

i/p space

Model
Complex
relations

o/p space

Model
Complex
relations

Simple linear
kernels

Graphical model

Complex graph
or first order
kernels

First order prob
models

$$K(o_1, o_2) \in \mathbb{R}$$

2 classes: Study ways of generating interpretable "features" for o_1 & o_2 s.t. the clustering/classification/regression problem at hand can be solved effectively.

Unsupervised feature mining

Transaction ID	Items
1	{D,O,N,K,E,Y}
2	{M,O,N,K,E,Y}
3	{M,A,K,E}
4	{M,U,C,K,Y}
5	{C,O,O,K,I,E}

$$\Sigma = \{A, C, D, E, I, K, M, N, O, U, Y\}$$

Figure 1: Dataset \mathcal{D}

- Let us say you are provided the data set \mathcal{D} as in Figure 1. Propose an efficient algorithm to find all subsets of Σ that are subsets of more than m transactions (for some fixed m). (This is the idea behind the classic apriori algorithm).

eg: Find all $x \subseteq \Sigma$ s.t. x covers at least 2 transactions

$$\text{covers}(x, t) = 1 \quad \text{iff } x \subseteq t$$
$$\text{covers}(x, t) = 0 \quad \text{o/w}$$

$$\{C\} = \{4, 5\} \quad \{E\} = \{1, 2, 3, 5\} \quad \{M\} = \{2, 3, 4\}$$

$$\{C, E\} = \{5\}$$

$$\{E, M\} = \{2, 3\}$$

$$\{D\} = \{1\}$$

$$\{Y\} = \{1, 2, 4\}$$

Because $|D|=1$
 $\therefore |\{D, Y\}| \leq 1$

$$\{D, Y\} = \{1\}$$

S_i = set of all subsets of Σ of size i st each element of S_i covers at least m transactions

Algo

① Find all $s_i \in \Sigma$ st $|\{t \mid \text{covers}(s_i, t) = 1\}| \geq m$

S_1 = union of all such s_i

② Given S_1, S_2, \dots, S_{i-1}

how can I construct S_i ?

ⓐ $S_{i-k} \times S_k$

ⓑ $S_{i-1} \times S_{i-1}$

③ Find subset of S_i that is actually frequent & go back to

Different options for step ②

①: $S_i = S_{i-k} \times S_k$ for any $k \in [1, i-1]$
 such that if S_{i-k} & S_k are disjoint then $S_i = S_{i-k} \cup S_k$

Choice of k may depend on how small S_{i-k} or S_k can be

eg: $\{N, K, E, \gamma\} \Rightarrow \{1, 2\}$

$\{N, K\} \cup \{E, \gamma\}$

$\{N, E\} \cup \{K, \gamma\}$

$\{N, \gamma\} \cup \{K, E\}$

$S_2 \cup S_2'$ ($k=2$)
 $S \in S_2 \cap S_2' = \emptyset$

Selectivity (size) of S_2 is more

$\{N, K, E\} \cup \{\gamma\}$

$\{N, K, \gamma\} \cup \{E\}$

$\{K, E, \gamma\} \cup \{N\}$

$\{N, E, \gamma\} \cup \{K\}$

$S_3 \cup S_1$
 $S \in S_3 \cap S_1 = \emptyset$

Selectivity (size) of S_3 is small

S_3 than that of
but less than
 S_1

Antimonotonicity property:

If S_n covers t (denoted by $S_n \sqsubseteq t$) then $\forall k \leq n$, s.t.
 $S_k \subseteq S_n$, S_k covers t ($S_k \sqsubseteq t$)

The covers relation \sqsubