succeq D$ iff C subsumes D and $C \succ D$ iff C properly subsumes D . Clearly, the \succeq relation on clauses is reflexive and

transitive. Thus it imposes a quasi-order on the set of clauses. Since \succeq is a quasi-order, we know a partial ordering must result from the partition of the set of clauses \mathcal{C} into a set of equivalence classes¹ \mathcal{C}_E .

An example of subsumption ordering on clause-sets is produced below:

$$\begin{array}{c} \{mem(A, [A|B]) \leftarrow, mem(A, [B, A|C]) \leftarrow\} \\ \succeq \\ \{mem(1, [1, 2]) \leftarrow, mem(2, [1, 2]) \leftarrow\} \end{array}$$

A clause is always subsume-equivalent with a clause that does not contain literals more than once. So for example $P(x) \vee Q(a) \vee P(x)$ is obviously subsume-equivalent with $P(x) \vee Q(a)$. Similarly, the order of literals in a clause does not matter much. For instance, $P(a) \vee P(b) \equiv P(b) \vee P(a)$. This amounts to treating a clause as a set of literals, instead of a disjunction of literals. Thus we may use the set $\{P(a), Q(x)\}$ to represent the clauses $\{Q(x) \vee P(a), P(a) \vee Q(x) \vee P(a), Q(x) \vee P(a) \vee Q(x)\}$, etc.

However, the above condition does not give a sufficient condition for subsume-equivalence. While two atoms are subsume-equivalent iff they are variants, this is not true for clauses in general. For instance, $C = \{P(x, x)\} \equiv \{P(x, x), P(x, y)\} = D$, since $C \subseteq D$ and $D\{y/x\} \equiv C$, yet C and D are not variants. In fact, the subsume-equivalence class of this C contains an infinite number of clauses which are not variants. For example, for each n , the clause $D_n = \{P(x, x), P(x, x_1), P(x_1, x_2), \dots, P(x_{n_1}, x_n)\}$ is subsume-equivalent with $C = \{P(x, x)\}$.

While subsume-equivalent clauses need not be variants, *reduced* subsume-equivalent clauses are however variants. A clause C is said to be *reduced* if there is no proper subset D of C ($D \subset C$) such that $C \equiv D$. Equivalently, a clause C is *reduced* if there is no substitution θ such that $C\theta$ is a proper subset of C . A reduced clause D such that $C \equiv D$ and $D \subseteq C$ is called a *reduction* of C . Thus, $C = \{P(x, y), P(y, x)\}$ is reduced whereas $D = \{P(x, x), P(x, y), P(y, x)\}$ is not reduced, since $D' = \{P(x, x)\} \subset D$ and $D \equiv D'$ is a reduction of D . Note that a subset of a reduced clause need not be reduced itself; if $C = \{\neg Q(x, a), \neg Q(y, a)\}$ and $D = \{P(x, y), \neg Q(x, a), \neg Q(y, a)\}$ then D is reduced, while $C \subset D$ is not reduced, since $C\{x/y\}$ is a proper subset of C . Thus, a reduced clause is a canonical member of the subsume-equivalence class. It can be obtained using the algorithm outlined in Figure 2.1. The algorithm itself is based on the following theorem:

Theorem 26 *Let C be a clause. If for some θ , $C\theta \subseteq C$, then there is a reduced clause $D \subseteq C\theta$ such that $C \equiv D$.*

Proof: Let $C_1 = C\theta$. Since $C \models C\theta$ and $C\theta \models C$, clearly $C \equiv C_1$. If C_1 is reduced, then let $D = C_1$, and we are done. Otherwise, there is a substitution

¹Note that applying the same substitution θ to two subsume-equivalent clauses may yield two clauses which are no longer subsume-equivalent. For example, if $C = \{P(x, y), P(z, u)\}$ and $D = \{P(x, y)\}$, then $C \equiv D$. Let $\theta = \{y/f(x), z/f(x), u/x\}$. Then $C\theta = \{P(x, f(x)), P(J(x), x)\}$ and $D\theta = \{P(x, f(x))\}$, which are no longer equivalent.

θ_1 such that $C_2 = C_1\theta_1 \subset C_1$. So C_2 is a proper subset of C_1 which is subsume-equivalent to C_1 . Since $C_1 \equiv C$, we also have $C_2 \equiv C$, in fact $C\theta\theta_l = C_2 \subset C$. If C_2 is still not reduced, we can go on defining $C_3 = C_2\theta_2 \subset C_2$, etc. Since C only contains a finite number of literals, this cannot go on indefinitely. Hence we must arrive at a $D = C_n$ such that D is reduced and $C \equiv D$. \square

INPUT: A clause C .
OUTPUT: A reduction D of C .
Set $D = C$, $\theta =$;
repeat
 Set D to $D\theta$;
 Find a literal $\mathbf{l} \in D$ and a substitution θ such that $D\theta \subseteq D \setminus \{\mathbf{l}\}$;
until Such a (\mathbf{l}, θ) does not exist;
return D .

Figure 2.1: Plotkin's reduction algorithm.

2.1.1 Lattice Structure of Clauses

It can be proved that every finite set \mathcal{S} of (general or horn) clauses has a greatest lower bound under subsumption (also called greatest specialization under subsumption or GSS in the ILP literature), as well as a least upper bound under subsumption (also called least generalization under subsumption or LGG in the ILP literature). This holds both for clausal languages \mathcal{C} , and for Horn languages \mathcal{H} .

Greatest Specialization Under Subsumption

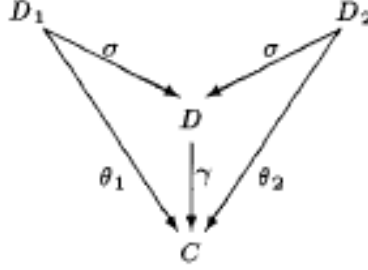
It is straightforward to show that the GSS of some finite set \mathcal{S} of clauses in \mathcal{C} is simply the union of all clauses in \mathcal{S} after they are standardized apart (EXERCISE). Proving the existence of a GSS of every finite set of Horn clauses in \mathcal{H} requires a little more work. We will assume that \mathcal{H} contains an artificial bottom element² \perp , such that $C \succeq \perp$ for every $C \in \mathcal{H}$, and $\perp \succeq C$ for every $C \in \mathcal{H}$.

Theorem 27 *Let \mathcal{H} be the Horn language \mathcal{H} , with an additional bottom element $\perp \in \mathcal{H}$. Then for every finite non-empty $\mathcal{S} \subseteq \mathcal{H}$, there exists a GSS (glb) of \mathcal{S} in \mathcal{H} .*

Proof: Suppose $\mathcal{S} = \{D_1, \dots, D_n\} \subseteq \mathcal{H}$. Without loss of generality we assume the clauses in \mathcal{S} are standardized apart, D_1, \dots, D_k are the definite program clauses in \mathcal{S} , and D_{k+1}, \dots, D_n are the definite goals in \mathcal{S} .

1. If $k = 0$ (i.e., if \mathcal{S} only contains goals), then it is easy to show that $D_1 \cup \dots \cup D_n$ is a GSS of \mathcal{S} in \mathcal{H} .

²Note that \perp is not subsume-equivalent with other tautologies. Two tautologies need not be subsume-equivalent either.

Figure 2.2: D is a GSS of D_1 and D_2 .

2. If $k \geq 1$ and the set $\{D_1^+, \dots, D_k^+\}$ (the set of heads of clauses in S), is not unifiable, then \perp is a GSS of S in \mathcal{H} .
3. Otherwise, let σ be an mgu for $\{D_1^+, \dots, D_k^+\}$, and let $D = D_1\sigma \cup \dots \cup D_n\sigma$ (note that actually $D_i\sigma = D_i$ for $k+1 \leq i \leq n$, since the clauses in S are standardized apart). Since D has exactly one literal in its head, it is a definite program clause. Furthermore, we have $D_i \succeq D$ for every $1 \leq i \leq n$, since $D_i\sigma \subseteq D$.

To show that D is a GSS of S in \mathcal{H} , suppose $C \in \mathcal{H}$ is some clause such that $D_i \succeq C$ for every $1 \leq i \leq n$. For every $1 \leq i \leq n$, let θ_i be such that $D_i\theta_i \subseteq C$, and θ_i only acts on variables in D_i . Let $\theta = \theta_1 \cup \dots \cup \theta_n$. For every $1 \leq i \leq k$, $D_i^+\theta = D_i^+\theta_i = C^+$, so θ is a unifier for $\{D_1^+, \dots, D_k^+\}$. But σ is an mgu for this set, so there is a γ such that $\theta = \sigma\gamma$. Now $D\gamma = D_1\sigma\gamma \cup \dots \cup D_n\sigma\gamma = D_1\theta \cup \dots \cup D_n\theta = D_1\theta_1 \cup \dots \cup D_n\theta_n \subseteq C$. Hence $D \succeq C$, so D is a GSS of S in \mathcal{H} . See Figure 2.2 for illustration of the case where $n = 2$.

For example, $D = P(a) \leftarrow P(f(a)), Q(y)$ is a GSS of $D_1 : P(x) \leftarrow P(f(x))$ and $D_2 : P(a) \leftarrow Q(y)$. Note that D can be obtained by applying $\sigma = \{x/a\}$ (the mgu for the heads of D_1 and D_2) to $D_1 \cup D_2$, the GSS of D_1 and D_2 in C .

□

Least Generalization Under Subsumption

Proving the existence of the least generalization LGG (lub) is a little harder. Let C and D be clauses. A *selection* of C and D is a pair of compatible literals (\mathbf{l}, \mathbf{m}) , such that $\mathbf{l} \in C$, $\mathbf{m} \in D$. Two literals are said to be *compatible* if they have the same sign and predicate symbol. For example, $C = \{P(x), P(y), \dots, P(a)\}$ and $D = \{Q(b), P(a), \dots, P(b)\}$ have three selections: $(P(x), P(a))$, $(P(y), P(a))$, and $(\neg P(a), \neg P(b))$. Given two clauses C and D , there is only a finite number of selections. Suppose C and D have a total of n selections. Then we can

order these in a sequence $S = (\mathbf{l}_1, \mathbf{m}_1), (\mathbf{l}_2, \mathbf{m}_2), \dots, (\mathbf{l}_n, \mathbf{m}_n)$, and construct two compatible ordered clauses $C_S = \mathbf{l}_1 \vee \dots \vee \mathbf{l}_n$ and $D_S = \mathbf{m}_1 \vee \dots \vee \mathbf{m}_n$. Treating $C_S = \vee(\mathbf{l}_1, \dots, \mathbf{l}_n)$ and $D_S = \vee(\mathbf{m}_1, \dots, \mathbf{m}_n)$ as atoms, it can be shown that the lub of C_S and D_S constructed using the anti-unification algorithm on page 84 is also an LGS of C and D . As an example, if

- $G = \{\mathbf{l}_1, \mathbf{l}_2, \mathbf{l}_3\}$, for $\mathbf{l}_1 = P(f(a), f(x))$, $\mathbf{l}_2 = P(f(x), g(a))$, $\mathbf{l}_3 = Q(a)$
- $d = \{\mathbf{m}_1, \mathbf{m}_2\}$, for $\mathbf{m}_1 = P(f(b), x)$, $\mathbf{m}_2 = P(y, g(b))$.

then, the set of all selections of C and D can be ordered in the following sequence:

$$S = (\mathbf{l}_1, \mathbf{m}_1), (\mathbf{l}_1, \mathbf{m}_2), (\mathbf{l}_2, \mathbf{m}_1), (\mathbf{l}_2, \mathbf{m}_2)$$

leading to the constructing of the following clauses:

$$C_S = \vee(\mathbf{l}_1, \mathbf{l}_1, \mathbf{l}_2, \mathbf{l}_2)$$

$$D_S = \vee(\mathbf{m}_1, \mathbf{m}_1, \mathbf{m}_2, \mathbf{m}_2)$$

Note that the order and duplication of literals is not ignored in the atomic order. The clauses C_S and D_S are compatible, and have the following lub:

$$E_S = P(f(z_1), z_2) \vee P(z_3, z_4) \vee P(f(z_5), z_6) \vee P(z_7, g(z_1))$$

This lub of C_S and D_S can be shown to be also an LGS of C and D and can be reduced to the LGS $E = \{P(f(z_1), z_2), P(z_7, g(z_1))\}$. Note that the predicate Q does not appear in the lub or LGS.

Theorem 28 *Let \mathcal{C} be a clausal language. Let $C, D \in \mathcal{C}$ be clauses, and S be a sequence of all selections of C and D . Then an $\text{lub}(C_S, D_S)$ is an LGS of $\{C, D\}$.*

Proof: First of all if C and D are clauses, and S a sequence of (not necessarily all) selections of C and D . Then $\text{lub}(C_S, D_S) \succeq C$ and $\text{lub}(C_S, D_S) \succeq D$. This is easy to see. If $E_S = \text{lub}(C_S, D_S)$, then $E \succeq C_S$. Also, $C_S \succeq C$, since the literals in C_S form a subset of C . Hence $E \succeq C$, by the transitivity of \succeq . Similarly $E \succeq D$.

Let $F = \{\mathbf{l}_1, \dots, \mathbf{l}_m\}$ be a clause such that $F \succeq C$ and $F \succeq D$. In order to establish that E is an LGS of $\{C, D\}$, we need to prove $F \succeq E$. Since $F \succeq C$ and $F \succeq D$, there are θ_1 and θ_2 , and $\mathbf{l}_1, \dots, \mathbf{l}_m \in C$ and $\mathbf{m}_1, \dots, \mathbf{m}_m \in D$, such that $\mathbf{l}_i \theta_1 = \mathbf{l}_i$, and $\mathbf{l}_i \theta_2 = \mathbf{m}_i$, for every $1 \leq i \leq m$. Then $S' = (\mathbf{l}_1, \mathbf{m}_1), \dots, (\mathbf{l}_m, \mathbf{m}_m)$ is a sequence of selections of C and D . Let $C_{S'} = \vee(\mathbf{l}_1, \dots, \mathbf{l}_m)$, $D_{S'} = \vee(\mathbf{m}_1, \dots, \mathbf{m}_m)$, let $G = \vee(\mathbf{k}_1, \dots, \mathbf{k}_m)$ be $\text{lub}(C_{S'}, D_{S'})$, and σ_1 and σ_2 be such that $G\sigma_1 = C_{S'}$ and $G\sigma_2 = D_{S'}$. Since $\vee(\mathbf{l}_1, \dots, \mathbf{l}_m)\theta_1 = C_{S'}$ and $\vee(\mathbf{l}_1, \dots, \mathbf{l}_m)\theta_2 = D_{S'}$, there must be a γ such that $\vee(\mathbf{l}_1, \dots, \mathbf{l}_m)\gamma = \vee(\mathbf{k}_1, \dots, \mathbf{k}_m)$. Thus we have the situation given in Figure 2.3.

$\vee(\mathbf{l}_1, \dots, \mathbf{l}_m)\gamma = G$, so we have $F \succeq G$. Since every selection in S' also occurs in S , we can prove³ that $G \succeq E$. Hence $F \succ E$. \square

Thus the LGS of any two clauses exists, and can be computed by the method made explicit in the algorithm in Figure 2.4

³We leave this an EXERCISE. You will need to iron out duplications from S and S' and permute them so that S' is a prefix of S .

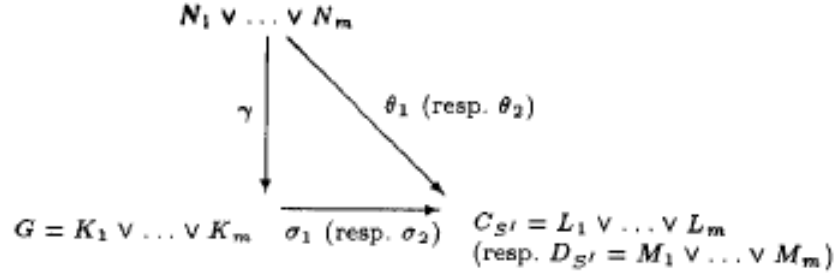


Figure 2.3: Illustration of the proof of theorem 2.3.

INPUT: Two clauses C and D ;
OUTPUT: An LGS of $\{C, D\}$;
 Let $(\mathbf{l}_1, \mathbf{m}_1), \dots, (\mathbf{l}_n, \mathbf{m}_n)$ be a sequence of all selections of C and D ;
 Obtain $\text{lub}(\vee(\mathbf{l}_1, \dots, \mathbf{l}_n), \vee(\mathbf{m}_1, \dots, \mathbf{m}_n)) = \vee(\mathbf{l}_1, \dots, \mathbf{l}_n)$ from the Anti-Unification Algorithm;
return $\{\mathbf{l}_1, \dots, \mathbf{l}_n\}$;

Figure 2.4: The LGS Algorithm.

The LGS of any finite set of clauses can be computed by repeatedly applying this algorithm. Notice that if two clauses C and D have no selections for instance, when they have no predicates in common - then their LGS is the empty clause \square . Thus \square can play the role of top element here, which means that we do not need to add an artificial top element \top to the language \mathcal{C} .

Note that if all literals in C and D have the same sign and predicate symbol, then C and D have $|C| \cdot |D|$ selections. Accordingly, the LGS of C and D that can be obtained from these selections may also contain $|C| \cdot |D|$ distinct literals. Thus the number of literals in an LGS may increase quite rapidly.

Since we have now proved the existence of a GSS and LGS of every two clauses, it follows that a clausal language ordered by subsumption has a lattice-structure (we do not need an artificial bottom element \perp for this).

Theorem 29 *Let \mathcal{C} be a clausal language. Then $\langle \mathcal{C}_E, \succeq \rangle$ is a lattice. Further, if \mathcal{H} be the horn clause language, including an addition symbol \perp , then $\langle \mathcal{H}_E, \succeq \rangle$ is also a lattice.*

Proof: The proof for (\mathcal{C}_E, \succeq) being a lattice follows naturally from the fact that lub (LGS) and glb (GSS) exist for every pair of clauses from \mathcal{C} .

As for \mathcal{H}_E , since there is at most one selection possible from the heads of a set of Horn clauses \mathcal{H} , the LGS of \mathcal{H} has at most one positive literal, and hence is itself also a Horn clause. Therefore (\mathcal{H}_E, \succeq) is a lattice. Here we need the bottom element \perp to guarantee the existence of a GSS of two definite program clauses with different predicate symbols in their head.

Thus, the p.o. set of equivalence classes of atoms $\mathcal{C}_E(\mathcal{H}_E)$ is a lattice with the binary operations \sqcap and \sqcup defined on elements of $\mathcal{C}_E(\mathcal{H}_E)$ as follows (note that $\perp \in \mathcal{C}$ is the empty set \emptyset):

- $[\perp] \sqcap [C] = [\perp]$, and $[\top] \sqcap [C] = [C]$
- If $C_1, C_2 \in \mathcal{C}$ ($C_1, C_2 \in \mathcal{H}$) have a *GSS* (glb) D then $[C_1] \sqcap [C_2] = [D] = [C_2\theta]$ otherwise $[C_1] \sqcap [C_2] = [\perp]$
- $[\perp] \sqcup [C] = [C]$, and $[\top] \sqcup [C] = [\top]$
- If C_1 and C_2 have *LGS* (lub) M then $[C_1] \sqcup [C_2] = [M]$ otherwise $[C_1] \sqcup [C_2] = [\top]$

For \mathcal{C} ,

- the *GSS* is simply the union $C_1 \cup C_2$ (union translates to ‘or’ing the two clauses)
- the *LGS* of C_1 and C_2 is obtained using the LGS algorithm outlined in Figure 2.4

whereas, for \mathcal{H} ,

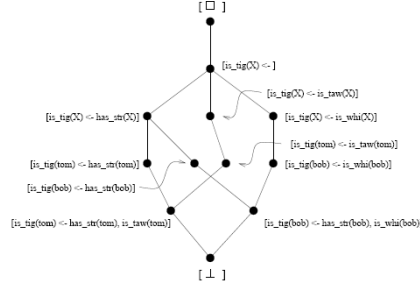
- the *GSS* (or meet operation) is given as follows: If the set of positive literals (heads) in $C_1 \cup C_2$ have an mgu θ , then $GSS(C_1, C_2) = (C_1 \cup C_2)\theta$. Otherwise $GSS(C_1, C_2) = \perp$
- the *LGS* of C_1 and C_2 is obtained using the LGS algorithm outlined in Figure 2.4.

□

As an example, if

$$\begin{aligned} \mathcal{H} = \{ & \square, \perp, \\ & is_tiger(tom) \leftarrow has_stripes(tom), is_tawny(tom) , \\ & is_tiger(bob) \leftarrow has_stripes(bob), is_white(bob) , \\ & is_tiger(tom) \leftarrow has_stripes(tom) , \\ & is_tiger(tom) \leftarrow is_tawny(tom) , \\ & is_tiger(bob) \leftarrow has_stripes(bob) , \\ & is_tiger(bob) \leftarrow is_white(tom) , \\ & is_tiger(X) \leftarrow has_stripes(X) , \\ & is_tiger(X) \leftarrow is_tawny(X) , \\ & is_tiger(X) \leftarrow is_white(X) , \\ & is_tiger(X) \leftarrow \} \end{aligned}$$

then the diagram of p.o. set \mathcal{H}_E is as outlined in Figure 2.5

Figure 2.5: The lattice of \mathcal{H}_E for the *is_tiger* relation.

2.1.2 Covers in the Subsume Order

The least generalization and the greatest specialization respectively concern generalizing or specializing a set of clauses to a single clause. What about generalizing and specializing single clauses? We will address this issue by investigating *covers* of clauses in the subsumption order.

It can be shown that there exist clauses which have no complete set of upward covers in the subsumption order. In fact, there are clauses which have no upward covers at all. For example, the clause $C = \{P(x_1, x_1)\}$ has no upward cover in $\langle C, \succeq \rangle$. This can be proved by defining

$$C_n = \{P(x_i, x_j) \mid i \neq j \text{ and } 1 \leq i, j \leq n\}, \quad n \geq 2$$

and noting that each C_n properly subsumes C_{n+1} and C . That is $C_2 \succ C_3 \succ \dots C_n \succ \dots C$. From this result, we know that C has no complete set of upward covers.

Dually, for the downward cover there exists a clause C which has no finite complete set of downward covers⁴. So if this particular C does have a complete set of downward covers, this set must be infinite. This result is sufficient to later prove the negative result that an ideal downward refinement operator does not exist for the subsumption order. One such C is

$$C = \{P(x_1, x_2), P(x_2, x_1)\}$$

If we define⁵

$$E_n = \{P(y_1, y_2), P(y_2, y_3), \dots, P(y_{n-1}, y_n), P(y_n, y_1)\}, \quad n \geq 2$$

and

$$C_n = C \cup E_n, \quad n \geq 3$$

⁴Whether all clauses have a (sometimes infinite) complete set of downward covers remains an open question.

⁵The clauses E_n have a special structure called a *cycle*.

then it can be shown that for any $n = 3^k$ ($k \geq 1$), $C \succ C_n$. Further, by contradiction, one can show that there is no downward cover D of C , such that $D \succeq C_{3^k}$ for every $k \geq 1$ and subsequently, that, C has no finite complete set of downward covers in $\langle \mathcal{C}, \succeq \rangle$.

2.2 The Implication Order

It is easy to see that implication is reflexive and transitive, and hence a quasi order. Implication between (Horn) clauses is undecidable, but using implication as a generality order is more desirable than using subsumption, for some important reasons:

1. It is better able to deal with recursive clauses. A clause C which implies another clause D , need not subsume this D . For instance, take

$$\begin{aligned} C &= P(f(x)) \leftarrow P(x) \\ D &= P(f^2(x)) \leftarrow P(x) \end{aligned}$$

Then $C \models D$, but $C \not\subseteq D$. Subsumption is too weak in this case. A further sign of this weakness is the fact that two tautologies need not be subsume equivalent, even though they are logically equivalent.

2. For the construction of least generalizations, subsumption is again not fully satisfactory and can over-generalize. For example,
 - If S consists of the clauses $D_1 = P(f^2(a)) \leftarrow P(a)$ and $D_2 = P(f(b)) \leftarrow P(b)$, then the LGS of S is $P(f(y)) \leftarrow P(x)$. the other hand, the clause $P(f(x)) \leftarrow P(x)$ seems more appropriate as a least generalization of S , since it implies D_1 and D_2 , and is implied by the LGS. However, it does not subsume D_1 .
 - Even for clauses without function symbols, the subsumption order may still be unsatisfactory. Consider $D_1 = P(x, y, z) \leftarrow P(y, z, x)$ and $D_2 = P(x, y, z) \leftarrow P(z, x, y)$. The clause D_1 is a resolvent of D_2 with D_2 , and D_2 is a resolvent of D_1 with D_1 ; so D_1 and D_2 are logically equivalent. This means that D_1 is a least generalization under implication (LGI) of the set $\{D_1, D_2\}$. Yet the LGS of these two clauses is $P(x, y, z) \leftarrow P(u, v, w)$, which is clearly an over-generalization.

As these examples also show, the subsumption order is particularly unsatisfactory if we consider recursive clauses: clauses where the same predicate symbol occurs both in a positive and a negative literal. Thus it is desirable to make the step from the subsumption order to the more powerful implication order.

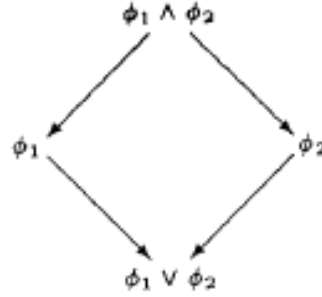


Figure 2.6: Least generalization and greatest specialization in first-order logic.

3. A further advantage of the implication order is that one can easily compare a set of clauses (a theory) with another theory or clause. For example, if $\Sigma = \{(P \leftarrow Q), (Q \leftarrow R)\}$ and $C = P \leftarrow R$, then we have $\Sigma \models C$. On the other hand, subsumption cannot be used here to compare the generality of Σ and C , because neither member of Σ subsumes C .

The main results on the implication order can be summed up as follows:

1. If α_1 and α_2 are two arbitrary first-order formulas, then it can be easily shown that their least generalization under implication (LGI) is just $\alpha_1 \wedge \alpha_2$, and their greatest specialization under implication (GSI) is just $\alpha_1 \vee \alpha_2$. See Figure 2.6. However, if α_1 and α_2 are clauses, $\alpha_1 \wedge \alpha_2$ is not a clause. Instead of using $\alpha_1 \wedge \alpha_2$, we have to find some least clause as LGI, which implies both clauses α_1 and α_2 . Such a clause appears quite hard to find sometimes, as we will see subsequently. However, $\alpha_1 \vee \alpha_2$ remains a clause and remains the GSI under the subsumption order on clauses.
2. Every finite set of clauses Σ which contains at least one function-free nontautologous clause, has a computable least generalization (LGI) under implication in \mathcal{C} . A related observation is that the problem whether $\Sigma \models C$, where Σ is a finite set of ground clauses and C is a ground clause, is decidable. This can be proved as follows: Let $C = L_1 \vee \dots \vee L_n$ and \mathcal{A} be the finite set of all ground atoms occurring in Σ and C . Now:

$$\begin{aligned}
 \Sigma \models C & \text{ iff } \Sigma \cup \{\neg L_1, \dots, \neg L_n\} \text{ is unsatisfiable} \\
 & \text{ iff } \Sigma \cup \{\neg L_1, \dots, L_n\} \text{ has no Herbrand model} \\
 & \text{ iff no subset of } \mathcal{A} \text{ is a Herbrand model of } \Sigma \cup \{\neg L_1, \dots, L_n\}
 \end{aligned}$$

Since \mathcal{A} is finite, the last statement is decidable. What follows from this is that the problem whether $\Sigma \models C$, where Σ is a finite set of function-free clauses and C is a clause, is decidable. The algorithm for computing LGI is outlined in Figure 2.7.

INPUT: A finite set \mathcal{S} of clauses, containing at least one non-tautologous function-free clause;

OUTPUT: An LGI of \mathcal{S} in \mathcal{C} ;

Remove all tautologies from \mathcal{S} , call the remaining set \mathcal{S}' ;

Let m be the number of distinct terms (including subterms) in \mathcal{S}' , let $V = \{x_1, \dots, x_m\}$;

Let \mathcal{G} be the (finite) set of all clauses which can be constructed from predicate symbols and constants in \mathcal{S}' and variables in V ;

Let $\{U_1, \dots, U_n\}$ be the set of all subsets of \mathcal{G} ;

Let H_i be an LGS (computed using algorithm in Figure 2.4) of U_i , for every $1 \leq i \leq n$;

Remove from $\{H_1, \dots, H_n\}$ all clauses which do not imply \mathcal{S}' (since each H_i is function-free, this implication is decidable), and standardize the remaining clauses $\{H_1, \dots, H_q\}$ apart. ;

return $H = H_1 \cup \dots \cup H_q$;

Figure 2.7: The LGI Algorithm.

3. Every finite set of clauses has a greatest specialization (GSI) under implication in \mathcal{C} . A GSI of a finite set \mathcal{S} is the same as the GSS of \mathcal{S} (refer Section 2.1.1), namely the union of the clauses in \mathcal{S} after these are standardized apart.

Finding least generalizations (under implication) of sets of clauses is common practice in ILP. On the other hand, the greatest specialization, which is the dual of the least generalization, is used hardly ever. Nevertheless, the GSI of two clauses D_1 and D_2 might be useful. For example, suppose we have one positive example e^+ , and two negative examples e_1^- and e_2^- , and suppose that D_1 implies e^+ and e_1^- , while D_2 implies e^+ and e_2^- . Then it might very well be that the GSI of D_1 and D_2 still implies e^+ , but is consistent with respect to $\{e_1^-, e_2^-\}$. Thus we could obtain a correct specialization by taking the GSI of D_1 and D_2 .

4. As a corollary, if \mathcal{C} is a function-free clausal language, then $\langle \mathcal{C}, \models \rangle$ is a lattice.
5. There exist pairs of Horn clauses which do not have an LGI in \mathcal{H} .
6. Similarly, there exist pairs of Horn clauses which do not have a GSI in \mathcal{H} .
7. For general clauses which all contain function symbols, the LGI-question is still open.
8. For covers, the negative results from the subsumption order carryover to the implication order.

- Some clauses, such as $\{P(x_1, x_2)\}$, have no upward covers.

- Some clauses, such as $\{P(x_1, x_2), P(x_2, x_1)\}$, have no finite complete set of downward covers.

2.3 Inverse Reduction

Plotkin's reduction algorithm (Figure 2.1) finds a reduction D of C . In this section we present an algorithm which does the inverse: given a reduced clause D , the algorithm constructs (possibly non-reduced) members C of the subsume-equivalence class of D . This will be useful in our treatment of refinement operators. Since the subsume-equivalence class of D is infinite, we have to limit the scope of the algorithm. This is done by restricting the number of literals in C .

We will next state a lemma, which will find use subsequently.

Theorem 30 *Let C and D be clauses. If D is a reduction of C , then there is a substitution θ such that $C\theta = D$ and $L\theta = L$ for every $L \in D$.*

Proof: Suppose D is a reduction of C , then there is a σ such that $C\sigma \subseteq D$. Then also $D\sigma \subseteq C\sigma \subseteq D$, since $D \subseteq C$. If $D\sigma \neq D$, then D would not be reduced, hence $D\sigma = C\sigma = D$.

Since σ is injective, for all $L_1, L_2 \in D$, if $L_1\sigma \neq L_2\sigma$ then $L_1 \neq L_2$, for otherwise $|D\sigma| < |D|$. Hence if $L_1\sigma = L_2\sigma$, then $L_1 = L_2$. For each $L \in D$, consider the following infinite sequence:

$$L, L\sigma, L\sigma^2, L\sigma^3, \dots$$

Since $D\sigma = D$, each literal in this sequence is a member of D . D contains only a finite number of literals, so for some $i < j$ we must have $L\sigma^i = L\sigma^j$. Then from the injectivity, also $L = L\sigma^{j-i}$. For this L , define $n(L) = j - i$. Notice that $L = L\sigma^m$ if m is a multiple of $n(L)$. Let k be the least common multiple of all $n(L)$. Then $L\sigma^k = L$ for every $L \in D$.

Define $\theta = \sigma^k$. Then since $C\sigma = D$ and $D\sigma = D$, we have $C\theta = D$. \square

From theorem 30, we know that for every non-reduced C such that $D \subset C$ and $D \equiv C$, we can find a θ such that $C\theta = D$ and θ only acts on variables not appearing in D . Thus C can be reduced by mapping $E = C \setminus D$ to literals in D . **In the inverse direction, we can find C by adding a set E to D , such that $E\theta \subseteq D$, where θ does not act on variables in D .** This is the idea used in the algorithm in Figure 2.8. If D is a reduced clause and m is some positive integer, then the algorithm finds a variant of every non-reduced C with m or less literals in the subsume-equivalence class of D .

As an illustration of the algorithm in Figure 2.8, if $D = \{P(x, x)\}$ and $m = 3$, then upto variants, $M_1 \in \{P(x, y), P(y, y)\}$ and $M_2 = \{P(y, z), P(x, z), P(y, x), P(z, z)\}$.

2.4 Generality order and ILP

The subsumption and implication orders discussed thus far are important for learning for the following reasons:

```

INPUT: A reduced clause  $D$  and an integer  $m$ ;
OUTPUT: Variants of every  $C$  such that  $D \equiv C$  and  $|C| \leq m$ ;
Set  $k = 0$ ;
If  $|D| \leq m$ , then output  $D$ ;
while  $l < (m - |D|)$  do
  Set  $l$  to  $l + 1$ ;
  for Every sequence  $L_1, \dots, L_l$  such that each  $L_i \in D$ , but the  $L_i$ 's are
  not necessarily distinct do
    Find every (up to variants) set  $E = \{M_1, \dots, M_t\}$  such that
    1. Every  $M_i$  contains at least one new variable not in  $D$ , and
    2. If  $x_1, \dots, x_n$  are all those new variables, then there is a  $\theta =$ 
        $x_1/t_1, \dots, x_n/t_n$ , such that  $M_i\theta = L_i$ , for  $i = 1, \dots, t$ .
    ;
  return  $D \cup E$ ;
end for
end while

```

Figure 2.8: The Inverse Reduction Algorithm: It finds a *variant* of every non-reduced equivalent clause (equivalent to D) with m or less literals.

- They provide a generality ordering for hypotheses, thus structuring the hypothesis space.
- They can be used to prune large parts of the search space.
 1. When generalizing C to C' , $C' \succ C$, all the examples covered by C will also be covered by C' (since if $\mathcal{B} \cup \{C\} \models e$ (e being an example) holds then also $\mathcal{B} \cup \{C'\} \models e$ holds). This property is used to prune the search of more general clauses when e **is a negative example: if e is inconsistent (covers a negative example) then all its generalizations will also be inconsistent**. Hence, the generalizations of C do not need to be considered.
 2. When specializing C to C' , $C \succ C'$, an example not covered by C will not be covered by any of its specializations either (since if $\mathcal{B} \cup \{C\} \not\models e$ holds then also $\mathcal{B} \cup \{C'\} \not\models e$ holds). This property is used to prune the search of more specific clauses when e **is an uncovered positive example: if C does not cover a positive example none of its specializations will**. Hence, the specializations of C do not need to be considered.
- The generality orderings provide the basis for two important ILP techniques:
 1. bottom-up building of least general generalizations from training examples, relative to background knowledge, and

2. top-down searching of refinement graphs.

The principal generality orderings of interest are subsumption (\succeq_θ) and implication (\succeq_{\models}). For clauses C, D , subsumption is *not* equivalent to implication

$$\text{if } C \succeq_\theta D \text{ then } C \succeq_{\models} D$$

but

not vice – versa

For example, the above holds for:

$$C : \text{natural}(s(X)) \leftarrow \text{natural}(X)$$

$$D : \text{natural}(s(s(X))) \leftarrow \text{natural}(X)$$

Recall the subsumption theorem; it is a key theorem linking subsumption and implication:

If Σ is a set of clauses and D is a clause, then $\Sigma \models D$ iff D is a tautology, or there exists a clause $D' \succeq_\theta D$ which can be derived from Σ using some form of resolution.

When Σ contains a single clause C then the only clauses that can be derived are the result of *self-resolutions* of C . Thus, the difference between $C \succeq_{\models} D$ and $C \succeq_\theta D$ arises when C is self-recursive or D is tautological

Is there a principled approach for comparing generality orderings? Fortunately, the answer is yes.

Given a set of clauses S , clauses $C, D \in S$ and quasi-orders \succeq_1 and \succeq_2 on S , then \succeq_1 is *stronger* than \succeq_2 if $C \succeq_2 D$ implies $C \succeq_1 D$. If also for some $C, D \in S$ $C \not\succeq_2 D$ and $C \succeq_1 D$ then \succeq_1 is *strictly stronger* than \succeq_2

It can be seen that as per above definition, the implication ordering is strictly stronger than the subsumption ordering. Quasi-orders that are increasingly weaker can be devised from stronger ones. Here are some other generality orderings, listed in decreasing order of strength.

- $C \succeq_{\models} D$ iff $C \models D$
- $C \succeq_\theta D$ iff there is a substitution θ s.t. $C \subseteq D$
- $C \succeq_{\theta'} D$ iff every literal in D is *compatible* to a literal in C and $C \succeq_\theta D$.
- $C \succeq_{\theta''} D$ iff $|C| \geq |D|$ and $C \succeq_{\theta'} D$

What generality ordering should we choose for the ILP search procedure? We would like the strongest ordering that is practical in terms of tractability, logical implication between clauses is undecidable (even for Horn clauses). Subsumption is decidable but NP-complete (even for Horn clauses). Restrictions to the form of clauses can make subsumption efficient. Here are two example restrictions that make subsumption efficient.

- Determinate Horn clauses. There exists an ordering of literals in C and exactly one substitution θ s.t. $C\theta \subseteq D$.
- k -local Horn clauses. Partition a Horn clause into k “disjoint” sub-parts and perform k independent subsumption tests.

The strongest quasi-order that is practical appears to be subsumption. Even that often requires restrictions (such as the two listed above) on the clauses being compared.

In summary, the subsumption order on clausal languages is used most often in ILP as the generality order (in contrast to the implication order), owing to its following properties:

1. Further, subsumption is more tractable and efficiently implementable. Subsumption between clauses is a decidable relation, whereas implication is not. The flip side is that subsumption is a weaker relation.
2. Equivalence classes under subsumption can be represented by a single reduced clause. Reduction can be undone by inverse reduction (*c.f.* Section 2.3).
3. Every finite set of clauses (function free or not) has a least generalization (LGS) and greatest specialization (GSS) under subsumption in \mathcal{C} . Hence $\langle \mathcal{C}, \succeq \rangle$ is a lattice. The same is not true for the implication quasi-ordering \succeq_{\models} (for restricted languages *lubs* for \succeq_{\models} may well exist).

Order	<i>lub</i>	<i>glb</i>
\succeq_{θ}	✓	✓
\succeq_{\models}	×	✓

4. Every finite set of Horn clauses has a least generalization (LGS) and greatest specialization (GSS) under subsumption in \mathcal{H} . Hence $\langle \mathcal{H}, \succeq \rangle$ is a lattice. The same does not hold for the implication order.
5. The negative results for covers hold for subsumption as well as implication.

2.5 Incorporating Background Knowledge

Why does background knowledge matter? The answer is that combining the examples With what we already know often allows for the construction of a more satisfactory theory than can be glanced from the examples by themselves. To illustrate this, we consider [Bun88] the following two clauses as positive examples (not just ground atoms as examples):

$$\begin{aligned}
 D_1 &= \text{CuddlyPet}(x) \leftarrow \text{Small}(x), \text{Fluffy}(x), \text{Dog}(x) \\
 D_2 &= \text{CuddlyPet}(x) \leftarrow \text{Fluffy}(x), \text{Cat}(x)
 \end{aligned}$$

Given only these clauses, the most obvious way to generalize them is to take their LGS or LGI, which is the rather general clause $C = \text{CuddlyPet}(x) \leftarrow \text{Fluffy}(x)$.

However, suppose we have the following definite program \mathcal{B} which expresses our background knowledge.

$$\begin{aligned} B_1 \quad & \text{Pet}(x) \leftarrow \text{Cat}(x) \\ B_2 \quad & \text{Pet}(x) \leftarrow \text{Dog}(x) \\ B_3 \quad & \text{Small}(x) \leftarrow \text{Cat}(x) \end{aligned}$$

Given \mathcal{B} , we may also use the following clause as generalization:

$$D = \text{CuddlyPet}(x) \leftarrow \text{Small}(x), \text{Fluffy}(x), \text{Pet}(x).$$

since D together with \mathcal{B} implies both examples. Note that without the background knowledge \mathcal{B} , our clause D neither subsumes nor implies the examples. If we interpret this example in human terms, the generalization D is much more satisfactory than C . After all, not every fluffy object (such as a teddy bear) is a cuddly pet. Thus the use of background knowledge allows us to find a better theory. Given the usefulness of background knowledge, it is important to outline a formalized way to reckon with it in our generality order.

There are three generality orders which are able to take background knowledge into account: (i) Plotkin's relative subsumption ($\succeq_{\mathcal{B}}$), (ii) Relative implication ($\models_{\mathcal{B}}$) and Buntine's (iii) Generalized subsumption ($\geq_{\mathcal{B}}$). Relative subsumption and relative implication apply to arbitrary clauses and the background knowledge \mathcal{B} may be an arbitrary finite set of clauses. Generalized subsumption only applies to definite program clauses and the background knowledge should be a definite program. We will state the existence of least generalizations in each of these orders, both in case we are dealing with a Horn language \mathcal{H} , and for a general clausal language \mathcal{C} . Further, each of the three orders can be related to some kind of deduction. Since empty background knowledge reduces relative and generalized subsumption to ordinary subsumption, and relative implication to ordinary implication, the negative results on covers for subsumption and implication carry over to the three orders: some clauses do not have finite complete sets of upward or downward covers in these orders. As to the existence and non-existence of least generalizations or greatest specializations in the three orders, we will only pay attention to least generalizations, since these are used much more often than their dual. In general, for all three, least generalizations do not always exist in the presence of background knowledge. However, certain restrictions on the background knowledge guarantee the existence of a least generalization for each of the three.

2.5.1 Plotkin's relative subsumption ($\succeq_{\mathcal{B}}$)

Definition 1 Let C and D be clauses, and \mathcal{B} be a set of clauses. We say $C \succeq_{\mathcal{B}} D$, if there is a substitution θ such that $\mathcal{B} \models \forall(C\theta \rightarrow D)$ (note that $C \rightarrow D$ need not be a clause).

With respect to the last example, we can verify that $B_3 \succeq_{\{B_1, E\}} F$, where $E = \text{CuddlyPet}(x) \leftarrow \text{Small}(x), \text{Fluffy}(x), \text{Pet}(x)$ and $F = \text{CuddlyPet}(x) \leftarrow \text{Fluffy}(x), \text{Cat}(x)$. Informally, this can be seen as follows: suppose for every x it holds that if x is a cat, then x is small (*i.e.*, C is true). Using B_1 , we also have that if x is a cat, then x is a pet. Thus if x is a fluffy cat, then x is small, fluffy, and a pet, and by E , x is a cuddly pet (*i.e.*, D is true). Following are some properties of relative subsumption:

1. Reflexivity and transitivity are easily proved, so relative subsumption is a quasi-order on clauses. Note that each set of clauses \mathcal{B} induces its own quasi-order: the quasi-orders induced by $\mathcal{B} = \{P(a)\}$ and $\mathcal{B} = \{P(a), P(b)\}$ are different. Note also that if D is a tautology, then $C \succeq_{\mathcal{B}} D$ for any C and \mathcal{B} . Furthermore, it is also easy to see that if $C \succeq_{\mathcal{B}} D$ and $\mathcal{B} \subseteq \mathcal{B}'$, then $C \succeq_{\mathcal{B}'} D$.
2. Relative subsumption is strictly stronger than subsumption.
 - (a) Firstly, it is easy to see that subsumption implies relative subsumption. If $C \succeq D$, then $C \succeq_{\mathcal{B}} D$, for some \mathcal{B} . But then $\forall(C\theta \succeq D)$ is a tautology, and $\mathcal{B} \models \forall(C\theta \rightarrow D)$ for any \mathcal{B} . Hence if $C \succeq D$, then $C \succeq_{\mathcal{B}} D$.
 - (b) Now consider propositional atoms P , Q , and R . Let $C = P$, $D = Q$, and $\mathcal{B} = \{Q \leftarrow P\}$. Then $\mathcal{B} \models (C \rightarrow D)$, so $C \succeq_{\mathcal{B}} D$. Since $C \not\succeq D$, we see that relative subsumption does not imply subsumption. This even holds for the case where \mathcal{B} is empty and D is a tautology: if $\mathcal{B} = \emptyset$, $C = Q$, and $D = P \leftarrow P$, then $C \succeq_{\mathcal{B}} D$, but $C \not\succeq D$.
3. If C and D are non-tautologous clauses, and \mathcal{B} a finite set of ground literals such that $\mathcal{B} \cap D = \emptyset$, then⁶ $C \succeq_{\mathcal{B}} D$ iff $C \succeq (D \vee \overline{\mathcal{B}})$.
4. Relative subsumption coincides with ordinary subsumption for non-tautologous clauses and empty background knowledge⁷. That is, if C and D are non-tautologous clauses, then $C \succeq_{\mathcal{B}} D$ iff $C \succ D$.
5. $C \succeq_{\mathcal{B}} D$ iff there exists a deduction of D from $\{C\} \cup \mathcal{B}$ in which C occurs at most once as a leaf. The proof of this long and windy and will not be dealt with here.
6. Least generalizations under relative subsumption (abbreviated to LGRSs) need not exist in the general case. The following counter example, adapted from [Nlb88] shows the non-existence of LGRSs both for the case of a clausal language \mathcal{C} , and for a Horn language \mathcal{H} .

Let

$$\begin{aligned} D_1 &= Q(a) & D_2 &= Q(b) \\ \mathcal{B} &= \{P(a, y), P(b, y)\} \end{aligned}$$

⁶Prove: EXERCISE

⁷Prove: EXERCISE

It can be shown there is no LGRS of $\{D_1, D_2\}$ relative to \mathcal{B} . Consider the following infinite sequence of clauses:

$$\begin{aligned} C_1 &= Q(x) \leftarrow P(x, f(x)) \\ C_2 &= Q(x) \leftarrow P(x, f(x)), P(x, f^2(x)) \\ C_3 &= Q(x) \leftarrow P(x, f(x)), P(x, f^2(x)), P(x, f^3(x)) \end{aligned}$$

It is easy to see that $C_i \succ_{\mathcal{B}} C_{i+1}$ for every $i \geq 1$. We also have $C_i \succ_{\mathcal{B}} D_1$ and $C_i \succ_{\mathcal{B}} D_2$, for every $i \geq 1$. Suppose some clause D is an LGRS of $\{D_1, D_2\}$ relative to \mathcal{B} , then we should have $C \succeq_{\mathcal{B}} D$, for every $i \geq 1$. Then by point 5 above, for every $i \geq 1$, there exists a deduction of D from $\{C_i\} \cup \mathcal{B}$ in which C_i occurs at most once as a leaf. C_i cannot occur zero times as a leaf, because then a clause from \mathcal{B} would simply subsume D , which is impossible. Thus for every $i \geq 1$, there exists a deduction of D from $\{C_i\} \cup \mathcal{B}$, in which C_i occurs once as a leaf. It cannot be that every C_i subsumes D , for then D would contain an instance of the term $f^i(x)$, for every $i \geq 1$. Thus for some j the deduction of D from $\{C_j\} \cup \mathcal{B}$ involves at least one resolution step. Since the members of \mathcal{B} are atoms and C_j only occurs once as a leaf, the parent clauses in the first resolution step must be C_j and a member of \mathcal{B} . Suppose this member of \mathcal{B} is $P(a, y)$. Then $P(a, y)$ must be unified with an atom of the form $P(x, f''(x))$ in the body of C_j . Then the head of the clause C_j is instantiated to $Q(a)$. This head will not be changed anymore in later resolution steps, so D would have $Q(a)$ as head - but then there is no deduction of $D_2 = Q(b)$ from $\{D\} \cup \mathcal{B}$.

Similarly, if the member of \mathcal{B} in the resolution step had been $P(b, y)$ instead of $P(a, y)$, D would have had $Q(b)$ as head, and there would be no deduction of $D_1 = Q(a)$ from $\{D\} \cup \mathcal{B}$. Either way, the assumption that D is an LGRS of $\{D_1, D_2\}$ relative to \mathcal{B} leads to a contradiction.

7. We however can identify a restriction on the background knowledge which guarantees existence (and computability) of an LGRS of any finite set of (horn) clauses. If $\mathcal{C}(\mathcal{H})$ is a (horn) clausal language and $\mathcal{B} \subseteq \mathcal{C}$ ($\mathcal{B} \subseteq \mathcal{H}$) is a finite set of ground literals, then every finite non-empty set $\mathcal{S} \subseteq \mathcal{C}$ ($\mathcal{S} \subseteq \mathcal{H}$) of clauses has an LGRS in $\mathcal{C}(\mathcal{H})$. This can be proved as follows: If a clause D is a tautology or $\mathcal{B} \cap D \neq \emptyset$, then $\mathcal{B} \models D$, hence for any clause C we have $C \succeq_{\mathcal{B}} D$. Remove from \mathcal{S} all tautologies and all D for which $\mathcal{B} \cap D \neq \emptyset$, call the remaining set \mathcal{S}' . If \mathcal{S}' is empty, any tautology is an LGRS of \mathcal{S} . If $\mathcal{S}' = \{D_1, \dots, D_n\}$ is non-empty, then it follows easily from point 3) above that an LGS of $\{(D_1 \vee \overline{\mathcal{B}}), \dots, (D_n \vee \overline{\mathcal{B}})\}$ (or equivalently of $\{(D_1 \leftarrow \mathcal{B}), \dots, (D_n \leftarrow \mathcal{B})\}$) in \mathcal{C} is an LGRS of \mathcal{S}' in \mathcal{C} , and hence also of \mathcal{S} . The existence of such an LGS follows from Theorem 28. Thus if the background knowledge \mathcal{B} is a finite set of ground literals, then we can construct an LGRS of a set $\mathcal{S} = \{D_1, \dots, D_n\}$ simply by constructing an LGS of $T = \{(D_1 \vee \overline{\mathcal{B}}), \dots, (D_n \vee \overline{\mathcal{B}})\}$. Additionally, if all clauses in \mathcal{S} are

horn clauses and each literal in \mathcal{B} is a ground atom, then each $D_i \vee \overline{\mathcal{B}}$ is also a horn clause and so will the LGS of T . So the result follows for \mathcal{H} . This result has been exploited in the GOLEM system.

8. The non-existence of finite chains of covers, in lattices of (Horn) clauses ordered by subsumption carries over to the lattice of clauses ordered by relative subsumption.

Relative Subsumption and ILP

We will develop further on point (7) discussed above. When presented with an example e , candidate hypothesis in the form of a clause C and background knowledge \mathcal{B} , what does it mean for clause C to “relatively subsume” example e . Recall normal subsumption: $C \succeq e$ means $\exists \theta$ s.t. $C\theta \subseteq e$. This also means $C\theta \models e$ or $\models (e \leftarrow C\theta)$.

$$\begin{array}{ll}
 e : & gfather(henry, john) \leftarrow \\
 B : & father(henry, jane) \leftarrow \\
 & father(henry, joe) \leftarrow \\
 & parent(jane, john) \leftarrow \\
 & parent(joe, robert) \leftarrow \\
 C : & gfather(X, Y) \leftarrow father(X, Z), parent(Z, Y)
 \end{array}$$

For this B, C, e with $\theta = \{X/henry, Y/john, Z/jane\}$, $B \cup \{C\theta\} \models e$ That is: $C \succeq_B e$ or equivalently, $B \models (e \leftarrow C\theta)$ (by subsumption theorem). However, note that $C \not\succeq e$. Clearly if $B = \emptyset$ normal subsumption between clauses results. Using the Deduction Theorem

$$\begin{aligned}
 \mathcal{B} \models (e \leftarrow C\theta) &\equiv \mathcal{B} \cup \{C\theta\} \models e \\
 &\equiv \mathcal{B} \cup \overline{e} \models \overline{C\theta} \\
 &\equiv \{C\theta\} \models \overline{\mathcal{B} \cup \overline{e}} \\
 &\equiv \models (\overline{\mathcal{B} \cup \overline{e}} \leftarrow C\theta)
 \end{aligned}$$

That is, $C \succeq_B e$ means $C \succeq \overline{\mathcal{B} \cup \overline{e}}$ or $C \models \overline{\mathcal{B} \cup \overline{e}}$. Recall that if $C_1 \succeq C_2$ then $C_1 \models C_2$. In fact, if $C_{1,2}$ are not self-recursive, then $C_1 \succeq C_2 \equiv C_1 \models C_2$

Let $a_1 \wedge a_2 \dots$ be the ground literals true in all models of $\mathcal{B} \cup \overline{e}$, that is these are all the members of the minimal model $MM(\mathcal{G}(\mathcal{B} \cup \overline{e}))$ (c.f. theorem 12 as well as theorem 17). Then

$$\begin{array}{l}
 \mathcal{B} \cup \overline{e} \models a_1 \wedge a_2 \dots \\
 \overline{a_1 \wedge a_2 \wedge \dots} \models \overline{\mathcal{B} \cup \overline{e}} \equiv e \leftarrow \mathcal{B}
 \end{array}$$

Let $\perp(\mathcal{B}, e) = \overline{a_1 \wedge a_2 \wedge \dots}$. Note that $\perp(\mathcal{B}, e)$ (the ‘most specific clause’ given \mathcal{B}, e) may not always be finite. If \mathcal{B} and e are both ground, it follows that $\perp(\mathcal{B}, e) \equiv (e \leftarrow \mathcal{B})$.

Now, if $D \succeq \perp(\mathcal{B}, e)$ then $D \models \perp(\mathcal{B}, e)$ and therefore $D \models \overline{\mathcal{B} \cup \bar{e}}$. Now, by transitivity of \models , if $D \models \perp(\mathcal{B}, e)$ then because $\perp(\mathcal{B}, e) \models e \leftarrow \mathcal{B}$, it should follow that $D \models e \leftarrow \mathcal{B}$. It can be further shown that if D, e are not self-recursive and $D \succeq \perp(\mathcal{B}, e)$ then $D \succeq \overline{\mathcal{B} \cup \bar{e}}$ (that is, $D \succeq (e \leftarrow \mathcal{B})$ or $D \succeq_B e$). In fact, if C, e are non self-recursive then

$$\begin{aligned} C \succeq \perp(\mathcal{B}, e) &\equiv \\ C \models \overline{\mathcal{B} \cup \bar{e}} &\equiv \\ C \succeq \overline{\mathcal{B} \cup \bar{e}} &\equiv \\ C \succeq e \leftarrow \mathcal{B} & \end{aligned}$$

An example of finding \perp is:

\mathcal{B} :

```
gfather(X,Y) ← father(X,Z), parent(Z,Y)
father(henry,jane) ←
mother(jane,john) ←
mother(jane,alice) ←
```

e_i :

```
gfather(henry,john) ←
```

Conjunction of ground atoms provable from $\mathcal{B} \cup \bar{e}_i$:

```
¬parent(jane,john) ∧
father(henry,jane) ∧
mother(jane,john) ∧
mother(jane,alice) ∧
¬gfather(henry,john)
```

$\perp(\mathcal{B}, e_i)$:

```
gfather(henry,john) ∨ parent(jane,john) ←
father(henry,jane),
mother(jane,john),
mother(jane,alice)
```

D_i :

```
parent(X,Y) ← mother(X,Y)
```

Mode declarations

Finding a clause D_i that subsumes $\perp(\mathcal{B}, e_i)$ is hampered by the fact that $\perp(\mathcal{B}, e_i)$ may be infinite! One workaround is to use a constrained subset of definite clauses to construct finite most-specific clauses. This can be enabled using **mode declarations**. An example set of mode declarations for the problem just considered is:

```
modeh(*,gfather(+person,-person))
modeh(*,parent(+person,-person))
modeb(*,father(+person,-person))
```

$modeb(*,parent(+person,-person))$
 $modeb(*,mother(+person,-person))$

A *definite mode language* is defined as follows.

Definition 2 Let $C : h \leftarrow b_1, \dots, b_n$ be a definite clause with an ordering over literals. Let M be a set of mode declarations. C is in the definite mode language $\mathcal{L}(M)$ iff

1. h is the atom of a modeh declaration in M with every place-marker of $+type$ and $-type$ replaced with variables, and every place marker of $\#type$ replaced by a ground term.
2. Every atom b_i in body of C is an atom in a modeb declaration in M with $+, -, \#$ places being replaced as above.
3. Every variable of $+type$ in b_i is either of $+type$ in h or of $-type$ in a b_j ($1 \leq j < i$)

Given a set of mode declarations M it is always possible to decide if a clause C is in $\mathcal{L}(M)$. Next, we define depth of variables.

Definition 3 Let C be a definite clause, v be a variable in an atom in C , and U_v all other variables in body atoms of C that contain v

$$d(v) = \begin{cases} 0 & \text{if } v \text{ in head of } C \\ (\max_{u \in U_v} d(u)) + 1 & \text{otherwise} \end{cases}$$

For example, if $C : gfather(X, Y) \leftarrow father(X, Z), parent(Z, Y)$

Then $d(X) = d(Y) = 0$, $d(Z) = 1$. Putting together the definitions of mode language and depth, we next define depth bounded definite mode language.

Definition 4 Let C be a definite clause with an ordering over literals. Let M be a set of mode declarations. C is in the depth-bounded definite mode language $\mathcal{L}_d(M)$ iff all variables in C have depth at most d

As an example, the clause for *gfather/2* earlier is in $\mathcal{L}_2(M)$.

We will state some properties of $\perp(\mathcal{B}, e_i)$ without their proofs. For every $\perp(\mathcal{B}, e_i)$ it is the case that:

There is a $\perp_d(\mathcal{B}, e_i)$ in $\mathcal{L}_d(M)$ s.t. $\perp_d(\mathcal{B}, e_i) \succeq \perp(\mathcal{B}, e_i)$

$\perp_d(\mathcal{B}, e_i)$ is finite

If $C \succeq \perp_d(\mathcal{B}, e_i)$ then $C \succeq \perp(\mathcal{B}, e_i)$

$\perp(\mathcal{B}, e_i)$ may not be Horn

$\perp(\mathcal{B}, e_i)$ may not be finite

Least upper bound of Horn clauses e_1, e_2 is:

$$lgg_B(e_1, e_2) = lgg(\perp(B, e_1), \perp(B, e_2))$$

Greatest lower bound of Horn clauses e_1, e_2 is:

$$glb_b(e_1, e_2) = glb(\perp(B, e_1), \perp(B, e_2))$$

Let us take an example of finding \perp_i :

$\perp(B, e_i)$:

```

gfather(henry, john)  $\vee$  parent(jane, john)  $\leftarrow$ 
    father(henry, jane),
    mother(jane, john),
    mother(jane, alice)

```

modes:

```

modeh(*, parent(+person, -person))
modeb(*, mother(+person, -person))
modeb(*, father(+person, -person))

```

$\perp_0(B, e_i)$:

```

parent(X, Y)  $\leftarrow$ 

```

$\perp_1(B, e_i)$:

```

parent(X, Y)  $\leftarrow$ 
    mother(X, Y),
    mother(X, Z)

```

2.5.2 Relative implication ($\models_{\mathcal{B}}$)

Definition 5 Let C and D be clauses, and \mathcal{B} be a set of clauses. C (logically) implies D relative to \mathcal{B} , denoted $C \models_{\mathcal{B}} D$, if $\{C\} \cup \mathcal{B} \models_{\mathcal{B}} D$.

For example, if $\mathcal{B} = \{P(a)\}$, $C = P(f(x)) \leftarrow P(x)$ and $D = P(f^2(a))$, then $C \models_{\mathcal{B}} D$, because there is a deduction of D from $\{C\} \cup \mathcal{B}$. However, we have $C \not\models_{\mathcal{B}} D$ because C has to be used more than once in the deduction of D . Following are some properties of relative subsumption:

1. Relative implication is perhaps the most obvious way to take background knowledge into account.
2. Obviously relative implication is reflexive and transitive so it can serve as a quasi-order on a set of clauses.

3. It is also obvious that if $C \models D$ then $C \models_{\mathcal{B}} D$. The converse need not hold though. Consider $C = P(a) \leftarrow P(b)$, $D = P(a)$, and $\mathcal{B} = \{P(b)\}$: then $C \models_{\mathcal{B}} D$, but $C \not\models D$.
4. Relative subsumption implies relative implication but not conversely. If $C \succeq_{\mathcal{B}} D$, then $\mathcal{B} \models \forall(C\theta \rightarrow D)$ (for some θ). If $C \models_{\mathcal{B}} D$ then $\{C\} \cup \mathcal{B} \models D$, which is equivalent to $\mathcal{B} \models \forall(C) \rightarrow \forall(D)$ by the deduction theorem.
5. $C \models_{\mathcal{B}} D$ iff there exists a deduction of D from $\{C\} \cup \mathcal{B}$ (contrast this with that for relative subsumption). This is another testimony to the fact that relative implication is a strictly stronger quasi-order than relative subsumption.
6. The negative results for existence of least generalizations under implication carry over to relative implication (LGRI), since ordinary logical implication is just a special case of relative implication. Further, the only positive result for LGI that was stated on page 110 in point (2) is also negative for LGRI. As an example, consider

$$\begin{aligned} D_1 &= P(a) & D_2 &= P(b) \\ \mathcal{B} &= \{(P(a) \vee \neg Q(x)), (P(b) \vee \neg Q(x))\} \end{aligned}$$

It can be shown that the set $\mathcal{S} = \{D_1, D_2\}$ has no LGRI relative to \mathcal{B} in \mathcal{C} . Suppose indeed D is an LGRI of \mathcal{S} relative to \mathcal{B} . Note that if D contains the literal $P(a)$, then the Herbrand interpretation which makes $P(a)$ true, and which makes all other ground atoms false, would be a model of $\mathcal{B} \cup \{D\}$ but not of D_2 , so then we would have $D \not\models_{\mathcal{B}} D_2$. Similarly, if D contains $P(b)$ then $D \not\models_{\mathcal{B}} D_1$. Hence D cannot contain $P(a)$ or $P(b)$. Now let d be a constant not appearing in D . Let $C = P(x) \vee Q(d)$, then $C \models_{\mathcal{B}} \mathcal{S}$. By the definition of an LGRI, we should have $C \models_{\mathcal{B}} D$. Then by the Subsumption Theorem, there must be a derivation from $\mathcal{B} \cup \{C\}$ of a clause E , which subsumes D . The set of all clauses which can be derived (in 0 or more resolution steps) from $\mathcal{B} \cup \{C\}$ is $\mathcal{B} \cup \{C\} \cup \{(P(a) \vee P(x)), (P(b) \vee P(x))\}$. But none of these clauses subsumes D , because D does not contain the constant d , nor the literals $P(a)$ or $P(b)$. Hence $C \not\models_{\mathcal{B}} D$, contradicting the assumption that D is an LGRI of \mathcal{S} relative to \mathcal{B} in \mathcal{C} . Thus, in general an LGRI of \mathcal{S} relative to \mathcal{B} need not exist.

7. We can identify a special case in which the existence of an LGRI is guaranteed. Let C and D be clauses, and \mathcal{B} be a finite set of function-free ground literals. Then $C \models_{\mathcal{B}} D$ iff $C \models (D \cup \bar{\mathcal{B}})$. Further, if $\mathcal{S} \subseteq \mathcal{C}$ is a finite set of clauses, containing at least one D for which $D \cup \bar{\mathcal{B}}$ is non-tautologous and function-free, then \mathcal{S} has an LGRI in \mathcal{C} . As a special case of this, if \mathcal{C} is itself function free clausal language, then for any finite $\mathcal{S} \subseteq \mathcal{C}$, \mathcal{S} has an LGRI in \mathcal{C} .

This can be proved as follows: Suppose $C \models_{\mathcal{B}} D$ i.e., $\{C\} \cup \mathcal{B} \models D$. Let M be a model of C . Then we need to show that M is also a model of

$D \cup \bar{\mathcal{B}}$. If M is not a model of \mathcal{B} , then it is a model of at least one literal in $\bar{\mathcal{B}}$, and hence of the clause $D \cup \bar{\mathcal{B}}$. If on the other hand, M is a model of \mathcal{B} then M is also a model of D , because $\{C\} \cup \mathcal{B} \models D$. Then M is also a model of $D \cup \bar{\mathcal{B}}$. This prove the ‘only if’ part.

Now suppose $C \models (D \cup \bar{\mathcal{B}})$. Let M be a model of $\{C\} \cup \mathcal{B}$. Then we need to show that M is also a model of D . Now, M is a model of C and hence of the clause $D \cup \bar{\mathcal{B}}$. But M is also a model of \mathcal{B} , and hence not a model of $\bar{\mathcal{B}}$. Therefore, M must be a model of D . This proves the ‘if’ part.

Now if $\mathcal{S} = \{D_1, D_2, \dots, D_n\}$, it follows that an LGI in \mathcal{C} of $T = \{(D_1 \cup \bar{\mathcal{B}}), \dots, (D_n \cup \bar{\mathcal{B}})\}$ is also an LGRI of \mathcal{S} in \mathcal{C} (which exists according to point 2 on page 110).

2.5.3 Buntine’s generalized subsumption ($\geq_{\mathcal{B}}$)

Definition 6 Let C and D be definite program clauses and \mathcal{B} be a definite program, comprising the background knowledge. We say that $C \geq_{\mathcal{B}} D$ (g-subsumes), if for every Herbrand model M of \mathcal{B} and every ground atom A such that D covers A under M , we have that C covers A under M . Definite clause C is said to ‘cover’ atom A under M if there exists a ground substitution θ for C (that is, $C\theta$ is ground), such that M is a model for $C^{-}\theta$ and $C^{+}\theta = A$.

Generalized subsumption applies only to definite program clauses. Loosely speaking, generalized subsumption says that for definite clauses C and D , C is more general than D , in any situation consistent with what we already know (through \mathcal{B}), C can be used to prove at least as many results as D . For example, if \mathcal{B} consists of the background knowledge on page 116 and C and D are defined somewhat similar to D_1 and D_2 as

$$\begin{aligned} C &= \text{CuddlyPet}(x) \leftarrow \text{Small}(x), \text{Pet}(x) \\ D &= \text{CuddlyPet}(x) \leftarrow \text{Cat}(x) \end{aligned}$$

then, we can show that $C \geq_{\mathcal{B}} D$. For suppose M is a Herbrand model of \mathcal{B} and D covers some ground atom $A = \text{CuddlyPet}(t)$ under M , then for $\theta = \{x/t\}$, $D^{-}\theta = \text{Cat}(t)$ is true under M . Since M is a model of \mathcal{B} , in particular of B_1 and B_3 , $\text{Pet}(t)$ and $\text{Small}(t)$ must be true under M as well. Then $C^{-}\theta = \text{Small}(t) \wedge \text{Pet}(t)$ is true under M , and $C^{+}\theta = A$, so C also covers A under M . Hence $C \geq_{\mathcal{B}} D$. In natural language, this can be rephrased as

**If small pets are cuddly pets, then cats are cuddly pets,
since we already know that cats are small pets.**

Some of the properties of $\geq_{\mathcal{B}}$ are:

1. Just as subsumption implies relative subsumption, subsumption also implies g-subsumption. Suppose $C \succeq D$. Then $C\theta \subseteq D$, so $C^{+}\theta = D^{+}$ and $C^{-}\theta \subseteq D^{-}\theta$ for some θ . If D covers some A under some I , there is a γ such that $D^{-}\gamma$ is true under I and $D^{+}\gamma = A$. But then $C^{+}\theta\gamma = D^{+}\gamma = A$,

and $C^- \gamma \subseteq D^- \gamma$, is true under I , so C covers A under I as well. Hence, if $C \succeq D$, then $C \geq_{\mathcal{B}} D$.

The converse need not hold: if $\mathcal{B} = \{P(a)\}$, $C = Q(a) \leftarrow P(a)$ and $D = Q(a)$, then $C \geq_{\mathcal{B}} D$ but $C \not\succeq D$.

2. Generalized subsumption is reflexive and transitive with respect to some definite program \mathcal{B} , so it can serve as a quasi-order on a set of definite clauses.
3. $C \geq_{\mathcal{B}} D$ iff there exists an SLD-deduction of D , with C as top clause and members of \mathcal{B} as input clauses. It follows from this that g-subsumption reduces to ordinary subsumption in the presence of empty background knowledge, as was the case for relative subsumption. That is, if C and D are definite program clauses, then $C \geq_{\emptyset} D$ iff $C \succeq D$.
4. From the previous point, it follows that if C and D are definite program clauses, and \mathcal{B} is a definite program, then if $C \geq_{\mathcal{B}} D$, it follows that $C \succeq_{\mathcal{B}} D$. That is, g-subsumption implies relative subsumption. But the converse does not hold. This is sufficient to prove that if an LGG (lub) does not exist under relative subsumption (such as for the example
5. An LGGS of a finite set \mathcal{S} (of definite program clauses which all have the same predicate symbol in their respective heads) always exists either if
 - (a) All clauses in \mathcal{S} are atoms, and the background knowledge \mathcal{B} implies only a finite number of ground atoms (i.e., $M_{\mathcal{B}}$ is finite)
 - (b) \mathcal{S} and \mathcal{B} are all function-free (see [Bun88]).
 - (c) \mathcal{B} is ground. This case differs from the first, because \mathcal{B} may imply only a finite number of ground examples, and still be non-ground itself. For example, $\mathcal{B} = \{P(a), (Q(x) \leftarrow P(x))\}$.

Actually, these three cases are special cases of the following theorem:

Theorem 31 *Let \mathcal{H} be a Horn language and \mathcal{B} be a definite program. Let $\mathcal{S} = \{D_1, \dots, D_n\} \subseteq \mathcal{H}$ be a finite non-empty set of definite program clauses, such that all D_i , have the same predicate symbol in their head. Furthermore, for every $1 \leq i \leq n$, let σ_i be a Skolem substitution for D_i , with respect to $\mathcal{B} \cup \mathcal{S}$, and M_i be the least Herbrand model of $\mathcal{B} \cup D_i^- \sigma_i$. If every M_i is finite, then there exists an LGGS of \mathcal{S} in \mathcal{H} .*

Thus if the least Herbrand models mentioned in the theorem are indeed finite, then we can find an LGGS of a set $\{D_1, \dots, D_n\}$ simply by constructing an LGS of $\{(\{D_1^+\} \cup \overline{M_1}), \dots, (\{D_n^+ \sigma_n\} \cup \overline{M_n})\}$.

2.6 Using Generalization and Specialization

The normal problem of inductive logic programming is to find a correct theory, a set of clauses which implies all given positive examples and which is consistent with respect to the given negative examples. Usually, it is not immediately obvious which set of clauses we should pick as our theory. Rather, we will have to search among the permitted clauses for a set of clauses with the right properties. If a positive example is implied by the theory, we should search for a more general theory. On the other hand, if the theory is not consistent with respect to the negative examples, we should search for a more specific theory - for instance, by replacing a clause in the theory by more specific clauses - such that the theory becomes consistent. Thus, the two most important operations in ILP are *generalization* and *specialization*. Repeated application of such generalization and specialization steps may finally yield a correct theory.

To systematically facilitate this search, it would be very handy if the set of clauses that has to be searched, is somehow structured. Fortunately, this is so as seen in chapter 1. We had seen several alternatives for what it means for some clause to be more general than another clause. Since generalization (or dually, specialization) can proceed along the lines of such a generality order, using such an order can direct the search for a correct theory. For example, least generalizations can be used to generalize given finite sets of examples.

2.7 Using the Cover Structure: Refinement Operators

One way to weaken an existing theory which is too strong is to find a false member of the theory Σ , and delete this clause from the theory. However, deleting a clause might make the theory in turn too weak. A way to strengthen the theory again, is to add weaker versions of previously deleted clauses. For instance, suppose the clause $P(x)$ is false under I , and has been deleted from Σ . It might be that $P(f(x))$, which is a “refinement” of $P(x)$, is true under I . Thus the theory might be strengthened by adding $P(f(x))$ to it.

The finite downward and upward cover chain algorithm provide the general direction of search in any search over the lattice structure over theories (atomic or otherwise). In *top-down search*, we want to find some unknown specialization \mathbf{l}_2 of \mathbf{l}_1 . Then we should use substitutions to try and find a chain of downward covers starting from \mathbf{l}_1 as in Algorithm 1.15. Since such finite chains always exist for atoms, we can restrict attention to downward covers of \mathbf{l}_1 , downward covers of downward covers of \mathbf{l}_1 , *etc.* Recall from the algorithm in Figure 1.15 that the progress from \mathbf{l}_i to \mathbf{l}_{i+1} is achieved by applying one of the following substitutions:

1. $\{X/f(X_1, \dots, X_k)\}$ where X is a variable in \mathbf{l}_i , X_1, \dots, X_k are distinct variables that do not appear in \mathbf{l}_i , and f is some k -ary function symbol in the language
2. $\{X/c\}$ where X is a variable in \mathbf{l}_i , and c is some constant in the language
3. $\{X/Y\}$ where X, Y are distinct variables in \mathbf{l}_i

In ILP, these 3 operations define a “downward refinement operator”

On the other hand, in *bottom-up search* we want to find some unknown generalization \mathbf{I}_1 of \mathbf{I}_2 . In that case, we should use inverse substitutions to find a chain of upward covers from \mathbf{I}_2 to \mathbf{I}_1 .

A systematic way to find refinements of clauses, is by using a refinement operator. There are two kinds of refinement operators: upward and downward ones. An upward refinement operator computes a set of generalizations of a given clause, a downward refinement operator computes a set of specializations. What constitutes a ‘specialization’ or ‘generalization’ of a clause, is determined by a generality order \succeq (such as subsumption, implication, relative subsumption, relative implication, *etc.*) on clauses. Then we can say that C is a generalization of D (dually: C is a specialization of D), if $C \succeq D$ holds. In each of these orders, the empty clause \square is the most general clause.

Downward refinement operators compute sets of specializations of a clause, upward ones compute sets of generalizations. Refinement operators are defined for a set of formulae S with a quasi-ordering \succeq . There are two refinement operators.

- ρ is a *downward refinement operator* if $\forall C \in S : \rho(C) \subseteq \{D \mid D \in S \text{ and } C \succeq D\}$
- δ is an *upward refinement operator* if $\forall C \in S : \delta(C) \subseteq \{D \mid D \in S \text{ and } D \succeq C\}$

For example, with an equality theory $=/2$, $D \in \rho(C)$ (downward refinement operator) if

$$D = \begin{cases} p(X_1, X_2, \dots, X_{n_p}) & \text{if } C = \square \text{ and } p/n_p \in \mathcal{L} \\ & \text{and the } X_i \text{ are distinct} \\ C \cup \{-l\} & \text{otherwise} \end{cases}$$

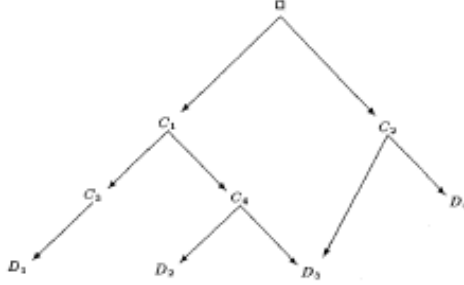
where

$$l = \begin{cases} V = W & \text{where } V, W \text{ occur in } C \\ V = f(X_1, X_2, \dots, X_{n_f}) & \text{where } V \text{ occurs in } C \\ & \text{and } f/n_f \in \mathcal{L} \text{ and} \\ & \text{the } X_i \text{ are distinct} \\ q(X_1, X_2, \dots, X_{n_q}) & \text{where } q/n_q \in \mathcal{L} \\ & \text{and the } X_i \text{ occur in } C \end{cases}$$

is an example of a refinement operator.

2.7.1 Example

We assume the set of clauses \mathcal{C}_h is ordered by such a generality order \succeq . Shapiro’s top-down approach only employs a downward refinement operator ρ , so $\rho(C)$ is a set of specializations of a given clause C . We start with $\Sigma = \{\square\}$. This is clearly too strong, since it implies any clause. Hence we want to find specialisations of \square , We use the set $\rho(\square)$ for this. If $\rho(\square)$ is still too strong, its false members in turn be replaced by their refinements, and so on. Thill allows us to search stepwisely through the generality order. This stepwise approach will

Figure 2.9: Paths through the refinement operator ρ .

only work if there is a path (a number of refinement steps) from Box to every clause in atleast one correct theory, For instance, suppose $\Sigma = \{D_2, D_3, D_4\}$ is a correct theory. Let the refinement operator ρ be such, that $\rho(\square) = \{C_1, C_2\}$, $\rho(C_1) = \{C_3, C_4\}$, $\rho(C_2) = \{D_3, D_4\}$, $\rho(C_3) = \{D_1\}$, and $\rho(C_4) = \{D_2, D_3\}$. Starting from \square , we can reach Σ by considering $\rho(\square)$, $\rho(C_1)$, $\rho(C_2)$, $\rho(C_4)$. See Figure 2.9

2.7.2 Ideal Refinement Operators

The sets of one-step refinements, n-step refinements, and refinements of some $C \in \mathcal{C}$, are respectively:

$$\rho^1(C) = \rho(C)$$

$$\rho^n(C) = \{D \mid \text{there is an } E \in \rho^{n-1}(C) \text{ such that } D \in \rho(E)\}, n \geq 2$$

$$\rho^*(C) = \rho^1(C) \cup \rho^2(C) \cup \rho^3(C) \dots$$

A ρ -chain from C to D is a sequence $C = C_0, C_1, \dots, C_n = D$, such that $C_{i-1} \in \rho(C_{i-1})$ for every $1 \leq i \leq n$. A refinement operator induces a refinement graph. This is a directed graph which has the members of \mathcal{C} as nodes (here variant clauses in \mathcal{C} can be viewed as the same node), and which contains an edge from C to D just in case $D \in \rho(C)$. This refinement graph is the space that is searched for candidates to include in the theory (we will later see refinements over theories).

Some desirable properties of ρ (and dually δ) are that they be:

1. **Locally Finite:** $\forall C \in \mathcal{C}$: $\rho(C)$ is finite and computable. Otherwise ρ will be of limited use in practice.
2. **Complete:** $\forall C \succ D$: $\exists E \in \rho^*(C)$ s.t. $E \sim D$. That is, every specialization should be reachable by a finite number of applications of the operator.
3. **Proper:** $\forall C \in \mathcal{C}$: $\rho(C) \subseteq \{D \mid D \in S \text{ and } C \succ D\}$. That is, it is better only to compute proper generalizations of a clause, for otherwise repeated

application of the operator might get stuck in a sequence of equivalent clauses (such as the application of inverse reduction in Section 2.3), without ever achieving any real specialization.

A refinement operator is *ideal* if it is locally finite, complete, and proper. An ideal refinement operator induces a refinement graph in which only a finite number of edges start from each node (local finiteness), in which there exists a path of finite length from C to a member of the equivalence class of D whenever $C \succ D$ (completeness), and which contains no cycles (by properness).

It can be shown that ideal downward $\rho_{\mathcal{A}}$ and upward $\delta_{\mathcal{A}}$ refinement operators exist for the simplest of our quasi-orders: the set of atoms ordered by subsumption.

Definition 7 *Let \mathcal{A} be the set of atoms in a language. The downward refinement operator $\rho_{\mathcal{A}}$ for \mathcal{A} is defined as follows:*

1. *For every variable z in an atom A and every n -ary function symbol f in the language, let x_1, \dots, x_n be distinct variables not appearing in A . Let $\rho_{\mathcal{A}}(A)$ contain $A\{z/f(x_1, \dots, x_n)\}$.*
2. *For every variable z in A and every constant a in the language let $\rho_{\mathcal{A}}(A)$ contain $A\{z/a\}$.*
3. *For every two distinct variables x and z in A , let $\rho_{\mathcal{A}}(A)$ contain $A\{z/x\}$.*

Note that $\rho_{\mathcal{A}}(A)$ may still contain variants. For instance, $\rho_{\mathcal{A}}(P(x, y))$ contains both $P(x, x)$ and $P(y, y)$. Clearly, in a practical application we can ignore any redundant variants. The three different kinds of atoms in $\rho(A)$ correspond exactly to the three kinds of downward covers that we discussed in theorem 25. It can be proved easily from the properties of sets of covers of atoms described on page 87 that if \mathcal{A} contains a finite number of constants function and predicate symbols, then $\rho_{\mathcal{A}}$ is ideal.

An ideal upward refinement operator $\delta_{\mathcal{A}}$ can be defined straightforwardly as follows:

Definition 8 *Let \mathcal{A} be the set of atoms in a language. The upward refinement operator $\delta_{\mathcal{A}}$ for \mathcal{A} is defined as follows:*

1. *For every $t = f(x_1, \dots, x_n)$ in A , for which x_1, \dots, x_n are distinct variables and each occurrence of some x_i in A is within an occurrence of t , $\delta_{\mathcal{A}}(A)$ contains an atom obtained by replacing all occurrences of t in A by some new variable z not in A .*
2. *For every constant a in A and every non-empty subset of the set of occurrences of a in A , $\delta_{\mathcal{A}}(A)$ contains an atom obtained by replacing those occurrences of a in A by some new variable z not in A .*
3. *For every variable x in A and every non-empty proper subset of the set of occurrences of x in A , $\delta_{\mathcal{A}}(A)$ contains an atom obtained by replacing*

those occurrences of x in A by some new variable z not in A . Note that in this last item, we cannot replace all occurrences of x by a new variable z , for then we would get a variant of A . For instance, $P(z, a, z)$ is a variant of $A = P(x, a, x)$.

As in the case of ρ_A , it easily follows that δ_A is locally finite, complete and proper. We do not even have to presuppose a finite number of constants function and predicate symbol for this, because when constructing $\text{JA}(A)$ we only have to deal with the finite number of symbols in A —there is no need to introduce new constants, functions or predicates.

2.7.3 Refinement Operators on Clauses for Subsumption

Unfortunately, the ideal conditions described for atoms cannot all be met at the same time for more complex orders. Ideal refinement operators do not exist for full clausal languages or Horn languages ordered by subsumption or by the stronger orders. This negative result is a consequence of the fact that finite complete sets of covers do not always exist⁸. We had seen (*c.f.* Section 2.1.2) that the existence of finite chains in lattices of atoms ordered by subsumption does *not* carry over to Horn clauses ordered by subsumption. This had followed from the observation that there are clauses which have no *finite* and complete set of downward covers. This makes it impossible to devise an ILP program that uses a refinement operator that is both complete and non-redundant. Therefore, there are no upward (downward) refinement operators that are locally finite as well as complete as well as proper for sets of clauses. That is, for clausal languages ordered by subsumption \succeq_θ or stronger orders, ideal refinement operators do not exist.

In order to define a refinement operator for full clausal languages, it will be required to drop one of the three properties of idealness. Of the three conditions of local finiteness, completeness and properness, finiteness seems indispensable: an infinite set $\rho(C)$ of refinements of a clause C cannot be handled well, because it would then be impossible to test all members of $\rho(C)$ in finite time. Furthermore, it is obvious that completeness is also a very valuable property, if you want to be able to guarantee that a solution will always be found whenever one exists. Of the three ideal properties, properness seems the least important and can be compromised. Ideal refinement operators can be approximated by

1. *Dropping the requirement of properness:* Locally finiteness and completeness are considered to be the two most important properties, so dropping the 'properness' is a common practice. We will define downward (ρ_L) and upward (δ_U) refinement operators that are locally finite and complete, but improper. For subsumption, such refinement operators exist, both for the downward and for the upward case. If C subsumes D , then $C\theta \subseteq D$ for some substitution θ . Thus specialization under subsumption can be

⁸In fact, it can be proved that if there exists an ideal downward (upward) refinement operator for $\langle \mathcal{C}, \succeq \rangle$, then every $C \in \mathcal{C}$ will have a finite complete set of downward (upward) covers

achieved by applying (elementary) substitutions and adding literals. In fact, when adding literals, it is sufficient to add only most general literals, since these can always be instantiated by a substitution later on to get the right literals. A literal $P(x_1, \dots, x_n)$ or $\neg P(x_1, \dots, x_n)$ is most general with respect to a clause C , if x_1, \dots, x_n are distinct variables not appearing in C .

Definition 9 *Let \mathcal{C} be a clausal language. The locally finite and complete downward refinement operator $\rho_{\mathcal{L}}$ for $\langle \mathcal{C}, \succeq \rangle$ is defined as follows:*

- (a) **Apply a substitution to a clause C** in one of the following ways:
 - i. For every variable z in a clause C and every n -ary function symbol f in the language, let x_1, \dots, x_n be distinct variables not appearing in C . Let $\rho_{\mathcal{L}}(C)$ contain $C\{z/f(x_1, \dots, x_n)\}$.
 - ii. For every variable z in C and every constant a in the language, let $\rho_{\mathcal{L}}(C)$ contain $C\{z/a\}$.
 - iii. For every two distinct variables x and z in C , let $\rho_{\mathcal{L}}(C)$ contain $C\{z/x\}$.
- (b) **Add a literal to a clause C :** For every n -ary predicate symbol P in the language, let x_1, \dots, x_n be distinct variables not appearing in C . Then $\rho_{\mathcal{L}}$ contains both $C \cup \{P(x_1, \dots, x_n)\}$ and $C \cup \{\neg P(x_1, \dots, x_n)\}$. Note that the literals $P(x_1, \dots, x_n)$ and $\neg P(x_1, \dots, x_n)$ that are added to C by the fourth item in the definition are most general with respect to C .

The proof of locally finiteness and completeness is straightforward. Since we already know that no ideal operators exist for this case, $\rho_{\mathcal{L}}$ cannot be proper. For instance, if $C = \{P(x)\}$ and $D = \{P(x), P(y)\}$, then $D \in \rho_{\mathcal{L}}(C)$ and $C \sim D$. However, this D is needed in a $\rho_{\mathcal{L}}$ -chain from C to $\{P(a)P(b)\}$ as follows: $\{P(x)\}, \{P(x), P(y)\}, \{P(a), P(y)\}, \{P(a), P(b)\}$. Notice that for every $\square \neq C$, we have $\square \succ C$, so $\rho^*(\square)$ contains a clause which is subsume-equivalent to C . In other words: If we start with the empty clause (as Shapiro's Model Inference Algorithm does), then for every $C \in \mathcal{C}$, a clause C' such that $C \sim C'$ can be reached by means of $\rho_{\mathcal{L}}$.

Definition 10 *Let \mathcal{C} be a clausal language. The locally finite and complete upward refinement operator δ_u for $\langle \mathcal{C}, \succeq \rangle$ is defined as follows:*

- (a) **Apply one of the following inverse substitution operations:**
 - i. For every $t = f(x_1, \dots, x_n)$ in C , for which all x_i are distinct variables and each occurrence of x_i in a clause C is within an occurrence of t , $\delta_u(C)$ contains the clause obtained by replacing all occurrences of t in C by some new variable z not previously in C .

ii. For every constant a in C and every non-empty subset of the set of occurrences of a in $\bar{C} = \text{dup}(C, a)$, if \bar{D} is the ordered clause obtained by replacing those occurrences of a in \bar{C} by the new variable z , then $\delta_u(C)$ contains D , where D is the set of literals in the ordered clause \bar{D} .

$\text{dup}(C, t)$ is defined as follows. If $C = \{L_1, \dots, L_n\}$ is a clause and t a term occurring in C and suppose t occurs k_1 times in L_1 , k_2 times in L_2 , etc. Then $\text{dup}(C, t) = \bar{C}$ is an ordered clause consisting of 2^{k_1} copies of L_1 , 2^{k_2} copies of L_2 , ... and 2^{k_n} copies of L_n . Note that if some $L \in C$ does not contain the term t , then \bar{C} contains L $2^0 = 1$ times, as it should.

iii. For every variable x in C and every non-empty proper subset of the set of occurrences of x in $\bar{C} = \text{dup}(C, x)$, if \bar{D} is the ordered clause obtained by replacing those occurrences of x in \bar{C} by the new variable z , then $\delta_u(C)$ contains D .

(b) **Remove a literal from the body of a clause:** If $C = D \cup \{L\}$ and L is a most general literal with respect to D , then $\delta_u(C)$ contains D .

2. *By bounding the language:* If we want to retain all three ideal properties, it seems that the only possibility is to restrict the search space. There exist ideal refinement operators for reduced finite clausal languages, ordered by subsumption. The fact that there always exists an ideal refinement operator for finite sets is mainly of theoretical interest. Thus in practice, we usually prefer more constructive-though possibly improper-refinement operators over such very elaborate ideal operators.

3. *Dropping the requirement of completeness:* Refinement operators are used very often in ILP systems, for instance in MIS [Sha81b], SIM [LD90, Lin92], FOIL [Qui90, QC93], CLAUDIEN [DB93], LINUS [LD94], and Progol [Mug95]. For reasons of efficiency, those operators are usually less general than the ones discussed here, and often incomplete. Nevertheless, the complete operators defined here form a good starting point for the construction of practical refinement operators.

An approximation to the ideal downward refinement operator, as adopted in Aleph is:

(a) Adding a literal drawn from \perp_i

$$p(X, Y) \leftarrow q(X) \text{ becomes } p(X, Y) \leftarrow q(X), r(Y)$$

(b) Equating two variables of the same type

$$p(X, Y) \leftarrow q(X) \text{ becomes } p(X, X) \leftarrow q(X)$$

(c) Instantiate a variable with a general functional term or constant

$$p(X, Y) \leftarrow q(X) \text{ becomes } p(3, Y) \leftarrow q(3)$$

Optimal cover-refinement operators do not exist for clausal languages ordered by subsumption.

2.7.4 Refinement Operators on Theories

Definition 11 Let \mathcal{C} be a clausal language, containing only a finite number of constants, function symbols, and predicate symbols. Let \mathcal{S} be the set of finite subsets of \mathcal{C} . The downward refinement operator $\rho_{\mathcal{I}}$ for $\langle \mathcal{S}, \models \rangle$ is defined as follows:

1. $(\Sigma \cup \{R \mid R \text{ is a resolvent of } C_1, C_2 \in \Sigma\}) \in \rho_{\mathcal{I}}(\Sigma)$.
2. If $\Sigma = \{C_1, \dots, C_n\}$, then $(\Sigma \cup \rho_{\mathcal{I}}(C_i)) \in \rho_{\mathcal{I}}(\Sigma)$, for each $1 \leq i \leq n$.
3. If $\Sigma = \{C_1, \dots, C_n\}$, then $(\Sigma \setminus \{C_i\}) \in \rho_{\mathcal{I}}(\Sigma)$, for each $1 \leq i \leq n$.

Note that every theory in $\rho_{\mathcal{I}}(\Sigma)$ that is specified by one of the first two items in the definition of $\rho_{\mathcal{I}}$ is logically equivalent to Σ . This shows that $\rho_{\mathcal{I}}$ is not proper. The completeness of $\rho_{\mathcal{I}}$ follows from the Subsumption Theorem, which tells us that logical implication can be implemented by a combination of resolution and subsumption.

2.8 Inverse Resolution

Since induction can be seen as the inverse of deduction, and resolution is our main tool for deduction, using inverse resolution for induction seems a sensible idea. Deduction moves from the general rules to the special case, while induction intends to find the general rules from special cases (examples). Muggleton and Buntine [MB88] introduced inverse resolution as a tool for induction. Their paper was followed by a wave of interest and research into the properties of inverse resolution [Wir89, HS91, Mug91b, Mug92b, Mug92c, RP89, RP90, Rou92, NF91, Ide93c, Ide92, Ide93b, Ide93a, LN92, SADB92, Tay93, SA93, BG93]. Inverting resolution is nowadays still a prominent *generalization operator* for bottom-up approaches to ILP. However, the theoretical foundation of this idea needs much more investigation. Moreover, in the application of inverse resolution, many indeterminacies arise: many different choices of literals, clauses and substitutions lay open. Accordingly, inverse resolution generates a very large search space of possibilities.

Muggleton and Buntine introduced two operators for this: the V-operator and the W-operator.

2.8.1 V-operator

The V-operator, outlined as a non-deterministic⁹ algorithm in Figure 2.10, generalizes two given clauses $\{C_1, R\}$ to $\{C_1, C_2\}$, such that R is an instance of a resolvent of C_1 and C_2 . The setting for the V-operator is pictorially depicted in Figure 2.11. The goal in inverse resolution is to construct a derivation of a positive example A (usually a ground atom) which hitherto was not implied by

⁹Since, in step 1, the algorithm has to choose one among many different possible θ_1 's, which all satisfy $C_1' \subseteq R$.

the theory. Using the algorithm for the V-operator, we can invert one resolution step, for given C_1 and R . By repeatedly applying the V-operator, we are able to invert any SLD-derivation.

INPUT: Horn clauses $C_1 = L_1 \vee C'_1$ and R , where $C'_1\theta_1 \subseteq R$ for some θ_1 .
OUTPUT: A Horn clause C_2 , such that R is an instance of a resolvent of C_1 and C_2 .
 Choose a substitution θ_1 such that $C'_1\theta_1 \subseteq R$.
 Choose an L_2 and C'_2 such that $L_1\theta_1 = \neg L_2\theta_2$ and $C'_2\theta_2 = R - C'_1\theta_1$, for some θ_2 .
return $C_2 = L_2 \vee C'_2$.

Figure 2.10: The V-operator.

The simplest situation is where $C_1 = L_1$, so where C'_1 is empty. Then for a given θ_1 , any C_2 and θ_2 with $C_2\theta_2 = \neg L_1\theta_1 \vee R$ will do. Since $L_1\theta_1 \vee R$ is an instance of any of these possible C_2 's, it is clear that $C_2 = \neg L_1\theta_1 \vee R$ is the "minimal" of all possible C_2 's, for a fixed θ_1 .

As an example application of inverse resolution, let $C_1 = P(x) \vee \neg Q(f(x))$, $L_1 = P(x)$, and $R = Q(g(y)) \vee \neg Q(f(g(y)))$. We assume $C_1\theta_1$ and $C_2\theta_2$ do not overlap.

1. Here only one θ_1 is possible, namely $\theta_1 = \{x/g(y)\}$.
2. L_2 and C'_2 should be such that, for some θ_2 , $L_1\theta_1 = P(g(y)) = \neg L_2\theta_2$ and $R - C'_1\theta_1 = Q(g(y)) = C'_2\theta_2$. Figure 2.12 shows all possible $C_2 = L_2 \vee C'_2$ (unique up to renaming of variables) of two literals, from top to bottom in decreasing order of generality. Small generalization and specialization steps are relevant for the V-operator. In this case, we are often interested in finding a 'minimal' C_2 , as in Figure 11.3, where $C_2 = \neg L_1\theta_1 \vee (R - C'_1\theta_1 = \neg P(g(y)) \vee Q(g(y)))$ is the minimal choice and can be obtained using the covers relation between clauses in the implication/subsumption

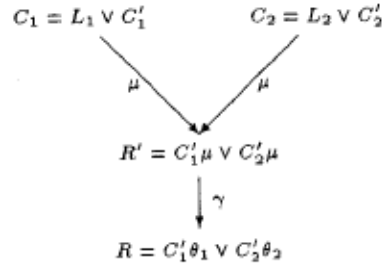
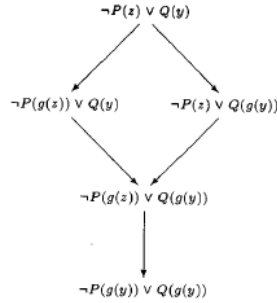


Figure 2.11: The setting for the V-operator.

Figure 2.12: C_2 derived using V-operator.

order studied in Section 2.1.2. The other C_2 's can then be found by taking small generalization steps starting from the minimal C_2 .

Note that for some C_2 , R itself is not a resolvent of C_1 and C_2 . For instance, if we let $C_2 = \neg P(z) \vee Q(y)$, then the resolvent of C_1 and C_2 is $\neg Q(f(z)) \vee Q(y)$, of which R is an instance.

We sometimes have to duplicate some literals in R before applying the V-operator, in order to be able to find the desired parent clauses, such that $C_1'\theta_1$ and $C_2'\theta_2$ overlap.

Note that there are many indeterminacies here, which make an unrestricted search through all possible invertible derivations very inefficient. Within the V-operator itself, many different choices for θ_1 , L_2 and C_2' are possible. And even before we can use the V-operator, we have to decide which clause from the old theory or the background knowledge to use as C_1 , and which literals to duplicate in R . Thus the total number of possibilities may become very large sometimes, which can make application of inverse resolution very inefficient.

In Figures 2.13 and 2.14, we contrast linear derivation against inverse linear derivation.

2.8.2 Predicate Invention: W-operator

One of the problems inductive learning algorithms have to face, is the fact that it is sometimes necessary to invent new predicates. For instance; suppose we want our algorithm to induce clauses from examples about family life. It would be very unfortunate if the system did not possess a predicate for the concept of 'parent'. If we have not given such a predicate to the system in advance, the system should be able to invent this predicate for itself. If we examine the V-operator carefully, it is clear that this operator cannot invent new predicates: all predicates appearing in any of the possible C_2 that we might construct already appear in C_1 or R . However, by putting two V-settings side-by-side, we get a W-shape. The W-operator combines two V-operators: it generalizes two given clauses $\{R_1, R_2\}$ to $\{C_1, C_2, C_3\}$, such that R_1 is an instance of a resolvent of

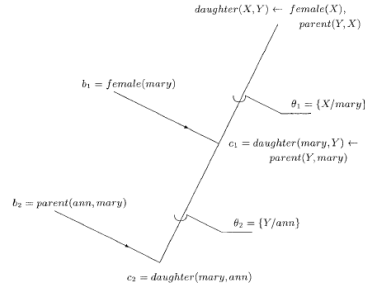


Figure 2.13: The linear derivation tree for $e_2 = Daughter(mary, ann)$ from background knowledge $\mathcal{B} = \{b_1, b_2\}$ where $b_1 = Female(mary)$ and $b_2 = Parent(ann, mary)$ and from a hypothesis $H = \{c\}$ where $c = Daughter(x, y) \leftarrow Female(x), Parent(y, x)$.

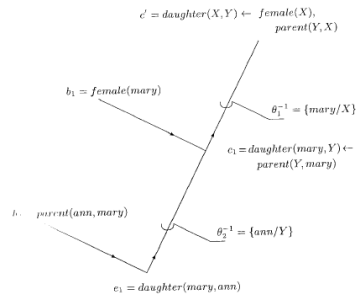


Figure 2.14: The inverse linear derivation tree for $H = \{c\}$ where $c = Daughter(x, y) \leftarrow Female(x), Parent(y, x)$ from background knowledge $\mathcal{B} = \{b_1, b_2\}$ where $b_1 = Female(mary)$ and $b_2 = Parent(ann, mary)$ and from an example $e_2 = Daughter(mary, ann)$.

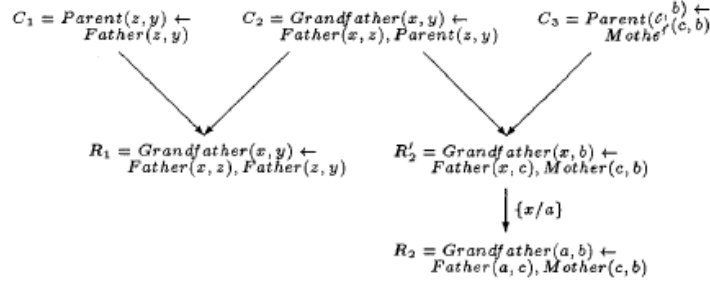


Figure 2.15: Generalization of $\{R_1, R_2\}$ to $\{C_1, C_2, C_3\}$ by the W-operator.

C_1 and C_2 , and R_2 is an instance of a resolvent of C_2 and C_3 . In addition, the W-operator is also able to invent new predicates.

We will present an example which shows the idea behind the W-operator. Suppose we have two Horn clauses $R_1 = \text{Grandfather}(x, y) \leftarrow \text{Father}(x, z), \text{Father}(z, y)$ and $R_2 = \text{Grandfather}(a, b) \leftarrow \text{Father}(a, c), \text{Mother}(c, b)$, and suppose we want to generalize these clauses. The Woperator constructs clauses C_1, C_2, C_3 , such that R_1 is an instance of a resolvent of C_1 and C_2 , and R_2 is an instance of a resolvent of C_2 and C_3 . Thus, the W-operator generalizes $\{R_1, R_2\}$ to $\{C_1, C_2, C_3\}$. In Figure 2.15, we give possible C_1, C_2, C_3 which can serve this purpose. The important point to notice about the figure is that the predicate *Parent*, which appears in C_1, C_2 and C_3 , did not appear in the clauses R_1 and R_2 we started with. Thus in generalizing $\{R_1, R_2\}$ to $\{C_1, C_2, C_3\}$, the W-operator has itself introduced a new predicate. The invention of this new predicate is quite useful, since it allows us to write out the definition of a ‘Grandfather’ in a very succinct way in C_2 : x is the grandfather of y , if x is the father of some z , and z is a parent of y . Note that any predicate name may be assigned the role of *Parent* here, including ‘old’ names such as *Grandfather* or *Mother*, since this predicate is resolved away in the two resolution steps anyway.

The general setting for the W-operator is pictured in Figure 2.16. Given R_1 and R_2 , the W-operator constructs C_1, C_2, C_3 , with the property that R_1 is an instance of a resolvent of C_1 and C_2 , and R_2 is an instance of a resolvent of C_2 and C_3 . What we want to find, are $C_1 = L_1 \vee C'_1, C_2 = L_2 \vee C - 2', C_3 = L_3 \vee C'_3, \theta_1, \theta_2, \sigma_1$ and σ_2 , such that $L_1\theta_1 = \neg L_2\theta_2, L_2\sigma_1 = \neg L_3\sigma_2, R_1 = C'_1\theta_1 \vee C'_2\theta_2$ and $R_2 = C'_2\sigma_1 \vee C'_3\sigma_2$. Thus L_1 and L_2 are resolved upon in deriving R_1 , while L_2 and L_3 are resolved upon in deriving R_2 . μ is an mgu for L_1 and $\neg L_2$, and ν is an mgu for $\neg L_2$ and L_3 . Hence L_1 and L_3 must either be both positive, or both negative. Note that L_1, L_2, L_3 do not appear in R_1 and R_2 , which gives the opportunity for inventing a new predicate.

The idea behind the construction of C_1, C_2 and C_3 using the W-operator is sketched below:

1. Given R_1 and R_2 , we first try to find a C'_2 such that $C'_2\theta_2 \subseteq R_1$ and

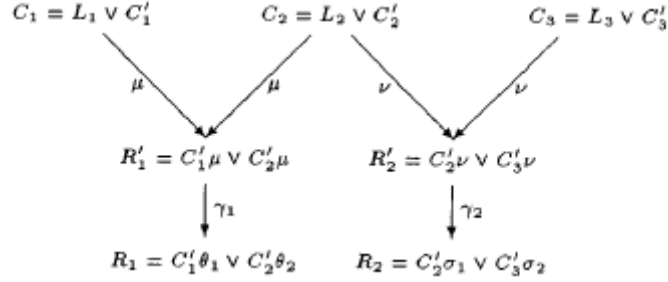


Figure 2.16: The setting for the W-operator.

$C'_2 \sigma_1 \subseteq R_2$, for some θ_2 and σ_1 . Let $D_1 = C'_2 \theta_2$ and $D_2 = C'_2 \sigma_1$. Clearly, many different C'_2 's can give the same D_1 and D_2 . These C'_2 's can be considered as generalizations of $\{D_1, D_2\}$. We would like to begin with a minimal C'_2 . This motivates the use of the notion of a 'least generalization' of clauses discussed in Section 2.1.1. If we have found a minimal C'_2 , the other possible C'_2 's can be found by taking small generalization steps, starting from the minimal C'_2 . This motivates the use of 'covers' (minimal generalizations or specializations of that clause) and consequently the refinement operator of the clause.

If such a C'_2 cannot be found - which means, intuitively, that R_1 and R_2 have 'nothing in common' - we should let C'_2 be empty.

2. If we have chosen an appropriate C'_2 , we can complete C_2 by choosing also L_2 . In principle, **any** L_2 will do.
3. Once we have decided which clause to take as C_2 , then a C_1 and C_3 can be *found independently*. C_1 can be constructed by the V-operator from C_2 and R_1 , and C_3 can be constructed by the V-operator from C_2 and R_2 .

Consider Figure 2.15 again. Given R_1 and R_2 , how did we find C_1 , C_2 and C_3 ? First we note that $(Grandfather(x, y) \vee \neg Father(x, z)) \in R_1$, and $(Grandfather(x, y) \vee \neg Father(x, z)) \{x/a, y/b, z/c\} \in R_2$. Hence $C'_2 = Grandfather(x, y) \vee \neg Father(x, z)$ is an appropriate choice. Secondly, we have to choose L_2 . Let us say we take $L_2 = \neg Parent(z, y)$. This gives $C_2 = Grandfather(x, y) \leftarrow \neg Parent(z, y)$. Thirdly, the V-operator can find $C_1 = Parent(z, y) \leftarrow \neg Father(x, z)$ from C_2 and R_1 . Similarly, it can construct the clause $C_3 = Parent(c, b) \leftarrow \neg Mother(c, b)$ from C_2 and R_2 . Thus the W-operator generalizes $\{R_1, R_2\}$ to $\{C_1, C_2, C_3\}$.

2.9 Summary

In this chapter, we discussed the following generalization operations:

1. Relative least general generalization (lgg/lub) with respect to subsumption (LGRS) as well as implication order (LGRI). This is utilized by systems such as GOLEM [Muggleton and Feng 92]. The progressive construction of lubs for the subsumption order, starting with terms and culminating in clauses is illustrated in Table 2.1.

GOLEM is a bottom-up, non-interactive, batch single-predicate learner, which employs LRGS/

2. The upward refinement operator δ_u .
3. Inverse resolution in terms of the V and W operators. This is utilized in ILP systems such as CIGOL (W-operator) [Muggleton and Buntine 88], MARVIN (V-operator) [Sammur and Banerji 86] and ITOU (V-operator) [Rouveirol 92].

CIGOL is a bottom-up, interactive, incremental, multiple predicate learner that uses the W-operator for predicate invention. The user is asked whether the induced clauses are true in the interpretation intended by the user and how invented predicates should be named.

The following specialization techniques were also discussed:

1. Relative greatest specialization with respect to subsumption as well as implication order. This is rarely used in practice.
2. The downward refinement operator $\rho_{\mathcal{L}}$. This is utilized in ILP systems such as MIS (model inference system) [Shapiro 83], FOIL [Quinlan and Cameron-Jone 93] and its successors mFOIL [Dzeroski 91], CLAUDIEN [L. De Raedt and Bruynooghe 93], PROGOL [Muggleton 95], MOBAL [Morik, Wrobel, Kietz and Emde 93], RDT [J-U Kietz and Wrobel 92], FOCL [Brunk and Pazzani 91], MARKUS [Grobelnik 92] and MPL [De Raedt et al 1993].

MIS is a top-down, interactive, incremental, multiple-predicate learner, restricted to Horn clauses. FOIL is a top-down, non-interactive, batch single-predicate learner, upgrading Quinlan's earlier decision tree learner ID3 [Quinlan 86] which learns function free normal programs using the (set) 'covering' approach. CLAUDIEN is a top-down, non-interactive batch learner which uses Herbrand interpretations as (positive only) examples. CLAUDIEN is one of the very few systems that induces full, rather than definite or normal program clauses. PROGOL is a top-down, non-interactive, batch, multi-predicate learner that learns clauses using an A^* -like heuristic algorithm to search top-down through a refinement graph, while restricting attention to clauses that subsume some bottom-clause $\perp_d(\mathcal{B}, e)$ (*c.f.* page 120). RDT is top-down, non-interactive, batch, multi-predicate

learner that learns function-free normal programs, using a set of ground literals as background knowledge (or reduced to that form as in GOLEM) and where language bias is introduced into the refinement operators using a rule schema. A *rule schema* is a clause with predicate-variables instead of ordinary predicate symbols. Only clauses that can be obtained by instantiating the predicate-variables in one of the given rule schemas to ordinary predicate symbols may be used in the theory. It also uses a *predicate topology* to further restrict the search for what predicate symbols can be added to the body, given a predicate symbol in the head. MOBAL makes use of RDT as one of its components.

3. There is another specialization technique called *unfolding*, which was not discussed so far. SPECTRE [Bostrom 95], IMPUT [Alexin, Gyimothy and Bostrom 96] and JIGSAW [Ade and Bostron 95] use unfolding. SPECTRE is a top-down, non-interactive, batch, single-predicate learner.

Systems such as CLINT make use of a hybrid of top-down (unfolding) and bottom-up (abduction) approaches.

Lub	Definition	Examples
lub of terms $lub(t_1, t_2)$	<ol style="list-style-type: none"> $lub(t, t) = t$, $lub(f(s_1, \dots, s_n), f(t_1, \dots, t_n)) = f(lub(s_1, t_1), \dots, lub(s_n, t_n))$, $lub(f(s_1, \dots, s_m), g(t_1, \dots, t_n)) = V$, where $f \neq g$, and V is a variable which represents $lub(f(s_1, \dots, s_m), g(t_1, \dots, t_n))$, $lub(s, t) = V$, where $s \neq t$ and at least one of s and t is a variable; in this case, V is a variable which represents $lub(s, t)$. 	<ul style="list-style-type: none"> $lub([a, b, c], [a, c, d]) = [a, X, Y]$. $lub(f(a, a), f(b, b)) = f(lub(a, b), lub(a, b)) = f(V, V)$ where V stands for $lub(a, b)$. When computing lggs one must be careful to use the same variable for multiple occurrences of the lubs of subterms, i.e., $lub(a, b)$ in this example. This holds for lubs of terms, atoms and clauses alike.
lub of atoms $lub(a_1, a_2)$	<ol style="list-style-type: none"> $lub(P(s_1, \dots, s_n), P(t_1, \dots, t_n)) = P(lub(s_1, t_1), \dots, lub(s_n, t_n))$, if atoms have the same predicate symbol P, $lub(P(s_1, \dots, s_m), Q(t_1, \dots, t_n))$ is undefined if $P \neq Q$. 	
lub of literals $lub(l_1, l_2)$	<ol style="list-style-type: none"> if l_1 and l_2 are atoms, then $lub(l_1, l_2)$ is computed as defined above, if both l_1 and l_2 are negative literals, $l_1 = \overline{a_1}$, $l_2 = \overline{a_2}$, then $lub(l_1, l_2) = lub(\overline{a_1}, \overline{a_2}) = \overline{lub(a_1, a_2)}$, if l_1 is a positive and l_2 is a negative literal, or vice versa, $lub(l_1, l_2)$ is undefined. 	<ul style="list-style-type: none"> $lub(Parent(ann, mary), Parent(ann, tom)) = Parent(ann, X)$. $lub(Parent(ann, mary), \overline{Parent(ann, tom)}) = undefined$. $lub(Parent(ann, X), Daughter(mary, ann)) = undefined$.
lub of clauses $lub(C_1, C_2)$	<ol style="list-style-type: none"> Let $C_1 = \{l_1, \dots, l_n\}$ and $C_2 = \{k_1, \dots, k_m\}$. Then $lub(C_1, C_2) = \{l_{ij} = lub(l_i, k_j) \mid l_i \in C_1, k_j \in C_2 \text{ and } lub(l_i, k_j) \text{ is defined}\}$. 	<ul style="list-style-type: none"> If $C_1 = Daughter(mary, ann) \leftarrow Female(mary), Parent(ann, mary)$ and $C_2 = Daughter(eve, tom) \leftarrow Female(eve), Parent(tom, eve)$, then $lub(C_1, C_2) = Daughter(X, Y) \leftarrow Female(X), Parent(Y, X)$ where X stands for $lub(mary, eve)$ and Y stands for $lub(ann, tom)$.
$rlgg(a_1, a_2)$	$lub(a_1 \leftarrow B, a_2 \leftarrow B)$	$rlgg(Daughter(mary, ann), Daughter(eve, tom)) = Daughter(X, Y) \leftarrow Female(X), Parent(Y, X)$ where B denotes the conjunction of the literals $Parent(ann, mary), Parent(ann, tom), Parent(tom, eve), Parent(tom, ian), Female(ann), Female(mary), Female(eve)$.

Table 2.1: Table showing progressive definitions of lubs, starting with terms and culminating in rlgg between clauses.

Chapter 3

Linear Algebra

Dixit algorizmi. Or, “So said al-Khwarizmi”, being the opening words of a 12th century Latin translation of a work on arithmetic by al-Khwarizmi (*ca.* 780–840).

3.1 Linear Equations

Elementary algebra, using the rules of completion and balancing developed by al-Khwarizmi, allows us to determine the value of an unknown variable x that satisfies an equation like the one below:

$$10x - 5 = 15 + 5x$$

An equation like this that only involves an unknown (like x) and not its higher powers (x^2 , x^3), along with additions (or subtractions) of the unknown multiplied by numbers (like $10x$ and $5x$) is called a *linear* equation. We now know, of course, that the equation above can be converted to a special form (“number multiplied by unknown equals number”, or $ax = b$, where a and b are numbers):

$$5x = 20$$

Once in this form, it becomes easy to see that $x = b/a = 4$. Linear algebra is, in essence, concerned with the solution of several linear equations in several unknowns. Here is a simple example of two equations and two unknowns x and y , written in a uniform way, with all unknowns (variables) to the left of the equality, and all numbers (constants) to the right:

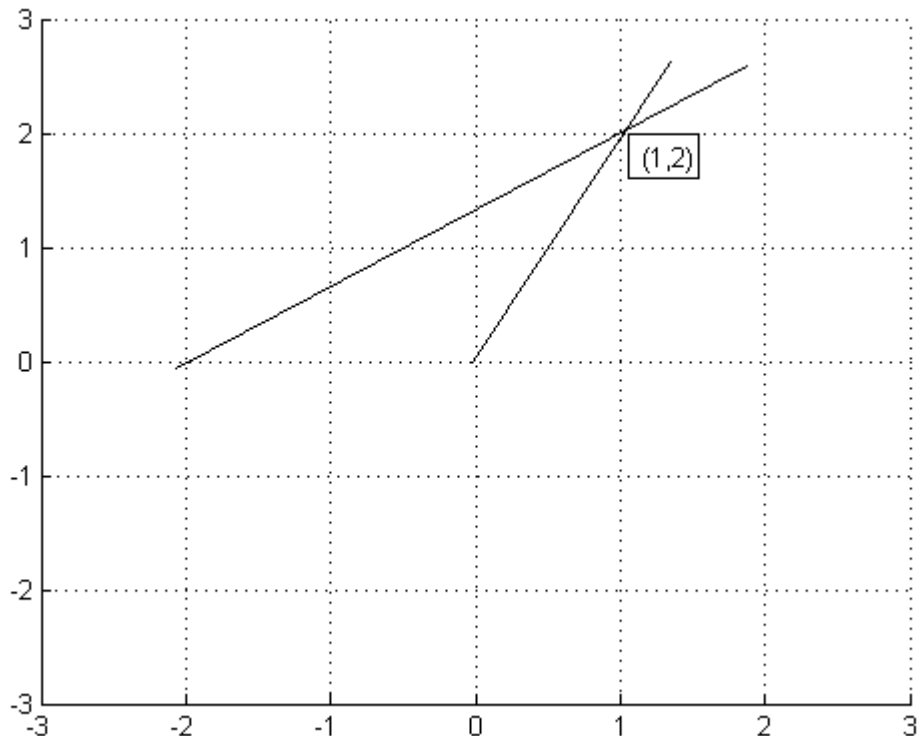


Figure 3.1: Solving linear equations: the geometric view.

$$2x - y = 0$$

$$-x + 2y = 3$$

We would like to find values of x and y for which these equations are true. School geometry tells us how to visualise this: each equation is a straight line in the xy plane, and since we want a value of x and y for which both equations are true, we are really asking for the values of x and y that lie on both lines (that is, the point of intersection of the two lines: see Fig. 3.1). Of course, if the lines do not meet at a point, then there are no values of x and y that satisfy the equations. And we can continue to solve problems like these geometrically: more unknowns means lines become higher-dimensional flat surfaces (“hyperplanes”), and more equations means we are looking for the single point of intersection of all these surfaces. Visually though, this is challenging for all but a small minority of us, geared as we are to live in a world of three spatial dimensions.

Linear algebra, an extension of elementary algebra, gives us a way of looking at the solution of any number of linear equations, with any number of variables without suffering from this visual overload. In effect, equations are once again converted to the simple form we just saw, that is, $Ax = b$, although A and b are no longer just numbers. In fact, we will see that A is a *matrix*, and that x and b are *vectors* (and in order not to confuse them with variables and numbers, we will from now on use the bold-face notation \mathbf{x} and \mathbf{b}). Linear algebra, shows us that solutions, if they exist, can be obtained in three different ways:

1. A direct solution, using techniques called elimination and back substitution.
2. A solution by “inverting” the matrix A , to give the solution $\mathbf{x} = A^{-1}\mathbf{b}$.
3. A vector space solution, by looking at notions called the column space and nullspace of A .

Understanding each of these requires a minimal understanding of vectors and matrices, which we give in a somewhat compressed form here.

3.2 Vectors and Matrices

It is easiest to think of a vector as a generalisation of a single number. A pair of numbers can be represented by a *two-dimensional vector*. Here is the two-dimensional vector representation of the pair $(2, -1)$:

$$\mathbf{u} = \begin{bmatrix} 2 \\ -1 \end{bmatrix}$$

This kind of vector is usually called a “column” vector. Geometrically, such a vector is often visualised by an arrow in the two-dimensional plane as shown on the left in Fig. ???. Multiplying such a vector by any particular number, say 2, multiplies each component by that number. That is, $2\mathbf{u}$ represents the pair $(4, -2)$. Geometrically, we can see that multiplication by a number—sometimes called *scalar multiplication*—simply makes gives a vector with a “longer” arrow as shown on the right in the figure (assuming, of course, that we are not dealing with zero-length vectors). In general, multiplication of a (non-zero) vector \mathbf{u} by different (non-zero) numbers a result in lines either in the direction of \mathbf{u} (if $a > 0$) or in the opposite direction

Suppose we now consider a second vector \mathbf{v} corresponding to the pair $(-1, 2)$, and ask: what is $\mathbf{u} + \mathbf{v}$. This simply adds the individual components. In our example:

$$\mathbf{u} = \begin{bmatrix} 2 \\ -1 \end{bmatrix} \quad \mathbf{v} = \begin{bmatrix} -1 \\ 2 \end{bmatrix} \quad \mathbf{u} + \mathbf{v} = \begin{bmatrix} 2 - 1 \\ -1 + 2 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

Geometrically, the addition of two vectors gives a third, which can be visualised as the diagonal of the parallelogram formed by \mathbf{u} and \mathbf{v} (Fig. ??, left). It should be straightforward to visualise that any point on the plane containing the vectors \mathbf{u} and \mathbf{v} can be obtained by some linear combination $a\mathbf{u} + b\mathbf{v}$, and that the space of all linear combinations is simply the full two-dimensional plane containing \mathbf{u} and \mathbf{v} (Fig. ??, right). For the two-dimensional example here, this plane is just the usual xy plane (we will see that this is the vector space \mathbb{R}^2).

Although we have so far only looked at vectors with two components, linear algebra is more general. It allows us to use the same operations with vectors of any size. Suppose our vectors \mathbf{u} and \mathbf{v} are three-dimensional. Linear combinations now still fill a plane containing the two vectors. But, this is no longer the xy plane, since the vectors generated by the linear combinations are points in three-dimensional space (we will see later, that is some “subspace” of the vector space \mathbb{R}^3). Addition of a third vector \mathbf{w} will also not necessarily result in a point on this plane, and the space of linear combinations $a\mathbf{u} + b\mathbf{v} + c\mathbf{w}$ could fill the entire three-dimensional space.

Let us return now to the two equations that we saw in the previous section:

$$\begin{aligned} 2x - y &= 0 \\ -x + 2y &= 3 \end{aligned}$$

It should be easy to see how these can be written in “vector” form:

$$x \begin{bmatrix} 2 \\ -1 \end{bmatrix} + y \begin{bmatrix} -1 \\ 2 \end{bmatrix} = \begin{bmatrix} 0 \\ 3 \end{bmatrix} \quad (3.1)$$

That is, we are asking if there is some linear combination of the column vectors $[2, -1]$ and $[-1, 2]$ that gives the column vector $[0, 3]$. And this is the point of departure with the usual geometric approach: we visualise solutions of equations not as points of intersections of surfaces, but as linear combination of vectors (of whatever dimension): see Fig. 3.2.

To get it into a form that is even more manageable, we need the concept of a “coefficient matrix”. A matrix is simply a rectangular array of numbers, and the coefficient matrix for the left hand side of the linear combination above is:

$$A = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}$$

This is a 2×2 (“two by two”) matrix, meaning it has 2 rows and 2 columns. You can see that the columns of the matrix are simply the column vectors of the linear combination. Let:

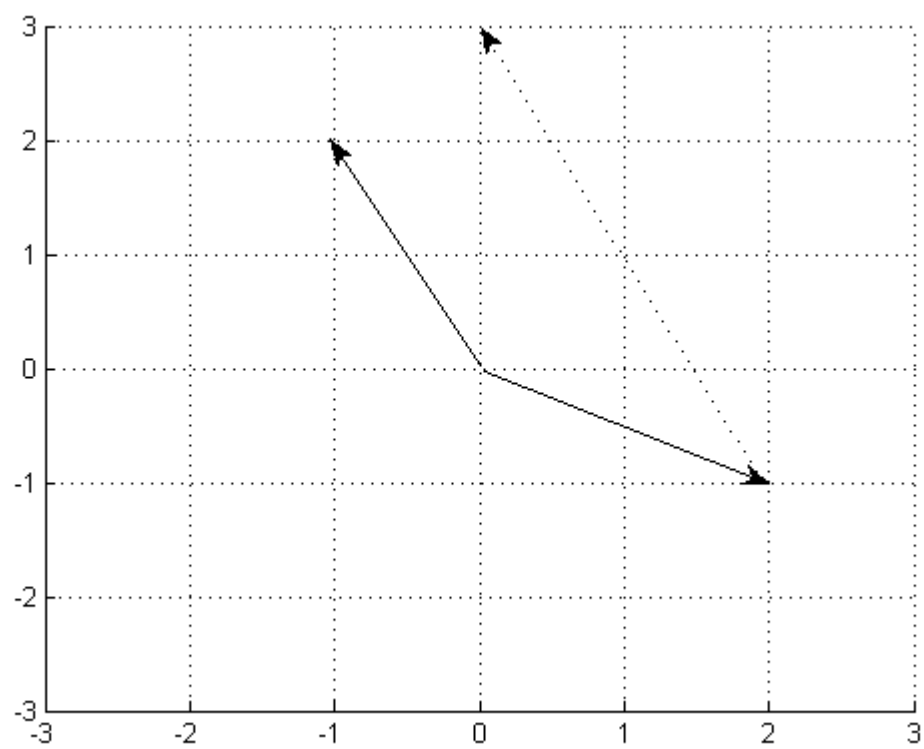


Figure 3.2: Solving linear equations: the geometric view from linear algebra.

$$\mathbf{x} = \begin{bmatrix} x \\ y \end{bmatrix} \text{ and } \mathbf{b} = \begin{bmatrix} 0 \\ 3 \end{bmatrix}$$

Then, the matrix equation representing the same linear combination is:

$$A\mathbf{x} = \mathbf{b} \tag{3.2}$$

This, as you can see, is just as simple, at least in form. as the very first equation we started with ($5x = 20$). We still need to know what $A\mathbf{x}$ means. Comparing Equations 3.2 and 3.2, $A\mathbf{x} = x$ (column 1 of A) + y (column 2 of A).

This extends easily enough to equations with more variables. Here are three linear equations in three unknowns:

$$\begin{aligned} 2x - y &= 0 \\ -x + 2y - z &= -1 \\ -3y + 4z &= 4 \end{aligned}$$

The coefficient matrix A is:

$$A = \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -3 & 4 \end{bmatrix}$$

The right hand side of the matrix equation is:

$$\mathbf{b} = \begin{bmatrix} 0 \\ -1 \\ 4 \end{bmatrix}$$

What we are trying to do is to find values of x, y and z such that:

$$x(\text{column 1 of } A) + y(\text{column 2 of } A) + z(\text{column 3 of } A) = \begin{bmatrix} 0 \\ -1 \\ 4 \end{bmatrix}$$

It is easy to see now that the solution we are after is $x = 0, y = 0, z = 1$. Or, in vector form, the solution to the matrix equation $A\mathbf{x} = \mathbf{b}$ is:

$$\mathbf{x} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

In general, things are not so obvious and it may be the case that for some values of A and \mathbf{b} , no values of x, y and z would solve $A\mathbf{x} = \mathbf{b}$. For example, \mathbf{b} may be a point in 3-dimensional space that could not be “reached” by any linear combinations of the vectors comprising the columns of A . Here is a simple example:

$$A = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 0 & 0 & 0 \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

Sequences of mathematical operations—algorithms, if you will—have been devised to check if solutions exist, and obtain these solutions mechanically when they exist. There are three such approaches we will look at: obtaining solutions by elimination (the simplest), obtaining solutions by matrix inversion (a bit more complex), and finally, a vector space solution (the hardest). We look at each of these in turn.

3.3 Solution of Linear Equations by Elimination

We will now examine a systematic method as “elimination”—first presented by Gauss, for solving a linear system of equations. The basic idea is to progressively *eliminate* variables from equations. For example, let us look once again at the two equations we saw earlier:

$$2x - y = 0$$

$$-x + 2y = 3$$

Elimination first multiplies both sides of the second equation by 2 (this clearly leaves it unchanged):

$$-2x + 4y = 6$$

We can also add equal amounts to the left and right sides without changing the equation. So, adding the left hand side of the first equation to the left hand

side of this new equation, and the right hand side of the first equation to the right hand side of this new equation also does not alter anything:

$$(-2x + 4y) + (2x - y) = 6 + 0 \quad \text{or} \quad 3y = 6$$

So, the two original equations are the same as:

$$\begin{aligned} 2x - y &= 0 \\ 3y &= 6 \end{aligned}$$

You can see that x has been “eliminated” from the second equation and the set of equations have been said to be transformed into an *upper triangular* form. In this form, it is easy to see that $y = 6/3 = 2$. The value of x can then be obtained by substituting back this value for y in the first equation, to give $2x - 2 = 0$ or $x = 1$. The different steps in the elimination process can be expressed clearly using matrices, which we do now. As a running example, we will use the following set of 3 equations:

$$\begin{aligned} x + 2y + z &= 2 \\ 3x + 8y + z &= 12 \\ 4y + z &= 2 \end{aligned}$$

We now know what the coefficient matrix for these equations is:

$$A = \begin{bmatrix} 1 & 2 & 1 \\ 3 & 8 & 1 \\ 0 & 4 & 1 \end{bmatrix}$$

A point of notation. The entry in row 1, column 1 of A will be denoted a_{11} ; row 1, column 2 will be a_{12} and so on. So, in the matrix above, $a_{11} = 1$, $a_{12} = 2$ etc.. In general, the entry in row i , column j will be denoted a_{ij} .

Before we plunge into the details of the matrix operations, let us just go through the procedure mechanically (taking on faith for the moment that the steps are indeed valid ones). Our first elimination step is to eliminate x from the second equation. We multiply the first equation by a multiplier and subtract it from the second equation with the goal of eliminating the x coefficient in the second equation. We will call it the (2,1) step. The first element of the first row a_{11} determines the value of the multiplier (3 in this case) and it is called a *pivot*. For reasons that will become clear, pivots should not be 0. The resultant coefficient matrix is:

$$A_1 = \begin{bmatrix} 1 & 2 & 1 \\ 0 & 2 & -2 \\ 0 & 4 & 1 \end{bmatrix}$$

The next step will be to get a 0 in the first column of the third row (a_{31}) of A_1 . Since this is already the case, we do not really need to do anything. But, just to be pedantic, let us take it as giving a coefficient matrix A_2 , which is just the same as A_1 :

$$A_2 = \begin{bmatrix} 1 & 2 & 1 \\ 0 & 2 & -2 \\ 0 & 4 & 1 \end{bmatrix}$$

We now move on to eliminating a_{32} in A_2 . Using a_{22} in A_2 as the next pivot, we subtract from the third row a multiple (2) of the second row. The resultant coefficient matrix is now:

$$A_3 = \begin{bmatrix} 1 & 2 & 1 \\ 0 & 2 & -2 \\ 0 & 0 & 5 \end{bmatrix}$$

A_3 is called an upper triangular matrix for obvious reasons (and sometimes denoted by U). We will see shortly that with the sequence of operations that we have just done, the left hand side of the original matrix equation $A\mathbf{x}$ is transformed into $A_3\mathbf{x}$ by progressively multiplying by a sequence of matrices called “elimination matrices”.

3.3.1 Elimination as Matrix Multiplication

Let us go back to the original matrix equation:

$$\begin{bmatrix} 1 & 2 & 1 \\ 3 & 8 & 1 \\ 0 & 4 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} 2 \\ 12 \\ 2 \end{bmatrix}$$

We take a step back and look again at the first elimination step (“multiply equation 1 by 3 and subtract from equation 2”). The effect of this step is to change the right-hand side second equation from 12 to $12 - 3 \times 2 = 6$ and leave the right-hand sides of all other equations unchanged. In matrix notation, the right hand side, after the first elimination step, is:

$$\mathbf{b}_1 = \begin{bmatrix} 2 \\ 6 \\ 2 \end{bmatrix}$$

A little calculation should be sufficient to convince yourself that \mathbf{b}_1 can be obtained by pre-multiplying \mathbf{b} by the matrix:

$$E = \begin{bmatrix} 1 & 0 & 0 \\ -3 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

That is, $\mathbf{b}_1 = E\mathbf{b}$. You can check this by doing the usual linear combination of the columns of E with the components of \mathbf{b} , but the following “row-by-column” view—which is simply the linear combination expanded out—may be even more helpful:

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} = \begin{bmatrix} a_{11}b_1 + a_{12}b_2 + a_{13}b_3 \\ a_{21}b_1 + a_{22}b_2 + a_{23}b_3 \\ a_{31}b_1 + a_{32}b_2 + a_{33}b_3 \end{bmatrix}$$

So, if the elimination step multiplies the left-hand side of the matrix equation $A\mathbf{x} = \mathbf{b}$ by the matrix E , then to make sure nothing is changed, we have to do the same to the left-hand side. That is, the elimination step changes the left-hand side to $E A \mathbf{x}$. But now we are stuck— $E A$ is a product of two matrices, which we have not come across before. What does this mean?

Well, we know what we would like $E A$ to mean. We would like $E A = A_1$. That is:

$$\begin{bmatrix} 1 & 0 & 0 \\ -3 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 & 1 \\ 3 & 8 & 1 \\ 0 & 4 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 2 & 1 \\ 0 & 2 & -2 \\ 0 & 4 & 1 \end{bmatrix} \quad (3.3)$$

Taking a vector as simply being a matrix with a single column, we would like to extend the old matrix-vector multiplication ($A\mathbf{x}$) idea to general matrix-matrix multiplication. Suppose B is a matrix comprised of the column vectors \mathbf{b}_1 , \mathbf{b}_2 and \mathbf{b}_3 . Then AB is a matrix that has columns $A\mathbf{b}_1$, $A\mathbf{b}_2$, and $A\mathbf{b}_3$. So, in the example above, $E A$ is a matrix that has columns $E\mathbf{a}_1$, $E\mathbf{a}_2$ and $E\mathbf{a}_3$ (where \mathbf{a}_1 , \mathbf{a}_2 and \mathbf{a}_3 are the columns of A). Let us work out what these are:

$$E\mathbf{a}_1 = \begin{bmatrix} 1 & 0 & 0 \\ -3 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 3 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \times 1 + 0 \times 3 + 0 \times 0 \\ -3 \times 1 + 1 \times 3 + 0 \times 0 \\ 0 \times 1 + 0 \times 3 + 1 \times 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

This is the first column of the matrix A_1 on the right-hand side of Equation 3.3.1. You can check that $E\mathbf{a}_2$ and $E\mathbf{a}_3$ do indeed give the other two columns of A_1 . Once again, there is a “row-by-column” view of multiplying two matrices that you will often find helpful:

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \begin{bmatrix} b_{11} & b_{12} & b_{13} \\ b_{21} & b_{22} & b_{23} \\ b_{31} & b_{32} & b_{33} \end{bmatrix} = \begin{bmatrix} a_{11}b_{11} + a_{12}b_{21} + a_{13}b_{31} & a_{11}b_{12} + a_{12}b_{22} + a_{13}b_{32} & a_{11}b_{13} + a_{12}b_{23} + a_{13}b_{33} \\ a_{21}b_{11} + a_{22}b_{21} + a_{23}b_{31} & a_{21}b_{12} + a_{22}b_{22} + a_{23}b_{32} & a_{21}b_{13} + a_{22}b_{23} + a_{23}b_{33} \\ a_{31}b_{11} + a_{32}b_{21} + a_{33}b_{31} & a_{31}b_{12} + a_{32}b_{22} + a_{33}b_{32} & a_{31}b_{13} + a_{32}b_{23} + a_{33}b_{33} \end{bmatrix}$$

At this point, it is important that you are aware of some properties of matrix multiplication. First, multiplying matrices A and B is only meaningful if the number of columns of A is the same as the number of rows of B . If A is an $m \times n$ matrix, and B is an $n \times k$ matrix, then AB is an $m \times k$ matrix. Second, just like with ordinary numbers, matrix multiplication is “associative”; that is, $(AB)C = A(BC)$ (with numbers, $(3 \times 4) \times 5 = 3 \times (4 \times 5)$). But, unlike ordinary numbers, matrix multiplication is not “commutative”. That is $AB \neq BA$ (but with numbers, $3 \times 4 = 4 \times 3$).

It is the associativity of matrix multiplication that allows us to build up a sequence of matrix operations representing elimination. Let us return once again to the matrix equation we started with:

$$\begin{bmatrix} 1 & 2 & 1 \\ 3 & 8 & 1 \\ 0 & 4 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} 2 \\ 12 \\ 2 \end{bmatrix}$$

We have seen, how, by multiplying both sides by the elimination matrix E (which we will now call E_{21} , for reasons that will be obvious), gives:

$$E_{21}(A\mathbf{x}) = (E_{21}A)\mathbf{x} = E_{21}\mathbf{b}$$

or:

$$A_1\mathbf{x} = E_{21}\mathbf{b}$$

where $A_1 = E_{21}A$. Without more elaboration, we now simply present the elimination matrices E_{31} and E_{32} that correspond to the remaining two elimination steps.

$$E_{31} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad E_{32} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -2 & 1 \end{bmatrix}$$

The general rule for constructing an elimination matrix is this. If we are looking at n equations in m unknowns, and an elimination step involves multiplying equation j by a number q and subtracting it from equation i , then the elimination matrix E_{ij} is simply the $n \times m$ “identity matrix” I , with $a_{ij} = 0$ in I replaced by $-q$. For example, with 3 equations in 3 unknowns, and an elimination step that “multiplies equation 2 by 2 and subtracts from equation 3”:

$$I = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad E_{32} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -2 & 1 \end{bmatrix}$$

Each elimination step therefore results in a multiplication of both sides of $A\mathbf{x} = \mathbf{b}$ by the corresponding elimination matrix. In our example, the three elimination steps give:

$$E_{32}E_{31}E_{21}(A\mathbf{x}) = E_{32}E_{31}E_{21}\mathbf{b}$$

which, using the property of associativity of matrix multiplication is:

$$(E_{32}(E_{31}(E_{21}A)))\mathbf{x} = (E_{32}E_{31}E_{21})\mathbf{b}$$

Or:

$$U\mathbf{x} = (E_{32}E_{31}E_{21})\mathbf{b} = \mathbf{c} \quad (3.4)$$

where U is the upper triangular matrix $E_{32}A_2 = E_{32}(E_{31}A_1 = E_{32}(E_{31}(E_{21}A)))$. Here:

$$U = \begin{bmatrix} 1 & 2 & 1 \\ 0 & 2 & -2 \\ 0 & 0 & 5 \end{bmatrix} \quad \mathbf{c} = \begin{bmatrix} 2 \\ 6 \\ -10 \end{bmatrix} \quad (3.5)$$

Before we leave this section, there is one aspect of elimination that we have not yet considered. Let us look at the same equations as before, but in the following order:

$$\begin{aligned}4y + z &= 2 \\x + 2y + z &= 2 \\3x + 8y + z &= 12\end{aligned}$$

The coefficient matrix A is then:

$$A = \begin{bmatrix} 0 & 4 & 1 \\ 1 & 2 & 1 \\ 3 & 8 & 1 \end{bmatrix} \quad (3.6)$$

Now clearly, no amount of elimination can get this matrix into an upper triangular form, since that would require a non-zero entry for a_{11} . Since there is no reason to keep the equations in this particular order, we can exchange their order until we reach the one we had in the previous section. Just as a single elimination step can be expressed as multiplication by an elimination matrix, exchange of a pair of equations can be expressed by multiplication by a *permutation* matrix.

The general rule for constructing a permutation matrix is this. If we are looking at m equations in n unknowns, and we want to exchange equations i and j , then the permutation matrix P_{ij} is simply the $m \times n$ “identity matrix” I , with rows i and j swapped:

$$I = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad P_{12} = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Multiplying a matrix A by P_{12} will swap rows 1 and 2 of A :

$$\begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & 4 & 1 \\ 1 & 2 & 1 \\ 3 & 8 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 2 & 1 \\ 0 & 4 & 1 \\ 3 & 8 & 1 \end{bmatrix}$$

What happens if, in spite of all exchanges, elimination still results in a 0 in any one of the pivot positions? Then we consider the process to have failed, and the equations do not have a solution. Assuming this does not happen, we will reach a point where the original equation $\mathbf{Ax} = \mathbf{b}$ is transformed into $U\mathbf{x} = \mathbf{c}$ (as we did in Equation 3.4). The final step is that of *back-substitution*, in which variables are progressively assigned values using the right-hand side of this transformed equation (in Equation 3.3.1, $z = -2$, back-substituted to give $y = 1$, which finally yields $x = 2$).

3.4 Solution of Linear Equations by Matrix Inversion

So, it is possible to represent the steps leading to the solution of a set of linear equations by elimination entirely as a sequence of matrix multiplications. We now look at obtaining the solution by considering matrix “inversion”. What we are trying to do is to really find a matrix analog for division with ordinary numbers. There, with an equation like $5x = 20$, we are able to get the answer immediately using division: $x = 20/5$. Can we not do the same with matrices? That is, given $A\mathbf{x} = \mathbf{b}$, can we not get $\mathbf{x} = \mathbf{b}/A$. Well, not quite. But we can get close: we find $\mathbf{x} = A^{-1}\mathbf{b}$, where A^{-1} is the matrix equivalent of $1/A$, and is called the *inverse* of the matrix.

3.4.1 Inverse Matrices

The starting point is just the same as with numbers. We know $a/a = aa^{-1} = 1$ for a non-zero number a . For matrices, we want to find A^{-1} such that $AA^{-1} = I$ where I is the identity matrix. Actually, with matrices, we can ask for inverses in two different ways: AA^{-1} and $A^{-1}A$, called for obvious reasons, right and left inverses of A (recall that since matrix multiplication does not necessarily commute, these could be different).

Let us start with $m \times n$ (“square”) matrices. Our definition of an inverse is simply this: if there exists a matrix A_L^{-1} such that $A_L^{-1}A = I$, where I is the $N \times N$ identity matrix, then A_L^{-1} is called the left inverse of A . On the other hand, if there exists a matrix A_R^{-1} such that $AA_R^{-1} = I$, then A_R^{-1} is called the right inverse of A . Now, for square matrices, it is easy to see that the left and right inverses are the same:

$$A_L^{-1}(AA_R^{-1}) = (AA_L^{-1})A_R^{-1}$$

Or,

$$A_L^{-1} = A_R^{-1}$$

So, for square matrices at least, we can simply talk about “the inverse” A^{-1} . The question that now concerns us is: do all square matrices have an inverse? The short answer is “no”. Here is a matrix that is not invertible:

$$A = \begin{bmatrix} 1 & 3 \\ 2 & 6 \end{bmatrix} \quad (3.7)$$

We can see the conditions under which an inverse exists by referring back to the matrix equation that formed the basis of solution by elimination:

$$A\mathbf{x} = \mathbf{b}$$

Let us assume that A^{-1} exists. Then, the solution reached by elimination would simply be:

$$\mathbf{x} = A^{-1}\mathbf{b} \quad (3.8)$$

Therefore, if the inverse exists, then elimination must produce an upper triangular matrix with non-zero pivots. In fact, the condition works both ways—if elimination produces non-zero pivots then the inverse exists (you can see very quickly that elimination applied to the matrix A in Equation 3.4.1 would give a row of 0s). Otherwise, the matrix is not invertible, or *singular*. Another way to look at this is that the matrix will be singular if its “determinant” is 0. We will look at what this means later (in Section 3.10), but it is related to the elimination producing non-zero pivots.

If the inverse exists, then the only solution to the matrix equation $A\mathbf{x} = \mathbf{b}$ is $\mathbf{x} = A^{-1}\mathbf{b}$. This gives another way to test for the singularity of a matrix: if there are solutions other than $\mathbf{x} = \mathbf{0}$ to $A\mathbf{x} = \mathbf{0}$. For example, with A in Equation 3.4.1, the vector $\mathbf{x} = [3, -1]$ is a solution to $A\mathbf{x} = \mathbf{0}$.

A final observation may be evident from the example in Equation 3.4.1. A matrix is singular if the columns (or rows) are not linearly independent.

Now let us consider a slight variant of the matrix A in Equation 3.4.1:

$$A = \begin{bmatrix} 1 & 3 \\ 2 & 7 \end{bmatrix}$$

We believe that this matrix is invertible. How can we determine it's inverse? Let the inverse be

$$A^{-1} = \begin{bmatrix} a & c \\ b & d \end{bmatrix} \quad (3.9)$$

The system of equations $AA^{-1} = I$ can be written as:

$$\begin{bmatrix} 1 & 3 \\ 2 & 7 \end{bmatrix} \begin{bmatrix} a & c \\ b & d \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

Again, recall the view of matrix multiplication in which each column on the right hand side is a linear combination of the columns of A :

$$\begin{bmatrix} 1 & 3 \\ 2 & 7 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

and

$$\begin{bmatrix} 1 & 3 \\ 2 & 7 \end{bmatrix} \begin{bmatrix} c \\ d \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

So, once we solve these two sets of linear equations, we can assemble A^{-1} from the values of a, b, c , and d . We are back, therefore, to solving linear systems of equations—the Gaussian elimination procedure for a single set of linear equations with a single column vector on the right-hand side has to be generalised. The process used is called the *Gauss-Jordan* procedure.

3.4.2 Gauss-Jordan Elimination

The Gauss-Jordan elimination method addresses the problem of solving several linear systems $A\mathbf{x}_i = \mathbf{b}_i$ ($1 \leq i \leq N$) at once, such that each linear system has the same coefficient matrix A but a different right hand side b_i .

From Section 3.3, we know that Gaussian elimination is nothing more than multiplication by elimination matrices, that transforms a coefficient matrix A into an upper-triangular matrix U :

$$U = E_{32}(E_{31}(E_{21}A)) = (E_{32}E_{31}E_{21})A$$

Here E_{ij} is an elimination matrix constructed as we described before (replacing the appropriate 0 in the identity matrix with a non-zero number). Of course, we might, in general, be required to perform row permutation operations and they will simply as appear as multiplication by permutation matrices. But, let us ignore this complication for the moment. Suppose now we applied further elimination steps until U was transformed into the identity matrix. This means multiplication by more matrices:

$$I = E_{13}(E_{12}(E_{23}(E_{32}(E_{31}(E_{21}A)))))) = (E_{13}E_{12}E_{23}E_{32}E_{31}E_{21})A = BA \quad (3.10)$$

By definition $B = (E_{13}E_{12}E_{23}E_{32}E_{31}E_{21})$ must be A^{-1} . And this is what Gauss-Jordan does: it simply runs the elimination steps further until the upper-triangular matrix is converted into the identity matrix. So, A^{-1} can be computed by applying the same sequence of elimination steps to the identity matrix. A standard technique for carrying out the same elimination steps on two matrices A and B is to create an augmented matrix $[A \ B]$ and carry out the elimination on this augmented matrix. Gauss-Jordan can therefore be summarised in a single line: perform elimination steps on the augmented matrix $[A \ I]$ (representing the equation $AB = I$) to give the augmented matrix $[I \ A^{-1}]$ (representing the equation $IB = A^{-1}$). Or, in matrix multiplication terms: We illustrate the process with the example matrix we looked at earlier:

$$\left[\begin{array}{cc|cc} 1 & 3 & 1 & 0 \\ 2 & 7 & 0 & 1 \end{array} \right] \xrightarrow{\text{Row}_2 - 2 \times \text{Row}_1} \left[\begin{array}{cc|cc} 1 & 3 & 1 & 0 \\ 0 & 1 & -2 & 1 \end{array} \right] \xrightarrow{\text{Row}_1 - 3 \times \text{Row}_2} \left[\begin{array}{cc|cc} 1 & 0 & 7 & -3 \\ 0 & 1 & -2 & 1 \end{array} \right]$$

One could verify that the inverse of A is given by

$$A^{-1} = \begin{bmatrix} 7 & -3 \\ -2 & 1 \end{bmatrix} \quad (3.11)$$

Gauss-Jordan therefore gives us a method to construct the inverse of a coefficient matrix A , and therefore directly solve $A\mathbf{x} = \mathbf{b}$ as $\mathbf{x} = A^{-1}\mathbf{b}$.

What if A is not a square matrix but rather a rectangular matrix of size $m \times n$, such that $m \neq n$. Does there exist a notion of A^{-1} ? The answer depends on the rank of A .

- If A is full row rank and $n > m$, then AA^T is a full rank $m \times m$ matrix and therefore $(AA^T)^{-1}$ exists. The matrix $A^T(AA^T)^{-1}$ yields the identity matrix when multiplied to A on its right, *i.e.*, $AA^T(AA^T)^{-1} = I$ and is called the right inverse of A . When the right inverse of A is multiplied on its left, we get the projection matrix $A^T(AA^T)^{-1}A$, which projects matrices onto the row space of A .
- If A is full column rank and $m > n$, then $A^T A$ is a full rank $n \times n$ matrix and therefore $(A^T A)^{-1}$ exists. The matrix $(A^T A)^{-1}A^T$ yields the identity matrix when multiplied to A on its left, *i.e.*, $(A^T A)^{-1}A^T A = I$ and is called the left inverse of A . When the left inverse of A is multiplied on its right, we get the projection matrix $A(A^T A)^{-1}A^T$, which projects matrices onto the column space of A .

What if A is neither full row rank nor full column rank? In Section 3.13, we define the pseudoinverse of any $m \times n$ matrix, without any restrictions on its rank.

3.5 Solution of Linear Equations using Vector Spaces

We now turn to the third approach for solving linear equations. This is, in some sense, the most abstract, and involves the idea a vector spaces. A vector space is a collection of vectors that, informally speaking, may be multiplied by a number and added. More formally, a vector space is a set of vectors on which two operations are defined: vector addition and scalar multiplication. Additionally, these two operations satisfy certain natural conditions which we will elaborate shortly. A well-known example is the vector space \mathbb{R}^2 . This consists of all 2 dimensional column vectors with real-valued components (this is also just the entire xy plane). Similarly, the space \mathbb{R}^n comprises all n dimensional column vectors with real-valued components.

More generally, if a set of vectors \mathcal{V} is to qualify as a “vector space” then two vector operations—addition and scalar multiplication—have to be defined, and they have to result in vectors within the set \mathcal{V} . The set, or space \mathcal{V} is then said to be “closed” under the operations of addition and scalar multiplication. Now, given vectors \mathbf{u} and \mathbf{v} in a vector space, all scalar multiples of vectors $a\mathbf{u}$ and $b\mathbf{v}$ are in the space, as is their sum $a\mathbf{u} + b\mathbf{v}$. That is, all linear combinations of elements in the space are also elements of the space (V is closed under linear combination). If a subset (V_S) of any such space is itself a vector space (that is, (V_S) is also closed under linear combination) then (V_S) is called a subspace of (V) . All this may seem a bit abstract, and some examples may help:

1. The set of vectors \mathbb{R}^2 consisting of all two-dimensional vectors with real-valued components is a vector space. Adding any two vectors in the set gives a vector that is also in the set. Scalar multiplication of any vector in the set is a vector also in \mathbb{R}^2 (you may find these easy to see by visualising the vectors in the xy plane).
2. The set of vectors $(\mathbb{R}^2)^+$ consisting of all two-dimensional vectors in the positive quadrant is *not* a vector space. Adding a pair of vectors in $(\mathbb{R}^2)^+$ results in a vector in $(\mathbb{R}^2)^+$. But, unfortunately, multiplying by a scalar may not. For example, every vector $-3\mathbf{v}$ ($\mathbf{v} \in (\mathbb{R}^2)^+$) does not belong to $(\mathbb{R}^2)^+$.
3. The subset \mathbb{R}_S^3 of \mathbb{R}^3 consisting of vectors of any length through the origin $(0, 0, 0)$ is a subspace of \mathbb{R}^3 . Adding vectors in \mathbb{R}_S^3 clearly results in an element of the set, as does multiplication by a scalar. It is important that the origin $(0, 0, 0)$ is included: otherwise, multiplication of a vector by 0 would result in a vector not in the subset.

3.5.1 Vector Spaces and Matrices

Our condition on vector spaces has nothing really to do with vectors: all we need is that a pair of operations, addition and scalar multiplication be defined

on a set of elements. So, we are now ready to go a further step, and drop the restriction that vector spaces consist only of vectors. We can, for example, talk of a “vector” space \mathcal{M} consisting of all 2×2 matrices. It is easy to check that this is indeed closed under (matrix) addition and scalar multiplication (we have not come across this before: it is simply the multiplication of every element of the matrix by the number comprising the scalar). Just as with vectors, a subspace of \mathcal{M} is then some subset that is also a vector space.

Vector spaces of matrices provide a novel way of looking at the solution of $A\mathbf{x} = \mathbf{b}$. Recall that $A\mathbf{x}$ is simply a linear combination of the columns of the matrix A . All possible linear combinations of the columns produce a set of all possible column vectors (in effect, all possible \mathbf{b} 's). This set is called the *column space* of A , or $C(A)$. Given \mathbf{b} , therefore, when we ask: is there a solution to $A\mathbf{x} = \mathbf{b}$, we are really asking if the particular \mathbf{b} we are given is in the column space of A . An example may help. Consider the matrix A :

$$A = \begin{bmatrix} 1 & 1 & 2 \\ 2 & 1 & 3 \\ 3 & 1 & 4 \\ 4 & 1 & 5 \end{bmatrix}$$

The column space $C(A)$ is a subspace of \mathbb{R}^4 (are you sure you understand why this is so?). We can ask an number of questions now. What is in this subspace? Obviously, each column of A is in $C(A)$. Additionally, $C(A)$ contains all linear combinations of the columns of A . Is $C(A)$ the entire 4-dimensional space \mathbb{R}^4 ? If not, how much smaller is $C(A)$ compared to \mathbb{R}^4 ?

Equivalently, we can pose this problem as follows. Consider the linear system $A\mathbf{x} = \mathbf{b}$. For which right hand sides \mathbf{b} does a solution \mathbf{x} always exist? A solution \mathbf{x} definitely does not exist for every right hand side \mathbf{b} , since there are 4 equations in 3 unknowns. Let us analyse this further by writing down the system of equations

$$A\mathbf{x} = \begin{bmatrix} 1 & 1 & 2 \\ 2 & 1 & 3 \\ 3 & 1 & 4 \\ 4 & 1 & 5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} \quad (3.12)$$

Our first point is that there are many vectors \mathbf{b} which cannot be expressed as the linear combination of the columns of A . That leads to the question, *which right hand side \mathbf{b} allows the equation to be solved*. One value of \mathbf{b} for which the equation can be solved is the zero vector, for which the corresponding solution is $\mathbf{x} = \mathbf{0}$. Three other trivial values of \mathbf{b} are the values corresponding to every column of A . In particular, we can solve $A\mathbf{x} = \mathbf{b}$ whenever \mathbf{b} is in the column

space $C(A)$. When \mathbf{b} is a combination of columns of A , the combination tells us what exactly \mathbf{x} must be.

Do all the columns of A contribute something ‘new’ to the space $C(A)$ ¹? In other words, can we get the same space $C(A)$ using less than three columns of A ? In this particular example, the third column of A is a linear combination of the first two columns of A . $C(A)$ is therefore a 2-dimensional subspace of \mathfrak{R}^4 . In general, if A is an $m \times n$ matrix, $C(A)$ is a subspace of \mathfrak{R}^m .

3.5.2 Null Space

The null space of a matrix A , referred to as $N(A)$, is the space of all solutions to the equation $A\mathbf{x} = 0$. The null space of an $m \times n$ matrix A is a subspace of \mathfrak{R}^n .

Cosider the example matrix A discussed in the previous section. Its null space is a subspace of \mathfrak{R}^3 . We will try to figure out the null space of the matrix A by observing the following system of equations:

$$A\mathbf{x} = \begin{bmatrix} 1 & 1 & 2 \\ 2 & 1 & 3 \\ 3 & 1 & 4 \\ 4 & 1 & 5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \quad (3.13)$$

One obvious solution to the system is the zero vector. The null space will always contain the zero vector. Making use of the observation that the columns of A are linearly dependent, we find a second solution to the system as:

$$\mathbf{x}^* = \begin{bmatrix} 1 \\ 1 \\ -1 \end{bmatrix} \quad (3.14)$$

Thus, \mathbf{x}^* is in the null space $N(A)$. Every multiple of \mathbf{x}^* is also in the null space. Each element of the null space $N(A)$ is of the form

$$c\mathbf{x}^* = \begin{bmatrix} c \\ c \\ -c \end{bmatrix} \quad (3.15)$$

¹In subsequent sections, we will refer to these columns as *pivot* columns.

where $c \in \mathfrak{R}$. Thus, the null space $N(A)$ is the line passing through the zero vector $[0\ 0\ 0]$ and $[1\ 1\ -1]$.

Do solutions to $Ax = 0$ always yield a vector space? The answer is yes and this can be proved by observing that if $A\mathbf{v} = 0$ and $A\mathbf{w} = 0$, then $A(\mathbf{v} + \mathbf{w}) = 0$ and $A(p\mathbf{v}) = 0$ where $p \in \mathfrak{R}$. In general, there are two equivalent ways of specifying a subspace.

1. The first way is to specify a bunch of vectors whose linear combinations will yield the subspace.
2. The second way is to specify a system of equations of the form $A\mathbf{x} = \mathbf{0}$ and any vector \mathbf{x} that satisfies the system is an element of the subspace.

What about the set of all solutions to the equation $A\mathbf{x} = \mathbf{b}$ - do elements of this set form a space? The answer is a *no*. An easy way to see this is that the zero vector is not a solution to this system (unless \mathbf{b} is the zero vector) and hence the solutions cannot form a space.

3.6 Elimination for Computing the Null Space ($A\mathbf{x} = \mathbf{0}$)

In the last section we defined the null space of a matrix A . In this section, we will turn the definition into an algorithm using the elimination technique discussed in Section 3.3. We will take as an example, the following rectangular matrix A

$$A = \begin{bmatrix} 1 & 2 & 2 & 2 \\ 2 & 4 & 6 & 8 \\ 3 & 6 & 8 & 10 \end{bmatrix} \quad (3.16)$$

3.6.1 Pivot Variables and Row Echelon Form

We note rightaway that column 2 of A is a multiple of column 1 - it is in the same direction as column 1 and is therefore not independent. We expect to discover that in the elimination process for computing the null space $N(A)$. In terms of rows, we observe that the third row is the sum of the first two rows. Thus, the rows are also not independent - and again we hope to discover that in the elimination process.

In essence, what elimination does is change the matrix A and consequently its column space, while leaving the null space of A intact. We first choose the element in position $[1, 1]$ as the pivot element and perform the steps $(2, 1)$ and $(3, 1)$ (recall the definition of a step from Section 3.3) to get the transformed matrix A_1 .

$$A_1 = \begin{bmatrix} [1] & 2 & 2 & 2 \\ 0 & 0 & 2 & 4 \\ 0 & 0 & 2 & 4 \end{bmatrix} \quad (3.17)$$

We have got the first column in the desirable form. So next, we try to use the element in position $[2, 2]$ as the pivot element. But unfortunately it is a 0. We look below it position $[3, 2]$ hoping for a non-zero element so that we can do a row exchange. But there is a zero below as well! That tells us that second column is dependent on the first column.

Since we have nothing to do with the second column, we move to the third column and choose the entry $[2, 3]$ as the pivot. We perform the next elimination step $(3, 2)$, and obtain a third row of zeros. We will denote the resultant matrix by U . Note that the pivots are marked in boxes.

$$U = \begin{bmatrix} [1] & 2 & 2 & 2 \\ 0 & 0 & [2] & 4 \\ 0 & 0 & 0 & 0 \end{bmatrix} \quad (3.18)$$

The matrix U is in the row echelon form. A matrix is said to be in row echelon form if it satisfies the following requirements:

1. All nonzero rows are above any rows of all zeroes.
2. The leading coefficient of a row is always strictly to the right of the leading coefficient of the row above it.

While reducing from the matrix A to U , we used only two pivots. In fact, we have already discovered the most important number about the matrix A . The number of pivots is 2 - which is also the *rank* of the matrix.

Fact: *The rank of a matrix is the number of pivots used while reducing it to the row echelon form using elimination.*

We can now solve $U\mathbf{x} = \mathbf{0}$, which has the same solution as $A\mathbf{x} = \mathbf{0}$ (since the elimination steps on zero vector always yield a zero vector). Thus, the null space of U is the same as that of A . How do we describe these solutions? We will first write down the equations corresponding to $U\mathbf{x} = \mathbf{0}$.

$$\begin{aligned} x_1 + 2x_2 + 2x_3 + 2x_4 &= 0 \\ 2x_3 + 4x_4 &= 0 \end{aligned}$$

We will describe the solution by first separating out the two columns containing the pivots, referred to as *pivot columns* and the remaining columns, referred

to as *free columns*. Variables corresponding to the free columns are called *free variables*, since they can be assigned any value. Variables corresponding to the pivot columns are called *pivot variables*, and their values can be determined based on the values assigned to the free variables. In our example, variables x_2 and x_4 are free variables while x_1 and x_3 are the pivot variables.

Let us say we use the following assignment of values to free variables: $x_2 = 1$, $x_4 = 0$. Then, by back substitution, we get the following values: $x_1 = -2$ and $x_3 = 0$. Thus, the following vector \mathbf{x}' is a solution to the system $U\mathbf{x} = \mathbf{0}$ (and consequently the solution to $A\mathbf{x} = \mathbf{0}$) and therefore lies in $N(A)$.

$$\mathbf{x}' = \begin{bmatrix} -2 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad (3.19)$$

This solution reiterates our first observation, *viz.*, that column 2 is a multiple of column 1.

We will find some more vectors in the null space. Any multiple $c\mathbf{x}'$, $c \in \mathfrak{R}$ is also in $N(A)$. Note that $c\mathbf{x}'$ is a line. Are these the only vectors in $N(A)$? Actually, no – we obtained this set of vectors by assigning only one set of values to the free variables x_2 and x_4 . We assign another set of values $x_2 = 0$, $x_4 = 1$, and obtain the values of x_1 and x_3 by back-substitution to get another vector \mathbf{x}'' in $N(A)$.

$$\mathbf{x}'' = \begin{bmatrix} 2 \\ 0 \\ -2 \\ 1 \end{bmatrix} \quad (3.20)$$

Now we are in a position to specify all vectors in $N(A)$. The null space will contain all linear combinations of the two special solutions \mathbf{x}' and \mathbf{x}'' . Every vector in $N(A)$ therefore has the form given in (3.21):

$$ax' + bx'', a \in \mathfrak{R}, b \in \mathfrak{R} \quad (3.21)$$

What is the number of special (linearly independent) solutions for $A\mathbf{x} = \mathbf{0}$ if A is an $m \times n$ matrix? As we saw earlier, the rank r of a matrix equals the number of pivots. The number of free variables is given by

$$\text{no. of free variables} = n - r$$

The number of special solutions is exactly equal to the number of free variables. In the above example, we had $n = 4$, $r = 2$ and therefore number of free variables was 2. The steps for characterizing the null space of a matrix A can be summarized as follows:

1. Reduce A to the row echelon form.
2. If r is the number of pivots, find the $k = n - r$ free variables.
3. Make k different assignments to the free variables. For each assignment, use backsubstitution (using the row echelon form) to find the values of the pivot variables. Each assignment to the free variables yields a vector in the null space.

3.6.2 Reduced Row Echelon Form

We will take a second look at the matrix U that we obtained by elimination.

$$U = \begin{bmatrix} [1] & 2 & 2 & 2 \\ 0 & 0 & [2] & 4 \\ 0 & 0 & 0 & 0 \end{bmatrix} \quad (3.22)$$

The last row of U is composed of zeroes. This is because row 3 of A was a linear combination of rows 1 and 2 and this fact was discovered by elimination. How can we clean U up further? We can do elimination upwards to get zeros above as well as below the pivots. Elimination step $(2, 1)$ on U yields the matrix U' .

$$U' = \begin{bmatrix} [1] & 2 & 0 & -2 \\ 0 & 0 & [2] & 4 \\ 0 & 0 & 0 & 0 \end{bmatrix} \quad (3.23)$$

Further, we will make all pivot elements equal to 1 by dividing the corresponding row by the pivot element. This yields the matrix R .

$$R = \begin{bmatrix} [1] & 2 & 0 & -2 \\ 0 & 0 & [1] & 2 \\ 0 & 0 & 0 & 0 \end{bmatrix} \quad (3.24)$$

The matrix R has zeroes above and below each pivot. This matrix is called the reduced row echelon form (rref) of A . Matlab has the function $rref(A)$ that returns the reduced row echelon form of A . The system of equations $R\mathbf{x} = 0$ is given as

$$\begin{aligned}x_1 + 2x_2 - 2x_4 &= 0 \\x_3 + 2x_4 &= 0\end{aligned}$$

The solution to this system is the same the solution to the original system of equations $Ax = 0$. By simple back-substitution, the vector x can be expressed as:

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} -2 & 2 \\ 1 & 0 \\ 0 & -2 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_2 \\ x_4 \end{bmatrix} \quad (3.25)$$

Note that the specification of the null space in equation 3.25 is the same as that in equation 3.21.

Let us suppose that we have already got a matrix A in the reduced row echelon form (rref) R . Further, let us pretend that the pivot columns I come before the free columns F . The matrix R can be written as

$$R = \begin{bmatrix} I & F \\ 0 & 0 \end{bmatrix} \quad (3.26)$$

This block matrix is a very typical rref. We can easily do all the special solutions at once. We will create a *null basis* N whose columns are the special solutions; *i.e.*, $RN = 0$. The following N satisfies this system:

$$N = \begin{bmatrix} -F \\ I \end{bmatrix} = \begin{bmatrix} -2 & 2 \\ 0 & -2 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} \quad (3.27)$$

In fact there is a Matlab command `null(A)` that returns the null basis of A . It first computes the rref of A and then composes N using the free columns of A as well as the identity matrix of size equal to the rank of A .

Next, we will illustrate the same sequence of steps on the transpose matrix A^t to obtain its row echelon form U and observe the pivots, rank and null space. The solution to $A^t \mathbf{x} = \mathbf{0}$ is a column vector of size 3. Notice that the rank of the transpose is again 2 and is the same as the number of pivot columns. There is a single free column corresponding to the free variable x_3 .

$$\begin{bmatrix} [1] & 2 & 3 \\ 2 & 4 & 6 \\ 2 & 6 & 8 \\ 2 & 8 & 10 \end{bmatrix} \xrightarrow{E_{2,1}, E_{3,1}, E_{4,1}} \begin{bmatrix} 1 & 2 & 3 \\ 0 & 0 & 0 \\ 0 & 2 & 2 \\ 0 & 4 & 4 \end{bmatrix} \xrightarrow{P_{2,3}} \begin{bmatrix} 1 & 2 & 3 \\ 0 & [2] & 2 \\ 0 & 0 & 0 \\ 0 & 4 & 4 \end{bmatrix} \xrightarrow{E_{4,2}} \begin{bmatrix} [1] & 2 & 3 \\ 0 & [2] & 2 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad (3.28)$$

Suppose we make the following assignment to the free variable $x_3 = -c$. Then the solution is given by

$$\begin{bmatrix} -c \\ -c \\ c \end{bmatrix} = c \begin{bmatrix} -1 \\ -1 \\ 1 \end{bmatrix} \quad (3.29)$$

Thus, the null space of A^t is a line. Taking the elimination steps forward, we can get the reduced row echelon form (as a block matrix) R for matrix A^t .

$$\begin{bmatrix} [1] & 2 & 3 \\ 0 & [2] & 2 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \xrightarrow{E_{1,2}} \begin{bmatrix} [1] & 0 & 1 \\ 0 & [2] & 2 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \xrightarrow{\left(\frac{Row2}{2}\right)} \begin{bmatrix} [1] & 0 & 1 \\ 0 & [1] & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \text{ of the form } \begin{bmatrix} I & F \\ 0 & 0 \end{bmatrix} \quad (3.30)$$

The null basis N is

$$N = \begin{bmatrix} -F \\ I \end{bmatrix} = \begin{bmatrix} -F \\ I \end{bmatrix} = \begin{bmatrix} -1 \\ -1 \\ 1 \end{bmatrix} \quad (3.31)$$

3.6.3 Solving $Ax = \mathbf{b}$

In this sub-section we will illustrate how to completely solve the system $Ax = \mathbf{b}$, if there is a solution. If there is no solution, we will need to identify this fact and if there is some solution, we need to identify how many solutions it has. Our running example will be a system of equations with the same coefficient matrix A that was considered in the previous section (3.6).

$$A\mathbf{x} = \begin{bmatrix} 1 & 2 & 2 & 2 \\ 2 & 4 & 6 & 8 \\ 3 & 6 & 8 & 10 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} \quad (3.32)$$

The third row is the sum of rows 1 and 2. In other words, if we add the left hand sides, we get the third left hand sides. So we can predict right away what elimination will discover about the right hand side. What is the condition that b_1 , b_2 and b_3 satisfy so that there is a solution? Since the sum of the first two left hand sides equals the third left hand side, we require that $b_1 + b_2 = b_3$.

Let us see how elimination discovers this fact. If some combination on the left hand side gives zeros, the same combination on the right hand side should give zeros. Tacking on the vector of \mathbf{b} 's as another column to the matrix A , we get the augmented matrix $[A \ \mathbf{b}]$. Applying the elimination steps $E_{2,1}$ and $E_{3,1}$ to the augmented matrix, followed by the elimination step $E_{3,2}$, we get:

$$[A \ \mathbf{b}] = \begin{bmatrix} 1 & 2 & 2 & 2 & b_1 \\ 2 & 4 & 6 & 8 & b_2 \\ 3 & 6 & 8 & 10 & b_3 \end{bmatrix} \xrightarrow{E_{2,1}, E_{3,1}} \begin{bmatrix} [1] & 2 & 2 & 2 & b_1 \\ 0 & 0 & [2] & 4 & b_2 - 2b_1 \\ 0 & 0 & 2 & 4 & b_3 - 3b_1 \end{bmatrix} \\ \xrightarrow{E_{3,2}} \begin{bmatrix} [1] & 2 & 2 & 2 & b_1 \\ 0 & 0 & [2] & 4 & b_2 - 2b_1 \\ 0 & 0 & 0 & 0 & b_3 - b_1 - b_2 \end{bmatrix} \quad (3.33)$$

The condition for solvability is therefore $b_3 - b_1 - b_2 = 0$. Thus, the system of equations will have a solution for $\mathbf{b} = [5 \ 1 \ 6]^T$.

We will now discuss the solvability conditions on the right hand side of a system of equations to ensure that the system of equations $A\mathbf{x} = \mathbf{b}$ is solvable. We will provide a definition in terms of the column space.

The system of equations $A\mathbf{x} = \mathbf{b}$ is solvable when \mathbf{b} is in the column space $C(A)$.

Another way of describing solvability is:

The system of equations $A\mathbf{x} = \mathbf{b}$ is solvable if a combination of the rows of A produces a zero row, the requirement on \mathbf{b} is that the same combination of the components of \mathbf{b} has to yield zero.

It is not immediately apparent that the two systems of equations are equivalent. We will come back to discuss this in a later part of the course. We will proceed to the case when the system of equations does have a solution.

Assuming that the systems of equations $A\mathbf{x} = \mathbf{b}$ is solvable, what is the algorithm (or sequence of steps) to find the complete solution? We will start by finding one particular solution.

1. $\mathbf{x}_{particular}$ ²: Set all free variables (corresponding to columns with no pivots) to 0. In the example above, we should set $x_2 = 0$ and $x_4 = 0$.
2. Solve $A\mathbf{x} = \mathbf{b}$ for pivot variables.

This leaves us with the equations

$$x_1 + 2x_3 = b_1 \quad 2x_3 = b_2 - 2b_1$$

Adopting the normal back substitution method, we get

$$x_3 = \frac{b_2 - 2b_1}{2}x_1 = b_2 + 3b_1 \quad (3.34)$$

Thus the particular solution is

$$\mathbf{x}_{particular} = \begin{bmatrix} b_2 + 3b_1 \\ 0 \\ \frac{b_2 - 2b_1}{2} \\ 0 \end{bmatrix}$$

For example, if we choose $\mathbf{b} = [5 \ 1 \ 6]^T$, we get

$$\mathbf{x}_{particular} = \begin{bmatrix} -2 \\ 0 \\ \frac{3}{2} \\ 0 \end{bmatrix}$$

The sequence of steps is (a) check for solvability conditions (b) substitute some values for the free variables and obtain values for pivot variables. How do we find the complete solution to $A\mathbf{x} = \mathbf{b}$? It is easy to see that any vector $\mathbf{x}_{nullspace}$ in the null space of the matrix A can be added to $\mathbf{x}_{particular}$ and the resultant vector will still remain a solution. Thus, a general solution to the system $A\mathbf{x} = \mathbf{b}$ is

$$\mathbf{x}_{complete} = \mathbf{x}_{particular} + \mathbf{x}_{nullspace} \quad (3.35)$$

Let us write out the complete solution to this example (recall the null space for this matrix from Equation 3.25).

²Since there are many solutions, one could have one's own way of finding one solution.

$$x_{complete} = \begin{bmatrix} -2 \\ 0 \\ \frac{3}{2} \\ 0 \end{bmatrix} + \begin{bmatrix} -2 & 2 \\ 1 & 0 \\ 0 & -2 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} \quad (3.36)$$

This pattern shows up in all of mathematics, whenever we have linear equations. The reason for this is that

$$A\mathbf{x}_{complete} = A(\mathbf{x}_{particular} + \mathbf{x}_{nullspace}) = \mathbf{b} + \mathbf{0} = \mathbf{b}$$

In words, this means that if we have one solution, we can add on anything in the null space. This gives us the ‘complete’ solution. Note that while the null vector can be scaled arbitrarily, the same does not hold for the *particular* solution.

Let us say we want to plot all solutions to this equation. The plot should be in \mathfrak{R}^4 because there are 4 unknowns. Does the set of solutions to $A\mathbf{x} = \mathbf{b}$ form a subspace? No, because the space of solutions to this system is not closed under the scaling operation. The null space is a 2-dimensional³ subspace inside \mathfrak{R}^4 . The set of solutions to $Ax = b$ does not however pass through the origin, because it must pass through $\mathbf{x}_{particular}$ and then onward. It is like a sub-space shifted away from the origin!

In summary, the algorithm is to go through elimination, find a particular solution and then a special solution. We will now visualize the bigger picture by answering some questions. Consider an $m \times n$ matrix A of rank r .

Q1. What is the relationship between m and r ? We know certainly that $r \leq m$ and $r \leq n$. This because, each row as well as column can contain only one pivot and therefore the number of pivots should be less than the number of rows as also less than the number of columns.

Q2. What happens when the rank r is as big as it can be? There are two possibilities here, depending on what the numbers m and n are.

Q3. In the case that A is full column rank, i.e., $r = n$, what can we infer about the null space and the complete solution? Full column rank implies that there is a pivot in every column, that is, there are n pivots. As a result, there are no free variables. The implication is that the null space will only have the $\mathbf{0}$ vector. Therefore, the complete solution is just $\mathbf{x}_{particular}$; there is just one solution, if there is one at all. Thus, the number of solutions is either 0 or 1. There are many applications in reality where the columns are independent and have nothing to look for in the null space, leading to just a particular solution.

We will illustrate by squeezing in an example.

³The dimension of the subspace corresponds to the number constants you can choose.

$$A = \begin{bmatrix} 1 & 3 \\ 2 & 1 \\ 6 & 1 \\ 5 & 1 \end{bmatrix} \quad (3.37)$$

The rank of this matrix is 2; elimination will yield exactly 2 pivots. Carrying out the elimination process to the end, we can get the following reduced row echelon form for this matrix:

$$A_{rref} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & -17 \\ 0 & -14 \end{bmatrix} \quad (3.38)$$

The first two rows are not independent, but the other rows are combinations of the first two rows. It is a case of full column rank. $A\mathbf{x} = \mathbf{b}$ is a system of four equations in two unknowns. If the right hand side is not consistent with the 4 equations, we will get zero solutions. The right hand side $b = [4 \ 3 \ 7 \ 6]^T$ is consistent with the equations and yields one solution. Similarly, the right hand side b which is the sum of the two independent columns of A also gives one unique solution $\mathbf{x} = [1 \ 1]^T$. We will maintain the natural symmetry of this discussion by next looking at *full row rank*.

Q4. In the case that A is full row rank, i.e., $r = m$, what can we infer about the null space and the complete solution? Elimination will lead to m pivots; every row will have a pivot. What is the implication on solvability, i.e., for which right hand sides will we have a solution to $A\mathbf{x} = \mathbf{b}$? Since we do not have any zero rows, we can solve the system for every right hand side \mathbf{b} . This resolves the question about the existence of a solution. How many free variables will we have? Since $n \geq r$, we will be left with $n - r = n - m$ free variables.

An easy way to obtain an example here (matrix B) is to transpose the above full column rank example matrix A .

$$B = A^T = \begin{bmatrix} 1 & 2 & 6 & 5 \\ 3 & 1 & 1 & 1 \end{bmatrix} \quad (3.39)$$

Elimination yields the following row reduced echelon form with two pivots:

$$\begin{bmatrix} 1 & 0 & 4 & 3 \\ 0 & 1 & -11 & -8 \end{bmatrix} \quad (3.40)$$

$r=m=n$	$r=m<n$	$r=n<m$
$R=I$	$R=[I \ F]$	$R=[I \ 0]^T$
Unique solution	Infinitely many solutions	0 or 1 solution

Figure 3.3: Summary of the properties of the solutions to the system of equations $Ax = b$.

The number of free variables is 2.

Q5. In the case that A is full rank, i.e., $r = m = n$, what can we infer about the null space and the complete solution?

This is the most important case of all. We will illustrate with an example.

$$A = \begin{bmatrix} 1 & 2 \\ 3 & 1 \end{bmatrix} \quad (3.41)$$

The reduced row echelon form for this matrix is

$$A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad (3.42)$$

The matrix is invertible; invertible matrices come out naturally in the rref which is the identity matrix. Which are the satisfiable right hand sides \mathbf{b} for the system $A\mathbf{x} = \mathbf{b}$? Since there are no zero rows, there are no constraints on the right hand side. What is the null space of A ? It is the zero vector only. Since the rank is also m , the only solution is the *particular solution*, and is therefore a unique solution.

Figure 3.3 summarizes the properties of the solutions to the system of equations $Ax = b$ for different inequality constraints between m , n and r . The rank summarizes the possible solutions, except the exact entries in the solutions.

3.7 Independence, Basis, and Dimension

In this section, we will develop the ideas of linear independence of vectors, the space vectors span, basis for vector spaces and finally the dimension of vector spaces. We will assign clear meanings to these terms.

To set the appropriate background, we will begin with a highly important fact which was mentioned earlier. Let us say we have the system $A\mathbf{x} = \mathbf{0}$, where A is an $m \times n$ matrix and $m < n$. That is, we have more unknowns than equations. The conclusion is that there are some non-zero vectors in the null space of A . If we perform elimination on A , we will get some pivots and some free columns that do not have pivots because there will be $n - m$ free variables. We can assign non-zero values to the free variables and automatically obtain values for the pivot variables. We will resume from this point.

3.7.1 Independence

Independence Vectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ are independent if no linear combination gives the zero vector, except the zero combination. That is, $\forall c_1, c_2, \dots, c_n \in \mathcal{R}$,

$$\text{such that not all of the } c_i \text{'s are simultaneously } 0, \sum_i^n c_i \mathbf{x}_i \neq \mathbf{0}.$$

For example, in a two dimensional space, a vector \mathbf{x} and twice the vector $2\mathbf{x}$ are dependent, because $(-2) \times \mathbf{x} + (1) \times 2\mathbf{x} = \mathbf{0}$. As another example, suppose we have the vectors \mathbf{v}_1 and a zero vector \mathbf{v}_2 , they are dependent because $(0) \times \mathbf{v}_1 + (100) \times \mathbf{v}_2 = \mathbf{0}$.

On the other hand, two non-zero vectors \mathbf{v}_1 and \mathbf{v}_2 in a two dimensional space that make an angle $0 < \theta < \frac{\pi}{2}$ with each other will be independent. If we however add a third vector \mathbf{v}_3 from the two dimensional space to the set, the three vectors will now be dependent. How do we determine the truth of the above two statements? We could do this as follows. We construct a matrix A with the vectors as three columns $A = [v_1 \ v_2 \ v_3]$. This matrix is a 2×3 matrix. Does there exist a non-zero solution to the following system?

$$A\mathbf{c} = \begin{bmatrix} \mathbf{v}_1 & \mathbf{v}_2 & \mathbf{v}_3 \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad (3.43)$$

It can be easily proved that a non-zero vector $[c_1 \ c_2 \ c_3]^T$ exists. We will restate the definition for independence in terms of the columns $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ of a matrix A .

Independence The columns $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ of a matrix A are independent if the null-space of A is the zero vector. The columns of A are dependent only if $A\mathbf{c} = \mathbf{0}$ for some $\mathbf{c} \neq \mathbf{0}$.

In other words, the rank of the matrix A , whose columns are independent is the number of columns n . And in the reduced echelon form, all columns will be pivot columns with no free variables. If the columns are dependent, the rank of A will be less than n , and there will be free variables.

What does it mean for a bunch of vectors to span a space? Recall that we could take all combinations of the columns of a matrix to yield the column space. This column space is the space spanned by the columns of the matrix.

Space spanned by vectors: *Vectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ span a space means that the space consists of all linear combinations of the vectors. Thus, the space spanned by the columns $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ of a matrix A , is the column space of A .*

3.7.2 Basis and Dimension

The vectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$, need not be independent in order to span a space. We are specifically interested in a set of vectors that span a space and are at the same time linearly independent. This set of vectors is in some sense, the right number of vectors; even without a single vector from this set, the space cannot be defined. This set is called the *basis*.

Basis for a space: *The basis for a space is a set of vectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ with two properties, viz., (1) The vectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ are independent and (2) These vectors span the space.*

The definition of basis is hinged on the preceding two definitions - the basis is the set of vectors that is necessary and sufficient for spanning the space. As an example, one basis for the four dimensional space \mathfrak{R}^4 is:

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \quad (3.44)$$

It is easy to verify that the above vectors are independent; if a combination of the vectors using the scalars in $[c_1, c_2, c_3, c_4]$ should yield the zero vector, we must have $c_1 = c_2 = c_3 = c_4 = 0$. Another way of proving this is by making the four vectors the columns of a matrix. The resultant matrix will be an identity matrix. The null space of an identity matrix is the zero vector. The above basis is often called the *standard basis* for \mathfrak{R}^4 .

This is not the only basis of \mathfrak{R}^4 . Consider the following three vectors

$$\begin{bmatrix} 2 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 2 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 2 \end{bmatrix} \quad (3.45)$$

These vectors are certainly independent. But they do not span \mathfrak{R}^4 . This can be proved by showing that the following vector in \mathfrak{R}^4 cannot be expressed as a linear combination of these vectors.

$$\begin{bmatrix} 0 \\ 2 \\ 0 \\ 0 \end{bmatrix} \quad (3.46)$$

In fact, if this vector is added to the set of three vectors in (3.45), together, they define another basis for \mathfrak{R}^4 . And this could be proved by introducing them as columns of a matrix A , subject A to row reduction and check if there are any free variables (or equivalently, whether all columns are pivot columns). If there are no free variables, we can conclude that the vectors form a basis for \mathfrak{R}^4 . This is also equivalent to the statement that *if the matrix A is invertible, its columns form a basis for its column space*. This statement can be generalized to \mathfrak{R}^n : *if an $n \times n$ matrix A is invertible, its columns form a basis for \mathfrak{R}^n* .

While there can be many bases for a space, a commonality between all the bases is that they have exactly the same number of vectors. This unique size of the basis is called the dimension of the space.

Dimension: *The number of vectors in any basis of a vector space is called the dimension of the space.*

Do the vectors in (3.45), form a basis for any space at all? The vectors are independent and therefore span the space of all linear combinations of the three vectors. The space spanned by these vectors is a hyperplane in \mathfrak{R}^4 . Let A be any matrix. By definition, the columns of A span the column space $C(A)$ of A . If there exists a $\mathbf{c} \neq \mathbf{0}$ such that, $A\mathbf{c} = \mathbf{0}$, then the columns of A are not linearly independent. For example, the columns of the matrix A given below are not linearly independent.

$$A = \begin{bmatrix} 1 & 2 & 3 & 1 \\ 2 & 3 & 5 & 2 \\ 3 & 4 & 7 & 3 \end{bmatrix} \quad (3.47)$$

A choice of $\mathbf{c} = [-1 \ 0 \ 0 \ 1]^T$ gives $A\mathbf{c} = \mathbf{0}$. Thus, the columns of A do not form a basis for its columns space. What is a basis for $C(A)$? A most natural choice is the first two columns of A ; the third column is the sum of the first and second columns, while the fourth column is the same as the first column. Also, column elimination⁴ on A yields pivots on the first two columns. Thus, a basis for $C(A)$ is

⁴Column elimination operations are very similar to row elimination operations.

$$\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \begin{bmatrix} 2 \\ 3 \\ 4 \end{bmatrix} \quad (3.48)$$

Another basis for $C(A)$ consists of the first and third columns. We note that the dimension of $C(A)$ is 2. We also note that the rank of A is the number of its pivots columns, which is exactly the dimension of $C(A)$. This gives us a nice theorem.

Theorem 32 *The rank of a matrix is the same as the dimension of its column space. That is, $\text{rank}(A) = \text{dimension}(C(A))$.*

What about the dimension of the null space? We already saw that $\mathbf{c} = [-1 \ 0 \ 0 \ 1]^T$ is in the null space. Another element of the null space is $\mathbf{c}' = [1 \ 1 \ -1 \ 0]^T$. These vectors in the null space specify combinations of the columns that yield zeroes. The two vectors \mathbf{c} and \mathbf{c}' are obviously independent. Do these two vectors span the entire null space? The dimension of the null space is the same as the number of free variables, which happens to be $4 - 2 = 2$ in this example. Thus the two vectors \mathbf{c} and \mathbf{c}' must indeed span the null space. In fact, it can be proved that the dimension of the null space of an $m \times n$ matrix A is $n - \text{rank}(A)$.

The space spanned by the rows of a matrix is called the *row space*. We can also define the row space of a matrix A as the column space of its transpose A^T . Thus the row space of A can be specified as $C(A^T)$. The null space of A , $N(A)$ is often called the *right null space* of A , while the null space of A^T , $N(A^T)$ is often referred to as its *left null space*. How do we visualize these four spaces? $N(A)$ and $C(A^T)$ of an $m \times n$ matrix A are in \mathfrak{R}^n , while $C(A)$ and $N(A^T)$ are in \mathfrak{R}^m . How can we construct bases for each of the four subspaces? We note that dimensions of $C(A)$ and the rank of $C(A^T)$ should be the same, since row rank of a matrix is its column rank. The bases of $C(A)$ can be obtained as the set of the pivot columns.

Let r be the rank of A . Recall that the null space is constructed by linear combinations of the special solutions of the null space (3.5.2) and there is one special solution for each assignment of the free variables. In fact, the number of special solutions exactly equals the number of free variables, which is $n - r$. Thus, the dimension of $N(A)$ will be $n - r$. Similarly, the dimension of $N(A^T)$ will be $m - r$.

Let us illustrate this on the sample matrix in (3.47).

$$\begin{bmatrix} 1 & 2 & 3 & 1 \\ 2 & 3 & 5 & 2 \\ 3 & 4 & 7 & 3 \end{bmatrix} \xrightarrow{E_{2,1}, E_{3,1}} \begin{bmatrix} 1 & 2 & 3 & 1 \\ 0 & -1 & -1 & 0 \\ 0 & -2 & -2 & 0 \end{bmatrix} \xrightarrow{E_{3,2}} (R =) \begin{bmatrix} 1 & 2 & 3 & 1 \\ 0 & -1 & -1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \quad (3.49)$$

The reduced matrix R has the same row space as A , by virtue of the nature of row reduction. In fact, the rows of A can be retrieved from the rows of R by reversing the linear operations involved in row elimination. The first two rows give a basis for the row space of A . The dimension of $C(A^T)$ is 2, which is also the rank of A . To find the left null space of A , we look at the system $\mathbf{y}^T A = 0$. Recall the Gauss-Jordan elimination method from Section 3.4.2 that augments A with an $m \times m$ identity matrix, and performs row elimination on the augmented matrix.

$$[A \ I_{m \times m}] \xrightarrow{\text{rref}} [R \ E_{m \times m}]$$

The rref will consist of the reduced matrix augmented with the elimination matrix reproduced on its right. For the example case in 3.49, we apply the same elimination steps to obtain the matrix E below:

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \xrightarrow{E_{2,1}, E_{3,1}} \begin{bmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ -3 & 0 & 1 \end{bmatrix} \xrightarrow{E_{3,2}} (E =) \begin{bmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ 1 & -2 & 1 \end{bmatrix} \quad (3.50)$$

Writing down $EA = R$,

$$\begin{bmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ 1 & -2 & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 & 1 \\ 2 & 3 & 5 & 2 \\ 3 & 4 & 7 & 3 \end{bmatrix} = \begin{bmatrix} 1 & 2 & 3 & 1 \\ 0 & -1 & -1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \quad (3.51)$$

We observe that the last row of E specifies a linear combination of the rows of A that yields a zero vector (corresponding to the last row of R). This is the only vector that yields a zero row in R and is therefore the only element in the basis of the left null space of A , that is, $N(A^T)$. The dimension of $N(A^T)$ is 1.

As another example, consider the space \mathcal{S} of vectors $\mathbf{v} \in \mathbb{R}^3$ where $\mathbf{v} = [v_1 \ v_2 \ v_3]^T$ such that $\mathbf{v}_1 + \mathbf{v}_2 + \mathbf{v}_3 = \mathbf{0}$. What is the dimension of this subspace? Note that this subspace is the right null space $N(A)$ of a 1×3 matrix $A = [1 \ 1 \ 1]$, since $A\mathbf{v} = 0$. The rank, $r = \text{rank}(A)$ is 1, implying that the dimension of the right null space is $n - r = 3 - 1 = 2$. One set of basis vectors for \mathcal{S} is $[-1 \ 1 \ 0]$, $[-1 \ 0 \ 1]$. The column space $C(A)$ is \mathbb{R}^1 with dimension 1. The left null space $N(A^T)$ is the singleton set $\{0\}$ and as expected, has a dimension of $m - r = 1 - 1 = 0$.

3.8 Matrix Spaces

We will extend the set of examples of vector spaces discussed in Section 3.5 with a new vector space, that of all $m \times n$ matrices with real entries, denoted by $\mathbb{R}^{m \times n}$.

It is easy to verify that the space of all matrices is closed under operations of addition and scalar multiplication. Additionally, there are interesting subspaces in the entire matrix space $\mathfrak{R}^{m \times n}$, viz.,

- set \mathcal{S} of all $n \times n$ symmetric matrices
- set \mathcal{U} of all $n \times n$ upper triangular matrices
- set \mathcal{L} of all $n \times n$ lower triangular matrices
- set \mathcal{D} of all $n \times n$ diagonal matrices

Let $\mathcal{M} = \mathfrak{R}^{3 \times 3}$ be the space of all 3×3 matrices. The dimension of \mathcal{M} is 9. Each element of this basis has a 1 in one of the 9 positions and the remaining entries as zeroes. Of these basis elements, three are symmetric (those having a 1 in any of the diagonal positions). These three matrices form the basis for the subspace of diagonal matrices. Six of the nine basis elements of \mathcal{M} form the basis of \mathcal{U} while six of them form the basis of \mathcal{L} .

The intersection of any two matrix spaces is also a matrix space. For example, $\mathcal{S} \cap \mathcal{U}$ is \mathcal{D} , the set of diagonal matrices. However the union of any two matrix spaces need not be a matrix space. For example, $\mathcal{S} \cup \mathcal{U}$ is not a matrix space; the sum $S + U$, $S \in \mathcal{S}$, $U \in \mathcal{U}$ need not belong to $\mathcal{S} \cup \mathcal{U}$. We will discuss a special set comprising all linear combinations of the elements of union of two vector spaces \mathcal{V}_1 and \mathcal{V}_2 (i.e., $\mathcal{V}_1 \cup \mathcal{V}_2$), and denote this set by $\mathcal{V}_1 \oplus \mathcal{V}_2$. By definition, this set is a vector space. For example, $\mathcal{S} + \mathcal{U} = \mathcal{M}$, which is a vector space.

A property fundamental to many properties of matrices is the expression for a rank 1 matrix. A rank 1 matrix can be expressed as the product of a column vector with a row vector (the row vector forming a basis for the matrix). Thus, any rank 1 matrix X can be expressed as

$$X_{m \times n} = u^T v = \begin{bmatrix} u_1 \\ u_2 \\ u_3 \\ \vdots \\ u_m \end{bmatrix} \begin{bmatrix} v_1 & v_2 & \dots & v_n \end{bmatrix} \quad (3.52)$$

Let $\mathcal{M}_{m \times n}$ be the set of all $m \times n$ matrices. Is the subset of $\mathcal{M}_{m \times n}$ matrices with rank k , a subspace? For $k = 1$, this space is obviously not a vector space as is evident from the sum of rank 1 matrices, A^1 and B^1 , which is not a rank 1 matrix. In fact, the subset of $\mathcal{M}_{m \times n}$ matrices with rank k is not a subspace.

$$A^1 + B^1 = \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 1 \\ 1 & 2 & 1 \end{bmatrix} + \begin{bmatrix} 4 & 4 & 2 \\ 2 & 2 & 1 \\ 4 & 4 & 2 \end{bmatrix} = \begin{bmatrix} 5 & 6 & 3 \\ 4 & 6 & 2 \\ 5 & 6 & 3 \end{bmatrix} \quad (3.53)$$

3.9 Orthogonality and Projection

In this section we will discuss the orthogonality of subspaces. Two vectors \mathbf{x} and \mathbf{y} are said to be orthogonal *iff*, their dot product is 0. In the euclidian space, the dot product of the two vectors is $\mathbf{x}^T \mathbf{y}$. The condition $\mathbf{x}^T \mathbf{y} = 0$ is equivalent to the pythagorous condition between the vectors \mathbf{x} and \mathbf{y} that form the perpendicular sides of a right triangle with the hypotenuse given by $\mathbf{x} + \mathbf{y}$. The *pythagorous condition* is $\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 = \|\mathbf{x} + \mathbf{y}\|^2$, where the norm is the euclidian norm, given by $\|\mathbf{x}\|^2 = \mathbf{x}^T \mathbf{x}$. This equivalence can be easily proved and is left to the reader as an exercise. By definition, the vector $\mathbf{0}$ is orthogonal to every other vector.

We will extend the definition of orthogonality to subspaces; a subspace \mathcal{U} is orthogonal to subspace \mathcal{V} *iff*, every vector in \mathcal{U} is orthogonal to every vector in \mathcal{V} . As an example:

Theorem 33 *The row space $C(A^T)$ of an $m \times n$ matrix A is orthogonal to its right null space $N(A)$.*

Proof: $A\mathbf{x} = \mathbf{0}$, $\forall \mathbf{x} \in N(A)$. On the other hand, $\forall \mathbf{y} \in C(A^T)$, $\exists \mathbf{z} \in \mathfrak{R}^m$, s.t., $\mathbf{y} = A^T \mathbf{z}$. Therefore, $\forall \mathbf{y} \in C(A^T)$, $\mathbf{x} \in N(A)$, $\mathbf{y}^T \mathbf{x} = \mathbf{z}^T A\mathbf{x} = \mathbf{z}^T \mathbf{0} = 0$. \square

Not only are $C(A^T)$ and the right null space $N(A)$ orthogonal to each other, but they are also *orthogonal complements* in \mathfrak{R}^n , that is, $N(A)$ contains all vectors that are orthogonal to some vector in $C(A^T)$.

Theorem 34 *The null space of A and its row space are orthogonal complements.*

Proof: We note, based on our discussion in Section 3.7.2 that the dimensions of the row space and the (right) null space add up to n , which is the number of columns of A . For any vector $\mathbf{y} \in C(A^T)$, we have $\exists \mathbf{z} \in \mathfrak{R}^m$, s.t., $\mathbf{y} = A^T \mathbf{z}$. Suppose $\forall \mathbf{y} \in C(A^T)$, $\mathbf{y}^T \mathbf{x} = 0$. That is, $\forall \mathbf{z} \in \mathfrak{R}^m$, $\mathbf{z}^T A\mathbf{x} = 0$. This is possible only if $A\mathbf{x} = \mathbf{0}$. Thus, necessarily, $\mathbf{x} \in N(A)$. \square

Along similar lines, we could prove that the column space $C(A)$ and the left null space $N(A^T)$ are orthogonal complements in \mathfrak{R}^m . Based on theorem 34, we prove that there is a one-to-one mapping between the elements of row space and column space.

Theorem 35 *If $\mathbf{x} \in C(A^T)$, $\mathbf{y} \in C(A^T)$ and $\mathbf{x} \neq \mathbf{y}$, then, $A\mathbf{x} \neq A\mathbf{y}$.*

Proof: Note that $A\mathbf{x}$ and $A\mathbf{y}$ are both elements of $C(A)$. Next, observe that $\mathbf{x} - \mathbf{y} \in C(A^T)$, which by theorem 34, implies that $\mathbf{x} - \mathbf{y} \notin N(A)$. Therefore, $A\mathbf{x} - A\mathbf{y} \neq \mathbf{0}$ or in other words, $A\mathbf{x} \neq A\mathbf{y}$. \square

Similarly, it can be proved that if $\mathbf{x} \in C(A)$, $\mathbf{y} \in C(A)$ and $\mathbf{x} \neq \mathbf{y}$, then, $A^T \mathbf{x} \neq A^T \mathbf{y}$. The two properties together imply a one-to-one mapping between the row and column spaces.

3.9.1 Projection Matrices

The projection of a vector \mathbf{t} on a vector \mathbf{s} is a vector $\mathbf{p} = c\mathbf{s}$, $c \in \mathfrak{R}$ (in the same direction as \mathbf{s}), such that $\mathbf{t} - c\mathbf{s}$ is orthogonal to \mathbf{s} . That is, $\mathbf{s}^T(\mathbf{t} - c\mathbf{s}) = 0$ or $\mathbf{s}^T\mathbf{t} = c\mathbf{s}^T\mathbf{s}$. Thus, the scaling factor c is given by $c = \frac{\mathbf{s}^T\mathbf{t}}{\mathbf{s}^T\mathbf{s}}$. The projection of the vector \mathbf{t} on a vector \mathbf{s} is then

$$\mathbf{p} = \mathbf{s} \frac{\mathbf{t}^T \mathbf{s}}{\mathbf{s}^T \mathbf{s}} \quad (3.54)$$

Using the associative property of matrix multiplication, the expression for \mathbf{p} can be re-written as

$$\mathbf{p} = P\mathbf{t} \quad (3.55)$$

where, $P = \mathbf{s}\mathbf{s}^T \frac{1}{\mathbf{s}^T\mathbf{s}}$ is called the *projection matrix*.

The rank of the projection matrix is 1 (since it is a column multiplied by a row). The projection matrix is symmetric and its column space is a line through \mathbf{s} . For any $d \in \mathfrak{R}$, $P(d\mathbf{s}) = d\mathbf{s}$, that is, the projection of any vector in the direction of \mathbf{s} is the same vector. Thus, $P^2 = P$.

3.9.2 Least Squares

In Section 3.6.3, we saw a method for solving the system $A\mathbf{x} = \mathbf{b}$ (A being an $m \times n$ matrix), when a solution exists. However, a solution may not exist, especially when $m > n$, that is when the number of equations is greater than the number of variables. In Section 3.6.3, we saw that the *rref* looks like $[I \ \mathbf{0}]^T$, where I is an $n \times n$ identity matrix. It could happen that the row reduction yields a zero submatrix in the lower part of A , but the corresponding elements in \mathbf{b} are not zeroes. In other words, \mathbf{b} may not be in the column space of A . In such cases, we are often interested in finding a ‘best fit’ for the system; a solution $\hat{\mathbf{x}}$ that satisfies $A\mathbf{x} = \mathbf{b}$ as well as possible.

We define the best fit in terms of a vector \mathbf{p} which is the projection of \mathbf{b} onto $C(A)$ and solve $A\hat{\mathbf{x}} = \mathbf{p}$. We require that $\mathbf{b} - \mathbf{p}$ is orthogonal to $C(A)$, which means

$$A^T(\mathbf{b} - A\hat{\mathbf{x}}) = \mathbf{0} \quad (3.56)$$

The vector $\mathbf{e} = \mathbf{b} - A\hat{\mathbf{x}}$ is the error vector and is in $N(A^T)$. The equation (3.9.2) can be rewritten as

$$(A^T A)\hat{\mathbf{x}} = A^T \mathbf{b} \quad (3.57)$$

A matrix that plays a key role in this problem is $A^T A$. It is an $n \times n$ symmetric matrix (since $(A^T A)^T = A^T A$). The right null space $N(A^T A)$ is the same as $N(A)$ ⁵. It naturally follows that the ranks of $A^T A$ and A are the same (since, the sum of the rank and dimension of null space equal n in either case). Thus, $A^T A$ is invertible exactly if $N(A)$ has dimension 0, or equivalently, A is a full column rank.

Theorem 36 *If A is a full column rank matrix (that is, its columns are independent), $A^T A$ is invertible.*

Proof: We will show that the null space of $A^T A$ is $\{0\}$, which implies that the square matrix $A^T A$ is full column (as well as row) rank is invertible. That is, if $A^T A\mathbf{x} = \mathbf{0}$, then $\mathbf{x} = \mathbf{0}$. Note that if $A^T A\mathbf{x} = \mathbf{0}$, then $\mathbf{x}^T A^T A\mathbf{x} = \|\mathbf{Ax}\|^2 = 0$ which implies that $\mathbf{Ax} = \mathbf{0}$. Since the columns of A are linearly independent, its null space is $\mathbf{0}$ and therefore, $\mathbf{x} = \mathbf{0}$. \square

Assuming that A is full column rank, the equation (3.9.2) can be rewritten as

$$\hat{\mathbf{x}} = (A^T A)^{-1} A^T \mathbf{b}. \quad (3.58)$$

Therefore the expression for the projection \mathbf{p} will be

$$\mathbf{p} = A(A^T A)^{-1} A^T \mathbf{b} \quad (3.59)$$

This expression is the n -dimensional equivalent of the one dimensional expression for projection in (3.9.1). The projection matrix in (3.59) is given by $P = A(A^T A)^{-1} A^T$. We will list the solution for some special cases:

- If A is an $n \times n$ square invertible matrix, its column space is the entire \mathcal{R}^n and the projection matrix will turn out to be the identity matrix.
- Also, if b is in the column space $C(A)$, then $\mathbf{b} = A\mathbf{t}$ for some t in \mathcal{R}^n and consequently, $P\mathbf{b} = A(A^T A)^{-1} (A^T A)\mathbf{t} = A\mathbf{t} = \mathbf{b}$.

⁵The proof is left as an exercise.

- On the other hand, if \mathbf{b} is orthogonal to $C(A)$, it will lie in $N(A^T)$, and therefore, $A^T \mathbf{b} = \mathbf{0}$, implying that $\mathbf{p} = \mathbf{0}$.

Another equivalent way of looking at the best fit solution $\hat{\mathbf{x}}$ is a solution that minimizes the square of the norm of the error vector

$$e(\hat{\mathbf{x}}) = \|A\hat{\mathbf{x}} - \mathbf{b}\|^2 \quad (3.60)$$

Setting $\frac{de(\hat{\mathbf{x}})}{d\hat{\mathbf{x}}} = \mathbf{0}$, we get the same expression for $\hat{\mathbf{x}}$ as in (3.9.2). The solution in 3.9.2 is therefore often called the *least squares solution*. Thus, we saw two views of finding a best fit; first was the view of projecting into the column space while the second concerned itself with minimizing the norm squared of the error vector.

We will take an example. Consider the data matrix A and the coefficient matrix \mathbf{b} as in (3.61).

$$Ax = \begin{bmatrix} 2 & -1 \\ -1 & 2 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} \hat{x}_1 \\ \hat{x}_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 3 \\ 3 \end{bmatrix} \quad (3.61)$$

The matrix A is full column rank and therefore $A^T A$ will be invertible. The matrix $A^T A$ is given as

$$A^T A = \begin{bmatrix} 6 & -3 \\ -3 & 6 \end{bmatrix}$$

Substituting the value of $A^T A$ in the system of equations (3.9.2), we get,

$$6\hat{x}_1 - 3\hat{x}_2 = 2 - 3\hat{x}_1 + 6\hat{x}_2 = 8$$

The solution of which is, $x_1 = \frac{4}{5}$, $x_2 = \frac{26}{15}$.

3.9.3 Orthonormal Vectors

A collection of vectors $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_n$ is said to be orthonormal *iff* the following condition holds $\forall i, j$:

$$\mathbf{q}_i^T \mathbf{q}_j \begin{cases} 0 & \text{if } i \neq j \\ 1 & \text{if } i = j \end{cases} \quad (3.62)$$

A large part of numerical linear algebra is built around working with orthonormal matrices, since they do not overflow or underflow. Let Q be a matrix comprising the columns \mathbf{q}_1 through \mathbf{q}_n . It can be easily shown that

$$Q^T Q = I_{n \times n}$$

When Q is square, $Q^{-1} = Q^T$. Some examples of matrices with orthonormal columns are:

$$Q_{rotation} = \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix}, \quad Q_{reflection} = \begin{bmatrix} \cos(\theta) & \sin(\theta) \\ \sin(\theta) & -\cos(\theta) \end{bmatrix},$$

$$Q_{Hadamard} = \frac{1}{2} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \end{bmatrix}, \quad Q_{rect} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \quad (3.63)$$

The matrix $Q_{rotation}$ when multiplied to a vector, rotates it by an angle θ , whereas $Q_{reflection}$ reflects the vector at an angle of $\theta/2$. These matrices present standard varieties of linear transformation, but in general, premultiplication by an $m \times n$ matrix transforms from an input space in \mathfrak{R}^m to an input space in \mathfrak{R}^n . The matrix $Q_{Hadamard}$ is an orthonormal matrix consisting of only 1's and -1's. Matrices of this form exist only for specific dimensions such as 2, 4, 8, 16, *etc.*, and are called *Hadamard matrices*⁶. The matrix Q_{rect} is an example rectangular matrix whose columns are orthonormal.

Suppose a matrix Q has orthonormal columns. What happens when we project any vector onto the column space of Q ? Substituting $A = Q$ in (3.59), we get⁷:

$$\mathbf{p} = Q(Q^T Q)^{-1} Q^T \mathbf{b} = Q Q^T \mathbf{b} \quad (3.64)$$

Making the same substitution in (3.9.2),

$$\hat{\mathbf{x}} = (A^T Q)^{-1} Q^T \mathbf{b} = Q^T \mathbf{b} \quad (3.65)$$

The i^{th} component of \mathbf{x} , is given by $x_i = q_i^T \mathbf{b}$.

Let Q_1 be one orthonormal basis and Q_2 be another orthonormal basis for the same space. Let A be the coefficient matrix for a set of points represented using Q_1 and B be the coefficient matrix for the same set of points represented using Q_2 . Then $Q_1 A = Q_2 B$, which implies that B can be computed as $B = Q_2^T Q_1 A$. This gives us the formula for changing basis.

⁶An exhaustive listing of different types of matrices can be found at http://en.wikipedia.org/wiki/List_of_matrices.

⁷Note that $Q^T Q = I$. However, $Q Q^T = I$ only if Q is a square matrix.

3.9.4 Gram-Schmidt Orthonormalization

The goal of the Gram-Schmidt orthonormalization process is to generate a set of orthonormal vectors $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_n$, given a set of independent vectors $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n$. The first step in this process is to generate a set of orthogonal vectors $\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_n$ from $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n$. To start with, \mathbf{t}_1 is chosen to be \mathbf{a}_1 . Next, the vector \mathbf{t}_2 is obtained by removing the projection of \mathbf{a}_2 on \mathbf{t}_1 , from \mathbf{a}_2 , based on (3.9.1). That is,

$$\mathbf{t}_2 = \mathbf{a}_2 - \frac{1}{\mathbf{a}_1^T \mathbf{a}_1} \mathbf{a}_1 \mathbf{a}_1^T \mathbf{a}_2 \quad (3.66)$$

This is carried out iteratively for $i = 1, 2, \dots, n$, using the expression below:

$$\mathbf{t}_i = \mathbf{a}_i - \frac{1}{\mathbf{t}_1^T \mathbf{t}_1} \mathbf{t}_1 \mathbf{t}_1^T \mathbf{a}_i - \frac{1}{\mathbf{t}_2^T \mathbf{t}_2} \mathbf{t}_2 \mathbf{t}_2^T \mathbf{a}_i - \dots - \frac{1}{\mathbf{t}_{i-1}^T \mathbf{t}_{i-1}} \mathbf{t}_{i-1} \mathbf{t}_{i-1}^T \mathbf{a}_i \quad (3.67)$$

This gives us the orthogonal vectors $\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_n$. Finally, the orthonormal vectors $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_n$ are obtained by the simple expression

$$\mathbf{q}_i = \frac{1}{\|\mathbf{t}_i\|} \mathbf{t}_i \quad (3.68)$$

Let A be the matrix with columns $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n$ and Q , the matrix with columns $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_n$. It can be proved that $C(V) = C(Q)$, that is, the matrices V and Q have the same column space. The vector \mathbf{a}_i can be expressed as

$$\mathbf{a}_i = \sum_{k=1}^n (\mathbf{a}_i^T \mathbf{q}_k) \mathbf{q}_k \quad (3.69)$$

The i^{th} column of A is a linear combination of the columns of Q , with the scalar coefficient $\mathbf{a}_i^T \mathbf{q}_k$ for the k^{th} column of Q . By the very construction procedure of the Gram-Schmidt orthonormalization process, \mathbf{a}_i is orthogonal to \mathbf{q}_k for all $k > i$. Therefore, (3.69) can be expressed more precisely as

$$\mathbf{a}_i = \sum_{k=1}^i (\mathbf{a}_i^T \mathbf{q}_k) \mathbf{q}_k \quad (3.70)$$

Therefore, matrix A can be decomposed into the product of Q with an upper triangular matrix R ; $A = QR$, with $R_{k,i} = \mathbf{a}_i^T \mathbf{q}_k$. Since $\mathbf{a}_i^T \mathbf{q}_k = 0, \forall k > i$, we can easily see that R is upper triangular.

3.9.5 Fourier and Wavelet Basis

The hermetian inner product of a complex vector \mathbf{x} with another complex vector \mathbf{y} is $\overline{\mathbf{x}}^T \mathbf{y}$, and is also denoted by $\mathbf{x}^H \mathbf{y}$. A complex matrix Q is called orthonormal if $\overline{Q}^T Q = I$. Consider the complex number $c = \cos(\frac{2\pi}{n}) + i\cos(\frac{2\pi}{n})$. Then, $w^n = 1$. The *fourier matrix* F_n is defined as

$$F_n = \frac{1}{n} \begin{bmatrix} 1 & 1 & \dots & 1 \\ 1 & c & \dots & c^{n-1} \\ \cdot & \cdot & \dots & \cdot \\ 1 & c^{k-1} & \dots & c^{(k-1)(n-1)} \\ \cdot & \cdot & \dots & \cdot \\ 1 & c^{n-1} & \dots & c^{(n-1)(n-1)} \end{bmatrix} \quad (3.71)$$

The (hermetian) inner products of distinct columns F_n are 0, while the inner product of a column with itself is 1. Therefore, the columns of F_n are orthonormal and form a basis for \mathfrak{R}^n . Consequently, the inverse of F_n is its conjugate transpose \overline{F}_n^T .

Further, $F_{2^k}, k \geq 1$ can expressed as

$$F_{2^k} == \begin{bmatrix} I & D \\ I & -D \end{bmatrix} \begin{bmatrix} F_{2^{k-1}} & 0_{2^{k-1}} \\ 0_{2^{k-1}} & F_{2^{k-1}} \end{bmatrix} \underbrace{\begin{bmatrix} 0 & 1 & 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & 0 & 1 & \dots & 0 & 0 \\ \cdot & \cdot & \cdot & \cdot & \dots & \cdot & \cdot \\ 0 & 0 & 0 & 0 & \dots & 0 & 1 \\ 1 & 0 & 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & 1 & 0 & \dots & 0 & 0 \\ \cdot & \cdot & \cdot & \cdot & \dots & \cdot & \cdot \\ 0 & 0 & 0 & 0 & \dots & 1 & 0 \end{bmatrix}}_P \quad (3.72)$$

where, $D = \text{diag}(1, c, c^2, \dots, c^{2^k})$ and 0_{2^k} is a $2^k \times 2^k$ matrix of 0's. This factorization, applied recursively, can reduce the time for computing F_{2^k} from $\mathcal{O}(n^2)$ to $\mathcal{O}(n \log n)$. This is the idea behind the *fast fourier transform algorithm*. Though the factorization discussed here, applies to only to powers of 2, there exist FFT algorithms [?] for any number (including primes).

An advantage of representing a vector in \mathfrak{R}^n (for example, a $\sqrt{n} \times \sqrt{n}$ sub-block of an image matrix) using the fourier basis is that certain basis components of the representation could be ignored to achieve minimally lossy compression of matrices such as image matrices. Another orthogonal basis that is used for

minimally lossy matrix compression is the wavelet basis. A sample wavelet basis matrix W for \mathfrak{R}^8 is

$$W = \begin{bmatrix} 1 & 1 & 1 & 0 & -1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & -1 & 0 & 0 & 0 \\ 1 & 1 & -1 & 0 & 0 & -1 & 0 & 0 \\ 1 & 1 & -1 & 0 & 0 & -1 & 0 & 0 \\ 1 & -1 & 0 & 1 & 0 & 0 & -1 & 0 \\ 1 & -1 & 0 & 1 & 0 & 0 & -1 & 0 \\ 1 & -1 & 0 & -1 & 0 & 0 & 0 & -1 \\ 1 & -1 & 0 & -1 & 0 & 0 & 0 & -1 \end{bmatrix} \quad (3.73)$$

The discrete wavelet transform can be computed efficiently using a *fast wavelet transform algorithm* which is less computationally complex, taking $\mathcal{O}(n)$ time as compared to $\mathcal{O}(n \log n)$ for the fast fourier transform.

3.10 Determinants

Every square matrix A has a real number associated with it, called its determinant and it is denoted by $\det(A)$. In this sequel, we will often refer to the following $n \times n$ matrix A :

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \cdot & \cdot & \dots & \cdot \\ a_{k1} & a_{k2} & \dots & a_{kn} \\ \cdot & \cdot & \dots & \cdot \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix} \quad (3.74)$$

We will describe four fundamental properties of the determinant, which essentially define the determinant.

1. The determinant of the identity matrix I is 1. That is, $\det(I) = 1$.
2. When two rows of A are permuted (*c.f.* Section 3.3.1), the sign of the determinant changes. That is $\det(\text{Perm}(A, j, k)) = -\det(A)$, where $\text{Perm}(A, j, k)$ returns a matrix formed by exchanging the j^{th} and k^{th} rows of A for any $1 \leq j, k \leq n$.
3. If any row of A is scaled by $t \in \mathfrak{R}$, the determinant also gets scaled by t . Thus, if

$$S(A, k, t) = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \cdot & \cdot & \dots & \cdot \\ ta_{k1} & ta_{k2} & \dots & ta_{kn} \\ \cdot & \cdot & \dots & \cdot \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix} \quad (3.75)$$

then

$$\det(S(A, k, t)) = t \times \det(A) \quad (3.76)$$

The function $S(A, k, t)$ returns a matrix exactly with all the entries of A , except for the k^{th} row, which is scaled by $t \in \mathfrak{R}$.

4. The sum of the determinants of two $n \times n$ matrices, A and B , with all $(n - 1)$ rows the same, except for the k^{th} row, $1 \leq k \leq n$, equals the determinant of an $n \times n$ matrix C that has the same $n - 1$ rows from A/B , but with the k^{th} row being the sum of the k^{th} rows of A and B . Thus, if

$$B = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \cdot & \cdot & \dots & \cdot \\ b_{k1} & b_{k2} & \dots & b_{kn} \\ \cdot & \cdot & \dots & \cdot \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix} \quad (3.77)$$

and,

$$C = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \cdot & \cdot & \dots & \cdot \\ a_{k1} + b_{k1} & a_{k2} + b_{k2} & \dots & a_{kn} + b_{kn} \\ \cdot & \cdot & \dots & \cdot \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix} \quad (3.78)$$

then

$$\det(C) = \det(A) + \det(B)$$

Using these basic properties of determinants, we infer some **derived properties**:

1. If a matrix A has two equal rows, its determinant must be 0.

Proof: Let B be the matrix obtained by permuting the two equal rows of A . By the second property of determinants, $\det(B) = -\det(A)$. Since the permuted rows are the same, $B = A$, which implies that $\det(B) = \det(A)$. The two equalities on determinants of A and B imply that $\det(A) = 0$. \square

2. The determinant of A obtained by subtracting $\rho \in \mathfrak{R}$ times the j^{th} row from the k^{th} row leaves the determinant unaltered. Therefore, if

$$E = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \cdot & \cdot & \cdots & \cdot \\ a_{k1} - \rho a_{j1} & a_{k2} - \rho a_{j2} & \cdots & a_{kn} - \rho a_{jn} \\ \cdot & \cdot & \cdots & \cdot \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix} \quad (3.79)$$

we will have

$$\det(E) = \det(A)$$

Proof:

$$\begin{aligned} \det(E) &= \det(A) + \det \left(\begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \cdot & \cdot & \cdots & \cdot \\ -\rho a_{j1} & -\rho a_{j2} & \cdots & -\rho a_{jn} \\ \cdot & \cdot & \cdots & \cdot \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix} \right) \\ &= \det(A) - \rho \times \det \left(\begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \cdot & \cdot & \cdots & \cdot \\ a_{j1} & a_{j2} & \cdots & a_{jn} \\ \cdot & \cdot & \cdots & \cdot \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix} \right) = \det(A) \quad (3.80) \end{aligned}$$

The first step follows from the fourth fundamental property of determinants, the second follows from the third fundamental property, while the third step is a consequence of the first derived property of determinants. \square

This property shows that the row elimination steps discussed in Section 3.3.1, leave the determinant unchanged.

3. If any k^{th} row of A is 0, then $\det(A) = 0$ *Proof:* Consider a matrix A' that has the same rows as A for all $1 \leq i \leq n$, except for the $i = k$. Let the k^{th} row of A' be the same as its j^{th} row, for some $1 \leq j \leq n$, such that $j \neq k$. Note that by the first derived property, $\det(A') = 0$. The matrix A can be obtained from A' by subtracting the j^{th} row of A' from its k^{th} row. By the second derived property, $\det(A) = \det(A')$. Thus, $\det(A) = 0$. Another simpler proof is that $S(A, k, 0) = A$ which implies that $\det(A) = \det(S(A, k, 0)) = 0 \times \det(A) = 0$ (by third fundamental property of determinants). \square
4. The determinant of an upper triangular matrix U is the product of its diagonal entries.

Proof: Row elimination operations can be performed on an upper triangular matrix U to yield a diagonal matrix D (c.f. Section 3.4.2 on Gauss-Jordan elimination), while neither performing any row exchanges nor altering the diagonal entries of U . By the second derived property of determinants, $\det(D) = \det(U)$. Using the first fundamental property of determinants, $\det(D)$ can be proved to be the product of its diagonal entries, which is also the product of the diagonal entries of U . \square

In fact, most mathematical softwares compute the determinant of a matrix A by first reducing it to an upper triangular matrix U by row elimination on A (which preserves the determinant, by virtue of the second derived property) and then compute the product of the diagonal entries (which also happen to be the pivots) of U . If some α row exchanges are performed during the reduction of A to U , the product of the diagonal entries, multiplied by $(-1)^\alpha$ yields the determinant of A . As an example, consider the 2×2 matrix A_2 :

$$A_2 = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \quad (3.81)$$

Using the derived property (2), the matrix A'_2 can be proved to have same determinant as A .

$$A'_2 = \begin{bmatrix} a_{11} & a_{12} \\ 0 & a_{22} - \frac{a_{21} * a_{12}}{a_{11}} \end{bmatrix} \quad (3.82)$$

A'_2 is an upper triangular matrix with $\det(A'_2)$ given by

$$\det(A'_2) = a_{11}a_{22} - a_{21} * a_{12} \quad (3.83)$$

Therefore,

$$\det(A) = a_{11}a_{22} - a_{21} * a_{12} \quad (3.84)$$

5. The determinant of a matrix A is 0 iff A is singular (or non-invertible).

Proof Sketch: We will consider the proof in two parts. When A is singular, elimination, with some possible permutations, yields (as discussed in Section 3.4.1) an upper triangular matrix with some diagonal (pivot) entries that are 0s. Therefore, by the derived property (4), $\det(A) = 0$. When A is non-singular, elimination yields an upper triangular matrix with no zero entries. Consequently, we will have $\det(A) \neq 0$. \square

6. The determinant of the product of two matrices A and B is the product of their determinants⁸, i.e., $\det(AB) = \det(A) \times \det(B)$.

A corollary of this property is that $\det(A^{-1}) = \frac{1}{\det(A)}$ because $\det(A^{-1})\det(A) = \det(I) = 1$. Similarly, $\det(A^n) = (\det(A))^n$ and $\det(A_1A_2 \dots A_n) = \det(A_1)\det(A_2) \dots \det(A_n)$.

In this context, it will be appropriate to point out that determinants also have relationship with volumes of solids. The determinant of matrix A in (3.74), is the volume of an n -dimensional parallelotope, with corners at $(0, 0, \dots, 0)$, $(a_{11}, a_{12}, \dots, a_{1n})$, \dots , $(a_{n1}, a_{n2}, \dots, a_{nn})$. The parallelotope corresponding to $I_{n \times n}$ is an n -dimensional unit hypercube in n dimensions and has a volume of 1. An orthonormal matrix Q represents a hypercube in n dimensions and has volume given by $\det(Q) = \sqrt{\det(I)} = 1$.

If Q is orthogonal (and not necessarily orthonormal), its volume is $\prod_{i=1}^n s_i$,

where s_i is the factor by which the i^{th} row of Q should be scaled, so that the row has unit norm. Determinants make easy the task of computing areas of parallelotopes. If the parallelotope does not have any corner at the origins, the coordinates of the corners can be computed relative any one of the corners and the area can be computed using determinants.

7. The determinant of a square matrix A equals the determinant of its transpose, i.e., $\det(A) = \det(A^T)$.

⁸Recall that determinant does not have the linear additive property.

Proof Sketch: We can decompose A as a LU , where L is a lower triangular matrix and U is an upper triangular matrix. That is $A = LU$. Consequently, $A^T = U^T L^T$. By derived property (6), $\det(A) = \det(L)\det(U)$ and $\det(A^T) = \det(U^T)\det(L^T)$. Since the diagonal entries of L^T and U^T are the same as the diagonal entries of L and U respectively, and since derived property (4) states that the determinants of L , L^T , U and U^T are just products of their respective diagonal entries, we have $\det(A) = \det(A^T)$. \square

By virtue of this property, all the properties of determinants discussed so far with respect to scaling or exchanging rows hold for similar manipulations on the columns, since column operations on A are row operations on A^T .

3.10.1 Formula for determinant

In (3.84), we showed the formula for the determinant of a 2×2 matrix A_2 . The formula can also be obtained by using the basic property (4), decomposing $\det(A_2)$ into the sum of determinants of 4 matrices, with one surviving element per row. We will use the notation $|\cdot|$ instead of $\det([\cdot])$ to denote the determinant of a matrix.

$$\begin{aligned} \det(A_2) &= \begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix} = \begin{vmatrix} a_{11} & 0 \\ a_{21} & 0 \end{vmatrix} + \begin{vmatrix} a_{11} & 0 \\ 0 & a_{22} \end{vmatrix} \\ &\quad + \begin{vmatrix} 0 & a_{12} \\ a_{21} & 0 \end{vmatrix} + \begin{vmatrix} 0 & a_{12} \\ 0 & a_{22} \end{vmatrix} = a_{11}a_{22} - a_{12}a_{21} \end{aligned} \quad (3.85)$$

Of these, there are only two nonzero terms; the terms with zero columns or zero rows have 0 determinant. The determinant of a 3×3 matrix can be similarly computed, by decomposing the determinant as the sum of $3 \times 3 \times 3 = 27$ determinants. However, many of the determinants in the sum turn out to be 0, either because of zero rows or zero columns. Each of the non-zero terms have exactly one entry for each row and each column. Thus, the determinant of a 3×3 matrix can be expressed as

$$\begin{aligned}
\det(A_3) &= \begin{vmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{vmatrix} = \begin{vmatrix} a_{11} & 0 & 0 \\ 0 & a_{22} & 0 \\ 0 & 0 & a_{33} \end{vmatrix} + \begin{vmatrix} a_{11} & 0 & 0 \\ 0 & 0 & a_{23} \\ 0 & a_{32} & 0 \end{vmatrix} \\
&\quad + \begin{vmatrix} 0 & a_{12} & 0 \\ 0 & 0 & a_{23} \\ a_{31} & 0 & 0 \end{vmatrix} + \begin{vmatrix} 0 & a_{12} & 0 \\ a_{21} & 0 & 0 \\ 0 & 0 & a_{33} \end{vmatrix} \\
&\quad + \begin{vmatrix} 0 & 0 & a_{13} \\ 0 & a_{22} & 0 \\ a_{31} & 0 & 0 \end{vmatrix} + \begin{vmatrix} 0 & 0 & a_{13} \\ a_{21} & 0 & 0 \\ 0 & a_{32} & 0 \end{vmatrix} \\
&= a_{11}a_{22}a_{33} - a_{11}a_{23}a_{32} + a_{12}a_{23}a_{31} - a_{12}a_{21}a_{33} - a_{13}a_{22}a_{31} + a_{13}a_{21}a_{32} \\
&= a_{11}(a_{22}a_{33} - a_{23}a_{32}) - a_{12}(a_{21}a_{33} - a_{23}a_{31}) + a_{13}(a_{21}a_{32} - a_{22}a_{31}) \\
&= a_{11} \underbrace{\begin{vmatrix} a_{22} & a_{23} \\ a_{32} & a_{33} \end{vmatrix}}_{\text{cofactor of } a_{11}} + a_{12}(-1) \underbrace{\begin{vmatrix} a_{21} & a_{23} \\ a_{31} & a_{33} \end{vmatrix}}_{\text{minor of } a_{12}} + a_{13} \underbrace{\begin{vmatrix} a_{21} & a_{22} \\ a_{31} & a_{32} \end{vmatrix}}_{\text{cofactor of } a_{13}} \quad (3.86)
\end{aligned}$$

In (3.86), the determinant of A_3 is decomposed into the sum of signed determinants of smaller 2×2 matrices called *co-factors*, each scaled by a corresponding *factor*. The sign of the co-factors depend on the number of row permutations required to get the matrix in a diagonal form; the sign is $(-1)^{\text{num perms}}$ which happens to be $(-1)^{i+j}$. In general, for an $n \times n$ matrix A , the *minor* of a term a_{ij} is the determinant of an $(n-1) \times (n-1)$ sub-matrix of A that has the row i and column j removed, while its *co-factor* is the minor multiplied by $(-1)^{i+j}$. The minor for a_{ij} is denoted by M_{ij} and its co-factor by C_{ij} . Minors of the form M_{ii} are called principal minors.

The general formula for the determinant of an $n \times n$ matrix contains $n!$ terms, corresponding to all permutations of the choice of the column index for the non-zero entries corresponding to each row index. That is,

$$\det(A) = \sum_{(p_1, p_2, \dots, p_n) \in \text{Perm}(1, 2, \dots, n)} a_{1p_1} a_{2p_2} \dots a_{p_n} \quad (3.87)$$

In terms of co-factors, the formula for determinant is

$$\det(A) = \sum_{k=1}^n a_{ik} C_{ik} \quad (3.88)$$

for any $1 \leq i \leq n$.

3.10.2 Formula for Inverse

Let A be an $n \times n$ invertible matrix. In Section 3.4.1, we saw an elegant algorithm for computing A^{-1} . In (3.89), we present a closed form expression for A^{-1} in terms of the co-factors of A , even though the expression is very expensive to compute.

$$A^{-1} = \frac{1}{\det(A)} \begin{bmatrix} C_{11} & C_{12} & \cdots & C_{1n} \\ C_{21} & C_{22} & \cdots & C_{2n} \\ \cdot & \cdot & \cdots & \cdot \\ C_{k1} & C_{k2} & \cdots & C_{kn} \\ \cdot & \cdot & \cdots & \cdot \\ C_{n1} & C_{n2} & \cdots & C_{nn} \end{bmatrix}^T = \frac{1}{\det(A)} C^T \quad (3.89)$$

We denote the matrix in (3.89) consisting of the co-factors of A , by C . It can be easily verified that the expression in (3.89) is indeed A^{-1} .

$$\begin{aligned} AA^{-1} &= \frac{1}{\det(A)} \begin{bmatrix} \sum_{j=1}^n a_{1j}C_{1j} & \sum_{j=1}^n a_{1j}C_{2j} & \cdots & \sum_{j=1}^n a_{1j}C_{nj} \\ \sum_{j=1}^n a_{2j}C_{1j} & \sum_{j=1}^n a_{2j}C_{2j} & \cdots & \sum_{j=1}^n a_{2j}C_{nj} \\ \cdot & \cdot & \cdots & \cdot \\ \cdot & \cdot & \sum_{j=1}^n a_{ij}C_{ij} & \cdot \\ \cdot & \cdot & \cdots & \cdot \\ \sum_{j=1}^n a_{nj}C_{1j} & \sum_{j=1}^n a_{nj}C_{2j} & \cdots & \sum_{j=1}^n a_{nj}C_{nj} \end{bmatrix} \\ &= \frac{1}{\det(A)} \begin{bmatrix} \det(A) & 0 & \cdots & 0 \\ 0 & \det(A) & \cdots & 0 \\ \cdot & \cdot & \cdots & \cdot \\ \cdot & \cdot & \det(A) & \cdot \\ \cdot & \cdot & \cdots & \cdot \\ 0 & 0 & \cdots & \det(A) \end{bmatrix} = I_{n \times n} \quad (3.90) \end{aligned}$$

Recall from (3.88) that $\det(A) = \sum_{j=1}^n a_{ij}C_{ij}$ for any $1 \leq i \leq n$. However,

$\sum_{j=1}^n a_{ij}C_{kj} = 0$, if $i \neq k$. This is because, for $i \neq k$, $\sum_{j=1}^n a_{ij}C_{kj}$ is the determinant

of a matrix that has identical i^{th} and k^{th} rows and hence equals 0, by the derived property (1) for matrices.

The formula (3.89) for matrix inverse (if it exists) can be substituted in (3.4.1) to yield the *Cramer's rule* for solving the system $A\mathbf{x} = \mathbf{b}$. The Cramer's rule is:

$$\begin{aligned}
 x = A^{-1}\mathbf{b} &= \frac{1}{\det(A)} \begin{bmatrix} \sum_{j=1}^n b_j C_{1j} \\ \sum_{j=1}^n b_j C_{2j} \\ \cdot \\ \sum_{j=1}^n b_j C_{nj} \end{bmatrix} \\
 &= \frac{1}{\det(A)} \begin{bmatrix} \det(B_1) \\ \det(B_2) \\ \cdot \\ \det(B_n) \end{bmatrix} \tag{3.91}
 \end{aligned}$$

B_i is a matrix obtained by replacing the i^{th} column of A with the vector \mathbf{b} and keeping all other columns unaltered. This rule is never used in practical computations; the explicit formula only helps in analysis or derivation.

3.11 Eigenvalues and Eigenvectors

Let A be an $n \times n$ square matrix. Consider the function $f : \mathfrak{R}^n \rightarrow \mathfrak{R}^n$, defined as $f(\mathbf{x}) = A\mathbf{x}$. Suppose we are interested in vectors \mathbf{x} , for which, f returns a vector in the same direction as \mathbf{x} . Such vectors are called *eigenvectors* of A .

Eigenvector: Vector $\mathbf{x} \in \mathfrak{R}^n$ is called an *eigenvector* of an $n \times n$ matrix A , iff

$$A\mathbf{x} = \lambda\mathbf{x}, \exists \lambda \in \mathfrak{R} \tag{3.92}$$

The scalar λ is called an **eigenvalue** of A , corresponding to the eigenvector \mathbf{x} .

We will consider some special examples of eigenvectors and eigenvalues.

- For the simple case of $\lambda = 0$, any $\mathbf{x} \in N(A)$ is an eigenvector of A . Thus, if A is singular, so that $N(A) \neq \{\}$, $\lambda = 0$ and $\mathbf{x} \in N(A)$ are a valid eigenvalue-eigenvector pair.
- If A happens to be a projection matrix, (c.f., Section 3.9.1), i.e.,

$$A = \mathbf{ss}^T \frac{1}{\mathbf{s}^T \mathbf{s}}$$

for some $\mathbf{s} \in \mathfrak{R}^n$. Recall that, $A\mathbf{x} = \mathbf{x}$ for any \mathbf{x} in the column space of A . Therefore, $\mathbf{x} = \rho\mathbf{s}$ is an eigenvector of A for any $\rho \in \mathfrak{R}$, with 1 as the corresponding eigenvalue. As discussed above, any $\mathbf{x} \in N(A)$ is also an eigenvector of A with a corresponding eigenvalue of 0. However, for any other $\mathbf{x} \notin C(A)$, $A\mathbf{x} = \mathbf{c}\mathbf{s}$ and therefore \mathbf{x} is not an eigenvector.

Consider the permutation matrix P_{23} in (3.93):

$$P_{23} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \quad (3.93)$$

By inspection, we find that P_{23} has at least three eigenvectors, viz., $x_1 = [1 \ 0 \ 0]$ with eigenvalue $\lambda_1 = 1$, $x_2 = [0 \ 1 \ 1]$ with eigenvalue $\lambda_2 = 1$, and $x_3 = [0 \ -1 \ 1]$ with eigenvalue $\lambda_3 = -1$. Does P_{23} have any more eigenvectors? The answer is no. It turns out that any $n \times n$ matrix has exactly n orthonormal eigenvectors. Moreover, the trace of a matrix (i.e., the sum of its diagonal entries) always equals the sum of the eigenvalues corresponding to the orthonormal eigenvectors.

$$\text{tr}(A) = \sum_{i=1}^n \lambda_i$$

Thus, if we knew $n - 1$ eigenvalues of a matrix, we could easily determine its n^{th} eigenvalue. We will defer this discussion to a later part of this chapter.

3.11.1 Solving for Eigenvalues

The equation (3.92) defining the criterion for an eigenvalue λ can be re-written as in (3.94).

$$(A - \lambda I)\mathbf{x} = \mathbf{0} \quad (3.94)$$

For a solution \mathbf{x} to exist, $A - \lambda I$ must be singular (*i.e.*, non-invertible) and \mathbf{x} must lie in the null space $N(A - \lambda I)$. Therefore, $\det(A - \lambda I) = 0$ is a necessary and sufficient condition for λ to be an eigenvalue. Once the eigenvalue λ is determined, the corresponding eigenvectors can be determined by computing $N(A - \lambda I)$, a procedure that has been already discussed in Section 3.5.2. We will therefore first discuss the procedure for computing the solution to

$$\det(A - \lambda I) = 0 \quad (3.95)$$

As an example, when we apply the criterion in (3.95), to the matrix P_{23} , we get solutions as shown in (3.96):

$$\begin{aligned} \det(P_{23} - \lambda I) &= (1 - \lambda)\lambda^2 = 0 \\ \Rightarrow \lambda &= 1 \text{ or } \lambda = -1 \end{aligned} \quad (3.96)$$

Substituting these two values into the system $(A - \lambda I)\mathbf{x} = \mathbf{0}$, we get one matrix for each possible value of λ . It can be verified that the basis for the null space of $(A - \lambda I)$ obtained using the elimination process discussed in Section 3.6 (particularly, equation 3.27) is indeed $[1 \ 0 \ 0]^T$ and $[0 \ 1 \ 1]^T$ for eigenvalue $\lambda_1 = 1$, and $[0 \ -1 \ 1]$ for eigenvalue $\lambda_3 = -1$.

3.11.2 Some Properties of Eigenvalues and Eigenvectors

How are the eigenvectors and eigenvalues of a matrix affected when transformations are performed on the matrix? Below, we list some properties of eigenvalues with respect to matrix transformations.

1. If $A\mathbf{x} = \lambda\mathbf{x}$, then $(A + kI)\mathbf{x} = (\lambda + k)\mathbf{x}$. That is, the eigenvalues of $A + \lambda I$ are the eigenvalues of A , incremented by k , without any change in corresponding eigenvectors.
2. Consider the matrix R in (3.97):

$$R = \begin{bmatrix} 0 & 3 \\ -2 & 0 \end{bmatrix} \quad (3.97)$$

The eigenvalues of R can be found as follows: $\det(R - \lambda I) = \lambda^2 + 6 = 0 \Rightarrow \lambda = \pm\sqrt{6}i$. The eigenvalues of a matrix could be complex numbers as this example illustrates. In fact, eigenvalues always appear as complex conjugates, as in this example.

3. Let λ be an eigenvalue of A and \mathbf{x} its corresponding eigenvector, *i.e.*, $A\mathbf{x} = \lambda\mathbf{x}$. It can be shown that the complex conjugates $\bar{\lambda}$ and $\bar{\mathbf{x}}$ also form an eigenvalue-eigenvector pair for \bar{A} . Thus, $\bar{A}\bar{\mathbf{x}} = \bar{\lambda}\bar{\mathbf{x}}$. If A happens to have only real entries, then, $A\bar{\mathbf{x}} = \bar{\lambda}\bar{\mathbf{x}}$.
4. The eigenvalues of upper and lower triangular matrices can be computed very easily. By derived property (4) of determinants, the determinant of an upper triangular matrix is the product of its diagonal entries. Let U be an $n \times n$ upper triangular matrix, with u_{ij} being the entry corresponding to the i^{th} row and j^{th} column. First we note that $U - \lambda I$ will also be upper triangular, since I is upper triangular and since the sum of upper triangular matrices is also upper triangular (the space of upper triangular matrices is a vector space, as shown in Section 3.8). Now, $\det(U - \lambda I) = \prod_{i=1}^n (u_{ii} - \lambda) = 0$. The eigenvalues correspond to solutions of this equation; they are $\lambda_i = u_{ii}$, $1 \leq i \leq n$. The eigenvectors can be computed by solving the systems $(U - \lambda_i I)\mathbf{x}_i = \mathbf{0}$ by simple back-substitutions, as illustrated in Section 3.3.1.
5. If \mathbf{x} is an eigenvector of A with a corresponding eigenvalue λ , we have $A\mathbf{x} = \lambda\mathbf{x}$. Therefore, $A^2\mathbf{x} = A(A\mathbf{x}) = \lambda A\mathbf{x} = \lambda^2\mathbf{x}$. Thus, \mathbf{x} is an eigenvector of A^2 as well, with a corresponding eigenvalue of λ^2 . This statement can be generalized: *If \mathbf{x} is an eigenvector of A with a corresponding eigenvalue λ , \mathbf{x} is also an eigenvector of A^k , with corresponding eigenvalue λ^k .*
6. The eigenvectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ of a matrix A are linearly independent if all its eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$ are different. This can be proved by contradiction⁹ However, the eigenvectors, could be independent even if eigenvalues are repeated; but it is not always true. For instance, any triangular matrix having some identical diagonal elements (as in the case of the identity matrix) has linearly independent eigenvectors, even though some eigenvalues are identical.
7. In many engineering problems, we are faced with the system of equations

$$\mathbf{b}_{i+1} = A\mathbf{b}_i, \quad \forall i \geq 0 \tag{3.98}$$

That is, $\mathbf{b}_i = A^i\mathbf{b}_0$. If A has n linearly independent eigenvectors (so that they span \mathfrak{R}^n), these systems can be solved efficiently, by expressing \mathbf{b}_0 as a linear combination of the eigenvectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ of A .

$$\mathbf{b}_0 = \sum_{k=1}^n c_k \mathbf{v}_k$$

⁹Exercise.

where, $c_k \in \Re, \forall 1 \leq k \leq n$. Consequently, any $\mathbf{b}_i, i \geq 0$ can be computed efficiently as

$$\mathbf{b}_i = \sum_{k=1}^n \lambda_k^i c_k \mathbf{v}_k \quad (3.99)$$

8. Consider the fibonacci sequence $f_{i+2} = f_i + f_{i+1}, i \geq 0$, with $f_0 = 0$ and $f_1 = 1$. The recurrence relation can be written as a linear system (3.100).

$$\underbrace{\begin{bmatrix} f_{i+2} \\ f_{i+1} \end{bmatrix}}_{\mathbf{b}_{i+1}} = \underbrace{\begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix}}_A = \underbrace{\begin{bmatrix} f_{i+1} \\ f_i \end{bmatrix}}_{\mathbf{b}_i} \quad (3.100)$$

Note that $\mathbf{b}_0 = [0 \ 1]^T$. The system of equations (3.100) is of the same form $\mathbf{b}_{i+1} = A\mathbf{b}_i, \forall 0 \leq i \leq n$ discussed above and therefore, the expression for \mathbf{b}_i can be derived using (3.99), after computing values of $\lambda_1, \mathbf{v}_1, c_1$ and λ_2, \mathbf{v}_2 and c_2 . The values of $\lambda_1 = \frac{1}{2}(1 + \sqrt{5}) = 1.6180$ and $\lambda_2 = \frac{1}{2}(1 - \sqrt{5}) = -0.6180$ can be computed by solving $\det(A - \lambda I) = 0$. Substituting these values of λ , eigenvectors \mathbf{v}_1 and \mathbf{v}_2 can be obtained as in (3.101).

$$\mathbf{v}_1 = \begin{bmatrix} -0.8507 \\ -0.5257 \end{bmatrix}, \mathbf{v}_2 = \begin{bmatrix} 0.5257 \\ -0.8507 \end{bmatrix} \quad (3.101)$$

A closed form expression is $\mathbf{b}_i = c_1(1.6180)^i[-0.8507 \ -0.5257]^T - c_2(0.6180)^i[0.5257 \ -0.8507]^T$.

Another application of this general technique is in differential equations. Let us say we are given the differential equation $x'' + a_1x' + a_2x = 0$. This equation can be equivalently expressed as

$$\mathbf{y}' = \underbrace{\begin{bmatrix} -a_1 & -a_2 \\ 1 & 0 \end{bmatrix}}_A \mathbf{y} \quad (3.102)$$

where, $\mathbf{y} = [x' \ x]^T$. The n^{th} derivative of x can be expressed in a closed form by determining the eigenvalues and eigenvectors of A .

9. If λ is an eigenvalue of a matrix A , then it is also an eigenvalue of A^T . This is because $\det(A - \lambda I) = \det((A - \lambda I)^T) = \det(A^T - \lambda I)$. The eigenvectors for the same eigenvalues could however differ between A and A^T .
10. Another general property of any square matrix is that the sum of its eigenvalues equals its trace. Additionally, the product of its eigenvalues equals its determinant. Consequently, for any 2×2 matrix, if the trace is negative and the determinant positive, the real parts of both its eigenvalues must be negative.
11. A Markov matrix¹⁰ M is an $n \times n$ matrix such that (1) all its entries are ≥ 0 and (2) the sum of all entries in each column is 1. An example Markov matrix M_3 is

$$M = \begin{bmatrix} 0.1 & 0.25 & 0.3 & 0.35 \\ 0.2 & 0.25 & 0.3 & 0.05 \\ 0.3 & 0.25 & 0.4 & 0.15 \\ 0.4 & 0.25 & 0 & 0.45 \end{bmatrix} \quad (3.103)$$

A very special property of markov matrices is that exactly one eigenvalue of any markov matrix equals 1 and the rest of its eigenvalues are strictly less than 0. For example, the M_3 has following eigenvalues: 1.0000, -0.2168, 0.3428, 0.0740. The first part of this property can be proved as follows. The matrix $M - I$ is singular, because the sum of the rows is a zero vector. Therefore, $\det(M - I) = 0$. Thus, $\lambda = 1$ must be an eigenvalue of M .

In probabilistic models, we often have systems of the form $\mathbf{p}_{i+1} = \mathbf{A}\mathbf{p}_i$, $\forall i \geq 0$, similar to equation (3.98). A closed form solution can be obtained using the idea of (3.99)

$$\mathbf{p}_i = \sum_{k=1}^n \lambda_k^i c_k \mathbf{v}_k$$

where, $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ are the eigenvectors of M and its eigenvalues are $\lambda_1, \lambda_2, \dots, \lambda_n$. If $\lambda_1 = 1$, then $\lambda_i < 1, \forall 2 \leq i \leq n$. Hence, as $i \rightarrow \infty$, $\mathbf{p}_i \rightarrow c_1 \mathbf{v}_1$.

12. If A is an $n \times n$ matrix with real valued entries and is symmetric, *i.e.*, $A = A^T$, then, its eigenvalues are real. Further, the eigenvectors of a symmetric matrix can be chosen to be orthogonal. In mathematics, this is called the *spectral theorem* while in mechanics it is called the *principal axis theorem*.

¹⁰The matrix entries of a markov entries represent probabilities of transitions within/between states.

Theorem 37 *If A is symmetric then (1) all its eigenvalues are real and (2) there exists an orthonormal basis Q of A , consisting of its eigenvectors.*

Proof for part (1): Let λ be an eigenvalue of A and \mathbf{x} be its corresponding eigenvector; $A\mathbf{x} = \lambda\mathbf{x}$. Then, premultiplying both sides by $\bar{\mathbf{x}}^T$, we get

$$\bar{\mathbf{x}}^T A\mathbf{x} = \lambda\bar{\mathbf{x}}^T \mathbf{x} \quad (3.104)$$

As mentioned earlier, the complex conjugates $\bar{\lambda}$ and \bar{x} also form an eigenvalue-eigenvector pair for a real matrix A ; $A\bar{\mathbf{x}} = \bar{\lambda}\bar{\mathbf{x}}$. This implies that $\bar{\mathbf{x}}^T A^T = \bar{\mathbf{x}}^T A = \bar{\lambda}\bar{\mathbf{x}}^T$ and therefore,

$$\bar{\mathbf{x}}^T A\mathbf{x} = \bar{\lambda}\bar{\mathbf{x}}^T \mathbf{x} \quad (3.105)$$

We note that the left hand sides of (3.104) and (3.105) are the same. Equating the right hand sides of these equations,

$$\lambda\bar{\mathbf{x}}^T \mathbf{x} = \bar{\lambda}\bar{\mathbf{x}}^T \mathbf{x} \quad (3.106)$$

$\bar{\mathbf{x}}^T \mathbf{x}$ is always real and non-negative. It is 0 only if $\mathbf{x} = \mathbf{0}$. Therefore, $\lambda = \bar{\lambda} \Rightarrow \lambda \in \mathfrak{R}$. \square

13. If A is a real symmetric matrix, the number of positive pivots and number of negative pivots are respectively equal to the number of positive and negative eigenvalues.
14. Two $n \times n$ matrices A and B are called *similar* if there exists an invertible $n \times n$ matrix M such that $M^{-1}BM = A$. A property of similar matrices is that they have same determinants, since $\det(A) = \det(M^{-1})\det(B)\det(M) = \frac{1}{\det(M)}\det(B)\det(M) = \det(B)$. A more fundamental property is that similar matrices have the same eigenvalues, though they could differ in their eigenvectors.

Theorem 38 *If A and B are similar matrices, they have the same eigenvalues.*

Proof: Let λ be an eigenvalue of A . Since A and B are similar, there exists an invertible matrix M such that, $M^{-1}BM = A$. $A\mathbf{x} = \lambda\mathbf{x} \Rightarrow (MAM^{-1})M\mathbf{x} = \lambda M\mathbf{x} \Rightarrow B(M\mathbf{x}) = \lambda(M\mathbf{x})$, that is, if λ is an eigenvalue

of A and \mathbf{x} is the corresponding eigenvector, then λ is an eigenvalue of B and $M\mathbf{x}$ is its corresponding eigenvector.

Similarly, $B\mathbf{x} = \lambda\mathbf{x} \Rightarrow (M^{-1}AM)M^{-1}\mathbf{x} = \lambda M^{-1}\mathbf{x} \Rightarrow B(M^{-1}\mathbf{x}) = \lambda(M^{-1}\mathbf{x})$, that is, if λ is an eigenvalue of B and \mathbf{x} is the corresponding eigenvector, then λ is an eigenvalue of A and $M^{-1}\mathbf{x}$ is its corresponding eigenvector. \square

At this point, we state the observation that matrices of the form $kI_{n \times n}$ are only similar to themselves, since, for any invertible matrix M , $M^{-1}(kI_{n \times n})M = kI_{n \times n}$.

3.11.3 Matrix Factorization using Eigenvectors

Let A be an $n \times n$ matrix, with n eigenvectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ and corresponding eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$. Let V be a matrix with the eigenvectors as columns. Postmultiplying A by V , we get

$$AV = [\lambda_1\mathbf{v}_1 \quad \lambda_2\mathbf{v}_2 \quad \dots \quad \lambda_n\mathbf{v}_n] = [\mathbf{v}_1 \quad \mathbf{v}_2 \quad \dots \quad \mathbf{v}_n] \underbrace{\begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \lambda_k & \cdot \\ \cdot & \cdot & \dots & \cdot \\ 0 & 0 & \dots & \lambda_n \end{bmatrix}}_{\text{Eigenvalue matrix } \Lambda} = V\Lambda$$

that is, $AV = V\Lambda$. The diagonal matrix Λ consists of eigenvalues along its diagonal and is called the *eigenvalue matrix*.

If the eigenvectors are linearly independent, V is invertible. Premultiplying AV by V^{-1} ,

$$V^{-1}AV = \Lambda$$

Another equivalent equation is

$$A = V\Lambda V^{-1} \tag{3.107}$$

This procedure of premultiplying a matrix by the inverse of its eigenvector matrix and post-multiplying it by the eigenvector matrix to obtain a diagonal matrix of its eigenvalues, is called *diagonalization*. Diagonalization can be generalized to powers of k :

$$A^k = V\Lambda^k V^{-1}$$

Thus, eigenvalues and eigenvectors provide a great way to understand the powers of a matrix. Further, if $|\lambda_i| < 1$, $\Lambda^k \rightarrow 0$, as $k \rightarrow \infty$. Therefore, if $|\lambda_i| < 1$,

$A^k \rightarrow 0$, as $k \rightarrow \infty$. As another example, if we define $e^{\rho A} = \sum_{n=0}^{\infty} \frac{1}{n!} (A\rho)^n$,

where $\rho \in \mathfrak{R}$, then using the above property, it can be shown that $e^{\rho A} = V e^{\rho \Lambda} V^{-1} = V \text{diag}(e^{\rho \lambda_1}, e^{\rho \lambda_2}, \dots, e^{\rho \lambda_n}) V^{-1}$, where $\text{diag}(c_1, c_2, \dots, c_n)$ returns an $n \times n$ diagonal matrix with the i^{th} diagonal entry as c_i .

If A is symmetric, the eigenvector matrix V could be chosen to be a matrix of orthonormal vectors, denoted by Q . Note that $Q^{-1} = Q^T$. Thus, for a symmetric A , the equation (3.107) can be re-written as:

$$A = Q \Lambda Q^T = \sum_{i=1}^n \lambda_i (\mathbf{q}_i \mathbf{q}_i^T) \quad (3.108)$$

From Section 3.9.1, we recall that $(\mathbf{q}_i \mathbf{q}_i^T)$ is a projection matrix. Moreover, if $i \neq j$, $(\mathbf{q}_i \mathbf{q}_i^T)$ is orthogonal to $(\mathbf{q}_j \mathbf{q}_j^T)$. This gives us another perspective of symmetric matrices - as a linear combination of orthogonal projection matrices. Also, since Q is of rank 1 and invertible, we can infer that A is similar to Λ . The diagonal matrix Λ can be thought of as a canonical form for the family of matrices similar to A . However, if A is not a full rank matrix, there exists an 'almost diagonal form', called the *Jordan form* [?], which is similar to A , containing the eigenvalues of A along its diagonal, with the only other non-zero entries being along the super-diagonal.

One more illustration of the utility of matrix factorization using eigenvectors is the interpretation of level sets involving the quadratic form $\mathbf{x}^T A \mathbf{x} = \mathbf{x}^T Q \Lambda Q^T \mathbf{x}$ for a symmetric matrix A . The *level set* of a real-valued function f of $\mathbf{x} \in \mathfrak{R}^n$ is a set of the form $\{\mathbf{x} | f(\mathbf{x}) = c\}$, where c is a constant. Using the eigenvalue factorization of matrices, the level set $\{\mathbf{x} | \mathbf{x}^T Q \Lambda Q^T \mathbf{x} = c\}$ can be interpreted as an ellipsoid in n dimensions, with each eigenvector-eigenvalue pair specifying the direction and the length respectively of an axis of the ellipsoid.

3.12 Positive Definite Matrices

Positive definite matrix: A *positive definite (p.d.) matrix* is a symmetric matrix with all positive eigenvalues. That M is a p.d. matrix is also denoted by $M > 0$.

By virtue of property of symmetric matrices, all the pivots in the rref of a p.d. matrix are also positive. Since the determinant of matrix equals the product of its eigenvalues, the determinant of a p.d. matrix is also positive; however, it is not necessary that a matrix with positive determinant is also p.d.

A matrix is called *positive semi-definite* (p.s.d.), if all its eigenvalues are non-negative. That M is p.s.d. is also denoted by $M \geq 0$.

3.12.1 Equivalent Conditions

We will list down some necessary and sufficient conditions for a matrix A to be positive definite or positive semi-definite:

1. A matrix A is p.d. iff all its principal minors (c.f. Section 3.10.1) are positive. As an example, if A is a 2×2 matrix, we must have $a_{11} > 0$ and $a_{11}a_{22} - a_{12}a_{21} > 0$ in order for A to be p.d. On the other hand, if all its principal minors are non-negative, the matrix is p.s.d.

2. Another equivalent definition for positive definiteness is: A matrix A is p.d. iff, $\forall \mathbf{x} \neq \mathbf{0}$, $\mathbf{x}^T A \mathbf{x} > 0$. This condition can be rewritten as $\forall \mathbf{x} \neq \mathbf{0}$, $\sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j > 0$. If $\forall \mathbf{x} \neq \mathbf{0}$, $\mathbf{x}^T A \mathbf{x} \geq 0$, A is p.s.d.

3. The condition $\forall \mathbf{x} \neq \mathbf{0}$, $\sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j > 0$ involves a quadratic expression. The expression is guaranteed to be greater than 0 $\forall \mathbf{x} \neq \mathbf{0}$ iff it can

be expressed as $\sum_{i=1}^n \lambda_i \left(\sum_{j=1}^{i-1} \beta_{ij} x_{ij} + x_{ii} \right)^2$, where $\lambda_i \geq 0$. This is possible

iff A can be expressed as LDL^T , where, L is a lower triangular matrix with 1 in each diagonal entry and D is a diagonal matrix of all positive diagonal entries. Or equivalently, it should be possible to factorize A as RR^T , where $R = LD^{1/2}$ is a lower triangular matrix. Note that any symmetric matrix A can be expressed as LDL^T , where L is a lower triangular matrix with 1 in each diagonal entry and D is a diagonal matrix; positive definiteness has only an additional requirement that the diagonal entries of D be positive. This gives another equivalent condition for positive definiteness: *Matrix A is p.d. if and only if, A can be uniquely factored as $A = RR^T$, where R is a lower triangular matrix with positive diagonal entries.* This factorization of a p.d. matrix is referred to as *Cholesky factorization*.

Recall that Gauss elimination on a matrix A yields its factorization as $A = LU$ and the diagonal entries of L are pivots. Therefore, if A is symmetric matrix such that Gauss elimination on it yields positive pivots, A is positive definite.

To illustrate the equivalence of the above definitions of positive definiteness, consider the matrix P below:

$$P = \begin{bmatrix} 1 & 1 & 2 & 1 \\ 1 & 10 & 14 & 4 \\ 2 & 14 & 21 & 9 \\ 1 & 4 & 9 & 20 \end{bmatrix} \quad (3.109)$$

The matrix is positive definite and this can be proved by showing any of the following properties:

1. All the eigenvalues of P , viz., $\lambda_1 = 0.1644$, $\lambda_2 = 0.9371$, $\lambda_3 = 14.4091$, $\lambda_4 = 36.4893$ are positive. and therefore $P > 0$.
2. The principal minors of P are 1, 9, 9 and 81. All the four principal minors are positive and thus $P > 0$.
3. Matrix P can be factorized as LL^T , where

$$L = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 3 & 0 & 0 \\ 2 & 4 & 1 & 0 \\ 1 & 1 & 3 & 3 \end{bmatrix} \quad (3.110)$$

Since L is lower triangular and since all its diagonal entries are positive, $P > 0$.

3.12.2 Some properties

We will list some properties of positive definite matrices, using an appropriate definition of positive definiteness as required.

1. If matrices $A > 0$ and $B > 0$, then $A + B > 0$. This follows from the fact that $\forall \mathbf{x} \neq \mathbf{0}$, $\mathbf{x}^T A \mathbf{x} > 0$ and $\forall \mathbf{x} \neq \mathbf{0}$, $\mathbf{x}^T B \mathbf{x} > 0$ implies that $\forall \mathbf{x} \neq \mathbf{0}$, $\mathbf{x}^T (A + B) \mathbf{x} > 0$. Similarly, $AB > 0$ and for any $c > 0$, $cA > 0$.
2. If $A > 0$, then $\forall \mathbf{x} \neq \mathbf{0}$, $\mathbf{x}^T A \mathbf{x} > 0$ implies $(\mathbf{x}^T A \mathbf{x})^T = \mathbf{x}^T A^T \mathbf{x} > 0$, that is, $A^T > 0$.
3. Let A be an $m \times n$ matrix. Recall from Section 3.9.2, the important matrix $A^T A$ which happened to be an $n \times n$ matrix. If A is full column rank, the only vector in its null space is $\mathbf{0}$. Note that $\forall \mathbf{x} \neq \mathbf{0}$, $\mathbf{x}^T A^T A \mathbf{x} = \|A \mathbf{x}\|^2 > 0$. Thus, $A^T A$ is always p.d. if A is non-singular.
4. Every p.d. matrix is invertible and its inverse is also p.d. That is, if $A > 0$ then A^{-1} exists and $A^{-1} > 0$.
5. If $A > 0$, the diagonal entries of A are real and positive. Consequently, the trace $tr(A)$ is also positive.

Testing for positive definiteness of a matrix arises in several applications, including optimization. Determining the local minimum of a function $f(\mathbf{x})$, $\mathbf{x} \in \mathcal{D}$, $\mathcal{D} \subseteq \Re^k$ involves determining points $\hat{\mathbf{x}}$ at which $\nabla f(\hat{\mathbf{x}}) = 0$ and $\nabla^2 f(\hat{\mathbf{x}}) > 0$ (positive curvature at $\hat{\mathbf{x}}$).

3.13 Singular Value Decomposition

In Section 3.11.3, we discussed that a full rank symmetric matrix can be factorized into $Q\Lambda Q^T$, where Q is an orthonormal matrix and Λ is a diagonal matrix. This factorization can be extended to any matrix and it is called *Singular Value Decomposition*, abbreviated as *SVD*. The singular value decomposition of any $m \times n$ matrix A is factorization of A as $U\Sigma V^T$, where Σ is a diagonal matrix and U and V are orthonormal matrices.

We will construct the matrices U and V as follows. Let r be the rank of A and let

- $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r$ be an orthonormal basis for the column space of A .
- $\mathbf{v}_{r+1}, \mathbf{v}_{r+2}, \dots, \mathbf{v}_n$ be an orthonormal basis for the null space of A .
- $\mathbf{u}_{r+1}, \mathbf{u}_{r+2}, \dots, \mathbf{u}_m$ be an orthonormal basis for the null space of A^T .
- $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r$ be such that $\mathbf{x}_i = A^T \mathbf{u}_i$ and $\mathbf{v}_i = \frac{1}{\|\mathbf{x}_i\|} \mathbf{x}_i$.

The relationship between \mathbf{u}_i and \mathbf{v}_i is therefore $A^T \mathbf{u}_i = \sigma_{ii} \mathbf{v}_i$, with

$$\sigma_{ii} = \begin{cases} \|A^T \mathbf{u}_i\| & \text{if } i \leq r \\ 0 & \text{if } i > r \end{cases} \quad (3.111)$$

This system of equations can be written in matrix form as

$$A^T U = V \Sigma \quad (3.112)$$

where, $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m$ are the columns of U and $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ are the columns of V . Σ is an $n \times n$ diagonal matrix with its ij^{th} entry given by σ_{ij} , such that

$$\sigma_{ij} = \begin{cases} 0 & \text{if } i \neq j \\ \|A^T \mathbf{u}_i\| & \text{if } i = j \text{ and } i \leq r \\ 0 & \text{if } i = j \text{ and } i > r \end{cases} \quad (3.113)$$

It can be shown that $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r$ are orthonormal and form a basis for the row space of A . Theorem 34 stated that the row space $C(A^T)$ and right null space $N(A)$ are orthogonal complements. Similarly, the column space $C(A)$ and left null space $N(A^T)$ are orthogonal complements. Therefore, $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m$ is an orthonormal basis for \mathfrak{R}^m , while $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ is an orthonormal basis for \mathfrak{R}^n .

Since $U^{-1} = U^T$, we can rewrite (3.112) as

$$A = U\Sigma V^T \quad (3.114)$$

Furthermore, $AA^T = U\Sigma^2U^T$ and $A^T A = V\Sigma^2V^T$, which are spectral decompositions, implying that the columns of U and V and eigenvectors of AA^T and $A^T A$ respectively and the diagonal entries of Σ are square roots of the eigenvalues of AA^T (or equivalently $A^T A$).

As an example, if P is the full rank, symmetric matrix in (3.109), the matrices U , Σ and V are

$$U = \begin{bmatrix} -165/2423 & 76/4167 & 637/688 & -892/2403 \\ -467/1012 & 373/992 & -577/1726 & -757/1036 \\ -367/508 & 48/133 & 318/1909 & 869/1536 \\ -172/337 & -407/477 & -329/5765 & -211/2328 \end{bmatrix} \quad (3.115)$$

$$\Sigma = \begin{bmatrix} 1715/47 & 0 & 0 & 0 \\ 0 & 11657/809 & 0 & 0 \\ 0 & 0 & 477/509 & 0 \\ 0 & 0 & 0 & 265/1612 \end{bmatrix} \quad (3.116)$$

$$V = \begin{bmatrix} -165/2423 & 76/4167 & 637/688 & -892/2403 \\ -467/1012 & 373/992 & -577/1726 & -757/1036 \\ -367/508 & 48/133 & 318/1909 & 869/1536 \\ -172/337 & -407/477 & -329/5765 & -211/2328 \end{bmatrix} \quad (3.117)$$

On the other hand, if P is a singular matrix of rank 2, given by

$$P = \begin{bmatrix} 1 & 3 & 1 \\ 2 & 3 & 1 \\ 3 & 6 & 2 \end{bmatrix} \quad (3.118)$$

then P can be decomposed into the following matrices:

$$U = \begin{bmatrix} -1301/3398 & 794/1101 & -780/1351 \\ -450/1039 & -715/1033 & -780/1351 \\ -337/413 & 203/6999 & 780/1351 \end{bmatrix} \quad (3.119)$$

$$\Sigma = \begin{bmatrix} 2565/299 & 0 & 0 \\ 0 & 687/1076 & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad (3.120)$$

$$V = \begin{bmatrix} -799/1854 & -647/717 & 0 \\ -1814/2119 & 453/1108 & -228/721 \\ -567/1987 & 151/1108 & 684/721 \end{bmatrix} \quad (3.121)$$

Notice that, since P is singular and of rank 2, its null space has dimension 1 and one of its eigenvalues is 0.

3.13.1 Pseudoinverse

The SVD of a matrix that is not full rank (such as P in (3.118)) can be used to compute its so-called *Moore-Penrose pseudoinverse*.

Pseudoinverse: *The pseudoinverse A^+ of an $m \times n$ matrix A is a unique $n \times m$ matrix, satisfying all the following criteria:*

1. $AA^+A = A$
2. $A^+AA^+ = A^+$
3. $\overline{(AA^+)}^T = AA^+$
4. $\overline{(A^+A)}^T = A^+A$

The pseudoinverse of a non-singular square matrix is the same as its inverse. A pseudoinverse of a rectangular matrix of full column rank is the left inverse, while a pseudoinverse of a rectangular matrix of full row rank is the right inverse (*c.f.* Section 3.4.2).

Consider an $n \times n$ diagonal matrix Σ having rank k .

$$\Sigma = \begin{bmatrix} \sigma_{11} & 0 & \dots & 0 & 0 & \dots & 0 \\ 0 & \sigma_{22} & \dots & 0 & 0 & \dots & 0 \\ \cdot & \cdot & \dots & \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \dots & \cdot & \cdot & \dots & \cdot \\ 0 & 0 & \dots & \sigma_{kk} & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 & 0 & \dots & 0 \\ \cdot & \cdot & \dots & \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \dots & \cdot & \cdot & \dots & \cdot \\ 0 & 0 & \dots & 0 & 0 & \dots & 0 \end{bmatrix} \quad (3.122)$$

The pseudoinverse Σ^+ of Σ is:

$$\Sigma^+ = \begin{bmatrix} \frac{1}{\sigma_{11}} & 0 & \dots & 0 & 0 & \dots & 0 \\ 0 & \frac{1}{\sigma_{22}} & \dots & 0 & 0 & \dots & 0 \\ \cdot & \cdot & \dots & \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \dots & \cdot & \cdot & \dots & \cdot \\ 0 & 0 & \dots & \frac{1}{\sigma_{kk}} & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 & 0 & \dots & 0 \\ \cdot & \cdot & \dots & \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \dots & \cdot & \cdot & \dots & \cdot \\ 0 & 0 & \dots & 0 & 0 & \dots & 0 \end{bmatrix} \quad (3.123)$$

The pseudoinverse P^+ of any non full rank matrix P can be computed using its singular value decomposition $U\Sigma V^T$ and the pseudoinverse Σ^+ of the diagonal matrix Σ as:

$$P^+ = V\Sigma^+U \quad (3.124)$$

Chapter 4

Convex Optimization

4.1 Introduction

4.1.1 Mathematical Optimization

The problem of mathematical optimization is to minimize a non-linear cost function $f_0(x)$ subject to inequality constraints $f_i(x) \leq 0, i = 1, \dots, m$ and equality constraints $h_i(x) = 0, i = 1, \dots, p$. $x = (x_1, \dots, x_n)$ is a vector of variables involved in the optimization problem. The general framework of a non-linear optimization problem is outlined in (4.1).

$$\begin{array}{ll} \text{minimize} & f_0(x) \\ \text{subject to} & f_i(x) \leq 0, \quad i = 1, \dots, m \\ & h_i(x) = 0, \quad i = 1, \dots, p \\ \text{variable } x = & (x_1, \dots, x_n) \end{array} \quad (4.1)$$

It is obviously very useful and arises throughout engineering, statistics, estimation and numerical analysis. In fact there is the tautology that ‘everything is an optimization problem’, though the tautology does not convey anything useful. The most important thing to note first is that the optimization problem is extremely hard in general. The solution and method is very much dependent on the property of the objective function as well as properties of the functions involved in the inequality and equality constraints. There are no good methods for solving the general non-linear optimization problem. In practice, you have to make some compromises, which usually translates to finding locally optimal solutions efficiently. But then you get only suboptimal solutions, unless you are willing to do global optimizations, which is for most applications too expensive.

There are important exceptions for which the situation is much better; the global optimum in some cases can be found efficiently and reliably. Three best known exceptions are

1. least-squares
2. linear programming
3. convex optimization problems - more or less the most general class of problems that can be solved efficiently.

Least squares and linear programming have been around for quite some time and are very special types of convex optimization problems. Convex programming was not appreciated very much until last 15 years. It has drawn attention more recently. In fact many combinatorial optimization problems have been identified to be convex optimization problems. There are also some exceptions besides convex optimization problems, such as singular value decomposition (which corresponds to the problem of finding the best rank- k approximation to a matrix, under the Frobenius norm) *etc.*, which has an exact global solution.

We will first introduce some general optimization principles. We will subsequently motivate the specific class of optimization problems called convex optimization problems and define convex sets and functions. Next, the theory of lagrange multipliers will be motivated and duality theory will be introduced. As two specific and well-studied examples of convex optimization, techniques for least squares and linear programming will be discussed to contrast them against generic convex optimization. Finally, we will dive into techniques for solving general convex optimization problems.

4.1.2 Some Topological Concepts in \mathfrak{R}^n

The definitions of some basic topological concepts in \mathfrak{R}^n could be helpful in the discussions that follow.

Definition 12 [Balls in \mathfrak{R}^n]: Consider a point $\mathbf{x} \in \mathfrak{R}^n$. Then the closed ball around \mathbf{x} of radius ϵ is defined as

$$\mathcal{B}[\mathbf{x}, \epsilon] = \{\mathbf{y} \in \mathfrak{R}^n \mid \|\mathbf{y} - \mathbf{x}\| \leq \epsilon\}$$

Likewise, the open ball around \mathbf{x} of radius ϵ is defined as

$$\mathcal{B}(\mathbf{x}, \epsilon) = \{\mathbf{y} \in \mathfrak{R}^n \mid \|\mathbf{y} - \mathbf{x}\| < \epsilon\}$$

For the 1-D case, open and closed balls degenerate to open and closed intervals respectively.

Definition 13 [Boundedness in \mathfrak{R}^n]: We say that a set $\mathcal{S} \subset \mathfrak{R}^n$ is bounded when there exists an $\epsilon > 0$ such that $\mathcal{S} \subseteq \mathcal{B}[0, \epsilon]$.

In other words, a set $\mathcal{S} \subseteq \mathfrak{R}^n$ is bounded means that there exists a number $\epsilon > 0$ such that for all $\mathbf{x} \in \mathcal{S}$, $\|\mathbf{x}\| \leq \epsilon$.

Definition 14 [Interior and Boundary points]: A point \mathbf{x} is called an interior point of a set \mathcal{S} if there exists an $\epsilon > 0$ such that $\mathcal{B}(\mathbf{x}, \epsilon) \subseteq \mathcal{S}$.

In other words, a point $\mathbf{x} \in \mathcal{S}$ is called an interior point of a set \mathcal{S} if there exists an open ball of non-zero radius around \mathbf{x} such that the ball is completely contained within \mathcal{S} .

Definition 15 [Interior of a set]: Let $\mathcal{S} \subseteq \mathbb{R}^n$. The set of all points lying in the interior of \mathcal{S} is denoted by $\text{int}(\mathcal{S})$ and is called the interior of \mathcal{S} . That is,

$$\text{int}(\mathcal{S}) = \{\mathbf{x} | \exists \epsilon > 0 \text{ s.t. } \mathcal{B}(\mathbf{x}, \epsilon) \subset \mathcal{S}\}$$

In the 1-D case, the open interval obtained by excluding endpoints from an interval \mathcal{I} is the interior of \mathcal{I} , denoted by $\text{int}(\mathcal{I})$. For example, $\text{int}([a, b]) = (a, b)$ and $\text{int}([0, \infty)) = (0, \infty)$.

Definition 16 [Boundary of a set]: Let $\mathcal{S} \subseteq \mathbb{R}^n$. The boundary of \mathcal{S} , denoted by $\text{bnd}(\mathcal{S})$ is defined as

$$\text{bnd}(\mathcal{S}) = \{\mathbf{y} | \forall \epsilon > 0, \mathcal{B}(\mathbf{y}, \epsilon) \cap \mathcal{S} \neq \emptyset \text{ and } \mathcal{B}(\mathbf{y}, \epsilon) \cap \mathcal{S}^C \neq \emptyset\}$$

For example, $\text{bnd}([a, b]) = \{a, b\}$.

Definition 17 [Open Set]: Let $\mathcal{S} \subseteq \mathbb{R}^n$. We say that \mathcal{S} is an open set when, for every $\mathbf{x} \in \mathcal{S}$, there exists an $\epsilon > 0$ such that $\mathcal{B}(\mathbf{x}, \epsilon) \subset \mathcal{S}$.

The simplest examples of an open set are the open ball, the empty set \emptyset and \mathbb{R}^n . Further, arbitrary union of opens sets is open. Also, finite intersection of open sets is open. The interior of any set is always open. It can be proved that a set \mathcal{S} is open if and only if $\text{int}(\mathcal{S}) = \mathcal{S}$.

The complement of an open set is the closed set.

Definition 18 [Closed Set]: Let $\mathcal{S} \subseteq \mathbb{R}^n$. We say that \mathcal{S} is a closed set when \mathcal{S}^C (that is the complement of \mathcal{S}) is an open set.

The closed ball, the empty set \emptyset and \mathbb{R}^n are three simple examples of closed sets. Arbitrary intersection of closed sets is closed. Furthermore, finite union of closed sets is closed.

Definition 19 [Closure of a Set]: Let $\mathcal{S} \subseteq \mathbb{R}^n$. The closure of \mathcal{S} , denoted by $\text{closure}(\mathcal{S})$ is given by

$$\text{closure}(\mathcal{S}) = \{\mathbf{y} \in \mathbb{R}^n | \forall \epsilon > 0, \mathcal{B}(\mathbf{y}, \epsilon) \cap \mathcal{S} \neq \emptyset\}$$

Loosely speaking, the closure of a set is the smallest closed set containing the set. The closure of a closed set is the set itself. In fact, a set \mathcal{S} is closed if and only if $\text{closure}(\mathcal{S}) = \mathcal{S}$. A bounded set can be defined in terms of a closed set; a set \mathcal{S} is bounded if and only if it is contained inside a closed set. A relationship between the interior, boundary and closure of a set \mathcal{S} is $\text{closure}(\mathcal{S}) = \text{int}(\mathcal{S}) \cup \text{bnd}(\mathcal{S})$.

4.1.3 Optimization Principles for Univariate Functions

Maximum and Minimum values of univariate functions

Let f be a function with domain \mathcal{D} . Then f has an *absolute maximum* (or global maximum) value at point $c \in \mathcal{D}$ if

$$f(x) \leq f(c), \quad \forall x \in \mathcal{D}$$

and an *absolute minimum* (or global minimum) value at $c \in \mathcal{D}$ if

$$f(x) \geq f(c), \quad \forall x \in \mathcal{D}$$

If there is an open interval \mathcal{I} containing c in which $f(c) \geq f(x)$, $\forall x \in \mathcal{I}$, then we say that $f(c)$ is a *local maximum value* of f . On the other hand, if there is an open interval \mathcal{I} containing c in which $f(c) \leq f(x)$, $\forall x \in \mathcal{I}$, then we say that $f(c)$ is a *local minimum value* of f . If $f(c)$ is either a local maximum or local minimum value of f in an open interval \mathcal{I} with $c \in \mathcal{I}$, the $f(c)$ is called a *local extreme value* of f .

The following theorem gives us the first derivative test for local extreme value of f , when f is differentiable at the extremum.

Theorem 39 *If $f(c)$ is a local extreme value and if f is differentiable at $x = c$, then $f'(c) = 0$.*

Proof: Suppose $f(c) \geq f(x)$ for all x in an open interval \mathcal{I} containing c and that $f'(c)$ exists. Then the difference quotient $\frac{f(c+h)-f(c)}{h} \leq 0$ for small $h \geq 0$ (so that $c+h \in \mathcal{I}$). This inequality remains true as $h \rightarrow 0$ from the right. In the limit, $f'(c) \leq 0$. Also, the difference quotient $\frac{f(c+h)-f(c)}{h} \geq 0$ for small $h \leq 0$ (so that $c+h \in \mathcal{I}$). This inequality remains true as $h \rightarrow 0$ from the left. In the limit, $f'(c) \geq 0$. Since $f'(c) \leq 0$ as well as $f'(c) \geq 0$, we must have $f'(c) = 0$ ¹. \square

The *extreme value theorem* is one of the most fundamental theorems in calculus concerning continuous functions on closed intervals. It can be stated as:

Theorem 40 *A continuous function $f(x)$ on a closed and bounded interval $[a, b]$ attains a minimum value $f(c)$ for some $c \in [a, b]$ and a maximum value $f(d)$ for some $d \in [a, b]$. That is, a continuous function on a closed, bounded interval attains a minimum and a maximum value.*

We must point out that either or both of the values c and d may be attained at the end points of the interval $[a, b]$. Based on theorem (39), the extreme value theorem can be extended as:

Theorem 41 *A continuous function $f(x)$ on a closed and bounded interval $[a, b]$ attains a minimum value $f(c)$ for some $c \in [a, b]$ and a maximum value $f(d)$ for some $d \in [a, b]$. If $a < c < b$ and $f'(c)$ exists, then $f'(c) = 0$. If $a < d < b$ and $f'(d)$ exists, then $f'(d) = 0$.*

¹By virtue of the *squeeze* or *sandwich theorem*

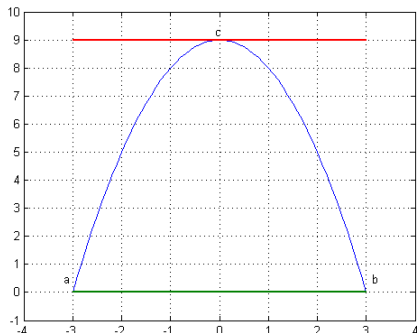


Figure 4.1: Illustration of Rolle's theorem with $f(x) = 9 - x^2$ on the interval $[-3, +3]$. We see that $f'(0) = 0$.

Next, we state the Rolle's theorem.

Theorem 42 *If f is continuous on $[a, b]$ and differentiable at all $x \in (a, b)$ and if $f(a) = f(b)$, then $f'(c) = 0$ for some $c \in (a, b)$.*

Figure 4.1 illustrates Rolle's theorem with an example function $f(x) = 9 - x^2$ on the interval $[-3, +3]$.

The *mean value theorem* is a generalization of the Rolle's theorem, though we will use the Rolle's theorem to prove it.

Theorem 43 *If f is continuous on $[a, b]$ and differentiable at all $x \in (a, b)$, then there is some $c \in (a, b)$ such that, $f'(c) = \frac{f(b) - f(a)}{b - a}$.*

Proof: Define $g(x) = f(x) - \frac{f(b) - f(a)}{b - a}(x - a)$ on $[a, b]$. We note rightaway that $g(a) = g(b)$ and $g'(x) = f'(x) - \frac{f(b) - f(a)}{b - a}$. Applying Rolle's theorem on $g(x)$, we know that there exists $c \in (a, b)$ such that $g'(c) = 0$. Which implies that $f'(c) = \frac{f(b) - f(a)}{b - a}$. \square

Figure 4.2 illustrates the mean value theorem for $f(x) = 9 - x^2$ on the interval $[-3, 0]$. We observe that the tangent at $x = -1$ is parallel to the secant joining -3 to 0 . One could think of the *mean value theorem* as a slanted version of Rolle's theorem. A natural corollary of the mean value theorem is as follows:

Corollary 44 *Let f be continuous on $[a, b]$ and differentiable on (a, b) with $m \leq f'(x) \leq M$, $\forall x \in (a, b)$. Then, $m(x - t) \leq f(x) - f(t) \leq M(x - t)$, if $a \leq t \leq x \leq b$.*

Let \mathcal{D} be the domain of function f . We define

1. the linear approximation of a differentiable function $f(x)$ as $L_a(x) = f(a) + f'(a)(x - a)$ for some $a \in \mathcal{D}$. We note that $L_a(x)$ and its first derivative at a agree with $f(a)$ and $f'(a)$ respectively.

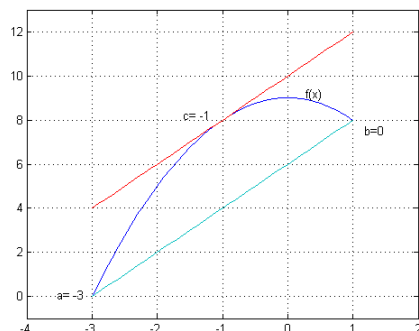


Figure 4.2: Illustration of mean value theorem with $f(x) = 9 - x^2$ on the interval $[-3, 0]$. We see that $f'(-1) = \frac{f(0) - f(-3)}{3}$.

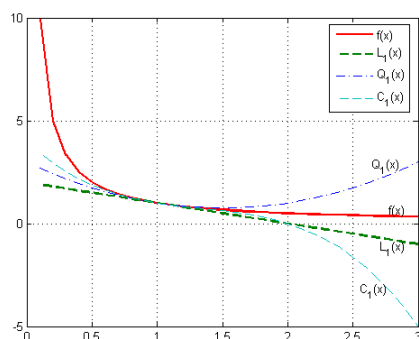


Figure 4.3: Plot of $f(x) = \frac{1}{x}$, and its linear, quadratic and cubic approximations.

2. the quadratic approximation of a twice differentiable function $f(x)$ as the parabola $Q_a(x) = f(a) + f'(a)(x - a) + \frac{1}{2}f''(a)(x - a)^2$. We note that $Q_a(x)$ and its first and second derivatives at a agree with $f(a)$, $f'(a)$ and $f''(a)$ respectively.
3. the cubic approximation of a thrice differentiable function $f(x)$ is $C_a(x) = f(a) + f'(a)(x - a) + \frac{1}{2}f''(a)(x - a)^2 + \frac{1}{6}f'''(a)(x - a)^3$. $C_a(x)$ and its first, second and third derivatives at a agree with $f(a)$, $f'(a)$, $f''(a)$ and $f'''(a)$ respectively.

The coefficient² of x^2 in $Q_a(x)$ is $\frac{1}{2}f''(a)$. Figure 4.3 illustrates the linear, quadratic and cubic approximations to the function $f(x) = \frac{1}{x}$ with $a = 1$.

²The parabola given by $Q_a(x)$ is strictly convex if $f''(a) > 0$ and is strictly concave if $f''(a) < 0$. Strict convexity for functions of single variable will be defined on page 222.

In general, an n^{th} degree polynomial approximation of a function can be found. Such an approximation will be used to prove a generalization of the mean value theorem, called the *Taylor's theorem*.

Theorem 45 *The Taylor's theorem states that if f and its first n derivatives $f', f'', \dots, f^{(n)}$ are continuous on the closed interval $[a, b]$, and differentiable on (a, b) , then there exists a number $c \in (a, b)$ such that*

$$f(b) = f(a) + f'(a)(b-a) + \frac{1}{2!}f''(a)(b-a)^2 + \dots + \frac{1}{n!}f^{(n)}(a)(b-a)^n + \frac{1}{(n+1)!}f^{(n+1)}(c)(b-a)^{n+1}$$

Proof: Define

$$p_n(x) = f(a) + f'(a)(x-a) + \frac{1}{2!}f''(a)(x-a)^2 + \dots + \frac{1}{n!}f^{(n)}(a)(x-a)^n$$

and

$$\phi_n(x) = p_n(x) + \Gamma(x-a)^{n+1}$$

The polynomials $p_n(x)$ as well as $\phi_n(x)$ and their first n derivatives match f and its first n derivatives at $x = a$. We will choose a value of Γ so that

$$f(b) = p_n(b) + \Gamma(b-a)^{n+1}$$

This requires that $\Gamma = \frac{f(b) - p_n(b)}{(b-a)^{n+1}}$. Define the function $g(x) = f(x) - \phi_n(x)$ that measures the difference between function f and the approximating function $\phi_n(x)$ for each $x \in [a, b]$.

- Since $g(a) = g(b) = 0$ and since g and g' are both continuous on $[a, b]$, we can apply the Rolle's theorem to conclude that there exists $c_1 \in [a, b]$ such that $g'(c_1) = 0$.
- Similarly, since $g'(a) = g'(c_1) = 0$, and since g' and g'' are continuous on $[a, c_1]$, we can apply the Rolle's theorem to conclude that there exists $c_2 \in [a, c_1]$ such that $g''(c_2) = 0$.
- In this way, Rolle's theorem can be applied successively to $g'', g''', \dots, g^{(n-1)}$ to imply the existence of $c_i \in (a, c_{i-1})$ such that $g^{(i)}(c_i) = 0$ for $i = 3, 4, \dots, n+1$. Note however that $g^{(n+1)}(x) = f^{(n+1)}(x) - 0 - (n+1)!\Gamma$ which gives us the value of Γ as $\frac{f^{(n+1)}(c_{n+1})}{(n+1)!}$.

Thus,

$$f(b) = f(a) + f'(a)(b-a) + \frac{1}{2!}f''(a)(b-a)^2 + \dots + \frac{1}{n!}f^{(n)}(a)(b-a)^n + \frac{f^{(n+1)}(c_{n+1})}{(n+1)!}(b-a)^{n+1}$$

□

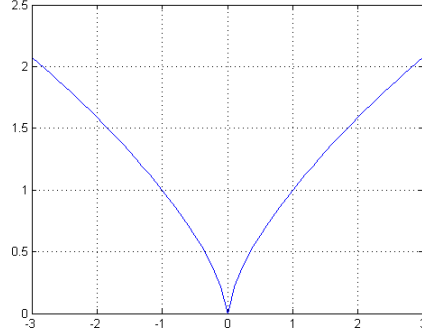


Figure 4.4: The mean value theorem can be violated if $f(x)$ is not differentiable at even a single point of the interval. Illustration on $f(x) = x^{2/3}$ with the interval $[-3, 3]$.

Note that if f fails to be differentiable at even one number in the interval, then the conclusion of the mean value theorem may be false. For example, if $f(x) = x^{2/3}$, then $f'(x) = \frac{2}{3\sqrt[3]{x}}$ and the theorem does not hold in the interval $[-3, 3]$, since f is not differentiable at 0 as can be seen in Figure 4.4.

We will introduce some definitions at this point:

- A function f is said to be *increasing* on an interval \mathcal{I} in its domain \mathcal{D} if $f(t) < f(x)$ whenever $t < x$.
- The function f is said to be *decreasing* on an interval $\mathcal{I} \in \mathcal{D}$ if $f(t) > f(x)$ whenever $t < x$.

These definitions help us derive the following theorem:

Theorem 46 *Let \mathcal{I} be an interval and suppose f is continuous on \mathcal{I} and differentiable on $\text{int}(\mathcal{I})$. Then:*

1. *if $f'(x) > 0$ for all $x \in \text{int}(\mathcal{I})$, then f is increasing on \mathcal{I} ;*
2. *if $f'(x) < 0$ for all $x \in \text{int}(\mathcal{I})$, then f is decreasing on \mathcal{I} ;*
3. *if $f'(x) = 0$ for all $x \in \text{int}(\mathcal{I})$, iff, f is constant on \mathcal{I} .*

Proof: Let $t \in \mathcal{I}$ and $x \in \mathcal{I}$ with $t < x$. By virtue of the mean value theorem, $\exists c \in (t, x)$ such that $f'(c) = \frac{f(x) - f(t)}{x - t}$.

- If $f'(x) > 0$ for all $x \in \text{int}(\mathcal{I})$, $f'(c) > 0$, which implies that $f(x) - f(t) > 0$ and we can conclude that f is increasing on \mathcal{I} .
- If $f'(x) < 0$ for all $x \in \text{int}(\mathcal{I})$, $f'(c) < 0$, which implies that $f(x) - f(t) < 0$ and we can conclude that f is decreasing on \mathcal{I} .

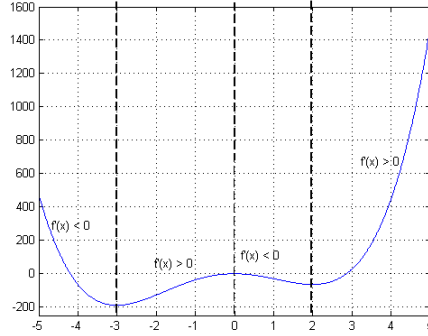


Figure 4.5: Illustration of the increasing and decreasing regions of a function $f(x) = 3x^4 + 4x^3 - 36x^2$

- If $f'(x) = 0$ for all $x \in \text{int}(\mathcal{I})$, $f'(c) = 0$, which implies that $f(x) - f(t) = 0$, and since x and t are arbitrary, we can conclude that f is constant on \mathcal{I} .

□

Figure 4.5 illustrates the intervals in $(-\infty, \infty)$ on which the function $f(x) = 3x^4 + 4x^3 - 36x^2$ is decreasing and increasing. First we note that $f(x)$ is differentiable everywhere on $(-\infty, \infty)$ and compute $f'(x) = 12(x^3 + x^2 - 6x) = 12(x - 2)(x + 3)x$, which is negative in the intervals $(-\infty, -3]$ and $[0, 2]$ and positive in the intervals $[-3, 0]$ and $[2, \infty)$. We observe that f is decreasing in the intervals $(-\infty, -3]$ and $[0, 2]$ and while it is increasing in the intervals $[-3, 0]$ and $[2, \infty)$.

There is a related sufficient condition for a function f to be increasing/decreasing on an interval \mathcal{I} , stated through the following theorem:

Theorem 47 *Let \mathcal{I} be an interval and suppose f is continuous on \mathcal{I} and differentiable on $\text{int}(\mathcal{I})$. Then:*

1. *if $f'(x) \geq 0$ for all $x \in \text{int}(\mathcal{I})$, and if $f'(x) = 0$ at only finitely many $x \in \mathcal{I}$, then f is increasing on \mathcal{I} ;*
2. *if $f'(x) \leq 0$ for all $x \in \text{int}(\mathcal{I})$, and if $f'(x) = 0$ at only finitely many $x \in \mathcal{I}$, then f is decreasing on \mathcal{I} .*

For example, the derivative of the function $f(x) = 6x^5 - 15x^4 + 10x^3$ vanishes at 0, and 1 and $f'(x) > 0$ elsewhere. So $f(x)$ is increasing on $(-\infty, \infty)$.

Are the sufficient conditions for increasing and decreasing properties of $f(x)$ in theorem 46 also necessary? It turns out that it is not the case. Figure 4.6 shows that for the function $f(x) = x^5$, though $f(x)$ is increasing in $(-\infty, \infty)$, $f'(0) = 0$.

In fact, we have a slightly different necessary condition for an increasing or decreasing function.

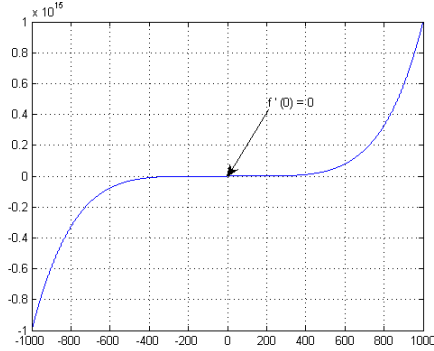


Figure 4.6: Plot of $f(x) = x^5$, illustrating that though the function is increasing on $(-\infty, \infty)$, $f'(0) = 0$.

Theorem 48 Let \mathcal{I} be an interval, and suppose f is continuous on \mathcal{I} and differentiable in $\text{int}(\mathcal{I})$. Then:

1. if f is increasing on \mathcal{I} , then $f'(x) \geq 0$ for all $x \in \text{int}(\mathcal{I})$;
2. if f is decreasing on \mathcal{I} , then $f'(x) \leq 0$ for all $x \in \text{int}(\mathcal{I})$.

Proof: Suppose f is increasing on \mathcal{I} , and let $x \in \text{int}(\mathcal{I})$. Then $\frac{f(x+h)-f(x)}{h} > 0$ for all h such that $x+h \in \text{int}(\mathcal{I})$. This implies that $f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h)-f(x)}{h} \geq 0$. For the case when f is decreasing on \mathcal{I} , it can be similarly proved that $f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h)-f(x)}{h} \leq 0$. \square

Next, we define the concept of *critical number*, which will help us derive the general condition for local extrema.

Definition 20 [Critical number]: A number c in the domain \mathcal{D} of f is called a *critical number* of f if either $f'(c) = 0$ or $f'(c)$ does not exist.

The general condition for local extrema is stated in the next theorem; it extends the result in theorem 39 to general non-differentiable functions.

Theorem 49 If $f(c)$ is a local extreme value, then c is a critical number of f .

That the converse of theorem 49 does not hold is illustrated in Figure 4.6; 0 is a critical number ($f'(0) = 0$), although $f(0)$ is not a local extreme value. Then, given a given critical number c , how do we discern whether $f(c)$ is a local extreme value? This can be answered using the *first derivative test*:

Procedure 1 [First derivative test]: Let c be an isolated critical number of f . Then,

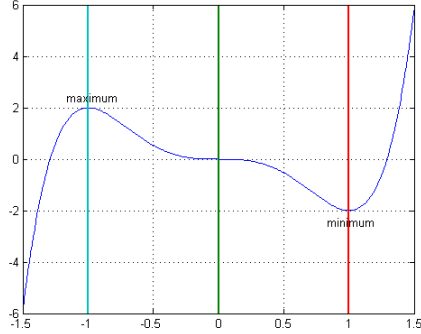


Figure 4.7: Example illustrating the derivative test for function $f(x) = 3x^5 - 5x^3$.

1. $f(c)$ is a local minimum if $f(x)$ is decreasing in an interval $[c - \epsilon_1, c]$ and increasing in an interval $[c, c + \epsilon_2]$ with $\epsilon_1, \epsilon_2 > 0$, or equivalently, the sign of $f'(x)$ changes from negative in $[c - \epsilon_1, c]$ to positive in $[c, c + \epsilon_2]$ with $\epsilon_1, \epsilon_2 > 0$.
2. $f(c)$ is a local maximum if $f(x)$ is increasing in an interval $[c - \epsilon_1, c]$ and decreasing in an interval $[c, c + \epsilon_2]$ with $\epsilon_1, \epsilon_2 > 0$, or equivalently, the sign of $f'(x)$ changes from positive in $[c - \epsilon_1, c]$ to negative in $[c, c + \epsilon_2]$ with $\epsilon_1, \epsilon_2 > 0$.
3. If $f'(x)$ is positive in an interval $[c - \epsilon_1, c]$ and also positive in an interval $[c, c - \epsilon_2]$, or $f'(x)$ is negative in an interval $[c - \epsilon_1, c]$ and also negative in an interval $[c, c - \epsilon_2]$ with $\epsilon_1, \epsilon_2 > 0$, then $f(c)$ is not a local extremum.

As an example, the function $f(x) = 3x^5 - 5x^3$ has the derivative $f'(x) = 15x^2(x+1)(x-1)$. The critical points are 0, 1 and -1 . Of the three, the sign of $f'(x)$ changes at 1 and -1 , which are local minimum and maximum respectively. The sign does not change at 0, which is therefore not a local supremum. This is pictorially depicted in Figure 4.7 As another example, consider the function

$$f(x) = \begin{cases} -x & \text{if } x \leq 0 \\ 1 & \text{if } x > 0 \end{cases}$$

Then,

$$f'(x) = \begin{cases} -1 & \text{if } x < 0 \\ 0 & \text{if } x > 0 \end{cases}$$

Note that $f(x)$ is discontinuous at $x = 0$, and therefore $f'(x)$ is not defined at $x = 0$. All numbers $x \geq 0$ are critical numbers. $f(0) = 0$ is a local minimum, whereas $f(x) = 1$ is a local minimum as well as a local maximum $\forall x > 0$.

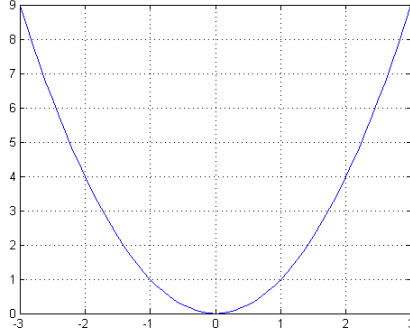


Figure 4.8: Plot for the strictly convex function $f(x) = x^2$ which has $f''(x) = 2 > 0, \forall x$.

Strict Convexity and Extremum

We define strictly convex and concave functions as follows:

1. A differentiable function f is said to be *strictly convex* (or *strictly concave up*) on an open interval \mathcal{I} , iff, $f'(x)$ is increasing on \mathcal{I} . Recall from theorem 46, the graphical interpretation of the first derivative $f'(x)$; $f'(x) > 0$ implies that $f(x)$ is increasing at x . Similarly, $f'(x)$ is increasing when $f''(x) > 0$. This gives us a sufficient condition for the strict convexity of a function:

Theorem 50 *If at all points in an open interval \mathcal{I} , $f(x)$ is doubly differentiable and if $f''(x) > 0, \forall x \in \mathcal{I}$, then the slope of the function is always increasing with x and the graph is strictly convex. This is illustrated in Figure 4.8.*

On the other hand, if the function is strictly convex and doubly differentiable in \mathcal{I} , then $f''(x) \geq 0, \forall x \in \mathcal{I}$.

There is also a slopeless interpretation of strict convexity as stated in the following theorem:

Theorem 51 *A differentiable function f is strictly convex on an open interval \mathcal{I} , iff*

$$f(ax_1 + (1-a)x_2) < af(x_1) + (1-a)f(x_2) \quad (4.2)$$

whenever $x_1, x_2 \in \mathcal{I}$, $x_1 \neq x_2$ and $0 < a < 1$.

Proof: First we will prove the necessity. Suppose f' is increasing on \mathcal{I} . Let $0 < a < 1$, $x_1, x_2 \in \mathcal{I}$ and $x_1 \neq x_2$. Without loss of generality assume that $x_1 < x_2$ ³. Then, $x_1 < ax_1 + (1-a)x_2 < x_2$ and therefore $ax_1 + (1-a)x_2 \in \mathcal{I}$. By the mean value theorem, there exist s and t with $x_1 < s < ax_1 + (1-a)x_2 < t < x_2$, such that $f(ax_1 + (1-a)x_2) - f(x_1) = f'(s)(x_2 - x_1)(1-a)$ and $f(x_2) - f(ax_1 + (1-a)x_2) = f'(t)(x_2 - x_1)a$. Therefore,

$$\begin{aligned} (1-a)f(x_1) - f(ax_1 + (1-a)x_2) + af(x_2) &= \\ a[f(x_2) - f(ax_1 + (1-a)x_2)] - (1-a)[f(ax_1 + (1-a)x_2) - f(x_1)] &= \\ a(1-a)(x_2 - x_1)[f'(t) - f'(s)] & \end{aligned}$$

Since $f(x)$ is strictly convex on \mathcal{I} , $f'(x)$ is increasing \mathcal{I} and therefore, $f'(t) - f'(s) > 0$. Moreover, $x_2 - x_1 > 0$ and $0 < a < 1$. This implies that $(1-a)f(x_1) - f(ax_1 + (1-a)x_2) + af(x_2) > 0$, or equivalently, $f(ax_1 + (1-a)x_2) < af(x_1) + (1-a)f(x_2)$, which is what we wanted to prove in 4.2.

Next, we prove the sufficiency. Suppose the inequality in 4.2 holds. Therefore,

$$\lim_{a \rightarrow 0} \frac{f(x_2 + a(x_1 - x_2)) - f(x_2)}{a} \leq f(x_1) - f(x_2)$$

that is,

$$f'(x_2)(x_1 - x_2) \leq f(x_1) - f(x_2) \quad (4.3)$$

Similarly, we can show that

$$f'(x_1)(x_2 - x_1) \leq f(x_2) - f(x_1) \quad (4.4)$$

Adding the left and right hand sides of inequalities in (4.3) and (4.4), and multiplying the resultant inequality by -1 gives us

$$(f'(x_2) - f'(x_1))(x_2 - x_1) \geq 0 \quad (4.5)$$

Using the mean value theorem, $\exists z = x_1 + t(x_2 - x_1)$ for $t \in (0, 1)$ such that

³For the case $x_2 < x_1$, the proof is very similar.

$$f(x_2) - f(x_1) = f'(z)(x_2 - x_1) \quad (4.6)$$

Since 4.5 holds for any $x_1, x_2 \in \mathcal{I}$, it also hold for $x_2 = z$. Therefore,

$$(f'(z) - f'(x_1))(x_2 - x_1) = \frac{1}{t}(f'(z) - f'(x_1))(z - x_1) \geq 0$$

Additionally using 4.6, we get

$$f(x_2) - f(x_1) = (f'(z) - f'(x_1))(x_2 - x_1) + f'(x_1)(x_2 - x_1) \geq f'(x_1)(x_2 - x_1) \quad (4.7)$$

Suppose equality holds in 4.5 for some $x_1 \neq x_2$. Then equality holds in 4.7 for the same x_1 and x_2 . That is,

$$f(x_2) - f(x_1) = f'(x_1)(x_2 - x_1) \quad (4.8)$$

Applying 4.7 we can conclude that

$$f(x_1) + af'(x_1)(x_2 - x_1) \leq f(x_1 + a(x_2 - x_1)) \quad (4.9)$$

From 4.2 and 4.8, we can derive that

$$f(x_1 + a(x_2 - x_1)) < (1 - a)f(x_1) + af(x_2) = f(x_1) + af'(x_1)(x_2 - x_1) \quad (4.10)$$

However, equations 4.9 and 4.10 contradict each other. Therefore, equality in 4.5 cannot hold for any $x_1 \neq x_2$, implying that

$$(f'(x_2) - f'(x_1))(x_2 - x_1) > 0$$

that is, $f'(x)$ is increasing and therefore f is convex on \mathcal{I} . \square

2. A differentiable function f is said to be *strictly concave* on an open interval \mathcal{I} , iff, $f'(x)$ is decreasing on \mathcal{I} . Recall from theorem 46, the graphical interpretation of the first derivative $f'(x)$; $f'(x) < 0$ implies that $f(x)$ is decreasing at x . Similarly, $f'(x)$ is monotonically decreasing when $f''(x) > 0$. This gives us a sufficient condition for the concavity of a function:

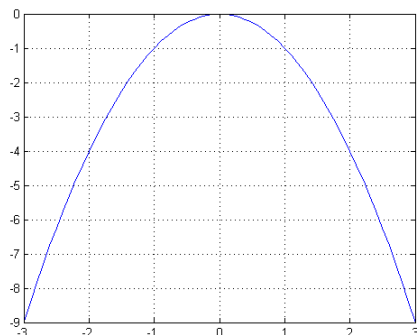


Figure 4.9: Plot for the strictly concave function $f(x) = -x^2$ which has $f''(x) = -2 < 0, \forall x$.

Theorem 52 *If at all points in an open interval \mathcal{I} , $f(x)$ is doubly differentiable and if $f''(x) < 0, \forall x \in \mathcal{I}$, then the slope of the function is always decreasing with x and the graph is strictly concave. This is illustrated in Figure 4.9.*

On the other hand, if the function is strictly concave and doubly differentiable in \mathcal{I} , then $f''(x) \leq 0, \forall x \in \mathcal{I}$.

There is also a slopeless interpretation of concavity as stated in the following theorem:

Theorem 53 *A differentiable function f is strictly concave on an open interval \mathcal{I} , iff*

$$f(ax_1 + (1-a)x_2) > af(x_1) + (1-a)f(x_2) \quad (4.11)$$

whenever $x_1, x_2 \in \mathcal{I}$, $x_1 \neq x_2$ and $0 < a < 1$.

The proof is similar to that for theorem 51.

Figure 4.10 illustrates a function $f(x) = x^3 - x + 2$, whose slope decreases as x increases to 0 ($f''(x) < 0$) and then the slope increases beyond $x = 0$ ($f''(x) > 0$). The point 0, where the $f''(x)$ changes sign is called the *inflection point*; the graph is strictly concave for $x < 0$ and strictly convex for $x > 0$. Along similar lines, we can diagnose the function $f(x) = \frac{1}{20}x^5 - \frac{7}{12}x^4 + \frac{7}{6}x^3 - \frac{15}{2}x^2$; it is strictly concave on $(-\infty, -1]$ and $[3, 5]$ and strictly convex on $[-1, 3]$ and $[5, \infty]$. The inflection points for this function are at $x = -1$, $x = 3$ and $x = 5$.

The *first derivative test* for local extrema can be restated in terms of strict convexity and concavity of functions.

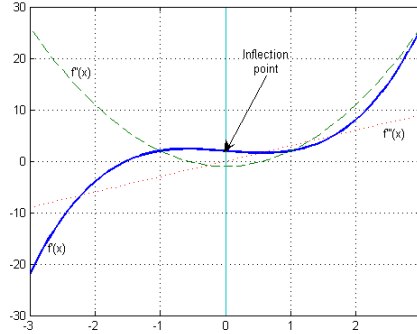


Figure 4.10: Plot for $f(x) = x^3 + x + 2$, which has an inflection point $x = 0$, along with plots for $f'(x)$ and $f''(x)$.

Procedure 2 [First derivative test in terms of strict convexity]: Let c be a critical number of f and $f'(c) = 0$. Then,

1. $f(c)$ is a local minimum if the graph of $f(x)$ is strictly convex on an open interval containing c .
2. $f(c)$ is a local maximum if the graph of $f(x)$ is strictly concave on an open interval containing c .

If the second derivative $f''(c)$ exists, then the strict convexity conditions for the critical number can be stated in terms of the sign of $f''(c)$, making use of theorems 50 and 52. This is called the *second derivative test*.

Procedure 3 [Second derivative test]: Let c be a critical number of f where $f'(c) = 0$ and $f''(c)$ exists.

1. If $f''(c) > 0$ then $f(c)$ is a local minimum.
2. If $f''(c) < 0$ then $f(c)$ is a local maximum.
3. If $f''(c) = 0$ then $f(c)$ could be a local maximum, a local minimum, neither or both. That is, the test fails.

For example,

- If $f(x) = x^4$, then $f'(0) = 0$ and $f''(0) = 0$ and we can see that $f(0)$ is a local minimum.
- If $f(x) = -x^4$, then $f'(0) = 0$ and $f''(0) = 0$ and we can see that $f(0)$ is a local maximum.
- If $f(x) = x^3$, then $f'(0) = 0$ and $f''(0) = 0$ and we can see that $f(0)$ is neither a local minimum nor a local maximum. $(0, 0)$ is an inflection point in this case.

- If $f(x) = x + 2 \sin x$, then $f'(x) = 1 + 2 \cos x$. $f'(x) = 0$ for $x = \frac{2\pi}{3}, \frac{4\pi}{3}$, which are the critical numbers. $f''\left(\frac{2\pi}{3}\right) = 2 \sin \frac{2\pi}{3} = -\sqrt{3} < 0 \Rightarrow f\left(\frac{2\pi}{3}\right) = \frac{2\pi}{3} + \sqrt{3}$ is a local maximum value. On the other hand, $f''\left(\frac{4\pi}{3}\right) = \sqrt{3} > 0 \Rightarrow f\left(\frac{4\pi}{3}\right) = \frac{4\pi}{3} - \sqrt{3}$ is a local minimum value.
- If $f(x) = x + \frac{1}{x}$, then $f'(x) = 1 + \frac{1}{x^2}$. The critical numbers are $x = \pm 1$. Note that $x = 0$ is not a critical number, even though $f'(0)$ does not exist, because 0 is not in the domain of f . $f''(x) = \frac{2}{x^3}$. $f''(-1) = -2 < 0$ and therefore $f(-1) = -2$ is a local maximum. $f''(1) = 2 > 0$ and therefore $f(1) = 2$ is a local minimum.

Global Extrema on Closed Intervals

Recall the extreme value theorem (theorem 40). An outcome of the extreme value theorem is that

- if either of c or d lies in (a, b) , then it is a critical number of f ;
- else each of c and d must lie on one of the boundaries of $[a, b]$.

This gives us a procedure for finding the maximum and minimum of a continuous function f on a closed bounded interval \mathcal{I} :

Procedure 4 [Finding extreme values on closed, bounded intervals]: 1.

Find the critical points in $\text{int}(\mathcal{I})$.

2. *Compute the values of f at the critical points and at the endpoints of the interval.*
3. *Select the least and greatest of the computed values.*

For example, to compute the maximum and minimum values of $f(x) = 4x^3 - 8x^2 + 5x$ on the interval $[0, 1]$, we first compute $f'(x) = 12x^2 - 16x + 5$ which is 0 at $x = \frac{1}{2}, \frac{5}{6}$. Values at the critical points are $f\left(\frac{1}{2}\right) = 1$, $f\left(\frac{5}{6}\right) = \frac{25}{27}$. The values at the end points are $f(0) = 0$ and $f(1) = 1$. Therefore, the minimum value is $f(0) = 0$ and the maximum value is $f(1) = f\left(\frac{1}{2}\right) = 1$.

In this context, it is relevant to discuss the one-sided derivatives of a function at the endpoints of the closed interval on which it is defined.

Definition 21 [One-sided derivatives at endpoints]: *Let f be defined on a closed bounded interval $[a, b]$. The (right-sided) derivative of f at $x = a$ is defined as*

$$f'(a) = \lim_{h \rightarrow 0^+} \frac{f(a+h) - f(a)}{h}$$

Similarly, the (left-sided) derivative of f at $x = b$ is defined as

$$f'(b) = \lim_{h \rightarrow 0^-} \frac{f(b+h) - f(b)}{h}$$

Essentially, each of the one-sided derivatives defines one-sided slopes at the endpoints. Based on these definitions, the following result can be derived.

Theorem 54 *If f is continuous on $[a, b]$ and $f'(a)$ exists as a real number or as $\pm\infty$, then we have the following necessary conditions for extremum at a .*

- *If $f(a)$ is the maximum value of f on $[a, b]$, then $f'(a) \leq 0$ or $f'(a) = -\infty$.*
- *If $f(a)$ is the minimum value of f on $[a, b]$, then $f'(a) \geq 0$ or $f'(a) = \infty$.*

If f is continuous on $[a, b]$ and $f'(b)$ exists as a real number or as $\pm\infty$, then we have the following necessary conditions for extremum at b .

- *If $f(b)$ is the maximum value of f on $[a, b]$, then $f'(b) \geq 0$ or $f'(b) = \infty$.*
- *If $f(b)$ is the minimum value of f on $[a, b]$, then $f'(b) \leq 0$ or $f'(b) = -\infty$.*

The following theorem gives a useful procedure for finding extrema on closed intervals.

Theorem 55 *If f is continuous on $[a, b]$ and $f''(x)$ exists for all $x \in (a, b)$. Then,*

- *If $f''(x) \leq 0, \forall x \in (a, b)$, then the minimum value of f on $[a, b]$ is either $f(a)$ or $f(b)$. If, in addition, f has a critical number $c \in (a, b)$, then $f(c)$ is the maximum value of f on $[a, b]$.*
- *If $f''(x) \geq 0, \forall x \in (a, b)$, then the maximum value of f on $[a, b]$ is either $f(a)$ or $f(b)$. If, in addition, f has a critical number $c \in (a, b)$, then $f(c)$ is the minimum value of f on $[a, b]$.*

The next theorem is very useful for finding global extrema values on open intervals.

Theorem 56 *Let \mathcal{I} be an open interval and let $f''(x)$ exist $\forall x \in \mathcal{I}$.*

- *If $f''(x) \geq 0, \forall x \in \mathcal{I}$, and if there is a number $c \in \mathcal{I}$ where $f'(c) = 0$, then $f(c)$ is the global minimum value of f on \mathcal{I} .*
- *If $f''(x) \leq 0, \forall x \in \mathcal{I}$, and if there is a number $c \in \mathcal{I}$ where $f'(c) = 0$, then $f(c)$ is the global maximum value of f on \mathcal{I} .*

For example, let $f(x) = \frac{2}{3}x - \sec x$ and $\mathcal{I} = (-\frac{\pi}{2}, \frac{\pi}{2})$. $f'(x) = \frac{2}{3} - \sec x \tan x = \frac{2}{3} - \frac{\sin x}{\cos^2 x} = 0 \Rightarrow x = \frac{\pi}{6}$. Further, $f''(x) = -\sec x(\tan^2 x + \sec^2 x) < 0$ on $(-\frac{\pi}{2}, \frac{\pi}{2})$. Therefore, f attains the maximum value $f(\frac{\pi}{6}) = \frac{\pi}{9} - \frac{2}{\sqrt{3}}$ on \mathcal{I} .

As another example, let us find the dimensions of the cone with minimum volume that can contain a sphere with radius R . Let h be the height of the cone and r the radius of its base. The objective to be minimized is the volume $f(r, h) = \frac{1}{3}\pi r^2 h$. The constraint between r and h is shown in Figure 4.11; the triangle AEF is similar to triangle ADB and therefore, $\frac{h-R}{R} = \frac{\sqrt{h^2+r^2}}{r}$. Our

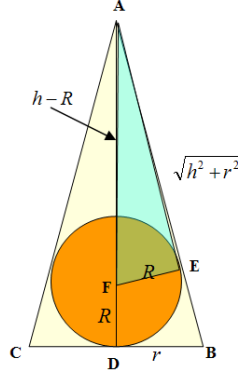


Figure 4.11: Illustrating the constraints for the optimization problem of finding the cone with minimum volume that can contain a sphere of radius R .

first step is to reduce the volume formula to involve only one of r^{2^4} or h . The algebra involved will be the simplest if we solved for h . The constraint gives us $r^2 = \frac{R^2 h}{h-2R}$. Substituting this expression for r^2 into the volume formula, we get $g(h) = \frac{\pi R^2}{3} \frac{h^2}{(h-2R)}$ with the domain given by $\mathcal{D} = \{h | 2R < h < \infty\}$. Note that \mathcal{D} is an open interval. $g' = \frac{\pi R^2}{3} \frac{2h(h-2R)-h^2}{(h-2R)^2} = \frac{\pi R^2}{3} \frac{h(h-4R)}{(h-2R)^2}$ which is 0 in its domain \mathcal{D} if and only if $h = 4R$. $g'' = \frac{\pi R^2}{3} \frac{2(h-2R)^3 - 2h(h-4R)(h-2R)^2}{(h-2R)^4} = \frac{\pi R^2}{3} \frac{2(h^2 - 4Rh + 4R^2 - h^2 + 4Rh)}{(h-2R)^3} = \frac{\pi R^2}{3} \frac{8R^2}{(h-2R)^3}$, which is greater than 0 in \mathcal{D} . Therefore, g (and consequently f) has a unique minimum at $h = 4R$ and correspondingly, $r^2 = \frac{R^2 h}{h-2R} = 2R^2$.

4.1.4 Optimization Principles for Multivariate Functions

Directional derivative and the gradient vector

Consider a function $f(\mathbf{x})$, with $\mathbf{x} \in \mathfrak{R}^n$. We start with the concept of the the direction at a point $\mathbf{x} \in \mathfrak{R}^n$. We will represent a vector by \mathbf{x} and the k^{th} component of \mathbf{x} by x_k . Let \mathbf{u}^k be a unit vector pointing along the k^{th} coordinate axis in \mathfrak{R}^n ; $u_k^k = 1$ and $u_j^k = 0$, $\forall j \neq k$. An arbitrary direction vector \mathbf{v} at \mathbf{x} is a vector in \mathfrak{R}^n with unit norm (i.e., $\|\mathbf{v}\| = 1$) and component v_k in the direction of \mathbf{u}^k . Let $f : \mathcal{D} \rightarrow \mathfrak{R}$, $\mathcal{D} \subseteq \mathfrak{R}^n$ be a function.

Definition 22 [Directional derivative]: *The directional derivative of $f(\mathbf{x})$ at \mathbf{x} in the direction of the unit vector \mathbf{v} is*

$$D_{\mathbf{v}} f(\mathbf{x}) = \lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\mathbf{v}) - f(\mathbf{x})}{h} \quad (4.12)$$

⁴Since r appears in the volume formula only in terms of r^2 .

provided the limit exists.

As a special case, when $\mathbf{v} = \mathbf{u}^k$ the directional derivative reduces to the partial derivative of f with respect to x_k .

$$D_{\mathbf{u}^k} f(\mathbf{x}) = \frac{\partial f(\mathbf{x})}{\partial x_k}$$

Theorem 57 *If $f(\mathbf{x})$ is a differentiable function of $\mathbf{x} \in \mathfrak{R}^n$, then f has a directional derivative in the direction of any unit vector \mathbf{v} , and*

$$D_{\mathbf{v}} f(\mathbf{x}) = \sum_{k=1}^n \frac{\partial f(\mathbf{x})}{\partial x_k} v_k \quad (4.13)$$

Proof: Define $g(h) = f(\mathbf{x} + \mathbf{v}h)$. Now:

- $g'(0) = \lim_{h \rightarrow 0} \frac{g(0+h) - g(0)}{h} = \lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\mathbf{v}) - f(\mathbf{x})}{h}$, which is the expression for the directional derivative defined in equation 4.12. Thus, $g'(0) = D_{\mathbf{v}} f(\mathbf{x})$.
- By definition of the chain rule for partial differentiation, we get another expression for $g'(0)$; $g'(0) = \sum_{k=1}^n \frac{\partial f(\mathbf{x})}{\partial x_k} v_k$

Therefore, $g'(0) = D_{\mathbf{v}} f(\mathbf{x}) = \sum_{k=1}^n \frac{\partial f(\mathbf{x})}{\partial x_k} v_k \quad \square$

The theorem works if the function is differentiable at the point, else it is not predictable. The above theorem leads us directly to the idea of the gradient. We can see that the right hand side of (4.13) can be realized as the dot product of two vectors, viz., $\left[\frac{\partial f(\mathbf{x})}{\partial x_1}, \frac{\partial f(\mathbf{x})}{\partial x_2}, \dots, \frac{\partial f(\mathbf{x})}{\partial x_n} \right]^T$ and \mathbf{v} . Let us denote $\frac{\partial f(\mathbf{x})}{\partial x_i}$ by $f_{x_i}(\mathbf{x})$. Then we assign a name to the special vector discovered above.

Definition 23 [Gradient Vector]: *If f is differentiable function of $\mathbf{x} \in \mathfrak{R}^n$, then the gradient of $f(\mathbf{x})$ is the vector function $\nabla f(\mathbf{x})$, defined as:*

$$\nabla f(\mathbf{x}) = [f_{x_1}(\mathbf{x}), f_{x_2}(\mathbf{x}), \dots, f_{x_n}(\mathbf{x})]$$

The directional derivative of a function f at a point \mathbf{x} in the direction of a unit vector \mathbf{v} can be now written as

$$D_{\mathbf{v}} f(\mathbf{x}) = \nabla^T f(\mathbf{x}) \cdot \mathbf{v} \quad (4.14)$$

What does the gradient $\nabla f(\mathbf{x})$ tell you about the function $f(\mathbf{x})$? We will illustrate with some examples. Consider the polynomial $f(x, y, z) = x^2y + z \sin xy$ and the unit vector $\mathbf{v}^T = \frac{1}{\sqrt{3}}[1, 1, 1]^T$. Consider the point $p_0 = (0, 1, 3)$. We will compute the directional derivative of f at p_0 in the direction of \mathbf{v} . To do this, we first compute the gradient of f in general: $\nabla f = [2xy + yz \cos xy, x^2 + xz \cos xy, \sin xy]^T$. Evaluating the gradient at a specific point p_0 , $\nabla f(0, 1, 3) = [3, 0, 0]^T$. The directional derivative at p_0 in the direction \mathbf{v} is $D_{\mathbf{v}}f(0, 1, 3) = [3, 0, 0] \cdot \frac{1}{\sqrt{3}}[1, 1, 1]^T = \sqrt{3}$. This directional derivative is the rate of change of f at p_0 in the direction \mathbf{v} ; it is positive indicating that the function f increases at p_0 in the direction \mathbf{v} . All our ideas about first and second derivative in the case of a single variable carry over to the directional derivative.

As another example, let us find the rate of change of $f(x, y, z) = e^{xyz}$ at $p_0 = (1, 2, 3)$ in the direction from $p_1 = (1, 2, 3)$ to $p_2 = (-4, 6, -1)$. We first construct a unit vector from p_1 to p_2 ; $\mathbf{v} = \frac{1}{\sqrt{57}}[-5, 4, -4]$. The gradient of f in general is $\nabla f = [yze^{xyz}, xze^{xyz}, xye^{xyz}] = e^{xyz}[yz, xz, xy]$. Evaluating the gradient at a specific point p_0 , $\nabla f(1, 2, 3) = e^6[6, 3, 2]^T$. The directional derivative at p_0 in the direction \mathbf{v} is $D_{\mathbf{v}}f(1, 2, 3) = e^6[6, 3, 2] \cdot \frac{1}{\sqrt{57}}[-5, 4, -4]^T = e^6 \frac{-26}{\sqrt{57}}$. This directional derivative is negative, indicating that the function f decreases at p_0 in the direction from p_1 to p_2 .

While there exist infinitely many direction vectors \mathbf{v} at any point \mathbf{x} , there is a unique gradient vector $\nabla f(\mathbf{x})$. Since we separated $D_{\mathbf{v}}f(\mathbf{x})$ as the dot product of $\nabla f(\mathbf{x})$ with \mathbf{v} , we can study $\nabla f(\mathbf{x})$ independently. What does the gradient vector tell us? We will state a theorem to answer this question.

Theorem 58 *Suppose f is a differentiable function of $\mathbf{x} \in \mathbb{R}^n$. The maximum value of the directional derivative $D_{\mathbf{v}}f(\mathbf{x})$ is $\|\nabla f(\mathbf{x})\|$ and it is so when \mathbf{v} has the same direction as the gradient vector $\nabla f(\mathbf{x})$.*

Proof: The *cauchy schwartz inequality* when applied in the euclidian space states that $|\mathbf{x}^T \cdot \mathbf{y}| \leq \|\mathbf{x}\| \cdot \|\mathbf{y}\|$ for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, with equality holding *iff* \mathbf{x} and \mathbf{y} are linearly dependent. The inequality gives upper and lower bounds on the dot product between two vectors; $-\|\mathbf{x}\| \cdot \|\mathbf{y}\| \leq \mathbf{x}^T \cdot \mathbf{y} \leq \|\mathbf{x}\| \cdot \|\mathbf{y}\|$. Applying these bounds to the right hand side of 4.14 and using the fact that $\|\mathbf{v}\| = 1$, we get

$$-\|\nabla f(\mathbf{x})\| \leq D_{\mathbf{v}}f(\mathbf{x}) = \nabla^T f(\mathbf{x}) \cdot \mathbf{v} \leq \|\nabla f(\mathbf{x})\|$$

with equality holding *iff* $\mathbf{v} = k\nabla f(\mathbf{x})$ for some $k \geq 0$. Since $\|\mathbf{v}\| = 1$, equality can hold *iff* $\mathbf{v} = \frac{\nabla f(\mathbf{x})}{\|\nabla f(\mathbf{x})\|}$. \square

The theorem implies that the maximum rate of change of f at a point \mathbf{x} is given by the norm of the gradient vector at \mathbf{x} . And the direction in which the rate of change of f is maximum is given by the unit vector $\frac{\nabla f(\mathbf{x})}{\|\nabla f(\mathbf{x})\|}$.

An associated fact is that the minimum value of the directional derivative $D_{\mathbf{v}}f(\mathbf{x})$ is $-\|\nabla f(\mathbf{x})\|$ and it occurs when \mathbf{v} has the opposite direction of the gradient vector, *i.e.*, $-\frac{\nabla f(\mathbf{x})}{\|\nabla f(\mathbf{x})\|}$. This fact is often used in numerical analysis when one is trying to minimize the value of very complex functions. The method

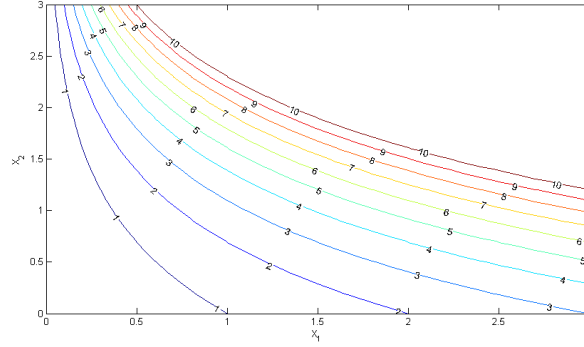


Figure 4.12: 10 level curves for the function $f(x_1, x_2) = x_1 e^{x_2}$.

of steepest descent uses this result to iteratively choose a new value of \mathbf{x} by traversing in the direction of $-\nabla f(\mathbf{x})$.

Consider the function $f(x_1, x_2) = x_1 e^{x_2}$. Figure 4.12 shows 10 level curves for this function, corresponding to $f(x_1, x_2) = c$ for $c = 1, 2, \dots, 10$. The idea behind a level curve is that as you change \mathbf{x} along any level curve, the function value remains unchanged, but as you move \mathbf{x} across level curves, the function value changes.

We will define the concept of a hyperplane next, since it will be repeatedly referred to in the sequel.

Definition 24 [Hyperplane]: A set of points $\mathcal{H} \subseteq \mathbb{R}^n$ is called a hyperplane if there exists a vector $\mathbf{v} \in \mathbb{R}^n$ and a point $\mathbf{q} \in \mathbb{R}^n$ such that

$$\forall \mathbf{p} \in \mathcal{H}, (\mathbf{p} - \mathbf{q})^T \mathbf{v} = 0$$

or in other words, $\forall \mathbf{p} \in \mathcal{H}, \mathbf{p}^T \mathbf{v} = \mathbf{q}^T \mathbf{v}$. This is the equation of a hyperplane orthogonal to vector \mathbf{v} and passing through point \mathbf{q} . The space spanned by vectors in the hyperplane \mathcal{H} which are orthogonal to vector \mathbf{v} , forms the orthogonal complement of the space spanned by \mathbf{v} .

Hyperplane \mathcal{H} can also be equivalently defined as the set of points \mathbf{p} such that $\mathbf{p}^T \mathbf{v} = c$ for some $c \in \mathbb{R}$ and some $\mathbf{v} \in \mathbb{R}^n$, with $c = \mathbf{q}^T \mathbf{v}$ in our definition. (This definition will be referred to at a later point.)

What if $D_{\mathbf{v}} f(\mathbf{x})$ turns out to be 0? What can we say about $\nabla f(\mathbf{x})$ and \mathbf{v} ? There is a useful theorem in this regard.

Theorem 59 Let $f : \mathcal{D} \rightarrow \mathbb{R}$ with $\mathcal{D} \subseteq \mathbb{R}^n$ be a differentiable function. The gradient ∇f evaluated at \mathbf{x}^* is orthogonal to the tangent hyperplane (tangent line in case $n = 2$) to the level surface of f passing through \mathbf{x}^* .

Proof: Let \mathcal{K} be the range of f and let $k \in \mathcal{K}$ such that $f(\mathbf{x}^*) = k$. Consider the level surface $f(\mathbf{x}) = k$. Let $\mathbf{r}(t) = [x_1(t), x_2(t), \dots, x_n(t)]$ be a curve on the level surface, parametrized by $t \in \mathfrak{R}$, with $\mathbf{r}(0) = \mathbf{x}^*$. Then, $f(x(t), y(t), z(t)) = k$. Applying the chain rule

$$\frac{df(\mathbf{r}(t))}{dt} = \sum_{i=1}^n \frac{\partial f}{\partial x_i} \frac{dx_i(t)}{dt} = \nabla^T f(\mathbf{x}(t)) \frac{d\mathbf{r}(t)}{dt} = 0$$

For $t = 0$, the equations become

$$\nabla^T f(\mathbf{x}^*) \frac{d\mathbf{r}(0)}{dt} = 0$$

Now, $\frac{d\mathbf{r}(t)}{dt}$ represents any tangent vector to the curve through $\mathbf{r}(t)$ which lies completely on the level surface. That is, the tangent line to any curve at \mathbf{x}^* on the level surface containing \mathbf{x}^* , is orthogonal to $\nabla f(\mathbf{x}^*)$. Since the tangent hyperplane to a surface at any point is the hyperplane containing all tangent vectors to curves on the surface passing through the point, the gradient is perpendicular to the tangent hyperplane to the level surface passing through that point. The equation of the tangent hyperplane is given by $(\mathbf{x} - \mathbf{x}^*)^T \nabla f(\mathbf{x}^*) = 0$. \square

Recall from elementary calculus, that the normal to a plane can be found by taking the cross product of any two vectors lying within the plane. The gradient vector at any point on the level surface of a function is normal to the tangent hyperplane (or tangent line in the case of two variables) to the surface at the same point, but can however be conveniently obtained using the partial derivatives of the function at that point.

We will use some illustrative examples to study these facts.

1. Consider the same plot as in Figure 4.12 with a gradient vector at $(2, 0)$ as shown in Figure 4.13. The gradient vector $[1, 2]^T$ is perpendicular to the tangent hyperplane to the level curve $x_1 e^{x_2} = 2$ at $(2, 0)$. The equation of the tangent hyperplane is $(x_1 - 2) + 2(x_2 - 0) = 0$ and it turns out to be a tangent line.
2. The level surfaces for $f(x_1, x_2, x_3) = x_1^2 + x_2^2 + x_3^2$ are shown in Figure 4.14. The gradient at $(1, 1, 1)$ is orthogonal to the tangent hyperplane to the level surface $f(x_1, x_2, x_3) = x_1^2 + x_2^2 + x_3^2 = 3$ at $(1, 1, 1)$. The gradient vector at $(1, 1, 1)$ is $[2, 2, 2]^T$ and the tangent hyperplane has the equation $2(x_1 - 1) + 2(x_2 - 1) + 2(x_3 - 1) = 0$, which is a plane in $3D$. On the other hand, the dotted line in Figure 4.15 is not orthogonal to the level surface, since it does not coincide with the gradient.
3. Let $f(x_1, x_2, x_3) = x_1^2 x_2^3 x_3^4$ and consider the point $\mathbf{x}^0 = (1, 2, 1)$. We will find the equation of the tangent plane to the level surface through \mathbf{x}^0 . The level surface through \mathbf{x}^0 is determined by setting f equal to its value evaluated at \mathbf{x}^0 ; that is, the level surface will have the equation $x_1^2 x_2^3 x_3^4 = 1^2 2^3 1^4 = 8$. The gradient vector (normal to tangent plane) at

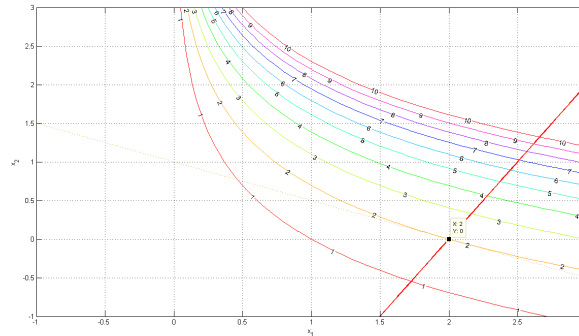


Figure 4.13: The level curves from Figure 4.12 along with the gradient vector at $(2, 0)$. Note that the gradient vector is perpendicular to the level curve $x_1 e^{x_2} = 2$ at $(2, 0)$.

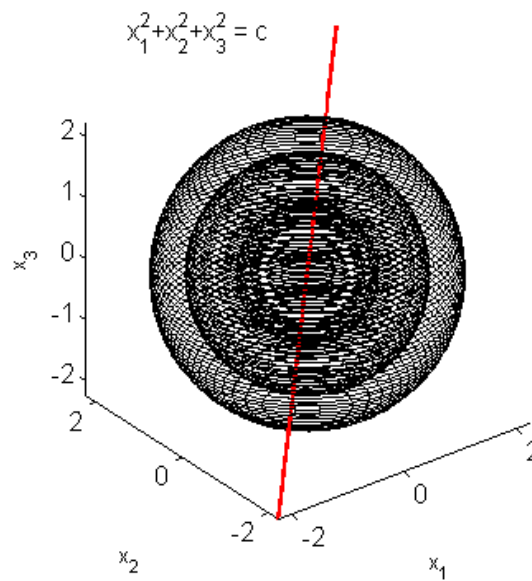


Figure 4.14: 3 level surfaces for the function $f(x_1, x_2, x_3) = x_1^2 + x_2^2 + x_3^2$ with $c = 1, 3, 5$. The gradient at $(1, 1, 1)$ is orthogonal to the level surface $f(x_1, x_2, x_3) = x_1^2 + x_2^2 + x_3^2 = 3$ at $(1, 1, 1)$.

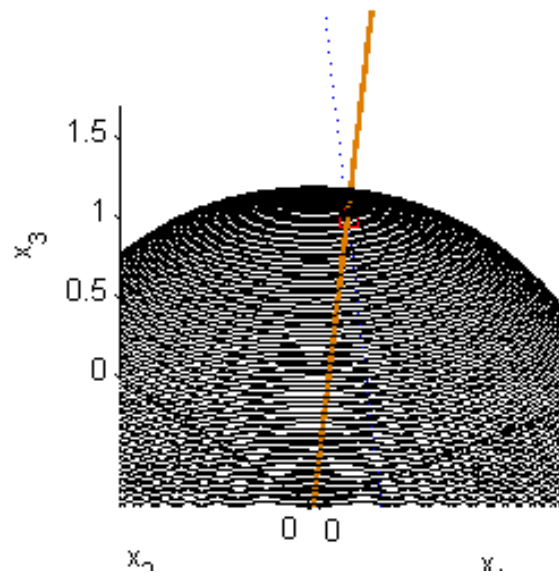


Figure 4.15: Level surface $f(x_1, x_2, x_3) = x_1^2 + x_2^2 + x_3^2 = 3$. The gradient at $(1, 1, 1)$, drawn as a bold line, is perpendicular to the tangent plane to the level surface at $(1, 1, 1)$, whereas, the dotted line, though passing through $(1, 1, 1)$ is not perpendicular to the same tangent plane.

$(1, 2, 1)$ is $\nabla f(x_1, x_2, x_3)|_{(1,2,1)} = [2x_1x_2^3x_3^4, 3x_1^2x_2^2x_3^4, 4x_1^2x_2^3x_3^3]^T|_{(1,2,1)} = [16, 12, 32]^T$. The equation of the tangent plane at \mathbf{x}^0 , given the normal vector $\nabla f(\mathbf{x}^0)$ can be easily written down: $\nabla f(\mathbf{x}^0)^T \cdot [\mathbf{x} - \mathbf{x}^0] = 0$ which turns out to be $16(x_1 - 1) + 12(x_2 - 2) + 32(x_3 - 1) = 0$, a plane in $3D$.

4. Consider the function $f(x, y, z) = \frac{x}{y+z}$. The directional derivative of f in the direction of the vector $\mathbf{v} = \frac{1}{\sqrt{14}}[1, 2, 3]$ at the point $\mathbf{x}^0 = (4, 1, 1)$ is $\nabla^T f|_{(4,1,1)} \cdot \frac{1}{\sqrt{14}}[1, 2, 3]^T = \left[\frac{1}{y+z}, -\frac{x}{(y+z)^2}, -\frac{x}{(y+z)^2} \right]|_{(4,1,1)} \cdot \frac{1}{\sqrt{14}}[1, 2, 3]^T = \left[\frac{1}{2}, -1, -1 \right] \cdot \frac{1}{\sqrt{14}}[1, 2, 3]^T = -\frac{9}{2\sqrt{14}}$. The directional derivative is negative, indicating that the function decreases along the direction of \mathbf{v} . Based on theorem 58, we know that the maximum rate of change of a function at a point \mathbf{x} is given by $\|\nabla f(\mathbf{x})\|$ and it is in the direction $\frac{\nabla f(\mathbf{x})}{\|\nabla f(\mathbf{x})\|}$. In the example under consideration, this maximum rate of change at \mathbf{x}^0 is $\frac{3}{2}$ and it is in the direction of the vector $\frac{2}{3} \left[\frac{1}{2}, -1, -1 \right]$.
5. Let us find the maximum rate of change of the function $f(x, y, z) = x^2y^3z^4$ at the point $\mathbf{x}^0 = (1, 1, 1)$ and the direction in which it occurs. The gradient at \mathbf{x}^0 is $\nabla^T f|_{(1,1,1)} = [2, 3, 4]$. The maximum rate of change at \mathbf{x}^0 is therefore $\sqrt{29}$ and the direction of the corresponding rate of change is $\frac{1}{\sqrt{29}}[2, 3, 4]$. The minimum rate of change is $-\sqrt{29}$ and the corresponding direction is $-\frac{1}{\sqrt{29}}[2, 3, 4]$.
6. Let us determine the equations of (a) the tangent plane to the paraboloid $\mathcal{P} : x_1 = x_2^2 + x_3^2 + 2$ at $(-1, 1, 0)$ and (b) the normal line to the tangent plane. To realize this as the level surface of a function of three variables, we define the function $f(x_1, x_2, x_3) = x_1 - x_2^2 - x_3^2$ and find that the paraboloid \mathcal{P} is the same as the level surface $f(x_1, x_2, x_3) = -2$. The normal to the tangent plane to \mathcal{P} at \mathbf{x}^0 is in the direction of the gradient vector $\nabla f(\mathbf{x}^0) = [1, -2, 0]^T$ and its parametric equation is $[x_1, x_2, x_3] = [-1+t, 1-2t, 0]$. The equation of the tangent plane is therefore $(x_1 + 1) - 2(x_2 - 1) = 0$.

We can embed the graph of a function of n variables as the 0-level surface of a function of $n + 1$ variables. More concretely, if $f : \mathcal{D} \rightarrow \mathfrak{R}$, $\mathcal{D} \subseteq \mathfrak{R}^n$ then we define $F : \mathcal{D}' \rightarrow \mathfrak{R}$, $\mathcal{D}' = \mathcal{D} \times \mathfrak{R}$ as $F(\mathbf{x}, z) = f(\mathbf{x}) - z$ with $\mathbf{x} \in \mathcal{D}'$. The function f then corresponds to a single level surface of F given by $F(\mathbf{x}, z) = 0$. In other words, the 0-level surface of F gives back the graph of f . The gradient of F at any point (\mathbf{x}, z) is simply, $\nabla F(\mathbf{x}, z) = [f_{x_1}, f_{x_2}, \dots, f_{x_n}, -1]$ with the first n components of $\nabla F(\mathbf{x}, z)$ given by the n components of $\nabla f(\mathbf{x})$. We note that the level surface of F passing through point $(\mathbf{x}^0, f(\mathbf{x}^0))$ is its 0-level surface, which is essentially the surface of the function $f(\mathbf{x})$. The equation of the tangent hyperplane to the 0-level surface of F at the point $(\mathbf{x}^0, f(\mathbf{x}^0))$ (that is, the tangent hyperplane to $f(\mathbf{x})$ at the point \mathbf{x}_0), is $\nabla F(\mathbf{x}^0, f(\mathbf{x}^0))^T \cdot [\mathbf{x} - \mathbf{x}^0, z - f(\mathbf{x}^0)]^T = 0$. Substituting appropriate expression for $\nabla F(\mathbf{x}^0)$, the equation of the tangent plane can be written as

$$\left(\sum_{i=1}^n f_{x_i}(\mathbf{x}^0)(x_i - x_i^0) \right) - (z - f(\mathbf{x}^0)) = 0$$

or equivalently as,

$$\left(\sum_{i=1}^n f_{x_i}(\mathbf{x}^0)(x_i - x_i^0) \right) + f(\mathbf{x}^0) = z$$

As an example, consider the paraboloid, $f(x_1, x_2) = 9 - x_1^2 - x_2^2$, the corresponding $F(x_1, x_2, z) = 9 - x_1^2 - x_2^2 - z$ and the point $x^0 = (\mathbf{x}^0, z) = (1, 1, 7)$ which lies on the 0-level surface of F . The gradient $\nabla F(x_1, x_2, z)$ is $[-2x_1, -2x_2, -1]$, which when evaluated at $x^0 = (1, 1, 7)$ is $[-2, -2, -1]$. The equation of the tangent plane to f at x^0 is therefore given by $-2(x_1 - 1) - 2(x_2 - 1) + 7 = z$.

Recall from theorem 39 that for functions of single variable, at local extreme points, the tangent to the curve is a line with a constant component in the direction of the function and is therefore parallel to the x -axis. If the function is differentiable at the extreme point, then the derivative must vanish. This idea can be extended to functions of multiple variables. The requirement in this case turns out to be that the tangent plane to the function at any extreme point must be parallel to the plane $z = 0$. This can happen if and only if the gradient ∇F is parallel to the z -axis at the extreme point, or equivalently, the gradient to the function f must be the zero vector at every extreme point.

We will formalize this discussion by first providing the definitions for local maximum and minimum as well as absolute maximum and minimum values of a function of n variables.

Definition 25 [Local maximum]: A function f of n variables has a local maximum at \mathbf{x}^0 if $\exists \epsilon > 0$ such that $\forall \|\mathbf{x} - \mathbf{x}^0\| < \epsilon$. $f(\mathbf{x}) \leq f(\mathbf{x}^0)$. In other words, $f(\mathbf{x}) \leq f(\mathbf{x}^0)$ whenever \mathbf{x} lies in some circular disk around \mathbf{x}^0 .

Definition 26 [Local minimum]: A function f of n variables has a local minimum at \mathbf{x}^0 if $\exists \epsilon > 0$ such that $\forall \|\mathbf{x} - \mathbf{x}^0\| < \epsilon$. $f(\mathbf{x}) \geq f(\mathbf{x}^0)$. In other words, $f(\mathbf{x}) \geq f(\mathbf{x}^0)$ whenever \mathbf{x} lies in some circular disk around \mathbf{x}^0 .

These definitions are exactly analogous to the definitions for a function of single variable. Figure 4.16 shows the plot of $f(x_1, x_2) = 3x_1^2 - x_1^3 - 2x_2^2 + x_2^4$. As can be seen in the plot, the function has several local maxima and minima.

We will next state a theorem fundamental to determining the locally extreme values of functions of multiple variables.

Theorem 60 If $f(\mathbf{x})$ defined on a domain $\mathcal{D} \subseteq \mathbb{R}^n$ has a local maximum or minimum at \mathbf{x}^* and if the first-order partial derivatives exist at \mathbf{x}^* , then $f_{x_i}(\mathbf{x}^*) = 0$ for all $1 \leq i \leq n$.

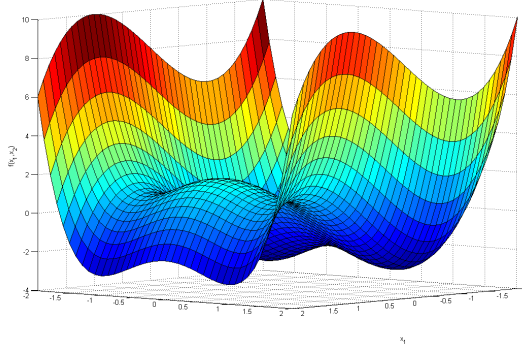


Figure 4.16: Plot of $f(x_1, x_2) = 3x_1^2 - x_1^3 - 2x_2^2 + x_2^4$, showing the various local maxima and minima of the function.

Proof: The idea behind this theorem can be stated as follows. The tangent hyperplane to the function at any extreme point must be parallel to the plane $z = 0$. This can happen if and only if the gradient $\nabla F = [\nabla^T f, -1]^T$ is parallel to the z -axis at the extreme point. Or equivalently, the gradient to the function f must be the zero vector at every extreme point, *i.e.*, $f_{x_i}(\mathbf{x}^*) = 0$ for $1 \leq i \leq n$.

To formally prove this theorem, consider the function $g_i(x_i) = f(x_1^*, x_2^*, \dots, x_{i-1}^*, x_i, x_{i+1}^*, \dots, x_n^*)$. If f has a local extremum at \mathbf{x}^* , then each function $g_i(x_i)$ must have a local extremum at x_i^* . Therefore $g_i'(x_i^*) = 0$ by theorem 39. Now $g_i'(x_i^*) = f_{x_i}(\mathbf{x}^*)$ so $f_{x_i}(\mathbf{x}^*) = 0$. \square

Applying theorem 60 to the function $f(x_1, x_2) = 9 - x_1^2 - x_2^2$, we require that at any extreme point $f_{x_1} = -2x_1 = 0 \Rightarrow x_1 = 0$ and $f_{x_2} = -2x_2 = 0 \Rightarrow x_2 = 0$. Thus, f indeed attains its maximum at the point $(0, 0)$ as shown in Figure 4.17.

Definition 27 [Critical point]: A point \mathbf{x}^* is called a critical point of a function $f(\mathbf{x})$ defined on $\mathcal{D} \subseteq \mathbb{R}^n$ if

1. If $f_{x_i}(\mathbf{x}^*) = 0$, for $1 \leq i \leq n$.
2. OR $f_{x_i}(\mathbf{x}^*)$ fails to exist for any $1 \leq i \leq n$.

A procedure for computing all critical points of a function f is:

1. Compute f_{x_i} for $1 \leq i \leq n$.
2. Determine if there are any points where any one of f_{x_i} fails to exist. Add such points (if any) to the list of critical points.
3. Solve the system of equations $f_{x_i} = 0$ simultaneously. Add the solution points to the list of saddle points.

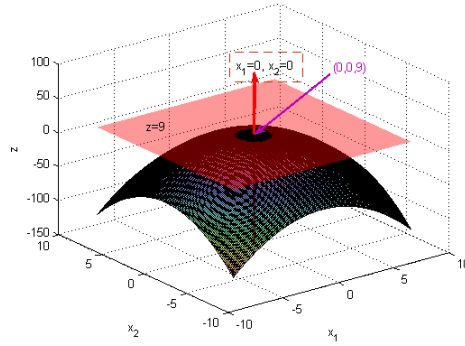


Figure 4.17: The paraboloid $f(x_1, x_2) = 9 - x_1^2 - x_2^2$ attains its maximum at $(0, 0)$. The tangent plane to the surface at $(0, 0, f(0, 0))$ is also shown, and so is the gradient vector ∇F at $(0, 0, f(0, 0))$.

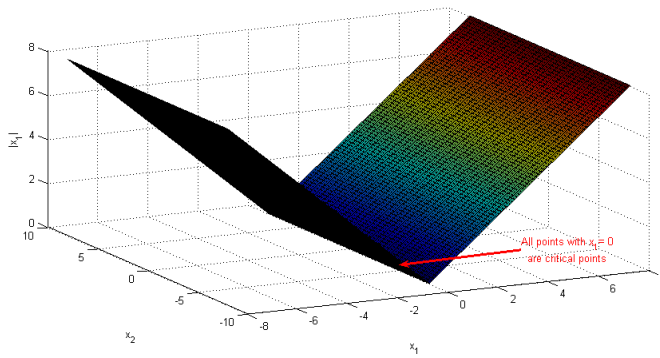


Figure 4.18: Plot illustrating critical points where derivative fails to exist.

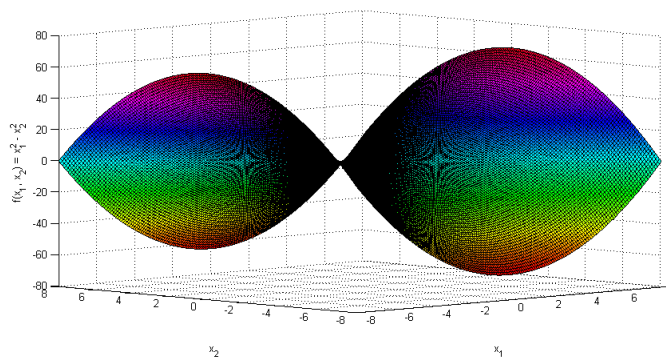


Figure 4.19: The hyperbolic paraboloid $f(x_1, x_2) = x_1^2 - x_2^2$, which has a saddle point at $(0, 0)$.

As an example, for the function $f(x_1, x_2) = |x_1|$, f_{x_1} does not exist for $(0, s)$ for any $s \in \Re$ and all of them are critical points. Figure 4.18 shows the corresponding 3-D plot.

Is the converse of theorem 60 true? That is, if you find an \mathbf{x}^* that satisfies $f_{x_i}(\mathbf{x}^*) = 0$ for all $1 \leq i \leq n$, is it necessary that \mathbf{x}^* is an extreme point? The answer is no. In fact, points that violate the converse of theorem 60 are called saddle points.

Definition 28 [Saddle point]: A point \mathbf{x}^* is called a saddle point of a function $f(\mathbf{x})$ defined on $\mathcal{D} \subseteq \Re^n$ if \mathbf{x}^* is a critical point of f but \mathbf{x}^* does not correspond to a local maximum or minimum of the function.

We saw the example of a saddle point in Figure 4.7, for the case $n = 1$. The *inflection point* for a function of single variable, that was discussed earlier, is the analogue of the saddle point for a function of multiple variables. An example for $n = 2$ is the hyperbolic paraboloid⁵ $f(x_1, x_2) = x_1^2 - x_2^2$, the graph of which is shown in Figure 4.19. The hyperbolic paraboloid opens up on x_1 -axis (Figure 4.20) and down on x_2 -axis (Figure 4.21) and has a saddle point at $(0, 0)$.

To get working on figuring out how to find the maximum and minimum of a function, we will take some examples. Let us find the critical points of $f(x_1, x_2) = x_1^2 + x_2^2 - 2x_1 - 6x_2 + 14$ and classify the critical point. This function is a polynomial function and is differentiable everywhere. It is a paraboloid that is shifted away from origin. To find its critical points, we will solve $f_{x_1} = 2x_1 - 2 = 0$ and $f_{x_2} = 2x_2 - 6 = 0$, which when solved simultaneously, yield a single critical point $(1, 3)$. For a simple example like this, the function f can be rewritten as $f(x_1, x_2) = (x_1 - 1)^2 + (x_2 - 3)^2 + 4$, which implies that $f(x_1, x_2) \geq 4 = f(1, 3)$. Therefore, $(1, 3)$ is indeed a local minimum (in fact a global minimum) of $f(x_1, x_2)$.

⁵The hyperbolic paraboloid is shaped like a *saddle* and can have a critical point called the saddle point.

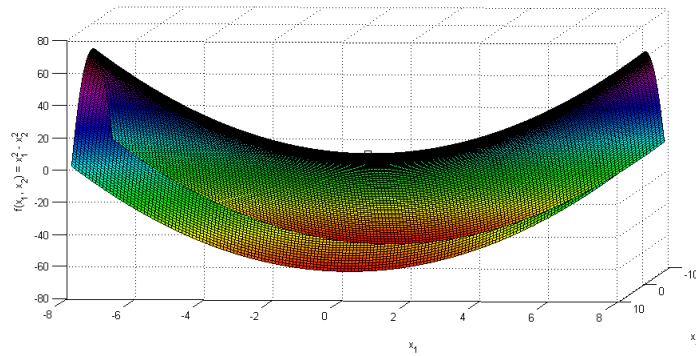


Figure 4.20: The hyperbolic paraboloid $f(x_1, x_2) = x_1^2 - x_2^2$, when viewed from the x_1 axis is concave up.

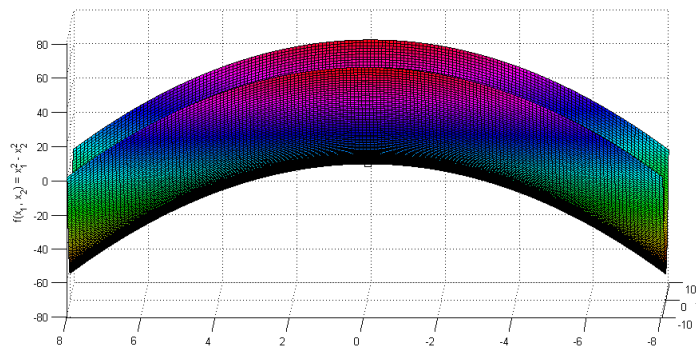


Figure 4.21: The hyperbolic paraboloid $f(x_1, x_2) = x_1^2 - x_2^2$, when viewed from the x_2 axis is concave down.

However, it is not always so easy to determine if a critical point is a point of local extreme value. To understand this, consider the function $f(x_1, x_2) = 2x_1^3 + x_1x_2^2 + 5x_1^2 + x_2^2$. The system of equations to be solved are $f_{x_1} = 6x_1^2 + x_2^2 + 10x_1 = 0$ and $f_{x_2} = 2x_1x_2 + 2x_2 = 0$. From the second equation, we get either $x_2 = 0$ or $x_1 = -1$. Using these values one at a time in the first equation, we get values for the other variables. The critical points are: $(0, 0)$, $(-\frac{5}{3}, 0)$, $(-1, 2)$ and $(-1, -2)$. Which of these critical points correspond to extreme values of the function? Since f does not have a quadratic form, it is not easy to find a lower bound on the function as in the previous example. However, we can make use of the Taylor series expansion for single variable to find polynomial expansions of functions of n variables. The following theorem gives a systematic method, similar to the second derivative test for functions of single variable, for finding maxima and minima of functions of multiple variables.

Theorem 61 *Let $f : \mathcal{D} \rightarrow \Re$ where $\mathcal{D} \subseteq \Re^n$. Let $f(\mathbf{x})$ have continuous partial derivatives and continuous mixed partial derivatives in an open ball \mathcal{R} containing a point \mathbf{x}^* where $\nabla f(\mathbf{x}^*) = 0$. Let $\nabla^2 f(\mathbf{x})$ denote an $n \times n$ matrix of mixed partial derivatives of f evaluated at the point \mathbf{x} , such that the ij^{th} entry of the matrix is $f_{x_i x_j}$. The matrix $\nabla^2 f(\mathbf{x})$ is called the Hessian matrix. The Hessian matrix is symmetric⁶. Then,*

- If $\nabla^2 f(\mathbf{x}^*)$ is positive definite, \mathbf{x}^* is a local minimum.
- If $\nabla^2 f(\mathbf{x}^*)$ is negative definite (that is if $-\nabla^2 f(\mathbf{x}^*)$ is positive definite), \mathbf{x}^* is a local maximum.

Proof: Since the mixed partial derivatives of f are continuous in an open ball containing \mathcal{R} containing \mathbf{x}^* and since $\nabla^2 f(\mathbf{x}^*) \succ 0$, it can be shown that there exists an $\epsilon > 0$, with $\mathcal{B}(\mathbf{x}^*, \epsilon) \subseteq \mathcal{R}$ such that for all $\|\mathbf{h}\| < \epsilon$, $\nabla^2 f(\mathbf{x}^* + \mathbf{h}) \succ 0$. Consider an increment vector \mathbf{h} such that $(\mathbf{x}^* + \mathbf{h}) \in \mathcal{B}(\mathbf{x}^*, \epsilon)$. Define $g(t) = f(\mathbf{x}^* + t\mathbf{h}) : [0, 1] \rightarrow \Re$. Using the chain rule,

$$g'(t) = \sum_{i=1}^n f_{x_i}(\mathbf{x}^* + t\mathbf{h}) \frac{dx_i}{dt} = \mathbf{h}^T \cdot \nabla f(\mathbf{x}^* + t\mathbf{h})$$

Since f has continuous partial and mixed partial derivatives, g' is a differentiable function of t and

$$g''(t) = \mathbf{h}^T \nabla^2 f(\mathbf{x}^* + t\mathbf{h}) \mathbf{h}$$

Since g and g' are continuous on $[0, 1]$ and g' is differentiable on $(0, 1)$, we can make use of the Taylor's theorem (45) with $n = 1$ and $a = 0$ to obtain:

$$g(1) = g(0) + g'(0) + \frac{1}{2}g''(c)$$

⁶By Clairaut's Theorem, if the partial and mixed derivatives of a function are continuous on an open region containing a point \mathbf{x}^* , then $f_{x_i x_j}(\mathbf{x}^*) = f_{x_j x_i}(\mathbf{x}^*)$, for all $i, j \in [1, n]$.

for some $c \in (0, 1)$. Writing this equation in terms of f gives

$$f(\mathbf{x}^* + \mathbf{h}) = f(\mathbf{x}^*) + \mathbf{h}^T \nabla f(\mathbf{x}^*) + \frac{1}{2} \mathbf{h}^T \nabla^2 f(\mathbf{x}^* + c\mathbf{h}) \mathbf{h}$$

We are given that $\nabla f(\mathbf{x}^*) = 0$. Therefore,

$$f(\mathbf{x}^* + \mathbf{h}) - f(\mathbf{x}^*) = \frac{1}{2} \mathbf{h}^T \nabla^2 f(\mathbf{x}^* + c\mathbf{h}) \mathbf{h}$$

The presence of an extremum of f at \mathbf{x}^* is determined by the sign of $f(\mathbf{x}^* + \mathbf{h}) - f(\mathbf{x}^*)$. By virtue of the above equation, this is the same as the sign of $H(c) = \mathbf{h}^T \nabla^2 f(\mathbf{x}^* + c\mathbf{h}) \mathbf{h}$. Because the partial derivatives of f are continuous in \mathcal{R} , if $H(0) \neq 0$, the sign of $H(c)$ will be the same as the sign of $H(0) = \mathbf{h}^T \nabla^2 f(\mathbf{x}^*) \mathbf{h}$ for \mathbf{h} with sufficiently small components. Therefore, if $\nabla^2 f(\mathbf{x}^*)$ is positive definite, we are guaranteed to have $H(0)$ positive, implying that f has a local minimum at \mathbf{x}^* . Similarly, if $-\nabla^2 f(\mathbf{x}^*)$ is positive definite, we are guaranteed to have $H(0)$ negative, implying that f has a local maximum at \mathbf{x}^* . \square

Theorem 61 gives sufficient conditions for local maxima and minima of functions of multiple variables. Along similar lines of the proof of theorem 61, we can prove necessary conditions for local extrema in theorem 62.

Theorem 62 *Let $f : \mathcal{D} \rightarrow \Re$ where $\mathcal{D} \subseteq \Re^n$. Let $f(\mathbf{x})$ have continuous partial derivatives and continuous mixed partial derivatives in an open region \mathcal{R} containing a point \mathbf{x}^* where $\nabla f(\mathbf{x}^*) = 0$. Then,*

- *If \mathbf{x}^* is a point of local minimum, $\nabla^2 f(\mathbf{x}^*)$ must be positive semi-definite.*
- *If \mathbf{x}^* is a point of local maximum, $\nabla^2 f(\mathbf{x}^*)$ must be negative semi-definite (that is, $-\nabla^2 f(\mathbf{x}^*)$ must be positive semi-definite).*

The following corollary of theorem 62 states a sufficient condition for a point to be a saddle point.

Corollary 63 *Let $f : \mathcal{D} \rightarrow \Re$ where $\mathcal{D} \subseteq \Re^n$. Let $f(\mathbf{x})$ have continuous partial derivatives and continuous mixed partial derivatives in an open region \mathcal{R} containing a point \mathbf{x}^* where $\nabla f(\mathbf{x}^*) = 0$. If $\nabla^2 f(\mathbf{x}^*)$ is neither positive semi-definite nor negative semi-definite (that is, some of its eigenvalues are positive and some negative), then \mathbf{x}^* is a saddle point.*

Thus, for a function of more than one variable, the second derivative test generalizes to a test based on the eigenvalues of the function's Hessian matrix at the stationary point. Based on theorem 61, we will derive the second derivative test for determining extreme values of a function of two variables.

Theorem 64 *Let the partial and second partial derivatives of $f(x_1, x_2)$ be continuous on a disk with center (a, b) and suppose $f_{x_1}(a, b) = 0$ and $f_{x_2}(a, b) = 0$ so that (a, b) is a critical point of f . Let $D(a, b) = f_{x_1 x_1}(a, b) f_{x_2 x_2}(a, b) - [f_{x_1 x_2}(a, b)]^2$. Then⁷,*

⁷ D here stands for the discriminant.

- If $D > 0$ and $f_{x_1x_1}(a, b) > 0$, then $f(a, b)$ is a local minimum.
- Else if $D > 0$ and $f_{x_1x_1}(a, b) < 0$, then $f(a, b)$ is a local maximum.
- Else if $D < 0$ then (a, b) is a saddle point.

Proof: Recall the definition of positive definiteness; a matrix is positive definite if all its eigenvalues are positive. For the 2×2 matrix $\nabla^2 f$ in this problem, the product of the eigenvalues is $\det(\nabla^2 f) = f_{x_1x_1}(a, b)f_{x_2x_2}(a, b) - [f_{x_1x_2}(a, b)]^2$ and the sum of the eigenvalues is $f_{x_1x_1}(a, b) + f_{x_2x_2}(a, b)$. Now:

- If $\det(\nabla^2 f(a, b)) > 0$ and if additionally $f_{x_1x_1}(a, b) > 0$ (or equivalently, $f_{x_2x_2}(a, b) > 0$), the product as well as the sum of eigenvalues will be positive, implying that the eigenvalues are positive and therefore $\nabla^2 f(a, b)$ is positive definite, According to theorem 61, this is a sufficient condition for $f(a, b)$ to be a local minimum.
- If $\det(\nabla^2 f(a, b)) > 0$ and if additionally $f_{x_1x_1}(a, b) < 0$ (or equivalently, $f_{x_2x_2}(a, b) < 0$), the product of the eigenvalue is positive whereas the sum is negative, implying that the eigenvalues are negative and therefore $\nabla^2 f(a, b)$ is negative definite, According to theorem 61, this is a sufficient condition for $f(a, b)$ to be a local maximum.
- If $\det(\nabla^2 f(a, b)) < 0$, the eigenvalues must have opposite signs, implying that the $\nabla^2 f(a, b)$ is neither positive semi-definite nor negative-semidefinite. By corollary 63, this is a sufficient condition for $f(a, b)$ to be a saddle point.

□

We saw earlier that the critical points for $f(x_1, x_2) = 2x_1^3 + x_1x_2^2 + 5x_1^2 + x_2^2$ are $(0, 0)$, $(-\frac{5}{3}, 0)$, $(-1, 2)$ and $(-1, -2)$. To determine which of these correspond to local extrema and which are saddle, we first compute the partial derivatives of f :

$$\begin{aligned} f_{x_1x_1}(x_1, x_2) &= 12x_1 + 10 \\ f_{x_2x_2}(x_1, x_2) &= 2x_1 + 2 \\ f_{x_1x_2}(x_1, x_2) &= 2x_2 \end{aligned}$$

Using theorem 64, we can verify that $(0, 0)$ corresponds to a local minimum, $(-\frac{5}{3}, 0)$ corresponds to a local maximum while $(-1, 2)$ and $(-1, -2)$ correspond to saddle points. Figure 4.22 shows the plot of the function while pointing out the four critical points.

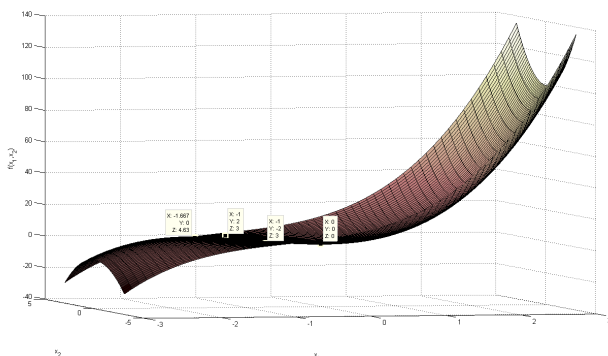


Figure 4.22: Plot of the function $2x_1^3 + x_1x_2^2 + 5x_1^2 + x_2^2$ showing the four critical points.

We will take some more examples:

1. Consider a significantly harder function $f(x, y) = 10x^2y - 5x^2 - 4y^2 - x^4 - 2y^4$. Let us find and classify its critical points. The gradient vector is $\nabla f(x, y) = [20xy - 10x - 4x^3, 10x^2 - 8y - 8y^3]$. The critical points correspond to solutions of the simultaneous set of equations

$$\begin{aligned} 20xy - 10x - 4x^3 &= 0 \\ 10x^2 - 8y - 8y^3 &= 0 \end{aligned} \quad (4.15)$$

One of the solutions corresponds to solving the system $-8y^3 + 42y - 25 = 0$ ⁸ and $10x^2 = 50y - 25$, which have four real solutions⁹, *viz.*, $(0.8567, 0.646772)$, $(-0.8567, 0.646772)$, $(2.6442, 1.898384)$, and $(-2.6442, 1.898384)$. Another real solution is $(0, 0)$. The mixed partial derivatives of the function are

$$\begin{aligned} f_{xx} &= 20y - 10 - 12x^2 \\ f_{xy} &= 20x \\ f_{yy} &= -8 - 24y^2 \end{aligned} \quad (4.16)$$

Using theorem 64, we can verify that $(2.6442, 1.898384)$ and $(-2.6442, 1.898384)$ correspond to local maxima whereas $(0.8567, 0.646772)$ and $(-0.8567, 0.646772)$ correspond to saddle points. This is illustrated in Figure 4.23.

⁸Solving this using matlab without proper scaling could give you complex values. With proper scaling of the equation, you should get $y = -2.545156$ or $y = 0.646772$ or $y = 1.898384$.

⁹The values of x corresponding to $y = -2.545156$ are complex

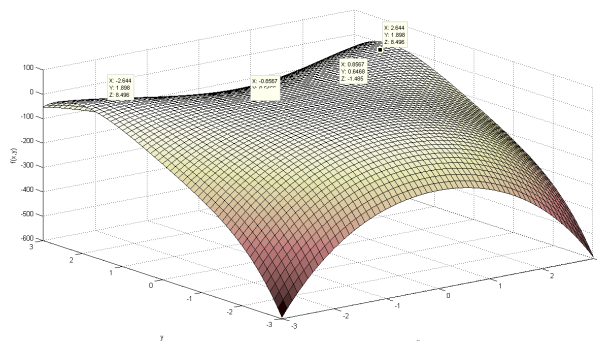


Figure 4.23: Plot of the function $10x^2y - 5x^2 - 4y^2 - x^4 - 2y^4$ showing the four critical points.

2. The function $f(x, y) = x \sin y$ has the gradient vector $[\sin y, x \cos y]$. The critical points correspond to the solutions to the simultaneous set of equations

$$\begin{aligned} \sin y &= 0 \\ x \cos y &= 0 \end{aligned} \tag{4.17}$$

The critical points are¹⁰ $(0, n\pi)$ for $n = 0, \pm 1, \pm 2, \dots$. The mixed partial derivatives of the function are

$$\begin{aligned} f_{xx} &= 0 \\ f_{xy} &= \cos y \\ f_{yy} &= -x \sin y \end{aligned} \tag{4.18}$$

which tell us that the discriminant function $D = -\cos^2 y$ is always negative. Therefore, all the critical points turn out to be saddle points. This is illustrated in Figure 4.24.

Along similar lines of the single variable case, we next define the global maximum and minimum.

Definition 29 [Global maximum]: A function f of n variables, with domain $\mathcal{D} \subseteq \mathbb{R}^n$ has an absolute or global maximum at \mathbf{x}^0 if $\forall \mathbf{x} \in \mathcal{D}, f(\mathbf{x}) \leq f(\mathbf{x}^0)$.

¹⁰Note that the *cosine* does not vanish wherever the *sine* vanishes.

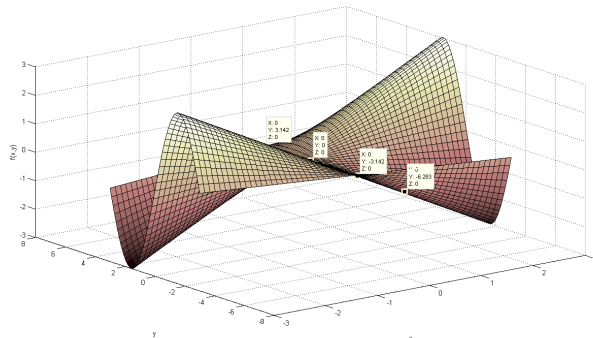


Figure 4.24: Plot of the function $x \sin y$ illustrating that all critical points are saddle points.

Definition 30 [Global minimum]: A function f of n variables, with domain $\mathcal{D} \subseteq \mathbb{R}^n$ has an absolute or global minimum at \mathbf{x}^0 if $\forall \mathbf{x} \in \mathcal{D}, f(\mathbf{x}) \geq f(\mathbf{x}^0)$.

We would like to find the absolute maximum and minimum values of a function of multiple variables in a closed interval, along similar lines of the method yielded by theorem 41 for functions of single variable. The procedure was to evaluate the value of the function at the critical points as well as the end points of the interval and determine the absolute maximum and minimum values by scanning this list. To generalize the idea to function of multiple variables, we point out that the analogue of finding the value of the function at the boundaries of closed interval in the single variable case is to find the function value along the boundary curve, which reduces the evaluation of a function of multiple variables to evaluating a function of a single variable. Recall from the definitions on page 212 that a closed set in \mathbb{R}^n is a set that contains its boundary points (analogous to closed interval in \mathbb{R}) while a bounded set in \mathbb{R}^n is a set that is contained inside a closed ball, $\mathcal{B}[\mathbf{0}, \epsilon]$. An example bounded set is $\{(x_1, x_2, x_3) | x_1^2 + x_2^2 + x_3^2 \leq 1\}$. An example unbounded set is $\{(x_1, x_2, x_3) | x_1 > 1, x_2 > 1, x_3 > 1\}$. Based on these definitions, we can state the extreme value theorem for a function of n variables.

Theorem 65 Let $f : \mathcal{D} \rightarrow \mathbb{R}$ where $\mathcal{D} \subseteq \mathbb{R}^n$ is a closed bounded set and f be continuous on \mathcal{D} . Then f attains an absolute maximum and absolute minimum at some points in \mathcal{D} .

The theorem implies that whenever a function of n variables is restricted to a bounded space, it has an absolute maximum and an absolute minimum. Following theorem 60, we note that the locally extreme values of a function occur at its critical points. By the very definition of local extremum, it cannot occur at the boundary point of \mathcal{D} . Since every absolute extremum is also a

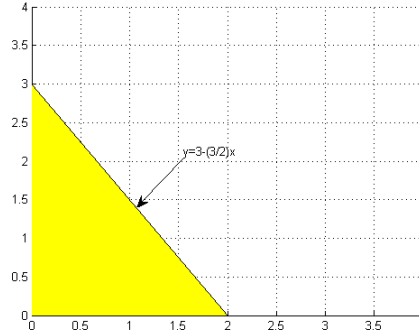


Figure 4.25: The region bounded by the points $(0, 3)$, $(2, 0)$, $(0, 0)$ on which we consider the maximum and minimum of the function $f(x, y) = 1 + 4x - 5y$.

local extremum, the absolute maximum and minimum of a function on a closed, bounded set will either happen at the critical points or at the boundary. The procedure for finding the absolute maximum and minimum of a function on a closed bounded set is outlined below and is similar to the procedure 4 for a function of single variable continuous on a closed and bounded interval:

Procedure 5 [Finding extreme values on closed, bounded sets]: *To find the absolute maximum and absolute minimum of a continuous function f on a closed bounded set \mathcal{D} ;*

- *evaluate f at the critical points of f on \mathcal{D}*
- *find the extreme values of f on the boundary of \mathcal{D}*
- *the largest of the values found above is the absolute maximum, and the smallest of them is the absolute minimum.*

We will take some examples to illustrate procedure 5.

1. Consider the function $f(x, y) = 1 + 4x - 5y$ defined on the region \mathcal{R} bounded by the points $(0, 3)$, $(2, 0)$, $(0, 0)$. The region \mathcal{R} is shown in Figure 4.25 and is bounded by three line segments

- \mathbf{B}_1 : $x = 0, 0 \leq y \leq 3$
- \mathbf{B}_2 : $y = 0, 0 \leq x \leq 2$
- and \mathbf{B}_3 : $y = 3 - \frac{3}{2}x, 0 \leq x \leq 2$.

The linear function $f(x, y) = 1 + 4x - 5y$ has no critical points, since $\nabla f(x, y) = [4, -5]^T$ is defined everywhere, though it cannot disappear at any point. In fact, linear functions have no critical points and the extreme values are always assumed at the boundaries; this forms the basis of linear programming. We will find the extreme values on the boundaries.

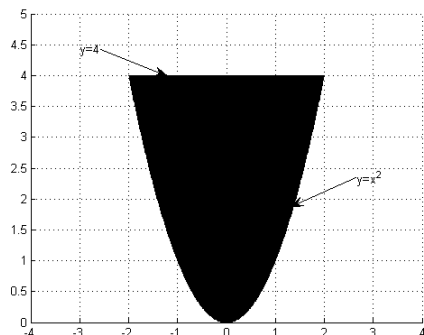


Figure 4.26: The region \mathcal{R} bounded by $y = x^2$ and $y = 4$ on which we consider the maximum and minimum of the function $f(x, y) = 1 - xy - x - y$.

- On \mathbf{B}_1 , $f(x, y) = f(0, y) = 1 - 5y$, for $y \in [0, 3]$. This is a single variable extreme value problem for a continuous function. Its largest value is assumed at $y = 0$ and equals 1 while the smallest value is assumed at $y = 3$ and equals -14 .
- On \mathbf{B}_2 , $f(x, y) = f(x, 0) = 1 + 4x$, for $x \in [0, 2]$. This is again a single variable extreme value problem for a continuous function. Its largest value is assumed at $x = 2$ and equals 9 while the smallest value is assumed at $x = 0$ and equals 1.
- On \mathbf{B}_3 , $f(x, y) = 1 + 4x - 5(3 - (3/2)x) = -14 + (23/2)x$, for $x \in [0, 2]$. This is also a single variable extreme value problem for a continuous function. Its largest value is assumed at $x = 2$ and equals 9 while the smallest value is assumed at $x = 0$ and equals -14 .

Thus, the absolute maximum is attained by f at $(2, 0)$ while the absolute minimum is attained at $(0, 3)$. Both extrema are at the vertices of the polygon (triangle) This example illustrates the general procedure for determining the absolute maximum and minimum of a function on a closed, bounded set. However, the problem can become very hard in practice as the function f gets complex.

2. Let us look at a harder problem. Let us find the absolute maximum and the absolute minimum of the function $f(x, y) = 1 - xy - x - y$ on the region \mathcal{R} bounded by $y = x^2$ and $y = 4$. This is not a linear function any longer. The region \mathcal{R} is shown in Figure 4.26 and is bounded by

- \mathbf{B}_1 : $y = x^2$, $-2 \leq x \leq 2$
- \mathbf{B}_2 : $y = 4$, $-2 \leq x \leq 2$

Since $f(x, y) = 1 - xy - x - y$ is differentiable everywhere, the critical point of f is characterized by $\nabla f(x, y) = [-y - 1, x - 1]^T = \mathbf{0}$, that is

$x = -1$, $y = -1$. However, this point does not lie in \mathcal{R} and hence, there are no critical points, in \mathcal{R} . Along similar lines of the previous problem, we will find the extreme values of f on the boundaries of \mathcal{R} .

- On \mathbf{B}_1 , $f(x, y) = 1 - x^3 - x - x^2$, for $x \in [-2, 2]$. This is a single variable extreme value problem for a continuous function. Its critical points correspond to solutions of $3x^2 + 2x + 1 = 0$. However, this equation has no real solutions¹¹ and therefore, the function's extreme values are only at the boundary points; the minimum value -13 is attained at $x = 2$ and the maximum value 7 is attained at $x = -2$.
- On \mathbf{B}_2 , $f(x, y) = 1 - 4x - x - 4 = -3 - 5x$, for $x \in [-2, 2]$. This is again a single variable extreme value problem for a continuous function. It has no critical points and extreme values correspond to the boundary points; its maximum value 7 is assumed at $x = -2$ while the minimum value -13 is assumed at $x = 2$.

Thus, the absolute maximum value 7 is attained by f at $(-2, 4)$ while the absolute minimum value -13 is attained at $(2, 4)$.

3. Consider the same problem as the previous one, with a slightly different objective function, $f(x, y) = 1 + xy - x - y$. The critical point of f is characterized by $\nabla f(x, y) = [y - 1, x - 1]^T = \mathbf{0}$, that is $x = 1$, $y = 1$. This lies within \mathcal{R} and f takes the value 0 at $(1, 1)$. Next, we find the extreme values of f on the boundaries of \mathcal{R} .

- On \mathbf{B}_1 , $f(x, y) = 1 + x^3 - x - x^2$, for $x \in [-2, 2]$. Its critical points correspond to solutions of $3x^2 - 2x - 1 = 0$. Its solutions are $x = 1$ and $x = -\frac{1}{3}$. The function values corresponding to these points are $f(1, 1) = 0$ and $f(-1/3, 1/9) = 32/27$. At the boundary points, the function assumes the values $f(-2, 4) = -9$ and $f(2, 4) = 3$. Thus, the maximum value on \mathbf{B}_1 is $f(2, 4) = 3$ and the minimum value is $f(-2, 4) = -9$.
- On \mathbf{B}_2 , $f(x, y) = 1 + 4x - x - 4 = -3 + 3x$, for $x \in [-2, 2]$. It has no critical points and extreme values correspond to the boundary points; At the boundary points, the function assumes the values $f(-2, 4) = -9$ and $f(2, 4) = 3$, which correspond to the minimum and maximum values respectively of f on \mathbf{B}_2 .

Thus, the absolute maximum value 3 is attained by f at $(2, 4)$ while the absolute minimum value -9 is attained at $(-2, 4)$.

4.1.5 Absolute extrema and Convexity

Theorem 61 specified a sufficient condition for the local minimum of a differentiable function with continuous partial and mixed partial derivatives, while

¹¹The complex solutions are $x = -\frac{1}{3} + i\frac{1}{3}\sqrt{2}$ and $x = -\frac{1}{3} - i\frac{1}{3}\sqrt{2}$.

theorem 62 specified a necessary condition for the same. Can these conditions be extended to globally optimal solutions? The answer is that the extensions to globally optimal solutions can be made for a specific class of optimization problems called convex optimization problems. In the next section we introduce the concept of convex sets and convex functions, enroute to discussing convex optimization.

4.2 Convex Optimization Problem

A function $f(\cdot)$ is called convex if its value at the scalar combination of two points x and y is less than the same scalar combination of the function at the two points. In other words, $f(\cdot)$ is convex if and only if:

$$\begin{aligned} f(\alpha x + \beta y) &\leq \alpha f(x) + \beta f(y) \\ \text{if } \alpha + \beta &= 1, \alpha \geq 0, \beta \geq 0 \end{aligned} \quad (4.19)$$

For a convex optimization problem, the objective function $f(x)$ as well as the inequality functions $g_i(x), i = 1, \dots, m$ are convex. The equality constraints are linear, *i.e.*, of the form, $Ax = b$.

$$\begin{aligned} \text{minimize} \quad & f(x) \\ \text{subject to} \quad & g_i(x) \leq 0, \quad i = 1, \dots, m \\ & Ax = b \end{aligned} \quad (4.20)$$

Least squares and linear programming are special cases of convex optimization problems. Like in the case of linear programming, there are no analytical solutions for convex optimization problems. But they can be solved reliably, efficiently and optimally. There are not many well developed software for the general class of convex optimization problems, though there are several software packages in matlab, C, *etc.*, and many free softwares as well. The computation time is polynomial but more complicated to be expressed exactly because the computation time depends on the cost of validating the function values and their derivatives. Modulo that, computation time for convex optimization problems is similar to that for linear programming problems.

To pose practical problems as convex optimization problems is more difficult than to recognize least squares and linear programs. There exist many techniques to reformulate problems in the convex form. However, surprisingly, many problems in practice can be solved via convex optimization.

4.2.1 Why Convex Optimization?

We will see in this sequel, that generic convex programs, under mild computability and boundedness assumptions, are computationally tractable. Many convex

programs admit theoretically and practically efficient solution methods. Convex optimization admits *duality theory*, which can be used to quantitatively establish the quality of an approximate solution. Even though duality may not yield a closed-form solution, it often facilitates nontrivial reformulations of the problem. Duality theory also comes handy in confirming if an approximate solution is optimal.

In contrast to this, rarely does it happen that a global solution can be efficiently found for nonconvex optimization programs¹². For most nonconvex programs, there are no sound techniques for certifying the global optimality of approximate solutions or estimating how non-optimal an approximate solution is.

4.2.2 History

Numerical optimization started in the 1940s with the development of the simplex method for linear programming. The next obvious extension to linear programming was by replacing the linear cost function with a quadratic cost function. The linear inequality constraints were however maintained. This first extension took place in the 1950s. We can expect that the next step would have been to replace the linear constraints with quadratic constraints. But it did not really happen that way. On the other hand, around the end of the 1960s, there was another non-linear, convex extension of linear programming called *geometric programming*. Geometric programming includes linear programming as a special case. Nothing more happened until the beginning of the 1990s. The beginning of the 1990s was marked by a big explosion of activities in the area of convex optimizations, and development really picked up. Researches formulated different and more general classes of convex optimization problems that are known as semidefinite programming, second-order cone programming, quadratically constrained quadratic programming, sum-of-squares programming, *etc.*

The same happened in terms of applications. Since 1990s, applications have been investigated in many different areas. One of the first application areas was control, and the optimization methods that were investigated included semi-definite programming for certain control problem. Geometric programming had been around since late 1960s and it was applied extensively to circuit design problems. Quadratic programming found application in machine learning problem formulations such as support vector machines. Semi-definite programming relaxations found use in combinatorial optimization. There were many other interesting applications in different areas such as image processing, quantum information, finance, signal processing, communications, *etc.*

This first look at the activities involving applications of optimization clearly indicates that a lot of development took place around the 1990s. Further, people extended interior-point methods (which were already known for linear

¹²Optimization problems such as singular value decomposition are some few exceptions to this.

programming since 1984¹³) to non-linear convex optimization problems. A highlight in this area was the work of Nesterov and Nemirovski who extended Karmarkar's work to polynomial-time interior-point methods for nonlinear convex programming in their book published in 1994, though the work actually took place in 1990. As a result, people started looking at non-linear convex optimization in a special way; instead of treating non-linear convex optimization as a special case of non-linear optimization, they looked at it as an extension of linear programming which can be solved almost with the same efficiency. Once people started looking at applications of non-linear convex optimization, they discovered many!

We will begin with a background on convex sets and functions. Convex sets and functions constitute the basic theory for the entire area of convex optimization. Next, we will discuss some standard problem classes and some recently studied problem classes such as semi-definite programming and cone programming. Finally, we will look at applications.

4.2.3 Affine Set

Definition 31 [Affine Set]: A set \mathcal{A} is called affine if the line connecting any two distinct points in the set is completely contained within \mathcal{A} . Mathematically, the set \mathcal{A} is called affine if

$$\forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}, \quad \theta \in \mathfrak{R} \quad \Rightarrow \quad \theta \mathbf{x}_1 + (1 - \theta) \mathbf{x}_2 \in \mathcal{A} \quad (4.21)$$

Theorem 66 The solution set of the system of linear equations $A\mathbf{x} = \mathbf{b}$ is an affine set.

Proof: Suppose \mathbf{x}_1 and \mathbf{x}_2 are solutions to the system $A\mathbf{x} = \mathbf{b}$ with $\mathbf{x}_1 \neq \mathbf{x}_2$. Then, $A(\theta \mathbf{x}_1 + (1 - \theta) \mathbf{x}_2) = \theta \mathbf{b} + (1 - \theta) \mathbf{b} = \mathbf{b}$. Thus, $\theta \mathbf{x}_1 + (1 - \theta) \mathbf{x}_2 \in \mathcal{A}$, implying that the solution set of the system $A\mathbf{x} = \mathbf{b}$ is an affine set. \square

In fact, converse of theorem 66 is also true; any affine set can be expressed as the solution set of a system of linear equations $A\mathbf{x} = \mathbf{b}$.

4.2.4 Convex Set

Definition 32 [Convex Set]: A set \mathcal{C} is called convex if the line segment connecting any two points in the set is completely contained within \mathcal{C} . Else \mathcal{C} is called concave. That is,

$$\forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{C} \quad 0 \leq \theta \leq 1 \quad \Rightarrow \quad \theta \mathbf{x}_1 + (1 - \theta) \mathbf{x}_2 \in \mathcal{C} \quad (4.22)$$

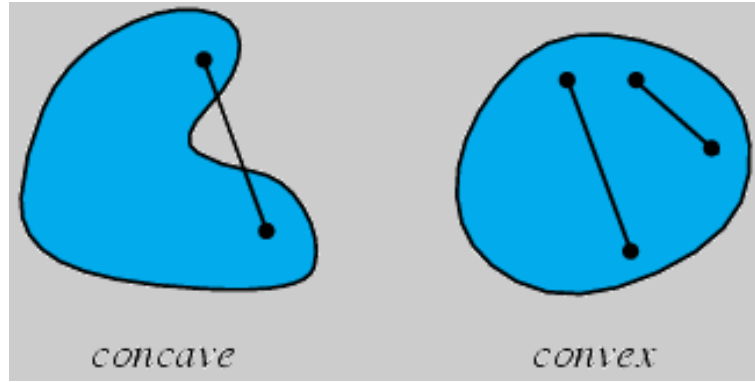


Figure 4.27: Examples of convex and non-convex sets.

Figure 4.27 shows examples of convex and non-convex (concave) sets. Since an affine set contains any line passing through two distinct points in the set, it also contains any line segment connecting two points in the set. Thus, an affine set is our first example of a convex set.

A set \mathcal{C} is a convex cone if it is convex and additionally, for every point $\mathbf{x} \in \mathcal{C}$, all non-negative multiples of \mathbf{x} are also in \mathcal{C} . In other words,

$$\forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{C} \quad \theta_1, \theta_2 \geq 0 \quad \Rightarrow \quad \theta_1 \mathbf{x}_1 + \theta_2 \mathbf{x}_2 \in \mathcal{C} \quad (4.23)$$

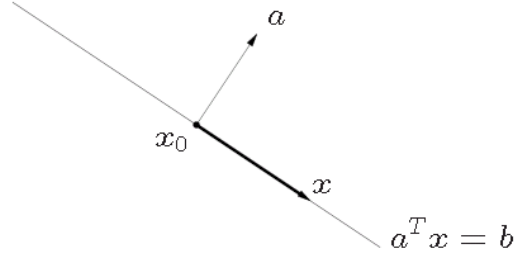
Combinations of the form $\theta_1 \mathbf{x}_1 + \theta_2 \mathbf{x}_2$ for $\theta_1 \geq 0, \theta_2 \geq 0$ are called conic combinations. We will state a related definition next - that of the convex hull of a set of points.

Definition 33 [Convex Hull]: A convex combination of the set of points $\mathcal{S} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k\}$ is any point \mathbf{x} of the form

$$\mathbf{x} = \sum_{i=1}^k \theta_i \mathbf{x}_i \quad \text{with} \quad \sum_{i=1}^k \theta_i = 1 \quad \text{and} \quad \theta_i \geq 0 \quad (4.24)$$

The convex hull $\text{conv}(\mathcal{S})$ of the set of points \mathcal{S} is the set of all convex combinations of points in \mathcal{S} . The convex hull of a convex set \mathcal{S} is \mathcal{S} itself.

¹³The first practical polynomial time algorithm for linear programming by Karmarkar (1984) involved interior-point methods.

Figure 4.28: Example of a hyperplane in \mathbb{R}^2 .

4.2.5 Examples of Convex Sets

We will look at simple but important examples of convex sets. We will also look at some operations that preserve convexity.

A *hyperplane* is the most common example of a convex set. A hyperplane is the set of solutions to a linear system of equations of the form $a^T \mathbf{x} = b$ with $a \neq 0$ and was defined earlier in definition 24. A *half space* is a solution set over the linear inequality $a^T \mathbf{x} \leq b, a \neq 0$. The hyperplane $a^T \mathbf{x} = b$ bounds the half-space from one side.

Formally,

Hyperplane: $\{\mathbf{x} | a^T \mathbf{x} = b, a \neq 0\}$. Figure 4.28 shows an example hyperplane in \mathbb{R}^2 . a is the normal vector.

Halfspace: $\{\mathbf{x} | a^T \mathbf{x} \leq b, a \neq 0\}$. Figure 4.29 shows an example half-space in \mathbb{R}^2 .

The hyperplane is convex and affine, whereas the halfspace is merely convex and not affine.

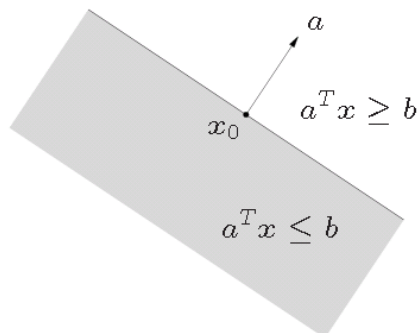
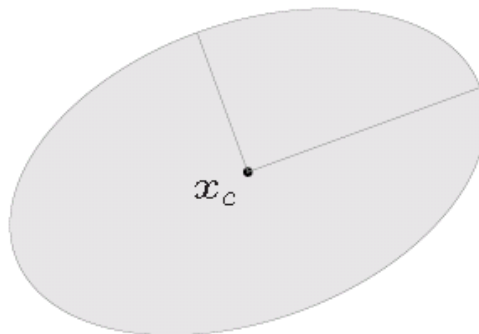
Another simple example of a convex set is a *closed ball* in \mathbb{R}^n with radius r and center \mathbf{x}_c which is an n -dimensional vector.

$$\mathcal{B}[\mathbf{x}_c, r] = \{\mathbf{x}_c + r\mathbf{u} \mid \|\mathbf{u}\|_2 \leq 1\}$$

where \mathbf{u} is a vector with norm less than or equal to 1. The open ball $\mathcal{B}(\mathbf{x}_c, r)$ is also convex. Replacing r with a non-singular square matrix A , we get an *ellipsoid* given by

$$\{\mathbf{x}_c + A\mathbf{u} \mid \|\mathbf{u}\|_2 \leq 1\}$$

which is also a convex set. Another equivalent representation of the ellipsoid can be obtained by observing that for any point \mathbf{x} in the ellipsoid, $\|A^{-1}(\mathbf{x} - \mathbf{x}_c)\|_2 \leq 1$, that is $(\mathbf{x} - \mathbf{x}_c)^T (A^{-1})^T A^{-1} (\mathbf{x} - \mathbf{x}_c) \leq 1$. Since $(A^{-1})^T = (A^T)^{-1}$ and

Figure 4.29: Example of a half-space in \mathbb{R}^2 .Figure 4.30: Example of a ellipsoid in \mathbb{R}^2 .

$A^{-1}B^{-1} = (BA)^{-1}$, the ellipsoid can be equivalently defined as $\{\mathbf{x} | (\mathbf{x} - \mathbf{x}_c)^T P^{-1} (\mathbf{x} - \mathbf{x}_c) \leq 1\}$ where $P = (AA^T)$ is a symmetric matrix. Furthermore, P is positive definite, since A is non-singular (*c.f.* page 205).

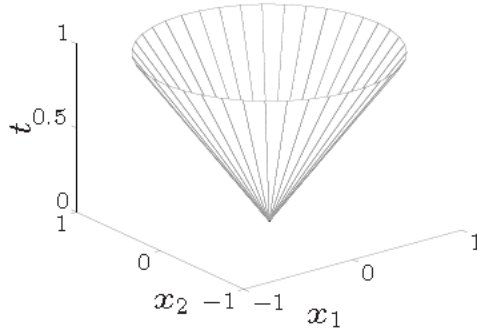
Matrix A determines the size of the ellipsoid; the eigenvalue λ_i of A determines the length of the i^{th} semi-axis of the ellipsoid (see page number 203). The ellipsoid is another example of a convex set and is a generalization of the euclidian ball. Figure 4.30 illustrates an ellipsoid in \mathbb{R}^2 .

A *norm ball* is a ball with an arbitrary norm. A norm ball with center \mathbf{x}_c and radius r is given by

$$\{\mathbf{x} \mid \|\mathbf{x} - \mathbf{x}_c\| \leq r\}$$

By the definition of the norm, a ball in that norm will be convex. The norm ball with the ∞ -norm corresponds to a square in \mathbb{R}^2 , while the norm ball with the 1-norm in \mathbb{R}^2 corresponds to the same square rotated by 45° . The norm ball is convex for all norms.

The definition of cone can be extended to any arbitrary norm to define a

Figure 4.31: Example of a cone in \mathfrak{R}^2 .

norm cone. The set of all pairs (\mathbf{x}, t) satisfying $\|\mathbf{x}\| \leq t$, *i.e.*,

$$\{(\mathbf{x}, t) \mid \|\mathbf{x}\| \leq t\}$$

is called a norm cone

When the norm is the euclidian norm, the cone (which looks like an ice-cream cone) is called the *second order cone*. Norm cones are always convex. Figure 4.31 shows a cone in \mathfrak{R}^2 . In general, the cross section of a norm cone has the shape of a norm ball with the same norm. The norm cone for the ∞ -norm is a square pyramid in \mathfrak{R}^3 and the cone for 1-norm in \mathfrak{R}^3 is the same square pyramid rotated by 45° .

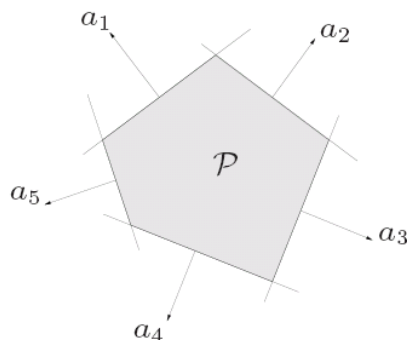
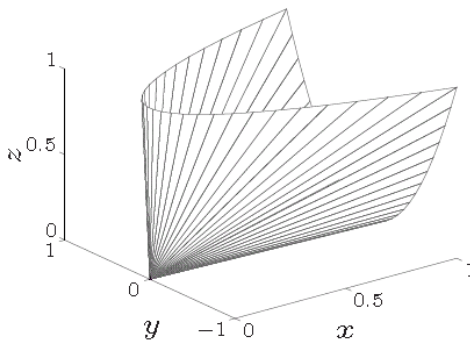
A *polyhedron* is another convex set which is given as the solution set of a finite set of linear equalities and inequalities. In matrix form, the inequalities can be stated as

$$\begin{aligned} A\mathbf{x} &\preceq \mathbf{b} & A &\in \mathfrak{R}^{m \times n} \\ C\mathbf{x} &= \mathbf{d} & C &\in \mathfrak{R}^{p \times n} \end{aligned} \tag{4.25}$$

where \preceq stands for component-wise inequality of the form \leq ¹⁴. A polyhedron can also be represented as the intersection of a finite number of halfspaces and hyperplanes. Figure 4.32 depicts a typical polyhedron in \mathfrak{R}^2 . An affine set is a special type of polyhedron.

A last simple example is the positive semi-definite cone. Let \mathcal{S}^n be the set of all symmetric $n \times n$ matrices and $\mathcal{S}_+^n \subset \mathcal{S}^n$ be the set of all positive semi-definite $n \times n$ matrices. The set \mathcal{S}_+^n is a convex cone and is called the *positive semi-definite cone*. Consider a positive semi-definite matrix S in \mathfrak{R}^2 . Then S must of the form

¹⁴The component-wise inequality corresponds to a generalized inequality \preceq_K with $K = \mathfrak{R}_+^n$.

Figure 4.32: Example of a polyhedron in \mathbb{R}^2 .Figure 4.33: Example of a positive semidefinite cone in \mathbb{R}^3 .

$$S = \begin{bmatrix} x & y \\ y & z \end{bmatrix} \quad (4.26)$$

We can represent the space of matrices \mathcal{S}_+^2 of the form $S \in \mathcal{S}_+^2$ as a three dimensional space with non-negative x , y and z coordinates and a non-negative determinant. This space corresponds to a cone as shown in Figure 4.33.

4.2.6 Convexity preserving operations

In practice if you want to establish the convexity of a set \mathcal{C} , you could either

1. prove it from first principles, *i.e.*, using the definition of convexity or
2. prove that \mathcal{C} can be built from simpler convex sets through some basic operations which preserve convexity.

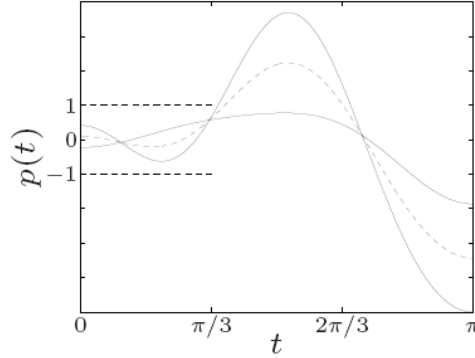


Figure 4.34: Plot for the function in (4.28)

Some of the important operations that preserve complexity are:

Intersection

The intersection of any number of convex sets is convex¹⁵. Consider the set \mathcal{S} :

$$\mathcal{S} = \left\{ \mathbf{x} \in \mathbb{R}^n \mid |p(t)| \leq 1 \text{ for } |t| \leq \frac{\pi}{3} \right\} \quad (4.27)$$

where

$$p(t) = x_1 \cos t + x_2 \cos 2t + \dots + x_m \cos mt \quad (4.28)$$

Any value of t that satisfies $|p(t)| \leq 1$, defines two regions, *viz.*,

$$\mathfrak{R}^{\leq}(t) = \{ \mathbf{x} \mid x_1 \cos t + x_2 \cos 2t + \dots + x_m \cos mt \leq 1 \}$$

and

$$\mathfrak{R}^{\geq}(t) = \{ \mathbf{x} \mid x_1 \cos t + x_2 \cos 2t + \dots + x_m \cos mt \geq -1 \}$$

Each of these regions is convex and for a given value of t , the set of points that may lie in \mathcal{S} is given by

$$\mathfrak{R}(t) = \mathfrak{R}^{\leq}(t) \cap \mathfrak{R}^{\geq}(t)$$

This set is also convex. However, not all the points in $\mathfrak{R}(t)$ lie in \mathcal{S} , since the points that lie in \mathcal{S} satisfy the inequalities for every value of t . Thus, \mathcal{S} can be given as:

$$\mathcal{S} = \bigcap_{|t| \leq \frac{\pi}{3}} \mathfrak{R}(t)$$

¹⁵Exercise: Prove.

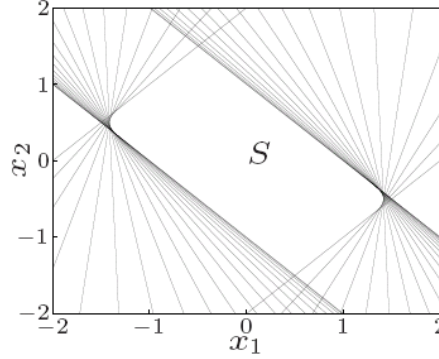


Figure 4.35: Illustration of the closure property for \mathcal{S} defined in (4.27), for $m = 2$.

Affine transform

An affine transform is one that preserves

- Collinearity between points, *i.e.*, three points which lie on a line continue to be collinear after the transformation.
- Ratios of distances along a line, *i.e.*, for distinct collinear points $\mathbf{p}_1, \mathbf{p}_2, \mathbf{p}_3$, $\frac{\|\mathbf{p}_2 - \mathbf{p}_1\|}{\|\mathbf{p}_3 - \mathbf{p}_2\|}$ is preserved.

An affine transformation or affine map between two vector spaces $f: \mathfrak{R}^n \rightarrow \mathfrak{R}^m$ consists of a linear transformation followed by a translation:

$$\mathbf{x} \mapsto A\mathbf{x} + \mathbf{b}$$

where $A \in \mathfrak{R}^{n \times m}$ and $\mathbf{b} \in \mathfrak{R}^m$. In the finite-dimensional case each affine transformation is given by a matrix A and a vector \mathbf{b} .

The image and pre-image of convex sets under an affine transformation defined as

$$f(\mathbf{x}) = \sum_i^n x_i a_i + b$$

yield convex sets¹⁶. Here a_i is the i^{th} row of A . The following are examples of convex sets that are either images or inverse images of convex sets under affine transformations:

1. the solution set of linear matrix inequality ($A_i, B \in S^m$)

$$\{\mathbf{x} \in \mathfrak{R}^n \mid x_1 A_1 + \dots + x_n A_n \preceq B\}$$

¹⁶Exercise: Prove.

is a convex set. Here $A \preceq B$ means $B - A$ is positive semi-definite¹⁷.

This set is the inverse image under an affine mapping of the positive semi-definite cone. That is, $f^{-1}(\text{cone}) = \{\mathbf{x} \in \Re^n \mid B - (x_1 A_1 + \dots + x_n A_n) \in \mathcal{S}_+^m\} = \{\mathbf{x} \in \Re^n \mid B \succeq (x_1 A_1 + \dots + x_n A_n)\}$.

2. hyperbolic cone ($P \in \mathcal{S}_+^n$), which is the inverse image of the norm cone $\mathcal{C}_{m+1} = \{(\mathbf{z}, u) \mid \|\mathbf{z}\| \leq u, u \geq 0, \mathbf{z} \in \Re^m\} = \{(\mathbf{z}, u) \mid \mathbf{z}^T \mathbf{z} - u^2 \leq 0, u \geq 0, \mathbf{z} \in \Re^m\}$ is a convex set. The inverse image is given by $f^{-1}(\mathcal{C}_{m+1}) = \{\mathbf{x} \in \Re^n \mid (A\mathbf{x}, \mathbf{c}^T \mathbf{x}) \in \mathcal{C}_{m+1}\} = \{\mathbf{x} \in \Re^n \mid \mathbf{x}^T A^T A \mathbf{x} - (\mathbf{c}^T \mathbf{x})^2 \leq 0\}$. Setting, $P = A^T A$, we get the equation of the hyperbolic cone:

$$\{\mathbf{x} \mid \mathbf{x}^T P \mathbf{x} \leq (\mathbf{c}^T \mathbf{x})^2, \mathbf{c}^T \mathbf{x} \geq 0\}$$

Perspective and linear-fractional functions

The perspective function $P : \Re^{n+1} \rightarrow \Re^n$ is defined as follows:

$$\begin{aligned} P : \Re^{n+1} &\rightarrow \Re^n \text{ such that} \\ P(x, t) &= x/t \qquad \text{dom } P = \{(x, t) \mid t > 0\} \end{aligned} \quad (4.29)$$

The linear-fractional function f is a generalization of the perspective function and is defined as: $\Re^n \rightarrow \Re^m$:

$$\begin{aligned} f : \Re^n &\rightarrow \Re^m \text{ such that} \\ f(\mathbf{x}) &= \frac{A\mathbf{x} + \mathbf{b}}{\mathbf{c}^T \mathbf{x} + d} \qquad \text{dom } f = \{\mathbf{x} \mid \mathbf{c}^T \mathbf{x} + d > 0\} \end{aligned} \quad (4.30)$$

The images and inverse images of convex sets under perspective and linear-fractional functions are convex¹⁸.

Consider the linear-fractional function $f = \frac{1}{x_1 + x_2 + 1}x$. Figure ?? shows an example convex set. Figure ?? shows the image of this convex set under the linear-fractional function f .

Supporting Hyperplane Theorem

On page 4.1.4, we introduced the concept of the hyperplane. For disjoint convex sets, we state the *separating hyperplane theorem*.

Theorem 67 *If \mathcal{C} and \mathcal{D} are disjoint convex sets, i.e., $\mathcal{C} \cap \mathcal{D} = \phi$, then there exists $\mathbf{a} \neq \mathbf{0}$, with a $b \in \Re$ such that*

$$\begin{aligned} \mathbf{a}^T \mathbf{x} &\leq b \text{ for } \mathbf{x} \in \mathcal{C}, \\ \mathbf{a}^T \mathbf{x} &\geq b \text{ for } \mathbf{x} \in \mathcal{D}. \end{aligned}$$

That is, the hyperplane $\{\mathbf{x} \mid \mathbf{a}^T \mathbf{x} = b\}$ separates \mathcal{C} and \mathcal{D} . The separating hyperplane need not be unique though.

¹⁷The inequality induced by positive semi-definiteness corresponds to a generalized inequality \preceq_K with $K = \mathcal{S}_+^n$.

¹⁸Exercise: Prove.

Proof: We first note that the set $\mathcal{S} = \{\mathbf{x} - \mathbf{y} | \mathbf{x} \in \mathcal{C}, \mathbf{y} \in \mathcal{D}\}$ is convex, since it is the sum of two convex sets. Since \mathcal{C} and \mathcal{D} are disjoint, $\mathbf{0} \notin \mathcal{S}$. Consider two cases:

1. Suppose $\mathbf{0} \notin \text{closure}(\mathcal{S})$. Let $\mathcal{E} = \{0\}$ and $\mathcal{F} = \text{closure}(\mathcal{S})$. Then, the euclidean distance between \mathcal{E} and \mathcal{F} , defined as

$$\text{dist}(\mathcal{E}; \mathcal{F}) = \inf \{ \|\mathbf{u} - \mathbf{v}\|_2 | \mathbf{u} \in \mathcal{E}, \mathbf{v} \in \mathcal{F} \}$$

is positive, and there exists a point $\mathbf{f} \in \mathcal{F}$ that achieves the minimum distance, i.e., $\|\mathbf{f}\|_2 = \text{dist}(\mathcal{E}, \mathcal{F})$. Define $\mathbf{a} = \mathbf{f}$, $\mathbf{b} = \|\mathbf{f}\|_2$. Then $\mathbf{a} \neq \mathbf{0}$ and the affine function $f(\mathbf{x}) = \mathbf{a}^T \mathbf{x} - b = \mathbf{f}^T (\mathbf{x} - \frac{1}{2} \mathbf{f})$ is nonpositive on \mathcal{E} and nonnegative on \mathcal{F} , i.e., that the hyperplane $\{\mathbf{x} | \mathbf{a}^T \mathbf{x} = b\}$ separates \mathcal{E} and \mathcal{F} . Thus, $\mathbf{a}^T (\mathbf{x} - \mathbf{y}) > 0$ for all $\mathbf{x} - \mathbf{y} \in \mathcal{S} \subseteq \text{closure}(\mathcal{S})$, which implies that, $\mathbf{a}^T \mathbf{x} \geq \mathbf{a}^T \mathbf{y}$ for all $\mathbf{x} \in \mathcal{C}$ and $\mathbf{y} \in \mathcal{D}$.

2. Suppose, $0 \in \text{closure}(\mathcal{S})$. Since $0 \notin \mathcal{S}$, it must be in the boundary of \mathcal{S} .

- If \mathcal{S} has empty interior, it must lie in an affine set of dimension less than n , and any hyperplane containing that affine set contains \mathcal{S} and is a hyperplane. In other words, \mathcal{S} is contained in a hyperplane $\{\mathbf{z} | \mathbf{a}^T \mathbf{z} = b\}$, which must include the origin and therefore $b = 0$. In other words, $\mathbf{a}^T \mathbf{x} = \mathbf{a}^T \mathbf{y}$ for all $\mathbf{x} \in \mathcal{C}$ and all $\mathbf{y} \in \mathcal{D}$ gives us a trivial separating hyperplane.

- If \mathcal{S} has a nonempty interior, consider the set

$$\mathcal{S}_{-\epsilon} = \{\mathbf{z} | B(\mathbf{z}, \epsilon) \subseteq \mathcal{S}\}$$

where $B(\mathbf{z}, \epsilon)$ is the Euclidean ball with center \mathbf{z} and radius $\epsilon > 0$. $\mathcal{S}_{-\epsilon}$ is the set \mathcal{S} , shrunk by ϵ . $\text{closure}(\mathcal{S}_{-\epsilon})$ is closed and convex, and does not contain $\mathbf{0}$, so as argued before, it is separated from $\{\mathbf{0}\}$ by at least one hyperplane with normal vector $\mathbf{a}(\epsilon)$ such that

$$\mathbf{a}(\epsilon)^T \mathbf{z} \geq 0 \text{ for all } \mathbf{z} \in \mathcal{S}_{-\epsilon}$$

Without loss of generality assume $\|\mathbf{a}(\epsilon)\|_2 = 1$. Let ϵ_k , for $k = 1, 2, \dots$ be a sequence of positive values of ϵ_k with $\lim_{k \rightarrow \infty} \epsilon_k = 0$. Since $\|\mathbf{a}(\epsilon_k)\|_2 = 1$ for all k , the sequence $\mathbf{a}(\epsilon_k)$ contains a convergent subsequence, and let $\bar{\mathbf{a}}$ be its limit. We have

$$\mathbf{a}(\epsilon_k)^T \mathbf{z} \geq 0 \text{ for all } \mathbf{z} \in \mathcal{S}_{-\epsilon_k}$$

and therefore $\bar{\mathbf{a}}^T \mathbf{z} \geq 0$ for all $\mathbf{z} \in \text{interior}(\mathcal{S})$, and $\bar{\mathbf{a}}^T \mathbf{z} \geq 0$ for all $\mathbf{z} \in \mathcal{S}$, which means

$$\bar{\mathbf{a}}^T \mathbf{x} \geq \bar{\mathbf{a}}^T \mathbf{y} \text{ for all } \mathbf{x} \in \mathcal{C}, \text{ and } \mathbf{y} \in \mathcal{D}.$$

□

Theorem 59 stated that the gradient evaluated at a point on a level set is orthogonal to the tangent hyperplane to the level set at that point. We now state the definition of a supporting hyperplane, which is special type of tangent hyperplane.

Definition 34 [Supporting Hyperplane]: *The supporting hyperplane to a set \mathcal{C} at a boundary point \mathbf{x}_0 is defined as $\{\mathbf{x} | \mathbf{a}^T \mathbf{x} = \mathbf{a}^T \mathbf{x}_0, \mathbf{a} \neq \mathbf{0}, \mathbf{a}^T \mathbf{y} \leq \mathbf{a}^T \mathbf{x}_0, \forall \mathbf{y} \in \mathcal{C}\}$*

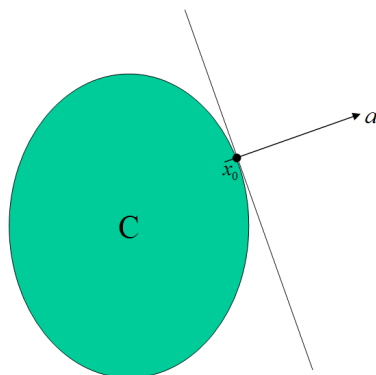


Figure 4.36: Example of a supporting hyperplane.

Figure 4.36 shows a supporting hyperplane at the point \mathbf{x}_0 on the boundary of a convex set \mathcal{C} .

For convex sets, there is an important theorem regarding supporting hyperplanes.

Theorem 68 *If the set \mathcal{C} is convex, then there exists a supporting hyperplane at every boundary point of \mathcal{C} . As in the case of the separating hyperplane, the supporting hyperplane need not be unique.*

Proof: If the interior of \mathcal{C} is nonempty, the result follows immediately by applying the separating hyperplane theorem to the sets $\{\mathbf{x}_0\}$ and $\text{interior}(\mathcal{C})$. If the interior of \mathcal{C} is empty, then \mathcal{C} must lie in an affine set of dimension less than n , and any hyperplane containing that affine set contains \mathcal{C} and \mathbf{x}_0 , and is a (trivial) supporting hyperplane. \square

4.2.7 Convex Functions

Definition 35 [Convex Function]: *A function $f : \mathcal{D} \rightarrow \Re$ is convex if \mathcal{D} is a convex set and*

$$f(\theta\mathbf{x} + (1 - \theta)\mathbf{y}) \leq \theta f(\mathbf{x}) + (1 - \theta)f(\mathbf{y}) \quad \forall \mathbf{x}, \mathbf{y} \in \mathcal{D} \quad 0 \leq \theta \leq 1 \quad (4.31)$$

Figure 4.37 illustrates an example convex function. A function $f : \mathcal{D} \rightarrow \Re$ is strictly convex if \mathcal{D} is convex and

$$f(\theta\mathbf{x} + (1 - \theta)\mathbf{y}) < \theta f(\mathbf{x}) + (1 - \theta)f(\mathbf{y}) \quad \forall \mathbf{x}, \mathbf{y} \in \mathcal{D} \quad 0 \leq \theta \leq 1 \quad (4.32)$$

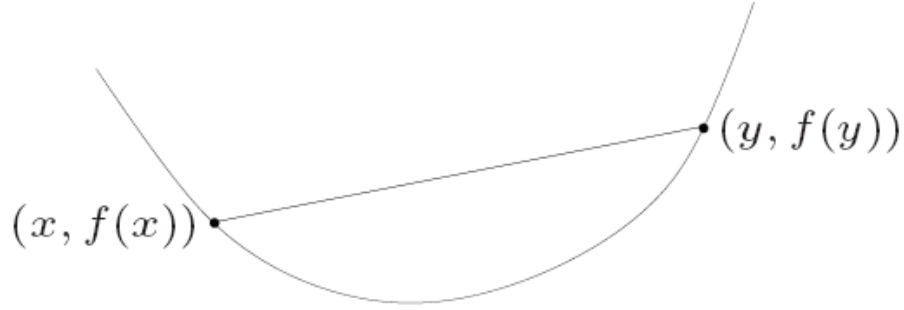


Figure 4.37: Example of convex function.

A function $f : \mathcal{D} \rightarrow \Re$ is called uniformly or strongly convex if \mathcal{D} is convex and there exists a constant $c > 0$ such that

$$f(\theta \mathbf{x} + (1 - \theta)\mathbf{y}) \leq \theta f(\mathbf{x}) + (1 - \theta)f(\mathbf{y}) - \frac{1}{2}c\theta(1 - \theta)\|\mathbf{x} - \mathbf{y}\| \quad \forall \mathbf{x}, \mathbf{y} \in \mathcal{D} \quad 0 \leq \theta \leq 1 \quad (4.33)$$

A function $f : \Re^n \rightarrow \Re$ is said to be concave if the function $-f$ is convex. Examples of convex functions on the set of reals \Re as well as on \Re^n and $\Re^{m \times n}$ are shown in Table 4.1. Examples of concave functions on the set of reals \Re are shown in Table 4.2. If a function is both convex and concave, it must be affine, as can be seen in the two tables.

Function type	Domain	Additional Constraints
The affine function: $ax + b$	\Re	Any $a, b \in \Re$
The exponential function: e^{ax}	\Re	Any $a \in \Re$
Powers: x^α	\Re_{++}	$\alpha \geq 1$ or $\alpha \leq 1$
Powers of absolute value: $ x ^p$	\Re	$p \geq 1$
Negative entropy: $x \log x$	\Re_{++}	
Affine functions of vectors: $\mathbf{a}^T \mathbf{x} + b$	\Re^n	
p-norms of vectors: $\ \mathbf{x}\ _p = \left(\sum_{i=1}^n x_i ^p \right)^{1/p}$	\Re^n	$p \geq 1$
inf norms of vectors: $\ \mathbf{x}\ _\infty = \max_k x_k $	\Re^n	
Affine functions of matrices: $\text{tr}(A^T X) + b = \sum_{i=1}^m \sum_{j=1}^n A_{ij} X_{ij} + b$	$\Re^{m \times n}$	
Spectral (maximum singular value) matrix norm: $\ X\ _2 = \sigma_{\max}(X) = (\lambda_{\max}(X^T X))^{1/2}$	$\Re^{m \times n}$	

Table 4.1: Examples of convex functions on \Re , \Re^n and $\Re^{m \times n}$.

4.2.8 Convexity and Global Minimum

One of the most fundamental and useful characteristics of convex functions is that any point of local minimum point for a convex function is also a point of global minimum.

Function type	Domain	Additional Constraints
The affine function: $ax + b$	\mathfrak{R}	Any $a, b \in \mathfrak{R}$
Powers: x^α	\mathfrak{R}_{++}	$0 \leq \alpha \leq 1$
logarithm: $\log x$	\mathfrak{R}_{++}	

Table 4.2: Examples of concave functions on \mathfrak{R} .

Theorem 69 *Let $f : \mathcal{D} \rightarrow \mathfrak{R}$ be a convex function on a convex domain \mathcal{D} . Any point of locally minimum solution for f is also a point of its globally minimum solution.*

Proof: Suppose $\mathbf{x} \in \mathcal{D}$ is a point of local minimum and let $\mathbf{y} \in \mathcal{D}$ be a point of global minimum. Thus, $f(\mathbf{y}) < f(\mathbf{x})$. Since \mathbf{x} corresponds to a local minimum, there exists an $\epsilon > 0$ such that

$$\forall \mathbf{z} \in \mathcal{D}, \|\mathbf{z} - \mathbf{x}\| \leq \epsilon \Rightarrow f(\mathbf{z}) \geq f(\mathbf{x})$$

Consider a point $\mathbf{z} = \theta\mathbf{y} + (1 - \theta)\mathbf{x}$ with $\theta = \frac{\epsilon}{2\|\mathbf{y} - \mathbf{x}\|}$. Since \mathbf{x} is a point of local minimum (in a ball of radius ϵ), and since $f(\mathbf{y}) < f(\mathbf{x})$, it must be that $\|\mathbf{y} - \mathbf{x}\| > \epsilon$. Thus, $0 < \theta < \frac{1}{2}$ and $\mathbf{z} \in \mathcal{D}$. Furthermore, $\|\mathbf{z} - \mathbf{x}\| = \frac{\epsilon}{2}$. Since f is a convex function

$$f(\mathbf{z}) \leq \theta f(\mathbf{x}) + (1 - \theta)f(\mathbf{y})$$

Since $f(\mathbf{y}) < f(\mathbf{x})$, we also have

$$\theta f(\mathbf{x}) + (1 - \theta)f(\mathbf{y}) < f(\mathbf{x})$$

The two equations imply that $f(\mathbf{z}) < f(\mathbf{x})$, which contradicts our assumption that \mathbf{x} corresponds to a point of local minimum. That is f cannot have a point of local minimum, which does not coincide with the point \mathbf{y} of global minimum. \square

Since any locally minimum point for a convex function also corresponds to its global minimum, we will drop the qualifiers ‘locally’ as well as ‘globally’ while referring to the points corresponding to minimum values of a convex function. For any strictly convex function, the point corresponding to the global minimum is also unique, as stated in the following theorem.

Theorem 70 *Let $f : \mathcal{D} \rightarrow \mathfrak{R}$ be a strictly convex function on a convex domain \mathcal{D} . Then f has a unique point corresponding to its global minimum.*

Proof: Suppose $\mathbf{x} \in \mathcal{D}$ and $\mathbf{y} \in \mathcal{D}$ with $\mathbf{y} \neq \mathbf{x}$ are two points of global minimum. That is $f(\mathbf{x}) = f(\mathbf{y})$ for $\mathbf{y} \neq \mathbf{x}$. The point $\frac{\mathbf{x} + \mathbf{y}}{2}$ also belongs to the convex set \mathcal{D} and since f is strictly convex, we must have

$$f\left(\frac{\mathbf{x} + \mathbf{y}}{2}\right) < \frac{1}{2}f(\mathbf{x}) + \frac{1}{2}f(\mathbf{y}) = f(\mathbf{x})$$

which is a contradiction. Thus, the point corresponding to the minimum of f must be unique. \square

In the following section, we state some important properties of convex functions, including relationships between convex functions and convex sets, and first and second order conditions for convexity. We will also draw relationships between the definitions of convexity and strict convexity stated here, with the definitions on page 222 for the single variable case.

4.2.9 Properties of Convex Functions

We will first extend the domain of a convex function to all \mathfrak{R}^n , while retaining its convexity and preserving its value in the domain.

Definition 36 [Extended value extension]: *If $f : \mathcal{D} \rightarrow \mathfrak{R}$, with $\mathcal{D} \subseteq \mathfrak{R}^n$ is a convex function, then we define its extended-valued extension $\tilde{f} : \mathfrak{R}^n \rightarrow \mathfrak{R}$ as*

$$\tilde{f}(\mathbf{x}) = \begin{cases} f(\mathbf{x}) & \text{if } \mathbf{x} \in \mathcal{D} \\ \infty & \text{if } \mathbf{x} \notin \mathcal{D} \end{cases} \quad (4.34)$$

In what follows, we will assume if necessary, that all convex functions are implicitly extended to the domain \mathfrak{R}^n . A useful technique for verifying the convexity of a function is to investigate its convexity, by restricting the function to a line and checking for the convexity of a function of single variable. This technique is hinged on the following theorem.

Theorem 71 *A function $f : \mathcal{D} \rightarrow \mathfrak{R}$ is (strictly) convex if and only if the function $\phi : \mathcal{D}_\phi \rightarrow \mathfrak{R}$ defined below, is (strictly) convex in t for every $\mathbf{x} \in \mathfrak{R}^n$ and for every $\mathbf{h} \in \mathfrak{R}^n$*

$$\phi(t) = f(\mathbf{x} + t\mathbf{h})$$

with the domain of ϕ given by $\mathcal{D}_\phi = \{t | \mathbf{x} + t\mathbf{h} \in \mathcal{D}\}$.

Proof: We will prove the necessity and sufficiency of the convexity of ϕ for a convex function f . The proof for necessity and sufficiency of the strict convexity of ϕ for a strictly convex f is very similar and is left as an exercise.

Proof of Necessity: Assume that f is convex. And we need to prove that $\phi(t) = f(\mathbf{x} + t\mathbf{h})$ is also convex. Let $t_1, t_2 \in \mathcal{D}_\phi$ and $\theta \in [0, 1]$. Then,

$$\begin{aligned} \phi(\theta t_1 + (1 - \theta)t_2) &= f(\theta(\mathbf{x} + t_1\mathbf{h}) + (1 - \theta)(\mathbf{x} + t_2\mathbf{h})) \\ &\leq \theta f(\mathbf{x} + t_1\mathbf{h}) + (1 - \theta)f(\mathbf{x} + t_2\mathbf{h}) = \theta\phi(t_1) + (1 - \theta)\phi(t_2) \end{aligned} \quad (4.35)$$

Thus, ϕ is convex.

Proof of Sufficiency: Assume that for every $\mathbf{h} \in \mathfrak{R}^n$ and every $\mathbf{x} \in \mathfrak{R}^n$, $\phi(t) = f(\mathbf{x} + t\mathbf{h})$ is convex. We will prove that f is convex. Let $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{D}$. Take, $\mathbf{x} = \mathbf{x}_1$ and $\mathbf{h} = \mathbf{x}_2 - \mathbf{x}_1$. We know that $\phi(t) = f(\mathbf{x}_1 + t(\mathbf{x}_2 - \mathbf{x}_1))$ is convex, with $\phi(1) = f(\mathbf{x}_2)$ and $\phi(0) = f(\mathbf{x}_1)$. Therefore, for any $\theta \in [0, 1]$

$$\begin{aligned} & f(\theta \mathbf{x}_2 + (1 - \theta)\mathbf{x}_1) = \phi(\theta) \\ & \leq \theta\phi(1) + (1 - \theta)\phi(0) \leq \theta f(\mathbf{x}_2) + (1 - \theta)f(\mathbf{x}_1) \end{aligned} \quad (4.36)$$

This implies that f is convex. \square

Next, we will draw the parallel between convex sets and convex functions by introducing the concept of the *epigraph* of a function.

Definition 37 [Epigraph]: Let $\mathcal{D} \subseteq \mathbb{R}^n$ be a nonempty set and $f : \mathcal{D} \rightarrow \mathbb{R}$. The set $\{(\mathbf{x}, f(\mathbf{x}) | \mathbf{x} \in \mathcal{D})\}$ is called graph of f and lies in \mathbb{R}^{n+1} . The epigraph of f is a subset of \mathbb{R}^{n+1} and is defined as

$$\text{epi}(f) = \{(\mathbf{x}, \alpha) | f(\mathbf{x}) \leq \alpha, \mathbf{x} \in \mathcal{D}, \alpha \in \mathbb{R}\} \quad (4.37)$$

In some sense, the epigraph is the set of points lying above the graph of f . Similarly, the hypograph of f is a subset of \mathbb{R}^{n+1} , lying above the graph of f and is defined by

$$\text{hyp}(f) = \{(\mathbf{x}, \alpha) | f(\mathbf{x}) \geq \alpha, \mathbf{x} \in \mathcal{D}, \alpha \in \mathbb{R}\} \quad (4.38)$$

There is a one to one correspondence between the convexity of function f and that of the set $\text{epi}(f)$, as stated in the following theorem.

Theorem 72 Let $\mathcal{D} \subseteq \mathbb{R}^n$ be a nonempty convex set, and $f : \mathcal{D} \rightarrow \mathbb{R}$. Then f is convex if and only if $\text{epi}(f)$ is a convex set.

Proof: Let f be convex. For any $(\mathbf{x}_1, \alpha_1) \in \text{epi}(f)$ and $(\mathbf{x}_2, \alpha_2) \in \text{epi}(f)$ and any $\theta \in (0, 1)$,

$$f(\theta \mathbf{x}_1 + (1 - \theta)\mathbf{x}_2) \leq \theta f(\mathbf{x}_1) + (1 - \theta)f(\mathbf{x}_2) \leq \theta \alpha_1 + (1 - \theta)\alpha_2$$

Since \mathcal{D} is convex, $\theta \mathbf{x}_1 + (1 - \theta)\mathbf{x}_2 \in \mathcal{D}$. Therefore, $(\theta \mathbf{x}_1 + (1 - \theta)\mathbf{x}_2, \theta \alpha_1 + (1 - \theta)\alpha_2) \in \text{epi}(f)$. Thus, $\text{epi}(f)$ is convex if f is convex. This proves the necessity part.

To prove sufficiency, assume that $\text{epi}(f)$ is convex. Let $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{D}$. So, $(\mathbf{x}_1, f(\mathbf{x}_1)) \in \text{epi}(f)$ and $(\mathbf{x}_2, f(\mathbf{x}_2)) \in \text{epi}(f)$. Since $\text{epi}(f)$ is convex, for $\theta \in (0, 1)$,

$$(\theta \mathbf{x}_1 + (1 - \theta)\mathbf{x}_2, \theta f(\mathbf{x}_1) + (1 - \theta)f(\mathbf{x}_2)) \in \text{epi}(f)$$

which implies that $f(\theta \mathbf{x}_1 + (1 - \theta)\mathbf{x}_2) \leq \theta f(\mathbf{x}_1) + (1 - \theta)f(\mathbf{x}_2)$ for any $\theta \in (0, 1)$. This proves the sufficiency. \square

There is also a correspondence between the convexity of a function and the convexity of its *sublevel sets*.

Definition 38 [Sublevel Sets]: Let $\mathcal{D} \subseteq \mathfrak{R}^n$ be a nonempty set and $f : \mathcal{D} \rightarrow \mathfrak{R}$. The set

$$L_\alpha(f) = \{\mathbf{x} | \mathbf{x} \in \mathcal{D}, f(\mathbf{x}) \leq \alpha\}$$

is called the α -sub-level set of f .

The correspondence between the convexity of f and its α -sub-level set is stated in the following theorem. Unlike the correspondence with the epigraph, the correspondence with the α -sub-level set is not one to one.

Theorem 73 Let $\mathcal{D} \subseteq \mathfrak{R}^n$ be a nonempty convex set, and $f : \mathcal{D} \rightarrow \mathfrak{R}$ be a convex function. Then $L_\alpha(f)$ is a convex set for any $\alpha \in \mathfrak{R}$.

Proof: Consider $\mathbf{x}_1, \mathbf{x}_2 \in L_\alpha(f)$. Then by definition of the level set, $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{D}$, $f(\mathbf{x}_1) \leq \alpha$ and $f(\mathbf{x}_2) \leq \alpha$. From convexity of \mathcal{D} it follows that for all $\theta \in (0, 1)$, $\mathbf{x} = \theta\mathbf{x}_1 + (1 - \theta)\mathbf{x}_2 \in \mathcal{D}$. Moreover, since f is also convex,

$$f(\mathbf{x}) \leq \theta f(\mathbf{x}_1) + (1 - \theta)f(\mathbf{x}_2) \leq \theta\alpha + (1 - \theta)\alpha = \alpha$$

which implies that $\mathbf{x} \in L_\alpha(f)$. Thus, $L_\alpha(f)$ is a convex set. \square The converse of this theorem does not hold. To illustrate this, consider the function $f(\mathbf{x}) = \frac{x_2}{1+2x_1^2}$. The 0-sublevel set of this function is $\{(x_1, x_2) \mid x_2 \leq 0\}$, which is convex. However, the function $f(\mathbf{x})$ itself is not convex.

An important property of a convex function is that it is continuous in the interior of its domain.

Theorem 74 Let $f : \mathcal{D} \rightarrow \mathfrak{R}$ be a convex function with $\mathcal{D} \subseteq \mathfrak{R}^n$ being a convex set. Let $\mathcal{S} \subset \mathcal{D}$ be an open convex set. Then f is continuous on \mathcal{S} .

Proof: Let us consider a point $\mathbf{x}_0 \in \mathcal{S}$. Since \mathcal{S} is an open convex set, we can find $n + 1$ points $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{n+1}$ such that the interior of the convex hull

$$\mathcal{C} = \left\{ \mathbf{x} | \mathbf{x} = \sum_{i=1}^{n+1} a_i \mathbf{x}_i, a_i \geq 0, \sum_1^{n+1} a_i = 1 \right\}$$

is not empty and $\mathbf{x}_0 \in \text{interior}(\mathcal{C})$. Let $M = \max_{1 \leq i \leq n+1} f(x_i)$. Then, for any

$$\mathbf{x} = \sum_{i=1}^{n+1} a_i \mathbf{x}_i \in \mathcal{C},$$

$$f(\mathbf{x}) = f\left(\sum_{i=1}^{n+1} a_i \mathbf{x}_i\right) \leq \sum_{i=1}^{n+1} a_i f(\mathbf{x}_i) \leq M$$

Since $\mathbf{x}_0 \in \mathcal{C}$, there exists a $\delta > 0$ such that $B(\mathbf{x}_0, \delta) \subset \mathcal{C}$, where, $B(\mathbf{x}_0, \delta) = \{\mathbf{x} | \|\mathbf{x} - \mathbf{x}_0\| \leq \delta\}$. Therefore, \mathbf{x}_0 can be expressed as a convex combination of $(\mathbf{x}_0 + \theta\mathbf{h}$ and $\mathbf{x}_0 - \mathbf{h}$ for some $\mathbf{h} \in B(\mathbf{x}_0, \delta)$ and some $\theta \in [0, 1]$.

$$\mathbf{x}_0 = \frac{1}{1+\theta}(\mathbf{x}_0 + \theta\mathbf{h}) + \frac{\theta}{1+\theta}(\mathbf{x}_0 - \mathbf{h})$$

Since f is convex on \mathcal{C} ,

$$f(\mathbf{x}_0) \leq \frac{1}{1+\theta}f(\mathbf{x}_0 + \theta\mathbf{h}) + \frac{\theta}{1+\theta}f(\mathbf{x}_0 - \mathbf{h})$$

From this, we can conclude that

$$f(\mathbf{x}_0 + \theta\mathbf{h}) - f(\mathbf{x}_0) \geq \theta(f(\mathbf{x}_0 - f(\mathbf{x}_0 - \mathbf{h})) \geq -\theta(M - f(\mathbf{x}_0)) \quad (4.39)$$

On the other hand,

$$f(\mathbf{x}_0 + \theta\mathbf{h}) \leq \theta f(\mathbf{x}_0 + \mathbf{h}) + (1 - \theta)f(\mathbf{x}_0)$$

which implies that

$$f(\mathbf{x}_0 + \theta\mathbf{h}) - f(\mathbf{x}_0) \leq \theta(f(\mathbf{x}_0 + \mathbf{h}) - f(\mathbf{x}_0)) \leq \theta(M - f(\mathbf{x}_0)) \quad (4.40)$$

From equations 4.39 and 4.40, we can infer that

$$|f(\mathbf{x}_0 + \theta\mathbf{h}) - f(\mathbf{x}_0)| \leq \theta|f(\mathbf{x}_0) - M|$$

For a given $\epsilon > 0$, select $\delta' \leq \delta$ such that $\delta'|f(\mathbf{x}_0) - M| \leq \epsilon\delta$. Then $\mathbf{d} = \theta\mathbf{h}$ with $\|\mathbf{h}\| = \delta$, implies that $\mathbf{d} \in B(\mathbf{x}_0, \delta)$ and $|f(\mathbf{x}_0 + \mathbf{d}) - f(\mathbf{x}_0)| \leq \epsilon$. This proves the theorem. \square

Analogous to the definition of increasing functions introduced on page number 218, we next introduce the concept of monotonic functions. This concept is very useful for characterization of a convex function.

Definition 39 Let $\mathbf{f} : \mathcal{D} \rightarrow \mathbb{R}^n$ and $\mathcal{D} \subseteq \mathbb{R}^n$. Then

1. \mathbf{f} is monotone on \mathcal{D} if for any $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{D}$,

$$(\mathbf{f}(\mathbf{x}_1) - \mathbf{f}(\mathbf{x}_2))^T (\mathbf{x}_1 - \mathbf{x}_2) \geq 0 \quad (4.41)$$

2. \mathbf{f} is strictly monotone on \mathcal{D} if for any $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{D}$ with $\mathbf{x}_1 \neq \mathbf{x}_2$,

$$(\mathbf{f}(\mathbf{x}_1) - \mathbf{f}(\mathbf{x}_2))^T (\mathbf{x}_1 - \mathbf{x}_2) > 0 \quad (4.42)$$

3. \mathbf{f} is uniformly or strongly monotone on \mathcal{D} if for any $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{D}$, there is a constant $c > 0$ such that

$$(\mathbf{f}(\mathbf{x}_1) - \mathbf{f}(\mathbf{x}_2))^T (\mathbf{x}_1 - \mathbf{x}_2) \geq c\|\mathbf{x}_1 - \mathbf{x}_2\|^2 \quad (4.43)$$

First-Order Convexity Conditions

The first order convexity condition for differentiable functions is provided by the following theorem:

Theorem 75 *Let $f : \mathcal{D} \rightarrow \Re$ be a differentiable convex function on an open convex set \mathcal{D} . Then:*

1. f is convex if and only if, for any $\mathbf{x}, \mathbf{y} \in \mathcal{D}$,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla^T f(\mathbf{x})(\mathbf{y} - \mathbf{x}) \quad (4.44)$$

2. f is strictly convex on \mathcal{D} if and only if, for any $\mathbf{x}, \mathbf{y} \in \mathcal{D}$, with $\mathbf{x} \neq \mathbf{y}$,

$$f(\mathbf{y}) > f(\mathbf{x}) + \nabla^T f(\mathbf{x})(\mathbf{y} - \mathbf{x}) \quad (4.45)$$

3. f is strongly convex on \mathcal{D} if and only if, for any $\mathbf{x}, \mathbf{y} \in \mathcal{D}$,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla^T f(\mathbf{x})(\mathbf{y} - \mathbf{x}) + \frac{1}{2}c\|\mathbf{y} - \mathbf{x}\|^2 \quad (4.46)$$

for some constant $c > 0$.

Proof:

Sufficiency: The proof of sufficiency is very similar for all the three statements of the theorem. So we will prove only for statement (4.44). Suppose (4.44) holds. Consider $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{D}$ and any $\theta \in (0, 1)$. Let $\mathbf{x} = \theta\mathbf{x}_1 + (1 - \theta)\mathbf{x}_2$. Then,

$$\begin{aligned} f(\mathbf{x}_1) &\geq f(\mathbf{x}) + \nabla^T f(\mathbf{x})(\mathbf{x}_1 - \mathbf{x}) \\ f(\mathbf{x}_2) &\geq f(\mathbf{x}) + \nabla^T f(\mathbf{x})(\mathbf{x}_2 - \mathbf{x}) \end{aligned} \quad (4.47)$$

Adding $(1 - \theta)$ times the second inequality to θ times the first, we get,

$$\theta f(\mathbf{x}_1) + (1 - \theta)f(\mathbf{x}_2) \geq f(\mathbf{x})$$

which proves that $f(\mathbf{x})$ is a convex function. In the case of strict convexity, strict inequality holds in (4.47) and it follows through. In the case of strong convexity, we need to additionally prove that

$$\theta \frac{1}{2}c\|\mathbf{x} - \mathbf{x}_1\|^2 + (1 - \theta) \frac{1}{2}c\|\mathbf{x} - \mathbf{x}_2\|^2 = \frac{1}{2}c\theta(1 - \theta)\|\mathbf{x}_2 - \mathbf{x}_1\|^2$$

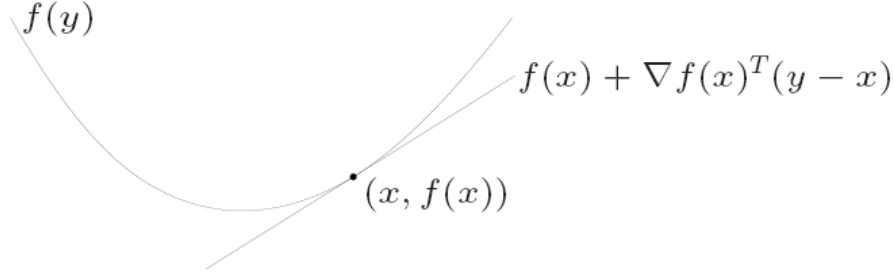


Figure 4.38: Figure illustrating Theorem 75.

Necessity: Suppose f is convex. Then for all $\theta \in (0, 1)$ and $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{D}$, we must have

$$f(\theta \mathbf{x}_2 + (1 - \theta) \mathbf{x}_1) \leq \theta f(\mathbf{x}_2) + (1 - \theta) f(\mathbf{x}_1)$$

Thus,

$$\nabla^T f(\mathbf{x}_1)(\mathbf{x}_2 - \mathbf{x}_1) = \lim_{\theta \rightarrow 0} \frac{f(\mathbf{x}_1 + \theta(\mathbf{x}_2 - \mathbf{x}_1)) - f(\mathbf{x}_1)}{\theta} \leq f(\mathbf{x}_2) - f(\mathbf{x}_1)$$

This proves necessity for (4.44). The necessity proofs for (4.45) and (4.46) are very similar, except for a small difference for the case of strict convexity; the strict inequality is not preserved when we take limits. Suppose equality does hold in the case of strict convexity, that is for a strictly convex function f , let

$$f(\mathbf{x}_2) = f(\mathbf{x}_1) + \nabla^T f(\mathbf{x}_1)(\mathbf{x}_2 - \mathbf{x}_1) \quad (4.48)$$

for some $\mathbf{x}_2 \neq \mathbf{x}_1$. Because f is strictly convex, for any $\theta \in (0, 1)$ we can write

$$f(\theta \mathbf{x}_1 + (1 - \theta) \mathbf{x}_2) = f(\mathbf{x}_2 + \theta(\mathbf{x}_1 - \mathbf{x}_2)) < \theta f(\mathbf{x}_1) + (1 - \theta) f(\mathbf{x}_2) \quad (4.49)$$

Since (4.44) is already proved for convex functions, we use it in conjunction with (4.48), and (4.49), to get

$$f(\mathbf{x}_2) + \theta \nabla^T f(\mathbf{x}_2)(\mathbf{x}_1 - \mathbf{x}_2) \leq f(\mathbf{x}_2 + \theta(\mathbf{x}_1 - \mathbf{x}_2)) < f(\mathbf{x}_2) + \theta \nabla^T f(\mathbf{x}_2)(\mathbf{x}_1 - \mathbf{x}_2)$$

which is a contradiction. Thus, equality can never hold in (4.44) for any $\mathbf{x}_1 \neq \mathbf{x}_2$. This proves the necessity of (4.45). \square

The geometrical interpretation of theorem 75 is that at any point, the linear approximation based on a local derivative gives a lower estimate of the function, *i.e.* the convex function always lies above the supporting hyperplane at that point. This is pictorially depicted in Figure 4.38. There are some implications of theorem 75 for strongly convex functions. We state them next.

Definition 40 [Some corollaries of theorem 75 for strongly convex functions]:

For a fixed \mathbf{x} , the right hand side of the inequality (4.46) is a convex quadratic function of \mathbf{y} . Thus, the critical point of the RHS should correspond to the minimum value that the RHS could take. This yields another lower bound on $f(\mathbf{y})$.

$$f(\mathbf{y}) \geq f(\mathbf{x}) - \frac{1}{2c} \|\nabla f(\mathbf{x})\|_2^2 \quad (4.50)$$

Since this holds for any $\mathbf{y} \in \mathcal{D}$, we have

$$\min_{\mathbf{y} \in \mathcal{D}} f(\mathbf{y}) \geq f(\mathbf{x}) - \frac{1}{2c} \|\nabla f(\mathbf{x})\|_2^2 \quad (4.51)$$

which can be used to bound the suboptimality of a point \mathbf{x} in terms of $\|\nabla f(\mathbf{x})\|_2$. This bound comes handy in theoretically understanding the convergence of gradient methods. If $\hat{\mathbf{y}} = \min_{\mathbf{y} \in \mathcal{D}} f(\mathbf{y})$, we can also derive a bound on the distance between any point $\mathbf{x} \in \mathcal{D}$ and the point of optimality $\hat{\mathbf{y}}$.

$$\|\mathbf{x} - \hat{\mathbf{y}}\|_2 \leq \frac{2}{c} \|\nabla f(\mathbf{x})\|_2 \quad (4.52)$$

Theorem 75 motivates the definition of the *subgradient* for non-differentiable convex functions, which has properties very similar to the gradient vector.

Definition 41 [Subgradient]: Let $f : \mathcal{D} \rightarrow \mathfrak{R}$ be a convex function defined on a convex set \mathcal{D} . A vector $\mathbf{h} \in \mathfrak{R}^n$ is said to be a subgradient of f at the point $\mathbf{x} \in \mathcal{D}$ if

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{h}^T(\mathbf{y} - \mathbf{x})$$

for all $\mathbf{y} \in \mathcal{D}$. The set of all such vectors is called the *subdifferential* of f at \mathbf{x} .

For a differentiable convex function, the gradient at point \mathbf{x} is the only subgradient at that point. Most properties of differentiable convex functions that hold in terms of the gradient also hold in terms of the subgradient for non-differentiable convex functions. Theorem 75 gives a very simple optimality criterion for a differentiable function f .

Theorem 76 Let $f : \mathcal{D} \rightarrow \mathfrak{R}$ be a convex function defined on a convex set \mathcal{D} . A point $\mathbf{x} \in \mathcal{D}$ corresponds to a minimum if and only if

$$\nabla^T f(\mathbf{x})(\mathbf{y} - \mathbf{x}) \geq 0$$

for all $\mathbf{y} \in \mathcal{D}$.

If $\nabla f(\mathbf{x})$ is nonzero, it defines a supporting hyperplane to \mathcal{D} at the point \mathbf{x} . Theorem 77 implies that for a differentiable convex function defined on an open set, every critical point must be a point of (global) minimum.

Theorem 77 *Let $f : \mathcal{D} \rightarrow \mathfrak{R}$ be differentiable and convex on an open convex domain $\mathcal{D} \subseteq \mathfrak{R}^n$. Then \mathbf{x} is a critical point of f if and only if it is a (global) minimum.*

Proof: If \mathbf{x} is a global minimum, it is a local minimum and by theorem 60, it must be a critical point and therefore $\nabla f(\mathbf{x}) = 0$. Conversely, let $\nabla f(\mathbf{x}) = 0$. By theorem 75, we know that for all $\mathbf{y} \in \mathcal{D}$,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla^T f(\mathbf{x})(\mathbf{y} - \mathbf{x})$$

Substituting $\nabla f(\mathbf{x}) = 0$ in this inequality, we get for all $\mathbf{y} \in \mathcal{D}$,

$$f(\mathbf{y}) \geq f(\mathbf{x})$$

That is, \mathbf{x} corresponds to a (global) minimum. \square

Based on the definition of monotonic functions in definition 39, we show the relationship between convexity of a function and monotonicity of its gradient in the next theorem.

Theorem 78 *Let $f : \mathcal{D} \rightarrow \mathfrak{R}$ with $\mathcal{D} \subseteq \mathfrak{R}^n$ be differentiable on the convex set \mathcal{D} . Then,*

1. *f is convex on \mathcal{D} if and only if its gradient ∇f is monotone. That is, for all $\mathbf{x}, \mathbf{y} \in \mathfrak{R}$*

$$(\nabla f(\mathbf{x}) - \nabla f(\mathbf{y}))^T (\mathbf{x} - \mathbf{y}) \geq 0 \quad (4.53)$$

2. *f is strictly convex on \mathcal{D} if and only if its gradient ∇f is strictly monotone. That is, for all $\mathbf{x}, \mathbf{y} \in \mathfrak{R}$ with $\mathbf{x} \neq \mathbf{y}$,*

$$(\nabla f(\mathbf{x}) - \nabla f(\mathbf{y}))^T (\mathbf{x} - \mathbf{y}) > 0 \quad (4.54)$$

3. *f is uniformly or strongly convex on \mathcal{D} if and only if its gradient ∇f is uniformly monotone. That is, for all $\mathbf{x}, \mathbf{y} \in \mathfrak{R}$,*

$$(\nabla f(\mathbf{x}) - \nabla f(\mathbf{y}))^T (\mathbf{x} - \mathbf{y}) \geq c\|\mathbf{x} - \mathbf{y}\|^2 \quad (4.55)$$

for some constant $c > 0$.

Proof:

Necessity: Suppose f is uniformly convex on \mathcal{D} . Then from theorem 75, we know that for any $\mathbf{x}, \mathbf{y} \in \mathcal{D}$,

$$\begin{aligned} f(\mathbf{y}) &\geq f(\mathbf{x}) + \nabla^T f(\mathbf{x})(\mathbf{y} - \mathbf{x}) - \frac{1}{2}c\|\mathbf{y} - \mathbf{x}\|^2 \\ f(\mathbf{x}) &\geq f(\mathbf{y}) + \nabla^T f(\mathbf{y})(\mathbf{x} - \mathbf{y}) - \frac{1}{2}c\|\mathbf{x} - \mathbf{y}\|^2 \end{aligned}$$

Adding the two inequalities, we get (4.55). If f is convex, the inequalities hold with $c = 0$, yielding (4.54). If f is strictly convex, the inequalities will be strict, yielding (4.54).

Sufficiency: Suppose ∇f is monotone. For any fixed $\mathbf{x}, \mathbf{y} \in \mathcal{D}$, consider the function $\phi(t) = f(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))$. By the mean value theorem applied to $\phi(t)$, we should have for some $t \in (0, 1)$,

$$\phi(1) - \phi(0) = \phi'(t) \tag{4.56}$$

Letting $\mathbf{z} = \mathbf{x} + t(\mathbf{y} - \mathbf{x})$, (4.56) translates to

$$f(\mathbf{y}) - f(\mathbf{x}) = \nabla^T f(\mathbf{z})(\mathbf{y} - \mathbf{x}) \tag{4.57}$$

Also, by definition of monotonicity of ∇f , (from (4.53)),

$$(\nabla f(\mathbf{z}) - \nabla f(\mathbf{x}))^T (\mathbf{y} - \mathbf{x}) = \frac{1}{t} (\nabla f(\mathbf{z}) - \nabla f(\mathbf{x}))^T (\mathbf{z} - \mathbf{x}) \geq 0 \tag{4.58}$$

Combining (4.57) with (4.58), we get,

$$\begin{aligned} f(\mathbf{y}) - f(\mathbf{x}) &= (\nabla f(\mathbf{z}) - \nabla f(\mathbf{x}))^T (\mathbf{y} - \mathbf{x}) + \nabla^T f(\mathbf{x})(\mathbf{y} - \mathbf{x}) \\ &\geq \nabla^T f(\mathbf{x})(\mathbf{y} - \mathbf{x}) \end{aligned} \tag{4.59}$$

By theorem 75, this inequality proves that f is convex. Strict convexity can be similarly proved by using the strict inequality in (4.58) inherited from strict monotonicity, and letting the strict inequality follow through to (4.59). For the case of strong convexity, from (4.55), we have

$$\begin{aligned} \phi'(t) - \phi'(0) &= (\nabla f(\mathbf{z}) - \nabla f(\mathbf{x}))^T (\mathbf{y} - \mathbf{x}) \\ &= \frac{1}{t} (\nabla f(\mathbf{z}) - \nabla f(\mathbf{x}))^T (\mathbf{z} - \mathbf{x}) \geq \frac{1}{t}c\|\mathbf{z} - \mathbf{x}\|^2 = ct\|\mathbf{y} - \mathbf{x}\|^2 \end{aligned} \tag{4.60}$$

Therefore,

$$\phi(1) - \phi(0) - \phi'(0) = \int_0^1 [\phi'(t) - \phi'(0)] dt \geq \frac{1}{2} c \|\mathbf{y} - \mathbf{x}\|^2 \quad (4.61)$$

which translates to

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla^T f(\mathbf{x})(\mathbf{y} - \mathbf{x}) + \frac{1}{2} c \|\mathbf{y} - \mathbf{x}\|^2$$

By theorem 75, f must be strongly convex. \square

Second Order Condition

For twice continuously differentiable convex functions the convexity condition can be characterized as follows.

Theorem 79 *A twice differential function $f : \mathcal{D} \rightarrow \Re$ for a nonempty open convex set \mathcal{D}*

1. *is convex if and only if its domain is convex and its Hessian matrix is positive semidefinite at each point in \mathcal{D} . That is*

$$\nabla^2 f(\mathbf{x}) \succeq 0 \quad \forall \mathbf{x} \in \mathcal{D} \quad (4.62)$$

2. *is strictly convex if its domain is convex and its Hessian matrix is positive definite at each point in \mathcal{D} . That is*

$$\nabla^2 f(\mathbf{x}) \succ 0 \quad \forall \mathbf{x} \in \mathcal{D} \quad (4.63)$$

3. *is uniformly convex if and only if its domain is convex and its Hessian matrix is uniformly positive definite at each point in \mathcal{D} . That is, for any $\mathbf{v} \in \Re^n$ and any $\mathbf{x} \in \mathcal{D}$, there exists a $c > 0$ such that*

$$\mathbf{v}^T \nabla^2 f(\mathbf{x}) \mathbf{v} \geq c \|\mathbf{v}\|^2 \quad (4.64)$$

In other words

$$\nabla^2 f(\mathbf{x}) \succeq c I_{n \times n}$$

where $I_{n \times n}$ is the $n \times n$ identity matrix and \succeq corresponds to the positive semidefinite inequality. That is, the function f is strongly convex iff $\nabla^2 f(\mathbf{x}) - c I_{n \times n}$ is positive semidefinite, for all $\mathbf{x} \in \mathcal{D}$ and for some constant $c > 0$, which corresponds to the positive minimum curvature of f .

Proof: We will prove only the first statement in the theorem; the other two statements are proved in a similar manner.

Necessity: Suppose f is a convex function, and consider a point $\mathbf{x} \in \mathcal{D}$. We will prove that for any $\mathbf{h} \in \Re^n$, $\mathbf{h}^T \nabla^2 f(\mathbf{x}) \mathbf{h} \geq 0$. Since f is convex, by theorem 75, we have

$$f(\mathbf{x} + t\mathbf{h}) \geq f(\mathbf{x}) + t\nabla^T f(\mathbf{x})\mathbf{h} \quad (4.65)$$

Consider the function $\phi(t) = f(\mathbf{x} + t\mathbf{h})$ considered in theorem 71, defined on the domain $\mathcal{D}_\phi = [0, 1]$. Using the chain rule,

$$\phi'(t) = \sum_{i=1}^n f_{x_i}(\mathbf{x} + t\mathbf{h}) \frac{dx_i}{dt} = \mathbf{h}^T \cdot \nabla f(\mathbf{x} + t\mathbf{h})$$

Since f has partial and mixed partial derivatives, ϕ' is a differentiable function of t on \mathcal{D}_ϕ and

$$\phi''(t) = \mathbf{h}^T \nabla^2 f(\mathbf{x} + t\mathbf{h}) \mathbf{h}$$

Since ϕ and ϕ' are continuous on \mathcal{D}_ϕ and ϕ' is differentiable on $\text{int}(\mathcal{D}_\phi)$, we can make use of the Taylor's theorem (45) with $n = 3$ to obtain:

$$\phi(t) = \phi(0) + t\phi'(0) + t^2 \cdot \frac{1}{2} \phi''(0) + O(t^3)$$

Writing this equation in terms of f gives

$$f(\mathbf{x} + t\mathbf{h}) = f(\mathbf{x}) + t\mathbf{h}^T \nabla f(\mathbf{x}) + t^2 \frac{1}{2} \mathbf{h}^T \nabla^2 f(\mathbf{x}) \mathbf{h} + O(t^3)$$

In conjunction with (4.65), the above equation implies that

$$\frac{t^2}{2} \mathbf{h}^T \nabla^2 f(\mathbf{x}) \mathbf{h} + O(t^3) \geq 0$$

Dividing by t^2 and taking limits as $t \rightarrow 0$, we get

$$\mathbf{h}^T \nabla^2 f(\mathbf{x}) \mathbf{h} \geq 0$$

Sufficiency: Suppose that the Hessian matrix is positive semidefinite at each point $\mathbf{x} \in \mathcal{D}$. Consider the same function $\phi(t)$ defined above with $\mathbf{h} = \mathbf{y} - \mathbf{x}$ for $\mathbf{y}, \mathbf{x} \in \mathcal{D}$. Applying Taylor's theorem (45) with $n = 2$ and $a = 0$, we obtain,

$$\phi(1) = \phi(0) + t\phi'(0) + t^2 \cdot \frac{1}{2} \phi''(c)$$

for some $c \in (0, 1)$. Writing this equation in terms of f gives

$$f(\mathbf{x}) = f(\mathbf{y}) + (\mathbf{x} - \mathbf{y})^T \nabla f(\mathbf{y}) + \frac{1}{2} (\mathbf{x} - \mathbf{y})^T \nabla^2 f(\mathbf{z}) (\mathbf{x} - \mathbf{y})$$

where $\mathbf{z} = \mathbf{y} + c(\mathbf{x} - \mathbf{y})$. Since \mathcal{D} is convex, $\mathbf{z} \in \mathcal{D}$. Thus, $\nabla^2 f(\mathbf{z}) \succeq 0$. It follows that

$$f(\mathbf{x}) \geq f(\mathbf{y}) + (\mathbf{x} - \mathbf{y})^T \nabla f(\mathbf{y})$$

By theorem 75, the function f is convex. \square

Examples of differentiable/twice differentiable convex functions, along with the value of their respective gradients/hessians are tabulated in Table 4.3.

Function type	Constraints	Gradient/Hessian
Quadratic : $\frac{1}{2}\mathbf{x}^T A\mathbf{x} + \mathbf{b}^T \mathbf{x} + c$	$A \succeq 0$	$\nabla^2 f(\mathbf{x}) = P$
Quadratic over linear: $\frac{x^2}{y} \geq 0$	$y > 0$	$\nabla^2 f(x, y) = \frac{2}{y^3} \begin{bmatrix} y^2 & -xy \\ -xy & x^2 \end{bmatrix}$
Log-sum-exp: $\log \sum_{k=1}^n \exp(x_k)$		$\nabla^2 f(\mathbf{x}) = \frac{1}{(\mathbf{1}^T \mathbf{z})^2} ((\mathbf{1}^T \mathbf{z}) \mathbf{diag}(\mathbf{z}) - \mathbf{z}\mathbf{z}^T)$ where $\mathbf{z} = [e^{x_1}, e^{x_2}, \dots, e^{x_n}]$
Negative Geometric mean: $-\left(\prod_{k=1}^n x_k\right)^{\frac{1}{n}}$	$\mathbf{x} \in \mathfrak{R}_{++}^n$	$\nabla^2 f(\mathbf{x}) = \frac{\prod_{i=1}^n nx_i^{1/n}}{n^2} \left(n \mathbf{diag}\left(\frac{1}{x_1^2}, \dots, \frac{1}{x_n^2}\right) - qq^T \right)$

Table 4.3: Examples of twice differentiable convex functions on \mathfrak{R} .

4.2.10 Convexity Preserving Operations on Functions

In practice if you want to establish the convexity of a function f , you could either

1. Prove it from first principles, i.e., using the definition of convexity or
2. If f is twice differentiable, show that $\nabla^2 f(x) \succeq 0$
3. Show that f is obtained from simple convex functions by operations that preserve convexity. Following are operations on functions that preserve convexity (proofs omitted, since they are trivial):

- **Nonnegative weighted sum:** $f = \sum_{i=1}^n \alpha_i f_i$ is convex if each f_i for $1 \leq i \leq n$ is convex and $\alpha_i \geq 0, 1 \leq i \leq n$.
- **Composition with affine function:** $f(Ax + b)$ is convex if f is convex. For example:
 - The log barrier for linear inequalities, $f(x) = -\sum_{i=1}^m \log(b_i - a_i^T x)$, is convex since $-\log(x)$ is convex.
 - Any norm of an affine function, $f(x) = \|Ax + b\|$, is convex.
- **Pointwise maximum:** If f_1, f_2, \dots, f_m are convex, then $f(x) = \max\{f_1(x), f_2(x), \dots, f_m(x)\}$ is also convex, For example:

- Sum of r largest components of $\mathbf{x} \in \Re^n$ $f(\mathbf{x}) = x_{[1]} + x_{[2]} + \dots + x_{[r]}$, where $x_{[i]}$ is the i^{th} largest component of \mathbf{x} , is a convex function.
- **Pointwise supremum:** If $f(x, y)$ is convex in x for every $y \in \mathcal{S}$, then $g(x) = \sup_{y \in \mathcal{S}} f(x, y)$ is convex. For example:
 - The function that returns the maximum eigenvalue of a symmetric matrix X , viz., $\lambda_{\max}(X) = \sup_{y \in \mathcal{S}} f(x, y)$ is a convex function of the symmetric matrix X .
- **Composition with functions:** Let $h : \Re^k \rightarrow \Re$ with $h(x) = \infty, \forall x \notin \text{dom } h$ and $g : \Re^n \rightarrow \Re^k$. Define $f(x) = h(g(x))$. f is convex if
 - g_i is convex, h is convex and nondecreasing in each argument
 - or g_i is concave, h is convex and nonincreasing in each argument

Some examples illustrating this property are:

- *exp* $g(x)$ is convex if g is convex
- $\sum_{i=1}^m \log g_i(x)$ is concave if g_i are concave and positive
- $\log \sum_{i=1}^m \exp g_i(x)$ is convex if g_i are convex
- $1/g(x)$ is convex if g is concave and positive
- **Infimum:** If $f(x, y)$ is convex in (x, y) and \mathcal{C} is a convex set, then $g(x) = \inf_{y \in \mathcal{C}} f(x, y)$ is convex. For example:
 - Let $f(x, \mathcal{S})$ that returns the distance of a point x to a convex set \mathcal{S} . That is $f(x, \mathcal{S}) = \inf_{y \in \mathcal{S}} \|x - y\|$. Then $f(x, \mathcal{S})$ is a convex.
- **Perspective Function:** The perspective of a function $f : \Re^n \rightarrow \Re$ is the function $g : \Re^n \times \Re \rightarrow \Re$, $g(x, t) = tf(x/t)$. Function g is convex if f is convex on $\text{dom } g = \{(x, t) | x/t \in \text{dom } f, t > 0\}$. For example,
 - The perspective of $f(x) = x^T x$ is (quadratic-over-linear) function $g(x, t) = \frac{x^T x}{t}$ and is convex.
 - The perspective of negative logarithm $f(x) = -\log x$ is the relative entropy function $g(x, t) = t \log t - t \log x$ and is convex.

4.3 Convex Optimization Problem

Formally, a convex program is defined as

$$\min_{\mathbf{x} \in \mathcal{X}} c^T \mathbf{x} \tag{4.66}$$

where $\mathcal{X} \subset \mathfrak{R}^n$ is a convex set and \mathbf{x} is a vector of n optimization or decision variables. In applications, convex optimization programs usually arise in the form:

$$\begin{array}{ll}
 \text{minimize} & f(\mathbf{x}) \\
 \text{subject to} & g_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m \\
 & A\mathbf{x} = \mathbf{b} \\
 \text{variable } \mathbf{x} & = (x_1, \dots, x_n)
 \end{array} \tag{4.67}$$

If it is given that the functions f, g_1, \dots, g_m are convex, by theorem 73, the feasible set \mathcal{X} of this problem, which is the intersection of a finite number of 0-sub-level sets of convex functions is also convex. Therefore, this problem can be posed as the following convex optimization problem:

$$\begin{array}{ll}
 \min_{x=(t,u) \in \mathcal{X}} & t \\
 \mathcal{X} = \{(t, u) \mid & f(u) \leq t, g_1(u) \leq 0, g_2(u) \leq 0, \dots, g_m(u) \leq 0\}
 \end{array} \tag{4.68}$$

The set \mathcal{X} is convex, and hence the problem in (4.68) is a convex optimization problem. Further, every locally optimal point is also globally optimal. The computation time of algorithms for solving convex optimization problems is roughly proportional to $\max(n^2, n^2m, C)$, where C is the cost of evaluating f , the g_i 's and their first and second derivatives. There are many reliable and efficient algorithms for solving convex optimization problems. However, it is often difficult to recognize convex optimization problems in practice.

Examples

Consider the optimization problem

$$\begin{array}{ll}
 \text{minimize} & f(\mathbf{x}) = x_1^2 + x_2^2 \\
 \text{subject to} & g_1(\mathbf{x}) = \frac{x_1}{1+x_2^2} \leq 0 \\
 & h(\mathbf{x}) = (x_1 + x_2)^2 = 0
 \end{array} \tag{4.69}$$

We note that the optimization problem above is not a convex problem according to our definition, since g_1 is not convex and h is not affine. However, we note that the feasible set $\{(x_1, x_2) \mid x_1 = -x_2, x_1 \leq 0\}$ is convex (recall that the converse of theorem 73 does not hold - the 0-sublevel set of a non convex function can be convex). This problem can be posed as an equivalent (but not identical) convex optimization problem:

$$\begin{aligned}
&\text{minimize} && f(\mathbf{x}) = x_1^2 + x_2^2 \\
&\text{subject to} && x_1 \leq 0 \\
&&& x_1 + x_2 = 0
\end{aligned} \tag{4.70}$$

4.4 Duality Theory

Duality is a very important component of nonlinear and linear optimization models. It has a wide spectrum of applications that are very popular. It arises in the basic form of linear programming as well as in interior point methods for linear programming. The duality in linear programming was first observed by Von Neumann, and later formalized by Tucker, Gale and Kuhn. In the first attempt at extending duality beyond linear programs, duals of quadratic programs were next developed. It was subsequently observed that you can always write a dual for any optimization problem and the modern Lagrange-based ‘constructive’¹⁹ duality theory followed in the late 1960s.

An extremely popular application of duality happens to be in the quadratic programming for Support Vector Machines. The primal and dual both happen to be convex optimization programs in this case. The Minimax theorem²⁰, a fundamental theorem of Game Theory, proved by John von Neumann in 1928, is but one instance of the general duality theory. In the consideration of equilibrium in electrical networks, current are ‘primal variables’ and the potential differences are the ‘dual variables’. In models of economic markets, the ‘primal’ variables are production and consumption levels while the ‘dual’ variables are prices (of goods, *etc.*). Dual price-based decomposition methods were developed by Danzig. In the case of thrust structures in mechanics, forces are primal variables and the displacements are the dual variables. Dual problems and their solutions are used for proving optimality of solutions, finding near-optimal solutions, analysing how sensitive the solution of the primal is to perturbations in the right hand side of constraints, analysing convergence of algorithms, *etc.*

4.4.1 Lagrange Multipliers

Consider the following quadratic function of $\mathbf{x} \in \mathfrak{R}^n$.

$$F(\mathbf{x}) = \mathbf{x}^T A \mathbf{x} - \mathbf{x}^T \mathbf{b} \tag{4.71}$$

where A is an $n \times n$ square matrix. Consider the unconstrained minimization problem

¹⁹As we will see, the theory helps us construct duals that are useful in practice.

²⁰The name Minimax was invented by Tucker.

$$\min_{\mathbf{x} \in \mathcal{D}} F(\mathbf{x}) \quad (4.72)$$

A locally optimum solution $\hat{\mathbf{x}}$ to this objective can be obtained by setting $\nabla F(\hat{\mathbf{x}}) = \mathbf{0}$. This condition translates to $A\hat{\mathbf{x}} = \mathbf{b}$. A sufficient condition for $\hat{\mathbf{x}}$ to be a point of local minimum is that $\nabla^2 F(\hat{\mathbf{x}}) \succ 0$. This condition holds *iff*, $A \succ 0$, that is, A is a positive definite matrix. Given that $A \succ 0$, A must be invertible (*c.f.* Section 3.12.2) and the unique solution is $\mathbf{x} = A^{-1}\mathbf{b}$.

Now suppose we have a constrained minimization problem

$$\begin{aligned} \min_{\mathbf{y} \in \mathfrak{R}^n} \quad & \frac{1}{2}\mathbf{y}^T B \mathbf{y} \\ \text{subject to} \quad & A^T \mathbf{y} = \mathbf{b} \end{aligned} \quad (4.73)$$

where $\mathbf{y} \in \mathfrak{R}^n$, A is an $n \times m$ matrix, B is an $n \times n$ matrix and \mathbf{b} is a vector of size m . To handle constrained minimization, let us consider minimization of the modified objective function $L(\mathbf{y}, \lambda) = \frac{1}{2}\mathbf{y}^T B \mathbf{y} + \lambda^T (A^T \mathbf{y} - \mathbf{b})$.

$$\min_{\mathbf{y} \in \mathfrak{R}^n, \lambda \in \mathfrak{R}^m} \quad \frac{1}{2}\mathbf{y}^T B \mathbf{y} + \lambda^T (A^T \mathbf{y} - \mathbf{b}) \quad (4.74)$$

The function $L(\mathbf{y}, \lambda)$ is called the lagrangian and involves the lagrange multiplier $\lambda \in \mathfrak{R}^m$. A sufficient condition for optimality of $L(\mathbf{y}, \lambda)$ at a point $L(\mathbf{y}^*, \lambda^*)$ is that $\nabla L(\mathbf{y}^*, \lambda^*) = \mathbf{0}$ and $\nabla^2 L(\mathbf{y}^*, \lambda^*) \succ 0$. For this particular problem:

$$\nabla L(\mathbf{y}^*, \lambda^*) = \begin{bmatrix} B\mathbf{y}^* + A\lambda^* \\ A^T \mathbf{y}^* - \mathbf{b} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

and

$$\nabla^2 L(\mathbf{y}^*, \lambda^*) = \begin{bmatrix} B & A \\ A^T & 0 \end{bmatrix} \succ 0$$

The point $(\mathbf{y}^*, \lambda^*)$ must therefore satisfy, $A^T \mathbf{y}^* = \mathbf{b}$ and $A\lambda^* = -B\mathbf{y}^*$. If B is taken to be the identity matrix, $n = 2$ and $m = 1$, the minimization problem (4.73) amounts to finding a point \mathbf{y}^* on a line $a_{11}y_1 + a_{12}y_2 = b$ that is closest to the origin. From geometry, we know that the point on a line closest to the origin is the point of intersection \mathbf{p}^* of a perpendicular from the origin to the line. On the other hand, the solution for the minimum of (4.74), for these conditions coincides with \mathbf{p}^* and is given by:

$$\begin{aligned} y_1 &= \frac{a_{11}b}{(a_{11})^2 + (a_{12})^2} \\ y_2 &= \frac{a_{12}b}{(a_{11})^2 + (a_{12})^2} \end{aligned}$$

That is, for $n = 2$ and $m = 1$, the solution to (4.74) is the same as the solution to (4.72). Can this construction be used to always find optimal solutions to a minimization problem? We will answer this question by first motivating the concept of lagrange multipliers and in Section 4.4.2, we will formalize the lagrangian dual.

Lagrange Multipliers with Equality Constraints

The concept of lagrange multipliers can be attributed to the mathematician Lagrange, who was born in the year 1736 in Turin. He largely worked on mechanics, the calculus of variations probability, group theory, and number theory. He was party to the choice of base 10 for the metric system (rather than 12). We will here give a brief introduction to lagrange multipliers; Section 4.4.2 will discuss the *Karush-Kuhn-Tucker conditions*, which are a generalization of lagrange multipliers.

Consider the equality constrained minimization problem (with $\mathcal{D} \subseteq \Re^n$)

$$\begin{aligned} \min_{\mathbf{x} \in \mathcal{D}} \quad & f(\mathbf{x}) \\ \text{subject to} \quad & g_i(\mathbf{x}) = 0 \quad i = 1, 2, \dots, m \end{aligned} \tag{4.75}$$

A direct approach to solving this problem is to find a parametrization of the constraints (as in the example on page 228) such that f is expressed in terms of the parameters, to give an unconstrained problem. For example if there is a single constraint of the form $\mathbf{x}^T \mathbf{A} \mathbf{x} = k$, and $A \succ 0$, then the coordinate system can be rotated and \mathbf{x} can be rescaled so that we get the constraint $\mathbf{y}' \mathbf{y} = k$. Further, we can substitute with parametrization of the y_i 's as

$$\begin{aligned} y_1 &= k \sin \theta_1 \sin \theta_2 \dots \sin \theta_{n-1} \\ y_2 &= k \sin \theta_1 \sin \theta_2 \dots \cos \theta_{n-1} \\ &\dots \dots \dots \end{aligned}$$

However, this is not possible for general constraints. The method of lagrange multipliers presents an indirect approach to solving this problem.

Consider a schematic representation of the problem in (4.75) with a single constraint, *i.e.*, $m = 1$ in Figure 4.39. The figure shows some level curves of the function f . The constraint function g_1 is also plotted with dotted lines in the same figure. The gradient of the constraint ∇g_1 is not parallel to the gradient ∇f of the function²¹ at $f = 10.4$; it is therefore possible to move along the constraint surface so as to further reduce f . However, as shown in Figure 4.39, ∇g_1 and ∇f are parallel at $f = 10.3$, and any motion along $g_1(\mathbf{x}) = 0$ will

²¹Note that the (negative) gradient at a point is orthogonal to the contour line going through that point. This was proved in Theorem 59.

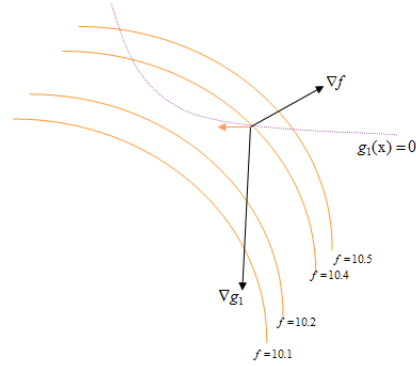


Figure 4.39: At any non-optimal and non-saddle point of the equality constrained problem, the gradient of the constraint will not be parallel to that of the function.

increase f , or leave it unchanged. Hence, at the solution \mathbf{x}^* , $\nabla f(\mathbf{x}^*)$ must be proportional to $-\nabla g_1(\mathbf{x}^*)$, yielding, $\nabla f(\mathbf{x}^*) = -\lambda \nabla g_1(\mathbf{x}^*)$, for some constant $\lambda \in \mathfrak{R}$; λ is called a *Lagrange multiplier*. In several problems, the value of λ itself need never be computed and therefore λ is often qualified as the *undetermined* lagrange multiplier.

The necessary condition for an optimum at \mathbf{x}^* for the optimization problem in (4.75) with $m = 1$ can be stated as in (4.76), where the gradient is now $n + 1$ dimensional with its last component being a partial derivative with respect to λ .

$$\nabla L(\mathbf{x}^*, \lambda^*) = \nabla f(\mathbf{x}^*) + \lambda^* \nabla g_1(\mathbf{x}^*) = 0 \quad (4.76)$$

The solutions to (4.76) are the stationary points of the lagrangian L ; they are not necessarily local extrema of L . L is unbounded: given a point \mathbf{x} that doesn't lie on the constraint, letting $\lambda \rightarrow \pm\infty$ makes L arbitrarily large or small. However, under certain stronger assumptions, as we shall see in Section 4.4.2, if the *strong Lagrangian principle* holds, the minima of f minimize the Lagrangian globally.

We will extend the necessary condition for optimality of a minimization problem with single constraint to minimization problems with multiple equality constraints (*i.e.*, $m > 1$. in (4.75)). Let \mathcal{S} be the subspace spanned by $\nabla g_i(\mathbf{x})$ at any point \mathbf{x} and let \mathcal{S}_\perp be its orthogonal complement. Let $(\nabla f)_\perp$ be the component of ∇f in the subspace \mathcal{S}_\perp . At any solution \mathbf{x}^* , it must be true that the gradient of f has $(\nabla f)_\perp = 0$ (*i.e.*, no components that are perpendicular to all of the ∇g_i), because otherwise you could move \mathbf{x}^* a little in that direction (or in the opposite direction) to increase (decrease) f without changing any of the g_i , *i.e.* without violating any constraints. Hence for multiple equality constraints, it must be true that at the solution \mathbf{x}^* , the space \mathcal{S} contains the

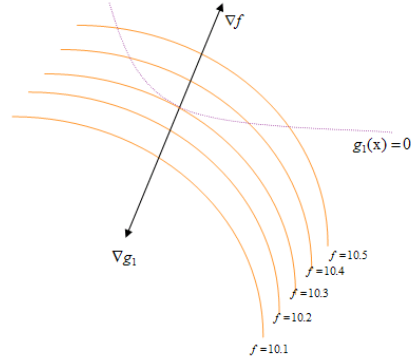


Figure 4.40: At the equality constrained optimum, the gradient of the constraint must be parallel to that of the function.

vector ∇f , *i.e.*, there are some constants λ_i such that $\nabla f(\mathbf{x}^*) = \lambda_i \nabla g_i(\mathbf{x}^*)$. We also need to impose that the solution is on the correct constraint surface (*i.e.*, $g_i = 0, \forall i$). In the same manner as in the case of $m = 1$, this can be encapsulated by introducing the Lagrangian $L(\mathbf{x}, \lambda) = f(\mathbf{x}) - \sum_{i=1}^m \lambda_i g_i(\mathbf{x})$, whose gradient with respect to both \mathbf{x} , and λ vanishes at the solution.

This gives us the following necessary condition for optimality of (4.75):

$$\nabla L(\mathbf{x}^*, \lambda^*) = \nabla \left(f(\mathbf{x}) - \sum_{i=1}^m \lambda_i g_i(\mathbf{x}) \right) = 0 \quad (4.77)$$

Lagrange Multipliers with Inequality Constraints

Instead of a single equality constraint $g_1(\mathbf{x}) = 0$, we could have a single inequality constraint $g_1(\mathbf{x}) \leq 0$. The entire region labeled $g_1(\mathbf{x}) \leq 0$ in Figure 4.41 then becomes feasible. At the solution \mathbf{x}^* , if $g_1(\mathbf{x}^*) = 0$, *i.e.*, if the constraint is active, we must have (as in the case of a single equality constraint) that ∇f is parallel to ∇g_1 , by the same argument as before. Additionally, it is necessary that the two gradients must point in opposite directions; otherwise a move away from the surface $g_1 = 0$ and into the feasible region would further reduce f . Since we are minimizing f , if the Lagrangian is written as $L = f + \lambda g_1$, we must have $\lambda \geq 0$. Therefore, with an inequality constraint, the sign of λ is important, and $\lambda \geq 0$ becomes a constraint.

However, if the constraint is not active at the solution $\nabla f(\mathbf{x}^*) = 0$, then removing g_1 makes no difference and we can drop it from $L = f + \lambda g_1$, which is equivalent to setting $\lambda = 0$. Thus, whether or not the constraints $g_i = 0$ are active, we can find the solution by requiring that the gradients of the Lagrangian vanish, and also requiring that $\lambda g_i(\mathbf{x}^*) = 0$. This latter condition is one of the

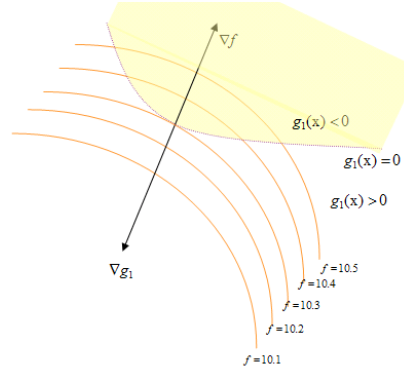


Figure 4.41: At the inequality constrained optimum, the gradient of the constraint must be parallel to that of the function.

important Karush-Kuhn-Tucker conditions of convex optimization theory that can facilitate the search for the solution and will be more formally discussed in Section 4.4.2.

Now consider the general inequality constrained minimization problem

$$\begin{aligned} \min_{\mathbf{x} \in \mathcal{D}} \quad & f(\mathbf{x}) \\ \text{subject to} \quad & g_i(\mathbf{x}) \leq 0 \quad i = 1, 2, \dots, m \end{aligned} \quad (4.78)$$

With multiple inequality constraints, for constraints that are active, as in the case of multiple equality constraints, ∇f must lie in the space spanned by the ∇g_i 's, and if the Lagrangian is $L = f + \sum_{i=1}^m \lambda_i g_i$, then we must additionally have $\lambda_i \geq 0, \forall i$ (since otherwise f could be reduced by moving into the feasible region). As for an inactive constraint g_j ($g_j < 0$), removing g_j from L makes no difference and we can drop ∇g_j from $\nabla f = -\sum_{i=1}^m \lambda_i \nabla g_i$ or equivalently set $\lambda_j = 0$. Thus, the above KKT condition generalizes to $\lambda_i g_i(\mathbf{x}^*) = 0, \forall i$. The necessary condition for optimality of (4.78) is summarily given as

$$\begin{aligned} \nabla L(\mathbf{x}^*, \lambda^*) = \nabla \left(f(\mathbf{x}) - \sum_{i=1}^m \lambda_i g_i(\mathbf{x}) \right) = 0 \\ \forall i \quad \lambda_i g_i(\mathbf{x}) = 0 \end{aligned} \quad (4.79)$$

A simple and often useful trick called the *free constraint gambit* is to solve ignoring one or more of the constraints, and then check that the solution satisfies those constraints, in which case you have solved the problem.

4.4.2 The Dual Theory for Constrained Optimization

Consider the general inequality constrained minimization problem in (4.78), restated below.

$$\begin{aligned} \min_{\mathbf{x} \in \mathcal{D}} \quad & f(\mathbf{x}) \\ \text{subject to} \quad & g_i(\mathbf{x}) \leq 0, \quad i = 1, 2, \dots, m \end{aligned} \quad (4.80)$$

There are three simple and straightforward steps in forming a dual problem.

1. The first step involves forming the lagrange function by associating a price λ_i , called a lagrange multiplier, with the constraint involving g_i .

$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \sum_{i=1}^n \lambda_i g_i(\mathbf{x}) = f(\mathbf{x}) + \lambda^T \mathbf{g}(\mathbf{x})$$

2. The second step is the construction of the dual function $L^*(\lambda)$ which is defined as:

$$L^*(\lambda) = \min_{\mathbf{x} \in \mathcal{D}} L(\mathbf{x}, \lambda) = \min_{\mathbf{x} \in \mathcal{D}} f(\mathbf{x}) + \lambda^T \mathbf{g}(\mathbf{x})$$

What makes the theory of duality constructive is when we can solve for L^* efficiently - either in a closed form or some other ‘simple’ mechanism. If L^* is not easy to evaluate, the duality theory will be less useful.

3. We finally define the dual problem:

$$\begin{aligned} \max_{\lambda \in \mathbb{R}^m} \quad & L^*(\lambda) \\ \text{subject to} \quad & \lambda \geq \mathbf{0} \end{aligned} \quad (4.81)$$

It can be immediately proved that the dual problem is a concave maximization problem.

Theorem 80 *The dual function $L^*(\lambda)$ is concave.*

Proof: Consider two values of the dual variables, *viz.*, $\lambda_1 \geq \mathbf{0}$ and $\lambda_2 \geq \mathbf{0}$. Let $\lambda = \theta\lambda_1 + (1 - \theta)\lambda_2$ for any $\theta \in [0, 1]$. Then,

$$\begin{aligned} L^*(\lambda) &= \min_{\mathbf{x} \in \mathcal{D}} f(\mathbf{x}) + \lambda^T g(\mathbf{x}) \\ &= \min_{\mathbf{x} \in \mathcal{D}} \theta [f(\mathbf{x}) + \lambda_1^T g(\mathbf{x})] + (1 - \theta) [f(\mathbf{x}) + \lambda_2^T g(\mathbf{x})] \\ &\geq \min_{\mathbf{x} \in \mathcal{D}} \theta [f(\mathbf{x}) + \lambda_1^T g(\mathbf{x})] + \min_{\mathbf{x} \in \mathcal{D}} (1 - \theta) [f(\mathbf{x}) + \lambda_2^T g(\mathbf{x})] \\ &= \theta L^*(\lambda_1) + (1 - \theta) L^*(\lambda_2) \end{aligned}$$

This proves that $L^*(\lambda)$ is a concave function. \square

The dual is concave (or the negative of the dual is convex) irrespective of the primal. Solving the dual is therefore always a convex programming problem. Thus, in some sense, the dual is better structured than the primal. However, the dual cannot be drastically simpler than the primal. For example, if the primal is not an LP, the dual cannot be an LP. Similarly, the dual can be quadratic only if the primal is quadratic.

A tricky thing in duality theory is to decide what we call the domain or *ground set* \mathcal{D} and what we call the constraints g_i 's. Based on whether constraints are explicitly stated or implicitly stated in the form of the ground set, the dual problem could be very different. Thus, many duals are possible for the given primal.

We will look at two examples to give a flavour of how the duality theory works.

1. We will first look at linear programming.

$$\begin{array}{ll} \min_{\mathbf{x} \in \mathbb{R}^n} & \mathbf{c}^T \mathbf{x} \\ \text{subject to} & -A\mathbf{x} + \mathbf{b} \leq \mathbf{0} \end{array}$$

The lagrangian for this problem is:

$$L(\mathbf{x}, \lambda) = \mathbf{c}^T \mathbf{x} + \lambda^T \mathbf{b} - \lambda^T A\mathbf{x} = \mathbf{b}^T \lambda + (\mathbf{c}^T - A^T \lambda) \mathbf{x}$$

The next step is to get L^* , which we obtain using the first derivative test:

$$L^*(\lambda) = \min_{\mathbf{x} \in \mathbb{R}^n} \mathbf{b}^T \lambda + (\mathbf{c}^T - A\lambda)^T \mathbf{x} = \begin{cases} \mathbf{b}^T \lambda & \text{if } A^T \lambda = \mathbf{c} \\ -\infty & \text{if } A^T \lambda \neq \mathbf{c} \end{cases}$$

The function L^* can be thought of as the extended value extension of the same function restricted to the domain $\{\lambda | A^T \lambda = \mathbf{c}\}$. Therefore, the dual problem can be formulated as:

$$\begin{array}{ll} \max_{\lambda \in \mathbb{R}^m} & \mathbf{b}^T \lambda \\ \text{subject to} & A^T \lambda = \mathbf{c} \\ & \lambda \geq \mathbf{0} \end{array} \quad (4.82)$$

This is the dual of the standard LP. What if the original LP was the following?

$$\begin{aligned} \min_{\mathbf{x} \in \mathfrak{R}^n} \quad & \mathbf{c}^T \mathbf{x} \\ \text{subject to} \quad & -A\mathbf{x} + \mathbf{b} \leq \mathbf{0} \quad \mathbf{x} \geq \mathbf{0} \end{aligned}$$

Now we have a variety of options based on what constraints are introduced into the ground set (or domain) and what are explicitly treated as constraints. Some working out will convince us that treating $\mathbf{x} \in \mathfrak{R}^n$ as the constraint and the explicit constraints as part of the ground set is a very bad idea. One dual for this problem is the same as (4.82).

2. Let us look at a modified version of the problem in (4.83).

$$\begin{aligned} \min_{\mathbf{x} \in \mathfrak{R}^n} \quad & \mathbf{c}^T \mathbf{x} - \sum_{i=1}^n \ln x_i \\ \text{subject to} \quad & -A\mathbf{x} + \mathbf{b} = \mathbf{0} \\ & \mathbf{x} > \mathbf{0} \end{aligned}$$

Typically, when we try to formulate a dual problem, we look for constraints that get in the way of conveniently solving the problem. We first formulate the lagrangian for this problem.

$$L(\mathbf{x}, \lambda) = \mathbf{c}^T \mathbf{x} - \sum_{i=1}^n \ln x_i + \lambda^T \mathbf{b} - \lambda^T A\mathbf{x} = \mathbf{b}^T \lambda + \mathbf{x}^T (\mathbf{c} - A^T \lambda) - \sum_{i=1}^n \ln x_i$$

The domain (or ground set) for this problem is $\mathbf{x} > \mathbf{0}$, which is open.

The expression for L^* can be obtained using the first derivative test, while keeping in mind that L can be made arbitrarily small (tending to $-\infty$) unless $(\mathbf{c} - A^T \lambda) > \mathbf{0}$. This is because, even if one component of $\mathbf{c} - A^T \lambda$ is less than or equal to zero, the value of L can be made arbitrarily small by decreasing the value of the corresponding component of \mathbf{x} in the $\sum_{i=1}^n \ln x_i$ part. Further, the sum $\mathbf{b}^T \lambda + (\mathbf{c} - A^T \lambda)^T \mathbf{x} - \sum_{i=1}^n \ln x_i$ can be separated out into the individual components of λ_i , and this can be exploited while determining the critical point of L .

$$L^*(\lambda) = \min_{\mathbf{x} > \mathbf{0}} \mathbf{b}^T \lambda + n + \sum_{i=1}^n \ln \frac{1}{(\mathbf{c} - A^T \lambda)_i} = \begin{cases} \mathbf{b}^T \lambda & \text{if } (\mathbf{c} - A^T \lambda) > \mathbf{0} \\ -\infty & \text{otherwise} \end{cases}$$

Finally, the dual will be

$$\begin{aligned} \max_{\lambda \in \mathfrak{R}^m} \quad & \mathbf{b}^T \lambda + n + \sum_{i=1}^n \ln \frac{1}{(\mathbf{c} - A^T \lambda)_i} \\ \text{subject to} \quad & -A^T \lambda + \mathbf{c} > \mathbf{0} \end{aligned}$$

As noted earlier, the theory of duality remains a theory unless the dual lends itself to some constructive evaluation; not always is the dual a useful form.

The following *Weak duality theorem* states an important relationship between solutions to the primal (4.80) and the dual (4.81) problems.

Theorem 81 *If $p^* \in \mathfrak{R}$ is the solution to the primal problem in (4.80) and $d^* \in \mathfrak{R}$ is the solution to the dual problem in (4.81), then*

$$p^* \geq d^*$$

In general, if $\hat{\mathbf{x}}$ is any feasible solution to the primal problem (4.80) and $\hat{\lambda}$ is a feasible solution to the dual problem (4.81), then

$$f(\hat{\mathbf{x}}) \geq L^*(\hat{\lambda})$$

Proof: If $\hat{\mathbf{x}}$ is a feasible solution to the primal problem (4.80) and $\hat{\lambda}$ is a feasible solution to the dual problem, then

$$f(\hat{\mathbf{x}}) \geq f(\hat{\mathbf{x}}) + \hat{\lambda}^T \mathbf{g}(\hat{\lambda}) \geq \min_{\mathbf{x} \in \mathcal{D}} f(\mathbf{x} + \hat{\lambda}^T \mathbf{g}(\lambda)) = L^*(\hat{\lambda})$$

This proves the second part of the theorem. A direct consequence of this is that

$$p^* = \min_{\mathbf{x} \in \mathcal{D}} f(\mathbf{x}) \geq \min_{\lambda \geq \mathbf{0}} L^*(\lambda) = d^*$$

□

The weak duality theorem has some important implications. If the primal problem is unbounded below, that is, $p^* = -\infty$, we must have $d^* = -\infty$, which means that the Lagrange dual problem is infeasible. Conversely, if the dual problem is unbounded above, that is, $d^* = \infty$, we must have $p^* = \infty$, which is equivalent to saying that the primal problem is infeasible. The difference, $p^* - d^*$ is called the duality gap.

In many hard combinatorial optimization problems with duality gaps, we get good dual solutions, which tell us that we are guaranteed of being some k % within the optimal solution to the primal, for some satisfactorily low values of k . This is one of the powerful uses of duality theory; constructing bounds for optimization problems.

Under what conditions can one assert that $d^* = p^*$? The condition $d^* = p^*$ is called *strong duality* and it does not hold in general. It usually holds for convex problems but there are exceptions to that - one of the most typical being that of the semi-definite optimization problem. The semi-definite program (SDP) is defined, with the linear matrix inequality constraint (*c.f.* page 260) as follows:

$$\begin{aligned} \min_{\mathbf{x} \in \mathfrak{R}^n} \quad & \mathbf{c}^T \mathbf{x} \\ \text{subject to} \quad & x_1 A_1 + \dots + x_n A_n + G \preceq 0 \\ & \mathbf{A}\mathbf{x} = \mathbf{b} \end{aligned} \tag{4.83}$$

Sufficient conditions for strong duality in convex problems are called *constraint qualifications*. One of the most useful sufficient conditions for strong duality is called the *Slaters constraint qualification*.

Definition 42 [Slaters constraint qualification]: For a convex problem

$$\begin{aligned} \min_{\mathbf{x} \in \mathcal{D}} \quad & f(\mathbf{x}) \\ \text{subject to} \quad & g_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m \\ & \mathbf{Ax} = \mathbf{b} \end{aligned} \quad (4.84)$$

variable $\mathbf{x} = (x_1, \dots, x_n)$

strong duality holds (that is $d^* = p^*$) if it is strictly feasible. That is,

$$\exists \mathbf{x} \in \text{int}(\mathcal{D}) : \quad g_i(\mathbf{x}) < 0 \quad i = 1, 2, \dots, m \quad \mathbf{Ax} = \mathbf{b}$$

However, if any of the g_i 's are linear, they do not need to hold with strict inequalities.

Table 4.4 summarizes some optimization problems, their duals and conditions for strong duality. Strong duality also holds for nonconvex problems

Problem type	Objective Function	Constraints	$L^*(\lambda)$	Dual constraints	Strong duality
Linear Program	$\mathbf{c}^T \mathbf{x}$	$\mathbf{Ax} \leq \mathbf{b}$	$-\mathbf{b}^T \lambda$	$\mathbf{A}^T \lambda + \mathbf{c} = \mathbf{0}$ $\lambda \geq \mathbf{0}$	Feasible primal and dual
Quadratic Program	$\frac{1}{2} \mathbf{x}^T Q \mathbf{x} + \mathbf{c}^T \mathbf{x}$ for $Q \in \mathcal{S}_{++}^n$	$\mathbf{Ax} \leq \mathbf{b}$	$-\frac{1}{2} (\mathbf{c} - \mathbf{A}^T \lambda)^T Q^{-1} (\mathbf{c} - \mathbf{A}^T \lambda) + \mathbf{b}^T \lambda$	$\lambda \geq \mathbf{0}$	Always
Entropy maximization	$x_i \sum_{i=1}^n \ln x_i$	$\mathbf{Ax} \leq \mathbf{b}$ $\mathbf{x}^T \mathbf{1} = 1$	$-\mathbf{b}^T \lambda - \mu - e^{-\mu-1} \sum_{i=1}^n e^{-\mathbf{a}_i^T \lambda}$ \mathbf{a}_i is the i^{th} column of \mathbf{A}	$\lambda \geq \mathbf{0}$	Primal constraints are satisfied.

Table 4.4: Examples of functions and their duals.

in extremely rare cases. One example of this is minimization of a nonconvex quadratic function over the unit ball.

4.4.3 Geometry of the Dual

We will study the geometry of the dual in the *column space* \Re^{m+1} . The column geometry of the dual will require definition of the following set:

$$\mathcal{I} = \{(\mathbf{s}, z) \mid \mathbf{s} \in \Re^m, z \in \Re, \exists \mathbf{x} \in \mathcal{D} \text{ with } g_i(\mathbf{x}) \leq s_i \quad \forall 1 \leq i \leq m, f(\mathbf{x}) \leq z\}$$

The set \mathcal{I} is a subset of \Re^{m+1} , where m is the number of constraints. Consider a plot in two dimensions, for $n = 1$, with s_1 along the x -axis and z along the y -axis. For every point, $\mathbf{x} \in \mathcal{D}$, we can identify all points (s_1, z) for $s_1 \geq g_1(\mathbf{x})$

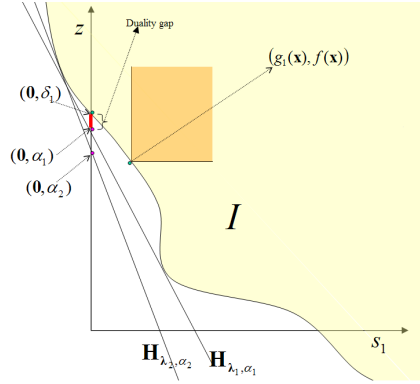


Figure 4.42: Example of the set \mathcal{I} for a single constraint (*i.e.*, for $n = 1$).

and $z \geq f(\mathbf{x})$ and these are points that lie to the left and above the point $(g_1(\mathbf{x}), f(\mathbf{x}))$. An example set \mathcal{I} is shown in Figure 4.42. It turns out that all the intuitions we need are in two dimensions, which makes it fairly convenient to understand the idea. It is straightforward to prove that if the objective function $f(\mathbf{x})$ is convex and each of the constraints $g_i(\mathbf{x})$, $1 \leq i \leq n$ is a convex function, then \mathcal{I} must be a convex set. Since the feasible region for the primal problem (4.78) is the region in \mathcal{I} with $\mathbf{s} \leq \mathbf{0}$, and since all points above and to the right of a point in \mathcal{I} also belong to \mathcal{I} , the solution to the primal problem corresponds to the point in \mathcal{I} with $\mathbf{s} = \mathbf{0}$ and least possible value of z . For example, in Figure 4.42, the solution to the primal corresponds to $(\mathbf{0}, \delta_1)$.

Let us define a hyperplane $\mathcal{H}_{\lambda, \alpha}$, parametrized by $\lambda \in \Re^m$ and $\alpha \in \Re$ as

$$\mathcal{H}_{\lambda, \alpha} = \{(\mathbf{s}, z) \mid \lambda^T \cdot \mathbf{s} + z = \alpha\}$$

Consider all hyperplanes that lie below \mathcal{I} . For example, in the Figure 4.42, both hyperplanes $\mathcal{H}_{\lambda_1, \alpha_1}$ and $\mathcal{H}_{\lambda_2, \alpha_2}$ lie below the set \mathcal{I} . Of all hyperplanes that lie below \mathcal{I} , consider the hyperplane whose intersection with the line $\mathbf{s} = \mathbf{0}$, corresponds to as high a value of z as possible. This hyperplane must be supporting hyperplane. Incidentally, $\mathcal{H}_{\lambda_1, \alpha_1}$ happens to be such a supporting hyperplane. Its point of intersection $(\mathbf{0}, \alpha_1)$ precisely corresponds to the solution to the dual problem. Let us derive this statement formally after setting up some more notation.

We will define two half-spaces corresponding to $\mathcal{H}_{\lambda, \alpha}$

$$\mathcal{H}_{\lambda, \alpha}^+ = \{(\mathbf{s}, z) \mid \lambda^T \cdot \mathbf{s} + z \geq \alpha\}$$

$$\mathcal{H}_{\lambda, \alpha}^- = \{(\mathbf{s}, z) \mid \lambda^T \cdot \mathbf{s} + z \leq \alpha\}$$

Let us define another set \mathcal{L} as

$$\mathcal{L} = \{(\mathbf{s}, z) \mid \mathbf{s} = \mathbf{0}\}$$

Note that \mathcal{L} is essentially the z or function axis. The intersection of $\mathcal{H}_{\lambda,\alpha}$ with \mathcal{L} is the point $(\mathbf{0}, \alpha)$. That is

$$(\mathbf{0}, \alpha) = \mathcal{L} \cap \mathcal{H}_{\lambda,\alpha}$$

We would like to manipulate λ and α so that the set \mathcal{I} lies in the half-space $\mathcal{H}_{\lambda,\alpha}^+$ as tightly as possible. Mathematically, we are interested in the problem of maximizing the height of the point of intersection of \mathcal{L} with $\mathcal{H}_{\lambda,\alpha}$ above the $\mathbf{s} = \mathbf{0}$ plane, while ensuring that \mathcal{I} remains a subset of $\mathcal{H}_{\lambda,\alpha}^+$.

$$\begin{array}{ll} \max & \alpha \\ \text{subject to} & \mathcal{H}_{\lambda,\alpha}^+ \supseteq \mathcal{I} \end{array}$$

By definitions of \mathcal{I} , $\mathcal{H}_{\lambda,\alpha}^+$ and the subset relation, this problem is equivalent to

$$\begin{array}{ll} \max & \alpha \\ \text{subject to} & \lambda^T \cdot \mathbf{s} + z \geq \alpha \quad \forall (\mathbf{s}, z) \in \mathcal{I} \end{array}$$

Now notice that if $(\mathbf{s}, z) \in \mathcal{I}$, then $(\mathbf{s}', z) \in \mathcal{I}$ for all $\mathbf{s}' \geq \mathbf{s}$. This was also illustrated in Figure 4.42. Thus, we cannot afford to have any component of λ negative; if any of the λ_i 's were negative, we could crank up s_i arbitrarily to violate the inequality $\lambda^T \cdot \mathbf{s} + z \geq \alpha$. Thus, we can add the constraint $\lambda \geq \mathbf{0}$ to the above problem without changing the solution.

$$\begin{array}{ll} \max & \alpha \\ \text{subject to} & \lambda^T \cdot \mathbf{s} + z \geq \alpha \quad \forall (\mathbf{s}, z) \in \mathcal{I} \\ & \lambda \geq \mathbf{0} \end{array}$$

Any equality constraint $h(\mathbf{x}) = 0$ can be expressed using two inequality constraints, *viz.*, $h(\mathbf{x}) \leq 0$ and $-h(\mathbf{x}) \leq 0$, implying that its corresponding lagrange multiplier should be exactly 0. This problem can again be proved to be equivalent to the following problem, using the fact that every point on $\partial\mathcal{I}$ must be of the form $(g_1(\mathbf{x}), g_2(\mathbf{x}), \dots, g_m(\mathbf{x}), f(\mathbf{x}))$ for some $\mathbf{x} \in \mathcal{D}$.

$$\begin{array}{ll} \max & \alpha \\ \text{subject to} & \lambda^T \cdot \mathbf{g}(\mathbf{x}) + f(\mathbf{x}) \geq \alpha \quad \forall \mathbf{x} \in \mathcal{D} \\ & \lambda \geq \mathbf{0} \end{array}$$

We will remind the reader at this point that $L(\mathbf{x}, \lambda) = \lambda^T \cdot \mathbf{g}(\mathbf{x}) + f(\mathbf{x})$. The above problem is therefore the same as

$$\begin{array}{ll} \max & \alpha \\ \text{subject to} & L(\mathbf{x}, \lambda) \geq \alpha \quad \forall \mathbf{x} \in \mathcal{D} \\ & \lambda \geq \mathbf{0} \end{array}$$

Since, $L^*(\lambda) = \min_{\mathbf{x} \in \mathcal{D}} L(\mathbf{x}, \lambda)$, we can deal with the equivalent problem

$$\begin{array}{ll} \max & \alpha \\ \text{subject to} & L^*(\lambda) \geq \alpha \\ & \lambda \geq \mathbf{0} \end{array}$$

This problem can be restated as

$$\begin{array}{ll} \max & L^*(\lambda) \\ \text{subject to} & \lambda \geq \mathbf{0} \end{array}$$

This is precisely the dual problem. We thus get a geometric interpretation of the dual.

Again referring to Figure 4.42, we note that if the set \mathcal{I} is not convex, there could be a *gap* between the z -intercept $(\mathbf{0}, \alpha_1)$ of the best supporting hyperplane $\mathcal{H}_{\lambda_1, \alpha_1}$ and the closest point $(\mathbf{0}, \delta_1)$ of \mathcal{I} on the z -axis, which corresponds to the solution to the primal. In fact, when the set \mathcal{I} is not convex, we can never prove that there will be no duality gap. And even when the set \mathcal{I} is convex, bizarre things can happen; for example, in the case of semi-definite programming, the set \mathcal{I} , though convex, is not at all well-behaved and this yields a large duality gap, as shown in Figure 4.43. In fact, the set \mathcal{I} is open from below (the dotted boundary) for a semi-definite program. We could create very simple problems with convex \mathcal{I} , for which there are duality gaps. For well-behaved convex functions (as in the case of linear programming), there are no duality gaps. Figure 4.44 illustrates the case of a well-behaved convex program.

4.4.4 Complementary slackness and KKT Conditions

We now state the conditions between the primal and dual optimal points for an arbitrary function. These conditions, called the *Karush-Kuhn-Tucker conditions* (abbreviated as KKT conditions) state a necessary condition for a solution to be optimal with zero duality gap. Consider the following general optimization problem.

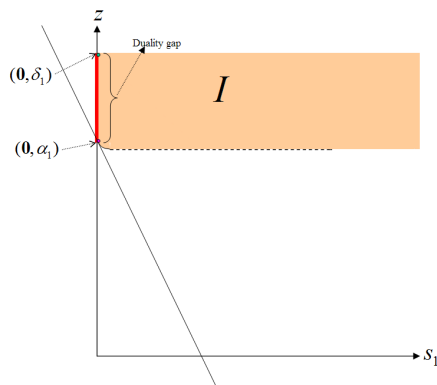


Figure 4.43: Example of the convex set \mathcal{I} for a single constrained semi-definite program.

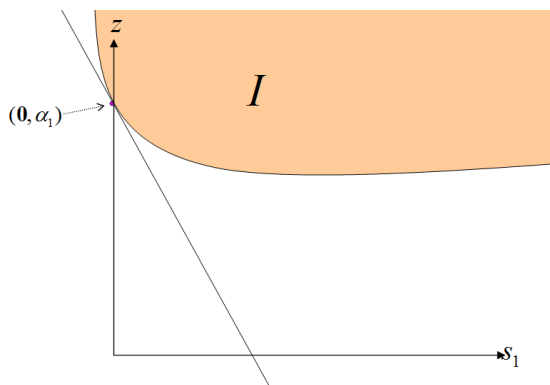


Figure 4.44: Example of the convex set \mathcal{I} for a single constrained well-behaved convex program.

$$\begin{aligned}
& \min_{\mathbf{x} \in \mathcal{D}} && f(\mathbf{x}) \\
& \text{subject to} && g_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m \\
& && h_j(\mathbf{x}) = 0, \quad j = 1, \dots, p \\
& \text{variable } \mathbf{x} = && (x_1, \dots, x_n)
\end{aligned} \tag{4.85}$$

Suppose that the primal and dual optimal values for the above problem are attained and equal, that is, strong duality holds. Let $\hat{\mathbf{x}}$ be a primal optimal and $(\hat{\lambda}, \hat{\mu})$ be a dual optimal point ($\hat{\lambda} \in \Re^m, \hat{\mu} \in \Re^p$). Thus,

$$\begin{aligned}
f(\hat{\mathbf{x}}) &= L^*(\hat{\lambda}, \hat{\mu}) \\
&= \min_{\mathbf{x} \in \mathcal{D}} f(\mathbf{x}) + \hat{\lambda}^T \mathbf{g}(\mathbf{x}) + \hat{\mu}^T \mathbf{h}(\mathbf{x}) \\
&\leq f(\hat{\mathbf{x}}) + \hat{\lambda}^T \mathbf{g}(\hat{\mathbf{x}}) + \hat{\mu}^T \mathbf{h}(\hat{\mathbf{x}}) \\
&\leq f(\hat{\mathbf{x}})
\end{aligned}$$

The last inequality follows from the fact that $\hat{\lambda} \geq \mathbf{0}$, $\mathbf{g}(\hat{\mathbf{x}}) \leq \mathbf{0}$, and $\mathbf{h}(\hat{\mathbf{x}}) = \mathbf{0}$. We can therefore conclude that the two inequalities in this chain must hold with equality. Some of the conclusions that we can draw from this chain of equalities are

1. That $\hat{\mathbf{x}}$ is a minimizer for $L(\mathbf{x}, \hat{\lambda}, \hat{\mu})$ over $\mathbf{x} \in \mathcal{D}$. In particular, if the functions f, g_1, g_2, \dots, g_m and h_1, h_2, \dots, h_p are differentiable (and therefore have open domains), the gradient of $L(\mathbf{x}, \hat{\lambda}, \hat{\mu})$ must vanish at $\hat{\mathbf{x}}$, since any point of global optimum must be a point of local optimum. That is,

$$\nabla f(\hat{\mathbf{x}}) + \sum_{i=1}^m \hat{\lambda}_i \nabla g_i(\hat{\mathbf{x}}) + \sum_{j=1}^p \hat{\mu}_j \nabla h_j(\hat{\mathbf{x}}) = \mathbf{0} \tag{4.86}$$

2. That

$$\hat{\lambda}^T \mathbf{g}(\hat{\mathbf{x}}) = \sum_{i=1}^m \hat{\lambda}_i g_i(\hat{\mathbf{x}}) = 0$$

Since each term in this sum is nonpositive, we conclude that

$$\hat{\lambda}_i g_i(\hat{\mathbf{x}}) = 0 \text{ for } i = 1, 2, \dots, m \tag{4.87}$$

This condition is called *complementary slackness* and is a necessary condition for strong duality. Complementary slackness implies that the i^{th}

optimal lagrange multiplier is 0 unless the i^{th} inequality constraint is active at the optimum. That is,

$$\begin{aligned}\widehat{\lambda}_i > 0 &\Rightarrow g_i(\widehat{\mathbf{x}}) = 0 \\ g_i(\widehat{\mathbf{x}}) < 0 &\Rightarrow \widehat{\lambda}_i = 0\end{aligned}$$

Let us further assume that the functions f, g_1, g_2, \dots, g_m and h_1, h_2, \dots, h_p are differentiable on open domains. As above, let $\widehat{\mathbf{x}}$ be a primal optimal and $(\widehat{\lambda}, \widehat{\mu})$ be a dual optimal point with zero duality gap. Putting together the conditions in (4.86), (4.87) along with the feasibility conditions for any primal solution and dual solution, we can state the following Karush-Kuhn-Tucker (KKT) necessary conditions for zero duality gap.

$$\begin{aligned}(1) \quad \nabla f(\widehat{\mathbf{x}}) + \sum_{i=1}^m \widehat{\lambda}_i \nabla g_i(\widehat{\mathbf{x}}) + \sum_{j=1}^p \widehat{\mu}_j \nabla h_j(\widehat{\mathbf{x}}) &= \mathbf{0} \\ (2) \quad g_i(\widehat{\mathbf{x}}) &\leq 0 \quad i = 1, 2, \dots, m \\ (3) \quad \widehat{\lambda}_i &\geq 0 \quad i = 1, 2, \dots, m \quad (4.88) \\ (4) \quad \widehat{\lambda}_i g_i(\widehat{\mathbf{x}}) &= 0 \quad i = 1, 2, \dots, m \\ (5) \quad h_j(\widehat{\mathbf{x}}) &= 0 \quad j = 1, 2, \dots, p\end{aligned}$$

When the primal problem is convex, the KKT conditions are also sufficient for the points to be primal and dual optimal with zero duality gap. If f is convex, g_i are convex and h_j are affine, the primal problem is convex and consequently, the KKT conditions are sufficient conditions for zero duality gap.

Theorem 82 *If the function f is convex, g_i are convex and h_j are affine, then KKT conditions in 4.88 are necessary and sufficient conditions for zero duality gap.*

Proof: The necessity part has already been proved; here we only prove the sufficiency part. The conditions (2) and (5) in (4.88) ensure that $\widehat{\mathbf{x}}$ is primal feasible. Since $\lambda \geq \mathbf{0}$, $L(\mathbf{x}, \widehat{\lambda}, \widehat{\mu})$ is convex in \mathbf{x} . Based on condition (1) in (4.88) and theorem 77, we can infer that $\widehat{\mathbf{x}}$ minimizes $L(\mathbf{x}, \widehat{\lambda}, \widehat{\mu})$. We can thus conclude that

$$\begin{aligned}L^*(\widehat{\lambda}, \widehat{\mu}) &= f(\widehat{\mathbf{x}}) + \widehat{\lambda}^T \mathbf{g}(\widehat{\mathbf{x}}) + \widehat{\mu}^T \mathbf{h}(\widehat{\mathbf{x}}) \\ &= f(\widehat{\mathbf{x}})\end{aligned}$$

In the equality above, we use $h_j(\widehat{\mathbf{x}}) = 0$ and $\widehat{\lambda}_i g_i(\widehat{\mathbf{x}}) = 0$. Further,

$$d^* \geq L^*(\widehat{\lambda}, \widehat{\mu}) = f(\widehat{\mathbf{x}}) \geq p^*$$

The duality theorem (theorem 81) however states that $p^* \geq d^*$. This implies that

$$d^* = L^*(\widehat{\lambda}, \widehat{\mu}) = f(\widehat{\mathbf{x}}) = p^*$$

This shows that $\widehat{\mathbf{x}}$ and $(\widehat{\lambda}, \widehat{\mu})$ correspond to the primal and dual optimals respectively and the problem therefore has zero duality gap. \square

In summary, for any convex optimization problem with differentiable objective and constraint functions, any points that satisfy the KKT conditions are primal and dual optimal, and have zero duality gap.

The KKT conditions play a very important role in optimization. In some rare cases, it is possible to solve the optimization problems by finding a solution to the KKT conditions analytically. Many algorithms for convex optimization are conceived as, or can be interpreted as, methods for solving the KKT conditions.

4.5 Algorithms for Unconstrained Minimization

We will now study some algorithms for solving convex problems. These techniques are relevant for most convex optimization problems that do not yield themselves to closed form solutions. We will start with unconstrained minimization.

Recall that the goal in unconstrained minimization is to solve the convex problem

$$\min_{\mathbf{x} \in \mathcal{D}} f(\mathbf{x})$$

We are not interested in f , whose solution can be obtained in closed form. For example, minimizing a quadratic is very simple and can be solved by linear equations, an example of which was discussed in Section 3.9.2. Let us denote the optimal solution of the minimization problem by p^* . We will assume that f is convex and twice continuously differentiable and that it attains a finite optimal value p^* . Most unconstrained minimization techniques produce a sequence of points $\mathbf{x}^{(k)} \in \mathcal{D}$, $k = 0, 1, \dots$ such that $f(\mathbf{x}^{(k)}) \rightarrow p^*$ as $k \rightarrow \infty$ or, $\nabla f(\mathbf{x}^{(k)}) \rightarrow \mathbf{0}$ as $k \rightarrow \infty$. Iterative techniques for optimization, further require a starting point $\mathbf{x}^{(0)} \in \mathcal{D}$ and sometimes that $\text{epi}(f)$ is closed. The $\text{epi}(f)$ can be inferred to be closed either if $\mathcal{D} = \mathfrak{R}^n$ or $f(\mathbf{x}) \rightarrow \infty$ as $\mathbf{x} \rightarrow \partial\mathcal{D}$. The function $f(x) = \frac{1}{x}$ for $x > 0$ is an example of a function whose $\text{epi}(f)$ is closed.

While there exist convergence proofs (including guarantees on number of optimization iterations) for many convex optimization algorithms, the proofs assume many conditions, many of which are either not verifiable or involve unknown constants (such as the Lipschitz constant). Thus, most convergence proofs for convex optimization problems are useless in practice, though it is good to know that there are conditions under which the algorithm converges. Since convergence proofs are only of theoretical importance, we will make the strongest possible assumption under which convergence can be proved easily, which is that the function f is strongly convex (*c.f.* Section 4.2.7 for definition of strong convexity) with the strong convexity constant $c > 0$ for which $\nabla^2 f(\mathbf{x}) \succeq cI \forall \mathbf{x} \in \mathcal{D}$.

Further, it can be proved that for a strongly convex function f , $\nabla^2 f(\mathbf{x}) \preceq DI$ for some constant $D \in \Re$. The ratio $\frac{D}{c}$ is an upper bound on the condition number of the matrix $\nabla^2 f(\mathbf{x})$.

4.5.1 Descent Methods

Descent methods for unconstrained optimization have been in use since the last 70 years or more. The general idea in descent methods is that the next iterate $\mathbf{x}^{(k+1)}$ is the current iterate $\mathbf{x}^{(k)}$ added with a descent or search direction $\Delta \mathbf{x}^{(k)}$ (a unit vector), which is multiplied by a scale factor $t^{(k)}$, called the step length.

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + t^{(k)} \Delta \mathbf{x}^{(k)}$$

The incremental step is determined while ensuring that $f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)})$. We assume that we are dealing with the extended value extension of the convex function f (c.f. definition 36), which returns ∞ for any point outside its domain. However, if we do so, we need to make sure that the initial point indeed lies in the domain \mathcal{D} .

A single iteration of the general descent algorithm (shown in Figure 4.45) consists of two main steps, *viz.*, determining a good descent direction $\Delta \mathbf{x}^{(k)}$, which is typically forced to have unit norm and determining the step size using some line search technique. If the function f is convex, and we require that $f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)})$ then, we must have $\nabla^T f(\mathbf{x}^{(k+1)})(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) < 0$. This can be seen from the necessary and sufficient condition for convexity stated in equation (4.44) within Section 4.2.9 and restated here for reference.

$$f(\mathbf{x}^{(k+1)}) \geq f(\mathbf{x}^{(k)}) + \nabla^T f(\mathbf{x}^{(k)})(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})$$

Since $t^{(k)} > 0$, we must have

$$\nabla^T f(\mathbf{x}^{(k)}) \Delta \mathbf{x}^{(k)} < 0$$

That is, the descent direction $\Delta \mathbf{x}^{(k)}$ must make an obtuse angle ($\theta \in (\frac{\pi}{2}, \frac{3\pi}{2})$) with the gradient vector.

Find a starting point $\mathbf{x}^{(0)} \in \mathcal{D}$
repeat
 1. Determine $\Delta \mathbf{x}^{(k)}$.
 2. Choose a step size $t^{(k)} > 0$ using ray^a search.
 3. Obtain $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + t^{(k)} \Delta \mathbf{x}^{(k)}$.
 4. Set $k = k + 1$.
until stopping criterion (such as $\|\nabla f(\mathbf{x}^{(k+1)})\| < \epsilon$) is satisfied

^aMany textbooks refer to this as line search, but we prefer to call it ray search, since the step must be positive.

Figure 4.45: The general descent algorithm.

There are many different empirical techniques for ray search, though it matters much less than the search for the descent direction. These techniques reduce the n -dimensional problem to a 1-dimensional problem, which can be easy to solve by use of plotting and eyeballing or even exact search.

1. **Exact ray search:** The exact ray search seeks a scaling factor t that satisfies

$$t = \underset{t>0}{\operatorname{argmin}} f(\mathbf{x} + t\Delta\mathbf{x}) \quad (4.89)$$

2. **Backtracking ray search:** The exact line search may not be feasible or could be expensive to compute for complex non-linear functions. A relatively simpler ray search iterates over values of step size starting from 1 and scaling it down by a factor of $\beta \in (0, \frac{1}{2})$ after every iteration till the following condition, called the *Armijo condition* is satisfied for some $0 < c_1 < 1$.

$$f(\mathbf{x} + t\Delta\mathbf{x}) < f(\mathbf{x}) + c_1 t \nabla^T f(\mathbf{x}) \Delta\mathbf{x} \quad (4.90)$$

Based on equation (4.44), it can be inferred that the Armijo inequality can never hold for $c_1 = 1$; for $c_1 = 1$, the right hand side of the Armijo condition gives a lower bound on the value of $f(\mathbf{x} + t\Delta\mathbf{x})$. The Armijo condition simply ensures that t decreases f sufficiently. Often, another condition is used for inexact line search in conjunction with the Armijo condition.

$$|\Delta\mathbf{x}^T \nabla f(\mathbf{x} + t\Delta\mathbf{x})| \leq c_2 |\Delta\mathbf{x}^T \nabla f(\mathbf{x})| \quad (4.91)$$

where $1 > c_1 > c_2 > 0$. This condition ensures that the slope of the function $f(\mathbf{x} + t\Delta\mathbf{x})$ at t is less than c_2 times that at $t = 0$. The conditions in (4.90) and (4.91) are together called the strong Wolfe conditions. These conditions are particularly very important for non-convex problems.

A finding that is borne out of plenty of empirical evidence is that exact ray search does better than empirical ray search in a few cases only. Further, the exact choice of the value of β and α seems to have little effect on the convergence of the overall descent method.

The trend of specific descent methods has been like a parabola - starting with simple steepest descent techniques, then accomodating the curvature hessian matrix through a more sophisticated Newton's method and finally, trying to simplify the Newton's method through approximations to the hessian inverse,

culminating in conjugate gradient techniques, that do away with any curvature matrix whatsoever, and form the internal combustion engine of many sophisticated optimization techniques today. We start the thread by describing the steepest descent methods.

Steepest Descent

Let $\mathbf{v} \in \mathfrak{R}^n$ be a unit vector under some norm. By theorem 75, for convex f ,

$$f(\mathbf{x}^{(k)}) - f(\mathbf{x}^{(k)} + \mathbf{v}) \leq -\nabla^T f(\mathbf{x}^{(k)})\mathbf{v}$$

For small \mathbf{v} , the inequality turns into approximate equality. The term $-\nabla^T f(\mathbf{x}^{(k)})\mathbf{v}$ can be thought of as (an upper-bound on) the first order prediction of decrease. The idea in the steepest descent method [?] is to choose a norm and then determine a descent direction such that for a unit step in that norm, the first order prediction of decrease is maximized. This choice of the descent direction can be stated as

$$\Delta \mathbf{x} = \operatorname{argmin} \{ \nabla^T f(\mathbf{x})\mathbf{v} \mid \|\mathbf{v}\| = 1 \}$$

The algorithm is outlined in Figure 4.46.

Find a starting point $\mathbf{x}^{(0)} \in \mathcal{D}$.
repeat
 1. Set $\Delta \mathbf{x}^{(k)} = \operatorname{argmin} \{ \nabla^T f(\mathbf{x}^{(k)})\mathbf{v} \mid \|\mathbf{v}\| = 1 \}$.
 2. Choose a step size $t^{(k)} > 0$ using exact or backtracking ray search.
 3. Obtain $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + t^{(k)} \Delta \mathbf{x}^{(k)}$.
 4. Set $k = k + 1$.
until stopping criterion (such as $\|\nabla f(\mathbf{x}^{(k+1)})\| \leq \epsilon$) is satisfied

Figure 4.46: The steepest descent algorithm.

The key to understanding the steepest descent method (and in fact many other iterative methods) is that it heavily depends on the choice of the norm. It has been empirically observed that if the norm chosen is aligned with the gross geometry of the sub-level sets²², the steepest descent method converges faster to the optimal solution. If the norm chosen is not aligned, it often amplifies the effect of oscillations. Two examples of the steepest descent method are the gradient descent method (for the euclidian or L_2 norm) and the coordinate-descent method (for the L_1 norm). One fact however is that no two norms should give exactly opposite steepest descent directions, though they may point in different directions.

Gradient Descent

A classic greedy algorithm for minimization is the gradient descent algorithm. This algorithm uses the negative of the gradient of the function at the current

²²The alignment can be determined by fitting, for instance, a quadratic to a sample of the points.

point \mathbf{x}^* as the descent direction $\Delta\mathbf{x}^*$. It turns out that this choice of $\Delta\mathbf{x}^*$ corresponds to the direction of steepest descent under the L_2 (eucledian) norm. This can be proved in a straightforward manner using theorem 58. The algorithm is outlined in Figure 4.47. The steepest descent method can be thought

Find a starting point $\mathbf{x}^{(0)} \in \mathcal{D}$
repeat
 1. Set $\Delta\mathbf{x}^{(k)} = -\nabla f(\mathbf{x}^{(k)})$.
 2. Choose a step size $t^{(k)} > 0$ using exact or backtracking ray search.
 3. Obtain $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + t^{(k)}\Delta\mathbf{x}^{(k)}$.
 4. Set $k = k + 1$.
until stopping criterion (such as $\|\nabla f(\mathbf{x}^{(k+1)})\|_2 \leq \epsilon$) is satisfied

Figure 4.47: The gradient descent algorithm.

of as changing the coordinate system in a particular way and then applying the gradient descent method in the changed coordinate system.

Coordinate-Descent Method

The co-ordinate descent method corresponds exactly to the choice of L_1 norm for the steepest descent method. The steepest descent direction using the L_1 norm is given by

$$\Delta\mathbf{x} = -\frac{\partial f(\mathbf{x})}{\partial x_i} \mathbf{u}^i$$

where,

$$\frac{\partial f(\mathbf{x})}{\partial x_i} = \|\nabla f(\mathbf{x})\|_\infty$$

and \mathbf{u}^i was defined on page 229 as the unit vector pointing along the i^{th} coordinate axis. Thus each iteration of the coordinate descent method involves optimizing over one component of the vector $\mathbf{x}^{(k)}$ and then updating the vector. The component chosen is the one having the largest absolute value in the gradient vector. The algorithm is outlined in Figure 4.48.

Convergence of Steepest Descent Method

For the gradient method, it can be proved that if f is strongly convex,

$$f(\mathbf{x}^{(k)}) - p^* \leq \rho^k \left(f(\mathbf{x}^{(0)}) - p^* \right) \quad (4.92)$$

The value of $\rho \in (0, 1)$ depends on the strong convexity constant c (*c.f.* equation (4.64) on page 275), the value of $\mathbf{x}^{(0)}$ and type of ray search employed. The suboptimality $f(\mathbf{x}^{(k)}) - p^*$ goes down by a factor $\rho < 1$ at every step and this is referred to as *linear convergence*²³. However, this is only of theoretical

²³A series s_1, s_2, \dots is said to have

Find a starting point $\mathbf{x}^{(0)} \in \mathcal{D}$.
Select an appropriate norm $\|\cdot\|$.
repeat
 1. Let $\frac{\partial f(\mathbf{x}^{(k)})}{\partial x_i^{(k)}} = \|\nabla f(\mathbf{x})\|_\infty$.
 2. Set $\Delta \mathbf{x}^{(k)} = -\frac{\partial f(\mathbf{x}^{(k)})}{\partial x_i^{(k)}} \mathbf{u}^i$.
 3. Choose a step size $t^{(k)} > 0$ using exact or backtracking ray search.
 4. Obtain $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + t^{(k)} \Delta \mathbf{x}^{(k)}$.
 5. Set $k = k + 1$.
until stopping criterion (such as $\|\nabla f(\mathbf{x}^{(k+1)})\|_\infty \leq \epsilon$) is satisfied

Figure 4.48: The coordinate descent algorithm.

importance, since this method is often very slow, indicated by values of ρ , very close to 1. Use of exact line search in conjunction with gradient descent also has the tendency to overshoot the next best iterate. It is therefore rarely used in practice. The convergence rate depends greatly on the condition number of the Hessian (which is upperbounded by $\frac{D}{c}$). It can be proved that the number of iterations required for the convergence of the gradient descent method is lower-bounded by the condition number of the hessian; large eigenvalues correspond to high curvature directions and small eigenvalues correspond to low curvature directions. Many methods (such as conjugate gradient) try to improve upon the gradient method by making the hessian better conditioned. Convergence can be very slow even for moderately well-conditioned problems, with condition number in the 100s, even though computation of the gradient at each step is only an $O(n)$ operation. The gradient descent method however works very well if the function is isotropic, that is if the level-curves are spherical or nearly spherical.

The convergence of the steepest descent method can be stated in the same form as in 4.92, using the fact that any norm can be bounded in terms of the Euclidean norm, *i.e.*, there exists a constant $\eta \in (0, 1]$ such that

$$\|\mathbf{x}\| \geq \eta \|\mathbf{x}\|_2$$

-
1. linear convergence to \bar{s} if $\lim_{i \rightarrow \infty} \frac{|s_{i+1} - \bar{s}|}{|s_i - \bar{s}|} = \delta \in (0, 1)$. For example, $s_i = (\gamma)^i$ has linear convergence to $\bar{s} = 0$ for any $\gamma < 1$. The rate of decrease is also sometimes called exponential or geometric. This is considered quite slow.
 2. superlinear convergence to \bar{s} if $\lim_{i \rightarrow \infty} \frac{|s_{i+1} - \bar{s}|}{|s_i - \bar{s}|} = 0$. For example, $s_i = \frac{1}{i!}$ has superlinear convergence. This is the most common.
 3. quadratic convergence to \bar{s} if $\lim_{i \rightarrow \infty} \frac{|s_{i+1} - \bar{s}|}{|s_i - \bar{s}|^2} = \delta \in (0, \infty)$. For example, $s_i = (\gamma)^{2^i}$ has quadratic convergence to $\bar{s} = 0$ for any $\gamma < 1$. This is considered very fast in practice.

4.5.2 Newton's Method

Newton's method [?] is based on approximating a function around the current iterate $\mathbf{x}^{(k)}$ using a second degree Taylor expansion.

$$f(\mathbf{x}) \approx \tilde{f}(\mathbf{x}) = f(\mathbf{x}^{(k)}) + \nabla^T f(\mathbf{x}^{(k)})(\mathbf{x} - \mathbf{x}^{(k)}) + \frac{1}{2}(\mathbf{x} - \mathbf{x}^{(k)})^T \nabla^2 f(\mathbf{x}^{(k)})(\mathbf{x} - \mathbf{x}^{(k)})$$

If the function f is convex, the quadratic approximation is also convex. Newton's method is based on solving it exactly by finding its critical point $\mathbf{x}^{(k+1)}$ as a function of $\mathbf{x}^{(k)}$. Setting the gradient of this quadratic approximation (with respect to \mathbf{x}) to $\mathbf{0}$ gives

$$\nabla^T f(\mathbf{x}^{(k)}) + \nabla^2 f(\mathbf{x}^{(k)})(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) = 0$$

solving which yields the next iterate as

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \left(\nabla^2 f(\mathbf{x}^{(k)}) \right)^{-1} \nabla f(\mathbf{x}^{(k)}) \quad (4.93)$$

assuming that the Hessian matrix is invertible. The term $\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}$ can be thought of as an update step. This leads to a simple descent algorithm, outlined in Figure 4.49 and is called the Newton's method. It relies on the invertibility of the hessian, which holds if the hessian is positive definite as in the case of a strictly convex function. In case the hessian is invertible, cholesky factorization (page 204) of the hessian can be used to solve the linear system (4.93). However, the Newton method may not even be properly defined if the hessian is not positive definite. In this case, the hessian could be changed to a nearby positive definite matrix whenever it is not. Or a line search could be added to seek a new point having a positive definite hessian.

This method uses a step size of 1. If instead, the stepsize is chosen using exact or backtracking ray search, the method is called the *damped Newton's method*. Each Newton's step takes $O(n^3)$ time (without using any fast matrix multiplication methods).

The Newton step can also be looked upon as another incarnation of the steepest descent rule, but with the quadratic norm defined by the (local) Hessian $\nabla^2 f(\mathbf{x}^{(k)})$ evaluated at the current iterate $\mathbf{x}^{(k)}$, *i.e.*,

$$\|\mathbf{u}\|_{\nabla^2 f(\mathbf{x}^{(k)})} = \left(\mathbf{u} \nabla^2 f(\mathbf{x}^{(k)}) \mathbf{u} \right)^{\frac{1}{2}}$$

The norm of the Newton step, in the quadratic norm defined by the Hessian at a point \mathbf{x} is called the *Newton decrement* at the point \mathbf{x} and is denoted by $\lambda(\mathbf{x})$. Thus,

$$\lambda(\mathbf{x}) = \|\Delta \mathbf{x}\|_{\nabla^2 f(\mathbf{x})} = \nabla^T f(\mathbf{x}) \left(\nabla^2 f(\mathbf{x}) \right)^{-1} \nabla f(\mathbf{x})$$

The Newton decrement gives an 'estimate' of the proximity of the current iterate \mathbf{x} to the optimal point \mathbf{x}^* obtained by measuring the proximity of \mathbf{x} to the

Find a starting point $\mathbf{x}^{(0)} \in \mathcal{D}$.
Select an appropriate tolerance $\epsilon > 0$.
repeat
 1. Set $\Delta \mathbf{x}^{(k)} = -(\nabla^2 f(\mathbf{x}^{(k)}))^{-1} \nabla f(\mathbf{x}^{(k)})$.
 2. Let $\lambda^2 = \nabla^T f(\mathbf{x}^{(k)}) (\nabla^2 f(\mathbf{x}^{(k)}))^{-1} \nabla f(\mathbf{x}^{(k)})$.
 3. If $\frac{\lambda^2}{2} \leq \epsilon$, **quit**.
 4. Set step size $t^{(k)} = 1$.
 5. Obtain $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + t^{(k)} \Delta \mathbf{x}^{(k)}$.
 6. Set $k = k + 1$.
until

Figure 4.49: The Newton's method.

minimum point of the quadratic approximation $\tilde{f}(\mathbf{x})$. The estimate is $\frac{1}{2}\lambda(\mathbf{x})^2$ and is given as

$$\frac{1}{2}\lambda(\mathbf{x})^2 = f(\mathbf{x}) - \min \tilde{f}(\mathbf{x})$$

Additionally, $\lambda(\mathbf{x})^2$ is also the directional derivative in the Newton direction.

$$\lambda(\mathbf{x})^2 = \nabla^T f(\mathbf{x}) \Delta \mathbf{x}$$

The estimate $\frac{1}{2}\lambda(\mathbf{x})^2$ is used to test the convergence of the Newton algorithm in Figure 4.49.

Next, we state an important property of the Newton's update rule.

Theorem 83 *If $\Delta \mathbf{x}^{(k)} = -(\nabla^2 f(\mathbf{x}^{(k)}))^{-1} \nabla f(\mathbf{x}^{(k)})$, $\nabla^2 f(\mathbf{x}^{(k)})$ is symmetric and positive definite and $\Delta \mathbf{x}^{(k)} \neq \mathbf{0}$, then $\Delta \mathbf{x}^{(k)}$ is a descent direction at $\mathbf{x}^{(k)}$, that is, $\nabla^T f(\mathbf{x}^{(k)}) \Delta \mathbf{x}^{(k)} < 0$.*

Proof: First of all, if $\nabla^2 f(\mathbf{x}^{(k)})$ is symmetric and positive definite, then it is invertible and its inverse is also symmetric and positive definite. Next, we see that

$$\nabla^T f(\mathbf{x}^{(k)}) \Delta \mathbf{x}^{(k)} = -\nabla^T f(\mathbf{x}^{(k)}) \left(\nabla^2 f(\mathbf{x}^{(k)}) \right)^{-1} \nabla f(\mathbf{x}^{(k)}) < 0$$

because $(\nabla^2 f(\mathbf{x}^{(k)}))^{-1}$ is symmetric and positive definite. \square

The Newton method is independent of affine changes of coordinates. That is, if optimizing a function $f(\mathbf{x})$ using the Newton's method with an initial estimate $\mathbf{x}^{(0)}$ involves the series of iterates $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(k)}, \dots$, then optimizing the same problem using the Newton's method with a change of coordinates given by $\mathbf{x} = A\mathbf{y}$ and the initial estimate $\mathbf{y}^{(0)}$ such that $\mathbf{x}^{(0)} = A\mathbf{y}^{(0)}$ yields the series of iterates $\mathbf{y}^{(1)}, \mathbf{y}^{(2)}, \dots, \mathbf{y}^{(k)}, \dots$, such that $\mathbf{x}^{(k)} = A\mathbf{y}^{(k)}$. This is a great advantage over the gradient method, whose convergence can be very sensitive to affine transformation.

Another well known feature of the Newton's method is that it converges very fast, if at all. The convergence is extremely fast in the vicinity of the point of

optimum. This can be loosely understood as follows. If \mathbf{x}^* is the critical point of a differentiable convex function f , defined on an open domain, the function is approximately equal to its second order Taylor approximation in the vicinity of \mathbf{x}^* . Further, $\nabla f(\mathbf{x}^*) = \mathbf{0}$. This gives the following approximation at any point \mathbf{x} in the vicinity of \mathbf{x}^* .

$$\begin{aligned} f(\mathbf{x}) &\approx f(\mathbf{x}^*) + \nabla^T f(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*) + \frac{1}{2}(\mathbf{x} - \mathbf{x}^*)^T \nabla^2 f(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*) \\ &= f(\mathbf{x}^*) + \frac{1}{2}(\mathbf{x} - \mathbf{x}^*)^T \nabla^2 f(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*) \end{aligned}$$

Thus, the level curves of a convex function are approximately ellipsoids near the point of minimum \mathbf{x}^* . Given this geometry near the minimum, it then makes sense to do steepest descent in the norm induced by the hessian, near the point of minimum (which is equivalent to doing a steepest descent after a rotation of the coordinate system using the hessian). This is exactly the Newton's step. Thus, the Newton's method²⁴ converges very fast in the vicinity of the solution.

This convergence analysis is formally stated in the following theorem and is due to Leonid Kantorovich.

Theorem 84 *Suppose $f(\mathbf{x}) : \mathcal{D} \rightarrow \Re$ is twice continuously differentiable on \mathcal{D} and \mathbf{x}^* is the point corresponding to the optimal value p^* (so that $\nabla f(\mathbf{x}^*) = \mathbf{0}$). Let f be strongly convex on \mathcal{D} with constant $c > 0$. Also, suppose $\nabla^2 f(\mathbf{x}^*)$ is Lipschitz continuous on \mathcal{D} with a constant $L > 0$ (which measures how well f can be approximated by a quadratic function or how fast the second derivative of f changes), that is*

$$\|\nabla^2 f(\mathbf{x}) - \nabla^2 f(\mathbf{y})\|_2 \leq L\|\mathbf{x} - \mathbf{y}\|_2$$

Then, there exist constants $\alpha \in (0, \frac{c^2}{L})$ and $\beta > 0$ such that

1. **Damped Newton Phase:** *If $\|\nabla^2 f(\mathbf{x})\|_2 \geq \alpha$, then $f(\mathbf{x}^{(k+1)}) - f(\mathbf{x}^{(k)}) \leq -\beta$. That is, at every step of the iteration in the damped Newton phase, the function value decreases by at least β and the phase ends after at most $\frac{f(\mathbf{x}^{(0)}) - p^*}{\beta}$ iterations, which is a finite number.*
2. **Quadratically Convergent Phase:** *If $\|\nabla^2 f(\mathbf{x})\|_2 < \alpha$, then $\frac{L}{2c^2} \|\nabla f(\mathbf{x}^{(k+1)})\|_2 \leq (\frac{L}{2c^2} \|\nabla f(\mathbf{x}^{(k)})\|_2)^2$. When applied recursively this inequality yields*

$$\frac{L}{2c^2} \|\nabla f(\mathbf{x}^{(k)})\|_2 \leq \left(\frac{1}{2}\right)^{2^{k-q}}$$

where q is iteration number, starting at which $\|\nabla^2 f(\mathbf{x}^{(q)})\|_2 < \alpha$. Using the result for strong convexity in equation (4.50) on page 271, we can derive

²⁴Newton originally presented his method for one-dimensional problems. Later on Raphson extended the method to multi-dimensional problems.

$$f(\mathbf{x}^{(k)}) - p^* \leq \frac{1}{2c} \|\nabla f(\mathbf{x}^{(k)})\|_2^2 \leq \frac{2c^3}{L^2} \left(\frac{1}{2}\right)^{2^{k-q+1}} \quad (4.94)$$

Also, using the result in equation (4.52) on page 271, we get a bound on the distance between the current iterate and the point \mathbf{x}^* corresponding to the optimum.

$$\|\mathbf{x}^{(k)} - \hat{\mathbf{x}}^*\|_2 \leq \frac{2}{c} \|\nabla f(\mathbf{x}^{(k)})\|_2 \leq \frac{c}{L} \left(\frac{1}{2}\right)^{2^{k-q}} \quad (4.95)$$

Inequality (4.94) shows that convergence is quadratic once the second condition is satisfied after a finite number of iterations. Roughly speaking, this means that, after a sufficiently large number of iterations, the number of correct digits doubles at each iteration²⁵. In practice, once in the quadratic phase, you do not even need to bother about any convergence criterion; it suffices to apply a fixed few number of Newton iterations to get a very accurate solution. Inequality (4.95) states that the sequence of iterates converges quadratically. The Lipschitz continuity condition states that if the second derivative of the function changes relatively slowly, applying Newton's method can be useful. Again, the inequalities are technical junk as far as practical application of Newton's method is concerned, since L , c and α are generally unknown, but it helps to understand the properties of the Newton's method, such as its two phases and identify them in problems. In practice, Newton's method converges very rapidly, if at all.

As an example, consider a one dimensional function $f(x) = 7x - \ln x$. Then $f'(x) = 7 - \frac{1}{x}$ and $f''(x) = \frac{1}{x^2}$. The Newton update rule at a point x is $x^{new} = x - x^2 \left(7 - \frac{1}{x}\right)$. Starting with $x^{(0)} = 0$ is really infeasible and useless, since the updates will always be 0. The unique global minimizer of this function is $x^* = \frac{1}{7}$. The range of quadratic convergence for Newton's method on this function is $x \in \left(0, \frac{2}{7}\right)$. However, if you start with an initial infeasible point $x^{(0)} = 0$, the function will quadratically tend to $-\infty$!

There are some classes of functions for which theorem 84 can be applied very constructively. They are

- $-\sum_{i=1}^m \ln x_i$
- $-\ln t^2 - \mathbf{x}^T \mathbf{x}$ for $t > 0$
- $-\ln \det(X)$

Further, theorem 84 also comes handy for linear combinations of these functions. These three functions are also at the heart of modern interior points method theory.

²⁵Linear convergence adds a constant number of digits of accuracy at each iteration.

4.5.3 Variants of Newton's Method

One important aspect of the algorithm in Figure 4.49 is the step (1), which involves solving a linear system $\nabla^2 f(\mathbf{x}^{(k)})\Delta\mathbf{x}^{(k)} = \nabla f(\mathbf{x}^{(k)})$. The system can be easy to solve if the Hessian is a 100×100 sparse matrix, but it can get hairy if it is a larger and denser matrix. Thus it can be unfair to claim that the Newton's method is faster than the gradient descent method on the grounds that it takes a fewer number of iterations to converge as compared to the gradient descent, since each iteration of the Newton's method involves inverting the hessian to solve a linear system, which can take time²⁶ $O(n^3)$ for dense systems. Further, the method assumes that the hessian is positive definite and therefore invertible, which might not always be so. Finally, the Newton's method might make huge-uncontrolled steps, especially when the hessian is positive semi-definite (for example, if the function is flat along the direction corresponding to a 0 or nearly 0 eigenvalue). Due to these disadvantages, most optimization packages do not use Newton's method.

There is a whole suite of methods called *Quasi-Newton methods* that use approximations of the hessian at each iteration in an attempt to either do less work per iteration or to handle singular hessian matrices. These methods fall in between gradient methods and Newton's method and were introduced in the 1960's. Work on quasi-Newton methods sprang from the belief that often, in a large linear system, most variables should not depend on most other variables (that is, the system is generally sparse).

We should however note that in some signal and image processing problems, the hessian has a nice structure, which allows one to solve the linear system $\nabla^2 f(\mathbf{x}^{(k)})\Delta\mathbf{x}^{(k)} = \nabla f(\mathbf{x}^{(k)})$ in time much less than $O(n^3)$ (often in time comparable to that required for quasi Newton methods), without having to explicitly store the entire hessian. We next discuss some optimization techniques that use specific approximations to the hessian $\nabla^2 f(\mathbf{x})$ for specific classes of problems, by reducing the time required for computing the second derivatives.

4.5.4 Gauss Newton Approximation

The Gauss Newton method decomposes the objective function (typically for a regression problem) as a composition of two functions²⁷ $f = l \circ \mathbf{m}$; (i) the vector valued model or regression function $\mathbf{m} : \mathbb{R}^n \rightarrow \mathbb{R}^p$ and (ii) the scalar-valued loss (such as the sum squared difference between predicted outputs and target outputs) function l . For example, if m_i is $y_i - r(\mathbf{t}_i, \mathbf{x})$, for parameter vector $\mathbf{x} \in \mathbb{R}^n$ and input instances (y_i, \mathbf{t}_i) for $i = 1, 2, \dots, p$, the function f can be written as

$$f(\mathbf{x}) = \frac{1}{2} \sum_{i=1}^p (y_i - r(\mathbf{t}_i, \mathbf{x}))^2$$

²⁶ $O(n^{2.7})$ to be precise.

²⁷Here, n is the number of weights.

An example of the function r is the linear regression function $r(\mathbf{t}_i, \mathbf{x}) = \mathbf{x}^T \mathbf{t}_i$. Logistic regression poses an example objective function, which involves a cross-entropy loss.

$$f(\mathbf{x}) = - \sum_{i=1}^p (y_i \log(\sigma(\mathbf{x}^T \mathbf{t}_i)) + (1 - y_i) \log(\sigma(-\mathbf{x}^T \mathbf{t}_i)))$$

where $\sigma(k) = \frac{1}{1+e^{-k}}$ is the logistic function.

The task of the loss function is typically to make the optimization work well and this gives freedom in choosing l . Many different objective functions share a common loss function. While the sum-squared loss function is used in many regression settings, cross-entropy loss is used in many classification problems. These loss functions arise from the problem of maximizing log-likelihoods in some reasonable way.

The Hessian $\nabla^2 f(\mathbf{x})$ can be expressed using a matrix version of the chain rule, as

$$\nabla^2 f(\mathbf{x}) = \underbrace{J_{\mathbf{m}}(\mathbf{x})^T \nabla^2 l(\mathbf{m}) J_{\mathbf{m}}(\mathbf{x})}_{G_f(\mathbf{x})} + \sum_{i=1}^p \nabla^2 m_i(\mathbf{x}) (\nabla l(\mathbf{m}))_i$$

where $J_{\mathbf{m}}$ is the jacobian²⁸ of the vector valued function \mathbf{m} . It can be shown that if $\nabla^2 l(\mathbf{m}) \succeq 0$, then $G_f(\mathbf{x}) \succeq 0$. The term $G_f(\mathbf{x})$ is called the Gauss-Newton approximation of the Hessian $\nabla^2 f(\mathbf{x})$. In many situations, $G_f(\mathbf{x})$ is the dominant part of $\nabla^2 f(\mathbf{x})$ and the approximation is therefore reasonable. For example, at the point of minimum (which will be the critical point for a convex function), $\nabla^2 f(\mathbf{x}) = G_f(\mathbf{x})$. Using the Gauss-Newton approximation to the hessian $\nabla^2 f(\mathbf{x})$, the Newton update rule can be expressed as

$$\Delta \mathbf{x} = (G_f(\mathbf{x}))^{-1} \nabla f(\mathbf{x}) = (G_f(\mathbf{x}))^{-1} J_{\mathbf{m}}^T(\mathbf{x}) \nabla l(\mathbf{m})$$

where we use the fact that $(\nabla f(\mathbf{x}))_i = \sum_{k=1}^p \frac{\partial l}{\partial m_k} \frac{\partial m_k}{\partial x_i}$, since the gradient of a composite function is a product of the jacobians.

For the cross entropy classification loss or the sum-squared regression loss l , the hessian is known to be positive semi-definite. For example, if the loss function is the sum of squared loss, the objective function is $f = \frac{1}{2} \sum_{i=1}^p m_i(\mathbf{x})^2$ and $\nabla^2 l(\mathbf{m}) = I$. The Newton update rule can be expressed as

$$\Delta \mathbf{x} = (J_{\mathbf{m}}(\mathbf{x})^T J_{\mathbf{m}}(\mathbf{x}))^{-1} J_{\mathbf{m}}(\mathbf{x})^T \mathbf{m}(\mathbf{x})$$

Recall that $(J_{\mathbf{m}}(\mathbf{x})^T J_{\mathbf{m}}(\mathbf{x}))^{-1} J_{\mathbf{m}}(\mathbf{x})^T$ is the Moore-Penrose pseudoinverse $J_{\mathbf{m}}(\mathbf{x})^+$ of $J_{\mathbf{m}}(\mathbf{x})$. The Gauss-Jordan method for the sum-squared loss can be interpreted as multiplying the gradient $\nabla l(\mathbf{m})$ by the pseudo-inverse of the jacobian of \mathbf{m}

²⁸The Jacobian is a $p \times n$ matrix of the first derivatives of a vector valued function, where p is arity of \mathbf{m} . The $(i, j)^{th}$ entry of the Jacobian is the derivative of the i^{th} output with respect to the j^{th} variable, that is $\frac{\partial m_i}{\partial x_j}$. For $m = 1$, the Jacobian is the gradient vector.

instead of its transpose (which is what the gradient descent method would do). Though the Gauss-Newton method has been traditionally used for non-linear least squared problems, recently it has also seen use for the cross entropy loss function. This method is a simple adoption of the Newton's method, with the advantage that second derivatives, which can be computationally expensive and challenging to compute, are not required.

4.5.5 Levenberg-Marquardt

Like the Gauss-Newton method, the Levenberg-Marquardt method has its main application in the least squares curve fitting problem (as also in the minimum cross-entropy problem). The Levenberg-Marquardt method interpolates between the Gauss-Newton algorithm and the method of gradient descent. The Levenberg-Marquardt algorithm is more robust than the Gauss Newton algorithm - it often finds a solution even if it starts very far off the final minimum. On the other hand, for well-behaved functions and reasonable starting parameters, this algorithm tends to be a bit slower than the Gauss Newton algorithm. The Levenberg-Marquardt method aims to reduce the uncontrolled step size often taken by the Newton's method and thus fix the stability issue of the Newton's method. The update rule is given by

$$\Delta \mathbf{x} = -(G_f(\mathbf{x}) + \lambda \text{diag}(G_f))^{-1} J_{\mathbf{m}}^T(\mathbf{x}) \nabla l(\mathbf{m})$$

where G_f is the Gauss-Newton approximation to $\nabla^2 f(\mathbf{x})$ and is assumed to be positive semi-definite. This method is one of the work-horses of modern optimization. The parameter $\lambda \geq 0$ adaptively controlled, limits steps to an elliptical model-trust region²⁹. This is achieved by adding λ to the smallest eigenvalues of G_f , thus restricting all eigenvalues of the matrix to be above λ so that the elliptical region has diagonals of shorter length that inversely vary as the eigenvalues (*c.f.* page 3.11.3). While this method fixes the stability issues in Newtons method, it still requires the $O(n^3)$ time required for matrix inversion.

4.5.6 BFGS

The Broyden-Fletcher-Goldfarb-Shanno³⁰ (BFGS) method uses linear algebra to iteratively update an estimate $B^{(k)}$ of $(\nabla^2 f(\mathbf{x}^{(k)}))^{-1}$ (the inverse of the curvature matrix), while ensuring that the approximation to the hessian inverse is symmetric and positive definite. Let $\Delta \mathbf{x}^{(k)}$ be the direction vector for the k^{th} step obtained as the solution to

$$\Delta \mathbf{x}^{(k)} = -B^{(k)} \nabla f(\mathbf{x}^{(k)})$$

The next point $\mathbf{x}^{(k+1)}$ is obtained as

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + t^{(k)} \Delta \mathbf{x}^{(k)}$$

²⁹Essentially the algorithm approximates only a certain region (the so-called trust region) of the objective function with a quadratic as opposed to the entire function.

³⁰The 4 authors wrote papers for exactly the same method at exactly at the same time.

where $t^{(k)}$ is the step size obtained by line search. Let $\Delta \mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k+1)}) - \nabla f(\mathbf{x}^{(k)})$. Then the BFGS update rule is derived by imposing the following logical conditions:

1. $\Delta \mathbf{x}^{(k)} = -B^{(k)} \nabla f(\mathbf{x}^{(k)})$ with $B^{(k)} \succ 0$. That is, $\Delta \mathbf{x}^{(k)}$ is the minimizer of the convex quadratic model

$$Q^{(k)}(\mathbf{p}) = f(\mathbf{x}^{(k)}) + \nabla^T f(\mathbf{x}^{(k)}) \mathbf{p} + \frac{1}{2} \mathbf{p}^T (B^{(k)})^{-1} \mathbf{p}$$

2. $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + t^{(k)} \Delta \mathbf{x}^{(k)}$, where $t^{(k)}$ is obtained by line search.
3. The gradient of the function $Q^{(k+1)} = f(\mathbf{x}^{(k+1)}) + \nabla^T f(\mathbf{x}^{(k+1)}) \mathbf{p} + \frac{1}{2} \mathbf{p}^T (B^{(k+1)})^{-1} \mathbf{p}$ at $\mathbf{p} = \mathbf{0}$ and $\mathbf{p} = -t^{(k)} \Delta \mathbf{x}^{(k)}$ agrees with gradient of f at $\mathbf{x}^{(k+1)}$ and $\mathbf{x}^{(k)}$ respectively. While the former condition is naturally satisfied, the latter need to be imposed. This quasi-Newton condition yields

$$(B^{(k+1)})^{-1} (\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) = \nabla f(\mathbf{x}^{(k+1)}) - \nabla f(\mathbf{x}^{(k)}).$$

This equation is called the *secant equation*.

4. Finally, among all symmetric matrices satisfying the secant equation, $B^{(k+1)}$ is closest to the current matrix $B^{(k)}$ in some norm. Different matrix norms give rise to different quasi-Newton methods. In particular, when the norm chosen is the Frobenius norm, we get the following BFGS update rule

$$B^{(k+1)} = B^{(k)} + R^{(k)} + S^{(k)}$$

where,

$$R^{(k)} = \frac{\Delta \mathbf{x}^{(k)} (\Delta \mathbf{x}^{(k)})^T}{(\Delta \mathbf{x}^{(k)})^T \Delta \mathbf{g}^{(k)}} - \frac{B^{(k)} \Delta \mathbf{g}^{(k)} (\Delta \mathbf{g}^{(k)})^T (B^{(k)})^T}{(\Delta \mathbf{g}^{(k)})^T B^{(k)} \Delta \mathbf{g}^{(k)}}$$

and

$$S^{(k)} = \mathbf{u} (\Delta \mathbf{x}^{(k)})^T B^{(k)} \Delta \mathbf{x}^{(k)} \mathbf{u}^T$$

with

$$\mathbf{u} = \frac{\Delta \mathbf{x}^{(k)}}{(\Delta \mathbf{x}^{(k)})^T \Delta \mathbf{g}^{(k)}} - \frac{B^{(k)} \Delta \mathbf{g}^{(k)}}{(\Delta \mathbf{g}^{(k)})^T B^{(k)} \Delta \mathbf{g}^{(k)}}$$

We have made use of the Sherman Morrison formula that determines how updates to a matrix relate to the updates to the inverse of the matrix.

The approximation to the Hessian is updated by analyzing successive gradient vectors and thus the Hessian matrix does not need to be computed at any stage. The initial estimate $B^{(0)}$ can be taken to be the identity matrix, so that the first step is equivalent to a gradient descent. The BFGS method has a reduced complexity of $O(n^2)$ time per iteration. The method is summarized

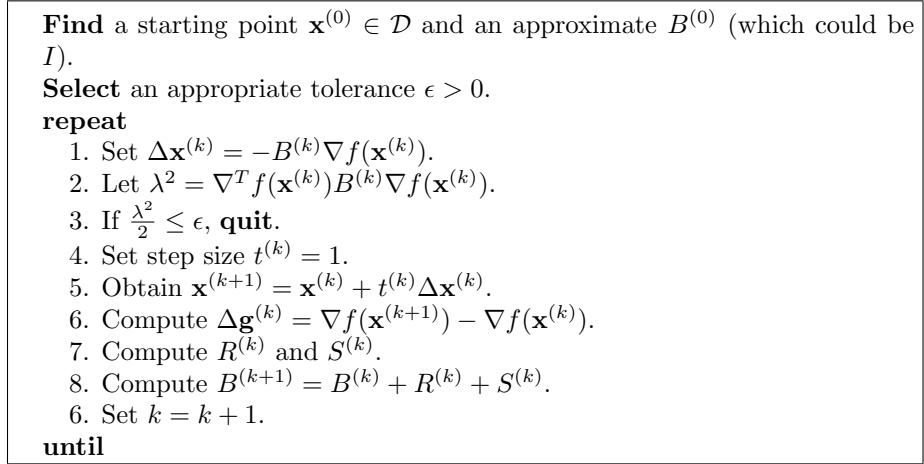


Figure 4.50: The BFGS method.

in Figure 4.50 The BFGS [?] method approaches the Newton's method in behaviour as the iterate approaches the solution. They are much faster than the Newton's method in practice. It has been proved that when BFGS is applied to a convex quadratic function with exact line search, it finds the minimizer within n steps. There is a variety of methods related to BFGS and collectively they are known as Quasi-Newton methods. They are preferred over the Newton's method or the Levenberg-Marquardt when it comes to speed. There is a variant of BFGS, called LBFGS [?], which stands for "Limited memory BFGS method". LBFGS employs a limited-memory quasi-Newton approximation that does not require much storage or computation. It limits the rank of the inverse of the hessian to some number $\gamma \in \Re$ so that only $n\gamma$ numbers have to be stored instead of n^2 numbers. For general non-convex problems, LBFGS may fail when the initial geometry (in the form of $B^{(0)}$) has been placed very close to a saddle point. Also, LBFGS is very sensitive to line search.

Recently, L-BFGS has been observed [?] to be the most effective parameter estimation method for Maximum Entropy model, much better than improved iterative scaling [?] (IIS) and generalized iterative scaling [?] (GIS).

4.5.7 Solving Systems Large Sparse Systems

In many convex optimization problems such as least squares, newton's method for optimization, *etc.*, one has to deal with solving linear systems involving large and sparse matrices. Elimination with ordering can be expensive in such cases. A lot of work has gone into solving such problems efficiently³¹ using iterative

³¹Packages such as LINPack (which is now renamed to LAPACK), EiSPACK, MINPACK, *etc.*, which can be found under the netlib repository, have focused on efficiently solving large linear systems under general conditions as well as specific conditions such as symmetry or positive definiteness of the coefficient matrix.

methods instead of direct elimination methods. An example iterative method is for solving a system $A\mathbf{x} = \mathbf{b}$ by repeated multiplication of a large and sparse matrix A by vectors to quickly get an answer $\hat{\mathbf{x}}$ that is sufficiently close to the optimal solution \mathbf{x}^* . Multiplication of an $n \times n$ sparse matrix A having k non-zero entries with a vector of dimension n takes $O(kn)$ time only, in contrast to $O(n^3)$ time for Gauss elimination. We will study three types of methods for solving systems with large and sparse matrices:

1. *Iterative Methods.*
2. *Multigrid Methods.*
3. *Krylov Methods.*

The most famous and successful amongst the Krylov methods has been the *conjugate gradient method*, which works for problems with positive definite matrices.

Iterative Methods

The central step in an iteration is

$$P\mathbf{x}_{k+1} = (P - A)\mathbf{x}_k + \mathbf{b}$$

where \mathbf{x}_k is the estimate of the solution at the k^{th} step, for $k = 0, 1, \dots$. If the iterations converge to the solution, that is, if $\mathbf{x}_{k+1} = \mathbf{x}_k$ one can immediately see that the solution is reached. The choice of matrix P , which is called the *preconditioner*, determines the rate of convergence of the solution sequence to the actual solution. The initial estimate \mathbf{x}_0 can be arbitrary for linear systems, but for non-linear systems, it is important to start with a good approximation. It is desirable to choose the matrix P reasonably close to A , though setting $P = A$ (which is referred to as perfect preconditioning) will entail solving the large system $A\mathbf{x} = \mathbf{b}$, which is undesirable as per our problem definition. If \mathbf{x}^* is the actual solution, the relationship between the errors \mathbf{e}_k and \mathbf{e}_{k+1} at the k^{th} and $(k + 1)^{\text{th}}$ steps respectively can be expressed as

$$P\mathbf{e}_{k+1} = (P - A)\mathbf{e}_k$$

where $\mathbf{e}_k = \mathbf{x}_k - \mathbf{x}^*$. This is called the *error equation*. Thus,

$$\mathbf{e}_{k+1} = (I - P^{-1}A)\mathbf{e}_k = M\mathbf{e}_k$$

Whether the solutions are convergent or not is controlled by the matrix M . The iterations are stationary (that is, the update is of the same form at every step). On the other hand, Multigrid and Krylov methods adapt themselves across iterations to enable faster convergence. The error after k steps is given by

$$\mathbf{e}_k = M^k \mathbf{e}_0 \tag{4.96}$$

Using the idea of eigenvector decomposition presented in (3.99), it can be proved that the error vector $\mathbf{e}_k \rightarrow \mathbf{0}$ if the absolute values of all the eigenvalues of M are less than 1. This is the *fundamental theorem of iteration*. In this case, the rate of convergence of \mathbf{e}_k to $\mathbf{0}$ is determined by the maximum absolute eigenvalue of M , called the spectral radius of M and denoted by $\rho(M)$.

Any iterative method should attempt to choose P so that it is easy to compute \mathbf{x}_{k+1} and at the same time, the matrix $M = I - P^{-1}A$ has small eigenvalues. Corresponding to various choices of the preconditioner P , there exist different iterative methods.

1. *Jacobi*: In the simplest setting, P can be chosen to be a diagonal matrix with its diagonal borrowed from A . This choice of A corresponds to the *Jacobi* method. The value of $\rho(M)$ is less than 1 for the Jacobi method, though it is often very close to 1. Thus, the Jacobi method does converge, but the convergence can be very slow in practice. While the residual $\hat{\mathbf{r}} = A\hat{\mathbf{x}} - \mathbf{b}$ converges rapidly, the error $\bar{\mathbf{x}} = \hat{\mathbf{x}} - \mathbf{x}^*$ decreases rapidly in the beginning, but the rate of decrease of $\bar{\mathbf{x}}$ reduces as iterations proceed. This happens because $\bar{\mathbf{x}} = A^{-1}\hat{\mathbf{r}}$ and A^{-1} happens to have large condition number for sparse matrices. In fact, it can be shown that Jacobi can take up to n^β iterations to reduce the error $\bar{\mathbf{x}}$ by a factor β .

We will take an example to illustrate the Jacobi method. Consider the following $n \times n$ tridiagonal matrix A .

$$A = \begin{bmatrix} 2 & -1 & 0 & \dots & 0 & 0 & 0 & \dots & 0 \\ -1 & 2 & -1 & \dots & 0 & 0 & 0 & \dots & 0 \\ \cdot & \cdot & \cdot & \dots & \cdot & \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \cdot & \dots & \cdot & \cdot & \cdot & \dots & \cdot \\ 0 & 0 & 0 & \dots & -1 & 2 & -1 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 & -1 & 2 & \dots & 0 \\ \cdot & \cdot & \cdot & \dots & \cdot & \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \cdot & \dots & \cdot & \cdot & \cdot & \dots & \cdot \\ 0 & 0 & 0 & \dots & 0 & 0 & 0 & \dots & 2 \end{bmatrix} \quad (4.97)$$

The absolute value of the i^{th} eigenvalue of M is $\cos \frac{j\pi}{n+1}$ and its spectral radius is $\rho(M) = \cos \frac{\pi}{n+1}$. For extremely large n , the spectral radius is approximately $1 - \frac{1}{2} \left(\frac{\pi}{n+1} \right)^2$, which is very close to 1. Thus, the Jacobi steps converge very slowly.

2. *Gauss-Seidel*: The second possibility is to choose P to be the lower-triangular part of A . The method for this choice is called the *Gauss-Seidel* method. For the example tridiagonal matrix A in (4.97), matrix

$P - A$ will be the strict but negated upper-triangular part of A . For the Gauss-Seidel technique, the components of \mathbf{x}_{k+1} can be determined from \mathbf{x}_k using back-substitution. The Gauss-sidel method provides only a constant factor improvement over the *Jacobi* method.

3. *Successive over-relaxation*: In this method, the preconditioner is obtained as a weighted composition of the preconditioners from the above two methods. It is abbreviated as *SOR*. In history, this was the first step of progress beyond *Jacobi* and *Gauss-Seidel*.
4. *Incomplete LU*: This method involves an incomplete elimination on the sparse matrix A . For a sparse matrix A , many entries in its *LU* decomposition will comprise of nearly 0 elements; the idea behind this method is to treat such entries as 0's. Thus, the L and U matrices are approximated based on the tolerance threshold; if the tolerance threshold is very high, the factors are exact. Else they are approximate.

Multigrid Methods

Multigrid methods come very handy in solving large sparse systems, especially differential equations using a hierarchy of discretizations. This approach often scales linearly with the number of unknowns n for a pre-specified accuracy threshold. The overall multi-grid algorithm for solving $A_h \mathbf{u}_h = \mathbf{b}_h$ with residual given by $\mathbf{r}_h = \mathbf{b} - A_h \mathbf{u}_h$ is

1. **Smoothing**: Perform a few (say 2-3) iterations on $A_h \mathbf{u} = \mathbf{b}_h$ using either *Jacobi* or *Gauss-sidel*. This will help remove high frequency components of the residual $\mathbf{r} = \mathbf{b} - A_h \mathbf{u}$. This step is really outside the core of the multi-grid method. Denote the solution obtained by \mathbf{u}_h . Let $\mathbf{r}_h = \mathbf{b} - A_h \mathbf{u}_h$.
2. **Restriction**: Restrict \mathbf{r}_h to coarse grid by setting $\mathbf{r}_{2h} = R\mathbf{r}_h$. That is, \mathbf{r}_h is downsampled to yield \mathbf{r}_{2h} . Let $k < n$ characterize the coarse grid. Then, the $k \times n$ matrix R is called the restriction matrix and it takes the residuals from a finer to a coarser grid. It is typically scaled to ensure that a vector of 1's on the fine mesh gets transformed to a vector of 1's on a coarse mesh. Calculations on the coarse grid are way faster than on the finer grid.
3. Solve $A_{2h} \mathbf{e}_{2h} = \mathbf{r}_{2h}$ with $A_{2h} = RA_h N$, which is a natural construction for the coarse mesh operation. This could be done by running few iterations of *Jacobi*, starting with $\mathbf{e}_{2h} = \mathbf{0}$.
4. **Interpolation/Prolongation**: This step involves interpolating the correction computed on a coarser grid to a finer grid. Interpolate back to $\mathbf{e}_h = N\mathbf{e}_{2h}$. Here N is a $k \times n$ interpolation matrix and it takes the residuals from a coarse to a fine grid. It is generally a good idea to connect N to R by setting $N = \alpha R^T$ for some scaling factor α . Add \mathbf{e}_h to \mathbf{u}_h . The

analytical expression for \mathbf{e}_h is

$$\mathbf{e}_h = N(A_{2h})^{-1}RA_h(\mathbf{u} - \mathbf{u}_h) = \underbrace{(N(RAN)^{-1}RA_h(\mathbf{u} - \mathbf{u}_h))}_S(\mathbf{u} - \mathbf{u}_h)$$

A property of the $n \times n$ matrix S is that $S^2 = S$. Thus, the only eigenvalues of S are 0 and 1. Since S is of rank $k < n$, k of its eigenvalues are 1 and $n - k$ are 0. Further, the eigenvectors for the 1 eigenvalues, which are in the null space of $I - S$ form the coarse mesh (and correspond to low frequency vectors) whereas the eigenvectors for the 0 eigenvalues, which are in the null space of S form the fine mesh (and correspond to high frequency vectors). We can easily derive that k eigenvalues of $I - S$ will be 0 and $n - k$ of them will be 1.

5. Finally as a post-smoothing step, iterate $A\mathbf{u}_h = \mathbf{b}_h$ starting from the improved $\mathbf{u}_h + \mathbf{e}_h$, using Jacobi or Gauss-Sidel.

Overall, the error \mathbf{e}^k after k steps will be of the form

$$\mathbf{e}_k = (M^t(I - S)M^t)\mathbf{e}_0 \quad (4.98)$$

where t is the number of Jacobi steps performed in (1) and (5). Typically t is 2 or 3. When you contrast (4.98) against (4.96), we discover that $\rho(M) \geq \rho(M^t(I - S)M^t)$. As t increases, $\rho(M^t(I - S)M^t)$ further decreases by a smaller proportion.

In general, you could have multiple levels of coarse grids corresponding to $2h$, $4h$, $8h$ and so on, in which case, steps (2), (3) and (4) would be repeated as many times with varying specifications of the coarseness. If A is an $n \times n$ matrix, multi-grid methods are known to run in $O(n^2)$ floating point operations (flops). The multi-grid method could be used an iterative method to solve a linear system. Alternatively, it could be used to obtain the preconditioner.

Linear Conjugate Gradient Method

The conjugate gradient method is one of the most popular Krylov methods. The Krylov matrix K_j , for the linear system $A\mathbf{u} = \mathbf{b}$ is given by

$$K_j = [\mathbf{b} \quad A\mathbf{b} \quad A^2\mathbf{b} \quad \dots \quad A^{j-1}\mathbf{b}]$$

The columns of K_j are easy to compute; each column is a result of a matrix multiplication A with the previous column. Assuming we are working with sparse matrices, (often symmetric matrices such as the Hessian) these computations will be inexpensive. The Krylov space \mathcal{K}_j is the column space of K_j . The columns of K_j are computed during the first j steps of an iterative method such as Jacobi. Most Krylov methods opt to choose vectors from \mathcal{K}_j instead of a fixed choice of the j^{th} column of K_j . A method such as *MinRes* chooses a vector

$\mathbf{u}_j \in \mathcal{K}_j$ that minimizes $\mathbf{b} - A\mathbf{u}_j$. One of the well-known Krylov methods is the *Conjugate gradient* method, which assumes that the matrix A is symmetric and positive definite and is faster than MinRes. In this method, the choice of u_j is made so that $\mathbf{b} - A\mathbf{u}_j \perp \mathcal{K}_j$. That is, the choice of \mathbf{u}_j is made so that the residual $\mathbf{r}_j = \mathbf{b} - A\mathbf{u}_j$ is orthogonal to the space \mathcal{K}_j . The conjugate gradient method gives an exact solution to the linear system if $j = n$ and that is how they were originally designed to be (and put aside subsequently). But later, they were found to give very good approximations for $j \ll n$.

The discussions that follow require the computation of a basis for \mathcal{K}_j . It is always preferred to have a basis matrix with low condition number³², and an orthonormal basis is a good choice, since it has a condition number of 1 (the basis consisting of the columns of K_j turns out to be not-so-good in practice). The *Arnoldi* method yields an orthonormal Krylov basis $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_j$ to get something that is numerically reasonable to work on. The method is summarized in Figure 4.51. Though the underlying idea is borrowed from Gram-Schmidt at every step, there is a difference; the vector \mathbf{t} is $\mathbf{t} = A\mathbf{Q}_j$ as against simply $\mathbf{t} = \mathbf{q}_j$. Will it be expensive to compute each \mathbf{t} ? Not if A is symmetric. First we note that by construction, $A\mathbf{Q} = \mathbf{Q}H$, where \mathbf{q}_j is the j^{th} column of \mathbf{Q} . Thus, $H = \mathbf{Q}^T A\mathbf{Q}$. If A is symmetric, then so is H . Further, since H has only one lower diagonal (by construction), it must have only one higher diagonal. Therefore, H must be symmetric and tridiagonal. If A is symmetric, it suffices to subtract only the components of \mathbf{t} in the direction of the last two vectors \mathbf{q}_{j-1} and \mathbf{q}_j from \mathbf{t} . Thus, for a symmetric A , the inner ‘for’ loop needs to iterate only over $i = j - 1$ and $i = j$.

Since A and H are similar matrices, they have exactly the same eigenvalues. Restricting the computation to a smaller number of orthonormal vectors (for some $k \ll n$), we can save time for computing \mathbf{Q}_k and H_k . The k eigenvalues of H_k are good approximations to the first k eigenvalues of H . This is called the *Arnoldi-Lanczos* method for finding the top k eigenvalues of a matrix.

As an example, consider the following matrix A .

$$A = \begin{bmatrix} 0.5344 & 1.0138 & 1.0806 & 1.8325 \\ 1.0138 & 1.4224 & 0.9595 & 0.8234 \\ 1.0806 & 0.9595 & 1.0412 & 1.0240 \\ 1.8325 & 0.8234 & 1.0240 & 0.7622 \end{bmatrix}$$

³²For any matrix A , the condition number $\kappa(A) = \frac{\sigma_{\max}(A)}{\sigma_{\min}(A)}$, where $\sigma_{\max}(A)$ and $\sigma_{\min}(A)$ are maximal and minimal singular values of A respectively. Recall from Section 3.13 that the i^{th} eigenvalue of $A^T A$ (the gram matrix) is the square of the i^{th} singular value of A . Further, if A is normal, $\kappa(A) = \left| \frac{\lambda_{\max}(A)}{\lambda_{\min}(A)} \right|$, where $\lambda_{\max}(A)$ and $\lambda_{\min}(A)$ are eigenvalues of A with maximal and minimal magnitudes respectively. All orthogonal, symmetric, and skew-symmetric matrices are normal. The condition number measures how much the columns/rows of a matrix are dependent on each other; higher the value of the condition number, more is the linear dependence. Condition number 1 means that the columns/rows of a matrix are linearly independent.

```

Set  $\mathbf{q}_1 = \frac{1}{\|\mathbf{b}\|} \mathbf{b}$ . //The first step in Gram schmidt.
for  $j = 1$  to  $n - 1$  do
   $\mathbf{t} = A\mathbf{q}_j$ .
  for  $i = 1$  to  $j$  do
    //If  $A$  is symmetric, it will be  $i = \max(1, j - 1)$  to  $j$ .
     $H_{i,j} = \mathbf{q}_i^T \mathbf{t}$ .
     $\mathbf{t} = \mathbf{t} - H_{i,j} \mathbf{q}_i$ .
  end for
   $H_{j+1,j} = \|\mathbf{t}\|$ .
   $\mathbf{q}_{j+1} = \frac{1}{\|\mathbf{t}\|} \mathbf{t}$ .
end for
 $\mathbf{t} = A\mathbf{q}_n$ .
for  $i = 1$  to  $n$  do
  //If  $A$  is symmetric, it will be  $i = n - 1$  to  $n$ .
   $H_{i,n} = \mathbf{q}_i^T \mathbf{t}$ .
   $\mathbf{t} = \mathbf{t} - H_{i,n} \mathbf{q}_i$ .
end for
 $H_{j+1,j} = \|\mathbf{t}\|$ .
 $\mathbf{q}_{j+1} = \frac{1}{\|\mathbf{t}\|} \mathbf{t}$ .

```

Figure 4.51: The Arnoldi algorithm for computing orthonormal basis.

and the vector \mathbf{b}

$$\mathbf{b} = \begin{bmatrix} 0.6382 & 0.3656 & 0.1124 & 0.5317 \end{bmatrix}^T$$

The matrix K_4 is

$$K_4 = \begin{bmatrix} 0.6382 & 1.8074 & 8.1892 & 34.6516 \\ 0.3656 & 1.7126 & 7.5403 & 32.7065 \\ 0.1124 & 1.7019 & 7.4070 & 31.9708 \\ 0.5317 & 1.9908 & 7.9822 & 34.8840 \end{bmatrix}$$

Its condition number is 1080.4.

The algorithm in Figure 4.51 computed the following basis for the matrix K_4 .

$$Q_4 = \begin{bmatrix} 0.6979 & -0.3493 & 0.5101 & -0.3616 \\ 0.3998 & 0.2688 & 0.2354 & 0.8441 \\ 0.1229 & 0.8965 & 0.1687 & -0.3908 \\ 0.5814 & 0.0449 & -0.8099 & -0.0638 \end{bmatrix}$$

The coefficient matrix H_4 is

$$H_4 = \begin{bmatrix} 3.6226 & 1.5793 & 0 & 0 \\ 1.5793 & 0.6466 & 0.5108 & 0 \\ 0 & 0.5108 & -0.8548 & 0.4869 \\ 0 & 0 & 0.4869 & 0.3459 \end{bmatrix}$$

and its eigenvalues are 4.3125, 0.5677, -1.2035 and 0.0835. On the other hand, the following matrix H_3 (obtained by restricting to K_3) has eigenvalues 4.3124, 0.1760 and -1.0741 .

The basic conjugate gradient method selects vectors in $\mathbf{x}_k \in \mathcal{K}_k$ that approach the exact solution to $A\mathbf{x} = \mathbf{b}$. Following are the main ideas in the conjugate gradient method.

1. The rule is to select an \mathbf{x}_k so that the new residual $\mathbf{r}_k = \mathbf{b} - A\mathbf{x}_k$ is orthogonal to all the previous residuals. Since $A\mathbf{x}_k \in \mathcal{K}_{k+1}$, we must have $\mathbf{r}_k \in \mathcal{K}_{k+1}$ and \mathbf{r}_k must be orthogonal to all vectors in \mathcal{K}_k . Thus, \mathbf{r}_k must be a multiple of \mathbf{q}_{k+1} . This holds for all k and implies that

$$\mathbf{r}_k^T \mathbf{r}_i = 0$$

for all $i < k$.

2. Consequently, the difference $\mathbf{r}_k - \mathbf{r}_{k-1}$, which is a linear combination of \mathbf{q}_{k+1} and \mathbf{q}_k , is orthogonal to each subspace \mathcal{K}_i for $i < k$.
3. Now, $\mathbf{x}_i - \mathbf{x}_{i-1}$ lies in the subspace \mathcal{K}_i . Thus, $\Delta\mathbf{r} = \mathbf{r}_k - \mathbf{r}_{k-1}$ is orthogonal to all the previous $\Delta\mathbf{x} = \mathbf{x}_i - \mathbf{x}_{i-1}$. Since $\mathbf{r}_k - \mathbf{r}_{k-1} = -A(\mathbf{x}_k - \mathbf{x}_{k-1})$, we get the following ‘conjugate directions’ condition for the updates

$$(\mathbf{x}_i - \mathbf{x}_{i-1})^T A(\mathbf{x}_k - \mathbf{x}_{k-1}) = 0$$

for all $i < k$. This is a necessary and sufficient condition for the orthogonality of the new residual to all the previous residuals. Note that while the residual updates are orthogonal in the usual inner product, the variable updates are orthogonal in the inner product with respect to A .

The basic conjugate gradient method consists of 5 steps. Each iteration of the algorithm involves a multiplication of vector \mathbf{d}_{k-1} by A and computation of two inner products. In addition, an iteration also involves around three vector updates. So each iteration should take time upto $(2+\theta)n$, where θ is determined by the sparsity of matrix A . The error \mathbf{e}_k after k iterations is bounded as follows.

$$\|\mathbf{e}_k\|_A = (\mathbf{x}_k - \mathbf{x})^T A(\mathbf{x}_k - \mathbf{x}) \leq 2 \left(\frac{\sqrt{\kappa(A)} - 1}{\sqrt{\kappa(A)} + 1} \right)^k \|\mathbf{e}_0\|$$

The ‘gradient’ part of the name *conjugate gradient* stems from the fact that solving the linear system $A\mathbf{x} = \mathbf{b}$ corresponds to finding the minimum value

```

x0 = 0, r0 = b, d0 = r0, k = 1.
repeat
  1.  $\alpha_k = \frac{\mathbf{r}_{k-1}^T \mathbf{r}_{k-1}}{\mathbf{d}_{k-1}^T A \mathbf{d}_{k-1}}$ . //Step length for next update. This corresponds to
  the entry  $H_{k,k}$ .
  2.  $\mathbf{x}_k = \mathbf{x}_{k-1} + \alpha_k \mathbf{d}_{k-1}$ .
  3.  $\mathbf{r}_k = \mathbf{r}_{k-1} - \alpha_k A \mathbf{d}_{k-1}$ . //New residual obtained using  $\mathbf{r}_k - \mathbf{r}_{k-1} =$ 
 $-A(\mathbf{x}_k - \mathbf{x}_{k-1})$ .
  4.  $\beta_k = \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{r}_{k-1}^T \mathbf{r}_{k-1}}$ . //Improvement over previous step. This corresponds
  to the entry  $H_{k,k+1}$ .
  5.  $\mathbf{d}_k = \mathbf{r}_k + \beta_k \mathbf{d}_{k-1}$ . //The next search direction, which should be
  orthogonal to the search direction just used.
  k = k + 1.
until  $\beta_k < \theta$ .

```

Figure 4.52: The conjugate gradient algorithm for solving $A\mathbf{x} = \mathbf{b}$ or equivalently, for minimizing $E(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T A \mathbf{x} - \mathbf{x}^T \mathbf{b}$.

of the convex (for positive definite A) energy function $E(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T A \mathbf{x} - \mathbf{b}^T \mathbf{x} = \mathbf{r}$ by setting its gradient $A\mathbf{x} - \mathbf{b}$ to the zero vector. The steepest descent method makes a move along at the direction of the residual \mathbf{r} at every step but it does not have a great convergence; we land up doing a lot of work to make a little progress. In contrast, as reflect in the step $\mathbf{d}_k = \mathbf{r}_k + \beta_k \mathbf{d}_{k-1}$, the conjugate gradient method makes a step in the direction of the residual, but only after removing any component β_k along the direction of the step it just took. Figures 4.53 and 4.54 depict the steps taken by the steepest descent and the conjugate descent techniques respectively, on the level-curves of the function $E(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T A \mathbf{x} - \mathbf{x}^T \mathbf{b}$, in two dimensions. It can be seen that while the steepest descent technique requires many iterations for convergence, owing to its oscillations, the conjugate gradient method takes steps that are orthogonal with respect to A (or are orthogonal in the transformed space obtained by multiplying with A), thus taking into account the geometry of the problem and taking a fewer number of steps. If the matrix A is a hessian, the steps taken by conjugate gradient are orthogonal in the local Mahalonobis metric induced by the curvature matrix A . Note that if $\mathbf{x}^{(0)} = \mathbf{0}$, the first step taken by both methods will be the same.

The conjugate gradient method is guaranteed to reach the minimum of the energy function E in exactly n steps. Further, if A has only r distinct eigenvalues, then the conjugate gradient method will terminate at the solution in at most r iterations.

4.5.8 Conjugate Gradient

We have seen that the Conjugate Gradient method in Figure 4.52 can be viewed as a minimization algorithm for the convex quadratic function $E(\mathbf{x}) =$

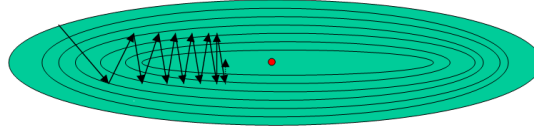


Figure 4.53: Illustration of the steepest descent technique on level curves of the function $E(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T A \mathbf{x} - \mathbf{x}^T \mathbf{b}$.

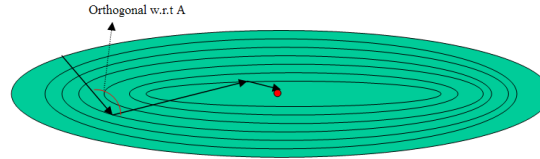


Figure 4.54: Illustration of the conjugate gradient technique on level curves of the function $E(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T A \mathbf{x} - \mathbf{x}^T \mathbf{b}$.

$\frac{1}{2}\mathbf{x}^T A \mathbf{x} - \mathbf{x}^T \mathbf{b}$. Can the approach be adapted to minimize general nonlinear convex functions? Nonlinear variants of the conjugate gradient are well studied [?] and have proved to be quite successful in practice. The general conjugate gradient method is essentially an incremental way of doing second order search.

Fletcher and Reeves showed how to extend the conjugate gradient method to nonlinear functions by making two simple changes³³ to the algorithm in Figure 4.52. First, in place of the exact line search formula in step (1) for the step length α_k , we need to perform a line search that identifies an approximate minimum of the nonlinear function f along $\mathbf{d}^{(k-1)}$. Second, the residual $\mathbf{r}^{(k)}$, which is simply the gradient of E (and which points in the direction of decreasing value of E), must be replaced by the gradient of the nonlinear objective f , which serves a similar purpose. These changes give rise to the algorithm for nonlinear optimization outlined in Figure 4.55. The search directions $\mathbf{d}^{(k)}$ are computed by Gram-Schmidt conjugation of the residuals as with linear conjugate gradient. The algorithm is very sensitive to the line minimization step and it generally requires a very good line minimization. Any line search procedure that yields an α_k satisfying the strong Wolfe conditions (see (4.90) and (4.91)) will ensure that all directions $\mathbf{d}^{(k)}$ are descent directions for the function f , otherwise, $\mathbf{d}^{(k)}$ may cease to remain a descent direction as iterations proceed. We note that each iteration of this method costs on $O(n)$, as against the Newton or quasi-newton methods which cost at least $O(n^2)$ owing to matrix operations. Most often, it yields optimal progress after $h \ll n$ iterations. Due to this property, the conjugate gradient method drives nearly all large-scale optimization today.

³³We note that in the algorithm in Figure 4.52, the residuals $\mathbf{r}^{(k)}$ in successive iterations (which are gradients of E) are orthogonal to each other, while the corresponding update directions are orthogonal with respect to A . While the former property is difficult to enforce for general non-linear functions, the latter condition can be enforced.

Select $\mathbf{x}^{(0)}$, Let $f_0 = f(\mathbf{x}^{(0)})$, $\mathbf{f}_0 = \nabla f(\mathbf{x}^{(0)})$, $\mathbf{d}^{(0)} = -\nabla f_0$, $k = 1$.

repeat

1. Compute α_k by line search.
2. Set $\mathbf{x}^{(k)} = \mathbf{x}^{(k-1)} + \alpha_k \mathbf{d}^{(k-1)}$.
3. Evaluate $\mathbf{f}^{(k)} = \nabla f(\mathbf{x}^{(k)})$.
4. $\beta_k = \frac{(\mathbf{f}^{(k)})^T \mathbf{f}^{(k)}}{(\mathbf{f}^{(k-1)})^T \mathbf{f}^{(k-1)}}$.
5. $\mathbf{d}_k = -\mathbf{f}^{(k)} + \beta_k \mathbf{d}^{(k-1)}$.

$k = k + 1$.

until $\frac{\|\mathbf{f}^{(k)}\|}{\mathbf{f}^{(0)}} < \theta$ OR $k > \text{maxIter}$.

Figure 4.55: The conjugate gradient algorithm for optimizing nonlinear convex function f .

It revolutionized optimization ever since it was invented in 1954.

Variants of the Fletcher-Reeves method use different choices of the parameter β_k . An important variant, proposed by Polak and Ribiere, defines β_k as

$$\beta_k^{PR} = \frac{(\mathbf{f}^{(k)})^T (\mathbf{f}^{(k)} - \mathbf{f}^{(k-1)})}{(\mathbf{f}^{(k)})^T \mathbf{f}^{(k)}}$$

The Fletcher-Reeves method converges if the starting point is sufficiently close to the desired minimum. However, convergence of the Polak-Ribiere method can be guaranteed by choosing

$$\beta_k = \max \{ \beta_k^{PR}, 0 \}$$

Using this value is equivalent to restarting³⁴ conjugate gradient if $\beta_k^{PR} < 0$. In practice, the Polak-Ribiere method converges much more quickly than the Fletcher-Reeves method. It is generally required to restart the conjugate gradient method after every n iterations, in order to get back conjugacy, *etc.*

If we choose f to be the strongly convex quadratic E and α_k to be the exact minimizer, this algorithm reduces to the linear conjugate gradient method. Unlike the linear conjugate gradient method, whose convergence properties are well understood and which is known to be optimal (see page 319), nonlinear conjugate gradient methods sometimes show bizarre convergence properties. It has been proved by Al-Baali that if the level set $\mathcal{L} = \{\mathbf{x} | f(\mathbf{x}) \leq f(\mathbf{x}^{(0)})\}$ of a convex function f is bounded and in some open neighborhood of \mathcal{L} , f is Lipschitz continuously differentiable and that the algorithm is implemented with a line search that satisfies the strong Wolfe conditions, with $0 < c_1 < c_2 < 1$, then

$$\liminf_{k \rightarrow \infty} \|\mathbf{f}^{(k)}\| = 0$$

³⁴Restarting conjugate gradient means forgetting the past search directions, and start it anew in the direction of steepest descent.

In summary, quasi-Newton methods are robust. But, they require $O(n^2)$ memory space to store the approximate Hessian inverse, and so they are not directly suited for large scale problems. Modifications of these methods called Limited Memory Quasi-Newton methods use $O(n)$ memory and they are suited for large scale problems. Conjugate gradient methods also work well and are well suited for large scale problems. However they need to be implemented carefully, with a carefully set line search. In some situations block coordinate descent methods (optimizing a selected subset of variables at a time) can be very much better suited than the above methods.

4.6 Algorithms for Constrained Minimization

The general form of constrained convex optimization problem was given in (4.20) and is restated below.

$$\begin{aligned} & \text{minimize} && f(x) \\ & \text{subject to} && g_i(x) \leq 0, \quad i = 1, \dots, m \\ & && Ax = b \end{aligned} \tag{4.99}$$

For example, when f is linear and g_i 's are polyhedral, the problem is a linear program, which was stated in (4.83) and whose dual was discussed on page 287. Linear programming is a typical example of constraint minimization problem and will form the subject matter for discussion in Section 4.7. As another example, when f is quadratic (of the form $\mathbf{x}^T Q \mathbf{x} + \mathbf{b}^T \mathbf{x}$) and g_i 's are polyhedral, the problem is called a quadratic programming problem. A special case of quadratic programming is the least squares problem, which we will take up in details in Section 4.8.

4.6.1 Equality Constrained Minimization

The simpler form of constrained convex optimization is when there is only the equality constrained in problem (4.99) and it turns out to be not much different from the unconstrained case. The equality constrained convex problem can be more explicitly stated as in (4.100).

$$\begin{aligned} & \text{minimize} && f(x) \\ & \text{subject to} && Ax = b \end{aligned} \tag{4.100}$$

where f is a convex and twice continuously differentiable function and $A \in \mathbb{R}^{p \times n}$ has rank p . We will assume that the finite primal optimal value p^* is attained by f at some point $\hat{\mathbf{x}}$. The following fundamental theorem for the equality constrained convex problem (4.100) can be derived using the KKT conditions stated

in Section 4.4.4 that were proved to be necessary and sufficiency conditions for optimality of a convex problem with differentiable objective and constraint functions.

Theorem 85 $\hat{\mathbf{x}}$ is optimal point for the primal iff there exists a $\hat{\boldsymbol{\mu}}$ such that the following conditions are satisfied.

$$\begin{aligned}\nabla f(\hat{\mathbf{x}}) + A^T \hat{\boldsymbol{\mu}} &= \mathbf{0} \\ A\hat{\mathbf{x}} &= \mathbf{b}\end{aligned}\tag{4.101}$$

The term $\nabla f(\hat{\mathbf{x}}) + A^T \hat{\boldsymbol{\mu}}$ is sometimes called the dual residual (r_d) while the term $A\hat{\mathbf{x}} - \mathbf{b}$ is referred to as the primal residual (r_p). The optimality condition basically states that both r_d and r_p should both be 0 and the success of this test is a certificate of optimality.

As an illustration of this theorem, consider the constrained quadratic problem

$$\begin{aligned}\text{minimize} \quad & \frac{1}{2} \mathbf{x}^T A \mathbf{x} + \mathbf{b}^T \mathbf{x} + c \\ \text{subject to} \quad & P \mathbf{x} = \mathbf{q}\end{aligned}\tag{4.102}$$

By theorem 85, the necessary and sufficient condition for optimality of a point $(\hat{\mathbf{x}}, \hat{\boldsymbol{\lambda}})$ is

$$\underbrace{\begin{bmatrix} A & P^T \\ P & 0 \end{bmatrix}}_{\text{KKT matrix}} \begin{bmatrix} \hat{\mathbf{x}} \\ \hat{\boldsymbol{\lambda}} \end{bmatrix} = \begin{bmatrix} -\mathbf{b} \\ \mathbf{q} \end{bmatrix}$$

The KKT matrix³⁵ is nonsingular iff, $P + A^T A \succ 0$. In such an event, the system of $n + p$ linear equations in $n + p$ unknowns will have a unique solution corresponding to the point of global minimum of (4.102). The linearly constrained least squared problem is a specific example of this and is discussed in Section 4.8.2.

Eliminating Equality Constraints

Figure 3.3 summarized the number of solutions to the system $A\mathbf{x} = \mathbf{b}$ under different conditions. In particular, when the rank of A is the number of its rows (p) and is less than the number of its columns (n), there are infinitely many solutions. This was logically derived in (3.35), and we restate it here for reference:

$$\mathbf{x}_{\text{complete}} = \mathbf{x}_{\text{particular}} + \mathbf{x}_{\text{nullspace}}$$

where the three vectors are defined with respect to the reduced row echelon form R of A (c.f. Section 3.6.2):

³⁵This matrix comes up very often in many areas such as optimization, mechanics, etc.

1. $\mathbf{x}_{complete}$: specifies any solution to $A\mathbf{x} = \mathbf{b}$
2. $\mathbf{x}_{particular}$: is obtained by setting all free variables (corresponding to columns with no pivots) to 0 and solving $A\mathbf{x} = \mathbf{b}$ for pivot variables.
3. $\mathbf{x}_{nullspace}$: is any vector in the null space of the matrix A , obtained as a linear combination of the basis vectors for $N(A)$.

Using formula (3.27) on page 167 to derive the null basis $N \in \mathbb{R}^{n \times n-p}$ (that is, $AN = 0$ and the columns of N span $N(A)$), we get the following free parameter expression for the solution set to $A\mathbf{x} = \mathbf{b}$:

$$\{\mathbf{x} | A\mathbf{x} = \mathbf{b}\} = \{N\mathbf{z} + \mathbf{x}_{particular} | \mathbf{z} \in \mathbb{R}^{n-p}\}$$

We can express the constrained problem in (4.100) in terms of the variables $\mathbf{z} \in \mathbb{R}^{n-p}$ (that is through an affine change of coordinates) to get the following equivalent problem:

$$\underset{\mathbf{z} \in \mathbb{R}^{n-p}}{\text{minimize}} \quad f(N\mathbf{z} + \mathbf{x}_{particular}) \quad (4.103)$$

This problem is equivalent to the original problem in (4.100), has no equality constraints and has p fewer variables. The optimal solutions $\hat{\mathbf{x}}$ and $\hat{\boldsymbol{\mu}}$ to the primal and dual of (4.100) respectively can be expressed in terms of the optimal solution $\hat{\mathbf{z}}$ to (??) as:

$$\begin{aligned} \hat{\mathbf{x}} &= N\hat{\mathbf{z}} + \mathbf{x}_{particular} \\ \hat{\boldsymbol{\mu}} &= -(AA^T)^{-1}A\nabla f(\hat{\mathbf{x}}) \end{aligned} \quad (4.104)$$

Any iterative algorithm that is applied to solve the problem (4.104) will ensure that all intermediate points are feasible, since for any $\mathbf{z} \in \mathbb{R}^{n-p}$, $\mathbf{x} = N\mathbf{z} + \mathbf{x}_{particular}$ is feasible, that is, $A\mathbf{x} = \mathbf{b}$. However, when the Newton's method is applied, the iterates are independent of the exact affine change of coordinates induced by the choice of the null basis N (c.f. page 304). The Newton update rule $\Delta\mathbf{z}^{(k)}$ for (4.103) is given by the solution to:

$$N\nabla^2 f(N\mathbf{z}^{(k)} + \mathbf{x}_{particular})N^T \Delta\mathbf{z}^{(k)} = N\nabla f(N\mathbf{z}^{(k)} + \mathbf{x}_{particular})$$

Due the affine invariance of Newton's method, if $\mathbf{z}^{(0)}$ is the starting iterate and $\mathbf{x}^{(0)} = N\mathbf{z}^{(0)} + \mathbf{x}_{particular}$, the k^{th} iterate $\mathbf{x}^{(k)} = N\mathbf{z}^{(k)} + \mathbf{x}_{particular}$ is independent of the choice of the null basis N . We therefore do not need separate convergence analysis. The algorithm for the Newton's method was outlined in Figure 4.49. Techniques for handling constrained optimization using Newton's method given an infeasible starting point $\mathbf{x}^{(0)}$ can be found in [?].

4.6.2 Inequality Constrained Minimization

The general inequality constrained convex minimization problem is

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && g_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m \\ & && A\mathbf{x} = b \end{aligned} \tag{4.105}$$

where f as well as the g_i 's are convex and twice continuously differentiable. As in the case of equality constrained optimization, we will assume that $A \in \mathbb{R}^{p \times n}$ and has rank p . Further, we will also assume that the finite primal optimal value p^* is attained by f at some point $\hat{\mathbf{x}}$. Finally, we will assume that the Slater's constraint qualification (*c.f.* page 290) conditions hold so that strong duality holds and the dual optimum is attained. Linear programs (LP), quadratically constrained quadratic programs (QCQP) (all listed in table 4.4 on page 290) and geometric programs³⁶ (GP) are some examples of convex optimization problems with inequality constraints. An example geometric program (in its convex form) is

$$\begin{aligned} & \text{minimize}_{\mathbf{y} \in \mathbb{R}^n} && \log \left(\sum_{k=1}^q e^{\mathbf{a}_k^T \mathbf{y} + b_k} \right) \\ & \text{subject to} && \log \left(\sum_{k=1}^r e^{\mathbf{c}_k^T \mathbf{y} + d_k} \right) \leq 0 \quad i = 1, 2, \dots, p \\ & && \mathbf{g}_i^T \mathbf{y} + h_i \quad i = 1, 2, \dots, m \end{aligned} \tag{4.106}$$

Semi-definite programs (SDPs) do not satisfy conditions such as zero duality gap, *etc.*, but can be handled by extensions of interior-point methods to problems having generalized inequalities.

Logarithmic Barrier

One idea for solving a minimization problem with inequalities is to replace the inequalities by a so-called barrier term. The barrier term is subtracted from the objective function with a weight μ on it. The solution to (4.105) is approximated by the solution to the following problem.

$$\begin{aligned} & \text{minimize} && B(\mathbf{x}, \mu) = f(\mathbf{x}) - \mu \sum_{i=1}^m \ln(-g_i(\mathbf{x})) \\ & \text{subject to} && A\mathbf{x} = b \end{aligned} \tag{4.107}$$

³⁶Although geometric programs are not convex in their natural form, they can, however, be transformed to convex optimization problems, by a change of variables and a transformation of the objective and constraint functions.

The objective function $B(\mathbf{x}, \mu)$ is called the *logarithmic barrier function*. This function is convex, which can be proved by invoking the composition rules described in Section 4.2.10. It is also twice continuously differentiable. The barrier term, as a function of \mathbf{x} approaches $+\infty$ as any feasible interior point \mathbf{x} approaches the boundary of the feasible region. Because we are minimizing, this property prevents the feasible iterates from crossing the boundary and becoming infeasible. We will denote the point of optimality $\hat{\mathbf{x}}(\mu)$ as a function of μ .

However, the optimal solution to the original problem (a typical example being the LP discussed in Section 4.7) is typically a point on the boundary of the feasible region (we will see this in the case of linear programming in Section 4.7). To obtain such a boundary point solution, it is necessary to keep decreasing the parameter μ of the barrier function to 0 in the limit. As a very simple example, consider the following inequality constrained optimization problem.

$$\begin{aligned} & \text{minimize} && x^2 \\ & \text{subject to} && x \geq 1 \end{aligned}$$

The logarithmic barrier formulation of this problem is

$$\text{minimize} \quad x^2 - \mu \ln(x - 1)$$

The unconstrained minimizer for this convex logarithmic barrier function is $\hat{\mathbf{x}}(\mu) = \frac{1}{2} + \frac{1}{2}\sqrt{1 + 2\mu}$. As $\mu \rightarrow 0$, the optimal point of the logarithmic barrier problem approaches the actual point of optimality $\hat{\mathbf{x}} = 1$ (which, as we can see, lies on the boundary of the feasible region). The generalized idea, that as $\mu \rightarrow 0$, $f(\hat{\mathbf{x}}) \rightarrow p^*$ (where p^* is the optimal for (4.105)) will be proved next.

Properties of the estimate $f(\hat{\mathbf{x}}(\mu))$

The following are necessary and sufficient conditions for $\hat{\mathbf{x}}(\mu)$ to be a solution to (4.107) for a fixed μ (see KKT conditions in (4.88)):

1. The point $\hat{\mathbf{x}}(\mu)$ must be strictly feasible. That is,

$$A\hat{\mathbf{x}}(\mu) = \mathbf{b}$$

and

$$g_i(\hat{\mathbf{x}}(\mu)) < 0$$

2. There must exist a $\eta \in \Re^p$ such that

$$\nabla f(\hat{\mathbf{x}}(\mu)) + \sum_{i=1}^m \frac{-\mu}{g_i(\hat{\mathbf{x}}(\mu))} \nabla g_i(\hat{\mathbf{x}}(\mu)) + A^T \hat{\eta} = \mathbf{0} \quad (4.108)$$

Define

$$\widehat{\lambda}_i(\mu) = \frac{-\mu}{g_i(\widehat{\mathbf{x}}(\mu))}$$

and

$$\widehat{\eta}(\mu) = \widehat{\eta}\mu$$

We claim that the pair $(\widehat{\lambda}(\mu), \widehat{\eta}(\mu))$ is dual feasible. The following steps prove our claim

1. Since $g_i(\widehat{\mathbf{x}}(\mu)) < 0$ for $i = 1, 2, \dots, m$, $\widehat{\lambda}(\mu) \succ \mathbf{0}$.
2. Based on the proof of theorem 82, we can infer that $L(\mathbf{x}, \lambda, \eta)$ is convex in \mathbf{x} .

$$L(\mathbf{x}, \lambda, \eta) = f(\mathbf{x}) + \sum_{i=1}^m \lambda_i g_i(\mathbf{x}) + \eta^T (A\mathbf{x} - \mathbf{b})$$

Since the lagrangian is convex in \mathbf{x} and since it is differentiable on its domain, from (4.108), we can conclude that $\widehat{\mathbf{x}}(\mu)$ is a critical point of $L(\mathbf{x}, \lambda, \eta)$ and therefore minimizes it for $(\widehat{\lambda}(\mu), \widehat{\eta}(\mu))$.

3. That is, the dual $L^*(\widehat{\lambda}(\mu), \widehat{\eta}(\mu))$ is defined and therefore, $(\widehat{\lambda}(\mu), \widehat{\eta}(\mu))$ is dual feasible.

$$L^*(\widehat{\lambda}(\mu), \widehat{\eta}(\mu)) = f(\widehat{\mathbf{x}}(\mu)) + \sum_{i=1}^m \widehat{\lambda}_i g_i(\widehat{\mathbf{x}}(\mu)) + \widehat{\eta}(\mu)^T (A\widehat{\mathbf{x}}(\mu) - \mathbf{b}) = f(\widehat{\mathbf{x}}(\mu)) - m\mu \quad (4.109)$$

From the weak duality theorem 81, we know that $d^* \leq p^*$, where d^* and p^* are the primal and dual optimals respectively, for (4.105). Since $L^*(\widehat{\lambda}(\mu), \widehat{\eta}(\mu)) \leq d^*$ (by definition), we will have from (4.109), $f(\widehat{\mathbf{x}}(\mu)) - m\mu \leq p^*$. Or equivalently,

$$f(\widehat{\mathbf{x}}(\mu)) - p^* \leq m\mu \quad (4.110)$$

The inequality in (4.110) forms the basis of the barrier method; it confirms the intuitive idea that $\widehat{\mathbf{x}}(\mu)$ converges to an optimal point as $\mu \rightarrow 0$. We will next discuss the barrier method.

The Barrier Method

The barrier method is a simple extension of the unconstrained minimization method to inequality constrained minimization. This method is based on the property in (4.110). This method solves a sequence of unconstrained (or linearly constrained) minimization problems, using the last point found as the starting point for the next unconstrained minimization problem. It computes $\widehat{\mathbf{x}}(\mu)$ for a

sequence of decreasing values of μ , until $m\mu \leq \epsilon$, which guarantees that we have an ϵ -suboptimal solution of the original problem. It was originally proposed as the sequential unconstrained minimization technique (SUMT) technique by Fiacco and McCormick in the 1960s. A simple version of the method is outlined in Figure 4.56.

Find a strictly feasible starting point $\hat{\mathbf{x}}$, $\mu = \mu^{(0)} > 0$, $\alpha > 0$.
Select an appropriate tolerance $\epsilon > 0$.
repeat
 1. **Centering Step:** Compute $\hat{\mathbf{x}}(\mu)$ by minimizing $B(\mathbf{x}, \mu)$ (optionally subject to $A\mathbf{x} = \mathbf{b}$) starting at \mathbf{x} .
 2. Update $\mathbf{x} = \hat{\mathbf{x}}(\mu)$.
 3. If $m\mu \leq \epsilon$, **quit**.
 4. Decrease μ : $\mu = \alpha\mu$.
until

Figure 4.56: The Barrier method.

The centering step (1) can be executed using any of the descent techniques discussed in Section 4.5. It can be proved [?] that the duality gap is $m\mu^{(0)}\alpha^k$ after k iterations. Therefore, the desired accuracy ϵ can be achieved by the

barrier method after exactly $\left\lceil \frac{\log\left(\frac{m\mu^{(0)}}{\epsilon}\right)}{-\log(\alpha)} \right\rceil$ steps.

Successive minima $\hat{\mathbf{x}}(\mu)$ of the Barrier function $B(\mathbf{x}, \mu)$ can be shown to have the following properties. Let $\bar{\mu} < \mu$ for sufficiently small μ , then

1. $B(\hat{\mathbf{x}}(\bar{\mu}), \bar{\mu}) < B(\hat{\mathbf{x}}(\mu), \mu)$
2. $f(\hat{\mathbf{x}}(\bar{\mu})) \leq f(\hat{\mathbf{x}}(\mu))$
3. $-\sum_{i=1}^m \ln(-g_i(\hat{\mathbf{x}}(\bar{\mu}))) \geq -\sum_{i=1}^m \ln(-g_i(\hat{\mathbf{x}}(\mu)))$

When a strictly feasible point $\hat{\mathbf{x}}$ is not known, the barrier method is preceded by a preliminary stage, called phase I, in which a strictly feasible point is computed (if it exists). The strictly feasible point found during phase I is then used as the starting point for the barrier method. This is discussed in greater details in [?].

4.7 Linear Programming

Linear programming has been widely used in the industry for maximizing profits, minimizing costs, *etc.* The word *linear* implies that the cost function is linear in the form of an inner product.

The inputs to the program are

1. \mathbf{c} , a cost vector of size n .

2. An $m \times n$ matrix A .
3. A vector \mathbf{b} of size m .

The unknown is a vector \mathbf{x} of size n , and this is what we will try to determine.

In linear programming (LP), the task is to minimize a linear objective function of the form $\sum_{j=1}^n c_j x_j$, subject to linear inequality constraints³⁷ of the form

$\sum_{j=1}^n a_{ij} x_j \geq b_i$, $i = 1, \dots, m$ and $x_i > 0$. The problem can be stated as in (4.111). In contrast to the LP specification on page 287, where the constraint $\mathbf{x} \geq \mathbf{0}$ was absorbed into the more general constraint $-\mathbf{A}\mathbf{x} + \mathbf{b} \leq \mathbf{0}$, here we choose to specify it as a separate constraint.

$$\begin{array}{ll} \min_{\mathbf{x} \in \mathbb{R}^n} & \mathbf{x}^T \mathbf{c} \\ \text{subject to} & -\mathbf{A}\mathbf{x} + \mathbf{b} \leq \mathbf{0} \quad \mathbf{x} \geq \mathbf{0} \end{array} \quad (4.111)$$

The flip side of this problem is that it has no analytical formula as a solution. However, that does not make a big difference in practice, because there exist reliable and efficient algorithms and software for linear programming. The computational time is roughly proportional to $n^2 m$, if $m \geq n$. This is basically the cost of one iteration in an interior point method.

Linear programming (*LP*) problems are harder to recognize in practice and often need reformulations to get into the standard form in (4.111). Minimizing a piecewise linear function of x is not an *LP*, though it can be written and solved as an *LP*. Other problems involving 1 or ∞ norms can also be written as linear programming problems.

The basis for linear programming was mentioned on page 248; linear functions have no critical points and therefore, by theorem 60, the extreme values are always assumed at the boundary of the feasible set. In the case of linear programs, the feasible set is itself defined by linear inequalities: $\{\mathbf{x} \mid -\mathbf{A}\mathbf{x} + \mathbf{b} \leq \mathbf{0}\}$. Applying the argument recursively, it can be proved that the extreme values for a linear program are assumed at some corners (*i.e.*, vertices) of the feasible set. A *corner* is the intersection point of n different planes, each given by a single equation. That is, a corner point is obtained by turning n of the $n+m$ inequalities into equalities and finding their intersection³⁸. An edge is the intersection of $n-1$ inequalities and connects two corners. Geometrically, it can be observed that when you maximize or minimize some linear function, as your progress in one direction in the search space, the objective will either increase monotonically or decrease monotonically. Therefore, the maximum and minimum will be found at the corners of the allowed region.

³⁷It is a rare feature to have linear inequality constraints.

³⁸In general, there are $\frac{(n+m)!}{n!m!}$ intersections.

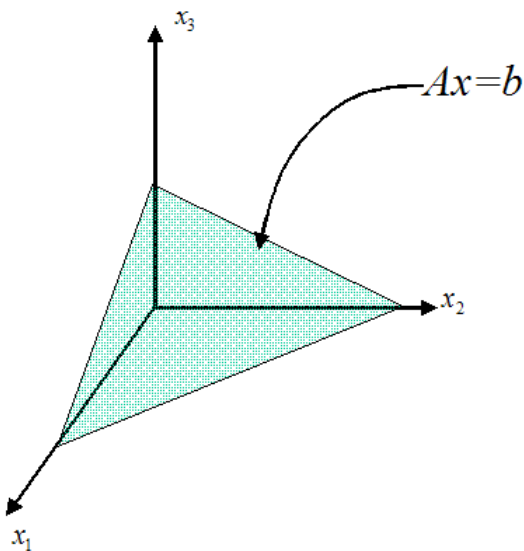


Figure 4.57: Example of the feasible set of a linear program for $n = 3$.

The feasible set is in the form of a finite interval in n dimensions. Figure 4.57 pictorially depicts a typical example of the feasible region for $n = 3$. The constraints $A\mathbf{x} \geq \mathbf{b}$ and $\mathbf{x} \geq \mathbf{0}$ would allow a tetrahedron or pyramid in the first (or completely positive) octant. If the constraint was an equality, $A\mathbf{x} = \mathbf{b}$, the feasible set would be the shaded triangle in the figure. In general for any n , the constraint $A\mathbf{x} \geq \mathbf{b}$, $\mathbf{x} \geq \mathbf{0}$ would yield as the feasible set, a polyhedron. The task of maximizing (or minimizing) the linear objective function $\mathbf{x}^T \mathbf{c} = \sum_{i=1}^n c_i x_i$ translates to finding a solution at one of the corners of the feasible region. Corners are points where some of the inequality constraints are tight or active, and others are not. At the corners, some of the inequality constraints translate to equalities. It is just a question of finding the right corner.

Why not just search all corners for the optimal answer? The trouble is that there are lots of corners. In n dimensions, with m constraints, the number of corners grows exponentially and there is no way to check all of them. There is an interesting competition between two quite different approaches for solving linear programs:

1. *The simplex method*
2. *Interior point barrier method*

4.7.1 Simplex Method

The simplex algorithm [?] is one of the fundamental methods for linear programming, developed in the late 1940s by Dantzig. It is the best established approach for solving linear problems. In the worst case, the algorithm takes a number of steps that is exponential in n ; but, in practice it is the most efficient method for solving linear programs.

The simplex method first constructs an admissible solution at a corner (which can be quite a bit of a job) of the polyhedron and then moves along its edges to vertices with successively higher values of the objective function until the optimum is reached. The movement along an edge originating at a vertex is performed by ‘loosening’ one of the inequalities that were tight at the vertex. The inequality chosen for ‘loosening’ is the one promising the fastest drop in the objective function $\mathbf{x}^T \mathbf{c}$. The rate of decrease along an edge can be measured using the gradient of the objective. This procedure is carried out iteratively, till the method encounters a vertex which has no edge (constraint) that is a promising descent direction (which means that the cost goes up along all edges incident at that vertex). Since an edge corresponding to decreasing value of the objective cannot correspond to its increasing value, no edge will be traversed twice in this process.

We will first rewrite the constraints $A\mathbf{x} \geq \mathbf{b}$ in the above LP as equations, by introducing a new non-negative “slack” variable s_j for the j^{th} constraint (for all j ’s) and subtracting it from the left-hand side of each inequality:

$$A\mathbf{x} - \mathbf{s} = \mathbf{b}$$

or equivalently in matrix notation

$$[-A \quad I] \begin{bmatrix} \mathbf{x} \\ \mathbf{s} \end{bmatrix} = -\mathbf{b}$$

We will treat the $m \times n + m$ matrix $M = [-A \quad I]$ as our new coefficient matrix and $\mathbf{y} = [\mathbf{x} \quad \mathbf{s}]^T$ as our new variable vector. With this, the above constraint can be rewritten as

$$M\mathbf{y} = -\mathbf{b}$$

The feasible set is now governed by these m equality constraints and the $n + m$ non-negativity constraints $\mathbf{x} \geq \mathbf{0}$ and $\mathbf{y} \geq \mathbf{0}$. The original cost vector \mathbf{c} is extended to a vector \mathbf{d} by appending m more zero components. This leaves us with the following problem, equivalent to the original LP (4.111).

$$\begin{aligned} \min_{\mathbf{y} \in \mathbb{R}^{n+m}} \quad & \mathbf{y}^T \mathbf{d} \\ \text{subject to} \quad & M\mathbf{y} = -\mathbf{b} \quad \mathbf{y} \geq \mathbf{0} \end{aligned} \tag{4.112}$$

We will assume that the matrix A (and therefore M) is of full row rank, that is of rank m . In practice, a preprocessing phase is applied to the user-supplied

data to remove some redundancies from the given constraints to get a full row rank matrix.

The following definitions and observations will set the platform for the simplex algorithm, which we will describe subsequently.

1. A vector \mathbf{y} is a *basic feasible point* if it is feasible and if there exists a subset \mathcal{B} of the index set $\{1, 2, \dots, n\}$ such that
 - (a) \mathcal{B} contains exactly m indices.
 - (b) $\mathbf{y} \geq \mathbf{0}$.
 - (c) $y_i \geq 0$ can be inactive (that is $y_i > 0$) only if $i \in \mathcal{B}$. In other words, $i \notin \mathcal{B} \Rightarrow y_i = 0$.
 - (d) If \mathbf{m}_i is the i^{th} column of M , the $m \times m$ matrix B defined as $B = [\mathbf{m}_i]_{i \in \mathcal{B}}$ is nonsingular.

A set \mathcal{B} satisfying these properties is called a *basis* for the problem (4.112). The corresponding matrix B is called the basis matrix. Any variable y_i for $i \in \mathcal{B}$ is called a *basic variable*, while any variable y_i for $i \notin \mathcal{B}$ is called a *free variable*.

2. It can be seen that all basic feasible points of (4.112) are corners of the feasible simplex $\mathcal{S} = \{\mathbf{x} | A\mathbf{x} \geq \mathbf{b}, \mathbf{x} \geq \mathbf{0}\}$ and vice versa. In other words, a corner of \mathcal{S} corresponds to a point \mathbf{y} in the new representation that has n components as zeroes.
3. Any two corners connected by an edge will have exactly $m - 1$ common basic variables. Each corner has n incident edges (corresponding to the addition of any one of n new basic variables and the corresponding drop of a basic variable).
4. Further, it can be proved that
 - (a) If (4.112) has a nonempty feasible region, then there is at least one basic feasible point
 - (b) If (4.112) has solutions, then at least one such solution is a basic optimal point
 - (c) If (4.112) is feasible and bounded, then it has an optimal solution.

This is known as the fundamental theorem of linear programming.

Using the ideas and notations presented above, the simplex algorithm can be outlined as follows.

1. Each iterate generated by the simplex algorithm is a *basic feasible point* of (4.112).

2. *Entering free variable:* The next iterate is determined by moving along an edge from one basic feasible solution to another. As discussed above, movement along an edge will mean that $m - 1$ variables will remain basic while one will become free. On the other hand, a new free variable will become basic. The real decision is which variable should be removed from the basis and which should be added. The idea in the simplex algorithm is to include that free variable y_k , which has the most negative component d_k (something like steepest descent in the L_1 norm).
3. *Leaving basic variable:* The basic variable from the current basis that will leave next is determined using a *pivot rule*. A commonly applied pivot rule is to determine the leaving basic variable through the constraint (say the j^{th} one) that has the smallest non-negative ratio of the right hand side b_j of the constraint to the coefficient m_{jk} of the entering variable y_k . If the coefficients of y_k are negative in all the constraints, it implies an unbounded case; the cost can be made $-\infty$ by arbitrarily increasing the value of y_k .
4. In order to facilitate the easy identification of leaving basic variables, we bring the equations into a form such that the basic variables stand by themselves. This is done by treating the new entering variable y_k as a 'pivot' in the j^{th} equation and substituting its value in terms of the other variables in the j^{th} equation into the other equations (as well as the cost function $\mathbf{y}^T \mathbf{d}$). In this form,
 - (a) the protocol is that variables corresponding to all columns of M that are in unit form are basic variables, while the rest are free variables.
 - (b) the choice of an equality in the step above automatically entails the choice of the leaving variable - the basic variable y_l corresponding to row j will be the next leaving variable.
5. In a large problem, it is possible for a leaving variable to reenter the basis at a later stage. Unless there is *degeneracy*, the costs keep going down and it can never happen that all of the m basic variables are the same as before. Thus, no corner is revisited and the method must end at the optimal corner or conclude that the cost is unbounded below. Degeneracy is said to occur if more than the usual n components of \mathbf{x} are 0 (in which case, cycling might occur but extremely rarely).

Since each simplex step involves decisions (choice of entering and leaving basic variables) and row operations (pivoting *etc.*), it is convenient to fit the data into a large matrix or *tableau*. The operations of the simplex method outlined above can be systematically translated to operations on the tableau.

1. The starting tableau is just a bigger $m + 1 \times m + n$ matrix

$$T = \begin{bmatrix} M & -\mathbf{b} \\ \mathbf{d} & \mathbf{0} \end{bmatrix}$$

2. Our first step is to get one basic variable alone on each row. Without loss of generality, we will renumber the variables and rearrange the corresponding coefficients of M so that at every iteration, y_1, y_2, \dots, y_m are the basic variables and the rest are free (*i.e.*, 0). The first m columns of A form an $m \times m$ square matrix B and the last n form an $m \times n$ matrix N . The cost vector \mathbf{d} can also be split as $\mathbf{d}^T = [\mathbf{d}_B^T \quad \mathbf{d}_N^T]$ and the variable vector can be split as $\mathbf{y}^T = [\mathbf{y}_B^T \quad \mathbf{y}_N^T]$ with $\mathbf{y}_N = \mathbf{0}$. To operate with the tableau, we will split it as

$$\begin{bmatrix} B & N & -\mathbf{b} \\ \mathbf{d}_B^T & \mathbf{d}_N^T & \mathbf{0} \end{bmatrix}$$

Performing Gauss Jordan elimination on the columns corresponding to basic variables, we get the equations into the form that will be preserved across iterations.

$$\begin{bmatrix} I & B^{-1}N & -B^{-1}\mathbf{b} \\ \mathbf{d}_B^T & \mathbf{d}_N^T & \mathbf{0} \end{bmatrix}$$

Further, we will ensure that all the columns corresponding to basic variables are in the unit form.

$$\begin{bmatrix} I & B^{-1}N & -B^{-1}\mathbf{b} \\ \mathbf{d}_B^T - \mathbf{d}_B^T I = \mathbf{0} & \mathbf{d}_N^T - \mathbf{d}_B^T B^{-1}N & \mathbf{d}_B^T B^{-1}\mathbf{b} \end{bmatrix}$$

This corresponds to a solution $\mathbf{y}_B = -B^{-1}\mathbf{b}$ with cost $\mathbf{d}^T \mathbf{y} = -\mathbf{d}_B^T B^{-1}\mathbf{b}$, which is the negative of the expression on the right hand bottom corner.

3. In the above tableau, the components of the expression $\mathbf{r} = \mathbf{d}_N^T - \mathbf{d}_B^T B^{-1}N$ are the *reduced costs* and capture what it costs to use the existing set of free variables; if the direct cost in \mathbf{d}_N is less than the saving due to use of the other basic variables, it will help to try a free variable. This guides us in the choice of the *entering variable*. If $\mathbf{r} = \mathbf{d}_N^T - \mathbf{d}_B^T B^{-1}N$ has any negative component, then the variable corresponding to the most negative component is picked up as the next entering variable and this choice corresponds to moving from a corner of the polytope \mathcal{S} to an adjacent corner with lower cost. Let y_k be the entering variable and d_k the corresponding cost.
4. As the entering component y_k is increased, to maintain $M\mathbf{y} = -\mathbf{b}$, the first component \mathbf{y}_j that decreases to 0 becomes the leaving variable and transforms from a basic to a free variable. The other components of \mathbf{y}_B would have moved around but would remain positive. Thus, the one that drops to zero should satisfy

$$j = \underset{t=1,2,\dots,m}{\operatorname{argmin}} \frac{(-B^{-1}\mathbf{b})_t}{(B^{-1}N)_{tk} > 0}$$

Note that the minimum is taken only over the positive components $(B^{-1}N)_{tk}$. If there are no positive components, the next corner is infinitely far away and the cost can be reduced forever to yield a minimum cost of $-\infty$.

- With the new choice of basic variables, steps (2)-(4) are repeated till the reduced cost is completely non-negative. The variables corresponding to the unit columns in the final tableau are the basic variables at the optimum.

What we have not discussed so far is how to obtain the initial basic feasible point. If $\mathbf{x} = 0$ satisfies $A\mathbf{x} \geq \mathbf{b}$, we can have an initial basic feasible point with the basic variables comprising of \mathbf{s} and \mathbf{x} constituting the free variables. This is illustrated through the following example. Consider the problem

$$\begin{array}{ll} \min_{x_1, x_2, x_3 \in \mathfrak{R}} & -15x_1 - 18x_2 - 20x_3 \\ \text{subject to} & -\frac{1}{6}x_1 - \frac{1}{4}x_2 - \frac{1}{2}x_3 \geq -60 \\ & -40x_1 - 50x_2 - 60x_3 \geq -2880 \\ & -25x_1 - 30x_2 - 40x_3 \geq -2400 \\ & x_1 \geq 0, x_2 \geq 0, x_3 \geq 0 \end{array}$$

The initial tableau is

$$\left[\begin{array}{cccccc|c} \frac{1}{6} & \frac{1}{4} & \frac{1}{2} & 1 & 0 & 0 & 60 \\ 40 & 50 & 60 & 0 & 1 & 0 & 2880 \\ 25 & 30 & 40 & 0 & 0 & 1 & 2400 \\ -15 & -18 & -20 & 0 & 0 & 0 & 0 \end{array} \right]$$

The most negative component of the reduced cost vector is for $k = 3$. The pivot row number is $2 = \operatorname{argmin}_{t=1,2,\dots,m} \frac{(B^{-1}\mathbf{b})_t}{(B^{-1}N)_{tk} > 0}$. Thus, the leaving basic variable is s_2 (the basic variable corresponding to the second row) while the entering free variable is x_3 . Performing Gauss elimination to obtain column $k = 3$ in the unit form, we get

$$\left[\begin{array}{cccccc|c} -\frac{1}{6} & -\frac{1}{6} & 0 & 1 & -\frac{1}{120} & 0 & 36 \\ \frac{2}{3} & \frac{5}{6} & 1 & 0 & \frac{1}{60} & 0 & 48 \\ -\frac{5}{3} & -\frac{10}{3} & 0 & 0 & -\frac{2}{3} & 1 & 480 \\ -\frac{5}{3} & -\frac{4}{3} & 0 & 0 & \frac{1}{3} & 0 & 960 \end{array} \right]$$

This tableau corresponds to the solution $x_1 = 0, x_2 = 0, x_3 = 48, s_1 = 0, s_2 = 36, s_3 = 480$ and cost $\mathbf{c}^T \mathbf{x} = -960$ (negative of the number on the right hand bottom corner). Since the reduced cost vector has still some negative components, it is possible to find a basic feasible solution with lower cost. Using the most negative component of the reduced cost vector, we select the next pivot

element to be $m_{21} = \frac{2}{3}$. Again performing Gaussian elimination, we obtain the tableau corresponding to the next iterate.

$$\left[\begin{array}{cccccc|c} 0 & \frac{1}{24} & \frac{1}{4} & 1 & -\frac{1}{240} & 0 & 48 \\ 1 & \frac{5}{4} & \frac{3}{2} & 0 & \frac{1}{40} & 0 & 72 \\ 0 & -\frac{5}{4} & \frac{5}{2} & 0 & -\frac{5}{8} & 1 & 600 \\ 0 & \frac{3}{4} & \frac{5}{2} & 0 & \frac{3}{8} & 0 & 1080 \end{array} \right]$$

Note that the optimal solution has been found, since the reduced cost vector is non-negative. The optimal solution is $x_1 = 72, x_2 = 0, x_3 = 0, s_1 = 48, s_2 = 0, s_3 = 600$ and cost $\mathbf{c}^T \mathbf{x} = -1080$

What if $\mathbf{x} = \mathbf{0}$ does not satisfy $A\mathbf{x} \geq \mathbf{b}$? The choice of \mathbf{s} as the basic variables and \mathbf{x} as the free variables will not be valid. As an example, consider the problem

$$\begin{array}{ll} \min_{x_1, x_2, x_3 \in \mathbb{R}} & 30x_1 + 60x_2 + 70x_3 \\ \text{subject to} & x_1 + 3x_2 + 4x_3 \geq 14 \\ & 2x_1 + 2x_2 + 3x_3 \geq 16 \\ & x_1 + 3x_2 + 2x_3 \geq 12 \\ & x_1 \geq 0, x_2 \geq 0, x_3 \geq 0 \end{array}$$

The initial tableau is

$$\left[\begin{array}{cccccc|c} -1 & -3 & -4 & 1 & 0 & 0 & -14 \\ -2 & -2 & -3 & 0 & 1 & 0 & -16 \\ -1 & -3 & -2 & 0 & 0 & 1 & -12 \\ 30 & 60 & 70 & 0 & 0 & 0 & 0 \end{array} \right]$$

With the choice of basic and free variables as above, we are not even in the feasible region to start off with. In general, if we have any negative number in the last column of the tableau, $\mathbf{x} = \mathbf{0}$ is not in the feasible region. Further, we have no negative numbers in the bottom row, which does not leave us with any choice of cost reducing free variable. But this is not of primary concern, since we first need to maneuver our way into the feasible region. We do this by moving from one basic point (that is, a point having not more than n zero components) to another till we land in the feasible region, which is indicated by all positive components in the extreme right hand column. This movement from one basic point to another is not driven by negative components in the cost vector, but rather by the negative components in the right hand column. The new rules for moving from one basic point to another are:

1. Pick³⁹ any negative number in the far right column (excluding the last row). Let this be in the q^{th} row for $q < m + 1$.

³⁹Note that there is no priority here.

2. In the q^{th} row, move⁴⁰ to left to a column number k where there is another negative number. The variable y_k will be the next entering variable.
3. Choose pivot element m_{jk} which gives the smallest positive ratio of an element in the j^{th} row of the last column to the element m_{jk} . The leaving variable will be y_j .
4. Once the pivot element is chosen, proceed as usual to convert the pivot element to 1 and the other elements in the pivot column to 0.
5. Repeat steps (1)-(4) on the modified tableau until there is no negative element in the right-most column.

Applying this procedure to the tableau above, we pick $m_{2,1} = -2$ as our first pivot element and do row elimination to get the first column in unit form.

$$\left[\begin{array}{cccccc|c} 0 & -2 & -\frac{5}{2} & 1 & -\frac{1}{2} & 0 & -6 \\ 1 & 1 & \frac{3}{2} & 0 & -\frac{1}{2} & 0 & 8 \\ 0 & -2 & -\frac{1}{2} & 0 & -\frac{1}{2} & 1 & -4 \\ 0 & 30 & 25 & 0 & 15 & 0 & -240 \end{array} \right]$$

We pick $m_{35} = -\frac{1}{2}$ as our next pivot element and do similar row elimination operations to obtain

$$\left[\begin{array}{cccccc|c} 0 & 0 & -2 & 1 & 0 & -1 & -2 \\ 1 & 3 & 2 & 0 & 0 & -1 & 12 \\ 0 & 4 & 1 & 0 & 1 & -2 & 8 \\ 0 & -30 & 10 & 0 & 0 & 30 & -360 \end{array} \right]$$

We have still not obtained a feasible basic point. We choose $m_{16} = -1$ as the next pivot and do row eliminations to get the next tableau.

$$\left[\begin{array}{cccccc|c} 0 & 0 & 2 & -1 & 0 & 1 & 2 \\ 1 & 3 & 4 & -1 & 0 & 0 & 14 \\ 0 & 4 & 5 & -2 & 1 & 0 & 12 \\ 0 & -30 & -50 & 30 & 0 & 0 & -420 \end{array} \right]$$

This tableau has not negative numbers in the last column and gives a basic feasible point $x_1 = 14, x_2 = 0, x_3 = 0$. Once we obtain the basic feasible point, we revert to the standard simplex procedure discussed earlier. The most negative component of the reduced cost vector is -50 and this leads to the pivot element $m_{13} = 2$. Row elimination yields

$$\left[\begin{array}{cccccc|c} 0 & 0 & 1 & -\frac{1}{2} & 0 & \frac{1}{2} & 1 \\ 1 & 3 & 0 & 1 & 0 & -2 & 10 \\ 0 & 4 & 0 & \frac{1}{2} & 1 & -\frac{5}{2} & 7 \\ 0 & -30 & 0 & 5 & 0 & 25 & -370 \end{array} \right]$$

⁴⁰Note that there is no priority here either.

Our next pivot element is $m_{32} = 4$. Row elimination yields

$$\left[\begin{array}{cccccc|c} 0 & 0 & 1 & -\frac{1}{2} & 0 & \frac{1}{2} & 1 \\ 1 & 0 & 0 & \frac{5}{8} & -\frac{3}{4} & -\frac{1}{8} & \frac{19}{4} \\ 0 & 1 & 0 & \frac{1}{8} & \frac{1}{4} & -\frac{5}{8} & \frac{7}{4} \\ 0 & 0 & 0 & \frac{35}{4} & \frac{15}{2} & \frac{25}{4} & -\frac{635}{2} \end{array} \right]$$

We are done! The reduced cost vector has no more negative components. The optimal basic feasible point is $x_1 = 4.75, x_2 = 1.75, x_3 = 1$ and the optimal cost is 317.5.

Revised Simplex Method

The simplex method illustrated above serves two purposes:

1. Doing all the eliminations completely makes the idea clear.
2. It is easier to follow the process when working out the solution by hand.

For computational purposes however, it is uncommon now to use the method as described earlier. This is because, once \mathbf{r} is computed, none of the columns above \mathbf{r} , (except for that corresponding to the leaving variable) are used. Therefore, computing them is a useless effort. Doing the eliminations completely at each step cannot be justified practically. Instead, the more efficient version of the simplex method, as outlined below, is used by software packages. It is called the revised simplex method and is essentially the simplex method itself, boiled down.

Compute the reduced costs $\mathbf{r} = \mathbf{d}_N - (\mathbf{d}_B)B^{-1}N$.
while $\mathbf{r} \not\geq \mathbf{0}$ **do**
 1. Let r_k be the most negative component of \mathbf{r} .
 2. Compute $\mathbf{v} = B^{-1}\mathbf{n}_i$, where \mathbf{n}_i is the i^{th} column of N .
 3. Let $j = \underset{t=1,2,\dots,m}{\operatorname{argmin}} \frac{(-B^{-1}\mathbf{b})_t}{(B^{-1}\mathbf{n}_i)_t}$.
 4. Update B (or B^{-1}) and $\mathbf{x}_B = B^{-1}\mathbf{b}$ to reflect the j^{th} leaving column and the k^{th} entering variable.
 Compute the new reduced costs $\mathbf{r} = \mathbf{d}_N - (\mathbf{d}_B)B^{-1}N$.
end while

Figure 4.58: The revised simplex method.

4.7.2 Interior point barrier method

Researchers have dreamt up pathological examples for the simplex method, for which the simplex method takes an exponential amount of time. In practice, however, the simplex method is one of the most efficient methods for a majority

of the LPs. Application of interior point methods to LP have led to a new competitor to the simplex method in the form of interior point methods for linear programming. In contrast to the simplex algorithm, which finds the optimal solution by progressing along points on the boundary of a polyhedral set, interior point methods traverse to the optimal through the interior of the feasible region (polyhedral set in the case of LPs). The first in this league was the iterative Karmarkar's method [?], developed by Narendra Karmarkar in 1984. Karmarkar also proved that the algorithm was polynomial time. This line of research was inspired by the *ellipsoid method* for linear programming, outlined by Leonid Khachiyan in 1979; the ellipsoid algorithm itself was introduced by Naum Z. Shor, *et. al.* in 1972 and used by Leonid Khachiyan [?] to prove the polynomial-time solvability of linear programs in linear programming, which was the first such algorithm known to have a polynomial running time.

This competitor to the simplex method takes a Newton's method-like approach through the interior of the feasible region. Newton steps are taken till the 'barrier' is encountered. It stirred up the world of optimization and inspired the whole class of barrier methods. Following this, a lot of interest was generated in the application of the erstwhile interior point methods for general non-linear constrained optimization problems. The Karmarkar's algorithm is now replaced by an improved logarithmic barrier method that makes use of the primal as well as the dual for solving an LP. Shanno and Bagchi [?] showed that Karmarkar's method is just a special case of the logarithmic barrier function method. We will restrict our discussion to the primal-dual barrier method [?].

The dual (4.113) for the linear program (4.111) can be derived in manner similar to the dual on page 287.

$$\begin{aligned} & \max_{\lambda \in \mathfrak{R}^m} && \lambda^T \mathbf{b} \\ & \text{subject to} && A^T \lambda \leq \mathbf{c} \\ & && \lambda \geq \mathbf{0} \end{aligned}$$

The weak duality theorem (theorem 81) states that the objective function value of the dual at any feasible solution is always less than or equal to the objective function value of the primal at any feasible solution. That is, for any primal feasible \mathbf{x} and any dual feasible λ ,

$$\mathbf{c}^T \mathbf{x} - \mathbf{b}^T \lambda \geq 0$$

For this specific case, the weak duality is easy to see: $\mathbf{b}^T \lambda \leq \mathbf{x}^T A^T \lambda \leq \mathbf{x}^T \mathbf{c}$.

Further, it can be proved using the *Farkas' lemma* that if the primal has an optimal solution \mathbf{x}^* (which is assumed to be bounded), then the dual also has an optimal solution⁴¹ λ^* , such that $\mathbf{c}^T \mathbf{x}^* = \mathbf{b}^T \lambda^*$.

The following steps will set the platform for the interior point method.

⁴¹For an LP and its dual D, there are only four possibilities:

1. (LP) is bounded and feasible and (D) is bounded and feasible.

1. As on page 4.7.1, we will rewrite the constraints $A\mathbf{x} \geq \mathbf{b}$ in the above LP as equations, by introducing a new non-negative "slack" variable s_j for the j^{th} constraint (for all j 's) and subtracting it from the left-hand side of each inequality. This gives us the constraint $M\mathbf{y} = -\mathbf{b}$, where $M = [-A \ I]$ and $\mathbf{y} = [\mathbf{x} \ \mathbf{s}]^T$. As before, the original cost vector \mathbf{c} is extended to a vector \mathbf{d} by appending m more zero components. This gives us the following problem, equivalent to the original LP (??).

$$\begin{aligned} \min_{\mathbf{y} \in \mathbb{R}^{n+m}} \quad & \mathbf{y}^T \mathbf{d} \\ \text{subject to} \quad & M\mathbf{y} = -\mathbf{b} \\ & \mathbf{y} \geq \mathbf{0} \end{aligned} \tag{4.113}$$

Its dual problem is given by

$$\begin{aligned} \max_{\lambda \in \mathbb{R}^m} \quad & -\lambda^T \mathbf{b} \\ \text{subject to} \quad & M^T \lambda \leq \mathbf{d} \end{aligned} \tag{4.114}$$

2. Next, we set up the barrier method formulation of the dual of the linear program. Letting $\mu > 0$ be a given fixed parameter, which is decreased during the course of the algorithm. We also insert slack variables $\xi = [\xi_1, \xi_2, \dots, \xi_n]^T \geq \mathbf{0}$. The barrier method formulation of the dual is then given by:

$$\begin{aligned} \max_{\lambda \in \mathbb{R}^m} \quad & -\lambda^T \mathbf{b} + \mu \sum_{i=1}^n \ln(\xi_i) \\ \text{subject to} \quad & M^T \lambda + \xi = \mathbf{d} \end{aligned} \tag{4.115}$$

The conditions $\xi_i \geq 0$ are no longer needed since $\ln(\xi_i) \rightarrow \infty$ as $\xi_i \rightarrow 0$, if $\xi_i > 0$. This latter property means that $\ln(\xi_i)$ serves as a barrier, discouraging ξ_i from going to 0.

3. To write the first-order necessary conditions for a minimum, we set the partial derivatives of the Lagrangian

$$L(\mathbf{y}, \lambda, \xi) = -\lambda^T \mathbf{b} + \mu \sum_{i=1}^n \ln(\xi_i) - \mathbf{y}^T (M^T \lambda + \xi - \mathbf{d})$$

-
2. (LP) is infeasible and (D) is unbounded and feasible.
 3. (LP) is unbounded and feasible and (D) is infeasible.
 4. (LP) is infeasible and (D) is infeasible.

This can be proved using the *Farkas' Lemma*.

with respect to \mathbf{y} , λ , ξ to zero. This results in the set of following three equations:

$$\begin{aligned} M^T \lambda + \xi &= \mathbf{d} \\ M\mathbf{y} &= -\mathbf{b} \\ \text{diag}(\xi) \text{diag}(\mathbf{y})\mathbf{1} &= \mu\mathbf{1} \end{aligned} \quad (4.116)$$

which include the dual and primal feasibility conditions excluding $\mathbf{y} \geq \mathbf{0}$ and $\xi \geq \mathbf{0}$.

4. We will assume that our current point $\mathbf{y}^{(k)}$ is primal feasible and the current point $(\lambda^{(k)}, \xi^{(k)})$ is dual feasible. We determine a new search direction $(\Delta\mathbf{y}^{(k)}, \Delta\lambda^{(k)}, \Delta\xi^{(k)})$ so that the new point $(\mathbf{y}^{(k)} + \Delta\mathbf{y}^{(k)}, \lambda^{(k)} + \Delta\lambda^{(k)}, \xi^{(k)} + \Delta\xi^{(k)})$ satisfies (4.116). This gives us the so-called Newton equations:

$$\begin{aligned} M^T \Delta\lambda + \Delta\xi &= \mathbf{0} \\ M\Delta\mathbf{y} &= \mathbf{0} \\ (y_i + \Delta y_i)(\xi_i + \Delta\xi_i) &= \mu \quad i = 1, 2, \dots, n \end{aligned} \quad (4.117)$$

Ignoring the second order term $\Delta y_i \Delta \xi_i$ in the third equation, and solving the system of equations in (4.117), we get the following update rules:

$$\begin{aligned} \Delta\lambda^{(k)} &= - (M \text{diag}(\mathbf{y}^{(k)}) \text{diag}(\xi^{(k)})^{-1} M^T)^{-1} M \text{diag}(\xi^{(k)})^{-1} (\mu\mathbf{1} - \text{diag}(\mathbf{y}^{(k)}) \text{diag}(\xi^{(k)})\mathbf{1}) \\ \Delta\xi^{(k)} &= -M^T \Delta\lambda^{(k)} \\ \Delta\mathbf{y}^{(k)} &= \text{diag}(\xi^{(k)})^{-1} (\mu\mathbf{1} - \text{diag}(\mathbf{y}^{(k)}) \text{diag}(\xi^{(k)})\mathbf{1}) - \text{diag}(\mathbf{y}^{(k)}) \text{diag}(\xi^{(k)})^{-1} \Delta\xi^{(k)} \end{aligned} \quad (4.118)$$

An affine variant of the algorithm can be developed by setting $\mu = 0$ in the equations (4.118).

5. $\Delta\mathbf{y}^{(k)}$ and $(\Delta\lambda^{(k)}, \Delta\xi^{(k)})$ correspond to partially constrained Newton steps, which might not honour the constraints $\mathbf{y}^{(k)} + \Delta\mathbf{y}^{(k)} \geq \mathbf{0}$ and $\xi^{(k)} + \Delta\xi^{(k)} \geq \mathbf{0}$. Since we have a separate search direction $\Delta\mathbf{y}^{(k)}$ in the primal space and a separate search direction $(\Delta\lambda^{(k)}, \Delta\xi^{(k)})$ in the dual space, we could compute the maximum step length $t_{(max,P)}^{(k)}$ that maintains the pending primal inequality $\mathbf{y}^{(k)} + t_{(max,P)}^{(k)} \Delta\mathbf{y}^{(k)} \geq \mathbf{0}$ and the maximum step length $t_{(max,D)}^{(k)}$ that maintains the pending dual inequality $\xi^{(k)} + t_{(max,D)}^{(k)} \Delta\xi^{(k)} \geq \mathbf{0}$.
6. Now we have a feasible primal solution $\mathbf{y}^{(k+1)}$ and feasible dual solution $(\lambda^{(k+1)}, \xi^{(k+1)})$ given by

$$\begin{aligned}
\mathbf{y}^{(k+1)} &= \mathbf{y}^{(k)} + t_{(max,P)}^{(k)} \Delta \mathbf{y}^{(k)} \\
\lambda^{(k+1)} &= \lambda^{(k)} + \Delta \lambda^{(k)} \\
\xi^{(k+1)} &= \xi^{(k)} + t_{(max,D)}^{(k)} \Delta \xi^{(k)}
\end{aligned} \tag{4.119}$$

7. For user specified small thresholds of $\epsilon_1 > 0$ and $\epsilon_2 > 0$, if the duality gap $\mathbf{d}^T \mathbf{y}^{(k+1)} + \mathbf{b}^T \lambda^{(k+1)}$ is not sufficiently close to 0, *i.e.*,

$$\mathbf{d}^T \mathbf{y}^{(k+1)} + \mathbf{b}^T \lambda^{(k+1)} > \epsilon_1$$

for a μ not yet sufficiently close to 0, *i.e.*,

$$\mu > \epsilon_2$$

we decrease μ by a user specified factor $\rho < 1$ (such as $\rho = 0.1$).

$$\mu = \mu \times \rho$$

8. Set $k = k + 1$. If μ was not modified in step (7), the duality gap is sufficiently small and the termination condition has been reached. So EXIT. Else, the last condition in (4.116) no longer holds with the modified value of μ . So steps (4)-(7) are re-executed.

In practice, the interior point method for LP gets down the duality gap to within 10^{-8} in just 20-80 steps (which is still slower than the simplex method for many problems), independent of the size of the problem specified through values of m and n .

4.8 Least Squares

Least squares was motivated in Section 3.9.2, based on the idea of projection. Least squares problems appear very frequently in practice. The objective for minimization in the case of least squares is the square of the euclidian norm of $A\mathbf{x} - \mathbf{b}$, where A is a $m \times n$ matrix, \mathbf{x} is a vector of n variables and \mathbf{b} is a vector of m knowns.

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|A\mathbf{x} - \mathbf{b}\|_2^2 \tag{4.120}$$

Very often one has a system of linear constraints on problem (4.120).

$$\begin{aligned}
&\min_{\mathbf{x} \in \mathbb{R}^n} \|A\mathbf{x} - \mathbf{b}\|_2^2 \\
&\text{subject to } C^T \mathbf{x} = \mathbf{0}
\end{aligned} \tag{4.121}$$

This problem is called the *least squares problem with linear constraints*.

In practice, incorporating the constraints $C^T \mathbf{x} = \mathbf{0}$ properly makes quite a difference. In lots of regularization problems, the least squares problem often comes with quadratic constraints in the following form.

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^n} \quad & \|A\mathbf{x} - \mathbf{b}\|_2^2 \\ \text{subject to} \quad & \|\mathbf{x}\|_2^2 = \alpha^2 \end{aligned} \quad (4.122)$$

This problem is termed as the *least squares problem with quadratic constraints*.

The classical statistical model assumes that all the error occurs in the vector \mathbf{b} . But sometimes, the data matrix A is itself not very well known, owing to errors in the variables. This is the model we have in the simplest version of the *total least squares problem*, which is stated as follows.

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^n, E \in \mathbb{R}^{m \times n}, \mathbf{r} \in \mathbb{R}^m} \quad & \|E\|_F^2 + \|\mathbf{r}\|_2^2 \\ \text{subject to} \quad & (A + E)\mathbf{x} = \mathbf{b} + \mathbf{r} \end{aligned} \quad (4.123)$$

While there is always a solution to the least squares problem (4.120), there is not always a solution to the total least squares problem (4.133). Finally, you can have a combination of linear and quadratic constraints in a least squares problem to yield a *least squares problem with linear and quadratic constraints*.

We will briefly discuss the problem of solving linear least squares problems and total least squares problems with linear or a quadratic constraint (due to regularization) The importance of lagrange multipliers will be introduced in the process. We will discuss stable numerical methods when the data matrix A is singular or near singular. We will also present iterative methods for large and sparse data matrices. There are many applications of least squares problems, which include statistical methods, image processing, data interpolation and surface fitting and finally geometrical problems.

4.8.1 Linear Least Squares

As a user of least squares in practice, one of the most important things to be known is that when A is of full column rank, it has an analytical solution given by \mathbf{x}^* (which was derived in Section 3.9.2 and gives a dual interpretation).

$$\text{Analytical solution: } \mathbf{x}^* = (A^T A)^{-1} A^T \mathbf{b} \quad (4.124)$$

This analytic solution can also be obtained by observing that

1. $\|\mathbf{y}\|_2^2$ is a convex function for $\mathbf{y} \in \mathbb{R}^m$.

2. Square of the convex euclidian norm function, applied to an affine transform is also convex. Thus $\|A\mathbf{x} - \mathbf{b}\|_2^2$ is convex.
3. Every critical point of a convex function defined on an open domain corresponds to its local minimum. The critical point \mathbf{x}^* of $\|A\mathbf{x} - \mathbf{b}\|_2^2$ should satisfy

$$\nabla(A\mathbf{x} - \mathbf{b})^T(A\mathbf{x} - \mathbf{b}) = 2A^T A\mathbf{x}^* - 2A^T \mathbf{b} = \mathbf{0}$$

Thus,

$$\mathbf{x}^* = (A^T A)^{-1} A^T \mathbf{b}$$

corresponds to a point of local minimum of (4.120) if $A^T A$ is invertible.

This is the classical way statisticians solve least squares problem. It can be solved very efficiently, and there exist many softwares that implement this solution. The computation time is linear in the number of rows of A and quadratic in the number of columns. For extremely large A , it can become important to look at the structure of A to solve it efficiently, but for most problems, it is efficient. In practice least-squares is very easy to recognize as an objective function. There are a few standard tricks to increase the flexibility. For example, constraints can be handled to a certain extent by adding weights. When the matrix A is not full column rank, the solution to (4.120) may not be unique.

We should note that while we get a closed form solution to the problem of minimizing the square of the euclidian norm, it is not so for most other norms such as the infinity norm. However, there exist iterative methods for solving least squares with infinity norm that yield a solution in as much time as is taken in computing the solution using the analytical formula in 4.124. Therefore, having a closed form solution is not always computationally helpful. In general, the method of solution to a least squares problem depends on the sparsity as well as the size of A and the degree of accuracy desired.

In practice, however, it is not recommended to solve least squares problem using the classical equation in 4.124 since the method is numerically unstable. Numerical linear algebra instead recommends the QR decomposition to accurately solve the least squares problem. This method is slower, but more numerically stable than the classical method. In theorem ??, we state a theory that compares the analytical solution (4.124) and the QR approach to the least squares problem.

Let A be an $m \times n$ matrix of either full row or full column rank. For the case of $n > m$, we saw on page ?? (summarised in Figure 3.3) that the system $A\mathbf{x} = \mathbf{b}$ will have at least one solution which means that minimum value of the objective function will be 0, corresponding to the solution. We are interested in the case $m \geq n$, for which there will either be no solution or a single solution to $A\mathbf{x} = \mathbf{b}$ and we are interested in one that minimizes $\|A\mathbf{x} - \mathbf{b}\|_2^2$.

1. We first decompose A into the product of an orthonormal $m \times m$ matrix Q with an upper triangular $m \times n$ matrix R , using the gram-schmidt or-

thonormalization process⁴² discussed in Section 3.9.4. The decomposition can also be performed using the Householder⁴³ transformation or Givens rotation. Householder transformation has the added advantage that new rows or columns can be introduced without requiring a complete redo of the decomposition process. The last $m - n$ rows of R will be zero rows. Since $Q^{-1} = Q^T$, the QR decomposition yields the system

$$Q^T A = \begin{bmatrix} R_1 \\ \mathbf{0} \end{bmatrix}$$

2. Applying the same orthogonal matrix to \mathbf{b} , we get

$$Q^T \mathbf{b} = \begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix}$$

where $\mathbf{d} \in \Re^{m-n}$.

3. The solution to the least squares problem is found by solving $R_1 \mathbf{x} = \mathbf{c}$. The solution to this can be found by simple back-substitution.

The next theorem examines how the least squares solution and its residual $\|A\mathbf{x} - \mathbf{b}\|$ are affected by changes in A and \mathbf{b} . Before stating the theorem, we will introduce the concept of the condition number.

Condition Number

The *condition number* associated with a problem is a measure of how numerically well-posed the problem is. A problem with a low condition number is said to be well-conditioned, while a problem with a high condition number is said to be ill-conditioned. For a linear system $A\mathbf{x} = \mathbf{b}$, the condition number is defined as maximum ratio of the relative error in \mathbf{x} (measured using any particular norm) divided by the relative error in \mathbf{b} . It can be proved (using the Cauchy Schwarz inequality) that the condition number equals $\|A^{-1}A\|$ and is independent of \mathbf{b} . It is denoted by $\kappa(A)$ and is also called the condition number of the matrix A .

$$\kappa(A) = \|A^{-1}A\|$$

If $\|\cdot\|_2$ is the L_2 norm, then

$$\kappa(A) = \frac{\sigma_{max}(A)}{\sigma_{min}(A)} = \|A\|_2 \|(A^T A)^{-1} A^T\|_2$$

⁴²The classical Gram-Schmidt method is often numerically unstable. Golub [?] suggests a modified Gram-Schmidt method that is numerically stable.

⁴³Householder was a numerical analyst. However, the first mention of the Householder transformation dates back to the 1930s in a book by Aikins, a statistician and a numerical analyst.

where $\sigma_{max}(A)$ and $\sigma_{min}(A)$ are maximal and minimal singular values of A respectively. For a real square matrix A , the square roots of the eigenvalues of $A^T A$, are called singular values. Further,

$$\kappa(A)^2 = \|A\|_2^2 \|(A^T A)^{-1}\|_2^2$$

Theorem 86 By $\|\cdot\|$, we will refer to the L_2 norm. Let

$$\mathbf{x}^* = \operatorname{argmin} \|A\mathbf{x} - \mathbf{b}\|$$

$$\hat{\mathbf{x}} = \operatorname{argmin} \|(A + \delta A)\mathbf{x} - (\mathbf{b} + \delta \mathbf{b})\|$$

where A and δA are in $\mathfrak{R}^{m \times n}$ with $m \geq n$. Let \mathbf{b} and $\delta \mathbf{b}$ be in \mathfrak{R}^m with $\mathbf{b} \neq \mathbf{0}$. Let us set

$$\mathbf{r}^* = \mathbf{b} - A\mathbf{x}^*$$

$$\hat{\mathbf{r}} = \mathbf{b} - A\hat{\mathbf{x}}$$

and

$$\rho^* = \|A\mathbf{x}^* - \mathbf{b}\|$$

If

$$\epsilon = \max \left\{ \frac{\|\delta A\|}{\|A\|}, \frac{\|\delta \mathbf{b}\|}{\|\mathbf{b}\|} \right\} < \frac{\sigma_n(A)}{\sigma_1(A)}$$

and

$$\sin \theta = \frac{\rho^*}{\|\mathbf{b}\|} \neq 1$$

then,

$$\frac{\|\hat{\mathbf{x}} - \mathbf{x}^*\|}{\|\mathbf{x}^*\|} \leq \epsilon \left\{ \frac{2\kappa(A)}{\cos \theta} + \tan \theta \kappa(A)^2 \right\} + O(\epsilon^2)$$

In this inequality most critical term for our discussion is $\kappa(A)^2$ and this is the term that can kill the analytical solution to least squares. Now matter how accurate an algorithm you use, you still have $\kappa(A)^2$, provided $\tan \theta$ is non-zero. Now $\tan \theta$ does not appear if you are solving a linear system, but if you solve a least squares problem this term appears, bringing along $\kappa(A)^2$. Thus, solving least squares problem is inherently more difficult and sensitive than linear equations. The perturbation theory for the residual vector depends just on the condition number $\kappa(A)$ (and not its square):

$$\frac{\|\hat{\mathbf{r}} - \mathbf{r}^*\|}{\|\mathbf{b}\|} \leq \epsilon \{1 + 2\kappa(A)\} \min\{1, m - n\} + O(\epsilon^2) + O(\epsilon^2)$$

However, having a small residual does not necessarily imply that you will have a good approximate solution.

The theorem implies that the sensitivity of the analytical solution \mathbf{x}^* for non-zero residual problems is measured by the square of the condition number. Whereas, sensitivity of the residual depends just linearly on $\kappa(A)$. We note that the QR method actually solves a nearby least squares problem.

Linear Least Squares for Singular Systems

To solve the linear least squares problem (4.120) for a matrix A that is of rank $r < \min\{m, n\}$, we can compute the pseudo-inverse (*c.f.* page 208) A^+ and obtain the least squares solution⁴⁴ as

$$\hat{\mathbf{x}} = A^+ \mathbf{b}$$

A^+ can be computed by first computing a singular orthogonal factorization

$$A = Q \begin{bmatrix} R & 0 \\ 0 & 0 \end{bmatrix} Z^T$$

where $Q^T Q = I_{m \times m}$ and $Z^T Z = I_{n \times n}$ and R is an $r \times r$ upper triangular matrix. A^+ can be computed in a straightforward manner as

$$A^+ = Z \begin{bmatrix} R^{-1} & 0 \\ 0 & 0 \end{bmatrix} Q^T$$

The above least squares solution can be justified as follows. Let

$$Q^T \mathbf{b} = \begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix}$$

and

$$Z^T \mathbf{x} = \begin{bmatrix} \mathbf{w} \\ \mathbf{y} \end{bmatrix}$$

Then

$$\|A\mathbf{x} - \mathbf{b}\|^2 = \|Q^T A Z Z^T \mathbf{x} - Q^T \mathbf{b}\|^2 = \|R\mathbf{w} - \mathbf{c}\|^2 + \|\mathbf{d}\|^2$$

The least squares solution is therefore given by

$$\hat{\mathbf{x}} = Z \begin{bmatrix} R^{-1} \mathbf{c} \\ 0 \end{bmatrix}$$

One particular decomposition that can be used is the singular value decomposition (*c.f.* Section 3.13) of A , with $Q^T \equiv U^T$ and $Z \equiv V$ and $U^T A V = \Sigma$. The pseudo-inverse A^+ has the following expression.

$$A^+ = V \Sigma^{-1} U^T$$

It can be shown that this A^+ is the unique minimal Frobenius norm solution to the following problem.

$$A^+ = \operatorname{argmin}_{X \in \mathbb{R}^{n \times m}} \|AX - I_{m \times m}\|$$

⁴⁴Note that this solution not only minimizes $\|A\mathbf{x} - \mathbf{b}\|$ but also minimizes $\|\mathbf{x}\|$. This may or may not be desirable.

This also shows that singular value decomposition can be looked upon as an optimization problem.

A greater problem is with systems that are nearly singular. Numerically and computationally it seldom happens that the rank of matrix is exactly r . A classical example is the following $n \times n$ matrix K , which has a determinant of 1.

$$K = \begin{bmatrix} 1 & -1 & \dots & -1 & -1 & \dots & -1 \\ 0 & 1 & \dots & -1 & -1 & \dots & -1 \\ \cdot & \cdot & \dots & \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \dots & \cdot & \cdot & \dots & \cdot \\ 0 & 0 & \dots & 1 & -1 & \dots & -1 \\ 0 & 0 & \dots & 0 & 1 & \dots & -1 \\ \cdot & \cdot & \dots & \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \dots & \cdot & \cdot & \dots & \cdot \\ 0 & 0 & \dots & 0 & 0 & \dots & 1 \end{bmatrix}$$

The eigenvalues of this matrix are also equal to 1, while its rank is n . However, a very small perturbation to this matrix can reduce its rank to $n - 1$; the rank of $K - 2^{-(n-1)}I_{n \times n}$ is $n - 1$! Such catastrophic problems occur very often when you do large computations. The solution using SVD is applicable for nearly singular systems as well.

4.8.2 Least Squares with Linear Constraints

We first reproduce the least squares problem with linear constraints that was stated earlier in (4.121).

$$\begin{aligned} \min_{\mathbf{x} \in \mathfrak{R}^n} \quad & \|A\mathbf{x} - \mathbf{b}\|^2 \\ \text{subject to} \quad & C^T\mathbf{x} = \mathbf{0} \end{aligned}$$

Let $C \in \mathfrak{R}^{n \times p}$, $A \in \mathfrak{R}^{m \times n}$ and $\mathbf{b} \in \mathfrak{R}^m$. We note that $\|A\mathbf{x} - \mathbf{b}\|^2$ is a convex function (since L_2 norm is convex and this function is the L_2 norm applied to an affine transform). We can thus solve this constrained problem by invoking the necessary and sufficient KKT conditions discussed in Section 4.4.4. The conditions can be worked out to yield

$$\begin{bmatrix} A^T A & C \\ C^T & 0 \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \lambda \end{bmatrix} = \begin{bmatrix} A^T \mathbf{b} \\ \mathbf{0} \end{bmatrix}$$

We need to now solve not only for the unknowns \mathbf{x} , but also for the lagrange multipliers; we have increased the dimensionality of the problem to $n + p$. If $\hat{\mathbf{x}} =$

$(A^T A)^{-1} A^T \mathbf{b}$ denotes the solution of the unconstrained least squares problem, then, using the first system of equality above, \mathbf{x} can be expressed as

$$\mathbf{x} = \hat{\mathbf{x}} - (A^T A)^{-1} C \lambda \quad (4.125)$$

In conjunction with the second system, this leads to

$$C^T (A^T A)^{-1} C \lambda = C^T \hat{\mathbf{x}} \quad (4.126)$$

The unconstrained least squares solution can be obtained using methods in Section 4.8.1. Next, the value of λ can be obtained by solving (4.126). If A is singular or nearly singular, we can use the singular value decomposition (or a similar decomposition) of A to determine $\hat{\mathbf{x}}$.

$$C^T R^{-1} (R^T)^{-1} C \lambda = C^T \hat{\mathbf{x}}$$

The QR factorization of $(R^T)^{-1} C$ can be efficiently used to determine λ . Finally, the value of λ can be substituted in (4.125) to solve for \mathbf{x} . This technique yields both the solutions, provided that both exist.

Another trick that is often employed when $A^T A$ is singular or nearly singular is to decrease its condition number by augmenting it in (4.125) with the ‘harmless’ $CW C^T$ and solve

$$\mathbf{x} = \hat{\mathbf{x}} - (A^T A + C W C^T)^{-1} C \lambda$$

The addition of $CW C^T$ is considered harmless, since $C^T \mathbf{x} = 0$ is to be imposed anyways. Matrix W can be chosen to be an identical or nearly identical matrix that chooses a few columns of C , just to make $A^T A + C W C^T$ non-singular.

If we use the following notation:

$$\mathcal{A}(W) = \begin{bmatrix} A^T A + C W C^T & C \\ C^T & 0 \end{bmatrix}$$

and

$$\mathcal{A} = \mathcal{A}(0) = \begin{bmatrix} A^T A & C \\ C^T & 0 \end{bmatrix}$$

and if \mathcal{A} and $\mathcal{A}(W)$ are invertible for $W \neq 0$, it can be proved that

$$\mathcal{A}^{-1}(W) = \mathcal{A}^{-1} - \begin{bmatrix} 0 & 0 \\ 0 & W \end{bmatrix}$$

Consequently

$$\kappa(\mathcal{A}(W)) \leq \kappa(\mathcal{A}) + \|W\|^2 \|C\|^2 + \alpha \|W\|$$

for some $\alpha > 0$. That is, the condition number of $\mathcal{A}(W)$ is bounded by the condition number of \mathcal{A} and some positive terms.

Another useful technique is to find an approximation to (4.121) by solving the following weighted unconstrained minimization problem.

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2 + \mu^2 \|C^T \mathbf{x}\|^2$$

For large values of μ , the solution $\widehat{\mathbf{x}}(\mu)$ of the unconstrained problem is a good approximation to the solution $\widehat{\mathbf{x}}$ of the constrained problem (4.121). We can use the generalized singular value decompositions of matrices A and C^T , that allows us to simultaneously diagonalize A and C^T .

$$U^T A X = \mathbf{diag}(\alpha_1, \dots, \alpha_m)$$

$$V^T C^T X = \mathbf{diag}(\gamma_1, \dots, \gamma_m)$$

where U and V are orthogonal matrices and X is some general matrix. The solution to the constrained problem can be expressed as

$$\widehat{\mathbf{x}} = \sum_{i=1}^p \frac{\mathbf{u}_i^T \mathbf{b}}{\alpha_i} \mathbf{x}_i$$

The analytical solution $\widehat{\mathbf{x}}(\mu)$ is then given as

$$\widehat{\mathbf{x}}(\mu) = \sum_{i=1}^p \frac{\alpha_i \mathbf{u}_i^T \mathbf{b}}{\alpha_i^2 + \mu^2 \gamma_i^2} \mathbf{x}_i + \widehat{\mathbf{x}}$$

It can be easily seen that as $\mu^2 \rightarrow \infty$, $\widehat{\mathbf{x}}(\mu) \rightarrow \widehat{\mathbf{x}}$.

Generally, if possible, it is better to eliminate the constraints, since this makes the problem better conditioned. We will discuss one final approach to solving the linearly constrained least squares problem (4.121), which reduces the dimensionality of the problem by eliminating the constraints. It is hinged on computing the QR factorization of C .

$$Q^T C = \begin{pmatrix} R \\ 0 \end{pmatrix} \begin{matrix} p \\ n-p \end{matrix} \quad (4.127)$$

This yields

$$A Q^T = \begin{pmatrix} A_1 & A_2 \end{pmatrix} \quad (4.128)$$

and

$$Q^T \mathbf{x} = \begin{pmatrix} \mathbf{y} \\ \mathbf{z} \end{pmatrix} \begin{matrix} p \\ n-p \end{matrix} \quad (4.129)$$

The constrained problem then becomes

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^n} \quad & \|\mathbf{b} - A_1\mathbf{y} - A_2\mathbf{z}\|^2 \\ \text{subject to} \quad & R^T\mathbf{y} = \mathbf{0} \end{aligned}$$

Since R is invertible, we must have $\mathbf{y} = \mathbf{0}$. Thus, the solution $\hat{\mathbf{x}}$ to the constrained least squares problem can be determined as

$$\hat{\mathbf{x}} = Q^T \begin{pmatrix} \mathbf{0} \\ \hat{\mathbf{z}} \end{pmatrix} \quad (4.130)$$

where

$$\hat{\mathbf{z}} = \underset{\mathbf{z}}{\operatorname{argmax}} \|\mathbf{b} - A_2\mathbf{z}\|^2$$

It can be proved that the matrix A_2 is at least as well-conditioned as the matrix A . Often, the original problem is singular and imposing the constraints makes it non-singular (and is reflected in a non-singular matrix A_2).

4.8.3 Least Squares with Quadratic Constraints

The quadratically constrained least squares problem is often encountered in regularization problems and can be stated as follows.

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^n} \quad & \|A\mathbf{x} - \mathbf{b}\|_2^2 \\ \text{subject to} \quad & \|\mathbf{x}\|_2^2 = \alpha^2 \end{aligned}$$

Since the objective function as well as the constraint function are convex, the KKT conditions (*c.f.* Section 4.4.4) are necessary and sufficient conditions for the optimality of the problem at the primal-dual variable pair given by $(\hat{\mathbf{x}}, \hat{\mu})$. The KKT conditions lead to the following equations

$$(A^T A + \mu I)\mathbf{x} = A^T \mathbf{b} \quad (4.131)$$

$$\mathbf{x}^T \mathbf{x} = \alpha^2 \quad (4.132)$$

The expression in (4.131) is the solution to what the statisticians sometimes refer to as the ridge regression problem. The solution to the problem under consideration has the additional constraint though, that the norm of the solution vector $\hat{\mathbf{x}}$ should equal $|\alpha|$. The two equations above yield the so-called *secular equation* stated below.

$$\mathbf{b}^T A(A^T A + \mu I)^{-2} A^T \mathbf{b} - \alpha^2 = 0$$

Further, the matrix A can be diagonalized using its singular value decomposition $A = U\Sigma V^T$ to obtain the following equation which is to be solved.

$$\sum_{i=1}^n \beta_i^2 \frac{\sigma_i^2}{(\sigma_i^2 + \mu)^2} - \alpha^2 = 0$$

4.8.4 Total Least Squares

The total least squares problem is stated as

$$\begin{array}{ll} \min_{\mathbf{x} \in \mathbb{R}^n, E \in \mathbb{R}^{m \times n}, \mathbf{r} \in \mathbb{R}^m} & \|E\|_F^2 + \|\mathbf{r}\|_2^2 \\ \text{subject to} & (A + E)\mathbf{x} = \mathbf{b} + \mathbf{r} \end{array}$$

Chapter 5

Statistics

Descriptive Statistics are used to describe the basic features of the data gathered from an experimental study. They provide *summary statistics* about the sample and the measures. In conjunction with *graphical descriptive statistics*, they form the basis of almost every quantitative analysis of data. Summary statistics include measures such as mean, median, mode, standard deviation, correlation (between two variables) *etc.* For example, the average height of the students in a room summarizes the heights of students in the room. Graphical descriptive statistics comprises of plots (such as scatter plots), charts (such as pi-charts, bar-charts), lines and histograms that describe some real world or simulated data without having you look at every single value. A graphical statistic in the form of a histogram or pi-chart of the heights of the students in a room, provides a broader summary.

In the later part of this chapter, we will introduce probability and sampling, with focus on application to data analysis. Some other topics to be discussed include confidence intervals, significance tests, hypothesis testing, *etc.* An example use of hypothesis testing in the drug industry is to decide statistically (based on past records) whether or not a specific drug is safe.

5.1 Variables

A *variable* is an attribute of any system which may change its value while it is under observation. For example, the number of people in the age group 75 – 79 in Canada is a variable. There are two basic types of variables, *viz.*, *qualitative* variables and *quantitative* variables. Qualitative variables take on values that are words. Examples of qualitative variables include marital status, nationality, color of skin, gender, *etc.* On the other hand, quantitative variables take on numeric values. Further, the quantitative variables could be either *discrete* or *continuous* based on whether the variable could take on only whole number values or any number respectively. Examples of discrete quantitative variables are number of people in a New York, number of cars in Germany or anything that

can be counted. Money, height, weight, age, *etc.*, are examples of continuous variables. The study of statistics is more concerned with qualitative variables than quantitative variables.

5.2 Descriptive Statistics

The histogram will be our first example of a (graphical) descriptive statistic. A histogram provides a picture of the data for single-variable data. We will thereafter discuss the scatter plot (or scatter diagram), which serves as a graphical description of data of two variables.

5.2.1 Histograms

A histogram is a graphical display of tabulated frequencies. A histogram shows what proportion of cases fall into each of several or many specified categories. In a histogram, it is the area of the bar that denotes the value, not the height. This is a crucial distinction to note, especially when the categories are not of uniform width. Table 5.1 shows a slightly modified estimate of the population of Canada by age group ¹ for the year 2007. The first column records the class intervals. *Class intervals* are ranges that the variable is divided into. Each class interval includes the left endpoint but not the right (by convention). The population (second column) is recorded in the thousands. The third column has a record of the percentage of the population that happens to fall in each age group.

There are three steps to be followed when plotting a histogram for tabulated frequencies as in Table 5.1.

1. Convert counts to percentages *percent* as shown in the third column of Table 5.1.

$$percentage = \frac{count}{total\ number\ of\ values}$$

2. Compute height for each class-interval as $height = \frac{percent}{width\ of\ range}$ as shown for the fourth column of Table 5.1.
3. Draw axes and label them. Label the class intervals (age groups in this case) along the x -axis and the heights along the y -axis.
4. Along each class interval on the x -axis, draw a rectangle of corresponding height and width as shown in Figure 5.1. This is precisely the histogram for the tabulated data in Table 5.1.

Figure 5.1 shows the histogram corresponding to Table 5.1. Note that the sum total area of all the bars is 100 (percent).

Histograms for discrete and continuous variables look slightly different. For histograms of continuous variables, class intervals (such as age ranges) are

¹Source: <http://www40.statcan.ca/101/cst01/demo10a.htm>

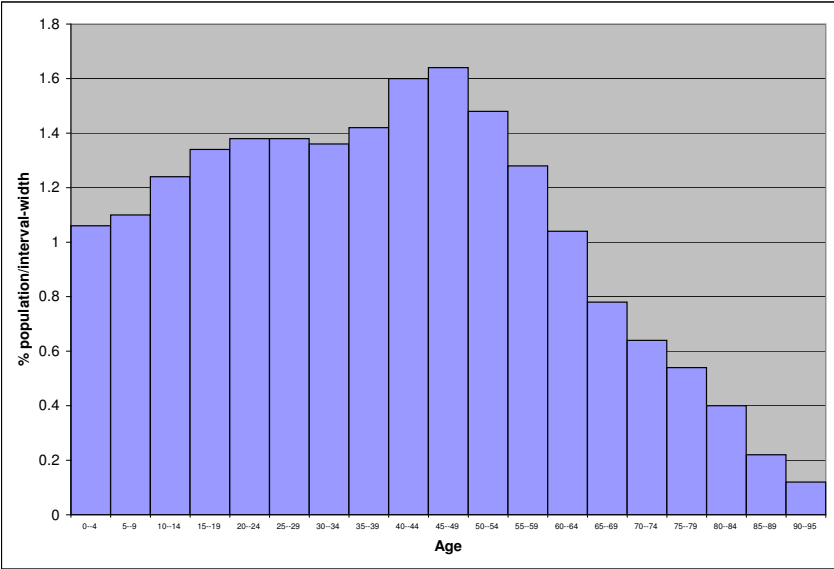


Figure 5.1: Histogram for Table 5.1.

Age group (class interval)	Persons (thousands) (count)	% of total for each group (area of histogram)	Height
0 to 4	1,740.20	5.3	1.06
5 to 9	1,812.40	5.5	1.1
10 to 14	2,060.50	6.2	1.24
15 to 19	2,197.70	6.7	1.34
20 to 24	2,271.60	6.9	1.38
25 to 29	2,273.30	6.9	1.38
30 to 34	2,242.00	6.8	1.36
35 to 39	2,354.60	7.1	1.42
40 to 44	2,640.10	8	1.6
45 to 49	2,711.60	8.2	1.64
50 to 54	2,441.30	7.4	1.48
55 to 59	2,108.80	6.4	1.28
60 to 64	1,698.60	5.2	1.04
65 to 69	1,274.60	3.9	0.78
70 to 74	1,047.90	3.2	0.64
75 to 79	894.7	2.7	0.54
80 to 84	650.8	2	0.4
85 to 89	369.3	1.1	0.22
90 to 95	186.2	0.6	0.12
Total	32,976.00	100	-

Table 5.1: A slightly modified estimate of the population of Canada by age group for the year 2007. The population (second column) is recorded in the thousands.

marked along the x -axis of the histogram and the width of the bars could be positive real number. On the other hand, histograms of discrete variables have generally a default width of 1 for every bar and different values assumed by the discrete variable are marked along the x -axis. Each bar must be centered on the corresponding value of the discrete variable and the height of every bar is the percentage value.

Density Scale

The height plotted along the y -axis of a histogram is often referred to as the *density scale*. It measures the ‘crowdedness’ of the histogram in units of ‘% per x unit’; taller the histogram bar, more is the density. In the last example, the

unit was census. Using Table 5.1 and the corresponding density estimates in Figure 5.1, one can estimate that the percentage of population aged between 75 and 77 years of age is around $\frac{2.7}{5} \times 3 = 1.62\%$. This is assuming that the density of population in the age group 75 – 79 is evenly distributed throughout the interval (that is the bar is really flat). But a close look at the bars surrounding that for 75–79 will suggest that the density in the interval 75–59 is probably not quite evenly distributed. While it would be accurate and lossless to have population counts corresponding to every age (instead of intervals), such data may not be as easy to digest as the population estimates based on intervals. There is a tradeoff between summarization and elaborate accounting or equivalently between wider bars and lots of bars.

5.2.2 Scatter Diagram

Suppose we are provided data comparing the marks (out of 100) obtained by some 500 students in the mathematics subjects in the year 1 and year 2 of a certain college. Consider a plot with ‘year 1 marks’ plotted along the x-axis and ‘year 2 marks’ plotted around the y-axis. The scatter diagram (or plot) for this marks data will consist of a point marked per student with its coordinates given by ‘(marks in year 1, marks in year 2)’. Figure 5.2 shows the scatter plot for some such hypothetical data. The dotted vertical and horizontal lines mark the average marks for year 1 and year 2 respectively. It can be seen from the plot that most students either performed well in both years or performed poorly in both years.

A point corresponding to an observation that is numerically far away from the rest of the points in a scatter plot is called an *outlier*. Statistics derived from data sets that include outliers can be misleading. Figure 5.3 shows the scatter plot of Figure 5.2 with an outlier introduced (in the form of a black point). The outlier results in a relatively drastic change in mean values of marks for years 1 and 2. While the mean value along the x-axis drops from 50.72 to 50.68, the mean value along the y-axis increases from 55.69 to 55.74.

5.3 Summary statistics

Summary statistics are numbers that summarize different features of a dataset. There are summary statistics such as mean, median, standard deviation for data of single variables and measures such as correlation coefficient for two variables.

5.4 Standard measures of location

The *arithmetic mean* μ (or simply the *mean*) is one of the most popular summary statistics and it is the average value assumed by the variable. If a variable V assumes values given by the set of data $\mathcal{V} = \{v_1, v_2, \dots, v_N\}$, then the mean μ

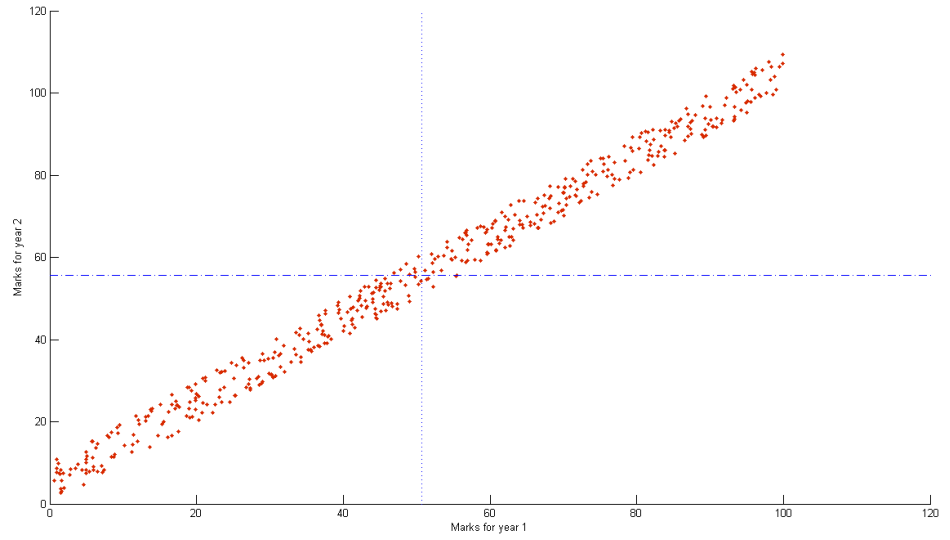


Figure 5.2: A sample scatter plot.

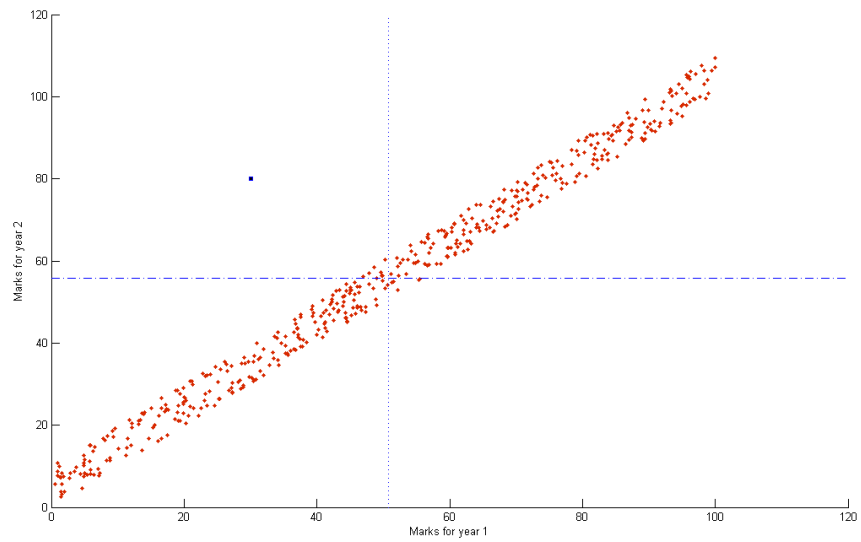


Figure 5.3: The sample scatter plot of Figure 5.2 with an outlier (30, 80) .

of V is computed as

$$\mu = \frac{\sum_{i=1}^N v_i}{N}$$

The mean is also the balancing point on a histogram; that is, if you think of the x -axis of the histogram as the beam of weight balance and weights to be proportional to the areas of the bars, then the fulcrum placed at the mean point will ensure that the beam stays horizontal.

Another measure of ‘center’ for a list of numbers is the *median*. The median is the number ν such that at least half the numbers in \mathcal{V} are less than or equal to ν and at least half the numbers are greater than or equal to ν . In order to determine the median, you need to sort the numbers (in either the ascending or descending order) and just pick the center. If more than one value qualifies as the middle value, their average (arithmetic mean) is taken to be the median. The median is the point on the x -axis of a histogram such that half of the total area of the bars lies to its left and half to its right.

As an example, the average of the set $\mathcal{V}' = \{1, 3, 4, 5, 7\}$ is 4, while its median is also 4. On the other hand, if the last number in this set is changed from 7 to 12 to yield the set $\mathcal{V}'' = \{1, 3, 4, 5, 12\}$, then the median remains 4, while the mean changes to 5. Thus, the median cares more for the number of values to its left and right rather than the actual values to its left and right. For the set $\mathcal{V}''' = \{1, 1, 3, 4, 5, 12\}$, the mean is $\frac{13}{3}$, whereas the median is the average of 3 and 4, which is 3.5. In general, for a symmetric histogram, the arithmetic mean equals the median. For a longer (shorter) right tail, the arithmetic mean is greater (smaller) than the median.

In most applications, mean is preferred over median as a measure of center. However, when the data is very skewed, median is preferred over mean as a measure of center. For example, while computing summary statistics for incomes, median is often preferred over mean, since you do not want a few very huge incomes to affect your measure of center. Similarly, median is preferred over mean as a measure of center of housing prices.

An often encountered measure is *percentile*. The k^{th} percentile of a set of values of a variable is the value (or score) of the variable below which k percent of the data points may be found. The 25th percentile is also known as the *first quartile*; the 50th percentile happens to be the median.

5.4.1 Standard measures of spread and association

The measures of center discussed thus far do not capture how spread out the data is. For example, if the average height of a class is 5 feet, 6 inches (5' 6"), it could be that everyone in the class has the same height or that someone in the class is just 4' 5" and the tallest student is 6' 3". The *interquartile range* is an illustrative but rarely used measure of the spread of data. It is defined as the distance between the 25th percentile and the 75th percentile. Generally, smaller the interquartile range, smaller will be the spread of the data.

An often used measure of spread is the *standard deviation* (SD), which measures the typical distance of a data point from the arithmetic mean. It is computed as the root mean square² (rms) of deviations from the arithmetic mean. That is, given a variable V that assumes values given by the set $\mathcal{V} = \{v_1, v_2, \dots, v_N\}$, if μ is the arithmetic mean of \mathcal{V} , then the standard deviation σ of \mathcal{V} is

$$\sigma = \sqrt{\frac{\sum_{i=1}^N v_i^2}{N}}$$

The SD of the set $\mathcal{V}' = \{1, 3, 4, 5, 7\}$ is 2, which is a typical distance of any number in \mathcal{V}' from the mean $\mu = 4$. Formulated by Galton in the late 1860s, the standard deviation remains the most common measure of statistical spread or dispersion. The square of standard deviation is called *variance*.

Effect of change of scale

What is the effect of modifying the data \mathcal{V} on the summary statistics of the data such as arithmetic mean, median and standard deviation? The effect of some data modifications have been studied in the past and are enumerated below.

1. *Adding a constant to every number of the data:* The effect is that arithmetic mean and median go up by that constant amount, while the standard deviation remains the same. This is fairly intuitive to see.
2. *Scaling the numbers in data by a positive constant:* The effect is that the arithmetic mean, the median and the standard deviation get scaled by the same positive constant.
3. *Multiplying numbers in data by -1 :* The average and the median get multiplied by -1 , whereas standard deviation remains the same.

5.4.2 Correlation coefficient

Correlation coefficient (r) measures the strength of the linear relationship between two variables. If the points with coordinates specified by the values of the two variables are close to some line, the points are said to be strongly correlated, else they are weakly correlated. More intuitively, the correlation coefficient measures how well the points are clustered around a line. Also called, *linear association*, the correlation coefficient r between sets of N values \mathcal{X} and \mathcal{Y} assumed by variables x and y respectively can be computed using the following three steps.

²The root mean square (rms) measures the typical ‘size’ or ‘magnitude’ of the numbers. It is the square root of the mean of the squares of the numbers. For example, the rms of the set $\{1, 3, 4, 5, 7\}$ is $\sqrt{20} = 4.47$.

1. Convert the values in \mathcal{X} and \mathcal{Y} into a set of values in standard unit, *viz.*, \mathcal{X}_{su} and \mathcal{Y}_{su} respectively. Recall that computing standard units requires knowledge about the mean and standard deviation and could therefore be an expensive step.
2. Let $\mathcal{P}_{su} = \{p_{su} = x_{su}y_{su} \mid x_{su} \in \mathcal{X}_{su}, y_{su} \in \mathcal{Y}_{su}\}$.
3. Let μ_{su} be the arithmetic mean of values in \mathcal{P}_{su} . The the correlation coefficient $r = \mu_{su}$.

Thus, if μ_x and μ_y are the means of x and y , and σ_x and σ_y are the respective standard deviations, then

$$r = \frac{\sum_{v_x \in \mathcal{X}, v_y \in \mathcal{Y}} \frac{v_x - \mu_x}{\sigma_x} \times \frac{v_y - \mu_y}{\sigma_y}}{N} \quad (5.1)$$

The sample scatter plot of Figure 5.2 is reproduced in Figure 5.4 with four regions marked, which are all bordered by the average lines. Points (v_x, v_y) in regions (1) and (3) contribute as positive quantities in the summation expression for r , whereas points (v_x, v_y) in regions (2) and (4) contribute as negative quantities. The correlation coefficient has no units and is always between -1 and $+1$; if r is $+1$ (-1), the points are on a line with positive (negative) slope. A simple case for which $r = 1$ is when all values assumed by y are scalar multiples of the corresponding values of x . If $r = 0$, the variables are uncorrelated. Two special but (statistically) uninteresting cases with $r = 0$, are when either of the variables always takes a constant value. Other interesting cases with $r = 0$ are when the scatter plot is symmetrical with respect to any horizontal or vertical line.

As an example, the correlation coefficient between the marks in years 1 and 2, for the data in Figure 5.2 is a positive quantity 0.9923. On the other hand, the correlation coefficient between the weight and mileage of cars is generally found to be negative. O-rings are one of the most common gaskets used in machine design. The failure of an O-ring seal was determined to be the cause of the Space Shuttle Challenger disaster on January 28, 1986. The material of the failed O-ring was a fluorinated elastomer called FKM, which is not a good material for cold temperature applications. When an O-ring is cooled below its glass transition temperature (Tg), it loses its elasticity and becomes brittle. In fact, the correlation coefficient between the extent of damage to the O-ring and temperature has been found to be negative.

There are some words of caution that one should exercise while interpreting and applying scatter plots:

1. Extrapolation is generally not a good idea for determining exact values. This is because, outside the range of values considered, the linear relationship might not hold.

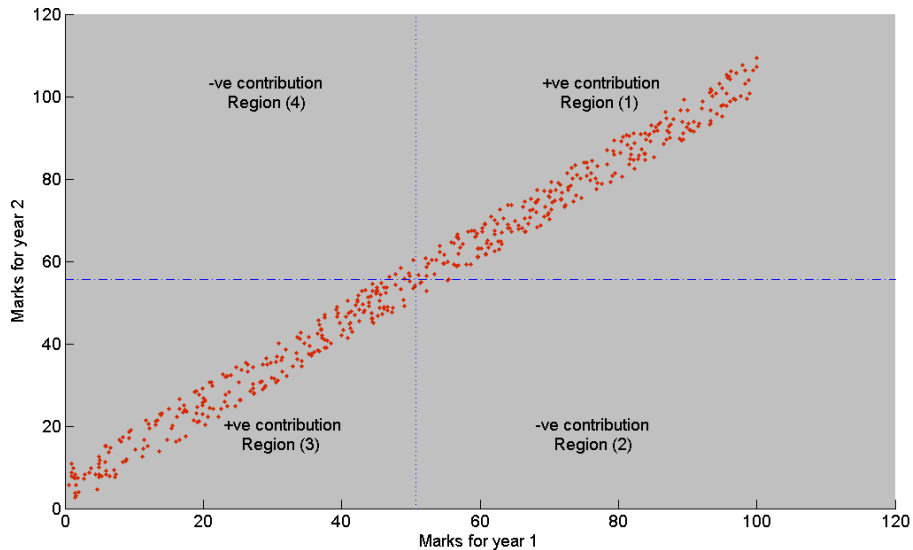


Figure 5.4: The sample scatter plot of Figure 5.2 reproduced with four regions marked based on positive and negative contributions to the correlation coefficient.

2. Even if you cannot do extrapolations, scatter plots can be informative and could give us hints about general trends (such as whether the value of one variable will increase with increase in value of the other variable).

While correlation measures only the strength of the linear association between two variables, the relationship could also be non-linear. In such cases, the scatter plot could show a strong pattern that is not linear (as an example, the scatter plot could assume the shape of a boomerang) and therefore, the quantity r is not as meaningful.

A word of caution before we move on; correlation does not imply causation, while it could definitely make one curious about a possible causal relationship. Just because the GPAs of students and their incomes are positively correlated, we cannot infer that high GPAs are caused by high incomes or vice versa. There could be latent cause of both observations, resulting in the positive correlation. For example, there is generally a positive correlation between the number of teachers and the number of failing students at a high school. But this mainly because generally, larger a school, larger is the number of teachers, greater is the student population and consequently, more are the number of failing students. Therefore, instead of treating the scatter diagram or the correlation measure as a proof of causation, these could be used as indicators that might possibly signal causation.

Effect of change of scale

The effects of change of scale on correlation are far simpler than they happened to be for arithmetic mean and standard deviation. We will list some effects on SD of changing a single variable. The effects of changing values of both variables can be derived by systematically considering effects produced by changing value of each variable.

1. When a constant is added or subtracted from every value of any single variable, the correlation coefficient stays the same, since such an operation involves translating all points and average lines by the same constant value along the corresponding axis. Consequently, the relative positions with respect to the closest line (or the standard units) remain the same.
2. When every value of any single variable is multiplied by the same constant, the correlation coefficient remains the same, since the standard units of the points remain the same (since the average and SD get scaled by the same amount as the values).
3. When every value of any single variable is multiplied by -1 , the signs of the values of each variable in standard units change (the value and mean change signs, whereas, the SD does not). Thus, r gets multiplied by -1 . However, if each value of both the variables is multiplied by -1 , the overall correlation coefficient will remain unchanged.
4. When the values of x are switched with the values of y , the correlation coefficient stays the same, since the terms within the summation expression for r in (5.1) remain the same.

5.4.3 Ecological Correlation

In contrast to a correlation between two variables that describe individuals, *ecological correlation* is a correlation calculated based on averages (or medians) of subgroups. The subgroups could be determined based on properties of the data and ecological correlation is just the correlation between means of the subgroups. For instance, the subgroups of students within a class could be determined by sections within the class or by zipcode of the residential area (which is indicative of the affluence) of the students. The ecological correlation between the incomes and the grades of students in a class could then be the standard correlation coefficient between the arithmetic means of the incomes and grades of students within each section or zipcode category. Some researchers suggest that the ecological correlation gives a better picture of the outcome of public policy actions [?]. However, what holds true for the group may not hold true for the individual and this discrepancy is often called the *ecological fallacy*. It is important to keep in mind that the ecological correlation captures the correlation between the values of the two variables across the subgroups (such as the zip code of residence) and not across individual students. The ecological correlation can help one draw a conclusion such as ‘*Students from*

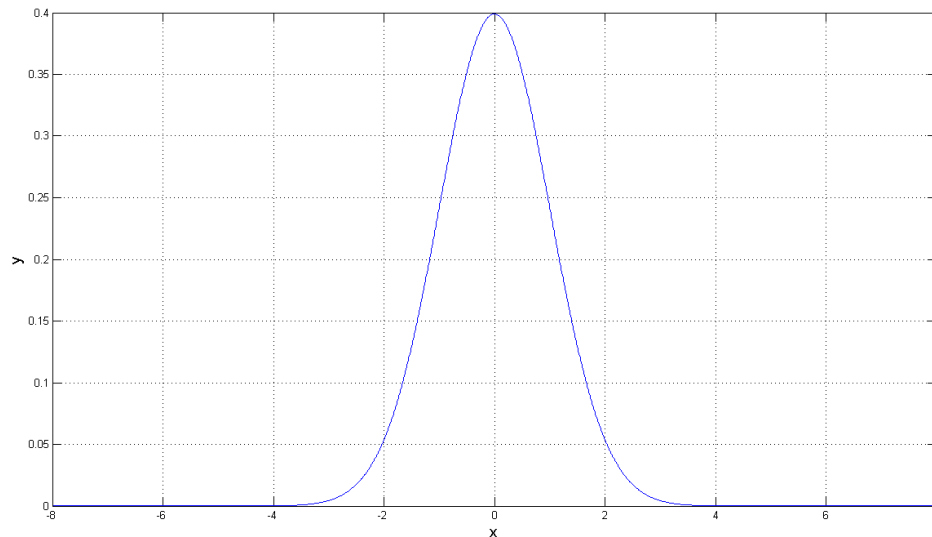


Figure 5.5: The standard normal curve.

wealthier zip codes have, on average, higher GPAs'. A recurring observation is that correlations for subgroup averages are usually larger than correlations for individuals.

5.5 Normal approximation

5.5.1 Standard Normal Curve

A standard normal curve is given by

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

where $2\pi \approx 3.14159\dots$ and $e \approx 2.71828\dots$. Figure 5.5 shows the standard normal curve.

To find the area under a normal curve between $-\infty$ and any x , there exists a standard (area) table of precomputed values, which you need to only lookup. The normal curve has some interesting properties.

1. The total area under the standard normal curve is 1. Therefore, when the rhs of the normal curve is multiplied by 100, the normal curve turns out to be a histogram - the total area under the curve of $f'(x) = \frac{100}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$ is 100.

2. The standard normal curve is symmetric about 0. This is revealed by the fact that $x^2 = (-x)^2$. Consequently, it follows from an earlier discussion that the normal curve has the same mean and median value of 0.
3. The standard deviation of the standard normal curve is 1.
4. The value of y decays very rapidly as x drifts away from 0. This is also reflected in the fact that an area of 0.998 is spanned by the normal curve lying between $x = -4.45$ and $x = 4.45$. The area table mentioned above, in fact ends at the value of $x = 3.009$. Further, the area between $x = -1$ and $x = 1$ is 68.2% (roughly two thirds) while that between -2 and 2 is 95.45%. That is, around 68% of the area is within one standard deviation of the average while roughly 95% of the area is within two standard deviations of the average.

Table 5.2 lists the value of $f(x)$ and the area between $-\infty$ and x for different values of x . Such a table is called the density table. The motivation behind discussing the normal curve is that many histograms look like the normal curve. The idea of much of mathematics and statistics is to make problems look like problems solved before.

5.5.2 General Normal Curve

The standard normal curve is centered around 0 and has a spread, though in practice, you will encounter normal curves that are not centered around 0. For example, the histogram of the heights (measured in inches) of people has the shape of a normal curve, but centered around a non-trivial height and having a spread different from 1. The general normal curve is specified by

$$f(x; \mu, \sigma) = \frac{x - \mu}{\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma}}$$

The mapping between a general normal curve and a standard normal curve is specified using the idea of standard units, which will be discussed next.

5.5.3 Standard Units

Standard units (SU) tell how many standard deviation units (SDs) above or below average a particular value on a general curve is. Thus a standard unit of 1 corresponds to distance equal to standard deviation σ from the average μ . The number of standard deviations sd of a value v is computed as

$$sd = \frac{v - \mu}{\sigma}$$

Thus, the SU of 78 for a mean value of 72 and a SD of 6 is 1 unit above average. Around $\frac{100-68}{2} = 16\%$ of the points on a normal curve are above the standard unit of 1 and therefore, the point 78 corresponds to a percentile rank of $100 - 16 = 84$ for the normal curve. On the other hand, the point

corresponding to percentile rank of 90 is obtained by looking up the standard unit as 1.30 corresponding to the approximate area value of 80 (*i.e.*, 80.64) in the density table, which translates to a value of $72 + 1.3 \times 6 = 79.8$ (for the mean and SD values of 72 and 6 respectively). In this way, the percentile rank tells you what percentage of the area under a normal curve is below a particular value.

5.6 Regression

Regression is a technique used for the modeling and analysis of numerical data consisting of values of a dependent variable y (response variable) and of a vector of independent variables \mathbf{x} (explanatory variables). The dependent variable in the regression equation is modeled as a function of the independent variables, corresponding parameters (“constants”), and an error term. However, the relationship between y and need not be causal (as in the case of correlation). Regression is used in several ways; one of the most often used ways is to estimate the average y value corresponding to a given \mathbf{x} value. As an example, you might want to guess the inflation next year, based on the inflation during the last three years. We will take three examples and discuss how ‘good’ the average estimate of the dependent variable could be.

1. In the case of the scatter diagram in Figure 5.6, one could make use of regression to estimate the average possible marks of a student in ‘year 2’, given her marks in ‘year 1’. In this case, the correlation coefficient (0.993) is very close to 1. The mean values for x and y are 60.26 and 60.33 respectively while the corresponding standard deviations are 10.48 and 10.54 respectively. The graph that shows the average y value corresponding to each x value is called the *graph of averages*. It can be seen in Figure 5.6, that the average values guessed for y corresponding to different x values are very close to a line and have very little variation across the line. Such a line is called the *regression line*. The regression line is a smoothed version of the graph of averages.
2. In contrast, if the x variable represented a student’s score in mathematics and the y variable represented her score in physics, the correlation might not be as strong. For a typical scatter plot as in Figure 5.7 for such data, the variance across every average estimate of y for every x value is larger and the correlation coefficient is further away from 1 (0.801 to be precise).
3. In the case of two uncorrelated variables (*i.e.*, $r \approx 0$) such a person’s height and his average scores in mathematics, the average value of y given any particular value of x will be more or less the same. The variance will also be large.

It turns out actually that if x is one SD above its average μ_x , the corresponding predicted value of y is r SDs above its average μ_y . The point (μ_x, μ_y) is in fact termed as the *point of averages*. In the third example discussed above,

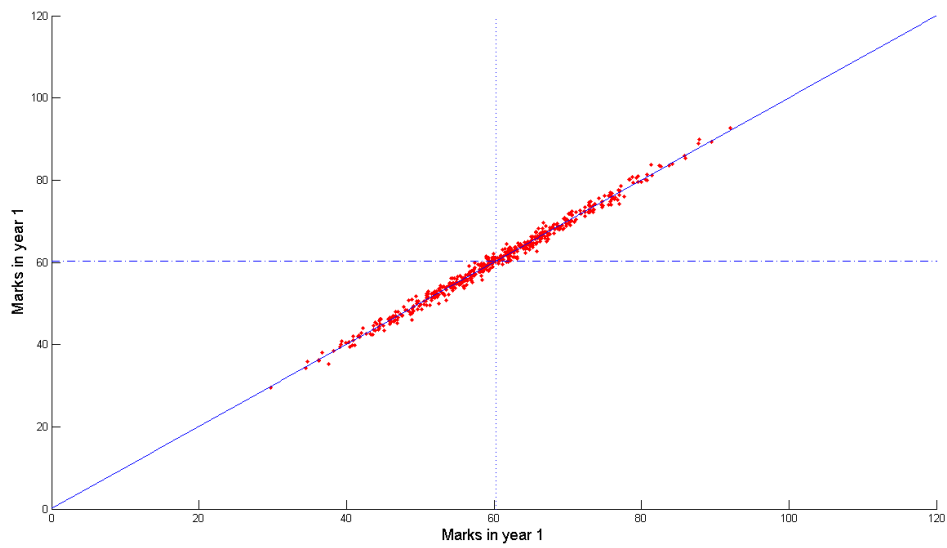


Figure 5.6: A scatter plot for data similar to that in Figure 5.2, along with the regression line. The mean values for x and y are 60.26 and 60.33 respectively while the corresponding standard deviations are 10.48 and 10.54 respectively. The correlation coefficient is 0.993.

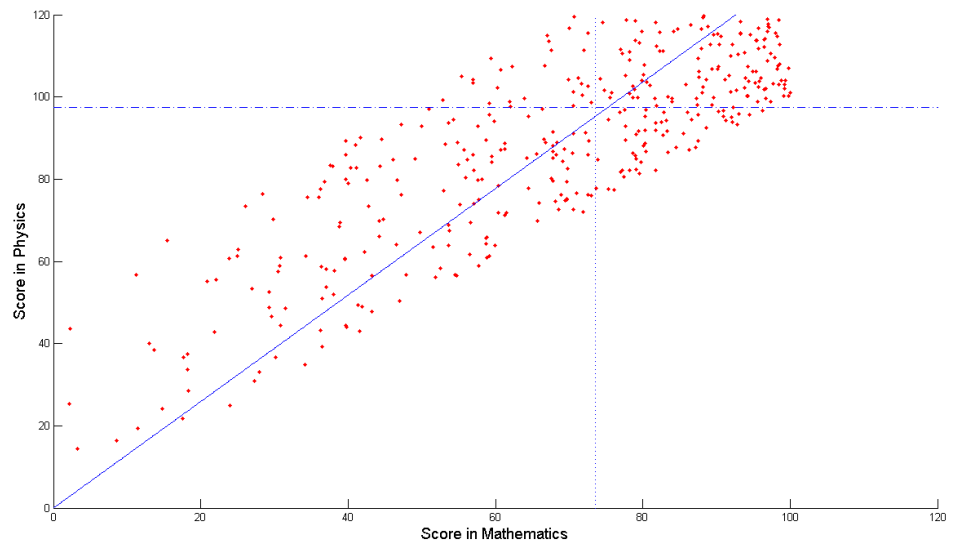


Figure 5.7: A scatter plot for the marks of students of some class in mathematics (x -axis) against physics (y -axis), along with the regression line. The mean values for x and y are 60.34 and 60.01 respectively while the corresponding standard deviations are 9.82 and 12.06 respectively. The correlation coefficient is 0.801.

with $r \approx 0$, the predicted y hard changes with the value of x . For the first example, the predicted y changes the most with the value of x .

The *linear regression method* concerns itself with finding points on the regression line, for a given table of data values with each row consisting of values of the independent variable x and the dependent variable y for a particular data point. The following three steps comprise the linear regression method.

1. Compute the correlation coefficient r between x and y . For the data in Figure 5.6, the correlation coefficient was 0.993.
2. Convert any given value of x into its standard units. Multiply this number by r to get the predicted value of y in its standard units. This step is essentially the *regression step*. For example, with respect to Figure 5.6, consider $x = 76$. This number in standard units will be $\frac{76-60.26}{10.48} = 1.50$. The predicted value of y in standard units will be 1.49.
3. Convert the value of y from standard units into normal units. The value of 0.993 for y in standard units will correspond to $0.993 \times 10.54 + 60.33 = 70.80$.

The regression line discussed above is for predicting value of y based on x (the usual notion of line). It is also possible to perform regression to predict value of x based on value of y , to yield another kind of regression line.

As an exercise, let us predict (using linear regression on the data in Figure 5.6) the marks for a student in year 2, given that she obtained 63 marks in year 1. The predicted value for marks in year 2 will be $60.33 + (0.993 \times \frac{63-60.26}{10.48}) \times 10.54 = 63.07$, which compares favourably with the actual value of $y = 63.53$.

As another exercise, with respect to the scatter plot in Figure 5.6, suppose a student is at 31st percentile in year 1. What is a good guess for this student's percentile rank in year 2? We note that the scatter diagram is '*football shaped*' (or tilted oval) which is an indication that the two variables roughly follow a normal distribution³ First we will determine the standard units for x corresponding to the 31st percentile by looking up the value in standard units corresponding to an area of $\frac{100-31}{100} = 0.69$ in standard normal Table 5.2. The value happens to be around 0.50. Multiplying this by the correlation coefficient $r = 0.993$, we get the standard units of the approximated y value to be 0.497. The corresponding percentile rank can be again obtained using the normal table to be $100 \times (1 - 0.6879) = 31.21$, *i.e.*, around 31 percentile. Note that we really did not make use of any of estimates for mean or standard deviation, since all computations dealt only with standard units.

What about the spread of values around the average predicted by linear regression? How accurate is regression? For each point in the scatter diagram, the actual value of y will not typically be on the regression line. The *prediction error* at a point x is defined as the difference between the actual and predicted values of y at that point. Analogous to the definition of the standard deviation,

³Later, we will discuss techniques that will help verify such assumptions.

the spread for linear regression is measured as the RMS of prediction errors and is called the *RMS error*. Linear regression should be suitable for datasets for which the RMS error is very close to sum of the deviations from average across all points. While RMS error can be used in the context of any prediction method, for the case of linear regression, the RMS error can be simplified to the following formula:

$$e_{rms} = \sigma_y \sqrt{1 - r^2} \quad (5.2)$$

where σ_y is the SD for y .

On the other hand, if the prediction method were to always predict the overall average μ_y of y regardless of the value of x , the RMS error would be the standard deviation σ_y of y (by definition). We can readily see that if $r \in (-1, 1)$, prediction using regression will always yield a lower value of e_{rms} than will a naive prediction using the average μ_y .

Mathematically, if $r = \pm 1$, all the points must be on a line (with slope r), resulting in 0 prediction errors and $e_{rms} = 0$. On the other hand, if $r = 0$ and $e_{rms} = \sigma_y$. This is because there is no linear relationship between x and y , implying that it does not help to know x while predicting y ; you can always guess y to be μ_y and could do no better. For any other $r \in (-1, 0) \cup (0, 1)$, $e_{rms} \in (0, \sigma_y)$. If e_{rms} is close to 0, it means that the predictions are very accurate. Whereas, if e_{rms} is close to σ_y , the predictions must be far from accurate.

What about answer to a question such as ‘*what percentage of all values are within one unit of RMS error of the regression line for a fixed value of x* ’? The answer will be around 68%, assuming that the points approximately follow a normal curve and that the average values of y for any given x approximately coincide with the point (x, rx) on the regression line, so that one unit of RMS error approximately corresponds to one standard deviation. For a football shaped scatter diagram, the normality condition can be assumed to hold for both x and y as well as for all points with a fixed value of x . Using this information, let us answer the following question: of all the students considered in Figure 5.6 who obtained a score of 65 in year 1, what percentage of them obtained over 67 in year 2? First of all, the average marks in year 2, for students who obtained 65 in year 1, predicted using linear regression will be $60.33 + (0.993 \times \frac{65 - 60.26}{10.48}) \times 10.54 = 65.06$. The RMS error can be computed using (5.2) to be $\sqrt{1 - (0.993)^2} \times 10.09 = 1.19$. The question then reduces to ‘what percentage of values on a normal curve with mean 65.06 and SD 1.19 are above 67?’ The answer to this can be easily obtained by looking up the standard normal table for the percentage area that is above 1.63 standard units, which is $100 \times (1 - 0.9484) = 5.16\%$.

Further, a football shaped scatter diagram is *homoscedastic*, *i.e.*, it has approximately the same spread around the regression line for every value of x . To cite a practical example, we would expect different divisions within a class to have the same mean and SD values for scores in a particular year, irrespective of their sizes. Of course, divisions with smaller sizes may have a smaller range

of scores than larger divisions.

5.6.1 The Regression Effect

The regression effect is also called ‘regression to the mean’. This effect is essentially the fact that for $r \in (-1, 1)$, and both x and y in standard units, the predicted (based on linear regression) y value will be closer to the mean than than the x value. This is because, if x is 1 SD above average, then the predicted value of y will be $r \in (-1, 1)$ SDs above average. Thus, if a student performs exceptionally well in the first year, the predicted performance of the student in the second year will not be as good. A corollary of this effect is that if y_1 is the predicted average value of y variable for a value x_1 of the x variable, then the predicted average value of x_2 of the x variable based on value y_1 will be less than x_1 .

5.6.2 SD Line

The SD line is a line that passes through the *point of averages* in a scatter plot and has slope equal to the ratio of standard deviation of y to that of x multiplied by $sgn(r)$. Points on the SD line are the same in standard units for both x and y . It can be proved that if $r = 1$, the SD and regression lines will coincide. On the other hand, further is r away from 1, further will

5.7 Probability and sampling

5.8 Binomial distribution

5.9 Interval estimation

5.10 Some standard significance tests

z	+0.00	+0.01	+0.02	+0.03	+0.04	+0.05	+0.06	+0.07	+0.08	+0.09
0.0	0.5000	0.5040	0.5080	0.5120	0.5160	0.5199	0.5239	0.5279	0.5319	0.5359
0.1	0.5398	0.5438	0.5478	0.5517	0.5557	0.5596	0.5636	0.5675	0.5714	0.5753
0.2	0.5793	0.5832	0.5871	0.5910	0.5948	0.5987	0.6026	0.6064	0.6103	0.6141
0.3	0.6179	0.6217	0.6255	0.6293	0.6331	0.6368	0.6406	0.6443	0.6480	0.6517
0.4	0.6554	0.6591	0.6628	0.6664	0.6700	0.6736	0.6772	0.6808	0.6844	0.6879
0.5	0.6915	0.6950	0.6985	0.7019	0.7054	0.7088	0.7123	0.7157	0.7190	0.7224
0.6	0.7257	0.7291	0.7324	0.7357	0.7389	0.7422	0.7454	0.7486	0.7517	0.7549
0.7	0.7580	0.7611	0.7642	0.7673	0.7704	0.7734	0.7764	0.7794	0.7823	0.7852
0.8	0.7881	0.7910	0.7939	0.7967	0.7995	0.8023	0.8051	0.8078	0.8106	0.8133
0.9	0.8159	0.8186	0.8212	0.8238	0.8264	0.8289	0.8315	0.8340	0.8365	0.8389
1.0	0.8413	0.8438	0.8461	0.8485	0.8508	0.8531	0.8554	0.8577	0.8599	0.8621
1.1	0.8643	0.8665	0.8686	0.8708	0.8729	0.8749	0.8770	0.8790	0.8810	0.8830
1.2	0.8849	0.8869	0.8888	0.8907	0.8925	0.8944	0.8962	0.8980	0.8997	0.9015
1.3	0.9032	0.9049	0.9066	0.9082	0.9099	0.9115	0.9131	0.9147	0.9162	0.9177
1.4	0.9192	0.9207	0.9222	0.9236	0.9251	0.9265	0.9279	0.9292	0.9306	0.9319
1.5	0.9332	0.9345	0.9357	0.9370	0.9382	0.9394	0.9406	0.9418	0.9429	0.9441
1.6	0.9452	0.9463	0.9474	0.9484	0.9495	0.9505	0.9515	0.9525	0.9535	0.9545
1.7	0.9554	0.9564	0.9573	0.9582	0.9591	0.9599	0.9608	0.9616	0.9625	0.9633
1.8	0.9641	0.9649	0.9656	0.9664	0.9671	0.9678	0.9686	0.9693	0.9699	0.9706
1.9	0.9713	0.9719	0.9726	0.9732	0.9738	0.9744	0.9750	0.9756	0.9761	0.9767
2.0	0.9772	0.9778	0.9783	0.9788	0.9793	0.9798	0.9803	0.9808	0.9812	0.9817
2.1	0.9821	0.9826	0.9830	0.9834	0.9838	0.9842	0.9846	0.9850	0.9854	0.9857
2.2	0.9861	0.9864	0.9868	0.9871	0.9875	0.9878	0.9881	0.9884	0.9887	0.9890
2.3	0.9893	0.9896	0.9898	0.9901	0.9904	0.9906	0.9909	0.9911	0.9913	0.9916
2.4	0.9918	0.9920	0.9922	0.9925	0.9927	0.9929	0.9931	0.9932	0.9934	0.9936
2.5	0.9938	0.9940	0.9941	0.9943	0.9945	0.9946	0.9948	0.9949	0.9951	0.9952
2.6	0.9953	0.9955	0.9956	0.9957	0.9959	0.9960	0.9961	0.9962	0.9963	0.9964
2.7	0.9965	0.9966	0.9967	0.9968	0.9969	0.9970	0.9971	0.9972	0.9973	0.9974
2.8	0.9974	0.9975	0.9976	0.9977	0.9977	0.9978	0.9979	0.9979	0.9980	0.9981
2.9	0.9981	0.9982	0.9982	0.9983	0.9984	0.9984	0.9985	0.9985	0.9986	0.9986
3.0	0.9987	0.9987	0.9987	0.9988	0.9988	0.9989	0.9989	0.9989	0.9990	0.9990

Table 5.2: Area under the standard normal curve between $-\infty$ and z for different values of z .

Chapter 6

Searching on Graphs

Consider the following problem:

Cab Driver Problem A cab driver needs to find his way in a city from one point (A) to another (B). Some routes are blocked in gray (probably because they are under construction). The task is to find the path(s) from (A) to (B). Figure 6.1 shows an instance of the problem..

Figure 6.2 shows the map of the united states. Suppose you are set to the following task

Map coloring problem Color the map such that states that share a boundary are colored differently..

A simple minded approach to solve this problem is to keep trying color combinations, facing dead ends, back tracking, etc. The program could run for a very very long time (estimated to be a lower bound of 10^{10} years) before it finds a suitable coloring. On the other hand, we could try another approach which will perform the task much faster and which we will discuss subsequently.

The second problem is actually isomorphic to the first problem and to problems of resource allocation in general. So if you want to allocate aeroplanes to routes, people to assignments, scarce resources to lots of tasks, the approaches we will show in this chapter will find use. We will first address search, then constraints and finally will bring them together.

6.1 Search

Consider the map of routes as in Figure 6.3. The numbers on the edges are distances. Consider the problem of organizing search for finding paths from S (start) to G (goal).

We humans can see the map geometrically and perhaps guess the answer. But the program sees the map only as a bunch of points and their distances. We will use the map as a template for the computer at every instance of time.

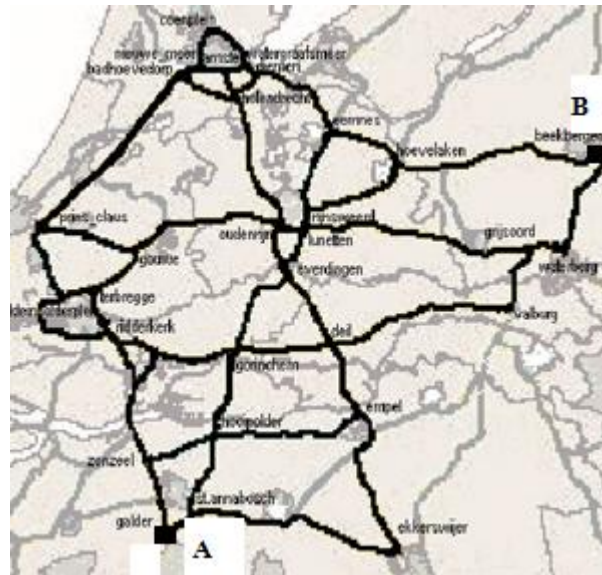


Figure 6.1: Routes in a city. The problem is to find the shortest path from point (A) to point (B).

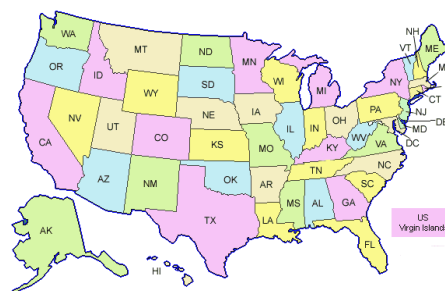


Figure 6.2: The map of the united states of America. The task is to color the map such that adjacent states have different colors.

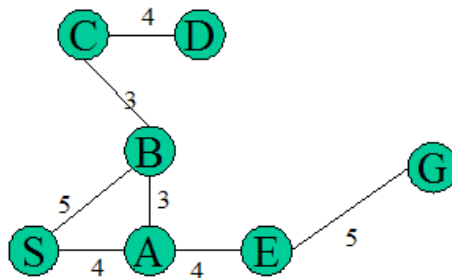


Figure 6.3: Map of routes. The numbers on the edges are distances.

At the beginning, the computer knows only about S . The search algorithm explores the possibilities for the next step originating from s . The possible next steps are A and B (thus excluding G). Further, from A , we could go to B or E . From B , we could go to A or C . From A , we could also go back to S . But we will never bother biting our own tail in this fashion; that is, we will never make a loop. Going ahead, from B , we could only go to C . From C , we can progress only to D and thereafter we are stuck on that path. From A , we could make a move to E and then from E to G . Figure 6.4 shows the exhaustive tree of possible paths (that do not bite their own tail) through this map. The process of finding all the paths is called the *British Museum Algorithm*¹.

But the British Museum algorithm can be very expensive. In fact, some of these searches can be exponential. If you had to look through a tree of chess moves for example, in the beginning it is essentially exponential, which is a bad news since it will imply 10^{10} years or so. We need a more organized approach for searching through these graphs. There exist many such organized methods. Some are better than others, depending on the graph. For example, a depth first search may be good for one problem but horrible for another problem. Words like *depth first search*, *breadth first search*, *beam search*, *A* search*, *etc.*, form the vocabulary of the search problem in Artificial Intelligence.

The representation we use will define the constraints (for example, the representation of the routes in Figure 6.3 defines the notion of proximity between nodes and also defines constraints on what sequences of vertices correspond to valid paths.

6.1.1 Depth First Search

This algorithm² boils down to the following method: Starting at the source, every time you get a choice for the next step, choose a next step and go ahead. We will have the convention that when we forge ahead, we will take the first

¹The British Museums are considered some of the largest in the world.

²Fortunately, the names given to these algorithms are representative of the algorithms actually do.

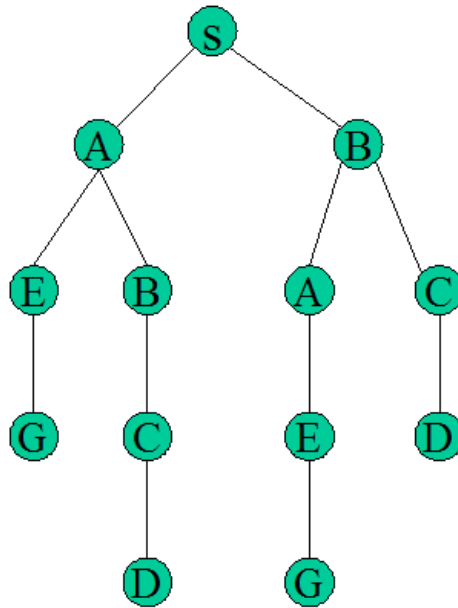


Figure 6.4: The tree of all possible paths from S to G through the map in Figure 6.3.

choice. Thus, a depth first search on the map in Figure 6.3 will tread along the path in Figure 6.5. In practice, what is used is a combination of depth first search and backup. What this means is that when a dead end is hit, the method goes back to the last state. In practice, this method could yield very complicated solutions.

6.1.2 Breadth First Search (BFS)

The way BFS organizes its examination of the tree is layer by layer; the algorithm first explores one layer of solutions (A or B), then forges ahead to another layer (B or E or A or C) and so on. Figure 6.6 illustrates the BFS traversal for the map in Figure 6.3. In practice, this search could expend a lot of time in useless parts of the graph. But it yields the shortest path in terms of number of streets covered.

6.1.3 Hill Climbing

Both BFS and DFS are completely uninformed about geography; neither information about distance nor direction is exploited. But often it helps if you have outside information about how good a place is to be in. For instance, in the map in Figure 6.3, between E and D , it is much better to be in E because that gets you closer to the goal. It makes less sense in general to head off to the right

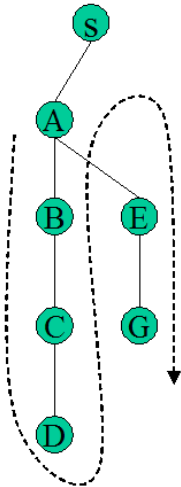


Figure 6.5: The DFS tree from S to G through the map in Figure 6.3.

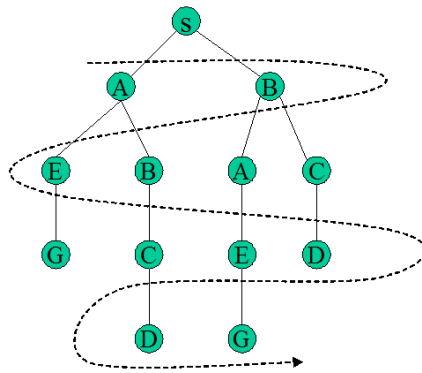


Figure 6.6: The BFS tree from S to G through the map in Figure 6.3.



Figure 6.7: The choices based on hill climbing for the route from S to G through the map in Figure 6.3.

if you know that the goal is on the left. It is heuristically good to be closer to the goal, though sometimes it may turn out not to be a good idea. So we could make use of heuristic information of this kind. And this forms the idea behind hill climbing. It is an idea drafted on top of depth first search. When initially at S , you could move to A or B . When you look at A and B , you find that one of them is closer to the goal G and you choose that for the next move. And B happens to be closer to G , which you pick for the next move; but this turns out to be a bad idea as we will see. Nevertheless, it looks good from the point of view of the short-sighted hill climbing algorithm. From B , the possible choices are A and C . And A is a natural choice for hill-climbing, owing to its proximity to G . From A onwards, the choices are straightforward - you move to E and then to G . Figure 6.7 shows the route from S to G as chalked out by hill climbing.

Hill climbing always is greedy because it plans only for one step at a time. The method requires a metric such as distance from the goal.

6.1.4 Beam Search

We drafted hill climbing on top of depth first search. Is there a similar thing we could do with breadth first search? And the answer is ‘yes’. And this approach

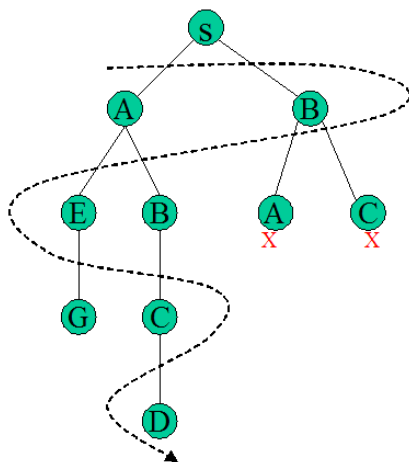


Figure 6.8: The choices based on beam search for the route from S to G through the map in Figure 6.3.

is called *beam search*. The problem with breadth first search is that it tends to be exponential. But what we could do to work around is to throw away at every level of BFS, all but the ‘most promising’ of the paths so far. The most promising step is defined as that step which will get us ‘closest’ to the goal. The only difference from the hill climbing approach is that beam search does not pick up just one path, it picks up some fixed number of paths to carry down. Let us try beam search on the map in Figure 6.3. For this simple example, let us keep track of two paths at every level. We will refer to the BFS plan in Figure 6.6. From S , we could move to A or B and we keep track of both possibilities. Further, from A , there are two possibilities, *viz.*, B and E , while from B there are two possibilities in the form of A and C . Of the four paths, which two should we retain? The two best (in terms of the heuristic measure of how far we are from the goal) are B and E . Carrying forward from these points, we arrive at G in a straight-forward manner as shown in Figure 6.8.

While the vanilla breadth first search is exponential in the number of levels, beam search has a fixed width and is therefore a constant in terms of number of levels. Beam search is however not guaranteed to find a solution (though the original BFS is guaranteed to find one). This is a price we pay for saving on time. However, the idea of backing up can be employed here; we could back up to the last unexplored path and try from there.

6.2 Optimal Search

For optimal search, we will start with the brute force method which is exponential and then slap heuristics on top to reduce the amount of work, while still assuring success.

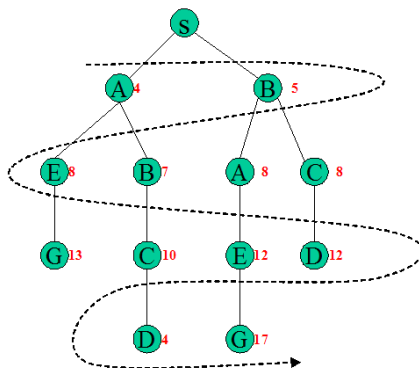


Figure 6.9: The tree of all possible paths from S to G through the map in Figure 6.3, with cumulative path lengths marked at each node.

6.2.1 Branch and Bound

Suppose the problem is to find the best possible path (in terms of distances) between S and G in the map as in Figure 6.3. An oracle suggests that the path $SAEG$ is the shortest one (its length is 13). Can we determine if $SAEG$ is indeed the shortest path? One method to answer this question is to verify if every other path is at least that long.

From S , we could go to A or B . The cumulative path length to B is 5. From B , you could either go to A or C . The cumulative path length upto A or C is 8. At this point, from A , you could go to E while from C you could move to D , with cumulative path lengths of 13 each and so on. Figure 6.9 shows the tree of possible paths from S to G , with cumulative path length (from S) marked at each node. It is evident from the figure that all paths to G have length greater than or equal to 13. We can therefore conclude that the shortest path to G from S is $SAEG$ and has length 13.

Most often we do not have any oracle suggesting the best path. So we have to think how to work without any oracle telling us what the best path is. One way is to find a path to the goal by some search technique (DFS or beam search) and use that as a reference (bound) to check if every other path is longer than that. Of course, our first search may not yield the best path, and hence we might have to change our mind about what the best path (bound) is as we keep trying to verify that the best one we got so far is in fact the best one by extending every other path to be longer than that. The intuition we work with is to always push paths that do not reach the goal until their length is greater than a path that does reach the goal. We might as well work only with shortest paths so far. Eventually, one of those will lead to the goal, with which we will almost be done, because all the other paths will be about that length too, following which we just have to keep pushing all other paths beyond the goal. This method is called branch and bound.

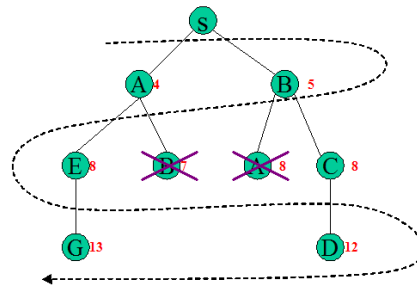


Figure 6.10: The pruned search tree for branch and bound search.

In the example in Figure 6.9, the shortest path upto B from S was of length 5 initially, which used as a bound, could eliminate the path SAB and therefore save the computation of the path $SABCD$. Similarly, the initial path to A , SA of length 4, sets a bound of 4 for the shortest path to A and can thus eliminate the consideration of the path SBA beyond A . Figure 6.10 illustrates the application of branch and bound to prune the BFS search tree of Figure 6.9. Crossed out nodes indicate search paths that are excluded because they yield paths that are longer than bounds (for the corresponding nodes). This crossing out using bounds is a variation on the theme of the dynamic programming principle.

6.2.2 A^* Search

The branch and bound algorithm can also be slapped on top of the hill climbing heuristic or the beam search heuristic. In each of these cases, the bound can be computed as the sum of the accumulated distance and the euclidian distance.

When the branch and bound technique is clubbed with shortest distance heuristic and dynamic programming principle, you get what is traditionally known as A^* search. Understanding A^* search is considered the culmination of optimal search. It is guaranteed to find the best possible path and is generally very fast.

6.3 Constraint Satisfaction

What we have discussed so far is mechanisms for finding the path to the goal (which may not be optimal). Rest of the discussion in this chapter will focus on resource allocation. Research allocation problems involve search. Too often people associate search only with maps. But maps are merely a convenient way³ to introduce the concept of search and the search problem is not restricted to maps. Very often, search does not involve maps at all. Let us again consider the

³Maps involve making a sequence of choices and therefore could involve search amongst the choices.

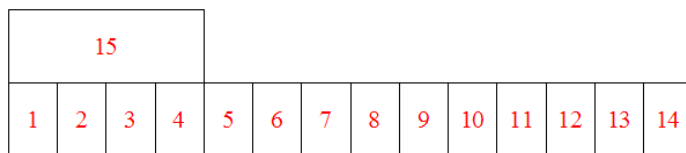


Figure 6.11: The sketch of a simple region for which we try to solve the map coloring problem.

“*Map coloring problem*” (refer to Figure 6.2). It involves searching for sequence of choices. Instead of the large city of USA, let us consider the simpler country of Simplea as in Figure 6.11. Let us say we need to pick a color from the set $\{R, G, B, Y\}$.

Let us say we number the regions arbitrarily. One approach is to try coloring them in the order of increasing numbers, which is a horrible way! This is particularly because the order chosen could be horrible. To color the region, let us start by considering the possible colors that could be assigned to each region. State number 1 could be assigned any of R, B, G or Y . Suppose we take a depth first search approach for this coloring job. Now we need to color the region number 2. Since we adopt depth first search, let us say we pick the color of R for region 1. And as we go from region to region, we keep rotating the set of colors.

The prescription for DFS is that you keep plunging head until you hit a dead end where you cannot do anything. But the problem is that we cannot infer we have hit a dead end till we assign colors to all 15 regions, to realise that the following constraint has been violated: *no two adjacent regions should have the same color*. And even if we backup a step or two, we will have to explore all 15 regions to infer if we have reached a dead end. At this rate, the projected time for successful coloring of the states in a US map is 10^{10} years⁴! That is no way to perform search. The culprit is the particular order in which we chose to color the regions. Figure 6.12 shows an example instance of random assignment of colors that leads to a dead-end (with no suitable color left for region number 15). We introduce the idea of *constraint checking* precisely to address this problem. It is essentially a (DFS) tree trimming method.

First of all, we will need some terminology.

1. *Variables (V):* V is a set of variables. Variables names will be referred to in upper case, whereas their specific instatiations will be referred to in lower case.
2. *Domain (D):* The domain D is the bag of values that the variables can take on. D_i is the domain of the $V_i \in V$. A good example of a domain

⁴The United states has 48 contiguous states. If the number of colors is 4, the branching factor of DFS will be 4 and the height will be 48. Thus, the number of computations will be 4^{48} or 2^{96} which is of the order of 10^{27} and since the length of a year is the order of 10^7 , it will take in the order of 10^{11} years assuming that a computation is performed in a nano second.

15 ?														
1	2	3	4	5	6	7	8	9	10	11	12	13	14	
R	G	B	Y	R	G	B	Y	R	G	B	Y	R	G	

Figure 6.12: A sequence of bad choices for the map coloring problem.

would be $\{R, B, G, Y\}$. In the map coloring problem, all variables have the same domain.

3. *Constraints (C)*: A constraint is a boolean function of two variables. C_{ij} represents the constraint between variables X_i and X_j . $C_{ij}(x_i, x_j) = true$ iff, the constraint is satisfied for the configuration $X_i = x_i, x_i \in D_i$ and $X_j = x_j, x_j \in D_j$. Else, $C_{ij}(x_i, x_j) = false$. The form of the constraint function could vary from problem to problem. An example of a constraint is: *no two adjacent regions should not have the same color*. Note that there need not be a C_{ij} for all i and j . For instance, in the map coloring problem, the only constraints are those involving adjacent regions, such as region 15 with the regions 1, 2, 3 and 4.

Suppose we are trying to explore if region number $i = 15$ is going to object to any coloring scheme. First we define the constraint $C_j, i = 15, j \in \{1, 2, 3, 4\}$ as follows:

$C_{ij}(x, y), i = 15, j \in \{1, 2, 3, 4\}$ is true iff $x \neq y$.

We can state the check as in Figure 6.13.

```

for  $x \in D_i$  ( $i = 15$ ) do
  for All constraints  $C_{ij}$  do
    Set  $X_i = x$  iff  $\exists y \in D_j$  such that  $C_{ij}(x, y) = true$ 
  end for
end for

```

Figure 6.13: The algorithm for constraint checking

6.3.1 Algorithms for map coloring

When assigning a color to each region, we could do different types of checks.

1. *Check all*: At one end of the spectrum, while assigning a color to a region, we check constraint satisfaction with respect to all regions (no matter how far they are from the region under consideration). This is an extreme case and potentially, we could consider all variants of what to check and what not to check. The check is only a speedup mechanism. It neither ensures nor prevents a solution. We could get a solution (if there is one) using depth first search and no search, though it might take 10^{10} years. If we

want to have minimum number of checks such that they are really helpful, our method for deciding what to check will be checks on regions that are in the immediate vicinity of the last region that had a color assigned.

2. *Check Neighbors*: Check for constraint satisfaction, only all regions in the immediate neighborhood of the last region that had a color assigned. This is at another end of the spectrum.
3. *Check Neighbors of Neighbors*: Check for constraint satisfaction, only all regions in the immediate neighborhood of the last region that had a color assigned as well as well as regions in the immediate neighborhood of the neighbors.
4. *Check Neighbors of 'unique' neighbors*: Check for constraint satisfaction, only all regions in the immediate neighborhood of the last region that had a color assigned as well as well as regions in the immediate neighborhood of the 'unique' neighbors. A neighbor is 'unique', if the constraint set for the region has been reduced to a single element (which means that it has a very tight constraint).

We will evaluate the above approaches on two fronts, *viz.*,

1. *Number of assignments performed*: An *assignment* corresponds to putting a color on a region. If we are very lucky, we might be able to color the map of the United states with just 48 assignments (that is we never had to backup). Constraint checking will enable the algorithm realise that it will soon hit a dead end and will have to back up. In practical situations, we could expect a bit of backup and also expect the constraints to show you some dead ends. So it could happen that the constraint checking reduces the number of assignments to some number slightly above the ideal (say 52).
2. *Number of checks performed*: A *check* corresponds to determining which color could be assigned to a region based on the color assigned so far to another regions.

A simple experiment for coloring the Unites States reveals the statistics in Table 6.1.

The message we can extract from Table 6.1 is that some of the constraint checking is essential, else it takes 10^{10} years. We can conclude that full propogation is probably not a worthy effort. Because full propogation, relative to propogation through unique values yielded the same number of assignments but relatively fewer checks. When we compare heuristic 2 against heuristic 4, we find that there is a tradeoff between number of assignments and number of checks. Checking all is simply a waste of time. The net result is that people tend to invoke either of heuristics 2 and 4. From the point of view of programming simplicity, one might just use heuristic number 2.

SrNo	Constraint Sat Heuristic	Assignments	Checks
1	Check all	54	7324
2	Check neighbors	78	417
3 of neighbors	Check neighbors	54	2806
4 of 'unique' neighbors	Check neighbors	54	894

Table 6.1: Number of assignments and checks for each of the constraint satisfaction heuristics for the problem of coloring the map of the USA using 4 colors.

The computational complexity of the discussed methods is still exponential⁵. But usually, the constraints try to suppress the exponential nature.

6.3.2 Resource Scheduling Problem

Let us forget about maps for a while and talk about airplanes. An upstart airline is trying to figure out how many airplanes they need. They have a schedule outlining when and where they need to fly planes. And the task is to figure out the minimum number of airplanes that need to be purchased. For every airplane saved, let us say the reward is saving on half the cost of the airplane.

F_1, F_2, \dots, F_n are the flights. Table 6.2 shows the schedule for the flights.

Flight No.	From	To	Dept. Time	Arr. Time
F_1	Boston	LGA	10:30	11:30
F_2	Boston	LGA	11:30	12:30
F_3	Boston	LGA	12:30	13:30
F_4	Boston	LGA	13:30	14:30
F_4	Boston	LAX	14:30	15:30
F_6	
...	
F_{15}	Boston	LAX	11:00	15:30

Table 6.2: Number of assignments and checks for each of the constraint satisfaction heuristics for the problem of coloring the map of the USA using 4 colors.

To fly this schedule, we have some number $m = 4$ of airplanes: P_1, P_2, P_3, P_4 .

⁵Note that these are NP complete problems. Any polynomial time algorithm will fetch field medals.

Of course, we do not want to fly any plane empty (dead head).

The assignment $P_1 \rightarrow F_1$, $P_2 \rightarrow F_2$, $P_3 \rightarrow F_3$ and $P_4 \rightarrow F_4$ will leave no airplane for flight F_5 . But if we could fly P_1 back from New York (LGA), the same plane could be used for flight F_4 , thus sparing P_4 for F_5 . We can draw the correspondence between this problem and the map coloring problem; the 4 airplanes correspond to 4 colors while the flights correspond to regions. The task is to assign planes (colors) to flights (regions), at all times honoring constraints between the flights. The constraints are slightly different from the ones we had in the US map. The constraint is that no airplane can be on two routes at the same time. This implies that there is a constraint between flights F_1 and F_2 . Similarly, there is a constraint between the pairs $\langle F_2, F_3 \rangle$, $\langle F_3, F_5 \rangle$, $\langle F_4, F_5 \rangle$, $\langle F_1, F_5 \rangle$, $\langle F_2, F_5 \rangle$ and $\langle F_3, F_5 \rangle$. We assume that the turn around duration for F_1 (which is of a one hour duration) will end by 14:30 hours, which sounds reasonable. Thus, there is no constraint $\langle F_1, F_4 \rangle$.

You schedule planes the same way you do map coloring. An assignment is tried and constraints for all unassigned flights are checked to ensure that there is at least one airplane that can fly each flight. There is one important difference between the flight scheduling problem and the map coloring problem: we know that we can color maps with 4 colors⁶, but we do not know how many airplanes it is going to take to fly a schedule. Hence, we will need to try the flight scheduling problem with different number of airplanes.

Let us say we over-resource the map coloring problem with 7 colors instead of 4. A sample run yields 48 assignments (that is no backup was required) and 274 checks. With 6 colors, you get 48 assignments and 259 checks. If on the extreme end, you used only 3 colors, you could never color Texas.

Frequently the problem is over-constrained and there is no solution with available resources (like say coloring the United States with 2 or 3 colors). In those circumstances, constraints can be turned into preferences so that some regions will not be allowed to be adjacent to regions of the same color. Or we might have to allow some 'dead-hit' flights. And on top of preferences, we could layer beam-search or some other search that tries to minimize the penalty cumulated or maximize the number of constraints that are satisfied.

⁶The 4 color theorem.

Chapter 7

Graphical Models

Graphical models [Lau96, CGH97, Pea88, Jor98, Jen01] provide a pictorial representation of the way a joint probability distribution factorizes into a product of factors. They have been widely used in the area of computer vision, control theory, bioinformatics, communication, signal processing, sensor networks, neurovision, *etc.*

There is a relation between the conditional independence properties of a joint distribution and its factorization properties. Each of these properties can be represented graphically. This allows us to

1. organize complicated mathematics through graph based algorithms for calculation and computation and thus save complicated computations (in many cases, map the pictorial representation onto computations directly)
2. gain new insights into existing models; it is a standard practice to modify an existing model by addition/deletion of nodes or links and adopt it to the problem at hand
3. motivate new models by simply amending pictures

However, everything that you could do in graphical models using pictures could also be done without pictures, by grinding through all the mathematics while consistently using the innocuous-looking sum (7.1) and product (7.2) rules of probability

$$\Pr(X) = \sum_{y \in \mathcal{Y}} \Pr(X = x, Y = y) \quad (7.1)$$

$$\Pr(X = x, Y = y) = \Pr(X = x | Y = y) \Pr(Y = y) \quad (7.2)$$

where X and Y are discrete random variables, assuming values $x \in \mathcal{X} = \{x_1, x_2, \dots, x_m\}$ and $y \in \mathcal{Y} = \{y_1, y_2, \dots, y_n\}$ respectively. A combination of the two rules yields the Bayes theorem:

$$\begin{aligned} \Pr(X = x_i, Y = y_j) &= \frac{\Pr(X = x_i | Y = y_j) \Pr(Y = y_j)}{\Pr(X = x_i)} \\ &= \frac{\Pr(X = x_i | Y = y_j) \Pr(Y = y_j)}{\sum_{y_j \in \mathcal{Y}} \Pr(X = x_i | Y = y_j) \Pr(Y = y_j)} \quad (7.3) \end{aligned}$$

The two main kinds of graphical models are *directed* and *undirected* models. The problems we will address in graphical models include

1. *Inference*: Broadly, there are two inference techniques for graphical models, *viz.*, exact and approximate inference. Exact inference is appropriate if the graphic is a tree, since it is a linear time algorithm. But for complex graphical models, exact inference may or may not be appropriate, since exact algorithms could be very slow. In such cases, approximate inference schemes are often resorted to. Markov chain monte carlo (which is exact if there were an infinite amount of computing resources and approximate otherwise) and variational inference (by approximating the analytical form for the posterior distribution) are two popular techniques. While variational techniques scale better, their other strengths and weaknesses are complementary to those of MCMC. An often adopted stepping stone for explaining variational inference is the expectation maximization algorithm (EM) and we will take the same route.

2. *Learning*:

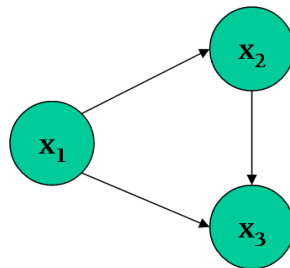
7.1 Semantics of Graphical Models

We will first discuss the semantics of graphical models, both directed and undirected. In the sections that follow, we will discuss the computational aspects of graphical models - in particular, inferencing and learning techniques.

7.1.1 Directed Graphical Models

We will start with the example of a directed graphical model. Consider an arbitrary joint distribution $\Pr(X_1 = x_1, X_2 = x_2, X_3 = x_3)$ over three discrete random variables X_1, X_2 and X_3 that assume values $x_1 \in \mathcal{X}_1, x_2 \in \mathcal{X}_2$ and $x_3 \in \mathcal{X}_3$ respectively. Denoting¹ $\Pr(X_i = x_i)$ by $p(x_i)$ and $\Pr(X_1 = x_1, X_2 = x_2, X_3 = x_3)$ by $p(x_1, x_2, x_3)$ and applying the product rule of probability successively, we obtain

¹As a convention, we will use capital letters to denote random variables and lower case letters to denote their realizations.



$$p(x_1, x_2, x_3) = p(x_1)p(x_2, x_3 \mid x_1) = p(x_1)p(x_2 \mid x_1)p(x_3 \mid x_1, x_2) \quad (7.4)$$

The successive application of the product rule is often termed as the *chain rule*. We should note that this rule applies even if x_1 , x_2 and x_3 happen to be continuous random variables (in which case p is a density function) or vectors of random variables. We should note that this factorization is quite non-symmetrical in the three random variables. This factorization can be represented in the form of the following directed graph: There is one node for each variable. We draw a directed edge between every conditioning variable (*i.e.* the corresponding node) to the conditioned variable. The way to go from a graph to the factorization of the joint distribution is to write down the product of the conditional distribution of every node (*i.e.*, the corresponding variable) conditioned on its parents within the graph. In the example above, x_1 had no parent, and therefore the term corresponds to its conditional $p(x_1)$ turned out to be its marginal distribution.

The factorization in the last example holds for any joint distribution over any three variables and the graph is therefore uninformative. In fact, any completely connected graph will be uninformative, as we will soon see. What interests us in graphical models is not the *presence* of edges but rather, the *absence* of edges. Since the graph in the previous example had no missing edges, it was uninteresting.

Definition 43 Let $\mathcal{R} = \{X_1, X_2, \dots, X_n\}$ be a set of random variables, with each X_i ($1 \leq i \leq n$) assuming values $x_i \in \mathcal{X}_i$. Let $\mathcal{X}_S = \{X_i \mid i \in S\}$ where $S \subseteq \{1, 2, \dots, n\}$. Let $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$ be a directed acyclic graph with vertices $\mathcal{V} = \{1, 2, \dots, n\}$ and $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ such that each edge $e = (i, j) \in \mathcal{E}$ is a directed edge. We will assume a one to one correspondence between the set of variables \mathcal{R} and the vertex set \mathcal{V} ; vertex i will correspond to random variable X_i . Let π_i be the set of vertices from which there is edge incident on vertex i . That is, $\pi_i = \{j \mid j \in \mathcal{V}, (j, i) \in \mathcal{E}\}$. Then, the family $\mathcal{F}(\mathcal{G})$ of joint distributions associated with the DAG² \mathcal{G} is specified by the factorization induced by \mathcal{G} as follows:

²As we saw earlier, the family of probability distributions specified by the related formalism of undirected graphical models is somewhat different.

$$\mathcal{F}(\mathcal{G}) = \left\{ p(\mathbf{x}) \left| p(\mathbf{x}) = \prod_{i=1}^n p(x_i | \mathbf{x}_{\pi_i}), p(x_i | \mathbf{x}_{\pi_i}) \forall 1 \leq i \leq n \geq 0, \sum_{x_i \in \mathcal{X}_i} p(x_i | \mathbf{x}_{\pi_i}) = 1 \right. \right\} \quad (7.5)$$

where, \mathbf{x} denotes the vector of values $[x_1, x_2, \dots, x_n]$ and x_i is the value assumed by random variable X_i and \mathbf{x}_{π_i} denotes the vector of values from \mathbf{x} , composed from positions in π_i .

For notational convenience with directed acyclic graphs, it is a common practice to assume a topological ordering on the indices of vertices in \mathcal{V} so that $\pi_i \subseteq \mu_{i-1} = \{1, 2, \dots, i-1\}$. Note that, by the chain rule, the following factorization always holds:

$$p(\mathbf{x}) = \prod_{i=1}^n p(x_i | \mathbf{x}_{\mu_{i-1}}) \quad (7.6)$$

Making use of the sum rule, in conjunction with (7.5), for any $p \in \mathcal{F}(\mathcal{G})$, we have

$$\begin{aligned} p(\mathbf{x}_{\mu_i}) &= \sum_{x_{i+1} \in \mathcal{X}_{i+1}, \dots, x_n \in \mathcal{X}_n} p(\mathbf{x}) \\ &= \sum_{x_{i+1} \in \mathcal{X}_{i+1}, \dots, x_n \in \mathcal{X}_n} p(x_1) p(x_2 | \mathbf{x}_{\pi_2}) \dots p(x_i | \mathbf{x}_{\pi_i}) \end{aligned} \quad (7.7)$$

Since the vertex indices are topologically ordered, it can be proved using the principle of induction (working backwards from n) on the basis of the sum rule in (7.7), that for any $p \in \mathcal{F}(\mathcal{G})$:

$$p(\mathbf{x}_{\mu_i}) = \prod_{j=1}^{i-1} p(x_j | \mathbf{x}_{\pi_j}) \quad (7.8)$$

Contrasting (7.6) against (7.5), we can think of the set of probability distribution $\mathcal{F}(\mathcal{G})$ as a sort of restricted class of distributions that arises from throwing away some of the dependencies. In particular, if $p \in \mathcal{F}(\mathcal{G})$ then

$$p(x_i | \mathbf{x}_{\mu_{i-1}}) = p(x_i | \mathbf{x}_{\pi_i})$$

that is, X_i is independent of $\mathbf{X}_{\mu_{i-1}}$, given \mathbf{X}_{π_i} . The independence is denoted by: $X_i \perp \mathbf{X}_{\mu_{i-1}} | \mathbf{X}_{\pi_i}$. This leads us to another approach to defining the class of probability distributions based on a DAG \mathcal{G} .

Definition 44 Let $\mathcal{R} = \{X_1, X_2, \dots, X_n\}$ be a set of random variables, with each X_i ($1 \leq i \leq n$) assuming values $x_i \in \mathcal{X}_i$. Let $\mathbf{X}_{\mathcal{S}} = \{X_i \mid i \in \mathcal{S}\}$ where $\mathcal{S} \subseteq \{1, 2, \dots, n\}$. Let $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$ be a directed acyclic graph with vertices $\mathcal{V} = \{1, 2, \dots, n\}$ and $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ such that each edge $e = (i, j)$ is a directed edge. Let $\pi_i = \{j \mid j \in \mathcal{V}, (j, i) \in \mathcal{E}\}$. Then, the family $\mathcal{C}(\mathcal{G})$ of joint distributions associated with the DAG \mathcal{G} is specified by the conditional independence induced by \mathcal{G} as follows:

$$\mathcal{C}(\mathcal{G}) = \left\{ p(\mathbf{x}) \mid X_i \perp \mathbf{X}_{\mu_{i-1}} \mid \mathbf{X}_{\pi_i} \quad \forall 1 \leq i \leq n, \sum_{\mathbf{x}} p(\mathbf{x}) = 1 \right\} \quad (7.9)$$

We will next show that the class $\mathcal{F}(\mathcal{G})$ defined in terms of factorizations is equivalent to the class $\mathcal{C}(\mathcal{G})$ defined in terms of independences. This is called the Hammersley Clifford theorem.

Theorem 87 The sets $\mathcal{F}(\mathcal{G})$ and $\mathcal{C}(\mathcal{G})$ are equal. That is $p \in \mathcal{F}(\mathcal{G})$ iff $p \in \mathcal{C}(\mathcal{G})$

Proof: \Leftarrow : We will first prove that $\mathcal{F}(\mathcal{G}) \subseteq \mathcal{C}(\mathcal{G})$. Let $p \in \mathcal{F}(\mathcal{G})$. We will prove that $p \in \mathcal{C}(\mathcal{G})$, that is, $p(x_i \mid \mathbf{x}_{\mu_{i-1}}, \mathbf{x}_{\pi_i}) = p(x_i \mid \mathbf{x}_{\pi_i})$. This trivially holds for $i = 1$, since $\mathbf{x}_{\pi_i} = \emptyset$. For $i = 2$:

$$p(x_1, x_2) = p(x_1)p(x_1 \mid x_2) = p(x_1)p(x_1 \mid \mathbf{x}_{\pi_2})$$

where, the first equality follows by chain rule, whereas the second equality follows by virtue of (7.8). Consequently,

$$p(x_1 \mid x_2) = p(x_1 \mid \mathbf{x}_{\pi_2})$$

Assume that $p(x_i \mid \mathbf{x}_{\mu_{i-1}}) = p(x_i \mid \mathbf{x}_{\pi_i})$ for $i \leq k$. For $i = k + 1$, it follows from chain rule and from (7.8) that

$$p(\mathbf{x}_{\mu_{k+1}}) = \prod_{i=1}^{k+1} p(x_i \mid \mathbf{x}_{\mu_{i-1}}) = \prod_{i=1}^{k+1} p(x_i \mid \mathbf{x}_{\pi_i})$$

Making use of the induction assumption for $i \leq k$ in the equation above, we can derive that

$$p(x_k \mid \mathbf{x}_{\mu_{k-1}}) = p(x_k \mid \mathbf{x}_{\pi_k})$$

By induction on i , we obtain that $p(x_i \mid \mathbf{x}_{\mu_{i-1}}) = p(x_i \mid \mathbf{x}_{\pi_i})$ for all i . That is, $p \in \mathcal{C}(\mathcal{G})$. Since this holds for any $p \in \mathcal{F}(\mathcal{G})$, we must have that $\mathcal{F}(\mathcal{G}) \subseteq \mathcal{C}(\mathcal{G})$.

\Rightarrow : Next we prove that $\mathcal{C}(\mathcal{G}) \subseteq \mathcal{F}(\mathcal{G})$. Let $p' \in \mathcal{C}(\mathcal{G})$ satisfy the conditional independence assertions. That is, for any $1 \leq i \leq n$, $p'(x_i \mid \mathbf{x}_{\mu_{i-1}}) = p'(x_i \mid \mathbf{x}_{\pi_i})$. Then by chain rule, we must have:

$$p'(\mathbf{x}_{\mu_n}) = \prod_{i=1}^n p'(x_i \mid \mathbf{x}_{\mu_{i-1}}) = \prod_{i=1}^{k+1} p'(x_i \mid \mathbf{x}_{\pi_i})$$

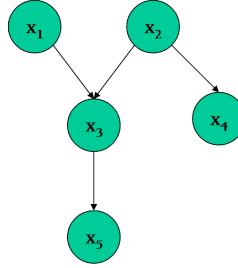


Figure 7.1: A directed graphical model.

which proves that $p' \in \mathcal{F}(\mathcal{G})$ and consequently that $\mathcal{C}(\mathcal{G}) \subseteq \mathcal{F}(\mathcal{G})$ \square

As an example, we will discuss the directed graphical model, as shown in Figure 7.1. Based on theorem 87, the following conclusions can be drawn from the graphical representation of a family of distributions represented by Figure 7.1.

1. Given the value of X_3 , the values of X_1, X_2 and X_4 will be completely uninformative about the value of X_5 . That is, $(X_5 \perp \{X_1, X_2, X_4\} \mid \{X_3\})$. Similarly, given the value of X_2 , the values of X_1 and X_3 will be completely uninformative about the value of X_4 . That is, $(X_4 \perp \{X_1, X_3\} \mid \{X_2\})$.
2. Secondly, since $p \in \mathcal{F}(\mathcal{G})$, we have

$$p(x_1, x_2, x_3, x_4, x_5) = p(x_1)p(x_2)p(x_3 \mid x_1, x_2)p(x_4 \mid x_2)p(x_5 \mid x_3)$$

What about other independence assertions? Is X_5 independent of X_4 given X_2 ? Is X_3 independent of X_4 , given X_2 ? The answer to both these questions happens to be yes. And these could be derived using either of the equivalent definitions of graphical models. In fact, some such additional conditional independence assertions can always follow from the two equivalent definitions of graphical models. Before delving into these properties, we will define an important concept called d-separation, introduced by Pearl [Pea88].

Definition 45 A set of nodes \mathcal{A} in a directed acyclic graph \mathcal{G} is d-separated from a set of nodes \mathcal{B} by a set of nodes \mathcal{C} , iff **every** undirected path from a vertex $A \in \mathcal{A}$ to a vertex $B \in \mathcal{B}$ is ‘blocked’. An undirected path between A and B is blocked by a node C either (i) if $C \in \mathcal{C}$ and both the edges (which might be the same) on the path through C are directed away from C (C is then called a tail-to-tail node) or (ii) if $C \in \mathcal{C}$ and of the two edges (which might be the same) on the path through C , one is directed toward C while the other is directed away from C (C is then called a head-to-tail node) or (iii) if $C \notin \mathcal{C}$ and both the edges (which might be the same) on the path through C are directed toward C (C is then called a head-to-head node).

In Figure 7.1, node X_2 blocks the only path between X_3 and X_4 , node X_3 blocks the path between X_2 and X_5 . Whereas, node X_3 does not block the

path between X_1 and X_2 . Consequently, $\{X_3\}$ and $\{X_4\}$ are d-separated by $\{X_2\}$, while $\{X_2\}$ and $\{X_5\}$ are d-separated by $\{X_3\}$. However, $\{X_1\}$ and $\{X_2\}$ are not d-separated by $\{X_3\}$, since X_3 is a head-to-head node. X_3 does not d-separate X_1 and X_2 even though it separates them. Thus, not every pair of (graph) separated nodes in the graph need be d-separated. We next define a family of probability distribution that have independences characterized by d-separation.

Definition 46 *The set of probability distributions $\mathcal{D}(\mathcal{G})$ for a DAG \mathcal{G} is defined as follows:*

$$\mathcal{D}(\mathcal{G}) = \{p(\mathbf{x}) \mid \mathbf{X}_{\mathcal{A}} \perp \mathbf{X}_{\mathcal{B}} \mid \mathbf{X}_{\mathcal{C}}, \text{ whenever } \mathcal{A} \text{ and } \mathcal{B} \text{ are d-separated by } \mathcal{C}\} \quad (7.10)$$

It can be proved that the notion of conditional independence is equivalent to the notion of d-separation in DAGs. That is,

Theorem 88 *For any directed acyclic graph \mathcal{G} , $\mathcal{D}(\mathcal{G}) = \mathcal{C}(\mathcal{G}) = \mathcal{F}(\mathcal{G})$.*

Thus, in Figure 7.1. $\{X_3\} \perp \{X_4\} \mid \{X_2\}$ and $\{X_2\} \perp \{X_5\} \mid \{X_3\}$. Whereas, $\{X_1\} \not\perp \{X_2\} \mid \{X_3\}$. We could think of X_1 and X_2 as completing explanations for X_3 . Thus, given a value of X_3 , any value of X_1 will, to some extent ‘explain away’ the value of X_3 , thus withdrawing the independence of X_2 . In terms of a real life example, if X_1 , X_2 and X_3 are discrete random variables corresponding to ‘the color of light’, ‘the surface color’ and ‘the image color’ respectively, then, given the value of X_3 (image color), any value of X_1 (color of light) will explain away the color of the image, thus constraining the values that X_3 (surface color) might take. On the other hand, $\{X_1\} \perp \{X_2\} \mid \{\}$. What about the independence of X_1 and X_2 given X_5 ? The path X_1, X_3, X_5, X_3, X_2 involves a head-to-head node X_5 and therefore, $X_1 \not\perp X_2 \mid \{X_5\}$. The Bayes ball algorithm provides a convenient algorithmic way for deciding if $\mathbf{X}_{\mathcal{A}} \perp \mathbf{X}_{\mathcal{B}} \mid \mathbf{X}_{\mathcal{C}}$, by using the d-separation property.

Bayesian Networks and Logic

The logical component of Bayesian networks essentially corresponds to a propositional logic program. This has already been observed by Haddawy [1994] and Langley [1995]. Langley, for instance, does not represent Bayesian networks graphically but rather uses the notation of propositional definite clause programs. Consider the following program. This program encodes the structure of the blood type Bayesian network in Figure 7.2. Observe that the random variables in this notation correspond to logical atoms. Furthermore, the direct influence relation in the Bayesian network corresponds to the immediate consequence operator.

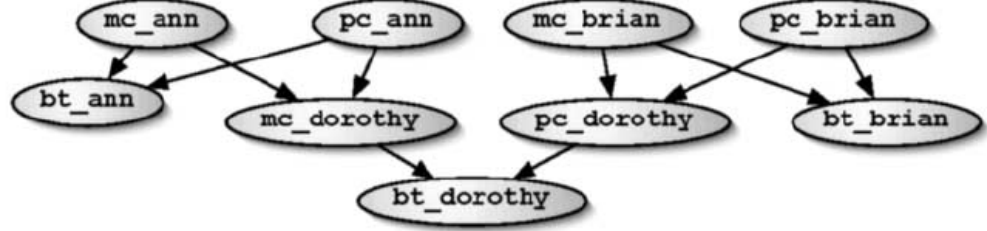


Figure 7.2: The graphical structure of a Bayesian network modeling the inheritance of blood types within a particular family.

$PC(ann)$.

$MC(ann)$.

$MC(dorothy) : -MC(ann), PC(ann)$.

$BT(ann) : -MC(ann), PC(ann)$.

$BT(dorothy) : -MC(dorothy), PC(dorothy)$.

$PC(brian)$.

$MC(brian)$.

$PC(dorothy) : -MC(brian), PC(brian)$.

$BT(brian) : -MC(brian), PC(brian)$.

7.1.2 Undirected Graphical Models

We will move on to an undirected graphical models (also known as Markov Random fields) primer, while drawing parallels with the directed counterpart. An undirected graph \mathcal{G} is a tuple $\langle \mathcal{V}, \mathcal{E} \rangle$ where $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ and such that each edge $e = (i, j) \in \mathcal{E}$ is a directed edge. While the conditional independence property for directed graphical models was tricky (involving concepts such as d-separation, *etc.*), the conditional independence property for undirected models is easier to state. On the other hand, the factorization property for directed graphical models simply involved local conditional probabilities as factors. It is however not as simple with undirected graphical models. Taking the easier route, we will first define the conditional independence property for undirected graphical models. Toward that, we introduce the notion of graph separation.

Definition 47 Given an undirected graph $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$, and $\mathcal{A}, \mathcal{B}, \mathcal{C} \subseteq \mathcal{V}$, we say that \mathcal{C} separates \mathcal{A} from \mathcal{B} in \mathcal{G} if every path from any node $A \in \mathcal{A}$ to any node $B \in \mathcal{B}$ passes through some node $C \in \mathcal{C}$. \mathcal{C} is also called a separator or a vertex cut set in \mathcal{G} .

In Figure 7.3, the set $\mathcal{A} = \{X_1, X_7, X_8\}$ is separated from $\mathcal{B} = \{X_3, X_4, X_5\}$ by the (vertex cut) set $\mathcal{C} = \{X_2, X_6\}$. It is easy to see that separation is symmetric in \mathcal{A} and \mathcal{B} . This simple notion of separation gives rise to a conditional independence assertion for undirected graphs. A random vector $\mathbf{X}_{\mathcal{C}}$ for $\mathcal{C} \subseteq \mathcal{V}$

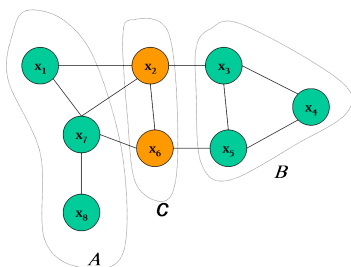


Figure 7.3: A directed graphical model.

is said to be *markov* with respect to a graph \mathcal{G} if $\mathbf{X}_A \perp \mathbf{X}_B \mid \mathbf{X}_C$ whenever \mathcal{C} separates \mathcal{A} from \mathcal{B} . That is, the random variables corresponding to the vertex cut set acts like a mediator between the values assumed by variables in \mathcal{A} and \mathcal{B} so that the variables in \mathcal{A} and \mathcal{B} are independent of each other, if we knew the values of variables in \mathcal{C} . It is straightforward to develop a ‘reachability’ algorithm (as with the bayes ball algorithm) for undirected graphs, to assess conditional independence assumptions. Based on the definition of markov random vector, we next define the family $\mathcal{M}(\mathcal{G})$ of distributions associated with an undirected graph \mathcal{G} .

Definition 48 Let $\mathcal{R} = \{X_1, X_2, \dots, X_n\}$ be a set of random variables, with each X_i ($1 \leq i \leq n$) assuming values $x_i \in \mathcal{X}_i$. Let $\mathbf{X}_S = \{X_i \mid i \in S\}$ where $S \subseteq \{1, 2, \dots, n\}$. Let $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$ be an undirected graph with vertices $\mathcal{V} = \{1, 2, \dots, n\}$ and $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ such that each edge $e = (i, j) \in \mathcal{E}$ is an undirected edge. Then, the family $\mathcal{M}(\mathcal{G})$ of joint distributions associated with \mathcal{G} is specified as follows:

$$\mathcal{M}(\mathcal{G}) = \{p(\mathbf{x}) \mid \mathbf{X}_A \perp \mathbf{X}_B \mid \mathbf{X}_C \forall \mathcal{A}, \mathcal{B}, \mathcal{C} \subseteq \mathcal{V}, \text{ whenever } \mathcal{C} \text{ separates } \mathcal{A} \text{ from } \mathcal{B}\} \quad (7.11)$$

As in the case of directed models, there is another family of probability distributions that could be specified for a given undirected graph \mathcal{G} based on a factorization assertion. The main difference is that, while the factorization for DAGs was obtained in terms of local conditional probabilities or local marginals, it turns out that this factorization is not possible for a general undirected graph (specifically when it has a cycle). Instead there is another notion of ‘local’ for undirected graphs – there should be no function involving any two variables X_i and X_j where $(i, j) \notin \mathcal{E}$ (otherwise, such a term will not break further, prohibiting assertions about conditional independences). Instead, we will have a function $\phi_{\mathcal{C}}(\mathbf{X}_{\mathcal{C}})$ for every clique $\mathcal{C} \subseteq \mathcal{V}$, since a clique is a subset of vertices that all ‘talk to’ one another. The most general version of factorization will be one for which there is a function corresponding to each *maximal clique*; all other

factorizations involving factors over smaller cliques will be specialized versions. These functions will be referred to as *compatibility* or *potential* functions. A *potential* or *compatibility* function on a clique \mathcal{C} is a non-negative real valued function $\phi_{\mathcal{C}}(\mathbf{x}_{\mathcal{C}})$ defined over all instantiations $\mathbf{x} \in \mathcal{X}_1 \times \mathcal{X}_2 \times \dots \times \mathcal{X}_n$ of \mathbf{X} . Other than these restrictions, the potential function can be arbitrary.

Definition 49 Let $\mathcal{R} = \{X_1, X_2, \dots, X_n\}$ be a set of random variables, with each X_i ($1 \leq i \leq n$) assuming values $x_i \in \mathcal{X}_i$. Let $\mathbf{X}_{\mathcal{S}} = \{X_i \mid i \in \mathcal{S}\}$ where $\mathcal{S} \subseteq \{1, 2, \dots, n\}$. Let $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$ be an undirected graph with vertices $\mathcal{V} = \{1, 2, \dots, n\}$ and $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ such that each edge $e = (i, j) \in \mathcal{E}$ is an undirected edge. Then, the family $\mathcal{M}(\mathcal{G})$ of joint distributions associated with \mathcal{G} is specified as follows:

$$\mathcal{F}(\mathcal{G}) = \left\{ p(\mathbf{x}) \mid p(\mathbf{x}) = \frac{1}{Z} \prod_{\mathcal{C} \in \Pi} \phi_{\mathcal{C}}(\mathbf{x}_{\mathcal{C}}), \text{ such that } \phi_{\mathcal{C}}(\mathbf{x}_{\mathcal{C}}) \in \mathbb{R}^+, \forall \mathcal{C} \in \Pi \right\} \quad (7.12)$$

where, Π is the set of all cliques in \mathcal{G} , \mathbf{x} denotes the vector of values $[x_1, x_2, \dots, x_n]$, $x_i \in \mathcal{X}_i$ is the value assumed by random variable X_i and where each $\phi_{\mathcal{C}}$ is a potential function defined over the clique $\mathcal{C} \subseteq \mathcal{V}$. Without loss of generality, we can assume that Π is the set of maximal cliques in \mathcal{G} . The normalization constant Z is called the partition function and is given by

$$Z = \sum_{\mathbf{x}_1 \in \mathcal{X}_1, \dots, \mathbf{x}_n \in \mathcal{X}_n} \prod_{\mathcal{C} \in \Pi} \phi_{\mathcal{C}}(\mathbf{x}_{\mathcal{C}}) \quad (7.13)$$

The potential functions are typically represented as tables – each row listing a unique assignment of values to the random variables in the clique and the corresponding potential. Thus, the value $\phi_{\mathcal{C}}(\mathbf{x}_{\mathcal{C}})$ can be obtained by a simple table lookup.

The form of the potential function can be chosen based on the particular application at hand. For instance, the clique potential can be decomposed into a product of potentials defined over each edge of the graph. When the domain of each random variable is the same, the form of the potential can be chosen to either encourage or discourage similar configurations (such as similar disease infection for patients who are related) at adjacent nodes. The potential function is often interpreted as an energy function in the modeling of crystals, protein folding, *etc.*, where a minimum energy configuration is desirable. Frequently, the energy function is assumed to have the form $\phi_{\mathcal{C}}(\mathbf{X}_{\mathcal{C}}) = \exp(-\theta_{\mathcal{C}}(\mathbf{X}_{\mathcal{C}}))$, which leads to the factorization as an exponential form distribution $p(\mathbf{x}) = \exp(-\sum_{\mathcal{C} \in \Pi} \theta_{\mathcal{C}}(\mathbf{x}_{\mathcal{C}}) - \log Z)$. The quantities $\theta_{\mathcal{C}}(\mathbf{X}_{\mathcal{C}})$ are called sufficient statistics.

From our discussion on the equivalence of directed and undirected trees in terms of conditional independencies, we may be tempted to conclude that the factors for the undirected tree can be the local conditional probabilities. This is easily established if we prove that for strictly positive distributions, the definition of an undirected graphical model in terms of conditional independencies is equivalent to the definition in terms of factorization, that is, $\mathcal{M}(\mathcal{G}) = \mathcal{F}(\mathcal{G})$.

Theorem 89 *For strictly positive distributions, $\mathcal{M}(\mathcal{G}) = \mathcal{F}(\mathcal{G})$. That is, $\mathcal{M}(\mathcal{G}) \cap \mathcal{D}^+ = \mathcal{F}(\mathcal{G}) \cap \mathcal{D}^+$, where, $\mathcal{D}^+ = \{p(\mathbf{x}) \mid p(\mathbf{x}) > 0, \forall \mathbf{x} \in \mathcal{X}_1 \times \mathcal{X}_2 \dots \times \mathcal{X}_n\}$.*

Proof: The proof that $\mathcal{M}(\mathcal{G}) \subseteq \mathcal{F}(\mathcal{G})$ is a bit involved and requires Mobius inversion. On the other hand, that $\mathcal{F}(\mathcal{G}) \subseteq \mathcal{M}(\mathcal{G})$ can be shown as follows. For any given $\mathcal{A}, \mathcal{B} \subseteq \mathcal{V}$ that are separated by $\mathcal{C} \subseteq \mathcal{V}$, consider the partitions $\mathcal{P}_1, \mathcal{P}_2$ and \mathcal{P}_3 of Π :

$$\begin{aligned}\mathcal{P}_1 &= \{\mathcal{K} \mid \mathcal{K} \in \Pi \text{ and } \mathcal{K} \subseteq \mathcal{A} \cap \mathcal{C} \text{ and } \mathcal{K} \not\subseteq \mathcal{C}\} \\ \mathcal{P}_2 &= \{\mathcal{K} \mid \mathcal{K} \in \Pi \text{ and } \mathcal{K} \subseteq \mathcal{B} \cap \mathcal{C} \text{ and } \mathcal{K} \not\subseteq \mathcal{C}\} \\ \mathcal{P}_3 &= \{\mathcal{K} \mid \mathcal{K} \in \Pi \text{ and } \mathcal{K} \subseteq \mathcal{C}\}\end{aligned}$$

Now, $p(\mathbf{x})$ can be factorized into factors involving cliques in $\mathcal{P}_1, \mathcal{P}_2$ and \mathcal{P}_3 . Consequently,

$$\frac{p(\mathbf{x}_A, \mathbf{x}_B, \mathbf{x}_C)}{p(\mathbf{x}_B, \mathbf{x}_C)} = \frac{\prod_{\mathcal{K} \in \mathcal{P}_1} \phi_{\mathcal{K}}(\mathbf{x}_{\mathcal{K}}) \prod_{\mathcal{K} \in \mathcal{P}_2} \phi_{\mathcal{K}}(\mathbf{x}_{\mathcal{K}}) \prod_{\mathcal{K} \in \mathcal{P}_3} \phi_{\mathcal{K}}(\mathbf{x}_{\mathcal{K}})}{\sum_{\mathbf{x}_A} \prod_{\mathcal{K} \in \mathcal{P}_1} \phi_{\mathcal{K}}(\mathbf{x}_{\mathcal{K}}) \prod_{\mathcal{K} \in \mathcal{P}_2} \phi_{\mathcal{K}}(\mathbf{x}_{\mathcal{K}}) \prod_{\mathcal{K} \in \mathcal{P}_3} \phi_{\mathcal{K}}(\mathbf{x}_{\mathcal{K}})} = \frac{\prod_{\mathcal{K} \subseteq \mathcal{A} \cup \mathcal{C}} \phi_{\mathcal{K}}(\mathbf{x}_{\mathcal{K}})}{\sum_{\mathbf{x}_A} \prod_{\mathcal{K} \subseteq \mathcal{A} \cup \mathcal{C}} \phi_{\mathcal{K}}(\mathbf{x}_{\mathcal{K}})} = p(\mathbf{x}_A \mid \mathbf{x}_C)$$

□

While conditional independence is useful for modeling purposes, factorization is more useful for computational purposes.

7.1.3 Comparison between directed and undirected graphical models

Is there any difference between the undirected and directed formalisms? Are they equally powerful? Or is one more powerful than the other. It turns out that there are families of probability distributions which can be represented using undirected models, whereas they have no directed counterparts. Figure 7.4 shows one such example. Imagine that random variables X_1 and X_3 represent hubs in a computer network, while X_2 and X_4 represent computers in the network. Computers do not interact directly, but only through hubs. Similarly, hubs interact only through computers. This leads to two independencies: (i) conditioned on X_1 and X_3 (the hubs), nodes X_2 and X_4 (the computers) become independent and (ii) conditioned on X_2 and X_4 , nodes X_1 and X_3 become independent. However, with a directed acyclic graph on four nodes, we will always have some head-to-head node and therefore, it is impossible to simultaneously satisfy both conditions (i) and (ii) using a DAG. Larger bipartite graphs

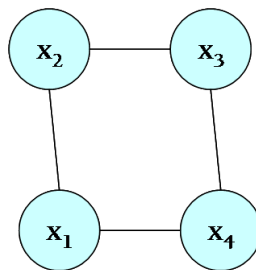


Figure 7.4: An undirected graphical model which has no equivalent directed model.

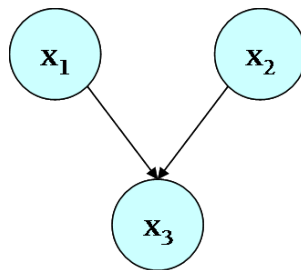


Figure 7.5: A directed graphical model which has no equivalent undirected model.

will have similar conditional independence assertions, which are inexpressible through DAGs.

Similarly, there are families of probability distributions which can be represented using directed models, whereas they have no undirected counterparts. Figure 7.5 shows one such example and corresponds to the ‘explaining away’ phenomenon, which we discussed earlier in this chapter. The node X_3 is blocked if not observed (so that $\{X_1\} \perp \{X_2\} \mid \emptyset$), whereas it is unblocked if its value is known (so that $\{X_1\} \not\perp \{X_2\} \mid \{X_3\}$). Can you get this behaviour with an undirected graph? The answer is no. This is because, with an undirected graph, there is no way of getting dependence between X_1 and X_2 if they were a priori independent.

On the other hand, for graphical models such as markov chains, dropping the arrows on the edges preserves the independencies, yielding an equivalent undirected graphical model. Similarly, directed trees are fundamentally no different from undirected trees.

An important point to note is that it is the absence of edges that characterizes a graphical model. For any graphical model, it is possible that the compatibility functions (or local conditional probabilities) assume a very special form so that there are more (conditional) independencies that hold than

what is indicated by the graphical model (which means that some of the edges in the graph could be redundant).

7.2 Inference

In this section, we discuss the problem of determining the marginal distribution $p(\mathbf{x}_{\mathcal{A}})$, the conditional distribution $p(\mathbf{x}_{\mathcal{A}} \mid \mathbf{x}_{\mathcal{B}})$ and the partition function Z , given a graphical model $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$ and for any $\mathcal{A}, \mathcal{B} \subseteq \mathcal{V}$. We will assume that the conditional probability tables (for directed models) or the potential function table (for undirected models) are already known. The sum and product rules of probability yield the following formulae³ for the marginal, the conditional⁴ and the partition function⁵ respectively:

$$p(\mathbf{x}_{\mathcal{A}}) = \sum_{\mathbf{x}_{\mathcal{V} \setminus \mathcal{A}}} p(\mathbf{x})$$

$$p(\mathbf{x}_{\mathcal{A}} \mid \mathbf{x}_{\mathcal{O}}) = \frac{p(\mathbf{x}_{\mathcal{A}}, \mathbf{x}_{\mathcal{O}})}{p(\mathbf{x}_{\mathcal{O}})}$$

$$Z = \sum_{\mathbf{x} \in \mathcal{X}_1 \times \mathcal{X}_2 \times \dots \times \mathcal{X}_n} \prod_{C \in \Pi} \phi_C(\mathbf{x}_C)$$

All these problems are somewhat similar, in that they involve summation over a very high dimensional space. Computing any of these quantities will involve number of computations that are atleast exponential in the size of $\mathcal{V} \setminus \mathcal{A}$. This is not feasible in many practical applications, where the number of nodes will run into the hundreds or the thousands. As we will see, there is ample redundancy in these computations. We will briefly discuss a simple and pedagogical algorithm called the *elimination algorithm* that provides the intuition as to how the structure of the graph could be exploited to answer some of the questions listed above. More clever algorithms such as the *sum-product algorithm* that captures redundancies more efficiently will be discussed subsequently.

A problem fundamentally different from the three listed above is that of *maximum a posteriori optimization* - determining the mode of a conditional distribution.

$$\hat{\mathbf{x}}_{\mathcal{A}} = \operatorname{argmax}_{\mathbf{x}_{\mathcal{A}}} p(\mathbf{x}_{\mathcal{A}} \mid \mathbf{x}_{\mathcal{O}})$$

For discrete problems, this is an integer linear programming problem. For general graphs, you cannot do much better than a brute force, whereas, for special graphs (such as trees), this mode can be computed efficiently.

³For continuous valued random variables, you can expect the summation to be replaced by integration in each formula.

⁴The problem of computing conditionals is not fundamentally different from the problem of computing the marginals, since every conditional is simply the ratio of two marginals.

⁵The partition function needs to be computed for parameter estimation.

7.2.1 Elimination Algorithm

The elimination algorithm provides a systematic framework for optimizing computations by rearranging sums and products, an instance of which we saw in the proof of theorem 89. Consider the simple directed graph in Figure 7.1. Let us say each random variable takes values from the set $\{v_1, v_2, \dots, v_k\}$. Brute force computation of $p(x_1 | x_5) = \frac{p(x_1, x_5)}{p(x_5)}$ will involve $k \times k^3 = k^4$ computations for the numerator and $k \times k^4 = k^5$ computations for the denominator. To take into account conditioning, we introduce a new potential function $\delta(x_i, x'_i)$ which is defined as follows:

$$\delta(x_i, x'_i) = \begin{cases} 1 & \text{if } x_i = x'_i \\ 0 & \text{otherwise} \end{cases}$$

and simply write the conditional as

$$p(x_1 | x_5) = \frac{p(x_1, x_5)}{p(x_5)} = \sum_{x_2, x_3, x_4, x'_5} p(x'_5 | x_3) \delta(x_5, x'_5) p(x_3 | x_1, x_2) p(x_4 | x_2) p(x_2)$$

That is, whenever a variable is observed, you imagine that you have imposed the indicator function for that observation into the joint distribution. Given this simplification, we will focus on efficiently computing $p(x_1)$, assuming that the δ function will be slapped onto the corresponding factor while computing conditionals.

Using the structure of the graphical model (implicit in the topological ordering over the indices), we can rearrange the sums and products for $p(x_1 | x_5)$

$$p(x_1 | x_5) = \left(\sum_{x_2} p(x_2) \left(\sum_{x_3} p(x_3 | x_1, x_2) \left(\sum_{x_4} p(x_4 | x_2) \right) \left(\sum_{x'_5} p(x'_5 | x_3) \delta(x_5, x'_5) \right) \right) \right) \quad (7.14)$$

where brackets have been placed at appropriate places to denote domains of summation.

Analysing this computational structure inside-out,

1. We find two innermost factors to be common across all the summations, *viz.*,

$$m_{x_4}(x_2) = \sum_{x_4} p(x_4 | x_2)$$

$$m_{x_5}(x_3) = \sum_{x'_5} p(x'_5 | x_3) \delta(x_5, x'_5)$$

where m_{x_4} is a *message* function of x_2 and is obtained using $k \times k = k^2$ computations and m_{x_5} is a *message* function of x_3 and similarly obtained using k^2 computations.

2. Further, we can decipher from (7.14), the following message function m_{x_3} of x_1 and x_2 which can be computed using k^3 operations:

$$m_{x_3}(x_1, x_2) = \sum_{x_3} p(x_3|x_1, x_2)m_{x_4}(x_2)m_{x_5'}(x_3)$$

3. Putting all the messages together, we get the message function m_{x_2} of x_1 , computable with k^2 operations.

$$m_{x_2}(x_1) = \sum_{x_2} m_{x_3}(x_1, x_2)$$

Thus, the summation over the factorization can be expressed as a flow of message from one node to another; when the message from a node passes on to another, the former gets stripped or *eliminated*. In the step (1), nodes x_4 and x_5 got stripped. In step (2), node x_3 was stripped and finally, in step (3), x_2 got stripped. This yields an elimination ordering⁶ $[x_4, x_5, x_3, x_2]$. The order of computation is thus brought down from $O(k^5)$ to $O(\max(k^2, k^3)) = O(k^3)$ computations. While some researchers could argue that this may not be a substantial decrease, for larger graphs the gain in speed using this procedure is always substantial.

More formally, consider a root to leaf ordering \mathcal{I} of the nodes, where r is the root node (equivalently, an ordering that corresponds to leaf stripping). Figure 7.1 shows a numbering of the nodes corresponding to such an ordering. We will define as the current active set $\mathcal{A}(k)$, a set of indices of general potential functions. At each step of the algorithm the potential functions are of three different types: (i) some of the local conditional probabilities $p(x_i|\mathbf{x}_{\pi_i})$, (ii) some of the indicators $\delta(x_j, x'_j)$ of the observed nodes and (iii) some messages (*c.f.* page 400) generated so far. More formally, the active set of potentials is given by $\{\Psi_\alpha(\mathbf{x}_\alpha)\}_{\alpha \in \mathcal{A}(k)}$, with α being a generic index that ranges over sets of nodes. $\mathcal{A}^{(0)}$ is initialized as the set of all cliques that are associated with potential functions in the graphical model. For example, in Figure 7.1, $\mathcal{A}^{(0)} = \{\{1\}, \{2\}, \{3, 2, 1\}, \{2, 4\}, \{3, 5\}\}$. Then, $\mathcal{A}^{(k)}$ can be computed using the algorithm presented in Figure 7.6.

When the active set construction is complete, the desired conditional/marginal probability can be obtained as

$$p(x_r | \mathbf{x}_o) = \frac{\Psi_{\{r\}x_r}}{\sum_{x_r} \Psi_{\{r\}x_r}}$$

A flip side of the elimination algorithm is that it requires a ‘good’ elimination order to be first determined. The number of elimination orderings is obviously a large value of $(n - 1)!$, where n is the number of nodes. Finding a good elimination ordering is an NP hard problem and heuristics have been the only recourse. We will not discuss the elimination algorithm any further, but instead jump to the more efficient sum-product algorithm.

⁶Since either x_4 and x_5 may get stripped first, $[x_5, x_4, x_3, x_2]$ is also a valid elimination ordering.

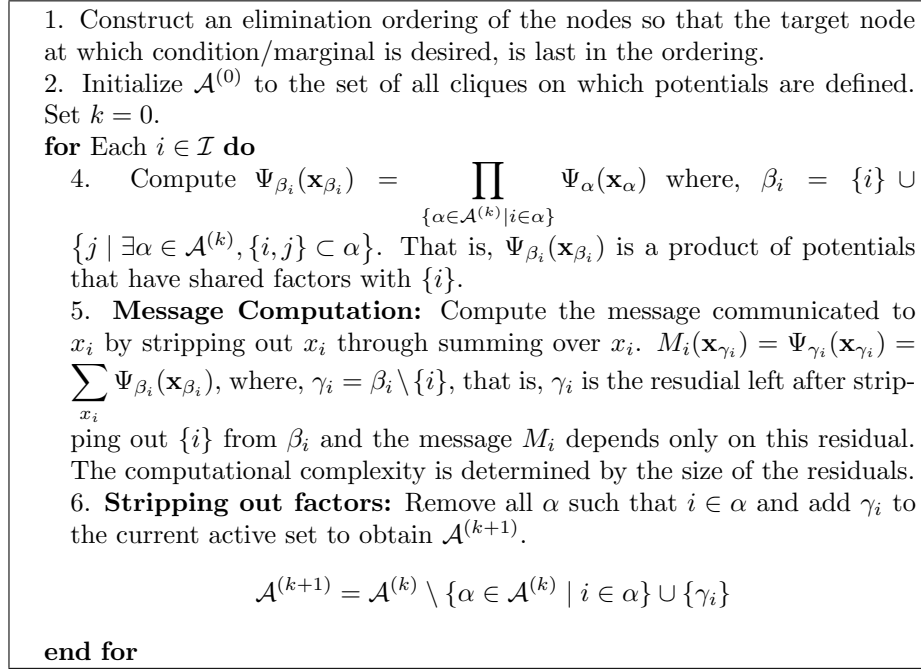


Figure 7.6: Procedure for constructing the active set, and the message at the desired target node t .

7.2.2 Sum-product Algorithm

The sum-product algorithm builds on the idea of ‘messages’ as motivated by the elimination algorithm. It involves local computations at nodes, to generate ‘messages’ which are related by nodes along their edges to their neighbors. This formalism enables simultaneous computation of marginals, conditionals as well as modes for several variable sets. This algorithm generalizes special algorithms such as viterbi, the forward-backward algorithm, kalman filtering, gaussian elimination as well as the fast fourier transform.

We will initially restrict our attention to undirected trees and will later generalize the algorithm. Figure 7.7 shows an example tree structured graphical model. Since the cliques consist of edges and individual nodes, the potential functions are basically either node potentials $\phi_p(x_p)$ or edge potentials $\phi_{p,q}(x_p, x_q)$. The joint distribution for the tree is

$$p(\mathbf{x}) = \frac{1}{Z} \prod_{p \in \mathcal{V}} \phi_p(x_p) \prod_{(p,q) \in \mathcal{E}} \phi_{p,q}(x_p, x_q) \quad (7.15)$$

Note that, all discussions that follow, hold equally well for directed trees, with the special parametrization $\phi_{p,q}(x_p, x_q) = p(x_p | x_q)$ if $x_p \in \pi_{x_q}$ and $\phi_p(x_p) =$

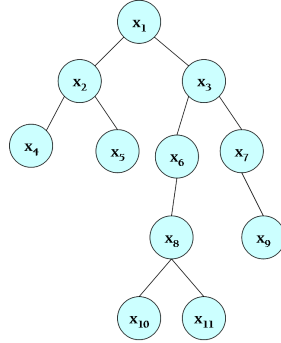


Figure 7.7: A tree graphical model.

$p(x_p)$.

We will deal with conditioning in a manner similar to the way we dealt with conditioning in the case of directed graphical models. For every observed variable $X_o = x_o$, we impose the indicator function $\delta(x_o, x'_o)$ onto $\phi_o(x'_o)$. Then, for a set of observed variables $\mathbf{X}_{\mathcal{O}}$, the conditional then takes the form:

$$p(\mathbf{x} | \mathbf{x}_{\mathcal{O}}) = \frac{1}{Z_{\mathcal{O}}} \prod_{p \in \mathcal{V} \setminus \mathcal{O}} \phi_p(x_p) \prod_{o \in \mathcal{O}} \phi_o(x'_o) \delta(x_o, x'_o) \prod_{(p,q) \in \mathcal{E}} \phi_{p,q}(x_p, x_q) \quad (7.16)$$

Thus, modifying the compatibility functions appropriately reduces the conditional problem to an instance of the base problem.

The crux of the sum-product algorithm is the following observation on the application of the leaf stripping and message construction procedure in Figure 7.6 to a tree-structured graphical model.

Theorem 90 *When a node i is eliminated, $\gamma_i = \{p\}$, where p is the unique parent of node i . This means, we can write the message as $M_{i \rightarrow p}$.*

Proof Sketch: This can be proved by induction on the number of steps in the leaf-stripping (elimination) order. You will need to show that at every step, $\beta_i = \{i, p\}$, where p is the unique parent of i . Consequently, $\gamma_i = \{p\}$ and we can derive the (recursive) expression for $M_{i \rightarrow p}(x_p)$ as

$$M_{i \rightarrow p}(x_p) = \Psi_{x_p}(x_p) = \underbrace{\sum_{x_i} \phi_i(x_i) \phi_{i,p}(x_i, x_p)}_{SUM} \underbrace{\prod_{q \in \mathcal{N}(i) \setminus p} M_{q \rightarrow i}(x_i)}_{PRODUCT} \quad (7.17)$$

Note that the proof of this statement will again involve induction. \square

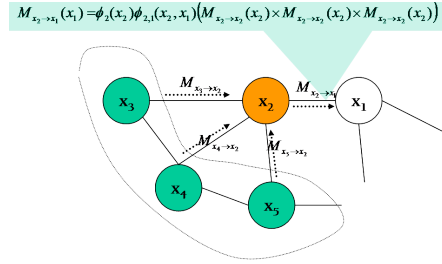


Figure 7.8: An illustration of (7.17) on a part of a tree structured (undirected) graphical model.

Figure 7.8 illustrates an application of (7.17) on a part of a tree structured graphical model.

Theorem 90 provides a very useful observation; the message is not a function of the vector of all nodes eliminated so far (which could have assumed k^m possible values for m eliminated nodes having k possible values each), but is instead is a function only of the child from which the message is being passed. This allows us to move from the elimination algorithm to the sum-product algorithm. In particular, the parts underlined as ‘SUM’ and ‘PRODUCT’ in (7.17) form the basis of the sum-product algorithm. In the sum-product algorithm, at every time step, the computation in (7.17) is performed at every node; each node does local computations and passes ‘updated’ messages to each of its neighbors. In this way, the computations are not tied to any particular elimination ordering.

Theorem 91 Consider an undirected tree structured graphical model $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with factorization given by (7.15). Let each random variable X_i , $i \in \mathcal{V}$ assume k possible values from $\mathcal{X} = \{\alpha_1, \alpha_2, \dots, \alpha_k\}$. For each edge (u, v) , define non-negative messages along both directions: $M_{v \rightarrow u}(\alpha_i)$ along $v \rightarrow u$ and $M_{u \rightarrow v}(\alpha_i)$ along $u \rightarrow v$ for each $\alpha_i \in \mathcal{X}$. If r is the iteration number, then the update rule

$$M_{u \rightarrow v}^{(r+1)}(x_v) = \frac{1}{Z_{u \rightarrow v}^{(r)}} \sum_{x_u} \underbrace{\phi_u(x_u) \phi_{u,v}(x_u, x_v) \prod_{q \in \mathcal{N}(u) \setminus v} M_{q \rightarrow u}^{(r)}(x_u)}_{M_{u \rightarrow v}^{(r+1)}(x_v)}$$

can be executed parallelly at every node u along every edge u, v originating at u and will converge, resulting in a fix-point for each $M_{u \rightarrow v}^{(r^*)}(x_v)$ for some r^* . That is, there exists an r^* such that $M_{u \rightarrow v}^{(r^*+1)}(x_v) = M_{u \rightarrow v}^{(r^*)}(x_v)$ for all $(u, v) \in \mathcal{E}$. At convergence,

$$p(x_v) = \phi(x_v) \prod_{u \in \mathcal{N}(v)} M_{u \rightarrow v}(x_v) \quad (7.18)$$

Note that $Z_{u \rightarrow v}^{(r)}$ is the normalization factor (required for ensuring numerical stability across iterations but not otherwise) and can be computed at each iteration as

$$Z_{u \rightarrow v}^{(r)} = \sum_{x_v} \widetilde{M_{u \rightarrow v}^{(r+1)}}(x_v)$$

Theorem 91 does away with the need for any elimination ordering on the nodes and lets computations at different nodes happen in parallel. It leads to the sum-product algorithm⁷ for trees, which is summarized in Figure 7.9. The so-called flooding⁸ schedule does a ‘for’ loop at each iteration of the algorithm, for each node in \mathcal{V} . By Theorem 91, the procedure will converge. In fact, it can be proved that the algorithm will converge after at most κ iterations, where κ is the diameter (length of the longest path) of \mathcal{G} . The intuition is that message passing needs the message from every node to reach every other node and this will take κ iterations for that to happen in the sum-product algorithm.

```

Initialize  $M_{u \rightarrow v}^{(0)}(x_v)$  for each  $(u, v) \in \mathcal{E}$  to some strictly positive random values.
Set  $r = 0$ .
repeat
  for Each  $u \in \mathcal{V}$  do
    for Each  $v \in \mathcal{N}(u)$  do
       $M_{u \rightarrow v}^{(r+1)}(x_v) = \frac{1}{Z_{u \rightarrow v}^{(r)}} \sum_{x_u} \phi_u(x_u) \phi_{u,v}(x_u, x_v) \prod_{q \in \mathcal{N}(u) \setminus v} M_{q \rightarrow u}^{(r)}(x_u)$ .
    end for
  end for
  Set  $r = r + 1$ .
until  $M_{u \rightarrow v}^{(r)}(x_v) = M_{u \rightarrow v}^{(r-1)}(x_v)$  for each  $(u, v) \in \mathcal{E}$  and each  $x_v \in \mathcal{X}$ 
Set  $r^* = r - 1$ .

```

Figure 7.9: The sum-product algorithm for a tree, using the flooding schedule. It converges for $r^* \leq \kappa$, κ being the tree diameter.

There have been modern, interesting uses of message passing techniques on graphs with cycles, such as in the field of sensor networks, where locally parallelizable algorithms such as message passing are highly desirable. While there is nothing that stops us in principle from applying the algorithm in Figure 7.9 to general graphs having cycles, the theory does not guarantee anything at all - neither in terms of convergence nor in terms of the number of iterations. However, in solving partial differential equations, the message passing algorithm is often used. The algorithm is widely used on certain interesting cyclic graphs in the

⁷SIMPLE EXERCISE : Implement the algorithm using Matlab.

⁸The flooding schedule somewhat mimics the flow of water or some liquid through a network; water flows in to a node along an edge and then spreads through the other edges incident on the node.

field of communications. For special problems such as solving linear systems, this method converges⁹ on graphs with cycles as well.

However, a schedule such as in Figure 7.9 will typically incur a heavy communication cost, owing to message transmissions. While the flooding schedule is conceptually easy in terms of parallelization, an alternative schedule called the serial schedule is often preferred when parallelization is not of paramount importance. The serial schedule minimizes the number of messages passed. In this schedule, a node u transmits message to node v only when it has received messages from all other nodes $q \in \mathcal{N}(u) \setminus v$. This algorithm will pass a message only once along every direction of each edge (although during different steps). Thus, the scheduling begins with each leaf passing a message to its immediate neighbor. An overhead involved in the serial schedule is that every node needs to keep track of the edges along which it has received messages thus far. For chip level design, the highly parallelizable flooding schedule is always preferred over the serial schedule.

A point to note is that while the algorithm in Figure 7.9 is guaranteed to converge within κ steps, in practice, you might want to run the algorithm for fewer steps, until the messages reasonably converge. This strategy is especially adopted in the belief propagation algorithm, which consists of the following steps at each node of a general graphical model, until some convergence criterion is met:

1. form product of incoming messages and local evidence
2. marginalize to give outgoing message
3. propagate one message in each direction across every link

The belief propagation algorithm will be discussed later.

7.2.3 Max Product Algorithm

The max product algorithm solves the problem of determining the mode or peak of a conditional distribution, specified by a graphical model \mathcal{G} , first addressed on page 399.

$$\hat{\mathbf{x}}_{\mathcal{A}} = \operatorname{argmax}_{\mathbf{x}_{\mathcal{A}}} p(\mathbf{x}_{\mathcal{A}} \mid \mathbf{x}_{\mathcal{O}})$$

where $\mathbf{X}_{\mathcal{O}}$ is the set of observed variables, having observed values $\mathbf{x}_{\mathcal{O}}$, and $\mathbf{X}_{\mathcal{A}}$ are the query variables.

Before looking at the procedure for finding the model of a distribution, we will take a peek at an algorithm for determining the maximum value of $p(\mathbf{x})$, assuming it to be a distribution over an undirected tree \mathcal{G} . This algorithm is closely related to the sum product algorithm and can easily be obtained from the sum-product algorithm by replacing the ‘SUM’ with a ‘MAX’ in (7.17). This is because, maximums can be pushed inside products and computations can be

⁹For solving linear systems, the message passing algorithm is more efficient than Jacobi, though less efficient than conjugate gradient.

structured in a similar manner as in the sum product algorithm. This places max-product in the league of message passing algorithms.

The sum-product and max-product algorithms are formally very similar. Both methods are based on the distributive law¹⁰:

- For sum-product: $ab + ac = a(b + c)$.
- For max-product (whenever $a \geq 0$): $\max \{ab, ac\} = a \times \max \{b, c\}$.

```

Initialize  $M_{u \rightarrow v}^{(0)}(x_v)$  for each  $(u, v) \in \mathcal{E}$  to some strictly positive random
values.
Set  $r = 0$ .
repeat
  for Each  $u \in \mathcal{V}$  do
    for Each  $v \in \mathcal{N}(u)$  do
       $M_{u \rightarrow v}^{(r+1)}(x_v) = \frac{1}{Z_{u \rightarrow v}^{(r)}} \max_{x_u} \phi_u(x_u) \phi_{u,v}(x_u, x_v) \prod_{q \in \mathcal{N}(u) \setminus v} M_{q \rightarrow u}^{(r)}(x_u)$ .
    end for
  end for
  Set  $r = r + 1$ .
until  $M_{u \rightarrow v}^{(r)}(x_v) = M_{u \rightarrow v}^{(r-1)}(x_v)$  for each  $(u, v) \in \mathcal{E}$  and each  $x_v \in \mathcal{X}$ 
Set  $r^* = r - 1$ .

```

Figure 7.10: The max-product algorithm for a tree, using the flooding schedule. It converges for $r^* \leq \kappa$, κ being the tree diameter.

The max-product and sum-product algorithms can be implemented in a common framework, with a simple flag to toggle between the two. The convergence of the max-product algorithm in Figure 7.10 is very similar to the proof of the convergence of the sum-product algorithm in Figure 7.9 (the proof is rooted in Theorem 91). After convergence of the max-product algorithm in Figure 7.10, the maximum value of $\Pr(\mathbf{X} = \mathbf{x})$ can be retrieved as

$$\max_{\mathbf{x}} \Pr(\mathbf{X} = \mathbf{x}) = \max_{x_r} \phi(x_r) \prod_{u \in \mathcal{N}(r)} M_{u \rightarrow r}(x_r) \quad (7.19)$$

where, we assume that x_r is the root of the tree. However, we need to go beyond the maximum value of a distribution; we need to find the point at which the maximum is attained. To traceback this, and in general to compute $\operatorname{argmax}_{\mathbf{x}_{\mathcal{A}}} p(\mathbf{x}_{\mathcal{A}} | \mathbf{x}_{\mathcal{O}})$ we will introduce some additional machinery.

As before (*c.f.* page 403), the conditioning can be obtained as in (7.16).

¹⁰See commutative semirings for the general algebraic framework. Also refer to work on generalized distributive laws.

Definition 50 We define the singleton marginal of a distribution $p(\mathbf{x})$ as

$$\mu_s(x_s) = \frac{1}{\alpha} \max_{\mathbf{x} \setminus \{x_s\}} p(\mathbf{x}) \quad (7.20)$$

and the pairwise marginal as

$$\nu_s(x_s, x_t) = \frac{1}{\alpha} \max_{\mathbf{x} \setminus \{x_s, x_t\}} p(\mathbf{x}) \quad (7.21)$$

where

$$\alpha = \max_{\mathbf{x}} p(\mathbf{x})$$

The max marginals are analogs of marginalization, but with the summation replaced by the max operator over all variables, except x_s in the former or (x_s, x_t) in the latter. While $\mu_s(x_s)$ gives a vector of max-marginals for the variable X_s , $\nu_s(x_s, x_t)$ corresponds to a matrix of values, for the pair of variables (X_s, X_t) . You can easily convince yourself that the maximum value of any max-marginal is 1 and it will be attained for atleast one value x_s for each variable X_s .

How can the marginals be tracedback efficiently? And how can they be useful? The marginals encode preferences in the form of sufficient statistics. For example:

$$\mu_1(x_1) = \frac{1}{\alpha} \max_{x_2, x_3, \dots, x_n} p(\mathbf{x})$$

In fact, we can look at the local maximal configurations, when they are unique and traceback the (unique) global maximal configuration. This is stated in the following powerful theorem.

Theorem 92 If $\operatorname{argmax}_{x_s} \mu_s(x_s) = \{x_s^*\} \forall s \in \mathcal{V}$ (that is, there is a unique value of variable X_s that maximizes $\mu_s(x_s)$ for every $s \in \mathcal{V}$), then $\mathbf{x}^* = \{x_1^*, x_2^*, \dots, x_n^*\} = \operatorname{argmax}_{\mathbf{x}} p(\mathbf{x})$ is the unique MAP configuration.

Proof Sketch: The theorem can be equivalently stated as follows: if $\mu_i(x_i^*) > \mu_i(x_i)$ for all $x_i \neq x_i^*$, and for all $i \in \mathcal{V}$, then $p(\mathbf{x}^*) \geq p(\mathbf{x})$ for all $\mathbf{x} \neq \mathbf{x}^*$. This can be proved by contradiction as follows. Suppose $\mathbf{x}' \in \operatorname{argmax}_{\mathbf{x}} p(\mathbf{x})$. Then, for any $s \in \mathcal{V}$,

$$\mu_s(x'_s) = \max_{\mathbf{x}} p(\mathbf{x}) = \max_{x_s} \max_{\mathbf{x} \setminus x_s} p(\mathbf{x}) > \mu_s(x_s) \forall x_s$$

But by definition

$$\max_{x_s} \max_{\mathbf{x} \setminus x_s} p(\mathbf{x}) = \max_{x_s} \mu_s(x_s) = \{x_s^*\}$$

Thus, $x'_s = x_s^*$. Since $s \in \mathcal{V}$ was arbitrary, we must have $\mathbf{x}' = \mathbf{x}^*$. \square

The singleton max marginal $\mu_s(x_s)$ can be directly obtained as an outcome of the max-product algorithm:

$$\mu_s(x_s) \propto \phi_s(x_s) \prod_{u \in \mathcal{N}(s)} M_{u \rightarrow s}(x_s) \quad (7.22)$$

What if $\{x_s^1, x_s^2\} \subseteq \mu_s(x_s)$? That is, if $\mu_s(x_s)$ violates the assumption in theorem 92? Then we have to start worrying about what is happening on the edges, through the medium of the max marginal $\nu_{s,t}(x_s, x_t)$. All we do is randomly sample from the set of configurations that have maximum probability, without caring which one we really get. We first randomly sample from $\operatorname{argmax}_{x_r} \mu_r(x_r)$ at the root r of the tree and then keep randomly sample for a configuration for a child s , given its parent t (*i.e.* for an edge) $\operatorname{argmax}_{x_s} \nu_{st}(x_s, x_t)$ which respects the pairwise coupling. The following theorem states the general traceback mechanism.

Theorem 93 *Given a set of singleton max-marginals, $\{\mu_s \mid s \in \mathcal{V}\}$ and a set of pairwise marginals: $\{\nu_{st} \mid (s, t) \in \mathcal{E}\}$ on a tree $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$, $\mathbf{x}^* = (x_1^*, \dots, x_n^*)$, constructed using the following procedure is a maximal configuration, that is $\mathbf{x}^* \in \operatorname{argmax}_{\mathbf{x}} p(\mathbf{x})$.*

1. Let r be the root. Let $x_r^* \in \operatorname{argmax}_{x_r} \mu_r(x_r)$.
2. In root to leaf order, choose $x_s^* \in \operatorname{argmax}_{x_s} \nu_{st}(x_s, x_t^*)$

Proof: We will first prove that \mathbf{x}^* is optimal for the root term. For any arbitrary \mathbf{x} , by step 1,

$$\mu_r(x_r) \leq \mu_r(x_r^*) \quad (7.23)$$

If $t \rightarrow s$ is an edge, then we have by definition of the singleton max-marginal

$$\operatorname{argmax}_{x_s} \nu_{st}(x_s, x_t) = \mu_t(x_t)$$

Thus,

$$\frac{\nu_{st}(x_s, x_t)}{\mu_t(x_t)} \leq 1$$

Since $x_s^* \in \operatorname{argmax}_{x_s} \nu_{st}(x_s, x_t^*)$ by step 2, we must have $\nu_{st}(x_s^*, x_t^*) = \mu_t(x_t^*)$ and therefore, the following upperbound

$$\frac{\nu_{st}(x_s, x_t)}{\mu_t(x_t)} \leq \frac{\nu_{st}(x_s^*, x_t^*)}{\mu_t(x_t^*)} \quad (7.24)$$

for all edges $t \rightarrow s$. With repeated application of step 2, we can stitch together the local optimality conditions (7.23) and (7.24) to get

$$\mu_r(x_r) \prod_{t \rightarrow s} \frac{\nu_{st}(x_s, x_t)}{\mu_t(x_t)} \leq \mu_r(x_r^*) \prod_{t \rightarrow s} \frac{\nu_{st}(x_s^*, x_t^*)}{\mu_t(x_t^*)} \quad (7.25)$$

Just as the singleton max marginals (7.22) can be expressed in terms of the messages from the max product algorithm, the edge max-marginals can also be expressed similarly, by restricting attention to the pair of variables (x_s, X_t) instead of the world of singletons X_s , and by avoiding accounting for the message $M_{s \rightarrow t}$ or $M_{t \rightarrow s}$ between the pair:

$$\nu_{st}(x_s, x_t) \propto \phi_s(x_s) \phi_t(x_t) \phi_{st}(x_s, x_t) \prod_{u \in \mathcal{N}(s) \setminus t} M_{u \rightarrow s}(x_s) \prod_{u \in \mathcal{N}(t) \setminus s} M_{u \rightarrow t}(x_t) \quad (7.26)$$

Combining (7.22) with (7.26), we get

$$\frac{\nu_{st}(x_s, x_t)}{\mu_s(x_s) \mu_t(x_t)} \propto \frac{\phi_{st}(x_s, x_t)}{M_{t \rightarrow s}(x_s) M_{s \rightarrow t}(x_t)} \quad (7.27)$$

Applying (7.27) and (7.25) in the factorization for $p(\mathbf{x})$, we get¹¹, we obtain

$$p(\mathbf{x}) \leq p(\mathbf{x}^*)$$

Since \mathbf{x} was arbitrary, we must have $\mathbf{x}^* \in \underset{\mathbf{x}}{\operatorname{argmax}} p(\mathbf{x})$. \square

An outcome of the proof above is that for trees, the factorization can be written in terms of max-marginals instead of the potential functions:

$$p(\mathbf{x}) \propto \mu_r(x_r) \prod_{t \rightarrow s} \frac{\nu_{ts}(x_s, x_t)}{\mu_t(x_t)}$$

The above form is a directed form of factorization and does not hold for general graphs that may contain cycles. For general graphs, it can be further proved that the following factorization holds

$$p(\mathbf{x}) \propto \prod_{s \in \mathcal{E}} \mu_s(x_s) \prod_{(s,t) \in \mathcal{E}} \frac{\nu_{ts}(x_s, x_t)}{\mu_s(x_s) \mu_t(x_t)}$$

¹¹EXERCISE: Prove.

7.2.4 Junction Tree Algorithm

In many applications, the graphs are not trees. We have not discussed any principled techniques for finding marginals and modes for graphs that are not trees, though we have discussed them for trees. The elimination algorithm discussed earlier is applicable for some general graphs, but the question of what elimination should be chosen, needs to be addressed. The junction tree algorithm is very much related to the elimination algorithm, but is a more principled approach for inferencing in directed acyclic graphs. Like the sum-product algorithm, the junction tree algorithm can ‘recycle’ computations. This is unlike the general elimination algorithm, whose computation was focused completely on a single node. The work on junction trees is primarily attributed to Lauritzen and Spiegelhalter (1998). The correspondence between the graph-theoretic aspect of locality and the algorithmic aspect of computational complexity is made explicit in the junction tree framework.

For a graph with nodes, it is an intuitive idea to consider clustering completely connected nodes, form a tree connecting these clusters and finally perform message passing the tree. This idea can be formalized using the concept of a *clique tree*.

Definition 51 Given a graph $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$ with a set $\Pi \subseteq 2^{\mathcal{V}}$ of maximal cliques, a clique tree $\mathcal{T}_{\mathcal{G}}$ is a tree, whose vertices correspond the maximal cliques Π of \mathcal{G} and such that there is an edge between two nodes in the tree only if¹² there is an edge in \mathcal{G} between two nodes across the corresponding maximal cliques.

Let us take some examples:

- For the acyclic graph \mathcal{G} with $\mathcal{V} = \{1, 2, 3, 4, 5, 6\}$, $\mathcal{E} = \{(1, 2), (1, 3), (2, 4), (2, 5), (3, 6)\}$, the (not so interesting) clique tree $\mathcal{T}_{\mathcal{G}} = \langle \Pi, \mathcal{E}_{\Pi} \rangle$ would be $\Pi = \{[1, 2], [1, 3], [2, 4], [2, 5], [3, 6]\}$ and $\mathcal{E}_{\Pi} = \{([1, 2], [2, 4]), ([1, 2], [2, 5]), ([1, 2], [1, 3]), ([1, 3], [3, 6])\}$.
- For the cyclic graph \mathcal{G} with $\mathcal{V} = \{1, 2, 3, 4\}$, $\mathcal{E} = \{(1, 2), (1, 3), (2, 4), (3, 4), (2, 3)\}$, the clique tree $\mathcal{T}_{\mathcal{G}} = \langle \Pi, \mathcal{E}_{\Pi} \rangle$ would have $\Pi = \{[1, 2, 3], [2, 3, 4]\}$ and $\mathcal{E}_{\Pi} = \{([1, 2, 3], [2, 3, 4])\}$. We will adopt the practice of labeling the edge connecting nodes corresponding to two maximal cliques \mathcal{C}_1 and \mathcal{C}_2 , with their intersection $\mathcal{C}_1 \cap \mathcal{C}_2$, which will be called the *separator set*. In the second example here, the separator set for vertices $[1, 2, 3]$ and $[2, 3, 4]$ is $[2, 3]$.
- For the slightly different cyclic graph \mathcal{G} with $\mathcal{V} = \{1, 2, 3, 4\}$, $\mathcal{E} = \{(1, 2), (1, 3), (2, 4), (3, 4)\}$, there are many possible clique trees. One possible clique tree $\mathcal{T}_{\mathcal{G}} = \langle \Pi, \mathcal{E}_{\Pi} \rangle$ would have $\Pi = \{[1, 2], [1, 3], [2, 4], [3, 4]\}$ and $\mathcal{E}_{\Pi} = \{([1, 2], [1, 3]), ([1, 2], [2, 4]), ([1, 3], [3, 4])\}$. It is a bit worrisome here that $[2, 4]$ and $[3, 4]$ are not connected, since a myopic or ‘local’ sum-product algorithm, running on the clique tree might make a wrong inference that $[2, 4]$ and $[3, 4]$ do not share anything in common. In this example, for instance, the message coming from $[2, 4]$,

¹²Since we are interested in a clique ‘tree’, it may be required to drop certain edges in the derived graph. This can be seen through the third example.

through $[1, 2]$ to $[3, 4]$ has marginalized over X_4 . But this is incorrect, since $[3, 4]$ include the variable X_4 .

With the last example in mind, we will refine the notion of a clique tree to a *junction tree*, in order to ensure that local computations are guaranteed to produce globally consistent answers; **that different copies of the same random variable have ways of communicating with each other.**

Definition 52 *A junction tree for a graph $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$ having a set $\Pi \subseteq 2^{\mathcal{V}}$ of maximal cliques, is a particular type of clique tree $\mathcal{T}_{\mathcal{G}}$ such that for any two $\mathcal{C}_1, \mathcal{C}_2 \in \Pi$, $\mathcal{C}_1 \cap \mathcal{C}_2$ is a subset of every separator set on the unique path from \mathcal{C}_1 to \mathcal{C}_2 in $\mathcal{T}_{\mathcal{G}}$. This property of junction trees is called the running intersection property.*

Based on this definition, the clique trees for the first two examples are junction trees, whereas that for the third is not a junction tree. In fact, there are no junction tree for the third example. Let us consider another example.

- Consider the cyclic graph \mathcal{G} with $\mathcal{V} = \{1, 2, 3, 4, 5\}$, $\mathcal{E} = \{(1, 2), (1, 3), (2, 3), (2, 4), (3, 4), (3, 5)\}$. There are many possible clique trees for \mathcal{G} . One possible clique tree $\mathcal{T}_{\mathcal{G}} = \langle \Pi, \mathcal{E}_{\Pi} \rangle$ has $\Pi = \{[1, 2, 3], [2, 3, 4], [3, 4, 5]\}$ and $\mathcal{E}_{\Pi} = \{([1, 2, 3], [2, 3, 4]), ([2, 3, 4], [3, 4, 5])\}$ and this happens to also be a junction tree. Another clique tree with same vertex set Π and $\mathcal{E}'_{\Pi} = \{([1, 2, 3], [3, 4, 5]), ([2, 3, 4], [3, 4, 5])\}$ is not a junction tree, since node 2 which is in the intersection of $[1, 2, 3]$ and $[2, 3, 4]$ is not on every separator on the path between these two nodes. This illustrates that how you generate your clique tree matters; some clique trees may happen to be junction trees, while some may not.
- For the cyclic graph \mathcal{G} with $\mathcal{V} = \{1, 2, 3, 4\}$, $\mathcal{E} = \{(1, 2), (1, 3), (2, 4), (3, 4)\}$, considered in the third example above, there are no junction trees possible.

The global picture on a graph is specified by the factorization property:

$$p(\mathbf{x}) \propto \prod_{\mathcal{C} \in \Pi} \phi_{\mathcal{C}}(\mathbf{x}_{\mathcal{C}}) \quad (7.28)$$

On the other hand, the local message passing algorithm should honor constraints set by the separator set; that the configuration of variables in the separator set are the same in both its neighbors. More precisely, if we define the set $\tilde{\mathbf{x}}_{\mathcal{C}} = \{\tilde{x}_{i,\mathcal{C}} | i \in \mathcal{C}\}$ for each $\mathcal{C} \in \Pi$, then, for each separator set $\mathcal{C}_1 \cap \mathcal{C}_2$, the separator sets define the constraints in the form

$$\prod_{i \in \mathcal{C}_1 \cap \mathcal{C}_2} \delta(\tilde{x}_{i,\mathcal{C}_1} = \tilde{x}_{i,\mathcal{C}_2}) \quad (7.29)$$

The factorization that message passing on $\mathcal{T}_{\mathcal{G}} = \langle \Pi, \mathcal{C}' \rangle$ should see with the set of additional constraints will involve multiple copies of each random variable, but will be tied together through the δ constraints in (7.29):

$$\tilde{p}(\tilde{x}_{i,\mathcal{C}}, \forall i \in \mathcal{C}, \forall \mathcal{C} \in \Pi) \propto \prod_{\mathcal{C} \in \Pi} \phi_{\mathcal{C}}(\tilde{\mathbf{x}}_{\mathcal{C}}) \prod_{(\mathcal{C}_1, \mathcal{C}_2) \in \mathcal{E}_{\Pi}} \prod_{i \in \mathcal{C}_1 \cap \mathcal{C}_2} \delta(\tilde{x}_{i,\mathcal{C}_1} = \tilde{x}_{i,\mathcal{C}_2}) \quad (7.30)$$

The (localized) sum-product algorithm will work precisely on the junction tree $\mathcal{T}_{\mathcal{G}}$ with factorization as specified in (7.30). The factorization in (7.30) is in fact equivalent to the factorization in (7.28). In (7.30), the multiple copies of x_i 's cancel out against the constraints δ constraint. This is called the junction tree property and is formally stated in the next proposition.

Theorem 94 *For a graph \mathcal{G} having distribution $p(\mathbf{x})$ as in (7.28) and for its junction tree $\mathcal{T}_{\mathcal{G}}$ having distribution \tilde{p} specified by (7.30), the \tilde{p} distribution satisfies the following property:*

$$\tilde{p}(\tilde{x}_{i,\mathcal{C}}, \forall i \in \mathcal{C}, \forall \mathcal{C} \in \Pi) = \begin{cases} p(\mathbf{x}) & \text{if } \{x_{i,\mathcal{C}_1} = x_{i,\mathcal{C}_2} \mid \forall \mathcal{C}_1, \mathcal{C}_2 \in \Pi, \forall i \in \mathcal{C}_1 \cap \mathcal{C}_2\} \\ 0 & \text{otherwise} \end{cases}$$

That is, the new distribution \tilde{p} is faithful to the original distribution. The transitivity due to the running intersection property of junction trees is exactly what you need for this desirable property to hold.

The proof of this theorem is trivial, given the junction tree assumption. As an exercise, you may verify the truth of this statement for all the junction tree examples considered so far. The message passing formalisms in Figures 7.9 and 7.10, when applied to the junction tree, will land up not accepting contributions from inconsistent configurations, owing to the δ constraint and will therefore discover the true marginal/mode.

Theorem 95 *Suppose that \mathcal{G} has a junction tree $\mathcal{T}_{\mathcal{G}} = \langle \Pi, \mathcal{E}_{\Pi} \rangle$. Then running the sum-product or max-product algorithm on the distribution \tilde{p} defined in (7.28) will output the correct marginals or modes respectively, for p defined for \mathcal{G} , in (7.30). The $\phi_{\mathcal{C}}(\tilde{\mathbf{x}}_{\mathcal{C}})$ can be thought of as node potentials for $\mathcal{T}_{\mathcal{G}}$, while $\delta(\tilde{x}_{i,\mathcal{C}_1} = \tilde{x}_{i,\mathcal{C}_2})$ are the edge potentials.*

Theorem 95 presents a transformation of the original problem to a problem on a right kind of tree on which the running intersection property holds so that marginals and modes are preserved. The proof of this theorem is also simple. In practice, the message passing algorithm need not create multiple copies of the shared variables; the sharing can be imposed implicitly.

7.2.5 Junction Tree Propagation

The *junction tree propagation algorithm* is the sum-product algorithm applied to the junction tree, with factorization specified by (7.30). It is due to Shafer and Shenoy [SS90]. Consider the message update rule from the algorithm in Figure 7.9.

$$M_{u \rightarrow v}^{(r+1)}(x_v) = \frac{1}{Z_{u \rightarrow v}^{(r)}} \sum_{x_u} \phi_u(x_u) \phi_{u,v}(x_u, x_v) \prod_{q \in \mathcal{N}(u) \setminus v} M_{q \rightarrow u}^{(r)}(x_u)$$

If neighbors u and v are replaced by neighboring cliques \mathcal{C}_1 and \mathcal{C}_2 respectively, the equation becomes

$$M_{\mathcal{C}_1 \rightarrow \mathcal{C}_2}^{(r+1)}(\mathbf{x}_{\mathcal{C}_2}) = \frac{1}{Z_{\mathcal{C}_1 \rightarrow \mathcal{C}_2}^{(r)}} \sum_{\mathbf{x}'_{\mathcal{C}_1}} \phi_{\mathcal{C}_1}(\mathbf{x}'_{\mathcal{C}_1}) \underbrace{\prod_{i \in \mathcal{C}_1 \cap \mathcal{C}_2} \delta(x'_{i, \mathcal{C}_1} = x_{i, \mathcal{C}_2})}_{\text{based on separator set } \mathcal{C}_1 \cap \mathcal{C}_2} \prod_{\mathcal{C}_3 \in \mathcal{N}(\mathcal{C}_1) \setminus \mathcal{C}_2} M_{\mathcal{C}_3 \rightarrow \mathcal{C}_1}^{(r)}(\mathbf{x}'_{\mathcal{C}_1})$$

The constraint based on the separator set ensures that configurations that are not consistent do not contribute to the outermost summation $\sum_{\mathbf{x}'_{\mathcal{C}_1}}$. The expression for the message can therefore be equivalently written as

$$M_{\mathcal{C}_1 \rightarrow \mathcal{C}_2}^{(r+1)}(\mathbf{x}_{\mathcal{C}_2}) = \frac{1}{Z_{\mathcal{C}_1 \rightarrow \mathcal{C}_2}^{(r)}} \sum_{\mathbf{x}'_{\mathcal{C}_1 \setminus \mathcal{C}_2}} \phi_{\mathcal{C}_1}(\mathbf{x}'_{\mathcal{C}_1 \setminus \mathcal{C}_2}, \mathbf{x}_{\mathcal{C}_2}) \prod_{\mathcal{C}_3 \in \mathcal{N}(\mathcal{C}_1) \setminus \mathcal{C}_2} M_{\mathcal{C}_3 \rightarrow \mathcal{C}_1}^{(r)}(\mathbf{x}'_{\mathcal{C}_1}) \quad (7.31)$$

Note that in (7.31), the constraints based on separator sets are implicitly captured in the summation $\sum_{\mathbf{x}'_{\mathcal{C}_1 \setminus \mathcal{C}_2}}$ over only a partial set of variables from $\mathbf{x}'_{\mathcal{C}_1}$.

Further, the message $M_{\mathcal{C}_1 \rightarrow \mathcal{C}_2}^{(r+1)}(\mathbf{x}_{\mathcal{C}_2})$ is not a function of the complete vector $\mathbf{x}_{\mathcal{C}_2}$ but is only a function of $\mathbf{x}_{\mathcal{C}_2 \cap \mathcal{C}_1}$. Rewriting (7.31) to reflect this finding, we have

$$M_{\mathcal{C}_1 \rightarrow \mathcal{C}_2}^{(r+1)}(\mathbf{x}_{\mathcal{C}_2 \cap \mathcal{C}_1}) = \frac{1}{Z_{\mathcal{C}_1 \rightarrow \mathcal{C}_2}^{(r)}} \sum_{\mathbf{x}'_{\mathcal{C}_1 \setminus \mathcal{C}_2}} \phi_{\mathcal{C}_1}(\mathbf{x}'_{\mathcal{C}_1 \setminus \mathcal{C}_2}, \mathbf{x}_{\mathcal{C}_2 \cap \mathcal{C}_1}) \prod_{\mathcal{C}_3 \in \mathcal{N}(\mathcal{C}_1) \setminus \mathcal{C}_2} M_{\mathcal{C}_3 \rightarrow \mathcal{C}_1}^{(r)}(\mathbf{x}'_{\mathcal{C}_1 \cap \mathcal{C}_3}) \quad (7.32)$$

This finding is important, since it helps reduce the computational complexity of the algorithm. You need to send messages whose sizes do not depend on the size of the cliques themselves but only on the size of the separator sets. Thus, if each variable was multinomial with k possible values, then the message size would be $k^{|\mathcal{C}_1 \cap \mathcal{C}_2|}$ instead of $k^{|\mathcal{C}_2|}$. Thus, the complexity of junction tree propagation is exponential in the size of the separator sets. Typically however, the size of separator sets are not much smaller than the cliques themselves.

The junction tree propagation will converge, by the convergence property of the sum-product algorithm. After convergence, the marginals can be obtained as

$$p(\mathbf{x}_C) = \phi_C(\mathbf{x}_C) \prod_{\mathcal{D} \in \mathcal{N}(C)} M_{\mathcal{D} \rightarrow C}(\mathbf{x}_{C \cap \mathcal{D}}) \quad (7.33)$$

7.2.6 Constructing Junction Trees

How do we obtain a junction tree for a graph. And what classes of graphs have junction trees? We already saw an example of a graph that did not have any junction tree; $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$, $\mathcal{V} = \{1, 2, 3, 4\}$, $\mathcal{E} = \{(1, 2), (1, 3), (2, 4), (3, 4)\}$. We also saw an example for which a particular clique tree was not a junction tree, though it had another clique tree that was a junction tree: \mathcal{G}' with $\Pi = \{1, 2, 3, 4, 5\}$, $\mathcal{E}_\Pi = \{(1, 2), (1, 3), (2, 3), (2, 4), (3, 4), (3, 5)\}$. The junction tree $\langle \mathcal{V}'_t, \mathcal{E}_t \rangle$ has $\mathcal{V}'_t = \{[1, 2, 3], [2, 3, 4], [3, 4, 5]\}$ and $\mathcal{E}'_t = \{([1, 2, 3], [2, 3, 4]), ([2, 3, 4], [3, 4, 5])\}$, which corresponds to the elimination ordering $[1, 3, 2, 4, 5]$. While the clique tree $\mathcal{V}''_t = \mathcal{V}'_t$ and $\mathcal{E}''_t = \{([1, 2, 3], [3, 4, 5]), ([2, 3, 4], [3, 4, 5])\}$ is not a junction tree and corresponds to the elimination ordering $[1, 2, 4, 5, 3]$. We can see that the construction of the junction tree really depends on the choice of a ‘nice’ elimination ordering. One difference between \mathcal{G} and \mathcal{G}' is that, while in the former, you cannot eliminate any node without adding additional edges, in the latter, you have an elimination ordering $[1, 3, 2, 4, 5]$ that does not need to add extra edges. For \mathcal{G}' , the elimination ordering $[1, 2, 3, 4, 5]$ will not yield a junction tree.

This leads to the definition of a triangulated graph, one of the key properties of any graph which can be transformed into a junction tree. In fact, a requirement will be that an elimination algorithm is ‘good’ for junction tree construction only if it leads to a triangulated graph.

Definition 53 *A cycle is chordless if no two non-adjacent vertices on the cycle are joined by an edge. A graph is triangulated if it has no chordless cycles.*

Thus, the graph \mathcal{G} with $\mathcal{V} = \{1, 2, 3, 4\}$, $\mathcal{E} = \{(1, 2), (1, 3), (2, 4), (3, 4)\}$ is not triangulated, whereas, the graph \mathcal{G}' with $\Pi' = \{1, 2, 3, 4, 5\}$, $\mathcal{E}'_\Pi = \{(1, 2), (1, 3), (2, 3), (2, 4), (3, 4), (3, 5)\}$ is triangulated. In fact, every triangulated graph has at least one junction tree. Another equivalent characterization of a triangulated graph is as a *decomposable graph*.

Definition 54 *A graph $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$ is decomposable either if it is complete or if \mathcal{V} can be recursively divided into three disjoint sets \mathcal{A} , \mathcal{B} and \mathcal{S} such that*

1. \mathcal{S} separates \mathcal{A} and \mathcal{B} and
2. \mathcal{S} is fully connected (i.e., a clique).
3. $\mathcal{A} \cup \mathcal{S}$ and $\mathcal{B} \cup \mathcal{S}$ are also decomposable.

Following are examples of decomposable graphs:

- $\mathcal{V} = \{1, 2, 3\}$, $\mathcal{E} = \{(1, 2), (2, 3)\}$.
- $\mathcal{V} = \{1, 2, 3, 4\}$, $\mathcal{E} = \{(1, 2), (2, 3), (3, 4), (4, 1), (2, 4)\}$. Application of elimination procedure on this graph, say starting with 3 should lead in 2 and 4 being connected together, which are already connected in this graph. This shows the connection between decomposable graphs and elimination.

However, for $\mathcal{V} = \{1, 2, 3, 4\}$ and $\mathcal{E} = \{(1, 2), (2, 3), (3, 4), (4, 1)\}$, \mathcal{G} is not decomposable.

Recall from Section 7.2.1, the elimination algorithm, which eliminated one node at a time, while connecting its immediate neighbors. Thus, for any undirected graph, by the very construction of the elimination algorithm, it is obvious that the reconstituted graph, output by the elimination algorithm is always triangulated. This can be proved by induction on the number of vertices in the graph¹³. The statement is trivial for the base case of a one node graph.

The following theorem is a *fundamental characterization* of graphs that have junction trees.

Theorem 96 *The following are equivalent ways of characterizing a graph \mathcal{G} :*

1. \mathcal{G} is decomposable.
 - *(This captures the ‘divide-and-conquer’ nature of message passing algorithms. In fact, the message passing algorithm exploited a divide and conquer strategy for computation on trees.)*
2. \mathcal{G} is triangulated.
 - *(Elimination can result in a triangulated graph.)*
3. \mathcal{G} has a junction tree.
 - *(If the graph is triangulated, it must have at least one junction tree. And junction tree is a good canonical data structure for conducting computations on general graphs.)*

Some practical impacts of this theorem are listed itemized in brackets by the side of each of the equivalent characterizations of \mathcal{G} . The equivalence of the first and second statements in the theorem can be proved very simply by induction on the number of nodes in the graph.

The first step in the junction tree algorithm is triangulating a graph. This might mean adding extra edges or increasing clique size. But this cannot be harmful¹⁴, since the potential function can be defined over a larger clique

¹³Prove: EXERCISE.

¹⁴What characterizes a graphical models is not the presence of edges, but the absence of edges. As an the extreme example, a complete graph potentially subsumes every graphical model.

as the product of potential functions over its sub-parts. Given a triangulated graph, we know that it must have a junction tree by virtue of theorem 96. How can a junction tree be constructed from a triangulated graph? The first step would be to isolate all its cliques. Going back to the second example on page 412, the triangulated graph \mathcal{G} with $\mathcal{V} = \{1, 2, 3, 4, 5\}$, $\mathcal{E} = \{(1, 2), (1, 3), (2, 3), (2, 4), (3, 4), (3, 5)\}$ has three maximal cliques: $\Pi = \{[1, 2, 3], [2, 3, 4], [3, 4, 5]\}$. There are different ways to connect up the cliques to form a tree. One possible clique tree $\mathcal{T}_{\mathcal{G}} = \langle \Pi, \mathcal{E}_{\Pi} \rangle$ has $\mathcal{E}_{\Pi} = \{([1, 2, 3], [2, 3, 4]), ([2, 3, 4], [3, 4, 5])\}$ and this happens to also be a junction tree. Another clique tree has $\mathcal{E}'_{\Pi} = \{([1, 2, 3], [3, 4, 5]), ([2, 3, 4], [3, 4, 5])\}$ is not a junction tree, since node 2 which is in the intersection of $[1, 2, 3]$ and $[2, 3, 4]$ is not on every separator on the path between these two nodes.

While you can discover a junction tree by exhaustive search, that is an infeasible idea. However, the search for a junction tree can be performed efficiently by making use of the following theorem.

Theorem 97 *Let $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$ be a triangulated graph with $\mathcal{V} = \{X_1, X_2, \dots, X_n\}$ and with Π being the set of maximal cliques. Let $\mathcal{T}_{\mathcal{G}}$ be a spanning tree for the clique graph of \mathcal{G} , having $\Pi = \{\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_m\}$ as the set of vertices and $\mathcal{E}_{\Pi} = \{e_1, e_2, \dots, e_{m-1}\}$ as the set of edges ($|\Pi| = m$ and $|\mathcal{E}_{\Pi}| = m - 1$). Let $\mathcal{S}(e)$ be the separator associated¹⁵ with any edge $e \in \mathcal{E}_{\Pi}$. We will define the weight¹⁶ of the spanning tree as*

$$w(\mathcal{T}_{\mathcal{G}}) \in \sum_{i=1}^{m-1} |\mathcal{S}(e_i)| = \sum_{i=1}^{m-1} \sum_{j=1}^n [X_j \in \mathcal{S}(e_i)] \quad (7.34)$$

where $[condition]$ is the indicator function that assumes value 1 if and only if condition is satisfied and is 0 otherwise. Then, if

$$\widehat{\mathcal{T}}_{\mathcal{G}} = \operatorname{argmax}_{\mathcal{T}_{\mathcal{G}}} w(\mathcal{T}_{\mathcal{G}})$$

$\widehat{\mathcal{T}}_{\mathcal{G}}$ must be a junction tree.

Proof: First we note that for any variable X_j , the number of separator sets in $\mathcal{T}_{\mathcal{G}}$ in which X_j appears is upper bounded by the number of cliques in which X_j appears.

$$\sum_{i=1}^{m-1} [X_j \in \mathcal{S}(e_i)] \leq \sum_{i=1}^m [X_j \in \mathcal{C}_i] \quad (7.35)$$

¹⁵Since any edge e in the clique graph is of the form $(\mathcal{C}_1, \mathcal{C}_2)$, where $\mathcal{C}_1, \mathcal{C}_2 \in \Pi$, the separator associated with e can be viewed as a function $\mathcal{S}(e) = \mathcal{C}_1 \cap \mathcal{C}_2$.

¹⁶For the example above with $\mathcal{E}_{\Pi} = \{([1, 2, 3], [2, 3, 4]), ([2, 3, 4], [3, 4, 5])\}$, the weight is $2 + 2 = 4$, and this also happens to be a junction tree. For $\mathcal{E}'_{\Pi} = \{([1, 2, 3], [3, 4, 5]), ([2, 3, 4], [3, 4, 5])\}$, the weight is $1 + 2 = 3$ and this is not a junction tree.

Equality will hold if and only if the running intersection property holds for X_j in $\mathcal{T}_{\mathcal{G}}$. By interchanging the order of summations in (7.34) and applying the inequality in (7.35), it follows that

$$w(\mathcal{T}_{\mathcal{G}}) = \sum_{i=1}^{m-1} \sum_{j=1}^n [X_j \in \mathcal{S}(e_i)] \leq \sum_{j=1}^n \left(\sum_{i=1}^m [X_j \in \mathcal{C}_i] - 1 \right)$$

Interchanging the summations in the rightmost term yields an upper-bound on $w(\mathcal{T}_{\mathcal{G}})$, which can be attained if and only if $\mathcal{T}_{\mathcal{G}}$ is a junction tree (that is, if and only if equality holds in (7.35) for all $1 \leq j \leq n$)

$$w(\mathcal{T}_{\mathcal{G}}) \leq \sum_{i=1}^m |\mathcal{C}_i| - n$$

We know from lemma 96 that if \mathcal{G} is triangulated, it must have a junction tree. Given that \mathcal{G} is triangulated, $\widehat{\mathcal{T}}_{\mathcal{G}} = \operatorname{argmax}_{\mathcal{T}_{\mathcal{G}}} w(\mathcal{T}_{\mathcal{G}})$ must be a junction tree and

$$\text{will satisfy } w(\widehat{\mathcal{T}}_{\mathcal{G}}) = \sum_{i=1}^m |\mathcal{C}_i| - n. \quad \square$$

The maximum weight spanning tree problem in (7.34) can be solved exactly by executing the following step $m - 1$ times, after initializing \mathcal{E}_{Π}^{all} to all possible ‘legal’ edges between nodes in Π and $\mathcal{E}_{\Pi} = \{\}$

1. For $i = 2$ to $m - 1$, if

$$\widehat{e} = \operatorname{argmax}_{e \in \text{acyclic}(\mathcal{E}_{\Pi}, \mathcal{E}_{\Pi}^{all})} |\mathcal{S}(e)|$$

then set $\mathcal{E}_{\Pi} = \mathcal{E}_{\Pi} \cup \{\widehat{e}\}$ and $\mathcal{E}_{\Pi}^{all} = \mathcal{E}_{\Pi}^{all} \setminus \{\widehat{e}\}$.

Here, $\text{acyclic}(\mathcal{E}_{\Pi}, \mathcal{E}_{\Pi}^{all}) = \{e \in \mathcal{E}_{\Pi}^{all} \mid \mathcal{E}_{\Pi} \cup \{e\} \text{ has no cycles}\}$. This can be efficiently implemented using Kruskal and Prim’s algorithm. The only additional requirement is that this problem requires specialized data structure to quickly check if $e \in \text{acyclic}(\mathcal{E}_{\Pi}, \mathcal{E}_{\Pi}^{all})$, that is, if addition of e to the current set of edges would induce any cycle. This discussion is also relevant for learning tree structured graphical models such that the structure maximizes some objective function on the data.

In Figure 7.11, we present the overall junction tree algorithm. Typically, junction tree propagation is the most expensive step (and is the main ‘online’ step) and has complexity $O(m|\mathcal{C}_{max}|^k)$, where k is the maximum number of states for any random variable and \mathcal{C}_{max} is the largest clique. The treewidth τ of a graph is defined as $\tau = \mathcal{C}_{max} - 1$ in the optimal triangulation. Thus, the junction tree propagation algorithm scales exponentially in the treewidth τ . There are many elimination orderings/triangulations. The best triangulation is the one that leads to smallest value of \mathcal{C}_{max} . The problem of finding the best elimination ordering or of finding the best junction tree is NP-hard. In

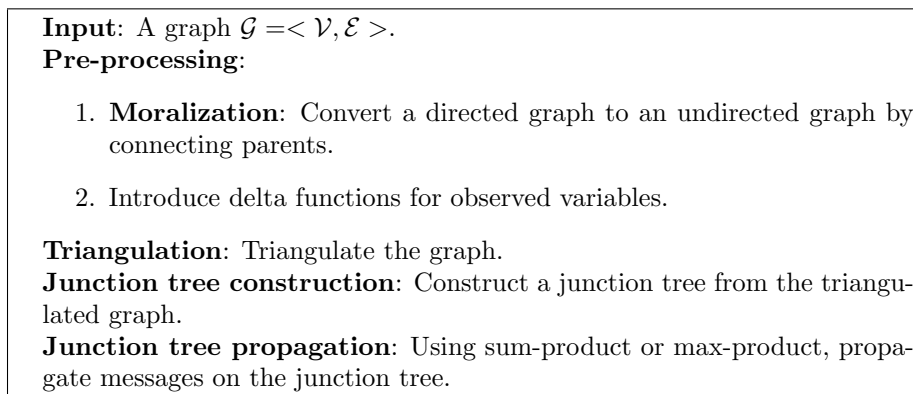


Figure 7.11: The high-level view of the Junction tree algorithm.

practice, there are many heuristics that work well. Though NP-hard in a worst case sense, this problem is much easier in the average case.

Many problems do not have bounded treewidth. The junction tree algorithm is limited in its application to families of graphical models that have bounded treewidth. Many common graphical models, such as grid structured graphical model that is commonplace in image processing have very high treewidth. The treewidth of an $n \times n$ grid (*i.e.*, n^2 nodes) scales as $O(n)$. Thus, junction tree becomes infeasible for grids as large as 500×500 , though it is applicable in theory.

7.2.7 Approximate Inference using Sampling

While the generic junction tree method is principled, it is limited to graphs with bounded treewidth. There are several approximate inference methods that could be considered as alternatives, in practice. One class of approximate inference techniques is the class of sampling methods.

Monte Carlo Methods

The general umbrella problem underlying Monte Carlo sampling methods is

$$E[f] = \int p(\mathbf{x})f(\mathbf{x})d\mathbf{x} \quad (7.36)$$

where $f(\mathbf{x})$ is some function defined on some possibly high dimensional space in \mathfrak{R}^n and p is a distribution defined on the same space. For $f(\mathbf{x}) = \mathbf{x}$, $E[f]$ turns out to be the mean. For $f(\mathbf{x}) = \delta(\mathbf{x} = \mathbf{x}')$, $E[f]$ becomes the marginal probability $p(\mathbf{X} = \mathbf{x}')$ or more general, for $f(\mathbf{x}) = \delta(\mathbf{x} \leq \mathbf{x}')$, $E[f]$ becomes the tail probability $p(\mathbf{x} \geq \mathbf{x}')$. The goal of sampling methods, such as Monte

Carlo methods is to approximate such integrals over such possibly high dimensional space using sampling techniques. If we could collect iid samples¹⁷ $\bar{\mathbf{x}} = (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)})$ from $p(\cdot)$, then

$$\hat{f} = \frac{1}{m} \sum_{i=1}^m f(\mathbf{x}^{(i)})$$

is a Monte Carlo estimate of $E[f]$. Since $\bar{\mathbf{x}}$ is a collection of random samples from $p(\cdot)$, \hat{f} is also a random variable. We could ask some questions about such an estimator:

- Is \hat{f} unbiased? That is, on an average, will the estimator produce the right quantity? The estimator is unbiased if

$$E[f] = E_{\bar{\mathbf{x}}}[\hat{f}]$$

Using the linearity property of expectation,

$$E_{\bar{\mathbf{x}}}[\hat{f}] = \frac{1}{m} \sum_{i=1}^m E[f(\mathbf{x}^{(i)})] = E[f]$$

that is, the estimator \hat{f} is indeed unbiased and gives the right answer on an average.

- The unbiased requirement on the part of the estimator only requires the sample mean to match the expected mean, on an average. It may also be desirable that the variance be stable. A related expectation is that as the number m of samples increases, the estimate should get closer to the actual answer. In fact, this is true for the estimator just discussed. It can be shown that if f has finite variance,

$$\text{var}(\hat{f}) = \frac{1}{m^2} \sum_{i=1}^m \text{var}(f(\mathbf{x}^{(i)})) = \sum_{i=1}^m \frac{\text{var}(f(\mathbf{x}^{(i)}))}{m} \quad (7.37)$$

or equivalently, the spread for $\text{var}(\hat{f})$ is $\frac{1}{\sqrt{m}}$ times the spread for \hat{f} . From this, we can infer that as $m \rightarrow \infty$, $\text{var}(\hat{f}) \rightarrow 0$.

The discussion thus far was centered around the assumption that we can draw samples from $p(\cdot)$. This area of sampling has warranted separate attention for research. Even generation of pseudo random numbers is not straightforward¹⁸. As another example, it is not straightforward to sample efficiently from

¹⁷The difference between $\bar{\mathbf{x}}$ here and in the case of maximum likelihood estimation is that in the latter case, $\bar{\mathbf{x}}$ is data provided, whereas here, we consider $\bar{\mathbf{x}}$ sampled from the model itself. However, the analytical form is similar.

¹⁸If a distribution function can be ‘inverted’ a common strategy is to pass a uniform distribution through it to generate samples.

a typically graphical model-like distribution (especially if it is multi-model) such as

$$p(x) = \frac{1}{Z} \underbrace{\exp \{ax^4 + bx^3 + cx^2 + dx + e\}}_{\varphi(x)}$$

where, the $\varphi(x)$ part is easy to compute, whereas $\frac{1}{Z}$ is hard to compute.

We will not look at a series of attempts at sampling from a distribution $p(\cdot)$.

Adopting Numeric Methods

One natural way out is to adopt standard numerical methods, such as the standard numerical recipe for evaluating an integral from first principles - creating discrete intervals and then letting the intervals shrink.

1. Discretize the space of x (such as the real line) into k discrete points, x_1, x_2, \dots, x_k

2. Compute an approximation to z as $\hat{z} = \sum_{i=1}^k \varphi(x_i)$.

3. The samples can then be drawn from one of k points based on the distribution p_d :

$$p_d(x_i) = \frac{\varphi(x_i)}{\hat{z}}$$

The new distribution has point masses at the samples x_1, x_2, \dots, x_k . As the number of grows larger, the approximation will get better

While a controllable approximation that works well in one dimension, how well does such a discretization method scale to higher dimensions. The number of discrete points scales exponentially¹⁹ in the number of dimensions as k^n (n being the dimensionality). Thus, discretization is not feasible in higher dimensions. Another factor in favour of Monte Carlo methods, vis-a-vis numerical techniques is that the statement (7.37) for the Monte Carlo estimator is really independent of the dimensionality n .

Rejection Sampling

Rejection sampling, which dates back to von Neumann in 1950's assumes that in addition to the decomposition $p(\mathbf{x}) = \frac{z}{\varphi(\mathbf{x})}$ with $\varphi(x)$ easy to compute and $\frac{1}{z}$ is hard to compute, we also have a proposal distribution $q(\mathbf{x})$ that is relatively easy to (exactly) sample from. q could be one of Gaussians, Cauchy, or some other member of the exponential family.

$$q(\mathbf{x}) = \frac{z_q}{\hat{q}(\mathbf{x})}$$

¹⁹This problem also goes under the name of the *curse of dimensionality* - the task that is easy in a single dimension becomes extremely complex at higher dimensions.

Rejection sampling also assumes that you have a constant M such that $\hat{q}(\mathbf{x})$ scaled by the constant yields an upper bound for the original distribution $\varphi(\mathbf{x})$.

$$\varphi(\mathbf{x}) \leq M\hat{q}(\mathbf{x})$$

The way rejection sampling works is:

1. First generate a sample $\mathbf{y} \sim q(\mathbf{x})$.
2. Secondly, sample u from a uniform distribution between 0 and $M\hat{q}(\mathbf{y})$:
 $u \sim U[0, M\hat{q}(\mathbf{y})]$.
 - If $u < \varphi(\mathbf{y})$, then accept \mathbf{y} as a realization of $\varphi(\mathbf{x})$.
 - Otherwise, reject \mathbf{y} .

We will see that this random procedure itself induces a distribution p_{rej} over \mathbf{x} that happens to be the same as $\varphi(\mathbf{x})$. For any \mathbf{y} that is an output of the rejection sampling procedure, its probability of being generated by the procedure is the product of the probability $q(\mathbf{y})$ of choosing \mathbf{y} and the probability $\frac{\varphi(\mathbf{y})}{M\hat{q}(\mathbf{y})}$ of accepting \mathbf{y} :

$$p_{gen}^{rej}(\mathbf{y}) = \frac{1}{z_{p_{gen}^{rej}}} q(\mathbf{y}) \left(\frac{\varphi(\mathbf{y})}{M\hat{q}(\mathbf{y})} \right) = \frac{1}{z_{p_{gen}^{rej}} z_q} \varphi(\mathbf{y})$$

where, the normalization constant $\frac{1}{z_{p_{gen}^{rej}}}$ is defined as

$$z_{p_{gen}^{rej}} = \int_{\mathbf{y}'} q(\mathbf{y}') \left(\frac{\varphi(\mathbf{y}')}{M\hat{q}(\mathbf{y}')} \right) d\mathbf{y}' = \int_{\mathbf{y}'} \frac{1}{z_q} \varphi(\mathbf{y}') d\mathbf{y}' = \frac{z_p}{z_q}$$

Combined, these two equalities mean that

$$p_{gen}^{rej}(\mathbf{y}) = \frac{1}{z_{p_{gen}^{rej}} z_q} \varphi(\mathbf{y}) = \frac{z_q}{z_{p_{gen}^{rej}} z_q z_p} \varphi(\mathbf{y}) = \frac{1}{z_p} \varphi(\mathbf{y}) = p(\mathbf{y})$$

In practice, it crucially matters how small you can make the reject region, since you would not like to spend too many sampling cycles to generate each sample \mathbf{y} . So it will be best to choose a \hat{q} that follows φ very closely. A measure of this ‘following closely’ is the ratios of the area A_{acc} under $\varphi(\mathbf{x})$ to the area A_{tot} under $M\hat{q}(\mathbf{x})$. This can be thought of as the acceptance probability p_{acc}^{rej} .

$$p_{acc}^{rej} = \frac{\int_{\mathbf{x}} \varphi(\mathbf{x}) d\mathbf{x}}{\int_{\mathbf{x}} M\hat{q}(\mathbf{x}) d\mathbf{x}} = \frac{A_{acc}}{M z_q} = \frac{A_{acc}}{A_{tot}}$$

One of the concerns with rejection sampling in high dimensions is that since $p_{acc}^{rej} \propto \frac{1}{M}$, the number of attempts before getting a single sample will scale as M . For instance, suppose \mathbf{X} and \mathbf{Y} are independent Gaussians with slightly different

variances - $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}, \mathbf{0}, \sigma_p^2 I)$ and $q(\mathbf{y}) = \mathcal{N}(\mathbf{y}, \mathbf{0}, \sigma_q^2 I)$. where $\rho \in (0, 1.1]$. For rejection sampling, we will need to choose an $M > \frac{\rho(\mathbf{x})}{\hat{q}(\mathbf{x})}$ for every $\mathbf{x} \in \mathfrak{R}^n$. For this Gaussian, we will require that $M > \frac{\rho(\mathbf{0})}{\hat{q}(\mathbf{0})} = \exp \left\{ n \log \frac{\sigma_q}{\sigma_p} \right\}$. Note that if σ_q is even slightly larger than σ_p , then M will increase exponentially with n . That is, we will have to do exponentially many rejection trials before getting a sample accepted.

In summary, while rejection sampling is useful in low dimensions, it breaks down in higher dimensions. Markov chain monte carlo (MCMC) builds on rejection sampling. While rejection sampling is memoryless - in the sense that rejected samples are naively abandoned, MCMC preserves rejected samples using memory.

Importance Sampling

Importance sampling is a more general form of Monte Carlo sampling method. Recall that Monte Carlo sampling was to solve the problem

$$E[f] = \int p(\mathbf{x})f(\mathbf{x})d\mathbf{x} \quad (7.38)$$

and the Monte Carlo estimate collects iid samples $\bar{\mathbf{x}} = (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)})$ from $p(\cdot)$ and computes the weighted sum

$$\hat{f} = \frac{1}{m} \sum_{i=1}^m f(\mathbf{x}^{(i)})$$

The idea of importance sampling is to generate samples $\bar{\mathbf{y}}$ from an alternative q which is somewhat easier to sample than p . However, how do we make use of $\bar{\mathbf{y}}$ in estimation? Assuming that $q(\mathbf{x}) > 0$ whenever $f(\mathbf{x})p(\mathbf{x}) > 0$, we rewrite the expression for $E[f]$ as

$$E[f] = \int \frac{p(\mathbf{x})f(\mathbf{x})}{q(\mathbf{x})}q(\mathbf{x})d\mathbf{x}$$

Defining $g(\mathbf{x}) = \frac{p(\mathbf{x})f(\mathbf{x})}{q(\mathbf{x})}$, we just re-express

$$E_p[f] = E_q[g]$$

This motivates the study of the Monte Carlo estimate \hat{g} based on samples $\bar{\mathbf{y}}$ from q .

$$\hat{g} = \frac{1}{m} \sum_{i=1}^m g(\mathbf{y}^{(i)}) = \frac{1}{m} \sum_{i=1}^m \frac{p(\mathbf{y}^{(i)})f(\mathbf{y}^{(i)})}{q(\mathbf{y}^{(i)})} \quad (7.39)$$

The estimate \hat{g} is called the importance sampling estimate. The terms $\frac{p(\mathbf{y}^{(i)})}{q(\mathbf{y}^{(i)})}$ are called *importance weights*. It can be shown that \hat{g} is unbiased, that is,

$$E_p[f] = E_{\hat{y}}[\hat{g}]$$

Like for the case of the Monte Carlo estimate, it can also be shown that if g has finite variance then,

$$\text{var}(\hat{g}) = \sum_{i=1}^m \frac{\text{var}(g(\mathbf{y}^{(i)}))}{m} \quad (7.40)$$

or equivalently, the spread for $\text{var}(\hat{g})$ is $\frac{1}{\sqrt{m}}$ times the spread for \hat{g} .

Importance sampling is useful when q is easier to sample than p . Backing off a bit, it may happen that q is itself of the form

$$q(\mathbf{y}) = \frac{1}{z_q} \hat{q}(\mathbf{y})$$

where $\hat{q}(\mathbf{y})$ is easy to compute while the normalization constant z_q is not. This is especially true in the case of graphical models, for which $\hat{q}(\mathbf{y})$ is simply the product of some compatibility functions or exponentiated weighted sum of sufficient statistics (exponential family). Similarly, it is often that $p(\mathbf{x}) = \frac{1}{z_p} \wp(\mathbf{x})$. In such a case, we can write

$$\hat{g} = \frac{1}{m} \sum_{i=1}^m \frac{p(\mathbf{y}^{(i)})f(\mathbf{y}^{(i)})}{q(\mathbf{y}^{(i)})} = \frac{1}{m} \sum_{i=1}^m \frac{\wp(\mathbf{y}^{(i)})f(\mathbf{y}^{(i)})}{\hat{q}(\mathbf{y}^{(i)})} \frac{z_q}{z_p}$$

For this formulation, $\frac{\wp(\mathbf{y}^{(i)})}{\hat{q}(\mathbf{y}^{(i)})}$ are called the importance weights $im(\mathbf{y}^{(i)})$. Also,

$$z_{p/q} = \frac{z_p}{z_q} = \frac{1}{z_q} \int \wp(\mathbf{x})d\mathbf{x} = \int \wp(\mathbf{x}) \frac{q(\mathbf{x})}{\hat{q}(\mathbf{x})} = \int \frac{\wp(\mathbf{x})}{\hat{q}(\mathbf{x})} q(\mathbf{x})d\mathbf{x}$$

which is again the expectation under q of $im(\mathbf{y})$. The Monte Carlo estimate $\hat{z}_{p/q}$ for $z_{p/q} = \frac{z_p}{z_q}$ is

$$\hat{z}_{p/q} = \frac{1}{m} \sum_{i=1}^m im(\mathbf{y}^{(i)}) = \frac{1}{m} \sum_{i=1}^m \frac{\wp(\mathbf{y}^{(i)})}{\hat{q}(\mathbf{y}^{(i)})}$$

This gives you a modified Monte Carlo estimate, which can be contrasted against (7.39).

$$\hat{g}' = \frac{\frac{1}{m} \sum_{i=1}^m \frac{p(\mathbf{y}^{(i)})f(\mathbf{y}^{(i)})}{q(\mathbf{y}^{(i)})}}{\frac{1}{m} \sum_{i=1}^m \frac{\wp(\mathbf{y}^{(i)})}{\hat{q}(\mathbf{y}^{(i)})}} = \frac{\sum_{i=1}^m f(\mathbf{y}^{(i)})im(\mathbf{y}^{(i)})}{\sum_{i=1}^m im(\mathbf{y}^{(i)})} \quad (7.41)$$

This modified Monte Carlo estimate uses samples $\bar{\mathbf{y}}$ from q and also does not require explicit computation of p . Rather it needs to compute only φ .

Importance sampling considerably widens the scope of Monte Carlo methods, since it provides flexibility of choosing q . In fact, even if you could sample from p , it can be useful to use a q , especially when f lies in the tail region(s) of p . Since q lets you reshape what you sample from, the modified g can help shift the mass over to regions of p such that getting information from f becomes much more likely. In practice, adapting q to the shape of f can also lead to reduction in the variance of \hat{g} .

Finally, importance sampling can be useful even when you can draw samples from $p(\cdot)$. For example, say $X_i \in \{0, 1\}$, for $1 \leq i \leq n$ are n binary random variables with $p(X_i = 1) = \epsilon$ for a very small ϵ , close to 0. Let the X_i 's be independent (though the example could hold for many non-independent X_i 's as well). Let us say, we are interested in estimating

$$p\left(\frac{\sum_{i=1}^n x_i}{n} \geq 0.5\right)$$

As can be seen, that $\frac{\sum_{i=1}^n x_i}{n} \geq 0.5$ is a very rare event. Importance sampling can be used to model such rare events which are otherwise computationally infeasible to simulate.

Like in the case of rejection sampling and numerical techniques, there are problems that importance sampling techniques have in high dimensional spaces.

Markov Chain Monte Carlo

Markov Chain Monte Carlo (MCMC) methods are a suite of methods based on setting up a markov chain to generate samples from a target distribution

$$p(\mathbf{x}) = \frac{1}{z} \varphi(\mathbf{x})$$

where z may not be easy to compute. The basic idea here is to set up a Markov chain $X_1 \rightarrow X_2 \rightarrow \dots \rightarrow X_n$ and a sequence of distributions q so that as $t \rightarrow \infty$, $q(\cdot | X^{(t)}) \rightarrow p$. While these methods have a flavour of rejection sampling, they are not memoryless; rejected samples are not entirely discarded. Introduced by physicists Metropolis, Rosenbluth, Teller in 1953, these methods were generalized by Hastings in 1970.

The Metropolis-Hastings algorithm is one of the more popular examples from this class. It is outlined in Figure 7.12. It builds on the rejection sampling technique, with two main differences. More importantly, a sample is not discarded if rejected; the value of X_{t+1} is set to X_t . Secondly, the acceptance probability determines if it is more likely to go $\mathbf{X}^{(t)} \rightarrow \mathbf{y}$ or $\mathbf{y} \rightarrow \mathbf{X}^{(t)}$.

Input: A target distribution $p(\mathbf{x}) = \frac{1}{z}\varphi(\mathbf{x})$.
Initialize: $\mathbf{X}^{(1)} \sim q(\mathbf{x})$ for some arbitrary q .
for $t = 1, 2, \dots$ **do**
 1. Sample \mathbf{y} from the conditional probability distribution $q(\cdot|\mathbf{X}^{(t)})$.
 2. Sample $u \sim Uniform[0, 1]$.
 3. Define the acceptance probability as

$$A(\mathbf{X}^{(t)}, \mathbf{y}) = \min \left\{ 1, \frac{\varphi(\mathbf{y})q(\mathbf{X}^{(t)}|\mathbf{y})}{\varphi(\mathbf{X}^{(t)})q(\mathbf{y}|\mathbf{X}^{(t)})} \right\}$$

4. Set

$$\mathbf{X}^{(t+1)} = \begin{cases} \mathbf{y}, & \text{if } u \leq A(\mathbf{X}^{(t)}, \mathbf{y}) \quad // \text{Accept} \\ \mathbf{X}^{(t)} & \text{otherwise} \quad // \text{Reject, but do not discard.} \end{cases}$$

end for

Figure 7.12: The Metropolis-Hastings algorithm.

In the special case of symmetry, that is if $q(\mathbf{y}|\mathbf{X}^{(t)}) = q(\mathbf{X}^{(t)}|\mathbf{y})$,

$$A(\mathbf{X}^{(t)}, \mathbf{y}) = \min \left\{ 1, \frac{\varphi(\mathbf{y})}{\varphi(\mathbf{X}^{(t)})} \right\}$$

the algorithm acts like a gradient algorithm with some randomness. To see this, note that if $\varphi(\mathbf{y}) \geq \varphi(\mathbf{X}^{(t)})$, the algorithm will always accept. Else it accepts with probability $\frac{\varphi(\mathbf{y})}{\varphi(\mathbf{X}^{(t)})} < 1$; the logic is that you do not always want to reject even if you were to go downhill, since you might want to wriggle out of local modes. So you reject, but probabilistically. Thus, the algorithm always tries to sample from regions of the space where the density is more.

We will eventually see that if $q(\mathbf{y}|\mathbf{X}^{(t)}) \rightarrow p(\mathbf{y})$ as $t \rightarrow \infty$. Enroute to proving this, we will require to understand the limiting behaviour of Markov chains as number of states goes to ∞ .

1. The Metropolis-Hastings algorithm in Figure 7.12 generates a first order Markov chain. Assume for simplicity that the Markov chain is finite state and that $\mathbf{X}^{(i)} \in \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k\}$. For many graphical models, k could be exponential in the number of nodes in the graphical model. For the Markov chain $\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \mathbf{X}^{(3)}, \dots$ generated by the algorithm in Figure 7.12, $\mathbf{X}^{(1)} \sim q(\cdot)$ while the transition $\mathbf{X}^{(t)} \rightarrow \mathbf{X}^{(t+1)}$ is specified by the homogeneous (*i.e.* fixed across time steps) conditional distribution

$$\Gamma(\mathbf{X}^{(t+1)}|\mathbf{X}^{(t)}) = A(\mathbf{X}^{(t)}, \mathbf{y})q(\mathbf{y}|\mathbf{X}^{(t)})$$

That is, the transition function for the algorithm is the combination of the original proposal distribution and the probability of acceptance in the

accept/reject step. Then, by the steps (1) and (4) of the algorithm in Figure 7.12,

$$q(\mathbf{X}^{(t+1)}) = \sum_{\mathbf{x}^{(t)}} q(\mathbf{X}^{(t)}) \Gamma(\mathbf{X}^{(t+1)} | \mathbf{X}^{(t)})$$

In matrix notation, this translates to

$$\mathbf{q}_{t+1} = \mathbf{q}_t^T T \quad (7.42)$$

where, $\Gamma[i, j] = \Gamma(\mathbf{X}^{(t+1)} = \mathbf{x}_j | \mathbf{X}^{(t)} = \mathbf{x}_i)$ and $\mathbf{q}_t[i] = q(\mathbf{X}^{(t)} = \mathbf{x}_i)$. By its very definition, Γ is row-stochastic, that is,

$$\Gamma \mathbf{1} = \mathbf{1}$$

We would like $\mathbf{q} \rightarrow \mathbf{p}$, where \mathbf{p} is the targetted probability vector. The following sequence of arguments will help arrive at conditions under which this will happen.

2. Let \mathbf{r} be the fix point of update equation (7.42). Then \mathbf{r} is also invariant²⁰ with respect to Γ . Thus, once the distribution \mathbf{q}_t hits an invariant \mathbf{r} , it stays in \mathbf{q}_t .
3. In order to have an invariant vector, the matrix must be non-negative ($\Gamma \geq 0$) and must be row-stochastic, which it is. The matrix Γ is not symmetric in general. The Perroon Forbenius theorem states that for any non-negative, row-stochastic matrix A , its spectral radius $\rho(A) = \max_{i=1,2,\dots,n} |\lambda_i(A)|$ satisfies the condition $\rho(A) = 1$ and that it has a left eigenvector $\mathbf{v} \geq \mathbf{0}$ such that $\mathbf{v}^T A = \mathbf{v}$. Since the matrix Γ is both non-negative and row-stochastic, a consequence of this theorem is that $\mathbf{r} = \frac{1}{\sum_{k=1}^k v_k} \mathbf{v}$ will be invariant with respect to Γ .
4. The above conditions and arguments state in principle that there is potentially a fix point for (7.42). In general, the fix point need not be unique; for a diagonal matrix, there are several invariant distributions. It can shown that an irreducible matrix Γ (*i.e.*, every state is reachable from every other state with strictly positive probability) will have a unique invariant distribution \mathbf{r} .

²⁰A vector $\mathbf{r} \in \mathbb{R}^k$ is invariant with respect to matrix Γ if \mathbf{r} represents a probability distribution (*i.e.*, $\mathbf{r}^T \mathbf{1} = \mathbf{1}$ and $\mathbf{r} \geq \mathbf{0}$) and is fixed under the updates of Γ , that is $\mathbf{r}^T \Gamma = \mathbf{r}$. As an example, you can verify that for random walks on a graph $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$ with 0 jump probability, \mathbf{r} such that $v_i = \frac{d_i}{\sum_{i \in \mathcal{V}} d_i}$ is invariant, d_i being the degree of vertex i .

5. But will the algorithm in Figure 7.12 converge to the fix point? To enable convergence, another necessary condition is that the Markov chain should be aperiodic (so that q_t does not keep toggling between values) or equivalently, have a period of 1.
6. With all the armour discussed so far, we state the following theorem, central to the correct convergence of the algorithm in Figure 7.12:

Theorem 98 *For any finite-chain irreducible aperiodic Markov Chain, the sequence $\mathbf{q}_{t+1} = \mathbf{q}_t^T \Gamma$ converges, for an arbitrary \mathbf{q}_1 , to the unique invariant distribution \mathbf{r} of the chain.*

The next few steps are dedicated to proving that (i) the Metropolis-Hastings algorithm satisfies the pre-requisites for this theorem under certain conditions and (ii) that the target distribution is indeed invariant with respect to the Markov chain in the Metropolis-Hastings algorithm.

7. The next step is to establish that the matrix Γ defined as

$$\Gamma(\mathbf{X}^{(t+1)}|\mathbf{X}^{(t)}) = A(\mathbf{X}^{(t)}, \mathbf{y})q(\mathbf{y}|\mathbf{X}^{(t)})$$

is irreducible and aperiodic. The Metropolis-Hastings algorithm is very general, allowing fairly arbitrary proposal distribution. Both these properties can be established if $q(\mathbf{y}|\mathbf{X}^{(t)}) > 0$ for all \mathbf{y} , so that $\Gamma(\mathbf{X}^{(t+1)}|\mathbf{X}^{(t)}) > 0$ for all \mathbf{y} . For complex models such as graphical models, the q should be designed so as to make that task of sampling from q efficient. Gibbs sampling is one such example.

8. The final step is in showing that the target distribution \mathbf{p} is invariant for the Markov chain Γ . That is, for every \mathbf{y} ,

$$\sum_{\mathbf{x}} p(\mathbf{x})\Gamma(\mathbf{y}|\mathbf{x}) = p(\mathbf{y})$$

This can be shown by first noting that the Metropolis Hastings Markov chain Γ satisfies the *detailed balance condition* with respect to \mathbf{p} , which means

$$p(\mathbf{x})\Gamma(\mathbf{y}|\mathbf{x}) = p(\mathbf{y})\Gamma(\mathbf{x}|\mathbf{y})$$

This can be proved by simply substituting for $\Gamma(\mathbf{y}|\mathbf{x})$. This leads to the desired result

$$\sum_{\mathbf{x}} p(\mathbf{x})\Gamma(\mathbf{y}|\mathbf{x}) = \sum_{\mathbf{x}} p(\mathbf{y})\Gamma(\mathbf{x}|\mathbf{y}) = p(\mathbf{y})$$

since Γ is row-stochastic.

While the algorithm in Figure 7.12 works in principle, it can be very slow, taking small steps into valleys. In practice, many refinements are made to the algorithm to make it work fast.

Gibbs Sampling

Gibbs sampling is a simple subclass of Metropolis-Hastings algorithm in which the transitions Γ are intuitive and easy to draw from. It is especially relevant for graphical models, which have a large state space of size c^n if each random variable can take c values and the graph has n nodes. Let $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$ be a graph with n nodes. The Gibbs sampling procedure involves two main steps, as outlined in Figure 7.13. It is very natural to think of Gibbs sampling in terms of

```

Input: A target distribution  $p(\mathbf{x}) = \frac{1}{z} \phi(\mathbf{x})$  over a graph  $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$ .
Initialize:  $\mathbf{X}^{(1)} = \mathbf{x}^{(1)}$  for some arbitrary  $q$ .
for step  $t = 1, 2, \dots$  do
  for  $i = 1, \dots, n$  do
    Sample:  $x_i^{(t+1)} \sim p(x_i \mid \mathbf{x}_{1,i-1}^{(t+1)} \cup \mathbf{x}_{i+1,n}^{(t)})$ .
  end for
end for

```

Figure 7.13: The Gibbs sampling algorithm.

the graphical model; the sampling is very much simplified given the conditional independence property - that a node is independent of all other nodes, given its Markov blanket. Thus, each distribution in the sampling step is the probability of a node given its Markov blanket. This procedure depends on the ordering of nodes, though. Though a very popular inference procedure for graphical models, it could take an extremely long time to converge.

For discrete random variables, the graph of configurations is in general a hypercube. While MCMC potentially allows jumps from one node to another in the hypercube, Gibbs sampling restricts every step to be along an edge of the hypercube. This can lead to poor convergence. There are many smarter algorithms that take larger, but calculated steps in the hypercube, leading to better convergence, without incurring excessive computational cost for individual steps.

7.3 Factor Graphs

We have seen that both directed and undirected models can be specified in one of two equivalent ways each; conditional independence and factorization. The semantics for directed and undirected models are however different. While the two specifications are equivalent, factorization can be specified at different levels of granularity or could be defined using different forms of potential functions and is therefore richer. For instance, factorization could be over maximal cliques or over smaller cliques or even over edges, while all these specifications could map to the same conditional independence properties demanding specialization in the conditional independence assertions.

Consider a triangular graph, given by the adjacency matrix

$$A = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

The factorization for this graph could be presented as

$$p(x_1, x_2, x_3) \propto \phi_{123}(x_1, x_2, x_3)$$

However, the following factorization is also a valid one for the graph

$$p(x_1, x_2, x_3) \propto \phi_{12}(x_1, x_2)\phi_{23}(x_2, x_3)\phi_{13}(x_1, x_3)$$

Such problems interest statisticians a lot, since existence of interactions between variables influence diagnosis and/or treatments in the medical domain, for example. Is it possible to graphically differentiate between the two factorizations? Factor graphs precisely serve this purpose. In some sense, there is more information being embedded in the second factorization; the bonafide triplet interaction has been decomposed into pairwise interactions. And this information is represented as ‘factor nodes’ in factor graphs. Corresponding to each factor in the factorization, factor graphs host a node. The nodes in the original graphical model are also retained. But edges are changed; there will be an edge between the factor node and each node whose random variable the factor is a function of.

Definition 55 Given a graphical model $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$, with a factorization $p(\mathbf{x}) = \prod_{a \in \mathcal{F}} \phi_a(\mathbf{x}_a)$ with $\mathcal{F} \subseteq 2^{\mathcal{V}}$, that is compatible with the independence assumptions asserted by the edges \mathcal{E} , the factor graph $\mathcal{G}_f = \langle \mathcal{V} \cup \mathcal{F}, \mathcal{E}_f \rangle$ is defined by $\mathcal{E}_f = \{(q, a) \mid q \in \mathcal{V}, a \in \mathcal{F}, q \in a\}$. The factor graph \mathcal{G}_f is said to encode the factorization of \mathcal{G} .

The set of factor nodes \mathcal{F} is a set of placeholders for particular terms in the factorization. Factor graphs can be looked upon as a ‘graphical way’ of representing hypergraphs (which can have a single edge spanning multiple vertices), where each hyperedge can be thought of as spanning all vertices that figure together in some potential function in the factorization.

As an example, the factor graph for $p(x_1, x_2, x_3) \propto \phi_{123}(x_1, x_2, x_3)$ will be given by $\mathcal{V} = \{1, 2, 3\}$, $\mathcal{F} = \{a\}$ and $\mathcal{E}_f = \{(1, a), (2, a), (3, a)\}$. Whereas, the factor graph for $p(x_1, x_2, x_3) \propto \phi_{12}(x_1, x_2)\phi_{13}(x_1, x_3)\phi_{23}(x_2, x_3)$ will be given by $\mathcal{V} = \{1, 2, 3\}$, $\mathcal{F} = \{a, b, c\}$ and $\mathcal{E}_f = \{(1, a), (2, a), (1, b), (3, b), (2, c), (3, c)\}$.

Any undirected/directed graph without any specialized factorization specified can also be converted into a factor graph. Thus, if $\mathcal{V} = \{1, 2, 3, 4, 5, 6, 7, 8\}$, $\mathcal{E} = \{(1, 2), (2, 3), (3, 4), (2, 5), (3, 4), (4, 5), (4, 7), (5, 7), (5, 6), (6, 8), (7, 8)\}$, then \mathcal{G}_f can be read off the maximal cliques; $\mathcal{F} = \{a, b, c, d, e\}$ and $\mathcal{E}_f = \{(1, a), (2, a), (3, a), (2, b), (5, b), (3, c), (4, d), (7, d), (6, e), (8, e)\}$. As another example, consider a hidden markov model with $\mathcal{V} = \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}$

and $\mathcal{E} = \{(1, 2), (1, 6), (2, 3), (2, 7), (3, 4), (3, 8), (4, 5), (4, 9), (5, 10)\}$. Then the factor graph will have $\mathcal{F} = \{a, b, c, d, p, q, r, s\}$ and $\mathcal{E}_f = \{(1, a), (2, a), (2, b), (3, b), (3, c), (4, c), (4, d), (5, d), (1, p), (6, p)\}$. As a final example, consider a markov decision process, which is a purely directed model and which has five ‘control’ variables in addition to the 10 for the HMM described above (decision process because there is a control that helps you decide the behaviour of the evolution across states in the HMM). It will be specified by $\mathcal{V} = \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15\}$ and $\mathcal{E} = \{(1, 2), (1, 6), (11, 1), (2, 3), (2, 7), (12, 2), (3, 4), (3, 8), (13, 3), (4, 5), (4, 9), (14, 4), (5, 10), (15, 5)\}$. What will be the factor graph corresponding to this graph? Note that, even though there are no directed triangles in this graph, nodes 1 through 5 have two parents each and the corresponding factors are conditional probabilities of the form $p(x_1|x_6, x_{11}), p(x_2|x_7, x_{12}), \text{ etc.}$ Thus, the factor graph will have a node for each such factor connected to three nodes each.

All the factor graph examples considered thus far are factor trees. The sum-product and max-product algorithms are exact on factor trees. Though they assume slightly different forms, the essential ideas are the same. There is no consistent reason for the factor graph being a better representation than the graphical model representation.

7.4 Exponential Family

The exponential family captures many common discrete²¹ and continuous²² graphical model formulations at an abstract level. It provides a rich (though not exhaustive) toolbox of models. The multinomial distribution which models random variables taking k discrete values is what we have looked at so far. Gaussian is one of the most widely used continuous distributions. The poisson distribution helps model distribution over random variables that can take any integral value (as against just k discrete values).

Definition 56 For a given vector of functions $\mathbf{f}(\mathbf{x}) = [f_1(\mathbf{x}), \dots, f_k(\mathbf{x})]$ and a parameter vector $\eta \in \mathbb{R}^k$, the exponential family of distributions is defined as

$$p(\mathbf{x}, \eta) = h(\mathbf{x}) \exp \{ \eta^T \mathbf{f}(\mathbf{x}) - A(\eta) \} \quad (7.43)$$

where the $h(\mathbf{x})$ is a conventional reference function and $A(\eta)$ is the log normalization constant²³ designed as

$$A(\eta) = \log \left[\int_{\mathbf{x} \in \text{Range}(\mathbf{X})} \exp \{ \eta^T \mathbf{f}(\mathbf{x}) \} h(\mathbf{x}) d\mathbf{x} \right]$$

The domain of the parameters η will be restricted to $\text{Domain}(\eta) = \{ \eta \in \mathbb{R}^k \mid A(\eta) < +\infty \}$. $A(\eta)$ plays a fundamental role in hypothesis testing, parameter estimation, etc.,

²¹Density with respect to the Counting measure.

²²Density with respect to the Lebesgue measure.

²³The same as $\log Z$ for graphical models.

though often not very easy to estimate. The central component here is the log-linear form. The parametrization $\eta \in \mathfrak{R}^k$, is also called the canonical parametrization. The function $\mathbf{f}(\mathbf{x})$ is a sufficient statistic function.

As an example, the Gaussian density function can be expressed in the exponential form. If $X \sim \mathcal{N}(\mu, \sigma^2)$ is a univariate Gaussian distribution, then its normal parametrization is

$$p(x, \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{1}{2\sigma^2}(x - \mu)^2 \right\}$$

This density can be manipulated and shown to be a member of the exponential family

$$p(x, \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}} \exp \left\{ \frac{\mu}{\sigma^2}x - \frac{1}{2\sigma^2}x^2 - \frac{1}{2\sigma^2}\mu^2 - \log \sigma \right\} = p(x, \eta)$$

The parameter vector for the exponential form is $\eta = \left[\frac{\mu}{\sigma^2}, -\frac{1}{2\sigma^2} \right]$ while the feature function vector is $\mathbf{f}(x) = [x, x^2]$. The log normalization constant is $A(\eta) = \frac{1}{2\sigma^2}\mu^2 + \log \sigma \equiv -\frac{\eta_1^2}{4\eta_2} + \frac{1}{2} \log(-2\eta_2)$, which determines $\text{Domain}(\eta) = \{\eta \in \mathfrak{R}^2 \mid \eta_2 < 0\}$. Finally $h(\mathbf{x}) = \frac{1}{\sqrt{2\pi}}$ for this example. The number of degrees of freedom (reflected through two parameters in the moment parametrization) will precisely be the number of canonical parameters. The canonical parametrization extended²⁴ to the multivariate Gaussian counterpart will be discussed in Section 7.5.

As a second example, consider the bernoulli distribution. A bernoulli random variable $X \in \{0, 1\}$ can model the outcome of any coin-flipping experiment. It can be expressed as

$$p(x, \mu) = \mu^x(1 - \mu)^{1-x}$$

where μ is the probability of $X = 1$. Rearranging the terms, we get the exponential form for bernoulli distribution as

$$p(x, \mu) = \exp \left\{ \left(\log \frac{\mu}{1 - \mu} \right) x + \log(1 - \mu) \right\} = p(x, \eta)$$

with parameter vector specified as $\eta = \left[\log \frac{\mu}{1 - \mu} \right]$ which gives μ as a logistic function (log of likelihood ratio or log-odds ratio) of η_1 and $\mathbf{f}(x) = [x]$. The log normalization constant is $A(\eta) = -\log(1 - \mu) \equiv \log\{1 + e^{\eta_1}\}$ ²⁵ while $h(x) = 1$. Also, $\text{Domain}(\eta) = \mathfrak{R}$. In general, if you started coupling together multiple distributions and try expressing them in exponential form, determining η could become a very hard problem.

As a third example, consider the poisson²⁶ random variable $X \in \mathcal{N}$ (set of natural numbers). Its density function is given by

$$p(x, \mu) = \frac{\mu^x e^{-\mu}}{x!}$$

²⁴EXERCISE.

²⁵EXERCISE.

²⁶When London was being bombed in World war 2, experts were trying to model where the bomb would next fall using a 2D poisson distribution.

Poisson distributions are often used to model events, such as the ringing of a telephone. The mean parameter μ controls the density or frequency with which the event is happening. The canonical parametrization for this distribution can be obtained using a routine rewrite as

$$p(x, \mu) = \frac{1}{x!} \exp\{(\log \mu)x - \mu\}$$

where $h(x) = \frac{1}{x!}$, $\eta = [\log \mu]$, $\mathbf{f}(x) = [x]$, $A(\eta) = \mu = e^{\eta}$ with $\text{Domain}(\eta) = \Re$.

In all examples considered so far, we algebraically converted the density function from a moment parametrization to a canonical representation as a member of the exponential family. How about going backwards from a canonical form to moment parametrization? The vector of moment parameters is defined as $\mu = [\mu_1, \mu_2, \dots, \mu_k] = [E[f_1(\mathbf{X})], E[f_2(\mathbf{X})], \dots, E[f_k(\mathbf{X})]]$. That is, we can get the moment parameters by taking expectations of the sufficient statistics $f_i(\mathbf{X})$ that sit in the exponent of the canonical form. And the expectation of the i^{th} component of this function can be proved to be the same as $\frac{\partial A}{\partial \eta_i}$. That is

$$\nabla A(\eta)_i = \frac{\partial A}{\partial \eta_i} = E[\mathbf{f}(\mathbf{X})] = \mu_i \quad (7.44)$$

Further, it can be proved that

$$\nabla^2 A(\eta)_{ij} = \frac{\partial^2 A}{\partial \eta_i \partial \eta_j} = \text{cov}\{f_i(\mathbf{X}), f_j(\mathbf{X})\} = E[(f_i(\mathbf{X}) - \mu_i)(f_j(\mathbf{X}) - \mu_j)] = \mu_{ij}$$

The proof of the two above statements are straightforward and use the property that $A(\eta)$ is infinitely differentiable. To illustrate this, we will revisit the canonical parametrization for the Gaussian example.

$$p(x, \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}} \exp\left\{\frac{\mu}{\sigma^2}x - \frac{1}{2\sigma^2}x^2 - \frac{1}{2\sigma^2}\mu^2 - \log \sigma\right\} = p(x, \eta)$$

The moments can be derived to me $\mu_1 = E[x] = \mu$ and $\mu_2 = E[x^2] = \sigma^2 + \mu$. Similarly, for the poisson distribution, the moment parameter is simply $[\mu]$ as also $\text{var}(X) = \mu$. For the bernoulli distribution, $E(X) = \mu$ and $\text{var}(X) = (1 - \mu)\mu$.

7.5 A Look at Modeling Continuous Random Variables

A normal or gaussian distribution is one of the most widely (ab)used probability distributions. They come up in the central limit theorem in the sums of independent or weakly dependent random variables, whose normalized sum converges to a gaussian (justifying their wide use). They are very often used

to model noise. Example graphical models that use gaussian distribution are Gaussian graphical models, Gauss-Markov time series, *etc.*

Another reason for their wide use is computational. In the case of continuous random variables, messages are functions rather than vectors. In general, this can often force us to quantize the messages. But in the case of gaussians (and in general for exponential models which we will shortly see), the message passing can be performed in terms of sufficient statistics. Kalman filters are a special case of message passing that pass sufficient statistics as messages.

The form of the density function for Gaussian distribution, with $\mathbf{x}, \mu \in \Re^n$ and $\Sigma \in \Re^{n \times n}$ is

$$p(\mathbf{x}, \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{1}{2}} \sqrt{\det(\Sigma)}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu) \right\} \quad (7.45)$$

where μ is the mean vector and $\Sigma \succ \mathbf{0}$ is the covariance matrix²⁷. It can be verified that μ is indeed the mean vector; $\mu = E[\mathbf{X}] = \int_{\Re^n} \mathbf{x} p(\mathbf{x}, \mu, \Sigma) d\mathbf{x}$. Similarly, it can be verified that $\Sigma = E[(\mathbf{X} - \mu)(\mathbf{X} - \mu)^T]$. The quantity $\frac{1}{(2\pi)^{\frac{1}{2}} \sqrt{\det(\Sigma)}}$ is the normalization constant and can be computed as $\int_{\Re^n} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu) \right\} d\mathbf{x}$. If $n = 1$, the integral is easy to compute. For computing integrals for multivariate problem, it is a good idea to reduce the matrix Σ by diagonalization.

What sits in the exponent of the density function is a quadratic term. The contours of constant probability are ellipsoids, with shape determined by Σ^{-1} and center as μ (which does not affect the shape of the ellipsoid). For example, if Σ^{-1} is diagonal, the axes of the ellipsoid will be aligned along the coordinate axes.

The parametrization (μ, Σ) is called the *moment parametrization* of the Gaussian; μ is the first order moment of \mathbf{x} while Σ is the matrix of second order centered moment. There is another *canonical parametrization* of the Gaussian which is related to the sum-product algorithm. The parameters are a new vector and a new matrix:

$$\eta = \Sigma^{-1} \mu$$

and a new matrix

$$\Omega = \Sigma^{-1}$$

With a bit of algebra, the Gaussian density can be re-expressed in terms of the canonical parameters as:

$$p(\mathbf{x}, \eta, \Omega) = \exp \left\{ \eta^T \mathbf{x} - \frac{1}{2} \mathbf{x}^T \Omega \mathbf{x} + \mathbf{x} \right\} \quad (7.46)$$

²⁷Recall that a matrix is positive definite if all its eigenvalues are strictly positive or equivalently, $\forall \mathbf{z} \neq \mathbf{0}, \mathbf{z}^T \Sigma \mathbf{z} > 0$.

where

$$\mathbf{x} = -\frac{1}{2} \{n \log(2\pi) - \log |\Omega| + \eta^T \Omega^{-1} \eta\}$$

is a constant, analogous to the normalization constant in the moment parametrization of the Gaussian density in (7.45). The parametrization for (7.45) is more naturally suited for graphical models, because, as we will see, the canonical parameters can be directly related to the compatibility functions when you start breaking up the Gaussian into a graphical model-like factorization. Typically, graphical model factorizations invoke canonical parameters. The above parametrization is often referred to as $\mathcal{N}(\eta, \Omega)$.

Equation (7.46) gives a quadratic form in the exponent and is a special case of the exponential family representation, which happens to be a very general concept. Note that all logarithms are to the base e . The quadratic term in the exponent can be rewritten as a trace:

$$\mathbf{x}^T \Omega \mathbf{x} = \text{trace}(\Omega \mathbf{x} \mathbf{x}^T)$$

where $\text{trace}(A)$ is the sum of the diagonal entries of A and happens to be linear. The Gaussian density using this transformation is obtained in its log-normal form, which has an exponent linear in its features of \mathbf{x} and $\mathbf{x} \mathbf{x}^T$.

It can be shown that linear functions of Gaussian random vectors are Gaussian; $\sum_{i=0}^n Y_i \sim \mathcal{N}(0, (n+1)\sigma^2)$ if each $Y_i \sim \mathcal{N}(0, \sigma^2)$. Also, if $X_i = \sum_{j=0}^n Y_j$ then $p(x_{i+1} | x_i) = \exp\{-\frac{1}{2\sigma^2}(x_{i+1} - x_i)^2\}$ and the random variables X_i form a Markov chain with factorization

$$p(\mathbf{x}) = \exp\left\{-\frac{1}{2\sigma^2}x_0^2\right\} \prod_{i=1}^n \exp\left\{-\frac{1}{2\sigma^2}(x_{i+1} - x_i)^2\right\}$$

We can verify that for this Markov chain, $E[X_1] = E_{X_0}[E_{X_1|X_0}[X_1 | X_0]] = E_{X_0}[X_0 + E[Y_0]] = 0$ and in general $E[X_i] = 0$. Also, $E[X_1 X_2] = 2\sigma^2$.

The canonical parameters for $p(\mathbf{x})$ can be read off the factorization as $\eta = \mathbf{0}$ and a tridiagonal, sparse

$$\Omega = \begin{bmatrix} 2 & -1 & 0 & \dots & 0 \\ -1 & 2 & -1 & \dots & 0 \\ 0 & -1 & 2 & \dots & 0 \\ 0 & 0 & -1 & \dots & 0 \\ \cdot & \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \cdot & \dots & -1 \\ 0 & 0 & 0 & \dots & 2 \end{bmatrix}$$

The matrix Ω is sparse, because the graph is sparse (just a markov chain). Missing edges translate to zeros; for any graph, if there is no edge between X_i and X_j , $\Omega_{ij} = 0$.

Factor analysis is another example application of gaussian density functions. Principal component analysis and Karhunen Loeve transform are limiting cases of factor analysis. The motivation here is dimensionality reduction for cases when $\mathbf{Y} \in \mathfrak{R}^n$ and n is very large, such as the size of an image for a naive $1 - d$ raster scanned pixel vector representation. The complete vector need not really capture the intrinsic dimensionality of the object being described (such as a face). Factor analysis is useful if you think that the data should be lying in some lower ($d \ll n$) dimensional space²⁸. Let $\omega_1, \omega_2, \dots, \omega_d \in \mathfrak{R}^n$ be d (linearly independent) basis vectors. Consider the linear subspace $M = \left\{ \mathbf{y} \in \mathfrak{R}^n \mid \mathbf{y} = \sum_{i=1}^d x_i \omega_i, x_i \in \mathfrak{R}, \omega_i \in \mathfrak{R}^n \right\}$. Factor analysis tries to induce a probabilistic distribution over the subspace M , by defining $\mathbf{X} \sim \mathcal{N}(\mathbf{0}_d, I_{d \times d})$ and

$$\mathcal{Y}_{n \times 1} = \mu_{n \times 1} + \Omega_{n \times d} \mathbf{X}_{d \times 1} + \mathbf{K}_{n \times 1}$$

where the i^{th} column of Ω is ω_i , $\mathbf{K} \sim \mathcal{N}(\mathbf{0}, D)$ is gaussian noise (capturing the assumption that the model may not be exactly correct) and $\mu \in \mathfrak{R}^d$ is the shift (in subspace \mathfrak{R}^n) vector. Though very naive, it has not stopped researchers and especially practitioners from using it successfully. The components ω_i 's are extracted from a large database using eigenvalue analysis (such as eigenface analysis in image processing literature).

The graphical model representation for factor analysis is very simple; $\mathcal{V} = \{X_1, X_2, \dots, X_d, \mathbf{Y}\}$ and $\mathcal{E} = \{(X_1, \mathbf{Y}), (X_2, \mathbf{Y}), \dots, (X_d, \mathbf{Y})\}$. The X_i 's are marginally independent as indicated by \mathcal{G} . Further, we can infer that $E[\mathbf{Y}] = \mu$ and

$$\Sigma = E[(\mathbf{Y} - \mu)(\mathbf{Y} - \mu)^T] = E[(\Omega \mathbf{X} + \mathbf{K})(\Omega \mathbf{X} + \mathbf{K})^T] = \Omega \Omega^T + D$$

Finally, you can show that

$$\begin{pmatrix} \mathbf{X}_{d \times 1} \\ \mathbf{Y}_{n \times 1} \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \mathbf{0}_{d \times 1} \\ \mu_{n \times 1} \end{pmatrix}, \begin{pmatrix} I_{d \times d} & \Omega_{d \times n}^T \\ \Omega_{n \times d} & (\Omega \Omega^T + D)_{n \times n} \end{pmatrix} \right)$$

In practice, it is useful to infer a distribution for \mathbf{X} (distribution on weights for different eigenfaces) used to synthesize a particular observation $\mathbf{Y} = \hat{\mathbf{y}}$ (such as a face image). That is, it can be required to infer $p(\mathbf{X} \mid \hat{\mathbf{y}})$. Fortunately²⁹, this turns out to be a Gaussian and this can be inferred using the Bayes rule and

²⁸There is a lot of interest these days in identifying the manifold (especially non-linear subspaces, such as spheres) in which the data lies. In the classical technique of factor analysis, we make a very restrictive assumption that the data lies in some *linear subspace*.

²⁹The conditional distribution need not always stay in the same family. But for Gaussians, the conditional distribution stays in the same family.

algebraic operations on matrices. In particular, the following property turns out to be useful. If

$$\begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \right)$$

then, $\mu_{\mathbf{X}|\mathbf{y}} = \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(\mathbf{y} - \mu_2)$, $Var(\mathbf{X} | \mathbf{y}) = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}$ and $Var(\mathbf{X} | \mathbf{y}) \preceq \Sigma_{11} = Var(\mathbf{X})$. The last expression indicates that observation over \mathbf{Y} should reduce uncertainty over \mathbf{X} .

7.6 Exponential Family and Graphical Models

We will now discuss the exponential family in the setting of graphical models. Particularly, we will discuss how to stitch together exponential models on graphs. Many, if not all graphical models take an exponential form. We will mainly discuss undirected graphs. Recall from (7.12), the standard factorization for an undirected graphical model

$$p(\mathbf{x}) = \frac{1}{Z} \prod_{\mathcal{C} \in \Pi} \phi_{\mathcal{C}}(\mathbf{x}_{\mathcal{C}})$$

where $\phi_{\mathcal{C}}$ is a compatibility or potential function and Π is the set of all maximal cliques. All we do in exponential form is rewrite it in the log form to get an additive decomposition. This does not always hold, since we will need that all $\phi_{\mathcal{C}}$ are non-negative in order to take their logs.

$$p(\mathbf{x}) = \exp \left\{ \sum_{\mathcal{C} \in \Pi} \log \phi_{\mathcal{C}}(\mathbf{x}_{\mathcal{C}}) - \log Z \right\}$$

The above additive decomposition is essentially in exponential form, though not quite in the exponential family form (we still need to specify the canonical parameters). As an example, consider a Gaussian graphical model with $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \Omega)$ (canonical parametrization) defined on a tree $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$ with noisy observations $\mathcal{Y} = C\mathcal{X} + \mathcal{V}$, $\mathcal{V} \sim \mathcal{N}(\mathbf{0}, R)$. Then the conditional distribution $p(\mathbf{X} | \mathbf{Y} = \mathbf{y})$ is also Gaussian, which can be decomposed according to the tree structure of \mathcal{G} .

$$p(\mathbf{x}|\mathbf{y}) \propto \prod_{s \in \mathcal{V}} p(y_s|x_s) \prod_{(s,t) \in \mathcal{E}} \exp \left\{ -\frac{1}{2} [x_s \ x_t] J_{st} [x_s \ x_t]^T \right\}$$

where,

$$J_{st} = \begin{bmatrix} \Omega_{s(t)} & \Omega_{ts} \\ \Omega_{st} & \Omega_{t(s)} \end{bmatrix}$$

with $\Omega_{s(t)}$ specified so that $\Omega_{ss} = \sum_{t \in \mathcal{N}(s)} \Omega_{s(t)}$. In this representation, it is possible to do message passing using the canonical parameters in the sum/product updates (the basis of Kalman filters).

You could also have bernoulli random variables at every node (for example, to model the spread of a contagious disease in a social network of people) and one might be interested in building a probability distribution over the social network. Then, with every node $X_s \sim \text{Ber}(\lambda_s) \in \{0, 1\}$ (canonical parametrization) and choosing³⁰ $\phi_{st} = \exp\{\lambda_{st} f_{st}(x_s, x_t)\}$ where, $f_{st} = x_s x_t + (1 - x_s)(1 - x_t)$ defined so that it takes value 1 iff $x_s = x_t$ (both are diseased or healthy) and is 0 otherwise, we have the following factorization:

$$p(\mathbf{x}) = \prod_{s \in \mathcal{V}} \exp\{\lambda_s x_s\} \prod_{(s,t) \in \mathcal{E}} \exp\{\lambda_{st} f_{st}(x_s, x_t)\} \quad (7.47)$$

If $\lambda_{st} = 0$, it indicates no coupling between the related individuals. For a reasonably contagious disease, we expect λ_{st} to be non-negative atleast. The value of θ_s indicates the belief that an individual X_s was initially diseased; higher the value, more will be the belief. Though not directly reflecting the marginal probability, λ_s is an indicator of the log-odds ratio.

Let us consider another instance - a problem of counting the vertex coverings in a graph using the sum-product formalism. We will associate a random variable $X_i \in \{0, 1\}$ with each node in the graph. Then, the following exponential form for $p(\mathbf{x})$ will serve our purpose:

$$p(\mathbf{x}, \eta) = \exp\left\{\eta \sum_{i=1}^n x_i - A(\eta)\right\} \prod_{(s,t) \in \mathcal{E}} (1 - (1 - x_s)(1 - x_t))$$

where $h(\mathbf{x}) = \prod_{(s,t) \in \mathcal{E}} (1 - (1 - x_s)(1 - x_t))$ ensures that we indeed have a vertex cover and $\mathbf{f}(\mathbf{x}) = \sum_{i=1}^n x_i$. It can be proved that as $\eta \rightarrow -\infty$ (which means you are increasingly penalizing larger coverings), the distribution goes to the minimum cardinality vertex covering (*i.e.*, $\sum_{i=1}^n x_i$).

If a graphical model were to be used to represent a language model for spell-checking or validity, *etc.* of an input string \mathbf{x} of length upto n , you can have a substantially fewer number of parameters than if you were to have a naive potential over all 26 characters in an alphabet leading to a table of size 26^n . This parameter reduction can be achieved by having indicator feature functions

³⁰Remember that in a graphical model, we are free to choose the form of the potential functions so that they match the semantics. Here we pick edge potentials that reflect the coupling between health of related individuals, thereby shaping the distribution.

$\mathbf{f}(\mathbf{x})$ corresponding to interesting (and not all) substrings such as ‘ion’, ‘ing’. To model such apriori information about ‘striking patterns’, it is useful to think of graphical models in the following *reduced parametrization* form, where feature function f_α can represent some such information as a feature function:

$$p(\mathbf{x}, \eta) = h(\mathbf{x}) \exp \left\{ \sum_{\mathcal{C} \in \Pi} \sum_{\alpha \in I_{\mathcal{C}}} \eta_\alpha f_\alpha(\mathbf{x}_{\mathcal{C}}) - A(\eta) \right\} \quad (7.48)$$

where $I_{\mathcal{C}}$ is the index for clique \mathcal{C} so that $f_\alpha(\mathbf{x}_{\mathcal{C}})$ corresponds to only those interesting features that are local in terms of clique \mathcal{C} .

7.6.1 Exponential Models and Maximum Entropy

The data fitting principle of maximum entropy, which is suitable for learning models from data leads naturally to graphical models in exponential form and also gives nice semantics to the weights in the exponential family. There are many problems with constraints on distributions, but where the information is not complete. For instance, if we knew that for two binary random variables $X, Y \in \{0, 1\}$ and for their paired observations, $2/3$ times X takes value 0, $3/4$ times, Y takes value 1 and if you are asked, to specify the fraction of times $(X, Y) = (0, 0)$, what would you answer with the insufficient specification? Or going back to our diseased network model, if you are given observations on healths of similarly networked people and were to translate these observations into constraints on the moments (and/or the joint quantities), how would you determine the parameters of the probability distribution? There could be many distributions/parameters that are consistent with the observations. What principle should be adopted for making a good choice of the parameters/distribution?

More concretely, say you are given some (feature) functions $f_\alpha(\mathbf{x})$ with $\alpha \in I$ and are also given some observations $\hat{\mu}_\alpha$ on the expected values of the functions, called *empirical moments*. The observations could come either from data or from physical constraints. We are interested in the family of distributions that are consistent with these constraints.

$$\mathcal{P}(\hat{\mu}) = \left\{ p(\cdot) \left| \sum_{\mathbf{x}} p(\mathbf{x}) f_\alpha(\mathbf{x}) = \hat{\mu}_\alpha \quad \forall \alpha \in I \right. \right\}$$

Note that the family could be an empty set if the set of constraints are inconsistent. However, we will assume that the constraints are consistent. We are interested in choosing a particular \hat{p} from $\mathcal{P}(\mu)$ in a principled manner. Maximum entropy is one such principled method. Entropy is, roughly speaking, a measure of uncertainty or randomness. It has played a vital role in physics, chemistry, statistics, computer science, communication³¹, etc.

³¹Entropy plays a fundamental role in deciding how far you could compress a sequence of bits.

Definition 57 The entropy $H(p)$ of the distribution p on a random variable (vector) \mathbf{X} is given by the expected value of the log of the distribution. Depending on whether the random variable (vector) is continuous or discrete, we will have two different definitions of expectation (and therefore for the entropy). For discrete random variable (vector) \mathbf{X} ,

$$H(p) = - \sum_{\mathbf{x}} p(\mathbf{x}) \log p(\mathbf{x})$$

whereas, for continuous valued random variable \mathbf{X} ,

$$H(p) = - \int_{\mathbf{x}} p(\mathbf{x}) \log p(\mathbf{x}) d\mathbf{x}$$

It can be easily proved that $H(p) \geq 0$ (convention being that $0 \log 0 = 0$). $H(p) = 0$ if and only if X is deterministically always some a , making $p(a) = 1$ and $p(x) = 0, \forall x \neq a$. $H(p)$ is maximal for the uniform distribution.

The principle of maximum entropy can be now defined:

Definition 58 Given $\mathcal{P}(\hat{\mu}) = \left\{ p(\cdot) \left| \sum_{\mathbf{x}} p(\mathbf{x}) f_{\alpha}(\mathbf{x}) = \hat{\mu}_{\alpha} \quad \forall \alpha \in I \right. \right\}$, choose

$$\hat{p} = \operatorname{argmax}_{p \in \mathcal{P}(\hat{\mu})} H(p)$$

The intuition here is to balance across the constraints. It is also called the *principle of least commitment* since the goal is to simultaneously respect the data and not commit to anything more.

Theorem 99 The maximum entropy solution exists and is unique. The unique solution takes the form

$$p(\mathbf{x}) \propto \exp \left\{ \sum_{\alpha \in I} \eta_{\alpha} f_{\alpha}(\mathbf{x}) \right\} \quad (7.49)$$

Proof: The first thing to note is that $\mathcal{P}(\hat{\mu})$ is a convex³² (linear), closed³³ and bounded set of distributions. Further, $H(p)$ is continuous, and this, in conjunction with the nature of $\mathcal{P}(\hat{\mu})$, guarantees that a maximum will be attained. Also, $H(p)$ is strictly concave³⁴, implying that the maximum will be unique.

The canonical parameters are actually the Lagrange multipliers. Consider the Lagrangian $L(p, \eta, \lambda)$:

$$L(p, \eta, \lambda) = H(p) + \sum_{\alpha \in I} \eta_{\alpha} \left[\hat{\mu}_{\alpha} - \sum_{\mathbf{x}} p(\mathbf{x}) f_{\alpha}(\mathbf{x}) \right] + \lambda \left[1 - \sum_{\mathbf{x}} p(\mathbf{x}) \right]$$

³²What if the set is empty?

³³Which means the optimum cannot escape to the boundary.

³⁴Since $\frac{\partial^2 H}{\partial p^2} = -\frac{1}{p} < 0$.

The KKT necessary and sufficient conditions (see (4.88) on page 296) for optimality of $(\widehat{p}(\mathbf{x}), \widehat{\eta}, \widehat{\lambda})$ yield:

1. $\frac{\partial L}{\partial p} = 0 \Rightarrow \log \widehat{p}(\mathbf{x}) - \sum_{\alpha \in I} \widehat{\eta}_\alpha f_\alpha(\mathbf{x}) - \widehat{\lambda} = 0$. That is,

$$\widehat{p}(\mathbf{x}) = \exp \left\{ \sum_{\alpha \in I} \widehat{\eta}_\alpha f_\alpha(\mathbf{x}) \right\} e^{\widehat{\lambda}} \quad (7.50)$$

2. $\sum_{\mathbf{x}} \widehat{p}(\mathbf{x}) = 1$. Substituting for $p(\mathbf{x})$ from (7.50), we get

$$e^{\widehat{\lambda}} = \frac{1}{\sum_{\mathbf{x}} \exp \left\{ \sum_{\alpha \in I} \widehat{\eta}_\alpha f_\alpha(\mathbf{x}) \right\}}$$

which is a constant.

This proves that $p(\mathbf{x}) \propto \exp \left\{ \sum_{\alpha \in I} \eta_\alpha f_\alpha(\mathbf{x}) \right\}$. \square

This result shows how exponential families can arise in a fairly natural way. It can be further shown that a slight generalization of the maximum entropy principle, called the *Kullback-Leibler divergence* very naturally leads to the maximum likelihood principle.

Definition 59 *The Kullback-Leibler (KL) divergence $D(p||q)$ between two distributions $p(\cdot)$ and $q(\cdot)$ is given by the expectation over $p(\cdot)$ of the log-likelihood ratio of $p(\cdot)$ over $q(\cdot)$.*

$$D(p||q) = \sum_{\mathbf{x}} p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})}$$

For distributions $p(\cdot)$ and $q(\cdot)$ over continuous valued random variables

$$D(p||q) = \int_{\mathbf{x}} p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}$$

Like any distance, the KL divergence is always non-negative. Also, $p \equiv q$ if and only if $D(p||q) = 0$. However, by inspection, we can infer that $D(p||q)$ is asymmetric, unlike metrics. That is $D(p||q) \neq D(q||p)$. If $q(\cdot)$ were the uniform distribution $D(p||q) = H(p) + c$, where c is some constant.

We will next define a slight generalization of the maximum entropy principle. Recall the definition of $\mathcal{P}(\widehat{\mu})$. Now let q be some additional prior or reference distribution. The generalized definition goes as:

Definition 60 Given $\mathcal{P}(\widehat{\mu}) = \left\{ p(\cdot) \left| \sum_{\mathbf{x}} p(\mathbf{x}) f_{\alpha}(\mathbf{x}) = \widehat{\mu}_{\alpha} \quad \forall \alpha \in I \right. \right\}$ and a reference distribution $q(\mathbf{x})$, choose

$$\widehat{p} = \underset{p \in \mathcal{P}(\widehat{\mu})}{\operatorname{argmin}} D(p||q)$$

The interpretation here is: get as close as possible to the reference distribution, while respecting the data in the form of the constraint set $\mathcal{P}(\widehat{\mu})$. For the maximum entropy principle, the reference distribution happened to be the uniform distribution. Note that we have an ‘argmin’ here instead of an ‘argmax’.

Theorem 100 The solution to the minimum KL divergence problem in definition 60 exists and is unique. The unique solution takes the form

$$p(\mathbf{x}) \propto q(\mathbf{x}) \exp \left\{ \sum_{\alpha \in I} \eta_{\alpha} f_{\alpha}(\mathbf{x}) \right\}$$

The reference distribution $q(\mathbf{x})$ plays the role of the function $h(\mathbf{x})$ in the exponential form.

Proof: The proof of existence and uniqueness here are the same as that for the counterpart in the proof of theorem 99. The only change will be the KKT necessary and sufficient conditions for optimality of $(\widehat{p}(\mathbf{x}), \widehat{\eta}, \widehat{\lambda})$ (see (4.88) on page 296).

Consider the Lagrangian $L(p, \eta, \lambda)$:

$$L(p, \eta, \lambda) = D(p||q) + \sum_{\alpha \in I} \eta_{\alpha} \left[\widehat{\mu}_{\alpha} - \sum_{\mathbf{x}} p(\mathbf{x}) f_{\alpha}(\mathbf{x}) \right] + \lambda [1 - \sum_{\mathbf{x}} p(\mathbf{x})]$$

The KKT necessary and sufficient conditions yield:

1. $\frac{\partial L}{\partial p} = 0 \Rightarrow \log \widehat{p}(\mathbf{x}) - \sum_{\alpha \in I} \widehat{\eta}_{\alpha} f_{\alpha}(\mathbf{x}) - \widehat{\lambda} - \log q(\mathbf{x}) = 0$. That is,

$$\widehat{p}(\mathbf{x}) = q(\mathbf{x}) \exp \left\{ \sum_{\alpha \in I} \widehat{\eta}_{\alpha} f_{\alpha}(\mathbf{x}) \right\} e^{\widehat{\lambda}} \quad (7.51)$$

2. $\sum_{\mathbf{x}} \widehat{p}(\mathbf{x}) = 1$. Substituting for $p(\mathbf{x})$ from (7.50), we get

$$e^{\widehat{\lambda}} = \frac{1}{\sum_{\mathbf{x}} q(\mathbf{x}) \exp \left\{ \sum_{\alpha \in I} \widehat{\eta}_{\alpha} f_{\alpha}(\mathbf{x}) \right\}}$$

which is a constant.

This proves that $p(\mathbf{x}) \propto q(\mathbf{x}) \exp \left\{ \sum_{\alpha \in I} \eta_{\alpha} f_{\alpha}(\mathbf{x}) \right\}$. \square

7.7 Model fitting

Thus far we have discussed graphical models and inference in these models. But how do we learn the potential functions associated with a model or the canonical parameters associated with the exponential form, given some data from the real world that describes the phenomena? This opens up a new interesting fundamental question: ‘how to infer models from data?’. We will also need a measure of how ‘good’ a model is and this is often driven by the end goal. These discussions form the topic of discussion for this section.

Consider a coin tossing experiment, in which $X \sim \text{Ber}(\mu)$ where $\mu = \Pr(X = 1)$. Let us say we observe a sequence of coin tosses: $[x^{(1)}, x^{(2)}, \dots, x^{(m)}]$. Can we use this data to infer something about μ ? There are two primary ways to do this.

1. The Bayesian school of thought views the Bayes rule as the fundamental method for inferring model from the data:

$$p(\theta | \bar{\mathbf{x}}) = \frac{p(\bar{\mathbf{x}} | \theta)p(\theta)}{p(\bar{\mathbf{x}})}$$

where θ is the underlying model parameter, $p(\bar{\mathbf{x}} | \theta)$ is the likelihood term, $p(\theta | \bar{\mathbf{x}})$ is the posterior that summarizes the remaining uncertainty in θ , given that you have observed the data $\bar{\mathbf{x}}$. This simple statement has consequences on the way you might view parameters; by imposing a distribution over θ , you are encoding the belief that the parameter itself is a random quantity. This requires viewing all parameters as random variables. The Bayesian technique views μ as a random quantity following a $\text{Beta}(a, b)$ distribution³⁵ on $[0, 1]$.

$$\mu \sim \text{Beta}(a, b) \propto \mu^{a-1}(1 - \mu)^{b-1}$$

The Bayesian tries to model the inherent uncertainty in the statistician about the value of the parameter (μ). The Bayesian thus has a more or less fixed procedure for folding in the data into the model.

2. The frequentist’s way of measuring probabilities is as numbers measured as outcomes over repeated trials as against the subjective notion of probability adopted by the Bayesians. The frequentist would object to the Bayesian view on several counts: (i) the subjectivity of the procedure; is there a sound justification for the choice of $\mu \sim \text{Beta}(a, b)$ and for the particular choice of a and b ? (ii) that different results for the same data set could be obtained by different statisticians adhering to different choices of the prior. The frequentist belief is that one should be completely objective with the data; the data should speak so that you get the same model, no matter who builds the model. It does not make sense for the probability

³⁵The shape of the Beta distribution depends on the value of the parameters - it could be either uniform or bell-shaped.

of a coin throwing up heads to be a random variable. However, the frequentist does not have any fixed procedure for inventing a model from the data. The frequentist thus considers several estimators for the parameters. One estimator is the sample mean of the data. Any estimator is assessed for its properties such as bias, variability, *etc.* Variability asks: How much would the estimated value vary if the data sample changed? How will the estimate behave if we repeated the experiment several times?

Both Bayesian and frequentist views are compatible with graphical models. Also, the two views agree that there could be different models for different (coin-tossing) experiments. The Bayesian view often gives you a good way of generating good procedures (that is, procedures with good properties) whereas the frequentist view often gives you a very good way of evaluating a procedure.

7.7.1 Sufficient Statistics

Let \mathbf{X} be a random variable. In the present discussion, we will often assume it to correspond to a vector of observations (that is data). Any function $\tau(\mathbf{X})$ is called a statistic. From the Bayesian point of view, a statistic $\tau(\mathbf{X})$ is *sufficient* for the parameter variable θ if $\theta \perp \mathbf{X} \mid \tau(\mathbf{X})$. This relation can be expressed graphically as a markov chain with $\mathcal{V} = \{\mathbf{X}, \tau(\mathbf{X}), \theta\}$, $\mathcal{E} = \{(\mathbf{X}, \tau(\mathbf{X})), (\tau(\mathbf{X}), \theta)$. Intuitively, this means that the function of the observations $\tau(\mathbf{X})$ is what is essential in data to explain the characteristics of θ and that all the dependence between \mathbf{X} and θ is being mediated by $\tau(\mathbf{X})$. $\tau(\mathbf{X})$ is typically much smaller than \mathbf{X} itself.

From the frequentist point of view, θ is an unknown fixed parametrization that we would like to estimate. The sufficiency for a frequentist is that

$$p(\mathbf{x} \mid \tau(\mathbf{x}); \theta) = p(\mathbf{x} \mid \tau(\mathbf{x}))$$

That is, the family of distributions specified by the parametrization θ becomes independent of θ when conditioned on the sufficient statistic $\tau(\mathbf{x})$; it indicates that we have conditioned on sufficient information to capture any information from θ .

A view of sufficiency unified across the Bayesian and frequentist perspectives can be obtained by treating it as a statement about factorization:

$$p(\mathbf{x}, \tau(\mathbf{x}), \theta) = \phi_1(\tau(\mathbf{x}), \theta) \phi_2(\tau(\mathbf{x}), \mathbf{x})$$

It can be proved that this is indeed a unified view of sufficiency, consistent with the Bayesian and frequentist views. Though there is a difference in interpretations, it amounts to the same factorization. Note the similarity of the factorization here with that for general graphical models. As we will see, for the graphical model family, sufficiency properties can be read out, just on the basis of their factorization. Further, the sufficiency property stands out clearly for exponential family. Given iid data³⁶ $\bar{\mathbf{x}} = (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)})$ each from some

³⁶This assumption is common-place because it makes computations very simple and decomposable. It is often not the case, since there could be dependence in the sampling.

exponential family, and given an η (which corresponds to the θ being discussed thus far in the context of model estimation), we can easily infer that

$$p(\bar{\mathbf{x}}_{1,m}; \eta) = \prod_{i=1}^m p(\mathbf{x}^{(i)}; \eta) = \underbrace{\left(\prod_{i=1}^m h(\mathbf{x}^{(i)}) \right)}_{\phi_2(\tau(\mathbf{x}), \mathbf{x})} \underbrace{\exp \left\{ \eta^T \left(\sum_{i=1}^m \mathbf{f}(\mathbf{x}^{(i)}) \right) - mA(\eta) \right\}}_{\phi_1(\tau(\mathbf{x}), \eta)}$$

This algebra tells us that the quantity $\tau(\bar{\mathbf{x}}) = \sum_{i=1}^m \mathbf{f}(\mathbf{x}^{(i)})$ is the sufficient statistic. Note that ϕ_1 can potentially depend on $\tau(\mathbf{x})$, though it does not depend in this case. The key to estimating η are the sufficient statistics $\tau(\bar{\mathbf{x}})$. We will also realize that the prefactor $h(\bar{\mathbf{x}})$ does not play a significant role here.

The quantity

$$\bar{\boldsymbol{\mu}} = \frac{1}{m} \sum_{i=1}^m \mathbf{f}(\mathbf{x}^{(i)}) = \sum_{\mathbf{x}} \wp(\mathbf{x}) \mathbf{f}(\mathbf{x})$$

is called the *empirical moment parameter* vector, where the empirical distribution $\wp(\mathbf{x})$ is obtained by placing point mass at data points.

$$\wp(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m \delta(\mathbf{x} - \mathbf{x}^{(i)})$$

The empirical moment parameters are a special type of empirical moments discussed on page 439. Note that the empirical moment parameters are simply the sufficient statistics scaled down by the number of data points. The empirical moment parameter vector can be contrasted against the moment parameter vector $\boldsymbol{\mu} = [E[f_1(\mathbf{X})], E[f_2(\mathbf{X})], \dots, E[f_k(\mathbf{X})]]$ for exponential families, as defined on page 433. The two differ in the probability distributions with respect to which they compute their expectations; while the former computes with respect to the empirical distribution, the latter computes with respect to the true distribution.

As an example, consider the Ising model, defined on a graph $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$, which is a member of the exponential family. Its distribution is very similar to that of the disease network model (7.47), only difference is that only the first term $x_s x_t$ is retained in the definition of $f_{st}(\mathbf{x})$.

$$p(\mathbf{x}, \eta) \propto \exp \left\{ \sum_{s \in \mathcal{V}} \eta_s x_s + \sum_{(s,t) \in \mathcal{E}} \eta_{st} x_s x_t \right\}$$

By appropriately identifying the feature function vector \mathbf{f} and η of size $|\mathcal{V}| + |\mathcal{E}|$ each, we can determine the empirical moments to be

$$\hat{\boldsymbol{\mu}}_s = \frac{1}{m} \sum_{i=1}^m x_s^i$$

and

$$\hat{\mu}_{st} = \frac{1}{m} \sum_{i=1}^m x_s^i x_t^i$$

which are just the marginal counts.

What about empirical moments for the reduced parametrization discussed earlier in (7.48) on page 439?

$$p(\mathbf{x}, \eta) = h(\mathbf{x}) \exp \left\{ \sum_{\mathcal{C} \in \Pi} \sum_{\alpha \in I_{\mathcal{C}}} \eta_{\alpha} f_{\alpha}(\mathbf{x}_{\mathcal{C}}) - A(\eta) \right\}$$

Given iid data $\bar{\mathbf{x}} = (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)})$, the following factorization holds if each $\mathbf{x}^{(i)}$ follows the distribution specified by the reduced parametrization. Since the a priori information is captured using indicative feature functions, you can have a reduced set of sufficient statistics.

$$p(\bar{\mathbf{x}}_{1,m}; \eta) = \prod_{i=1}^m p(\mathbf{x}^{(i)}; \eta) = \underbrace{\left(\prod_{i=1}^m h(\mathbf{x}^{(i)}) \right)}_{\phi_2(\tau(\mathbf{x}), \mathbf{x})} \exp \left\{ \underbrace{\sum_{\alpha} \eta_{\alpha} \left(\sum_{i=1}^m \sum_{\mathcal{C} \in \Pi | \alpha \in I_{\mathcal{C}}} f_{\alpha}(\mathbf{x}^{(i)}) \right)}_{\phi_1(\tau(\mathbf{x}), \eta)} - mA(\eta) \right\}$$

The form of the empirical moments in this case pools over all cliques that are relevant for a particular feature type f_{α} (for instance, in a grid, some features may be commonly defined/parametrized across all vertical edges while others across all horizontal edges) and then pools over all data:

$$\hat{\mu}_{\alpha} = \sum_{i=1}^m \sum_{\mathcal{C} \in \Pi | \alpha \in I_{\mathcal{C}}} f_{\alpha}(\mathbf{x}^{(i)})$$

The reduced parametrization assumes some kind of homogeneity in the model. More precisely, it makes the statistical assumption that the parameter η_{α} is independent of the exact position in \mathcal{G} and is determined by the local graphical structure. This parametrization assuming homogeneity is often desirable since it gets a more stable estimator with less variance even using relatively less data. In fact, the reduced parametrization yields a new exponential family with much lower dimensions, in which the parameters η_{α} enter in a purely linear way³⁷. The issues of when and how to pool data are important ones in modeling.

7.7.2 Maximum Likelihood

The likelihood function interests both Bayesian and frequentists alike. Recall that likelihood $p(\bar{\mathbf{x}} | \theta)$ was one part of the Bayes rule. The frequentist interpretation of this quantity is as ‘the likelihood of a fixed $\bar{\mathbf{x}}$ as θ varies’ whereas the Bayesian interpretation of this quantity is as ‘the conditional probability of

³⁷A widely used different class of exponential families is called *curved exponential families*, in which the parameters η_{α} enter in non-linear ways.

a varying $\bar{\mathbf{x}}$ given a fixed θ .³⁸ The frequentist thus views likelihood as a function of theta $L(\theta)$ that measures, in some sense, the goodness of θ as it is varied for the particular sample of data in hand. It is often useful to talk about the log-likelihood instead³⁸.

Definition 61 Given a sample $\bar{\mathbf{x}} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)}\}$, the log-likelihood (LL) is defined as

$$LL(\theta; \bar{\mathbf{x}}) = \frac{1}{m} \log p(\bar{\mathbf{x}} | \theta)$$

Note that the $\frac{1}{m}$ term does not really change the optimization and its usage is conventionally just to normalize the log-likelihood.

The maximum likelihood estimate (MLE) is defined next.

Definition 62 Given a sample $\bar{\mathbf{x}} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)}\}$, the maximum likelihood estimate (MLE) $\hat{\theta}_{MLE}$ of θ is defined as

$$\hat{\theta}_{MLE} = \underset{\theta}{\operatorname{argmax}} LL(\theta; \bar{\mathbf{x}})$$

As the name suggests, this principle treats log-likelihood as a measure of goodness of the parameter and picks that value which maximizes the LL. Though not against Bayesian principles, MLE has been of greater interest to the frequentists. The estimate $\hat{\theta}_{MLE}$ is called a point estimate, since the method gives you just a single point.

Let us try to fit (that is, adjust the parameters of) a scalar Gaussian $\mathcal{N}(\mu, \sigma^2)$ to some data $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ using the MLE principle. We will assume that the data is iid. This will mean that the joint distribution over the data will factorize into individual data instances, given the parameters.

$$LL(\mu; \bar{\mathbf{x}}) = \frac{1}{m} \sum_{i=1}^m \log p(x^{(i)} | \mu) = -\frac{1}{2m} \log 2\pi - \frac{1}{2m} \sum_{i=1}^m (x^{(i)} - \mu)^2$$

The LL is a quadratic in μ (that is θ), which achieves its maximum at

$$\hat{\mu}_{MLE} = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

as expected. The data determines the shape of the likelihood and MLE looks for the point $\hat{\mu}_{MLE}$ that maximizes the LL.

This example of parameter estimation for the simple Gaussian is one of the motivations for using MLE; MLE corresponds to an intuitively appealing estimation for the mean parameter. Secondly MLE has good asymptotic properties.

³⁸Since it is a monotonic transformation between the likelihood and log-likelihood, it does not really matter much whether we look at the former or the latter. It is usually more convenient to look at the log-likelihood (LL).

A frequentist takes more interest in the asymptotic properties of the estimator. We can look upon $\hat{\theta}_{MLE}$ as a random variable, since the data $\bar{\mathbf{x}}$ was itself random and $\hat{\theta}_{MLE}$ is a function of the data. As the amount of data grows, any reasonable estimator should behave ‘well’. One yardstick of good behaviour of the estimator is *consistency*, that is, if there was some true underlying θ^* , we would like

$$\hat{\theta}_{MLE} \xrightarrow{m \rightarrow \infty} \theta^*$$

Under most conditions, MLE estimators are consistent. Generally, it can be worrisome if the estimator does not converge to the true answer even with infinite amount of data.

Another yardstick of good behaviour is that asymptotically, the estimator’s spread around the true value of the parameter must be Gaussian-like. MLE estimators honor this feature as well; with increasing m , the distribution of $\hat{\theta}_{MLE}$ tends to be Gaussian with the true mean θ^* and the spread Σ given by the Fisher information matrix³⁹.

In small sample regimes (that is, if you do not have too much data), MLE does not behave well. In such settings, the frequentist adds a penalty term, called the *regularization* term (such as μ^2 in the above example) to the likelihood objective to somehow prevent the estimator from drifting away based on a small m .

7.7.3 What the Bayesians Do

The Bayesians focus more on the posterior $p(\theta \mid \bar{\mathbf{x}})$, which is the likelihood multiplied by a prior.

$$p(\theta \mid \bar{\mathbf{x}}) \propto p(\bar{\mathbf{x}} \mid \theta)p(\theta)$$

For example, let $x^{(i)} \sim \mathcal{N}(\mu, 1)$ for $1 \leq i \leq m$ and let $\mu \sim \mathcal{N}(\nu, \sigma^2)$. The parameter μ is a random variable, whereas the *hyperparameters* ν and σ are fixed. The choice of the values ν and σ could significantly influence the posterior. There is a class of models⁴⁰ called the ‘hierarchical Bayes models’ that put a prior on the hyper-parameters, priors again on the hyper-hyper-parameters and so on. It stops (for pragmatic reasons) when the abstraction has been sufficiently performed so that it matters little what the fixed values of the hyperⁿ-parameters are.

The posterior for this example, takes the form

$$p(\theta \mid \bar{\mathbf{x}}_{1,m}) \propto \exp \left\{ -\frac{1}{2} \sum_{i=1}^m (x^{(i)} - \mu)^2 \right\} \exp \left\{ -\frac{1}{2\sigma^2} (\mu - \nu)^2 \right\}$$

The Bayesian is not interested just in the point estimate, but moreso in the entire posterior distribution. In this case (a special instance of Kalman filters) the posterior again turns out to be a Gaussian; $\mu_{\bar{\mathbf{x}}_{1,m}} \sim \mathcal{N}(E(\mu \mid \bar{\mathbf{x}}_{1,m}), Var(\mu \mid \bar{\mathbf{x}}_{1,m}))$.

³⁹Related to this is the Cramer-Rao lower bound.

⁴⁰The broader area is called *Robust Bayesian statistics*.

A little calculation should convince you that $\mu_{\bar{\mathbf{x}}_{1,m}} \sim \mathcal{N}\left(\frac{m\sigma^2}{m\sigma^2+1}\hat{\mu}_{MLE} + \frac{1}{m\sigma^2+1}\nu, \text{Var}(\mu | \bar{\mathbf{x}}_{1,m})\right)$. How does the posterior behave as $m \rightarrow \infty$? We can easily see that with increasing amount of data, the likelihood will completely swamp the prior. That is, $E(\mu | \bar{\mathbf{x}}_{1,m}) \xrightarrow{m \rightarrow \infty} \hat{\mu}_{MLE}$. So the choice of the hyper-parameters is irrelevant as the size of the data goes to ∞ . Thus, MLE behaves well asymptotically even from the Bayesian point of view. Unlike the frequentist, the Bayesian does not need to explicitly handle the small sample regimes.

In general, the posterior might not have a finite parametrization. In such cases, the Bayesians do compute point estimates. Examples of point estimates computed by Bayesians are:

1. *Bayes estimate*: This is the mean of the posterior:

$$\hat{\theta}_{Bayes} = \int \theta p(\theta | \bar{\mathbf{x}}_{1,m}) d\theta$$

2. *MAP estimate*: This is the mode of the posterior:

$$\hat{\theta}_{MAP} = \underset{\theta}{\operatorname{argmax}} \frac{1}{m} \log \{p(\theta | \bar{\mathbf{x}}_{1,m})\}$$

The MAP estimate is very much related to the MLE. Invoking the Bayes rule and ignoring the term involving $\log p(\bar{\mathbf{x}}_{1,m})$,

$$\begin{aligned} \hat{\theta}_{MAP} &= \underset{\theta}{\operatorname{argmax}} \frac{1}{m} \log \{p(\bar{\mathbf{x}}_{1,m} | \theta)\} + \frac{1}{m} \log \{p(\theta)\} & (7.52) \\ &\stackrel{\text{iid}}{=} \underset{\theta}{\operatorname{argmax}} \frac{1}{m} \sum_{i=1}^m \log \{p(x^{(i)} | \theta)\} + \frac{1}{m} \log \{p(\theta)\} \end{aligned}$$

The first term in the decomposition is the data likelihood term, while the second term is the prior term. Let us study this decomposition.

- (a) This decomposition suggests that if the prior is uniform, $\hat{\theta}_{MAP} = \hat{\theta}_{MLE}$.
- (b) Secondly, in the asymptotic case, as $m \rightarrow \infty$, the prior component fades away (since it has no dependence on m in the numerator) in contrast to the likelihood term⁴¹. In fact, as $m \rightarrow \infty$, $\hat{\theta}_{MAP} \rightarrow \hat{\theta}_{MLE}$.
- (c) If the prior is assumed to be Gaussian, then the prior term will assume the form of a quadratic penalty. This is the same as a quadratic regularization term for MLE. With different choices of priors on θ , you get different regularizations. For example, a Laplacian prior on θ results in $L1$ regularization. For example, Lasso is $L1$ regularization for a regression problem. How the choice of regularization affects these estimators is an area of active research.

⁴¹One could use the central limit theorem or law of large numbers to study the behaviour of the likelihood term as $m \rightarrow \infty$.

7.7.4 Model Fitting for Exponential Family

We will initiate the discussion of model fitting in exponential families by looking at maximum likelihood (ML) estimation. Just as for most other properties of exponential models, there is something crisper that we can state about the ML properties of exponential models. The optimization conditions will turn out to take special and easily interpretable forms.

Recall the specification of any member of the exponential family from (7.43).

$$p(\mathbf{x} \mid \eta) = h(\mathbf{x}) \exp \{ \eta^T \mathbf{f}(\mathbf{x}) - A(\eta) \}$$

We will now see how the empirical moment parameters, defined on page 445 and discussed subsequently become relevant for parameter estimation. Given iid observations $\bar{\mathbf{x}} = \{ \mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)} \}$, the normalized LL for this distribution decouples as a sum

$$LL(\eta; \bar{\mathbf{x}}) = \frac{1}{m} \log p(\bar{\mathbf{x}} \mid \eta) \tag{7.53}$$

$$= \frac{1}{m} \left\{ \eta^T \sum_{i=1}^m \mathbf{f}(\mathbf{x}^{(i)}) - mA(\eta) \right\} + \log h(\bar{\mathbf{x}}) \tag{7.54}$$

$$= \left\{ \frac{1}{m} \eta^T \sum_{i=1}^m \mathbf{f}(\mathbf{x}^{(i)}) \right\} - A(\eta) + \log h(\bar{\mathbf{x}}) \tag{7.55}$$

$$= \{ \eta^T \bar{\boldsymbol{\mu}}(\bar{\mathbf{x}}) \} - A(\eta) + \log h(\bar{\mathbf{x}}) \tag{7.56}$$

where $\bar{\boldsymbol{\mu}}(\bar{\mathbf{x}}) = \frac{1}{m} \sum_{i=1}^m \mathbf{f}(\mathbf{x}^{(i)})$ is the vector of empirical moment parameters. Since $h(\bar{\mathbf{x}})$ is not a function of the parameter (and instead, is a constant offset), it follows that the MLE for the exponential family takes the form

$$\hat{\boldsymbol{\mu}}_{ML} \in \underset{\eta \in \text{Domain}(\eta)}{\text{argmax}} \{ \eta^T \bar{\boldsymbol{\mu}}(\bar{\mathbf{x}}) - A(\eta) \}$$

This is a very crisp formulation of sufficiency; given data, all you need to compute are $\bar{\boldsymbol{\mu}}(\bar{\mathbf{x}})$, forget the data $\bar{\mathbf{x}}$ and focus on solving the optimization problem. Since $A(\eta)$ is infinitely differentiable and since the remainder of the objective is linear, we could make use of the first order necessary optimality conditions from (4.44) on page 270:

$$\nabla LL(\eta; \bar{\mathbf{x}}) = 0$$

That is,

$$\bar{\boldsymbol{\mu}}(\bar{\mathbf{x}}) - \nabla A(\eta) = \bar{\boldsymbol{\mu}}(\bar{\mathbf{x}}) - E(\mathbf{f}(\mathbf{X})) = 0$$

where, we may recall from (7.44) on page 433 that $\nabla A(\eta)$ is the vector of moments predicted by the model. This important condition

$$\bar{\mu}(\bar{\mathbf{x}}) = E(\mathbf{f}(\mathbf{X})) \quad (7.57)$$

is called the *moment matching condition*⁴² and it states that for maximizing the likelihood, we need to match the empirical moments to the model moments. In general, these coupled d equations are complex and non-linear ones, that are not very easy to solve, though there are some exceptions. The necessary conditions suggest that the parameters should be learnt in such a way that if you drew samples from the model, they should look like the given data.

As example, let X be a bernoulli random variable. Recollect the specification of the bernoulli distribution as a member of exponential family from page 7.4: $\eta = \log \frac{\mu}{1-\mu}$, $f(x) = x$ and $A(\eta) = \log \{1 + e^\eta\}$. This leads to the following moment matching conditions

$$\bar{\mu} = \frac{1}{m} \sum_{i=1}^m x^{(i)} = A'(\eta) = \frac{e^\eta}{1 + e^\eta}$$

The empirical moments is the fraction of heads that come up in the m coin tosses of the experiment. The expression for η will be

$$\eta = \log \left(\frac{\bar{\mu}}{1 - \bar{\mu}} \right)$$

which corresponds to the logs-odd ratio. When does this have a solution? It has a solution if $\bar{\mu} \in (0, 1)$ but not if $i = 1$ or $i = 0$, which can happen if all the coin tosses landed up with heads or tails⁴³. If $\bar{\mu} \in (0, \frac{1}{2})$, $\eta < 0$ and if $\bar{\mu} \in (\frac{1}{2}, 1)$, $\eta > 0$. So if $\bar{\mu} = 1$, strictly speaking the MLE does not exist. Note that if we were to use the original formulation of the bernoulli distribution on page 432, moment matching would have yielded the MLE as $\mu = \bar{\mu}$. For large amounts of data, using the weak law of large numbers, you can be convinced that this loss can be recovered using even the exponential formulation.

As another example, let us revisit the multivariate normal distribution in the canonical form first uncovered in (7.46):

$$p(\mathbf{x}, \eta, \Omega) = \exp \left\{ \eta^T \mathbf{x} - \frac{1}{2} \text{trace}(\Omega \mathbf{x} \mathbf{x}^T) + \mathbf{x} \right\}$$

where

$$\eta = \Sigma^{-1} \mu$$

and

$$\Omega = \Sigma^{-1}$$

⁴²This is related to moment matching rooted in classical statistics. Though not related directly to maximum likelihood, these two boil down to the same criteria for the exponential family.

⁴³This is not impossible, especially if the tossing is done by some magician such as Persi Warren Diaconis or by a machine built by him.

Moment matching will yield the model first order and second order moments in terms of the empirical mean and the empirical second order moment matrix respectively. That is,

$$E[\mathbf{x}] = \mu = \frac{1}{m} \sum_{i=1}^m \mathbf{x}^{(i)} = \bar{\mu}$$

and

$$E[\mathbf{x}\mathbf{x}^T] = \Sigma + \mu\mu^T = \frac{1}{m} \sum_{i=1}^m \mathbf{x}^{(i)} \left(\mathbf{x}^{(i)}\right)^T$$

Solving this system, we get Ω_{ML} as the inverse of the sample covariance matrix

$$\Omega_{ML} = \left(\frac{1}{m} \sum_{i=1}^m \mathbf{x}^{(i)} \left(\mathbf{x}^{(i)}\right)^T - \bar{\mu}\bar{\mu}^T \right)^{-1}$$

(which exists if the sample covariance matrix has full rank) and

$$\eta_{ML} = \Omega_{ML}\bar{\mu}$$

In practice (such as in text mining problems), it can happen that many features are never observed and consequently, the corresponding empirical moment parameters will be 0. Similarly, for the multivariate Gaussian example, if you do not have sufficient data, the sample covariance matrix may not have full rank. More precisely, we need to pay heed to $Domain(\eta)$, since the optimal values could reside on the boundary. This is addressed using regularization through a penalized log-likelihood objective.

7.7.5 Maximum Likelihood for Graphical Models

For general large graphical models, the moment matching is not at all easy to solve. It has therefore been a practice to often resort to iterative algorithms for solving the set of moment matching equations. We will illustrate this difficulty through the Ising model (*c.f.* page 445), which is a model on an undirected graph $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$.

$$p(\mathbf{x}, \eta) = \exp \left\{ \sum_{s \in \mathcal{V}} \eta_s x_s + \sum_{(s,t) \in \mathcal{E}} \eta_{st} x_s x_t - A(\eta) \right\}$$

Let us say we have iid data $\bar{\mathbf{x}} = (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)})$. The empirical moments are

$$\hat{\mu}_s = \frac{1}{m} \sum_{i=1}^m x_s^i$$

and

$$\hat{\mu}_{st} = \frac{1}{m} \sum_{i=1}^m x_s^i x_t^i$$

The $d = |\mathcal{V}| + |\mathcal{E}|$ moment matching equations will be

$$\sum_{\mathbf{x} \in \{0,1\}^{|\mathcal{V}|}} \exp \left\{ \sum_{s \in \mathcal{V}} \eta_s x_s + \sum_{(s,t) \in \mathcal{E}} \eta_{st} x_s x_t - A(\eta) \right\} x_s = \hat{\mu}_s$$

and

$$\sum_{\mathbf{x} \in \{0,1\}^{|\mathcal{V}|}} \exp \left\{ \sum_{s \in \mathcal{V}} \eta_s x_s + \sum_{(s,t) \in \mathcal{E}} \eta_{st} x_s x_t - A(\eta) \right\} x_s x_t = \hat{\mu}_{st}$$

The task is to estimate the values of the $|\mathcal{V}| + |\mathcal{E}|$ sized vector η . The log-normalization constant $A(\eta)$ is in turn very complex and involves

$$A(\eta) = \sum_{\mathbf{x} \in \{0,1\}^{|\mathcal{V}|}} \exp \left\{ \sum_{s \in \mathcal{V}} \eta_s x_s + \sum_{(s,t) \in \mathcal{E}} \eta_{st} x_s x_t - A(\eta) \right\}$$

In general, though the equations are intuitive (asserting that parameters be picked to match the data), solving them is very very complex. What if the graph were a tree? Then the solution can be obtained efficiently, as we will see subsequently. Also, $A(\eta)$ as well as the marginal $p(\mathbf{x}; \eta)$ could be computed in polynomial time leading to efficient inference as well. Thus, solving the moment matching equations is closely linked to the graph structure. In general, solving the estimation problem inevitably involves solution to the marginalization problem, which can at least be performed efficiently for tree structured models.

Let us consider some graphical models, whose estimation problems are easy. Consider a markov chain with $\mathcal{V} = \{X_1, X_2, X_3, X_4\}$, $\mathcal{E} = \{(X_1, X_2), (X_2, X_3), (X_3, X_4)\}$. Though it can be thought of as a directed model, we will choose an exponential factorization that will link it to the general exponential machinery. Let each X_i have k possible states. Also, let

$$\lambda_{st} = \sum_{i=1}^k \sum_{j=1}^k \lambda_{st,ij} (\delta(x_s, i) \delta(x_t, j))$$

where $\delta(x, i) \delta(y, j)$ is an indicator feature function for each fixed (i, j) pair and is expressible as a $k \times k$ matrix. $\lambda_{st,ij}$'s are the canonical parameters. We can then adopt the exponential form

$$p(\mathbf{x}; \lambda) = \exp \{ \lambda_{12}(x_1, x_2) + \lambda_{23}(x_2, x_3) + \lambda_{34}(x_3, x_4) \}$$

Given data $\bar{\mathbf{x}} = (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)})$, the empirical moment parameters (sufficient statistics scaled down by m) can be computed to be

$$\bar{\mu}_{st}(x_s, x_t) = \frac{1}{m} \sum_{i=1}^m \delta(x_s = x_s^{(i)}) \delta(x_t = x_t^{(i)})$$

for $(s, t) \in \{(1, 2), (2, 3), (3, 4)\}$ and

$$\bar{\mu}_s(x_s) = \frac{1}{m} \sum_{i=1}^m \delta(x_s = x_s^{(i)})$$

for $s \in \{1, 2, 3, 4\}$.

These simply correspond to the empirical marginal probability $\wp(x_s, x_t)$. The moment matching conditions can be satisfied⁴⁴ by the following assignments to the canonical parameters:

$$\begin{aligned}\hat{\lambda}_{12,12} &= \log \bar{\mu}_{12}(x_1, x_2) \\ \hat{\lambda}_{23,23} &= \log \frac{\bar{\mu}_{23}(x_2, x_3)}{\bar{\mu}_2(x_2)} \\ \hat{\lambda}_{34,34} &= \log \frac{\bar{\mu}_{34}(x_3, x_4)}{\bar{\mu}_3(x_3)}\end{aligned}$$

Thus, the canonical parameters can be computed in closed form. The task is similarly simple for trees - the basic intuition has been captured in this markov chain example. However, the simple solution suggested in this four node example does not hold for a cyclic graph having 4 nodes with edges defined as $\mathcal{E} = \{(X_1, X_2), (X_2, X_3), (X_3, X_4), (X_4, X_1)\}$.

The task is also equally simple, with closed form solutions for the class of decomposable graphs (which could be graphs with cycles), defined on page 415 in definition 54.

7.7.6 Iterative Algorithms

We will now move to more general graphs and describe iterative algorithms for solving fixed point equations. We will assume that the potential functions are defined on maximal cliques.

$$p(\mathbf{x}, \phi) \propto \prod_{C \in \Pi} \phi_C(\mathbf{x}_C) \equiv \exp \left\{ \sum_{C \in \Pi} \theta_C(\mathbf{x}_C) \right\}$$

The general setting for maximum likelihood optimization is that we have to solve a set of fixed point equations

$$\mathbf{F}(\eta) = E_\eta[\mathbf{f}(\mathbf{X})] - \bar{\mu} = \mathbf{0}$$

(such as $\hat{\mu}_\alpha = \sum_{i=1}^m \sum_{C \in \Pi | \alpha \in I_C} f_\alpha(\mathbf{x}^{(i)})$ in the case of reduced parametrization - *c.f.* page 446). For this system, we will investigate an iterative solution of the form

$$\eta^{(t+1)} = \eta^{(t)} + f(\eta^{(t)})$$

so that at the fixed point $\eta^{(t^*)}$, $f(\eta^{(t^*)}) = 0$ and therefore $\eta^{(t^*+1)} = \eta^{(t^*)}$. t is the time step.

This general algorithm is called Iterative Proportional Fitting (IPF) and goes as follows for a general undirected graph $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$:

⁴⁴EXERCISE: Prove.

1. Start with some initial factorization (which could be uniform)

$$p(\mathbf{x}) \propto \prod_{\mathcal{C} \in \Pi} \phi_{\mathcal{C}}^{(0)}(\mathbf{x}_{\mathcal{C}}) \equiv \exp \left\{ \sum_{\mathcal{C} \in \Pi} \theta_{\mathcal{C}}^{(0)}(\mathbf{x}_{\mathcal{C}}) \right\}$$

2. For $t = 1$ onwards, let \mathcal{C}' range stepwise over the cliques in Π . Update

$$\phi_{\mathcal{C}'}^{(t+1)}(\mathbf{x}_{\mathcal{C}'}') = \phi_{\mathcal{C}'}^{(t)}(\mathbf{x}_{\mathcal{C}'}') \frac{\bar{\mu}_{\mathcal{C}'}(\mathbf{x}_{\mathcal{C}'}')}{\mu_{\mathcal{C}'}^{(t)}(\mathbf{x}_{\mathcal{C}'}')} \quad (7.58)$$

where, $\mu_{\mathcal{C}'}^{(t)}(\mathbf{x}_{\mathcal{C}'}')$ is the current model marginal, computed as⁴⁵

$$\mu_{\mathcal{C}'}^{(t)}(\mathbf{x}_{\mathcal{C}'}') = \frac{1}{Z_{\phi^{(t)}}} \sum_{\mathbf{x}_{\mathcal{V} \setminus \mathcal{C}'}} p(\mathbf{x}; \phi^{(t)})$$

While the model marginals can be computed efficiently for tree structured or even decomposable graphs using message passing formalisms, we may have to resort to approximate inferencing for general graphs. $\bar{\mu}_{\mathcal{C}'}$ are the empirical marginals that are precomputed from data $\bar{\mathbf{x}}$ as

$$\bar{\mu}_{\mathcal{C}'}(\mathbf{x}_{\mathcal{C}'}') = \frac{1}{m} \sum_{i=1}^m \delta(\mathbf{x}_{\mathcal{C}'}', \mathbf{x}_{\mathcal{C}'}^{(i)})$$

Equivalently one may also use the following update rule:

$$\theta_{\mathcal{C}'}^{(t+1)}(\mathbf{x}_{\mathcal{C}'}') = \theta_{\mathcal{C}'}^{(t)}(\mathbf{x}_{\mathcal{C}'}') + \log \frac{\bar{\mu}_{\mathcal{C}'}(\mathbf{x}_{\mathcal{C}'}')}{\mu_{\mathcal{C}'}^{(t)}(\mathbf{x}_{\mathcal{C}'}')} \quad (7.59)$$

The algorithm is called *iterative proportional fitting* because at every step, the potentials are updated proportional to how well the empirical moment parameters fit the model moments for a clique \mathcal{C}' .

Convergence of the iterative algorithm

It is easy to see that at the fix point $t = t^*$, the algorithm will yield the MLE $\phi^{(t^*)}$ since, for all $\mathcal{C} \in \Pi$,

$$\mathbf{F}(\phi^{(t^*)}) = \mu_{\mathcal{C}}^{(t^*)}(\mathbf{x}_{\mathcal{C}}) - \bar{\mu}_{\mathcal{C}}(\mathbf{x}_{\mathcal{C}}) = 0$$

⁴⁵You could formulate this problem using a set of $\delta_{\mathcal{C}}$ feature functions, and $\lambda_{\mathcal{C}, \mathbf{v}_{\mathcal{C}}}$ canonical parameters, as was adopted in the example on page 453. Here, $\mathbf{v}_{\mathcal{C}}$ is a configuration vector that could be assumed by variables on clique \mathcal{C} .

But we also need to show that the updates will always converge to the fix point. We first observe that after updating using step (7.59) (or equivalently, (7.58)), the moment matching will be achieved for clique \mathcal{C}' .

$$\mu_{\mathcal{C}'}^{(t+1)}(\mathbf{x}_{\mathcal{C}'}) = \bar{\mu}_{\mathcal{C}'}(\mathbf{x}_{\mathcal{C}'})$$

Why is this so? By definition,

$$\mu_{\mathcal{C}'}^{(t+1)}(\mathbf{x}_{\mathcal{C}'}) = \frac{1}{Z_{\theta^{(t+1)}}} \sum_{\mathbf{x}_{\mathcal{V} \setminus \mathcal{C}'}} p(\mathbf{x}; \theta^{(t+1)}) = \frac{1}{Z_{\theta^{(t+1)}}} \sum_{\mathbf{x}_{\mathcal{V} \setminus \mathcal{C}'}} \exp \left\{ \theta_{\mathcal{C}'}^{(t)}(\mathbf{x}_{\mathcal{C}'}) + \sum_{\mathcal{C} \in \Pi, \mathcal{C} \neq \mathcal{C}'} \theta_{\mathcal{C}}^{(t)}(\mathbf{x}_{\mathcal{C}}) \right\}$$

That is, every factor that is not in \mathcal{C}' is not changing with respect to the previous step. Expanding upon this, using the update equation (7.59),

$$\mu_{\mathcal{C}'}^{(t+1)}(\mathbf{x}_{\mathcal{C}'}) = \frac{1}{Z_{\theta^{(t+1)}}} \sum_{\mathbf{x}_{\mathcal{V} \setminus \mathcal{C}'}} \exp \left\{ \theta_{\mathcal{C}'}^{(t)}(\mathbf{x}_{\mathcal{C}'}) + \log \frac{\bar{\mu}_{\mathcal{C}'}(\mathbf{x}_{\mathcal{C}'})}{\mu_{\mathcal{C}'}^{(t)}(\mathbf{x}_{\mathcal{C}'})} + \sum_{\mathcal{C} \in \Pi, \mathcal{C} \neq \mathcal{C}'} \theta_{\mathcal{C}}^{(t)}(\mathbf{x}_{\mathcal{C}}) \right\} = \frac{1}{Z_{\theta^{(t+1)}}} \frac{\bar{\mu}_{\mathcal{C}'}(\mathbf{x}_{\mathcal{C}'})}{\mu_{\mathcal{C}'}^{(t)}(\mathbf{x}_{\mathcal{C}'})} \sum_{\mathbf{x}_{\mathcal{V} \setminus \mathcal{C}'}} \exp \left\{ \theta_{\mathcal{C}'}^{(t)}(\mathbf{x}_{\mathcal{C}'}) + \sum_{\mathcal{C} \in \Pi, \mathcal{C} \neq \mathcal{C}'} \theta_{\mathcal{C}}^{(t)}(\mathbf{x}_{\mathcal{C}}) \right\}$$

since $\bar{\mu}_{\mathcal{C}'}(\mathbf{x}_{\mathcal{C}'})$ and $\mu_{\mathcal{C}'}^{(t)}(\mathbf{x}_{\mathcal{C}'})$ are independent of the inner summation. This leads to

$$\mu_{\mathcal{C}'}^{(t+1)}(\mathbf{x}_{\mathcal{C}'}) = \frac{Z_{\theta^{(t)}}}{Z_{\theta^{(t+1)}}} \frac{\bar{\mu}_{\mathcal{C}'}(\mathbf{x}_{\mathcal{C}'})}{\mu_{\mathcal{C}'}^{(t)}(\mathbf{x}_{\mathcal{C}'})} \sum_{\mathbf{x}_{\mathcal{V} \setminus \mathcal{C}'}} p(\mathbf{x}_{\mathcal{C}'}; \theta^{(t)})$$

By definition,

$$\mu_{\mathcal{C}'}^{(t)}(\mathbf{x}_{\mathcal{C}'}) = \sum_{\mathbf{x}_{\mathcal{V} \setminus \mathcal{C}'}} p(\mathbf{x}_{\mathcal{C}'}; \theta^{(t)})$$

Thus,

$$\mu_{\mathcal{C}'}^{(t+1)}(\mathbf{x}_{\mathcal{C}'}) = \frac{Z_{\theta^{(t)}}}{Z_{\theta^{(t+1)}}} \bar{\mu}_{\mathcal{C}'}(\mathbf{x}_{\mathcal{C}'})$$

Since both empirical and model moments are expected to sum to 1 across all cliques in the graph, we should have

$$1 = \sum_{\mathcal{C}' \in \Pi} \mu_{\mathcal{C}'}^{(t+1)}(\mathbf{x}_{\mathcal{C}'}) = \sum_{\mathcal{C}' \in \Pi} \frac{Z_{\theta^{(t)}}}{Z_{\theta^{(t+1)}}} \bar{\mu}_{\mathcal{C}'}(\mathbf{x}_{\mathcal{C}'}) = \frac{Z_{\theta^{(t)}}}{Z_{\theta^{(t+1)}}}$$

Consequently, we will have that once the update step (7.59) is applied, moment matching will hold on \mathcal{C}' .

$$\mu_{\mathcal{C}'}^{(t+1)}(\mathbf{x}_{\mathcal{C}'}) = \bar{\mu}_{\mathcal{C}'}(\mathbf{x}_{\mathcal{C}'})$$

This takes us one step forward. But are these step updates guaranteed to lead to convergence? To see that convergence holds, we will note that the IPF algorithm is an instance of the coordinate-wise ascent algorithm (*c.f.* page 301). This becomes obvious by stating the log-likelihood objective

$$LL(\theta; \bar{\mathbf{x}}) = \frac{1}{m} \left(\sum_{\mathcal{C} \in \Pi} \theta_{\mathcal{C}} \bar{\mu}_{\mathcal{C}}(\mathbf{x}_{\mathcal{C}}) \right) - A(\{\theta_{\mathcal{C}} \mid \mathcal{C} \in \Pi\})$$

and viewing it as a function of $\theta_{\mathcal{C}'}$

$$g(\theta_{\mathcal{C}'}) = LL\left(\theta_{\mathcal{C}'}; \left\{\theta_{\mathcal{C}} = \theta_{\mathcal{C}}^{(t)} \mid \mathcal{C} \neq \mathcal{C}'\right\}, \bar{\mathbf{x}}\right)$$

The first order necessary optimality condition for $g(\theta_{\mathcal{C}'})$ precisely yield the moment matching equation for \mathcal{C}' . Since $g(\theta_{\mathcal{C}'})$ can be shown to be concave, the first order necessary conditions are also sufficient. The fact that the application of coordinate descent in this setting will lead to convergence follows from the fact that $LL(\theta; \bar{\mathbf{x}})$ is strongly concave.

7.7.7 Maximum Entropy Revisted

The IPF algorithm exploited the idea of matching moments in order to maximize likelihood. Recall from Section 7.6.1, the maximum entropy principle that seeks to find $\varphi(\cdot) \in \mathcal{P}(\hat{\boldsymbol{\mu}})$ that maximizes $H(p)$. \mathcal{P} is the set of all distributions that satisfy empirical moment constraints $\hat{\boldsymbol{\mu}}_{\alpha}$:

$$\mathcal{P}(\hat{\boldsymbol{\mu}}) = \left\{ p(\cdot) \mid \sum_{\mathbf{x}} p(\mathbf{x}) f_{\alpha}(\mathbf{x}) = \hat{\boldsymbol{\mu}}_{\alpha} \quad \forall \alpha \in I \right\}$$

The problem of maximum entropy was also shown to be equivalent to the problem of minimizing the KL divergence between p and the reference uniform distribution u :

$$\hat{p}_{ME} = \operatorname{argmax}_{p \in \mathcal{P}(\hat{\boldsymbol{\mu}})} H(p) = \operatorname{argmin}_{p \in \mathcal{P}(\hat{\boldsymbol{\mu}})} D(p||u) \quad (7.60)$$

On the other hand, the problem of maximum likelihood estimation problem is specified as

$$\hat{p}_{MLE} = \operatorname{argmax}_{\eta} \frac{1}{m} \sum_{i=1}^m \log p(\mathbf{x}^{(i)}; \eta)$$

To help us establish the connection between the two, we will define the data $\bar{\mathbf{x}}$ driven empirical distribution as

$$\varphi_{\bar{\mathbf{x}}}(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m \delta(\mathbf{x} = \mathbf{x}^{(i)})$$

This definition yields an alternative expression for the MLE as the minimizer of the KL divergence between the empirical distribution and the model distribution.

$$\hat{p}_{MLE} = \operatorname{argmax}_{\eta} \left(\frac{1}{m} \sum_{i=1}^m \wp_{\bar{\mathbf{x}}} \log p(\mathbf{x}^{(i)}; \eta) \right) \quad (7.61)$$

$$= \operatorname{argmax}_{\eta} \left(\frac{1}{m} \sum_{i=1}^m \wp_{\bar{\mathbf{x}}} \log p(\mathbf{x}^{(i)}; \eta) \right) + \left(\frac{1}{m} \sum_{i=1}^m \wp_{\bar{\mathbf{x}}} \log \wp_{\bar{\mathbf{x}}} \right) \quad (7.62)$$

$$= \operatorname{argmax}_{\eta} - D(\wp_{\bar{\mathbf{x}}} \| p(\bar{\mathbf{x}}; \eta))$$

$$= \operatorname{argmin}_{\eta} D(\wp_{\bar{\mathbf{x}}} \| p(\bar{\mathbf{x}}; \eta))$$

$$= \operatorname{argmin}_{p \in \mathbf{E}(\mathbf{f})} D(\wp_{\bar{\mathbf{x}}} \| p(\cdot)) \quad (7.63)$$

where, $\mathbf{E}(\mathbf{f})$ is the exponential family of distributions having \mathbf{f} as the vector of feature functions that also figure in the specification of constraints in $\mathcal{P}(\bar{\mu})$:

$$\mathbf{E}(\mathbf{f}) = \{p(\cdot) \mid p(\mathbf{x}; \eta) \propto \exp \{ \eta^T \mathbf{f}(\mathbf{x}) \} \}$$

We will now show the equivalence of specifications in (7.60) and (7.63). The discussion so far will be relevant in our proof. On the one hand, we have seen in (7.49) in theorem 99 that the constraint in $\mathbf{E}(\mathbf{f})$ is satisfied by any solution to (7.60). While on the other hand, we know that the moment matching conditions for (7.63) in (7.57) are precisely the constraints in $\mathcal{P}(\bar{\mu})$.

Theorem 101 $\hat{p}_{MLE} = \hat{p}_{ML}$ for exponential families, where:

$$\hat{p}_{MLE} = \operatorname{argmin}_{p \in \mathbf{E}(\mathbf{f})} D(\wp_{\bar{\mathbf{x}}} \| p(\cdot))$$

and

$$\hat{p}_{ME} = \operatorname{argmin}_{p \in \mathcal{P}(\bar{\mu})} D(p \| u)$$

Two differences between these formulations are:

1. While \hat{p}_{MLE} involves optimization over the second argument in the KL divergence, \hat{p}_{ME} involves optimization over the first argument.
2. The entry point of the data is also toggled in the two; while \hat{p}_{ME} has data entering through constraint set \mathcal{P} , \hat{p}_{MLE} has data entering through the cost function.

This is characteristic of dual problems.

Proof: This can be proved by first showing that the problem in (7.60) is the dual of (7.63). Next we will need to apply duality theory in Theorem 82, which states that for convex cost function and convex inequality constraint set, the KKT conditions are necessary and sufficient conditions for zero duality gap. It can be proved that (7.60) is convex and that the parameters for \hat{p}_{MLE} are solutions to the KKT conditions.

The theorem can also be proved by invoking the so called *pythagorean theorem*⁴⁶ for general class of distance that includes distributions. In this particular case, it can be shown that for all $\hat{p} \in \mathcal{P}(\bar{\mu}) \cap \mathbf{E}(\mathbf{f})$ and for all $p_{\mathcal{P}} \in \mathcal{P}(\bar{\mu})$ and $p_{\mathcal{E}} \in \mathcal{E}(\mathbf{f})$,

$$D(p_{\mathcal{P}}||p_{\mathbf{E}}) = D(p_{\mathcal{P}}||\hat{p}) + D(\hat{p}||p_{\mathbf{E}})$$

If $\hat{q}, \hat{p} \in \mathcal{P}(\bar{\mu}) \cap \mathbf{E}(\mathbf{f})$, then it will turn out by simple application of the theorem that $D(\hat{q}||\hat{p}) + D(\hat{p}||\hat{q}) = 0$, which implies that $\hat{p} = \hat{q}$. That is, $\mathcal{P}(\bar{\mu}) \cap \mathbf{E}(\mathbf{f})$ is a singleton and \hat{p} should correspond to both \hat{p}_{MLE} and \hat{p}_{ML} . \square

7.8 Learning with Incomplete Observations

Thus far, we have focussed on learning parameters for graphical models using complete observations; the underlying model was $p(\mathbf{x}; \theta)$ and the observations (data) were $\bar{\mathbf{x}} = (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)})$. An example of such a learning task was presented in the case of Markov chains on page 453. Consider the hidden markov model from page 430 with $\mathcal{V} = \{X_1, X_2, X_3, X_4, X_5, Y_1, Y_2, Y_3, Y_4, Y_5\}$ and $\mathcal{E} = \{(X_1, X_2), (X_1, Y_1), (X_2, X_3), (X_2, Y_2), (X_3, X_4), (X_3, Y_3), (X_4, X_5), (X_4, Y_4), (X_5, Y_5)\}$. where nodes $\{X_1, X_2, X_3, X_4, X_5\}$ are hidden variables and nodes $\{Y_1, Y_2, Y_3, Y_4, Y_5\}$ are the observed variables.

Mixture models are another set of popular example, which feature hidden variables. Many real world problems are characterized by distinct classes/subpopulations that the data falls into and can be modeled using mixture models.

Definition 63 Let $Z \in \{z_1, z_2, \dots, z_k\}$ be a multinomial variable indicating mixture component. Let \mathbf{X} be a random variable (vector), whose distribution is specified, conditioned on different values z_i of Z as

$$p(\mathbf{x}|z_i; \theta_i) \sim f_i(\mathbf{x}; \theta_i)$$

Then the finite mixture model is defined as

$$p(\mathbf{x}) \sum_{i=1}^k p(z_i) f_i(\mathbf{x}; \theta_i)$$

with k being the number of mixture components, Z called the mixture indicator component, $f_i(\mathbf{x}; \theta_i)$ termed as the density of the i^{th} mixture component with parameters θ_i . The quantities $p(z_i) = \pi_i$ are also called mixing weights, representing the proportion of the population in subpopulation i . Thus, $\pi = [\pi_1, \pi_2, \dots, \pi_k]$ and $\theta = [\theta_1, \theta_2, \dots, \theta_k]$ are the parameters of a finite mixture model, with a fixed value of k . As a graphical model, the mixture model can be represented as a two node graph: $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$ with $\mathcal{V} = \{\mathbf{X}, Z\}$ and $\mathcal{E} = \{(\mathbf{X}, Z)\}$.

⁴⁶The right angle here is not the conventional one, but a notional one.

As an example, the density of each mixture component could be Gaussian with $\theta_i = (\mu_i, \Sigma_i)$.

$$f_i(\mathbf{x}; \theta_i) \sim \mathcal{N}(\mu_i, \Sigma_i)$$

The distribution $p(\mathbf{x})$ is then called a mixture of Gaussians. In general, it not Gaussian itself.

How can the parameters be learnt in the presence of incomplete data? In the case of the HMM example, we might be provided only with observations $\bar{\mathbf{y}}$ for \mathbf{Y} and expected to learn the parameters. Or in the case of mixture models, we might be presented only with instances of \mathbf{X} in the form of $\bar{\mathbf{x}}$ and required to learn parameters π and θ . We will illustrate parameter learning for the mixture model problem.

7.8.1 Parameter Estimation for Mixture Models

It can be shown that learning for mixture models is an easy problem if the data is fully observed in the form $(\bar{\mathbf{x}}, \bar{z}) = [(\mathbf{x}^{(1)}, z^{(1)}), (\mathbf{x}^{(2)}, z^{(2)}), \dots, (\mathbf{x}^{(m)}, z^{(m)})]$. The joint distribution can be decomposed as

$$p(\mathbf{x}, z; \theta) = p(z)p(\mathbf{x} | z, \theta)$$

If $p(\mathbf{x} | z, \theta)$ is Gaussian and since $p(z)$ is multinomial, the joint will be in exponential form with Gaussian and multinomial sufficient statistics. The maximum likelihood estimation will boil down to moment matching with respect to these sufficient statistics, leading to an easy estimation problem.

In the incomplete data setting, we are given only $\bar{\mathbf{x}}$ while observations \bar{z} on the mixture components are hidden. The likelihood can still be expressed and maximized:

$$LL(\pi, \theta; \bar{\mathbf{x}}) = \frac{1}{m} \sum_{i=1}^m \log p(\mathbf{x}^{(i)}; \theta) = \frac{1}{m} \sum_{i=1}^m \log \left[\sum_{j=1}^k \pi_j f_j(\mathbf{x}^{(i)}; \theta_j) \right]$$

subject to the constraints that $\pi_j \geq 0$ and $\sum_{j=1}^k \pi_j = 1$.

Unfortunately, log cannot be distributed over a summation and that creates the main bottleneck. In case the densities are Gaussians, the objective to be maximized will be

$$LL(\pi, \mu, \Sigma; \bar{\mathbf{x}}) = \frac{1}{m} \sum_{i=1}^m \log \left[\sum_{j=1}^k \pi_j \mathcal{N}(\mathbf{x}^{(i)}; \mu_j, \Sigma_j) \right]$$

subject to the constraints that $\pi_j \geq 0$ and $\sum_{j=1}^k \pi_j = 1$.

1. **M-step:** Writing down the KKT necessary and sufficient optimality conditions (see (4.88) on page 296) for this maximization problem, subject to its associated inequality and linear equality constraints yields:

(a) For μ_j

$$\mu_j = \frac{\sum_{i=1}^m p(z_j | \mathbf{x}^{(i)}, \mu, \Sigma) \mathbf{x}^{(i)}}{\sum_{i=1}^m \sum_{j=1}^k p(z_j | \mathbf{x}^{(i)}, \mu, \Sigma)} \quad (7.64)$$

(b) And for Σ_j

$$\Sigma_j' = \frac{\sum_{i=1}^m p(z_j | \mathbf{x}^{(i)}, \mu, \Sigma) (\mathbf{x}^{(i)} - \mu_j) (\mathbf{x}^{(i)} - \mu_j)^T}{\sum_{i=1}^m \sum_{j=1}^k p(z_j | \mathbf{x}^{(i)}, \mu, \Sigma)} \quad (7.65)$$

These steps are called the M – *steps* or the maximization steps, since they are obtained as necessary and sufficient conditions of optimality for a maximization problem.

2. **E-step:** The posterior $p(z_j | \mathbf{x}^{(i)}, \mu, \Sigma)$ and the prior π_j in (7.65) and (7.64) can be determined using Bayes rule as

(a)

$$p(z_j | \mathbf{x}^{(i)}, \mu, \Sigma) = \frac{\pi_j f(\mathbf{x}; \mu_j, \Sigma_j)}{\sum_{a=1}^k \pi_a f(\mathbf{x}; \mu_a, \Sigma_a)}$$

(b)

$$\pi_j = \frac{1}{m} \sum_{i=1}^m p(z_j | \mathbf{x}^{(i)}, \mu, \Sigma)$$

The problem is that we do not get a closed form solution here; what we obtain are a set of coupled, non-linear equations and need to iterate between these steps to arrive at the fix point. This is where the expectation maximization (EM) algorithm comes in. We now will specify the EM algorithm in a more general setting.

7.8.2 Expectation Maximization

Let \mathbf{X} be a set of observed variables and \mathbf{Z} be a set of hidden variables for some statistical model. Let $\bar{\mathbf{x}}$ be m observations on \mathbf{X} . In this general setting, we really need not assume that the samples in $\bar{\mathbf{x}}$ are iid (though you could). We will assume that the MLE problem would have been easy⁴⁷ if $\bar{\mathbf{z}}$ was observed data for the hidden variables \mathbf{Z} (such as in the case of the mixture model). The complete data log-likelihood would have been:

$$LL(\theta; \bar{\mathbf{x}}, \bar{\mathbf{z}}) = \frac{1}{m} \log p(\bar{\mathbf{x}}, \bar{\mathbf{z}}; \theta)$$

Given a predictive distribution $q(\mathbf{z}|\bar{\mathbf{x}})$, the expected complete data log-likelihood is a function of the observed $\bar{\mathbf{x}}$ and θ and is defined as

$$LL_E(\theta; \bar{\mathbf{x}}) = \sum_{\mathbf{z}} q(\mathbf{z}|\bar{\mathbf{x}}) \log p(\bar{\mathbf{x}}, \mathbf{z}; \theta) \quad (7.66)$$

The expected complete data log-likelihood is an auxiliary function that gives a lower bound on the actual log-likelihood we want to optimize for. The actual log-likelihood in this general setting will be:

$$LL(\theta; \bar{\mathbf{x}}) = \frac{1}{m} \log \left\{ \sum_{\mathbf{z}} p(\bar{\mathbf{x}}, \mathbf{z}; \theta) \right\}$$

For example, the actual log-likelihood with iid assumption will be:

$$LL(\theta; \bar{\mathbf{x}}) = \frac{1}{m} \sum_{i=1}^m \log \left\{ \sum_{\mathbf{z}} p(\mathbf{x}^{(i)}, \mathbf{z}; \theta) \right\}$$

Theorem 102 For all θ and every possible distribution $q(\mathbf{z}|\bar{\mathbf{x}})$, following holds:

$$LL(\theta; \bar{\mathbf{x}}) \geq LL_E(\theta; \bar{\mathbf{x}}) + \frac{1}{m} H(q) \quad (7.67)$$

Equality holds if and only if

$$q(\mathbf{z}|\bar{\mathbf{x}}) = p(\mathbf{z}|\bar{\mathbf{x}}; \theta)$$

Proof: First of all

$$LL(\theta; \bar{\mathbf{x}}) = \frac{1}{m} \log \left\{ \sum_{\mathbf{z}} q(\mathbf{z}|\bar{\mathbf{x}}) \frac{p(\bar{\mathbf{x}}, \mathbf{z}; \theta)}{q(\mathbf{z}|\bar{\mathbf{x}})} \right\}$$

⁴⁷The trick in such a setting is to identify the model, \mathbf{X} and \mathbf{Z} so that you make the MLE problem easy in the presence of complete data.

Using the Jensen's inequality (since \log is a strictly convex function),

$$LL(\theta; \bar{\mathbf{x}}) \geq \underbrace{\frac{1}{m} \sum_{\mathbf{z}} q(\mathbf{z}|\bar{\mathbf{x}}) \log p(\bar{\mathbf{x}}, \mathbf{z}; \theta)}_{LL_E(\theta; \bar{\mathbf{x}})} - \underbrace{\frac{1}{m} \sum_{\mathbf{z}} q(\mathbf{z}|\bar{\mathbf{x}}) \log q(\mathbf{z}|\bar{\mathbf{x}})}_{H(q)}$$

Equality holds if and only if $\frac{p(\bar{\mathbf{x}}, \mathbf{z}; \theta)}{q(\mathbf{z}|\bar{\mathbf{x}})}$ is a constant, that is,

$$q(\mathbf{z}|\bar{\mathbf{x}}) \propto p(\bar{\mathbf{x}}, \mathbf{z}; \theta) = p(\mathbf{z}|\bar{\mathbf{x}}; \theta)p(\bar{\mathbf{x}}; \theta) \propto p(\mathbf{z}|\bar{\mathbf{x}}; \theta)$$

This can happen if and only if $q(\mathbf{z}|\bar{\mathbf{x}}) = p(\mathbf{z}|\bar{\mathbf{x}}; \theta)$. \square

A consequence of theorem 102 is that

$$\max_{\theta} LL(\theta; \bar{\mathbf{x}}) = \max_{\theta} \max_q LL_E(\theta; \bar{\mathbf{x}}) + \frac{1}{m} H(q) \quad (7.68)$$

The EM algorithm is simply coordinate ascent on the auxiliary function $LL_E(\theta; \bar{\mathbf{x}}) + \frac{1}{m} H(q)$. The expectation and maximization steps at time instance t can be easily identified for the formulation in (7.68) as

1. Expectation Step:

$$q^{(t+1)} = \operatorname{argmax}_q LL_E(\theta^{(t)}; \bar{\mathbf{x}}) + \frac{1}{m} H(q) = \operatorname{argmax}_q -D\left(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}; \theta^{(t)})\right) + \log \left\{ \mathbf{x}; \theta^{(t)} \right\} \quad (7.69)$$

Since, $LL_E(\theta^{(t)}; \bar{\mathbf{x}}) + \frac{1}{m} H(q) \leq \log \left\{ \mathbf{x}; \theta^{(t)} \right\}$ by theorem 102, the maximum value is attained in (7.69) for $q(\mathbf{z}|\mathbf{x}) = p(\mathbf{z}|\mathbf{x}; \theta^{(t)})$. Thus, the E-step can be summarized by

$$q^{(t+1)}(\mathbf{x}|\mathbf{x}) = p(\mathbf{z}|\mathbf{x}; \theta^{(t)}) \quad (7.70)$$

The E-step can involve procedures such as sum-product for obtaining marginals and/or conditions, if the distribution is defined on a graphical model, to obtain $p(\mathbf{z}|\mathbf{x}; \theta^{(t)})$.

2. Maximization Step:

$$\theta^{(t+1)} = \operatorname{argmax}_{\theta} LL_E(\theta; \bar{\mathbf{x}}) + \frac{1}{m} H(q^{(t+1)})$$

Since the maximization is over θ and since $H(q)$ is independent of θ , we can rewrite the M-step to be

$$\theta^{(t+1)} = \operatorname{argmax}_{\theta} LL_E(\theta; \bar{\mathbf{x}}) = \operatorname{argmax}_{\theta} \sum_{\mathbf{z}} q(\mathbf{z}|\bar{\mathbf{x}}) \log p(\bar{\mathbf{x}}, \bar{\mathbf{z}}; \theta) \quad (7.71)$$

In essence, the M-step looks very much like an ordinary maximum likelihood estimation problem, but using predicted values of \mathbf{z} . The M-step may not have a closed form solution, in which case, it may be required to resort to iterative techniques such as IPF (7.7.6).

Let us take some examples to illustrate the generic EM procedure outlined here. It is particularly useful if the term $\log p(\bar{\mathbf{x}}, \bar{\mathbf{z}}; \theta)$ were to split up into smaller terms (such as sum of sufficient statistics in the case of exponential models). Consider the Gauss Markov process specified by $Z_{t+1} = \theta Z_t + W_t$, where $Z_0 \sim \mathcal{N}(0, 1)$ and $W_t \sim \mathcal{N}(0, 1)$. Let $\theta \in \Re$ be the parameter to be estimated. The graphical model representation is $\mathcal{V} = \{Z_1, Z_2, Z_2, \dots, Z_n, X_1, X_2, \dots, X_n\}$ and $\mathcal{E} = \{(Z_1, Z_2), (Z_1, X_1), (Z_2, Z_3), (Z_2, Z_2), \dots, (Z_{n-1}, Z_n), (Z_{n-1}, X_{n-1}), (Z_n, X_n)\}$.

Say what we observe are noisy, indirect observations $X_t = Z_t + V_t$, $V_t \sim \mathcal{N}(0, \sigma^2)$ being the observation noise. Let $\bar{\mathbf{x}}$ be a set of m iid observations for \mathbf{X} while \mathbf{Z} remains hidden. Both \mathbf{X} and \mathbf{Z} are vectors of random variables of size n each. Then,

$$\begin{aligned} LL(\theta; \bar{\mathbf{x}}) &= \frac{1}{m} \sum_{i=1}^m p(x^{(i)}; \theta) \\ &= \frac{1}{m} \sum_{i=1}^m \log \left\{ \int_{\mathbf{z}} \prod_{t=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{1}{2} \frac{(x_t^{(i)} - z_t)^2}{\sigma^2} p(\mathbf{z}; \theta) d\mathbf{z} \right\} \right\} \end{aligned} \quad (7.72)$$

which is a mess! In contrast, the lower-bound component LL_E allows us to move the integration outside the logarithm, enabling more simplification:

$$\begin{aligned} LL_E(\theta; \bar{\mathbf{x}}) &= \frac{1}{m} \int_{\mathbf{z}} q(\mathbf{z}|\bar{\mathbf{x}}) \log p(\bar{\mathbf{x}}\mathbf{z}; \theta) d\mathbf{z} \\ &= \frac{1}{m} \int_{\mathbf{z}} q(\mathbf{z}|\bar{\mathbf{x}}) \sum_{i=1}^m \log p(\mathbf{x}^{(i)}|\mathbf{z}; \theta) d\mathbf{z} \\ &= \frac{1}{m} \int_{\mathbf{z}} q(\mathbf{z}|\bar{\mathbf{x}}) \sum_{i=1}^m \left[\log p(\mathbf{z}, \theta) + \sum_{t=1}^n \log \left(\frac{1}{\sqrt{2\pi}\sigma^2} \right) - \frac{1}{2\sigma^2} (x_t^{(i)} - z_t)^2 \right] d\mathbf{z} \end{aligned}$$

As can be seen above, $p(\mathbf{x}^{(i)}|\mathbf{z}; \theta)$ is independent of θ and therefore, the term $q(\mathbf{z}|\bar{\mathbf{x}}) \log p(\mathbf{x}^{(i)}|\mathbf{z}; \theta)$ can be brought out as a separate constant (since that part of the integral will not change with θ). This leads to the following simplified expression for LL_E

$$\begin{aligned}
LL_E(\theta; \bar{\mathbf{x}}) &= \frac{1}{m} \int_{\mathbf{z}} q(\mathbf{z}|\bar{\mathbf{x}}) \log p(\bar{\mathbf{x}}\mathbf{z}; \theta) d\mathbf{z} \\
&= C + \frac{1}{m} \int_{\mathbf{z}} q(\mathbf{z}|\bar{\mathbf{x}}) [\log p(z_1) + \log p(z_2|z_1; \theta) + \dots + \log p(z_n|z_{n-1}; \theta)] d\mathbf{z} \\
&= C' + \frac{1}{m} \int_{\mathbf{z}} q(\mathbf{z}|\bar{\mathbf{x}}) [\log p(z_2|z_1; \theta) + \dots + \log p(z_n|z_{n-1}; \theta)] d\mathbf{z}
\end{aligned}$$

In the E-step, the Kalman filter can be used to compute $p(\mathbf{z}|\mathbf{x}; \theta^{(t+1)})$ in terms of $\theta^{(t)}$. In the M-step, first order necessary optimality conditions on LL_E will yield $\theta^{(t+1)}$.

Recall from Section 7.7.7, equation (7.63) that the likelihood maximization problem can be viewed as a problem of minimizing the KL divergence between the empirical distribution and the model distribution.

$$\hat{p}_{MLE} = \operatorname{argmin}_{p \in \mathbf{E}(\mathbf{f})} D(\varphi_{\bar{\mathbf{x}}} \| p(\mathbf{x}; \theta))$$

While the likelihood maximization perspective lead to a lower-bounding strategy in the form of EM, an alternative upper-bounding strategy can also be adopted to view EM, though it is only the older bound in disguise. Making use of theorem 102, we can prove that for all distributions $q(\mathbf{z}|\mathbf{x})$ and any parameter θ , the following always holds:

$$D(\varphi(\mathbf{x}) \| p(\mathbf{x}; \theta)) \leq D(\varphi(\mathbf{x})q(\mathbf{z}|\mathbf{x}) \| p(\mathbf{x}, \mathbf{z}; \theta)) \quad (7.73)$$

This statement says that the KL divergence between the empirical and model distributions that maximum likelihood tries to minimize is upperbounded by the KL divergence between the ‘completed’ empirical and model distributions. As before, it can also be shown that the bound is tight if and only if $q(\mathbf{z}|\mathbf{x}) = p(\mathbf{z}|\mathbf{x}; \theta)$. The E-step will remain the same as before. Only, the M-step will slightly change:

1. KL divergence perspective of E-Step:

$$q^{(t+1)}(\mathbf{z}|\mathbf{x}) = \operatorname{argmin}_q D(\varphi(\mathbf{x})q(\mathbf{z}|\mathbf{x}) \| p(\mathbf{x}, \mathbf{z}; \theta^{(t)}))$$

2. KL divergence perspective of M-Step:

$$\theta^{(t+1)} = \operatorname{argmin}_{\theta} D(\varphi(\mathbf{x})q^{(t+1)}(\mathbf{z}|\mathbf{x}) \| p(\mathbf{x}, \mathbf{z}; \theta))$$

These modified E and M steps correspond to coordinate descent in contrast to the earlier perspective of coordinate ascent.

7.9 Variational Methods for Inference

In contrast to the sampling methods, variational methods are deterministic and fast algorithms that generate good approximations to the problems of computing marginals and MAP configurations. They involve the reformulation of the quantity of interest (such as log-likelihood, first order moments, marginal distributions, *etc.*) as the solution of some optimization problem. They are useful in several ways:

- The variational formulation often leads to efficient algorithms for determining exact solutions. Many algorithms discussed thus far, could be discovered as efficient techniques for solving the variational optimization problem.
- For many quantities that are hard to compute, the variational perspective leads to approximate solutions. Mean field and loopy sum product algorithms can also be viewed as special cases of approximation through variational inference.
- In contrast to approximate sampling methods, these are faster, deterministic and cheaper (in terms of memory).

We will motivate variational methods using two examples.

1. The first is the mean field algorithm for the Ising model defined on a graph $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$ of binary (0/1) variables, with pairwise interactions between adjacent variables captured through their product $(x_s x_t)$

$$p(\mathbf{x}, \eta) \propto \exp \left\{ \sum_{s \in \mathcal{V}} \eta_s x_s + \sum_{(s,t) \in \mathcal{E}} \eta_{st} x_s x_t \right\}$$

The Gibbs sampling for the Ising model derives updates of the form

$$x_i^{(t+1)} = \begin{cases} 1 & \text{if } u \sim \text{uniform}[0, 1] \leq 1 + \exp \left\{ -\eta_i + \sum_{j \in \mathcal{N}(s)} \eta_{ij} x_j^{(t)} \right\} \\ 0 & \text{otherwise} \end{cases}$$

which correspond to a non-deterministic version of the message passing algorithm (owing to $u \sim \text{uniform}[0, 1]$). The updates are very simple and local, making this a good choice.

The mean field method has its roots in physics. It makes a deterministic to the Gibbs update by replacing each random variable X_i by a deterministic mean parameter $\mu_i \in [0, 1]$ (which can be thought of as the probability of $X_i = 1$) and updates μ_i using

$$\mu_i^{(t+1)} = \left(1 + \exp \left\{ -\eta_i + \sum_{j \in \mathcal{N}(s)} \eta_{ij} \mu_j^{(t)} \right\} \right)^{-1}$$

Thus, the mean field algorithm is exactly a message passing algorithm, but has some semantics related to sampling techniques. We will see that the mean field algorithm is a specific instance of variational methods and it can be formulated as coordinate descent on a certain optimization problem and subsequently be analysed for convergence, *etc.*

2. The *loopy sum-product* (also called the loopy belief propagation) method is another instance of variational methods. ‘Loopy’ here means on graphs with cycles. If the tree width were low, you could create a junction tree and perform message passing, but what if the tree width were large, such as with a grid. This algorithm is the most naive application of the sum-product updates (originally developed for trees in Section 7.2.2) and apply it to graphs with cycles. This naive procedure has had extraordinary success in the fields of signal processing, compute vision, bioinformatics and most importantly, in communication networks, *etc.*

The message passing algorithm, when applied on a tree, breaks it into subtrees and passes messages between subtrees, which are independent (share no variable). But the moment you add an edge connecting any two subtrees, they are not independent any more. Passing messages around in cycles can lead to over-counting (analogous to gossip spreading in a social network). Thus, the message passing algorithm ceases to remain an exact algorithm and does not even guarantee convergence.

What turns out to be actually very important is how long are the cycles on which messages are being propagated; for *long cycles, the effects of over-counting can be weakened*. More technically speaking, the behaviour will depend on

- (a) The girth of the graph (length of cycles): For larger girth, you could run the message passing for many iterations before you land up with overcounting.
 - (b) The strength of the potentials are, or in other words, how close to independence is the model. For example, in the Ising model itself, based on the coupling induced through terms of the form $x_s x_t$, if the coupling is weak, almost close to independence, the algorithm will be perfect, giving almost exact answers. There is a region of transition, based on strengthening of the coupling terms, beyond which the algorithm breaks down.
3. The idea behind variational methods is to represent the quantity of interest (such as the marginal or mode over a graphical model) as the solution to an optimization problem. For example, the solution to $A\mathbf{x} = \mathbf{b}$ (as in the case of inferencing for Kalman filters) with $A \succ 0$ and $\mathbf{b} \in \mathbb{R}^n$ can be represented as the solution to

$$\tilde{\mathbf{x}} = \operatorname{argmin}_{\mathbf{x}} \frac{1}{2} \mathbf{x}^T A \mathbf{x} - \mathbf{b}^T \mathbf{x}$$

This is precisely a *variational formulation* for the linear system $A\mathbf{x} = \mathbf{b}$. If the system of equations $A\mathbf{x} = \mathbf{b}$ is large⁴⁸ the solution $\tilde{\mathbf{x}} = A^{-1}\mathbf{b}$ may not be easy to compute, in the event of which, iterative (and sometimes approximate) solutions to the optimization problem can be helpful. One of the most successful techniques for solving such systems (without inverting matrices) is the conjugate gradient method, discussed in Section 4.5.8, that solves the system in exactly $O(n)$ steps.

4. As will be seen on page 432 in Section 7.4, the bernoulli distribution can be expressed as

$$p(\mathbf{x}, \eta) = \exp\{\eta x - A(\eta)\}$$

for $X \in \{0, 1\}$. $A(\eta) = \log(1 + e^\eta)$. We saw in Section 7.4 that the mean is given by

$$\bar{\mu} = E_\eta = \nabla A(\eta) = \frac{e^\eta}{1 + e^\eta} = (1 + e^{-\eta})^{-1}$$

The key is in noting that $\bar{\mu}$ corresponds to the ‘slope of’ a supporting hyperplane (see Figure 4.38 on page 271) for $\text{epi}(A(\eta))$ in a $(\eta, A(\eta))$ space. Thus, we are interested in all the hyperplanes that lie below $\text{epi}(A(\eta))$, with intercept C along the axis for $A(\eta)$:

$$\mu^T \eta - C \leq A(\eta)$$

and want to get as close to a supporting hyperplane as possible

$$C^* = \sup_{\eta} \{\mu^T \eta - A(\eta)\} = A^*(\mu)$$

Borrowing ideas from duality theory (*c.f.* Section 4.4), we call the function $\sup_{\eta} \{\mu^T \eta - A(\eta)\}$ as the dual function $A^*(\mu)$. C^* is the intercept for the supporting hyperplane. In the case of the bernoulli example, we have

$$A^*(\mu) = \sup_{\eta} \{\mu^T \eta - \log(1 + e^\eta)\} = \begin{cases} (1 - \mu) \log(1 - \mu) + \mu \log \mu & \text{if } \mu \in (0, 1) \\ \infty & \text{otherwise} \end{cases}$$

Under nice conditions (such as under *Slater's constraint qualification* discussed in definition 42 on page 290), the operation of taking duals is symmetric ($(A^*)^* = A$), that is,

$$A(\theta) = \sup_{\mu} \{\mu^T \eta - A^*(\mu)\}$$

Under the conditions of zero duality gap, it should happen that if

$$\hat{\mu}(\eta) = \operatorname{argmax}_{\mu} \{\mu^T \eta - A^*(\mu)\}$$

⁴⁸Of course, as we saw in Section 7.5, such a system comes up in the case of solution to Kalman filters, but can be computed efficiently by exploiting the tree structure of the graph in inverting the matrix. The discussion here is for general graphs.

is the primal optimum, then $\widehat{\mu}^T \eta - A^*(\widehat{\mu})$ is the supporting hyperplane to $A(\eta)$, meaning that $\widehat{\mu}$ is the mean we were seeking. This yields a variational representation for the original problem of finding the mean. The dual itself is also useful in determining the log-normalization constant for problems such as parameter estimation. We can confirm in the simple bernoulli case that indeed

$$\widehat{\mu}(\eta) = \bar{\mu}(\eta) = \frac{e^\eta}{1 + e^\eta}$$

We will now generalize the variational formulation of the bernoulli case to the exponential family.

$$p(\mathbf{x}; \theta) = \exp \{ \theta^T \mathbf{f}(\mathbf{x}) - A(\theta) \}$$

where, $\mathbf{x} \in \mathfrak{R}^n$ and $\mathbf{f} : \mathfrak{R}^n \rightarrow \mathfrak{R}^d$. Let us say we are interested in computing the first order moments

$$\mu = E[\mathbf{f}(\mathbf{x})] \quad (7.74)$$

Following a similar line of argument as for the case of the bernoulli distribution, we define the dual as

$$A^*(\mu) = \sup_{\theta} \{ \mu^T \theta - A(\theta) \}$$

The key ingredients in this calculation are to set the gradient with respect to θ to $\mathbf{0}$, as a necessary first order condition.

$$\mu - \nabla A(\theta) = \mathbf{0} \quad (7.75)$$

This looks like the moment matching problem that results from maximum likelihood estimation. The only difference is that μ need not come from data, it is the argument to $A^*(\eta)$. When will (7.75) have a solution? In the simple bernoulli case, we already saw that we will have a solution $\eta = -\log \mu - \log(1 - \mu)$ if $\mu \in (0, 1)$. As another example, for the univariate Gaussian, $\mathbf{f}(\mathbf{x}) = [x, x^2]$ and $A(\eta) = \frac{1}{2\sigma^2} \mu^2 + \log \sigma \equiv -\frac{\eta_1^2}{4\eta_2} + \frac{1}{2} \log(-2\eta_2)$ (as seen on page 432). The system (7.75) can be shown to have a solution only if $\mu_2 - \mu_1^2 > 0$, that is if the variance is positive. For two examples, the conditions under which solutions exist to (7.75) are extremely simple - it should be possible to generate data using the distribution.

Assuming that a solution $\theta(\mu)$ exists to (7.75), and exploiting the fact that $\theta(\mu)$ satisfies moment matching conditions (7.75)

$$\sum_{\mathbf{x}} p(\mathbf{x}; \theta(\mu)) \mathbf{f}(\mathbf{x}) = \mu \quad (7.76)$$

and using the property that

$$A(\mu) = \sum_{\mathbf{x}} A(\mu) p(\mathbf{x}; \theta(\mu))$$

the dual will have the form

$$\begin{aligned} A^*(\mu) &= \theta^T(\mu)\mu - A(\mu) \\ &= \theta^T(\mu) \left(\sum_{\mathbf{x}} p(\mathbf{x}; \theta(\mu)) \mathbf{f}(\mathbf{x}) \right) - A(\mu) \\ &= \sum_{\mathbf{x}} p(\mathbf{x}; \theta(\mu)) (\theta^T(\mu) \mathbf{f}(\mathbf{x}) - A(\mu)) \\ &= \sum_{\mathbf{x}} p(\mathbf{x}; \theta(\mu)) \log p(\mathbf{x}; \theta(\mu)) \\ &= -H(p(\mathbf{x}; \theta(\mu))) \end{aligned} \tag{7.77}$$

That is, the dual is precisely the negative entropy $-H(p(\mathbf{x}; \theta(\mu)))$ of the distribution whose parameters $\theta(\mu)$ are obtained by moment matching. The dual for the bernoulli case which resembled an entropy was by no means a coincidence! If \mathcal{M} is the set of all possible moments that make a solution to the system (7.75 or equivalently 7.76) feasible, that is

$$\mathcal{M} = \left\{ \mu \in \mathbb{R}^d \mid \sum_{\mathbf{x}} p(\mathbf{x}; \theta) \mathbf{f}(\mathbf{x}) = \mu \text{ for some } p(\cdot) \right\}$$

then the dual could also be expressed as

$$A^*(\mu) = \begin{cases} -\max H(p(\mathbf{x}; \theta(\mu))) & \text{such that } E[\mathbf{f}(\mathbf{x})] = \mu \text{ for } \mu \in \mathcal{M} \\ \infty & \text{otherwise} \end{cases}$$

Another way of characterizing \mathcal{M} is as the set of all first order moments that could be generated by $p(\mathbf{x}; \theta)$. In any case, \mathcal{M} is often very hard to characterize and loop belief propagation *etc.* are often used to approximate it.

Finally, we will write down the variational problem as a reformulation of (7.74), of which mean field, (loopy) sum-product, Gauss Seidel, Jacobi, etc can be found to be special cases. Writing down the dual of the dual of (7.77), and assuming zero duality gap, we get a reformulation of (7.74):

$$A(\theta) = \sup_{\mu \in \mathcal{M}} \{ \mu^T \theta - A^*(\mu) \} \tag{7.78}$$

Again, $A(\theta)$ is a very hard function to compute, mainly because \mathcal{M} is simple to characterize. This maximization problem is concave, since $A^*(\eta)$ is concave

and the constraint set \mathcal{M} is convex. Under zero duality gap conditions (which holds in this case), the optimal point will be achieved at $\hat{\mu}(\theta) = E[\mathbf{f}(\mathbf{x})]$.

We will use some examples to illustrate the use of variational techniques. For problems of estimating moments, \mathbf{f} could be the feature functions. For problems of estimating marginals, \mathbf{f} can be chosen to be the indicator function.

The simplest example is for a two node chain: $\mathcal{V} = \{X_1, X_2\}$, $\mathcal{E} = \{(X_1, X_2)\}$, $X_1, X_2 \in \{0, 1\}$ and

$$p(\mathbf{x}; \theta) \propto \exp \{ \theta_1 x_1 + \theta_2 x_2 + \theta_{12} x_1 x_2 \}$$

The moments are: $\mu_i = E[X_i] = p(X_i = 1)$ and $\mu_{12} = E[X_1 X_2] = p(X_1 = 1, X_2 = 1)$. The set \mathcal{M} is

$$\mathcal{M} = \left\{ \mu \in \mathbb{R}^3 \left| \sum_{\mathbf{x}} p(\mathbf{x}; \theta) \mathbf{f}(\mathbf{x}) = \mu \text{ for some } p(\cdot) \right. \right\} = \{ \mu_i \in [0, 1], 0 \leq \mu_{12} \leq \min(\mu_1, \mu_2), 1 + \mu_{12} - \mu_1 - \mu_2 \geq 0 \}$$

Let us next write down the dual in terms of the entropy of the distribution

$$\begin{aligned} A^*(\mu) = -H(p(\mathbf{x}; \mu)) &= \sum_{x_1, x_2} p(x_1, x_2) \log p(x_1, x_2) \\ &= \mu_{12} \log \mu_{12} + (\mu_1 - \mu_{12}) \log (\mu_1 - \mu_{12}) + (\mu_2 - \mu_{12}) \log (\mu_2 - \mu_{12}) \\ &\quad + (1 + \mu_{12} - \mu_1 - \mu_2) \log (1 + \mu_{12} - \mu_1 - \mu_2) \end{aligned} \quad (7.79)$$

The corresponding variational problem will be

$$A(\theta) = \max_{\mu \in \mathcal{M}} \{ \theta_1 \mu_1 + \theta_2 \mu_2 + \theta_{12} \mu_{12} - A^*(\mu) \}$$

Though this can be solved using the method of Lagrange multipliers (*c.f.*, Section 4.4.1), *etc.*, we expect the optimal solution to the variational problem to be

$$\hat{\mu}_1 = \frac{1}{z} \sum_{x_1 \in \{0,1\}, x_2 \in \{0,1\}} x_1 \exp \{ \theta_1 x_1 + \theta_2 x_2 + \theta_{12} x_1 x_2 \} = \frac{\exp \theta_1 + \exp \theta_1 + \theta_2 + \theta_{12}}{1 + \exp \theta_1 + \exp \theta_2 + \exp \theta_1 + \theta_2 + \theta_{12}}$$

There are many applications of variational inference to quantity estimation problems that have either no exact solutions, or that have solutions not computable in polynomial time. Further, variational principles can be used to study how the approximate algorithms behave; whether they have fix points, whether they converge. what answers do they give, *etc.*. For instance, in belief propagation from a variational perspective, the messages correspond to lagrange multipliers (for active constraints in \mathcal{M}) that are passed around.

In general, the sets \mathcal{M} are very complex. For example, with a complete 7 node graph, there will be $O(3 \times 10^8)$ constraints. For an n node tree, you will have $4(n - 1)$ constraints. The variational principle provides the foundation for many approximate methods such as the Naive mean field algorithm, which restricts optimization to a ‘tractable’ set of \mathcal{M} , such as one for which the joint distribution over a graphical model is factorized by treating the variables as completely independent.

Chapter 8

Statistical Relational Learning

One of the central open questions of artificial intelligence is concerned with combining (i) expressive knowledge representation formalisms such as relational and first-order logic with (ii) principled probabilistic and statistical approaches to inference and learning. Why? Here are some reasons:

1. The fields of knowledge representation and inductive logic programming stress the importance of relational and logical representations that provide the flexibility and modularity to model large domains. They also highlight the importance of making general statements, rather than making statements for every single aspect of the world separately.
2. The fields of statistical learning and uncertainty in artificial intelligence emphasize that agents that operate in the real world must deal with uncertainty. An agent typically receives only noisy or limited information about the world; actions are often non-deterministic; and an agent has to take care of unpredictable events. Probability theory provides a sound mathematical foundation for inference and learning under uncertainty.
3. Machine learning, in general, argues that an agent needs to be capable of improving its performance through experience.

Thus, the combination of expressive knowledge representation with probabilistic approaches to inference and learning is needed in order to face the challenges of real-world applications, which are complex and heterogeneous. Most traditional artificial intelligence and machine learning systems, however, are able to handle either uncertainty or rich relational structures but not both.

Statistical learning, reinforcement learning, and data mining methods have traditionally been developed for data in attribute-value form only; data is represented in matrix form: columns represent attributes, and rows represent exam-

ples. Indeed, matrices are simple and efficient matrix operations can be used. In turn, a matrix form makes it possible to devise efficient algorithms.

Many, if not most, real-world data sets, however, are not in matrix form. Applications contain several entities and relationships among them. Inductive logic programming and relational learning have been developed for coping with this type of data. They do not, however, handle uncertainty in a principled way.

It is therefore not surprising that there has been a significant interest in integrating statistical learning with first order logic and relational representations. This newly emerging research field is known under the name of statistical relational learning and probabilistic logic learning, and may be briefly defined as follows:

Definition 64 *Statistical relational learning deals with machine learning and data mining in relational domains where observations may be missing, partially observed, and/or noisy.*

Instead of giving a probabilistic characterization of logic programming such as [Ng and Subrahmanian, 1992], this line of research stresses the machine learning aspect.

8.1 Role of Logical Abstraction in SRL

Employing relational and logical abstraction within statistical learning has three advantages:

1. Variables, *i.e.*, placeholders for entities allow one to make abstraction of specific entities.
2. Unification allows one to share information among entities.
3. In many applications, there is a rich background theory available, which can efficiently and elegantly be represented as sets of general regularities.

Thus, instead of learning regularities for each single entity independently, statistical relational learning aims at finding general regularities among groups of entities. The learned knowledge is declarative and compact, which makes it much easier for people to understand and to validate. Although, the learned knowledge must be recombined at run time using some reasoning mechanism such as backward chaining or resolution, which bears additional computational costs, statistical relational models are more flexible, context-aware, and offer, in principle, the full power of logical reasoning. Further, background knowledge often improves the quality of learning as it focuses learning on relevant patterns, *i.e.*, restricts the search space. While learning, relational and logical abstraction allow one to reuse experience: learning about one entity improves the prediction for other entities; it might even generalize to objects, which have never been observed before. Thus, relational and logical abstraction can make statistical learning more robust and efficient. This has been proven beneficial in

many fascinating real-world applications in citation analysis, web mining, web navigation, web search, natural language processing, robotics, computer vision, social network analysis, bio-and chemo-informatics, electronic games, and activity recognition. For instance, Liao et al. [2005] applied Taskar et al.'s relational Markov networks to a problem related to the ride sharing service, namely learning and inferring transportation routines.

Whereas most of the existing works on statistical relational learning have started from a statistical and probabilistic learning perspective and extended probabilistic formalisms with relational aspects, several pieces of research take a different perspective, starting from inductive logic programming (ILP) and studying how inductive logic programming formalisms, settings and techniques can be extended to deal with probabilities.

8.2 Inductive Logic Programming

ILP is a research field at the intersection of machine learning and logic programming [Muggleton and De Raedt, 1994]. It is often also called multi-relational data mining (MRDM) [Dzeroski and Lavrac, 2001]. It aims at a formal framework as well as practical algorithms for inductively learning relational descriptions (in the form of logic programs) from examples and background knowledge. However, it does not explicitly deal with uncertainty such as missing or noisy information. Dealing explicitly with uncertainty makes probabilistic ILP more powerful than ILP and, in turn, than traditional attribute-value approaches. Moreover, there are several benefits of an ILP approach to statistical relational learning.

1. Classical ILP learning settings, as we will argue, naturally carry over to the probabilistic case. The probabilistic ILP settings make abstraction of specific probabilistic relational and first order logical representations and inference and learning algorithms yielding general statistical relational learning settings.
2. Many ILP concepts and techniques such as *more-general-than*, *refinement operators*, *least general generalization* (lub), and *greatest lower bound* (glb) can be reused. Therefore, many ILP learning algorithms such as Quinlan's FOIL and De Raedt and Dehaspes' Claudien can easily be adapted.
3. The ILP perspective highlights the importance of background knowledge within statistical relational learning. The research on ILP and on artificial intelligence in general has shown that background knowledge is the key to success in many applications.
4. An ILP approach should make statistical relational learning more intuitive to those coming from an ILP background and should cross-fertilize ideas developed in ILP and statistical learning.

Formally, ILP is concerned with finding a hypothesis H (a logic program, *i.e.* a definite clause program) from a set of positive and negative examples \mathcal{E}^+ and \mathcal{E}^- .

Consider learning a definition for the *Daughter/2* predicate, *i.e.*, a set of clauses with head predicates over *Daughter/2*, given the following facts as learning examples:

$$\begin{array}{l} \mathcal{E}^+ : \text{Daughter(dorothy, ann)}. \\ \qquad \text{Daughter(dorothy, brian)}. \\ \hline \mathcal{E}^- : \text{Daughter(rex, ann)}. \\ \qquad \text{Daughter(rex, brian)}. \end{array}$$

Additionally, we have some general knowledge called background knowledge \mathcal{B} , which describes the family relationships and sex of each person:

$$\begin{array}{lll} \text{Mother(ann, dorothy)}. & \text{Female(dorothy)}. & \text{Female(ann)}. \\ \text{Mother(ann, rex)}. & \text{Father(brian, dorothy)}. & \text{Father(brian, rex)}. \end{array}$$

From this information, we could induce the following H :

$$\begin{array}{l} \text{Daughter}(C, P) \quad : - \quad \text{Female}(C), \text{Mother}(P, C). \\ \text{Daughter}(C, P) \quad : - \quad \text{Female}(C), \text{Father}(P, C). \end{array}$$

which perfectly explains the examples in terms of the background knowledge, *i.e.*, \mathcal{E}^+ are entailed by H together with \mathcal{B} , but \mathcal{E}^- are not entailed.

Definition 65 *Given a set of positive and negative examples \mathcal{E}^+ and \mathcal{E}^- over some language \mathcal{L}_E , a background theory \mathcal{B} , in the form of a set of definite clauses, a hypothesis language \mathcal{L}_H , which specifies the clauses that are allowed in hypotheses, and a covers relation $\text{covers}(e, H, \mathcal{B}) \in \{0, 1\}$, which basically returns the classification of an example e with respect to \mathcal{H} and \mathcal{B} , find a hypothesis $H^* \in \mathcal{H}$ that covers (with respect to the background theory \mathcal{B}) all positive examples in \mathcal{E}^+ (completeness) and none of the negative examples in \mathcal{E}^- (consistency).*

In a noisy learning setting, the goal is to find a hypothesis $H^ \in \mathcal{H}$ that maximizes an objective score($\mathcal{E}, H, \mathcal{B}$).*

$$H^* = \underset{H}{\operatorname{argmax}} \text{score}(\mathcal{E}, H, \mathcal{B})$$

The (i) language \mathcal{L}_E chosen for representing the examples together with the (ii) covers relation determines the inductive logic programming setting De Raedt [1997]. There are three broad settings, *viz.*,

1. learning from entailment [Plotkin, 1970]
2. learning interpretations [Helft, 1989, De Raedt and Dzeroski, 1994]
3. an intermediate setting called learning from proofs [Shapiro 1983]

8.3 Probabilistic ILP Settings

The inductive logic programming settings can be extended to the probabilistic case. When working with probabilistic ILP representations, there are essentially two changes:

1. *Probabilistic Clauses:* Clauses in H and \mathcal{B} are annotated with probabilistic information, and Here, we use the following probability notations. With X , we denote a (random) variable. Furthermore, x denotes a state and \mathbf{X} (resp. \mathbf{x}) a set of variables (resp. states). We will use \Pr to denote a probability distribution, e.g., $\Pr(x)$, and p to denote a probability value, e.g., $p(X = x)$ and $p(X = x)$.
2. *Probabilistic Covers:* The covers relation becomes probabilistic. A probabilistic covers relation softens the hard covers relation employed in traditional ILP and is defined as the probability of an example given the hypothesis and the background theory. A probabilistic covers relation takes as arguments an example e , a hypothesis H and possibly the background theory \mathcal{B} , and returns the probability value $\Pr(e|H, \mathcal{B}) \in [0, 1]$ of the example e given H and \mathcal{B} , i.e., $\text{covers}(e, H, \mathcal{B}) = \Pr(e|H, \mathcal{B})$.

Using the probabilistic covers relation above, our first attempt at a definition of the probabilistic ILP learning problem is as follows:

Definition 66 *Given a probabilistic-logical language \mathcal{L}_H and a set \mathcal{E} of examples over some language \mathcal{L}_E , find the hypothesis H^* in \mathcal{L}_H that maximizes $\Pr(\mathcal{E}|H^*, \mathcal{B})$.*

Under the usual *i.i.d.* assumption, i.e., examples are sampled independently from identical distributions, this results in the maximization of

$$\Pr(\mathcal{E}|H^*, \mathcal{B}) = \prod_{e \in \mathcal{E}} P(e|H^*, \mathcal{B}) = \prod_{e \in \mathcal{E}} \text{covers}(e, H^*, \mathcal{B})$$

Similar to the ILP learning problem, the language \mathcal{L}_E selected for representing the examples together with the probabilistic covers relation determines different learning setting. Guided by Definition 66, we will introduce several probabilistic ILP settings for statistical relational learning. The main idea is to lift traditional ILP settings by associating probabilistic information with clauses and interpretations and by replacing ILPs deterministic covers relation by a probabilistic one. In the discussion, we will have one trivial but important observation:

Observation: Derivations might fail.

The probability of a failure is zero and, consequently, failures are never observable. Only succeeding derivations are observable, i.e., the probabilities of such derivations are greater than zero. As an extreme case, recall the negative

examples \mathcal{E}^- employed in the ILP learning problem on page 8.2. They are supposed to be not covered, *i.e.*,

$$p(\mathcal{E}^-|H, \mathcal{B}) = 0$$

In that example, *Rex* is a male person; he cannot be the daughter of *ann*. Thus, *Daughter(rex, ann)* was listed as a negative example. **Negative examples conflict with the usual view on learning examples in statistical learning.** In statistical learning, we seek to find that hypothesis H^* , which is most likely given the learning examples:

$$H^* = \underset{H}{\operatorname{argmax}} p(H|\mathcal{E}) = \underset{H}{\operatorname{argmax}} \frac{p(\mathcal{E}|H)p(H)}{p(\mathcal{E})} \text{ with } p(\mathcal{E}) > 0$$

Thus, examples \mathcal{E} are observable, *i.e.*, $p(\mathcal{E}) > 0$. Therefore, we refine the preliminary probabilistic ILP learning problem definition 66. In contrast to the purely logical case of ILP, we do not speak of positive and negative examples anymore but of possible and impossible ones.

Definition 67 *Given a set $\mathcal{E} = \mathcal{E}_p \cup \mathcal{E}_i$ of possible and impossible examples \mathcal{E}_p and \mathcal{E}_i (with $\mathcal{E}_p \cap \mathcal{E}_i = \emptyset$) over some example language \mathcal{L}_E , a probabilistic covers relation $\operatorname{covers}(e, H, \mathcal{B}) = P(e|H, \mathcal{B})$, a probabilistic logical language \mathcal{L}_H for hypotheses, and a background theory \mathcal{B} , find a hypothesis H^* in \mathcal{L}_H such that $H^* = \underset{H}{\operatorname{argmax}} \operatorname{score}(\mathcal{E}, H, \mathcal{B})$ and the following constraints hold:*

$$\forall e_p \in \mathcal{E}_p : \operatorname{covers}(e_p, H^*, \mathcal{B}) > 0 \quad \text{and} \quad \forall e_i \in \mathcal{E}_i : \operatorname{covers}(e_i, H^*, \mathcal{B}) = 0$$

The scoring function is some objective score, usually involving the probabilistic covers relation of the possible examples such as the observed likelihood

$$\prod_{e_p \in \mathcal{E}_p} \operatorname{covers}(e_p, H^*, \mathcal{B})$$

some penalized variant thereof.

The probabilistic ILP learning problem of Definition 67 unifies ILP and statistical learning in the following sense:

1. Using a deterministic covers relation (which is either 1 or 0) yields the classical ILP learning problem, see Definition 65.
2. On the other hand, sticking to propositional logic and learning from possible examples, *i.e.*, $P(\mathcal{E}) > 0$, only yields traditional statistical learning.
3. Definition 67 makes abstraction of many particular kinds of problems.
 - In density estimation, the joint probability distribution of some random variables is estimated

- Whereas, in classification and regression the dependency of a discrete respectively continuous target variable given the value of some other variables is estimated.
4. Furthermore, several types of learning can be distinguished.
 - In supervised learning, the training examples contain information about all variables including the target variable.
 - In reinforcement learning, the training examples contain only indirect target information such as the classifier did well or not (in the form of some reward).
 - Finally, in unsupervised learning, no values of the target variable are observed.
 5. Another important distinction is whether all the random variables variables are observed, or whether some of them are hidden, *e.g.*, they are specified in the background knowledge and never observed.
 6. We can also formulate the learning problem as either of the following:
 - As a ‘point estimation problem, *i.e.*, the goal is to find a single best hypothesis H^* .
 - As a ‘Bayesian learning problem’, where the goal is to return a posterior distribution over hypotheses.
 7. Finally, learning might refer to the structure, *i.e.*, the underlying logic program of the hypothesis, the parameters, or both. To come up with algorithms solving probabilistic ILP learning problems, say for density estimation, one typically distinguishes two subtasks because $H = (L, \lambda)$ is essentially a logic program L annotated with probabilistic parameters λ :
 - Parameter estimation where it is assumed that the underlying logic program L is fixed, and the learning task consists of estimating the parameters λ that maximize the likelihood.
 - Structure learning where both L and λ have to be learned from the data.

Similar to that for traditional ILP, there are three probabilistic ILP settings, which extend the purely logical ones. The three ILP settings and their probabilistic extensions are outlined in the following sections.

8.3.1 Probabilistic setting for justification

There are two important components in SRL:

Abduction. Process of hypothesis formation. The logical setting for abduction is what traditional ILP is mostly about.

Justification. The degree of belief assigned to an hypothesis given a certain amount of evidence.

Recall the Bayes' Theorem

$$p(h|E) = \frac{p(h) \cdot p(E|h)}{p(E)}$$

The best hypothesis in a set \mathcal{H} (ignoring ties)

$$H = \operatorname{argmax}_{h \in \mathcal{H}} p(h|E)$$

We will consider the following learning framework.

Let X be a countable set of instances (encodings of all objects of interest) and D_X be a probability measure on X

Let $\mathcal{C} \subseteq 2^X$ be a countable set of concepts and $D_{\mathcal{C}}$ be a probability measure on 2^X

Let \mathcal{H} be a countable set of hypotheses¹ and $D_{\mathcal{H}}$ be a probability measure (prior) over \mathcal{H}

Let the concept represented by $h \in \mathcal{H}$ be $c(h) \in \mathcal{C}$

Let \mathcal{C} and \mathcal{H} be such that

- for each $C \in \mathcal{C}$, there is an $h \in \mathcal{H}$ s.t. $C = c(h)$
- for each $C \in \mathcal{C}$, $D_{\mathcal{C}}(C) = \sum_{\{h \in \mathcal{H} | C = c(h)\}} P(h)$

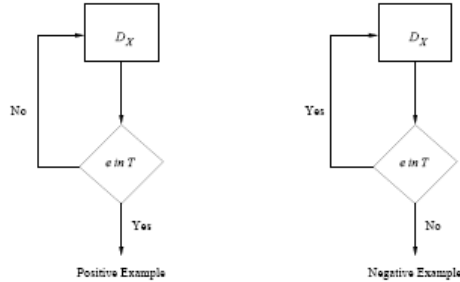
Target concept T is chosen using the distribution $D_{\mathcal{C}}$

Let $g(h)$ denote the proportion (w.r.t. the instance space) of the concept represented by a hypothesis $h \in \mathcal{H}$

- That is, $g(h) = \sum_{x \in c(h)} D_X(x)$
- $g(h)$ is a measure of the "generality" of h

Noise Free Data

Following is a model for noise free data.



¹Note that \mathcal{H} also been used earlier to denote the space of horn clauses. This time, it is used to denote the space of hypothesis, which anyways happens to be within the space of horn clauses.

Given $\mathcal{E} = \mathcal{E}^+ \cup \mathcal{E}^-$

$$p(h|E) \propto D_{\mathcal{H}}(h) \prod_{e \in E^+} p(e|h) \prod_{e \in E^-} p(e|h)$$

Or

$$P(h|E) \propto D_{\mathcal{H}}(h) \prod_{e \in E^+} \frac{D_X(e)}{g(h)} \prod_{e \in E^-} \frac{D_X(e)}{1-g(h)}$$

Assuming p positive and n negative examples

$$P(h|E) \propto D_{\mathcal{H}}(h) \left(\prod_{e \in E} D_X(e) \right) \left(\frac{1}{g(h)} \right)^p \left(\frac{1}{1-g(h)} \right)^n$$

Maximal $P(h|E)$ means finding the hypothesis that maximises

$$\log D_{\mathcal{H}}(h) + p \log \frac{1}{g(h)} + n \log \frac{1}{1-g(h)}$$

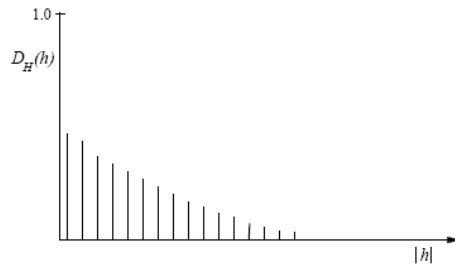
If there are no negative examples, then this becomes

$$\log D_{\mathcal{H}}(h) + p \log \frac{1}{g(h)}$$

Let us answer some questions at this point:

1. *What is the distribution $D_{\mathcal{H}}(h)$:*

A common assumption is that “larger” programs are less likely (in coding terminology, require more bits to encode)



As an example

$$D_{\mathcal{H}}(h) = 2^{-|h|}$$

That is

$$\log D_{\mathcal{H}}(h) = -|h|$$

2. What is generality function $g(h)$?

Recall that $g(h) = \sum_{x \in c(h)} D_X(x)$

- $c(h)$ may be infinite
- D_X is usually unknown (and is a mapping to the reals)

Have to be satisfied with approximate estimates of $g(h)$

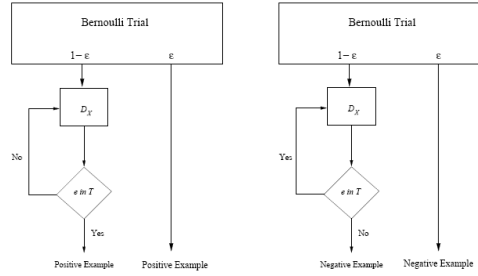
Estimation procedure

- (a) Randomly generate a finite sample of n instances using a known distribution (for eg. uniform)
- (b) Determine the number of these instances (say c) entailed by h
- (c) $g(h) \approx \frac{c+1}{n+2}$

3. What about noisy data? This will be addressed in the next subsection.

A Model for Noisy Data

Following is a model for justification in the presence of noisy data.



For any hypothesis h the examples $\mathcal{E} = \mathcal{E}^+ \cup \mathcal{E}^-$ can now be partitioned as follows

1. $TP = \{e | e \in E^+ \text{ and } e \in c(h)\}$ (true positives)
2. $FN = \{e | e \in E^+ \text{ and } e \notin c(h)\}$ (false negatives)
3. $FP = \{e | e \in E^- \text{ and } e \in c(h)\}$ (false positives)
4. $TN = \{e | e \in E^- \text{ and } e \notin c(h)\}$ (true negatives)

Recall

$$p(h|E) \propto D_{\mathcal{H}}(h) \prod_{e \in E^+} p(e|h) \prod_{e \in E^-} p(e|h)$$

Now

$$\prod_{e \in E^+} p(e|h) = \prod_{e \in TP} \left(\frac{D_X(e)(1-\epsilon)}{g(h)} + D_X(e)\epsilon \right) \prod_{e \in FN} D_X(e)\epsilon$$

$$\prod_{e \in E^-} p(e|h) = \prod_{e \in TN} \left(\frac{D_X(e)(1-\epsilon)}{1-g(h)} + D_X(e)\epsilon \right) \prod_{e \in FP} D_X(e)\epsilon$$

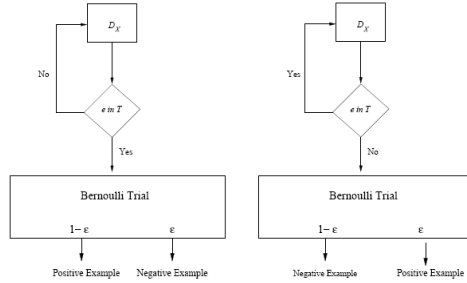
So, with $FPN = FP \cup FN$

$$p(h|E) \propto D_{\mathcal{H}}(h) \left(\prod_{e \in E} D_X(e) \right) \left(\frac{1-\epsilon}{g(h)} \right)^{|TP|} \left(\frac{1-\epsilon}{1-g(h)} \right)^{|TN|} \epsilon^{|FPN|}$$

Maximal $P(h|E)$ means finding the hypothesis that maximises

$$\log D_{\mathcal{H}}(h) + |TP| \log \frac{1-\epsilon}{g(h)} + |TN| \log \frac{1-\epsilon}{1-g(h)} + |FPN| \log \epsilon$$

Another model for noisy data is outlined below:



8.4 Learning from Entailment

Learning from entailment is by far the most popular ILP setting and it is addressed by a wide variety of well-known ILP systems such as FOIL [Quinlan and Cameron-Jones, 1995], Progol [Muggleton, 1995], and Aleph [Srinivasan, 1999]. When learning from entailment, the examples are definite clauses and a hypothesis H covers an example e with respect to the background theory \mathcal{B} if and only if $\mathcal{B} \cup H \models e$, *i.e.*, each model of $\mathcal{B} \cup H$ is also a model of e .

In many well-known systems, such as FOIL, one requires that the examples are ground facts. To illustrate the above setting, consider the following example inspired on the well-known mutagenicity application [Srinivasan et al., 1996].

Consider the following facts in the background theory \mathcal{B} , which describe part of molecule 225.

<i>Molecule</i> (225).	<i>Logmutag</i> (225, 0.64).
<i>Lumo</i> (225, -1.785).	<i>Logp</i> (225, 1.01).
<i>Nitro</i> (225, [<i>f1_4</i> , <i>f1_8</i> , <i>f1_10</i> , <i>f1_9</i>]).	<i>Atom</i> (225, <i>f1_1</i> , <i>c</i> , 21, 0.187).
<i>Atom</i> (225, <i>f1_2</i> , <i>c</i> , 21, -0.143).	<i>Atom</i> (225, <i>f1_3</i> , <i>c</i> , 21, -0.143).
<i>Atom</i> (225, <i>f1_4</i> , <i>c</i> , 21, -0.013).	<i>Atom</i> (225, <i>f1_5</i> , <i>o</i> , 52, -0.043).
<i>Bond</i> (225, <i>f1_1</i> , <i>f1_2</i> , 7).	<i>Bond</i> (225, <i>f1_2</i> , <i>f1_3</i> , 7).
<i>Bond</i> (225, <i>f1_3</i> , <i>f1_4</i> , 7).	<i>Bond</i> (225, <i>f1_4</i> , <i>f1_5</i> , 7).
<i>Bond</i> (225, <i>f1_5</i> , <i>f1_1</i> , 7).	<i>Bond</i> (225, <i>f1_8</i> , <i>f1_9</i> , 2).
<i>Bond</i> (225, <i>f1_8</i> , <i>f1_10</i> , 2).	<i>Bond</i> (225, <i>f1_11</i> , <i>f1_11</i> , 1).
<i>Ring_size_5</i> (225, [<i>f1_5</i> , <i>f1_1</i> , <i>f1_2</i> , <i>f1_3</i> , <i>f1_4</i>]).	<i>Bond</i> (225, <i>f1_11</i> , <i>f1_12</i> , 2).
<i>Hetero_aromatic_5_ring</i> (225, [<i>f1_5</i> , <i>f1_1</i> , <i>f1_2</i> , <i>f1_3</i> , <i>f1_4</i>]).	<i>Bond</i> (225, <i>f1_11</i> , <i>f1_13</i> , 1).
...	...

Consider now the positive example *Mutagenic*(225). It is covered by the following *H*:

$$\textit{Mutagenic}(M) \quad :- \quad \textit{Nitro}(M, R1), \textit{Logp}(M, C), C > 1.$$

together with the background knowledge \mathcal{B} , because $H \cup \mathcal{B}$ entails the example. To see this, we unify *Mutagenic*(225) with the clauses (*H*) head. This yields

$$\textit{Mutagenic}(225) \quad :- \quad \textit{Nitro}(225, R1), \textit{Logp}(225, C), C > 1.$$

Now, *Nitro*(225, *R1*) unifies with the third ground atom (left-hand side column) in \mathcal{B} , and *Logp*(225, *C*) with the second one on the right. Because $1.01 > 1$, we found a proof of *mutagenic*(225).

There are, broadly speaking, three types of ILP approaches for learning from entailment: *top-down approaches*, *bottom-up approaches* and *hybrid approaches*.

1. *Top-down approaches* start from short clauses, iteratively adding literals to their bodies as long as they do not become too general. Basically, in top-down approaches, hypotheses are generated in a pre-determined order, and then tested against the examples. More precisely:
 - (a) They start with the most general hypothesis, *i.e.*, clauses of the form $\textit{Daughter}(C, P) : -\textit{True}$, where all arguments are distinct variables.
 - (b) After seeing the first example that contradicts the hypothesis, *i.e.*, after seeing the first negative example, the hypothesis is specialized by specializing the general clause.
 - (c) The clause is specialized typically in two ways:
 - by applying a substitution,
 - by adding a literal, *i.e.*, an atom or its negation to the body, and

For instance, we can specialize $Daughter(C, P) : \text{True}$ and consider $Daughter(C, P) : \text{Female}(C)$ and $Daughter(C, P) : \text{Mother}(P, C)$ for further investigations. Several possibilities (and successive specializations) have to be tried before one finds a clause that covers some positive examples but no negative ones such as $Daughter(C, P) : \text{Female}(C), \text{Mother}(P, C)$.

- (d) If some positive examples are still not covered, these techniques typically add a new, maximally general clause to the hypothesis and essentially iterate the process in step 1c as before, until all positive examples are covered and no negative example is covered.
- (e) The background knowledge \mathcal{B} is typically viewed as a logic program (*i.e.* a definite clause program) that is provided to the inductive logic programming system and fixed during the learning process. The hypothesis H together with the background theory \mathcal{B} should cover all positive and none of the negative examples.

Top-down approaches are often employed by ILP systems that learn from entailment. More precisely, these systems often employ a separate-and-conquer rule-learning strategy [Furnkranz, 1999]. In an outer loop of the algorithm, they follow a set-covering approach [Mitchell, 1997] in which they repeatedly search for a rule covering many positive examples and none of the negative examples. They then delete the positive examples covered by the current clause and repeat this process until all positive examples have been covered. In the inner loop of the algorithm, they typically refine a clause by unifying variables, by instantiating variables to constants, and/or by adding literals to the clause.

Figure 8.1 presents the generic framework for top-down ILP systems while Figure 8.2 outlines one of the original ILP systems, MIS.

2. *Bottom-up approaches* start from long clauses, iteratively removing literals until they would become overly general. While top-down approaches successively specialize a very general starting hypothesis, bottom-up approaches successively generalize a very specific hypothesis. This is basically done
 - (a) by deleting literals (or clauses),
 - (b) by turning constants into variables and/or
 - (c) by turning bounded variables into new variables.

At each step, the theory is generalized by taking the least general generalization (under \cdot -subsumption) of pairwise clauses. Of course, care must be taken that the generalized theory does not cover negative examples.

3. *Hybrid approaches* that mix top-down and bottom-up searches. Hybrid approaches are usually employed for multiple predicate learning [De Raedt et al., 1993] and theory revision [Wrobel, 1996].

```

Initialize  $\mathcal{E}_{cur} = \mathcal{E}$ ;
Initialize  $\mathcal{H} = \emptyset$ ;
//Start covering
repeat
  Initialize  $C = T \leftarrow$ ;
  //Start specialization
  repeat
    Find the best refinement  $C_{best} \in \rho_{\mathcal{L}}(C)$ ;
    Assign  $C = C_{best}$ ;
  until Necessity stopping criterion is satisfied;
  Add  $C$  to  $\mathcal{H}$  to get new hypothesis  $\mathcal{H}' = \mathcal{H} \cup C$ ;
  Remove positive examples covered by  $C$  from  $\mathcal{E}_{cur}$  to get new training set
   $\mathcal{E}'_{cur} = \mathcal{E}_{cur} - \bigcup_{\{e \in \mathcal{E}_{cur}^+ | covers(B, \mathcal{H}', e) = 1\}} \{e\}$ ;
  Assign  $\mathcal{E}_{cur} = \mathcal{E}'_{cur}$ ,  $\mathcal{H} = \mathcal{H}'$ ;
until Sufficiency stopping criterion is satisfied;

```

Figure 8.1: A general strategy for top-down ILP approaches.

8.4.1 Constraining the ILP Search

It should be stressed that ILP is a difficult problem. Practical ILP systems fight the inherent complexity of the problem by imposing all sorts of constraints, mostly syntactic in nature. Such constraints include

1. *Language and search biases*: These are sometimes summarized as declarative biases, (see [Nedellec et al., 1996] for an overview). Essentially, the main source of complexity in ILP stems from the variables in the clauses. In top-down systems, the branching factor of the specialization operator increases with the number of variables in the clauses. Following are frequently adopted techniques for reducing this branching factor by introducing language and search bias.
 - (a) Introducing types for predicates can rule out main potential substitutions and unifications. As an example, the type definition $type(Father(person, person))$ specifies that both argument of atoms over $Father/2$ have to be persons.
 - (b) Refinement operators can also be used to encode a language bias, since they can be restricted to generate only a subset of the language \mathcal{L}_H . For instance, refinement operators can easily be modified to generate only constant-free and function-free clauses.
 - (c) Other methods use a kind of grammar construction to explicitly declare the range of acceptable clauses, see e.g. Cohen [1994].
 - (d) Lookaheads are an example of a search bias. In some cases, an atom might never be chosen by our algorithm because it will not, in itself,

```

Initialize hypothesis  $H$  to a (possibly empty) set of clauses in  $\mathcal{L}_H$ ;
loop
  Read the next (positive or negative) example;
  repeat
    if There exists an  $e \in \mathcal{E}^-$  covered by  $H$  then
      Delete incorrect clauses from  $H$ .
    end if
    if There exists an  $e \in \mathcal{E}^+$  not covered by  $H$  then
      With breadth-first search of the refinement graph develop a clause  $C$ 
      which covers  $e$  and add it to  $H$ ;
    end if
  until  $H$  is complete and consistent;
  return Hypothesis  $H$ ;
end loop

```

Figure 8.2: A simplified version of the MIS algorithm by [Shapiro 1983], which is the general strategy of top-down approaches.

result in a better score. However, such an atom, while not useful in itself, might introduce new variables that make a better coverage possible by adding another atoms later on [Quinlan, 1991] 7. It is usually solved by allowing the algorithm to look ahead in the search space. Instead of considering refinements with a single atom, one considers larger refinements consisting of multiple atoms [Blockeel and De Raedt, 1997].

2. *Bound on number of distinct variables:* One can also put a bound in the number of distinct variables that can occur in clauses.
3. *Mode declarations:* Mode declarations are another well-known ILP device. They are used to describe input-output behaviour of predicate definitions. For example, we might specify $mode(Daughter(+, -))$ and $mode(Father(-, +))$, meaning that the $+$ arguments must be instantiated, whereas the $-$ arguments will be bounded to the answer.

In general, a model can suffer from either underfitting or overfitting. A model that is not sufficiently complex can fail to fully detect the underlying rule of a complicated data set, leading to underfitting. A model that is too complex may fit the noise, not just the underlying rule, leading to overfitting and, for instance, wild predictions.

8.4.2 Example: Aleph

As an example for the learning from entailments setting, we will discuss Aleph. Given background knowledge \mathcal{B} and positive examples $\mathcal{E}^+ = e_1 \wedge e_2 \dots$, negative examples $\mathcal{E}^- = \overline{f_1} \wedge \overline{f_2} \wedge \dots$, the generic ILP system Aleph, is concerned

with finding a hypothesis $H = D_1 \wedge \dots$ that satisfies (note: \cup and \wedge used interchangeably)

Prior Satisfiability. $B \wedge \mathcal{E}^- \not\models \square$

Prior Necessity. $B \not\models \mathcal{E}^+$

Posterior Sufficiency. $\mathcal{B} \wedge H \models \mathcal{E}^+$ and $\mathcal{B} \wedge D_j \models e_1 \vee e_2 \vee \dots$

Posterior Satisfiability. $\mathcal{B} \wedge H \wedge \mathcal{E}^- \not\models \square$

If more than one H satisfies this, the one with highest posterior probability is chosen. The D_i can be found by examining clauses that “relatively subsume” at least one example, making use of plotkin’s relative subsumption discussed elaborately in Section 2.5.1. This gives the following sufficient implementation of Aleph, given \mathcal{B}, \mathcal{E} .

1. $H_0 = \mathcal{B}, i = 0, \mathcal{E}^+ = \{e_1, \dots, e_n\}$
2. repeat
 - (a) increment i
 - (b) Obtain the most specific clause $\perp(\mathcal{B}, e_i)$
 - (c) Find the clause D_i that: subsumes $\perp(\mathcal{B}, e_i)$; and is consistent with the negative examples;
 - (d) $H_i = H_{i-1} \cup \{D_i\}$
3. until $i > n$
4. return H_n

As discussed in Section 2.5.1, there are problems with this implementation. Particularly,

- $\perp(\mathcal{B}, e_i)$ may be infinite
- This implementation may perform a lot of redundant computation ($D_i \in H_{i-1}$)
- This implementation need not return the hypothesis with maximum posterior probability

So we next present a “Greedy” implementation for Aleph, given \mathcal{B}, \mathcal{E} .

1. $H_0 = \mathcal{B}, \mathcal{E}_0^+ = \mathcal{E}^+, i = 0$.
2. Repeat
 - (a) Increment i
 - (b) Randomly choose a positive example e_i from \mathcal{E}_{i-1}^+
 - (c) Obtain the most specific clause $\perp(\mathcal{B}, e_i)$

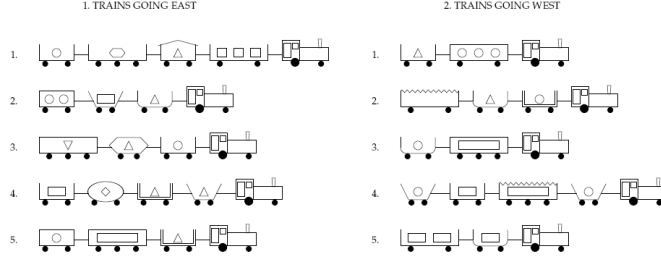


Figure 8.3: The train-spotting example for illustrating Aleph.

- (d) Find the clause D_i that: subsumes $\perp(\mathcal{B}, e_i)$; and is consistent with the negative examples; and maximises $p(H_{i-1} \cup \{D_i\} | e_i^+ \cup \mathcal{E}^-)$ where e_i^+ are the examples in \mathcal{E}^+ made redundant by $H_{i-1} \cup \{D_i\}$
 - (e) $H_i = H_{i-1} \cup \{D_i\}$
 - (f) $\mathcal{E}_i^+ = \mathcal{E}_{i-1}^+ \setminus e_i^+$
3. until $\mathcal{E}_i^+ = \emptyset$
 4. return H_i

However, this implementation does not address the problem that $\perp(\mathcal{B}, e_i)$ may be infinite and that it need not return in the hypothesis with maximum posterior probability. The problem of infinite $\perp(\mathcal{B}, e_i)$ may be addressed by making use of mode declarations introduced on page 120. This gives a revised “Greedy” implementation for Aleph, given $\mathcal{B}, \mathcal{E}, d$.

1. $H_0 = \mathcal{B}, \mathcal{E}_0^+ = \mathcal{E}^+, i = 0$
2. Repeat
 - (a) Increment i
 - (b) Randomly choose a positive example e_i from \mathcal{E}_{i-1}^+
 - (c) Obtain the most specific clause $\perp_d(\mathcal{B}, e_i)$
 - (d) Find the clause D_i that: subsumes $\perp(\mathcal{B}, e_i)$; and is consistent with the negative examples; and maximises $p(H_{i-1} \cup \{D_i\} | e_i^+ \cup \mathcal{E}^-)$ where e_i^+ are the examples in \mathcal{E}^+ made redundant by $H_{i-1} \cup \{D_i\}$
 - (e) $H_i = H_{i-1} \cup \{D_i\}$
 - (f) $\mathcal{E}_i^+ = \mathcal{E}_{i-1}^+ \setminus e_i^+$
3. until $\mathcal{E}_i^+ = \emptyset$
4. return H_i

As an example, let us consider the trainspotting problem (*c.f.* Figure 8.3) to illustrate the discussion thus far.

We will use the following mode declarations for the trainspotting example.


```

:- modeb(1, eastbound(+train)).
:- modeb(1, short(+car)).
:- modeb(1, closed(+car)).
:- modeb(1, long(+car)).
:- modeb(1, open_car(+car)).
:- modeb(1, double(+car)).
:- modeb(1, jagged(+car)).
:- modeb(1, shape(+car, #shape)).
:- modeb(1, load(+car, #shape, #int)).
:- modeb(1, wheels(+car, #int)).
:- modeb(*, has_car(+train, -car)).

```

Further, the examples \mathcal{E}^+ and \mathcal{E}^- for this example will be:

Positive	Negative
eastbound(east1).	eastbound(west6).
eastbound(east2).	eastbound(west7).
eastbound(east3).	eastbound(west8).
eastbound(east4).	eastbound(west9).
eastbound(east5).	eastbound(west10).

whereas, the background knowledge \mathcal{B} will comprise:

```

% type definitions
car(car_11). car(car_12). ...
car(car_21). car(car_22). ...
...

shape(ellipse). shape(hexagon). ...
...

% eastbound train 1
has_car(east1, car_11). has_car(east1, car_12). ...
shape(car_11, rectangle). shape(car_12, rectangle). ...
open_car(car_11). closed(car_12).
long(car_11). short(car_12). ...
...

% westbound train 6
has_car(west6, car_61). has_car(west6, car_62). ...
long(car_61). short(car_62).
shape(car_61, rectangle). shape(car_62, rectangle).
...

```

The outcome of a simple search in Aleph, for the trainspotting setting is outlined next:

```

eastbound(A) :-
    has_car(A,B).

```

```

[5/5]
eastbound(A) :-
    has_car(A,B), short(B).
[5/5]
eastbound(A) :-
    has_car(A,B), open_car(B).
[5/5]
eastbound(A) :-
    has_car(A,B), shape(B,rectangle).
[5/5]
...

[theory]

[Rule 1] [Pos cover = 5 Neg cover = 0]

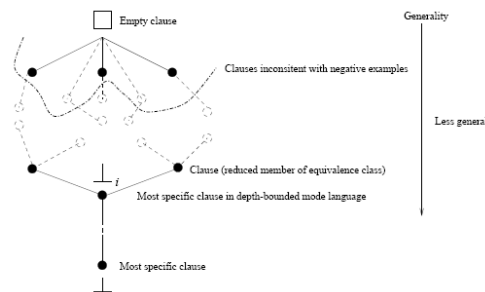
eastbound(A) :-
    has_car(A,B), short(B), closed(B).

[pos-neg] [5]

```

We will next discuss the search and redundancy aspects of Aleph. There are two stages in the clause-by-clause construction of hypothesis, which are summarized below:

1. Search



2. Remove redundant clauses once best clause is found

Aleph has provision for bottom-up as well as top-down search for moving about in the lattice. The following refinement steps can be used in Aleph.

General-to-specific search: start at \square , and move by

1. Adding a literal drawn from \perp_i

$$p(X, Y) \leftarrow q(X) \text{ becomes } p(X, Y) \leftarrow q(X), r(Y)$$

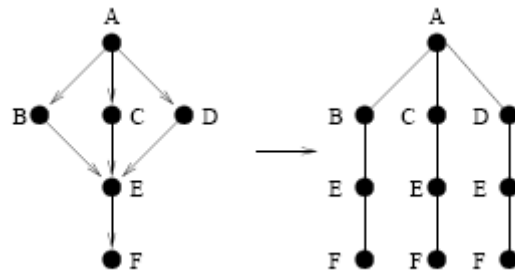


Figure 8.4: The conversion of (relative) subsumption lattice to tree.

2. Equating two variables of the same type

$$p(X, Y) \leftarrow q(X) \text{ becomes } p(X, X) \leftarrow q(X)$$

3. Instantiate a variable with a general functional term or constant

$$p(X, Y) \leftarrow q(X) \text{ becomes } p(3, Y) \leftarrow q(3)$$

Specific-to-general search: start at \perp_i

Dual operations to above (ie. remove literal etc.)

Progol does a general-to-specific search

Next, we discuss some search methods. The subsumption lattice can be represented as a directed acyclic graph. However, the DAG can be converted to a tree such that root is the first node (\square or \perp_i) and children of a node are refinements. The conversion of a lattice to a search tree is illustrated in Figure 8.4

Searching the lattice is therefore equivalent to searching a tree. Our goal is to find the node (goal node) that has greatest “compression”. There are two basic types of tree search: depth-first (DF) and breadth-first (BF) as discussed in chapter 6. DF and BF are “blind”. More guidance is desirable at any node s . The guidance can be in one or more of the following forms:

- g_s : cost of optimal path from root to s
- h_s : estimated cost of optimal path to goal from s

The different kinds of guided search can be summarized as follows:

Hill-climbing: DF with h_s

Best-first: BF with h_s

Best-cost: BF with g_s

A^* : BF with g_s and h_s

Progol does an A^* -like search – uses utilit instead of costs. For a clause C at a node

$$\begin{aligned}
 p_s &= |\{e : e \in E^+ \text{ and } B \wedge C \wedge \bar{e} \models \square\}| \\
 n_s &= |\{e : e \in E^- \text{ and } B \wedge C \wedge e \models \square\}| \\
 c_s &= |C| - 1 \\
 l_s &= \text{l.b. on literals needed} \\
 h_s &= -(n_s + c_s + l_s) \\
 g_s &= p_s \\
 f_s &= g_s + h_s
 \end{aligned}$$

Compression of a clause at node $s = p_s - n_s - c_s$ Progol seeks clauses with $c_s < c$ and $n_s = 0$ and guaranteed to have maximum $p_s - c_s - h_s$. Progol has several admissible pruning strategies to reduce search.

E.g. Let $children(s) = S$. If $f(s) > p_{s'} - c_{s'}$ for $s' \in S$ then $prune(s')$

We next outline an optimal search algorithm that uses branch-and-bound.

$bb(i, \rho, f)$: Given an initial element i from a discrete set S ; a successor function $\rho : S \rightarrow 2^S$; and a cost function $f : S \rightarrow \mathfrak{R}$, return $H \subseteq S$ such that H contains the set of cost-minimal models. That is for all $h_{i,j} \in H$, $f(h_i) = f(h_j) = f_{min}$ and for all $s' \in S \setminus H$ $f(s') > f_{min}$.

1. $Active := \langle i \rangle$.
2. $best := \infty$
3. $selected := \emptyset$
4. while $Active \neq \langle \rangle$
5. begin
 - (a) remove element k from $Active$
 - (b) $cost := f(k)$
 - (c) if $cost < best$
 - (d) begin
 - i. $best := cost$
 - ii. $selected := \{k\}$
 - iii. let $Prune_1 \subseteq Active$ s.t. for each $j \in Prune_1$, $\underline{f}(j) > best$ where $\underline{f}(j)$ is the lowest cost possible from j or its successors
 - iv. remove elements of $Prune_1$ from $Active$
 - (e) end
 - (f) elseif $cost = best$
 - i. $selected := selected \cup \{k\}$
 - (g) $Branch := \rho(k)$
 - (h) let $Prune_2 \subseteq Branch$ s.t. for each $j \in Prune_2$, $\underline{f}(j) > best$ where $\underline{f}(j)$ is the lowest cost possible from j or its successors
 - (i) $Bound := Branch \setminus Prune_2$
 - (j) add elements of $Bound$ to $Active$

6. end
7. return *selected*

Different search methods result from specific implementations of *Active*:

- Stack: depth-first search
- Queue: breadth-first search
- Prioritised Queue: best-first search

There are two types of redundancies in the hypothesis H :

1. Redundancy 1: Literal Redundancy:

Literal l is redundant in clause $C \vee l$ relative to background B iff

$$B \wedge (C \vee l) \equiv B \wedge C$$

It can be shown that the literal l is redundant in clause $C \vee l$ relative to the background B iff

$$B \wedge (C \vee l) \models C$$

The clause C is said to be reduced (*c.f.* Section 2.1, page 102) with respect to background knowledge B iff no literal in C is redundant.

2. Redundancy 2: Clause redundancy:

Clause C is redundant in the $B \wedge C$ iff $B \wedge C \equiv B$.

It can be shown that clause C is redundant in $B \wedge C$ iff

$$B \models C \equiv B \wedge \bar{C} \models \square$$

A set of clauses S is said to be reduced iff no clause in S is redundant. Progol uses this procedure to determine (and remove) examples made redundant by clause found in the search.

Example

e_j : $gfather(henry, john) \leftarrow$

B : $father(henry, jane) \leftarrow$
 $father(henry, joe) \leftarrow$
 $parent(jane, john) \leftarrow$
 $parent(joe, robert) \leftarrow$

D_j : $gfather(X, Y) \leftarrow father(X, Z), parent(Z, Y)$

e_j is redundant in $B \wedge D_j \wedge e_j$ since $B \wedge D_j \wedge \bar{e}_j \models \square$

There are several implementation issues that need further thought.

Question. Will the clause-by-clause search method yield the best set of clauses?
If no, why not?

Question. Is it possible to do a theory-by-theory search?

Question. Is it possible devise a complete search that is non-redundant? If no, why not?

8.5 Probabilistic Learning from Entailment

Probabilistic learning from entailment has been investigated for learning stochastic logic programs [Muggleton, 2000a,b, Cussens, 2001, Muggleton, 2002] and for parameter estimation of PRISM programs [Sato and Kameya, 2001, Kameya et al., 2004] from possible examples only.

In order to integrate probabilities in the entailment setting, we need to find a way to assign probabilities to clauses that are entailed by an annotated logic program. Since most ILP systems working under entailment employ ground facts for a single predicate as examples, we will restrict our attention to assign probabilities to facts for a single predicate².

More formally, let us annotate a logic program H consisting of a set of clauses of the form $q \leftarrow b_i$, where p is an atom of the form $p(V_1, \dots, V_n)$ with the V_i different variables, and the b_i are different bodies of clauses. Furthermore, we associate to each clause in H the probability values $p(b_i|q)$; they constitute the conditional probability distribution that for a random substitution θ for which $q\theta$ is ground and true (resp. false), the query $b_i\theta$ succeeds (resp. fails) in the knowledge base \mathcal{B} . Recall that the query q succeeds in \mathcal{B} if there is a substitution σ such that $\mathcal{B} \models q\sigma$. Furthermore, we assume the prior probability of q is given as $p(q)$, which denotes the probability that for a random substitution θ , $p\theta$ is true (resp. false). This can then be used to define the covers relation $p(q\theta | H, \mathcal{B})$ as follows (we delete the \mathcal{B} as it is fixed):

$$p(q\theta|H) = p(q\theta | b_1\theta, \dots, b_k\theta) = \frac{p(b_1\theta, \dots, b_k\theta|q\theta) \times p(q\theta)}{p(b_1\theta, \dots, b_k\theta)}$$

For instance, applying the naive Bayes assumption yields

$$p(q\theta|H) = \frac{\prod_i p(b_i\theta|q\theta) \times p(q\theta)}{p(b_1\theta, \dots, b_k\theta)} \quad (8.1)$$

²It remains an open question as how to formulate more general frameworks for working with entailment.

Finally, since $p(q\theta | H) + p(\neg q\theta | H) = 1$, we can compute $p(q\theta|H)$ without $p(b_1\theta, \dots, b_k\theta)$ through normalization.

As an example, consider again the mutagenicity domain and the following annotated logic program:

$$(0.01, 0.21) : \text{Mutagenetic}(M) \leftarrow \text{Atom}(M, -, -, 8, -)$$

$$(0.38, 0.99) : \text{Mutagenetic}(M) \leftarrow \text{Bond}(M, -, A, 1), \text{Atom}(M, A, c, 22, -), \text{Bond}(M, A, 2)$$

We denote the first clause by b_1 and the second one by b_2 . The vectors on the left-hand side of the clauses specify $p(b_i\theta = \text{true} | q\theta = \text{true})$ and $p(b_i\theta = \text{true} | q\theta = \text{false})$ respectively. The covers relation (assuming the Naive Bayes assumption) assigns probability 0.97 to example 225 because both features fail for $\theta = \{M/225\}$. Hence

$$\begin{aligned} p(\text{Mutagenetic}(225) = \text{true}, b_1\theta = \text{false}, b_2\theta = \text{false}) &= p(b_1\theta = \text{false} | \text{Mutagenetic}(225) = \text{true}) \\ &\times p(b_2\theta = \text{false} | \text{Mutagenetic}(225) = \text{true}) \\ &\times p(\text{Mutagenetic}(225) = \text{true}) \\ &= 0.99 \times 0.62 \times 0.31 \approx 0.19 \end{aligned}$$

and

$$p(\text{Mutagenetic}(225) = \text{false}, b_1\theta = \text{false}, b_2\theta = \text{false}) = 0.79 \times 0.01 \times 0.68 \approx 0.005$$

This yields

$$P(\text{Mutagenetic}(225) = \text{true} | b_1\theta = \text{false}, b_2\theta = \text{false}) = \frac{0.19}{0.19 + 0.005} \approx 0.97$$

8.5.1 Structure and Parameter Learning

Next, we outline typical solutions to the problems of structure and parameter learning in the setting of probabilistic learning from entailment (as well as interpretation).

1. The problem of parameter estimation is concerned with estimating the values of the parameters λ of a fixed probabilistic program $H = (L, \lambda)$ that best explains the examples \mathcal{E} . So, λ is a set of parameters and can be represented as a vector. To measure the extent to which a model fits the data, one usually employs the likelihood of the data, i.e. $P(\mathcal{E}|L, \lambda)$, though other scores or variants could be used as well.

When all examples are fully observable, maximum likelihood reduces to frequency counting. In the presence of missing data, however, the maximum likelihood estimate typically cannot be written in closed form. It is a numerical optimization problem, and all known algorithms involve nonlinear optimization. The most commonly adapted technique for probabilistic

logic learning is the Expectation-Maximization (EM) algorithm [Dempster et al., 1977, McLachlan and Krishnan, 1997]. EM is based on the observation that learning would be easy (*i.e.*, correspond to frequency counting), if the values of all the random variables would be known. Therefore, it estimates these values, maximizes the likelihood based on the estimates, and then iterates. More specifically, EM assumes that the parameters have been initialized (*e.g.*, at random) and then iteratively performs the following two steps until convergence:

- (a) **(E-Step):** On the basis of the observed data and the present parameters of the model, it computes a distribution over all possible completions of each partially observed data case.
 - (b) **(M-Step):** Treating each completion as a fully observed data case weighted by its probability, it computes the improved parameter values using (weighted) frequency counting. The frequencies over the completions are called the *expected counts*.
2. What if we need to learn both the structure L and the parameters λ of the probabilistic program $H = (L, \lambda)$ from the data? Often, further information is given as well. As in ILP, the additional knowledge can take various different forms, including a language bias that imposes restrictions on the syntax of L , and an initial hypothesis (L, λ) from which the learning process can start.

Nearly all (score-based) approaches to structure learning perform a heuristic search through the space of possible hypotheses. Typically, hill-climbing or beam-search is applied until the hypothesis satisfies the logical constraints and the $score(H, \mathcal{E})$ is no longer improving. The steps in the search-space are typically made using refinement operators, see Section 2.7.

Logical constraints often require that the possible examples are covered in the logical sense. For instance, when learning stochastic logic programs from entailment, the possible example clauses must be entailed by the logic program.

8.5.2 Example: Extending FOIL

Building on [Landwehr et al., 2005], we will next illustrate a promising, alternative approach with less computational complexity, which adapts FOIL [Quinlan and Cameron-Jones, 1995] with the conditional likelihood as described in Equation (8.1) as the scoring function $score(L, \lambda, \mathcal{E})$. This idea has been followed with nFOIL, see [Landwehr et al., 2005] for more details.

Given a training set \mathcal{E} containing positive and negative examples (*i.e.* true and false ground facts), this algorithm stays in the learning from possible examples only to induce a probabilistic logical model to distinguish between the positive and negative examples. It computes Horn clause features b_1, b_2, \dots in an outer loop. It terminates when no further improvements in the score are obtained, *i.e.*, when $score(\{b_1, \dots, b_i\}, \lambda_i, \mathcal{E}) < score(\{b_1, \dots, b_{i+1}\}, \lambda_{i+1}, \mathcal{E})$, where

λ denotes the maximum likelihood parameters. A major difference with FOIL is, however, that the covered positive examples are not removed. The inner loop is concerned with inducing the next feature b_{i+1} top-down, *i.e.*, from general to specific. To this aim it starts with a clause with an empty body, *e.g.*, $Mutagenic(M)$. This clause is then specialized by repeatedly adding atoms to the body, *e.g.*, $Mutagenic(M) \leftarrow Bond(M, A, 1)$, $Muta(M) \leftarrow Bond(M, A, 1)$, $Atom(M, A, c, 22, -)$, *etc.* For each refinement b'_{i+1} we then compute the maximum-likelihood parameters λ'_{i+1} and $score(\{b_1, \dots, b'_{i+1}\}, \mathcal{E})$. The refinement that scores best, say b_{i+1} , is then considered for further refinement and the refinement process terminates when

$$score(\{b_1, \dots, b_{i+1}\}, \lambda_{i+1}, \mathcal{E}) < score(\{b_1, \dots, b''_{i+1}\}, \lambda''_{i+1}, \mathcal{E})$$

It has been experimentally found that nFOIL performs well compared to other ILP systems on traditional ILP benchmark data sets. mFOIL and Aleph, two standard ILP systems, are never significantly better than nFOIL (paired sampled t-test, $p = 0.05$). nFOIL achieves significantly higher predictive accuracies than mFOIL on Alzheimer amine, toxic, and acetyl. Compared to Aleph, nFOIL achieves significantly higher accuracies on Alzheimer amine and acetyl (paired sampled t-test, $p = 0.05$). For more details, we refer to [Landwehr et al., 2005].

8.5.3 Example: Stochastic Logic Programming

Cussens [2001] and Sato and Kameya [2001], solve the parameter estimation problem for stochastic logic programs respectively PRISM programs, and Mugleton [2000a, 2002] presents an approach to structure learning of stochastic logic programs: adding one clause at a time to an existing stochastic logic program. The essentially use equation (8.4) as covers relation and, hence, employ the learning from entailment setting while making use of the semantics of probabilistic proofs. Here, the examples are ground atoms entailed by the target stochastic logic program. However, the Naive Bayes framework studied prior to this example has a much lower computational complexity. Also, learning stochastic logic programs from atoms only is much harder than learning them from proofs because atoms carry much less information than proofs.

8.6 Learning from Interpretations

The learning from interpretations setting [De Raedt and Dzeroski, 1994] upgrades boolean concept-learning in computational learning theory [Valiant, 1984]. When learning from interpretations, the examples are Herbrand interpretations³ and a hypothesis H covers an example e with respect to the background theory \mathcal{B} if and only if $\mathcal{B} \cup e$ is a model of H .

As an example, consider the interpretation I , which is the union of B and e :

³Recall that Herbrand interpretations are sets of true ground facts and they completely describe a possible situation.

$$\begin{aligned}
B &= \{Father(henry, bill), Father(alan, betsy), Father(alan, benny), Father(brian, bonnie), Father(bill, carl), \\
e &= \{Carrier(alan), Carrier(ann), Carrier(betsy).\}
\end{aligned}$$

The interpretation I is covered by the clause H :

$$Carrier(X) \quad :- \quad Mother(M, X), Carrier(M), Father(F, X), Carrier(F).$$

because I is a model of C , *i.e.*, for all substitutions θ such that $body(C)\theta \subseteq I$, it holds that $head(C)\theta \in I$.

The key difference between learning from interpretations and learning from entailment is that interpretations carry much more, even complete information. Indeed, when learning from entailment, an example can consist of a single fact, whereas when learning from interpretations, all facts that hold in the example are known. Therefore, learning from interpretations is typically easier and computationally more tractable than learning from entailment, *c.f.*, [De Raedt, 1997].

ILP systems that learn from interpretations work in a similar fashion as those that learn from entailment. There is, however, one crucial difference and it concerns the generality relationship: When learning from entailment, G is more general than S if and only if $G \models S$, whereas when learning from interpretations, when $S \models G$. Another difference is that learning from interpretations is well suited for learning from positive examples only. For this case, a complete search of the space ordered by θ -subsumption is performed until all clauses cover all examples [De Raedt and Dehaspe, 1997].

8.7 Learning from Probabilistic Interpretations

The large majority of statistical relational learning techniques proposed so far fall into the learning from interpretations setting including parameter estimation of probabilistic logic programs [Koller and Pfeffer, 1997], learning of probabilistic relational models [Getoor et al., 2002], parameter estimation of relational Markov models [Taskar et al. 2002], learning of object-oriented Bayesian networks [Bangs et al., 2001], learning relational dependency networks [Neville and Jensen, 2004], learning logic programs with annotated disjunctions [Vennekens et al., 2004, Riguzzi, 2004] and most recently, learning Bayesian logic programs [Kersting et. al. 2000,2005].

In order to integrate probabilities in the learning from interpretations setting, we need to find a way to assign probabilities to interpretations covered by an annotated logic program. In the past few years, this issue has received a lot of attention and various different approaches have been developed such as probabilistic-logic programs [Ngo and Haddawy, 1997], probabilistic relational models [Pfeffer, 2000], relational Bayesian networks Jager [1997], and Bayesian logic programs [Kersting, 2000, Kersting and De Raedt, 2001b]. Here, we focus on two methods:

1. *Undirected Graphical Model*: Domingos and Richardson [2004] Markov logic networks (MLNs) and
2. *Directed Graphical Model*: Kersting and De Raedt's [2001b] Bayesian logic programs.

8.7.1 Example: Markov Logic Networks

Markov logic networks combine Markov networks [Pearl, 1991], which represent probability distributions over propositional interpretations, with first order logic. The idea underlying Markov logic networks is to view logical formulas as soft constraints on the set of possible worlds, *i.e.*, as *soft constraints on interpretations*: if a world violates one formula, it is less probable but not necessarily impossible as in classical logic. The fewer formulas a world violates, the more probable it is. In a Markov logic network, this is realized by associating a weight with each formula that reflects how strong the constraint is. More precisely, a Markov logic network consists of weighted first-order predicate logic formulae $\Sigma = \{C_1, C_2, \dots, C_m\}$. The weights w_C of a formula C specify a bias for *ground instances to be true in a logical model*.

Consider the following example taken from [Richardson and Domingos, 2005]. Friends-smokers is a small Markov logic network that calculates the probability of a person P having lung cancer $Ca(P)$ based whether or not a person or her friends $Fr(P, P')$ smokes $Sm(P)$ respectively $Sm(P')$. This can be encoded using the following Markov logic formulas:

$$\begin{aligned} 1.5 : \quad \forall X : \quad & Sm(X) \Rightarrow Ca(X) \\ 1.1 : \quad \forall X, Y : \quad & Fr(X, Y) \Rightarrow (Sm(X) \iff Sm(Y)) \end{aligned}$$

For a given finite domain (roughly speaking a finite set of constants) $D = \{d_1, d_2, \dots, d_n\}$, the Markov logic network defines a probability distribution over interpretations I over domain D and the relations occurring in the Markov logic network via

$$P(I|H, \mathcal{B}) = \frac{1}{Z(I)} \prod_{C \in H \cup \mathcal{B}} e^{n_C(I) \cdot w_C} = \frac{1}{Z(I)} \prod_{C \in H \cup \mathcal{B}} \phi_C(I)^{n_C(I)} \quad (8.2)$$

where $n_C(I)$ is the number of true groundings of C in I , $\phi_C(I) = e^{w_C}$, and \mathcal{B} is a possible background theory.

Markov logic networks can be viewed as proving templates for constructing Markov networks. Given a set D constants,

- the nodes correspond to the ground atoms in the Herbrand base of the corresponding set of formulas Σ
- and there is an edge between two nodes if and only if the corresponding ground atoms appear together in at least one grounding of one formula $C_i \in \Sigma$.

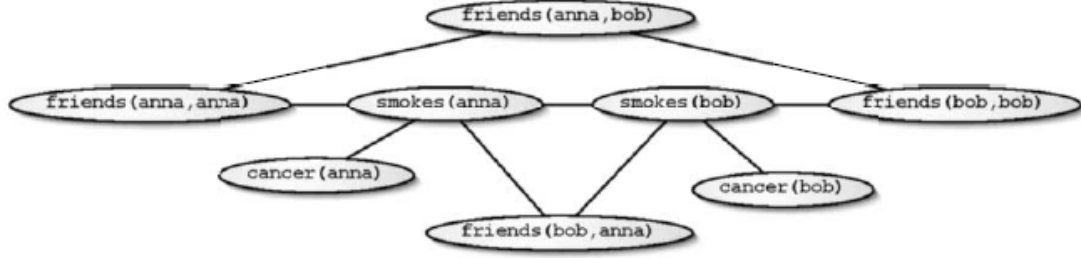


Figure 8.5: The Markov network induced by the friends-smoker Markov logic network assuming *anna* and *bob* as constants.

Assuming *anna* and *bob* as constants, the friends-smoker Markov logic network induces the Markov network in Figure 8.5.

Note that given different sets of constants, the Markov logic network will produce different Markov networks. From Equation (8.2), we can see that an example e consists of

- a *logical part*, which is a Herbrand interpretation of the annotated logic program, and
- a *probabilistic part*, which is a partial state assignment of the random variables occurring in the logical part.

To see this, consider that a possible example I in the friends-smokers domain is

$$\begin{aligned} Friends(anna, bob) &= true, & Friends(bob, anna) &= true, & Friends(anna, anna) &= ?, \\ Friends(bob, bob) &= true, & Smokes(anna) &= false, & Smokes(bob) &= ?, \\ Cancer(anna) &= ?, & Cancer(bob) &= false \end{aligned}$$

where ? denotes an unobserved state. The covers relation for e can now be computed using any Markov network inference engine based on Equation (8.2).

Nearly all (score-based) approaches to structure learning perform a heuristic search through the space of possible hypotheses. Typically, hill-climbing or beam-search is applied until the hypothesis satisfies the logical constraints and the $score(H, \mathcal{E})$ is no longer improving. The steps in the search-space are typically made using refinement operators, see Section 2.7.

As described in Section 8.5, logical constraints often require that the possible examples are covered in the logical sense. For instance, when learning Markov logic networks, the possible interpretations must be models of the underlying logic program. Thus, for a probabilistic program $H = (\mathcal{L}_H, \lambda_H)$ and a background theory $\mathcal{B} = (\mathcal{B}, \lambda_B)$ it holds that $\forall e_p \in \mathcal{E}_p : p(e|H, \mathcal{B}) > 0$ if and only if $covers(e, \mathcal{L}_H, \mathcal{B}) = 1$, where \mathcal{L}_H (respectively \mathcal{B}) is the underlying logic

program (logical background theory) and $\text{covers}(e, \mathcal{L}_H, \mathcal{B})$ is the purely logical covers relation, which is either 0 or 1.

Kok and Domingos [2005] proposed a beam-search based approach for learning clausal Markov logic networks from possible examples only. Recall that a clausal Markov logic program consists of weighted clauses, *i.e.*, disjunction of literals. The clauses without associated weights constitute a clausal program L , and the weights the parameters λ . Starting with some initial clausal Markov logic network $H = (L, \lambda)$, the parameters maximizing $\text{score}(L, \lambda, \mathcal{E})$ are computed. Then, refinement operators (Section 2.7) generalizing respectively specializing L are used to compute all neighbours of L in the hypothesis space. Literals are added and deleted, and signs of literals are flipped. Each neighbour is scored, yielding new hypotheses (L', λ') . To speed-up scoring, Kok and Domingos employ a variant of the pseudo-log-likelihood

$$\sum_{i=1}^n \log p(X_i = x_i | MB_x(X_i))$$

where x is the Herbrand base, x_i is the i^{th} ground atoms truth value, and $MB_x(X_i)$ is the state of X_i 's Markov blanket in the data, where, the Markov blanket of a node is its set of neighbouring nodes. The b best ones with $\text{score}(L', \lambda', \mathcal{E}) > \text{score}(L, \lambda, \mathcal{E})$ are kept. On these b best ones, the refining and scoring process is iteratively applied again until no new clauses improve the score or a maximal number of literals is reached. The clause with highest score in all iterations is added to H , and the process is continued until no improvement in score of the current best hypothesis is obtained.

8.7.2 Example: Bayesian Logic Programs

Recall from Chapter 7 that a Bayesian network specifies a joint probability distribution over a finite set of random variables and consists of two components: (1) a qualitative or logical one that encodes the local influences among the random variables using a directed acyclic graph, and (2) a quantitative one that encodes the probability densities over these local influences.

Despite these interesting properties, Bayesian networks also have a major limitation: they are essentially propositional representations.

For an illustrative example, imagine the task of modeling the localization of genes/proteins. When using a Bayesian network, every gene is a single random variable. There is no way of formulating general probabilistic regularities among the localizations of the genes such as the localization L of gene G is influenced by the localization L' of another gene G' that interacts with G . The propositional nature and limitation of Bayesian networks are similar to those of traditional attribute-value learning techniques, which have motivated work on upgrading these techniques within ILP. This in turn also explains the interest in upgrading Bayesian networks towards using first order logical representations.

Bayesian logic programs (BLPs) unify Bayesian networks with (definite clause) logic programming, which allows one to overcome the propositional character of Bayesian networks and the purely logical nature of logic programs. From a

knowledge representation point of view, BLPs can be distinguished from alternative frameworks by having both

1. *Logic programs* (i.e. definite clause programs, which are sometimes called pure Prolog programs) as well as
2. it Bayesian networks as an immediate special case.

This is realized through the use of a small but powerful set of primitives. Indeed, similar to that for markov logic networks, the underlying idea of BLPs is to establish a one-to-one mapping between

1. ground atoms in the least Herbrand model ($MM(\Sigma)$ and random variables, and
2. between the immediate consequence operator and the direct influence relation.

Therefore, BLPs can also handle domains involving structured terms as well as continuous random variables. We will briefly describe BLPs, their representation language, their semantics, and a query-answering process, and present the learning of BLPs from data based on [Kersting 200, Kersting and De Raedt 2005].

The least Herbrand model of a BLP together with its direct influence relation is viewed as a (possibly infinite) Bayesian network. BLPs inherit the advantages of both Bayesian networks and definite clause logic, including

1. the strict separation of qualitative and quantitative aspects and consequently,
2. the introduction of a graphical representation, which stays close to the graphical representation of Bayesian networks.

Indeed, BLPs can naturally model any type of Bayesian network (including those involving continuous variables) as well as any type of ‘pure’ Prolog program (including those involving functors). In fact, BLPs can model hidden Markov models and stochastic grammars, and can also be related to other first order extensions of Bayesian networks.

The framework for learning BLPs is an instance of the probabilistic learning from interpretations setting as described in Section 2.4.3. It is unifying as it combines traditional Bayesian network learning and ILP principles. Therefore, the BLP framework builds upon

- Many of the results for Bayesian network learning from the Uncertainty in AI community, see e.g. [Heckerman, 1995],
- Many of the results from the ILP community,
- The EM and gradient ascent algorithms for parameter estimation,

- The general structure learning mechanisms from the field of Bayesian networks, and
- The clausal discovery and learning from interpretations settings from ILP for probabilistic learning from interpretations.

Consider the family disease example from page 393. Now, imagine another totally separated family, which could be described by a similar Bayesian network. The graphical structure and associated conditional probability distribution for the two families are controlled by the same intensional regularities. But these overall regularities cannot be captured by a traditional Bayesian network. So, we need BLPs to represent these overall regularities. The central notion of BLPs is that of a Bayesian clause.

Definition 68 *A Bayesian (definite) clause C is an expression of the form $A|A_1, \dots, A_n$ where $n \geq 0$, the A, A_1, \dots, A_n are Bayesian atoms and all Bayesian atoms are (implicitly) universally quantified. When $n = 0$, C is called a Bayesian fact and expressed as A . The differences between a Bayesian clause and a logical clause are:*

1. *the atoms $P(t_1, \dots, t_k)$ and predicates P/k are Bayesian, which means that they have an associated finite set $S(P/k)$ of possible states (the ideas easily generalize to discrete and continuous random variables, modulo the well-known restrictions for Bayesian networks)*
2. *' \leftarrow ' is used instead of ' \vdash ' to highlight the conditional probability distribution.*

For instance, consider the Bayesian clause $C : BT(X) \mid MC(X), PC(X)$ where $S(BT/1) = \{a, b, ab, 0\}$ and $S(MC/1) = S(PC/1) = \{a, b, 0\}$. Intuitively, a Bayesian predicate P/k generically represents a set of random variables. More precisely, each Bayesian ground atom \mathbf{g} over P/k represents a random variable over the states $S(\mathbf{g}) = S(P/k)$. For example, $BT(ann)$ represents the blood type of a person named *Ann* as a random variable over the states $\{a, b, ab, 0\}$. Apart from that, most logical notions carry over to BLPs. So, we can speak of Bayesian predicates, terms, constants, substitutions, propositions, ground Bayesian clauses, Bayesian Herbrand interpretations *etc.* For the sake of simplicity we will sometimes omit the term Bayesian as long as no ambiguities arise. We will assume that all Bayesian clauses are range-restricted, i.e., $Var(head(C)) \subseteq Var(body(C))$. Range restriction is often imposed in the database literature; it allows one to avoid the derivation of non-ground true facts. As already indicated while discussing Figure 7.2, a set of Bayesian clauses encodes the qualitative or structural component of the BLPs. More precisely, ground atoms correspond to random variables, and the set of random variables encoded by a particular BLP corresponds to its *least Herbrand domain*. In addition, the direct influence relation corresponds to the immediate consequence.

Consider again, the following

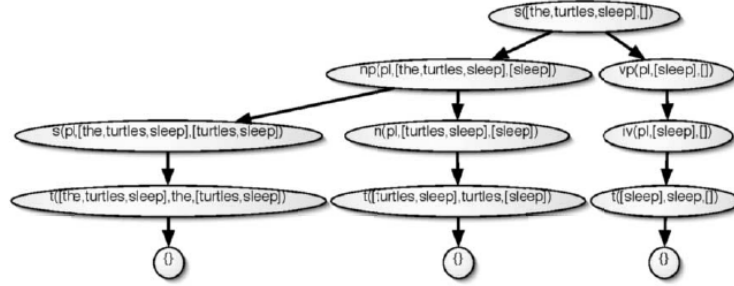


Figure 8.6: An example proof tree (which is covered by the definite clause grammar on page 8.8).

We assume that for each Bayesian predicate P/k there is a corresponding combining rule $cr(P/k)$, such as noisy or (see e.g. [Jensen, 2001]) or average. The latter assumes $n_1 = \dots = n_m$ and $S(A_{i_j}) = S(A_{k_j})$, and computes the average of the distributions over $S(A)$ for each joint state over $j \in S(A_{i_j})$.

8.8 Learning from Proofs

Because learning from entailment (with ground facts as examples) and interpretations occupy extreme positions with respect to the information the examples carry, it is interesting to investigate intermediate positions. Ehud Shapiro [1983] Model Inference System (MIS) fits nicely within the learning from entailment setting where examples are facts. However, to deal with missing information, Shapiro employs a clever strategy: MIS queries the users for missing information by asking them for the truth-value of facts. The answers to these queries allow MIS to reconstruct the trace or the proof of the positive examples.

When learning from proofs, the examples are ground proof-trees and an example e is covered by a hypothesis H with respect to the background theory \mathcal{B} if and only if e is a proof-tree for $H \cup \mathcal{B}$. At this point, there exist various possible forms of proof-trees. Here, we will, assume that the proof-tree is given in the form of a ground and-tree where the nodes contain ground atoms. More formally, a tree t is a proof-tree for a logic program Σ if and only if t is a rooted tree where for every node $n \in t$ with $children(n)$ satisfies the property that there exists a substitution θ and a clause $C \in \Sigma$ such that $n = head(C)\theta$. and $children(n) = body(C)$.

Consider the following definite clause grammar.

$Sentence(A, B)$: -	$Noun_phrase(C, A, D), Verb_phrase(C, D, B).$
$Noun_phrase(A, B, C)$: -	$Article(A, B, D), Noun(A, D, C).$
$Verb_phrase(A, B, C)$: -	$Intransitive_verb(A, B, C).$
$Article(singular, A, B)$: -	$Terminal(A, a, B).$
$Article(singular, A, B)$: -	$Terminal(A, the, B).$
$Article(plural, A, B)$: -	$Terminal(A, the, B).$
$Noun(singular, A, B)$: -	$Terminal(A, turtle, B).$
$Noun(plural, A, B)$: -	$Terminal(A, turtles, B).$
$Intransitive_verb(singular, A, B)$: -	$Terminal(A, sleeps, B).$
$Intransitive_verb(plural, A, B)$: -	$Terminal(A, sleep, B).$
$Terminal([A B], A, B).$		

It covers the proof tree shown in Figure 8.6. Proof-trees contain, as interpretations, a lot of information. Indeed, they contain instances of the clauses that were used in the proofs. Therefore, it may be hard for the user to provide these type of examples. Even though this is generally true, there exist specific situations for which this is feasible. Indeed, consider tree banks such as the UPenn Wall Street Journal corpus [Marcus et al., 1994], which contain parse trees. These trees directly correspond to the proof-trees we talk about. Another example is explanation-based learning (EBL) [Ellman, 1989, Mooney and Zelle, 1994]. It uses an existing domain theory to deductively explain an example (explanation step) in terms of a proof-tree and variablizes the explanation, *i.e.*, generalizes the proof as far as possible while maintaining its correctness (generalization step).

In the learning from proofs setting, one could turn all the proof-trees (corresponding to positive examples) into a set of ground clauses, which would constitute the initial theory. This theory can then be generalized by taking the least general generalization (under θ -subsumption) of pairwise clauses. Of course, care must be taken that the generalized theory does not cover negative examples.

8.9 Probabilistic Proofs

To define probabilities on proofs, ICL [Poole, 1993], PRISMs [Sato, 1995, Sato and Kameya, 2001], and stochastic logic programs [Eisele, 1994, Muggleton, 1996, Cussens, 2001] attach probabilities to facts (respectively clauses) and treat them as stochastic choices within resolution. Relational Markov models [Anderson et al., 2002] and logical hidden Markov models [Kersting 2000], can be viewed as a simple fragment of them, where heads and bodies of clauses are singletons only, so-called iterative clauses. We will illustrate probabilistic learning from proofs using stochastic logic programs. For a discussion of the close relationship among stochastic logic programs, ICL, and PRISM, we refer to [Cussens, 2005].

Stochastic logic programs are inspired on stochastic context free grammars [Abney, 1997, Manning and Schutze, 1999]. The analogy between context free grammars and logic programs is that

1. grammar rules correspond to definite clauses,
2. sentences (or strings) correspond to atoms, and
3. productions correspond to derivations.

Furthermore, in stochastic context-free grammars, the rules are annotated with probability labels in such a way that the sum of the probabilities associated to the rules defining a non-terminal is 1.0 (see this relation with the closed world assumption discussed in Section 1.3.7).

Eisele and Muggleton have exploited this analogy to define stochastic logic programs. These are essentially definite clause programs, where each clause C has an associated probability label p_C such that the sum of the probabilities associated to the rules defining any predicate is 1.0 (though Cussens [1999] considered less restricted versions as well). This framework allows ones to assign probabilities to proofs for a given predicate q given a stochastic logic program $H \cup B$ in the following manner. Let D_q denote the set of all possible ground proofs for atoms over the predicate q . For simplicity reasons, it will be assumed that there is a finite number of such proofs and that all proofs are finite (but again see [Cussens, 1999] for the more general case). Now associate to each proof $t_q \in D_q$ the probability

$$v_t = \prod_C p_C^{n_{C,t}}$$

where the product ranges over all clauses C and $n_{C,t}$ denotes the number of times clause C has been used in the proof t_q . For stochastic context free grammars, the values v_t correspond to the probabilities of the production.

However, the **difference** between context free grammars and logic programs is that in grammars two rules of the form $n \leftarrow q, n_1, \dots, n_m$ and $q \leftarrow q_1, \dots, q_k$ always resolve to give $n \leftarrow q_1, \dots, q_k, n_1, \dots, n_m$, **whereas resolution may fail due to unification**. Therefore, the probability of a proof tree t in D_q , *i.e.*, a successful derivation is

$$p(t|H, \mathcal{B}) = \frac{v_t}{\sum_{s \in D_q} v_s} \quad (8.3)$$

The probability of a ground atom a is then defined as the sum of all the probabilities of all the proofs for that ground atom.

$$p(a|H, \mathcal{B}) = \sum_{\substack{s \in D_q \\ s \text{ is a proof for } a}} v_s \quad (8.4)$$

As an example, consider a stochastic variant of the definite clause grammar in the example on page 507 with uniform probability values for each predicate. The value v_u of the proof (tree) u in the example on the page 507 is $p(v_u) = \frac{1}{3} \times \frac{1}{2} \times \frac{1}{2} = \frac{1}{12}$. The only other ground proofs s_1, s_2 of atoms over the predicate *Sentence* are those of $Sentence([a, turtle, sleeps], [])$ and $Sentence([the, turtle, sleeps], [])$. Both get the value $p(v_{s_1}) = p(v_{s_2}) = 1$. Because there is only one proof for each of the sentences,

$$p(Sentence([the, turtles, sleep], [])) = v_u = \frac{1}{3}$$

For stochastic logic programs, there are at least two natural learning settings.

1. Motivated by Equation (8.3), we can learn them from proofs. This makes structure learning for stochastic logic programs relatively easy, because proofs carry a lot of information about the structure of the underlying stochastic logic program. Furthermore, the learning setting can be considered as an extension of the work on learning stochastic grammars from proof-banks. It should therefore also be applicable to learning unification based grammars. We will present a probabilistic ILP approach within the learning from proofs setting subsequently in this section.
2. On the other hand, we can use Equation (8.4) as covers relation and, hence, employ the learning from entailment setting as mentioned in Section 8.5.3. Here, the examples are ground atoms entailed by the target stochastic logic program. Learning stochastic logic programs from atoms only is much harder than learning them from proofs because atoms carry much less information than proofs.

8.9.1 Example: Extending GOLEM

Given a training set \mathcal{E} containing ground proofs as examples, one possible approach to learning from possible proofs only combines ideas from the early ILP system Golem [Muggleton and Feng, 1992] that employs Plotkins [1970] least general generalization(LGG) with bottom-up generalization of grammars and hidden Markov models [Stolcke and Omohundro, 1993]. The resulting algorithm employs the likelihood of the proofs $score(L, \lambda, \mathcal{E})$ as the scoring function. It starts by taking as L_0 the set of ground clauses that have been used in the proofs in the training set and scores it to obtain λ_0 . After initialization, the algorithm will then repeatedly select a pair of clauses in L_i , and replace the pair by their LGG (lub) to yield a candidate L' . The candidate that scores best is then taken as $H_{i+1} = (L_{i+1}, \lambda_{i+1})$, and the process iterates until the score no longer improves.

One interesting issue is that strong logical constraints can be imposed on the LGG. These logical constraints directly follow from the fact that the example proofs should still be valid proofs for the logical component L of all hypotheses considered. Therefore, it makes sense to apply the LGG only to clauses that define the same predicate, that contain the same predicates, and whose (reduced)

LGG also has the same length as the original clauses. In general, the length (number of literals) of the LGG of m (ground) clauses of length at most n is n^m , see [Muggleton and Feng, 1992].

As an example, consider the following target stochastic logic program:

$$\begin{aligned}
1 : & \quad S(A, B) \leftarrow NP(\text{Number}, A, C), VP(\text{Number}, C, B). \\
\frac{1}{2} : & \quad NP(\text{Number}, A, B) \leftarrow \text{det}(A, C), n(\text{Number}, C, B). \\
\frac{1}{2} : & \quad NP(\text{Number}, A, B) \leftarrow \text{pronom}(\text{Number}, A, B). \\
\frac{1}{2} : & \quad VP(\text{Number}, A, B) \leftarrow V(\text{Number}, A, B). \\
\frac{1}{2} : & \quad VP(\text{Number}, A, B) \leftarrow V(\text{Number}, A, C), NP(D, C, B). \\
1 : & \quad \text{Det}(A, B) \leftarrow \text{Term}(A, \text{the}, B). \\
\frac{1}{4} : & \quad N(s, A, B) \leftarrow \text{Term}(A, \text{man}, B). \\
\frac{1}{4} : & \quad N(s, A, B) \leftarrow \text{Term}(A, \text{apple}, B). \\
\frac{1}{4} : & \quad N(pl, A, B) \leftarrow \text{Term}(A, \text{men}, B). \\
\frac{1}{4} : & \quad N(pl, A, B) \leftarrow \text{Term}(A, \text{apples}, B). \\
\frac{1}{4} : & \quad V(s, A, B) \leftarrow \text{Term}(A, \text{eats}, B). \\
\frac{1}{4} : & \quad V(s, A, B) \leftarrow \text{Term}(A, \text{sings}, B). \\
\frac{1}{4} : & \quad V(pl, A, B) \leftarrow \text{Term}(A, \text{eat}, B). \\
\frac{1}{4} : & \quad V(pl, A, B) \leftarrow \text{Term}(A, \text{sing}, B). \\
1 : & \quad \text{Pronom}(pl, A, B) \leftarrow \text{Term}(A, \text{you}, B). \\
1 : & \quad \text{Term}([A|B], A, B).
\end{aligned}$$

From this program, (independent) training sets of 50, 100, 200, and 500 proofs were generated. For each training set, 4 different random initial sets of parameters were tried. The learning algorithm was run on each data set starting from each of the initial sets of parameters. The algorithm stopped when a limit of 200 iterations was exceeded or a change in log-likelihood between two successive iterations was smaller than 0.0001. In all runs, the original structure was induced from the proof-trees. Moreover, already 50 proof-trees suffice to rediscover the structure of the original stochastic logic program. Further experiments with 20 and 10 samples respectively show that even 20 samples suffice to learn the given structure. Sampling 10 proofs, the original structure is rediscovered in one of five experiments. This supports that the learning from proof trees setting carries a lot information. Furthermore, the method scales well. Runs on two independently sampled sets of 1000 training proofs yield similar results: the original structure was learned in both cases. More details can be found in [De Raedt et al., 2005].

Other statistical relational learning frameworks that have been developed within the learning from proofs setting are relational Markov models [Anderson et al., 2002] and logical hidden Markov models [Kersting 2000?].

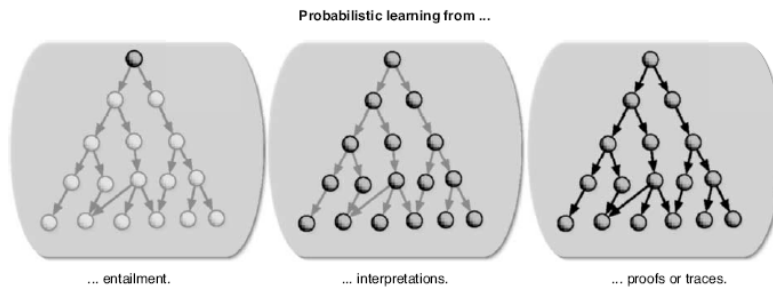


Figure 8.7: The level of information on the target probabilistic program provided by probabilistic ILP settings: shaded parts denote unobserved information. Learning from entailment provides the least information. Only roots of proof tree are observed. In contrast, learning from proofs or traces provides the most information. All ground clauses and atoms used in proofs are observed. Learning from interpretations provides an intermediate level of information. All ground atoms but not the clauses are observed.

8.10 Summary

The three learning settings are graphically compared in Figure 8.7.

References

- [CGH97] E. Castillo, J. M. Gutierrez, and A. S. Hadi. *Expert systems and probabilistic network models*. Springer-Verlag, 1997.
- [CKT91] Peter Cheeseman, Bob Kanefsky, and William M. Taylor. Where the Really Hard Problems Are. In *Proceedings of the Twelfth International Joint Conference on Artificial Intelligence, IJCAI-91, Sidney, Australia*, pages 331–337, 1991.
- [CM97] Stephen A. Cook and David G Mitchell. Finding hard instances of the satisfiability problem: A survey. In Du, Gu, and Pardalos, editors, *Satisfiability Problem: Theory and Applications*, volume 35, pages 1–17. American Mathematical Society, 1997.
- [Gol79] A. Goldberg. Average case complexity of the satisfiability problem. In *Proceedings of the 4th Workshop on Automated Deduction*, pages 1–6, Austin, Texas, USA, 1979.
- [Jen01] F. V. Jensen. *Bayesian Networks and Decision Graphs*. Springer, 2001.
- [Jor98] M. I. Jordan. *Learning in Graphical Models*. MIT Press, 1998.
- [Lau96] S. L. Lauritzen. *Graphical Models*. Clarendon Press, Oxford, 1996.
- [NCdW97] Shan Wei Nienhuys-Cheng and Ronald de Wolf. *Foundations of Inductive Logic Programming*. Springer Verlag, Germany, 1997.
- [Pea88] J. Pearl. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufmann, 1988.
- [SKC93] Bart Selman, Henry A. Kautz, and Bram Cohen. Local search strategies for satisfiability testing. In Michael Trick and David Stifler Johnson, editors, *Proceedings of the Second DIMACS Challenge on Cliques, Coloring, and Satisfiability*, Providence RI, 1993.
- [SS90] Glenn Shafer and Prakash P. Shenoy. Probability propagation. *Ann. Math. Artif. Intell.*, 2:327–351, 1990.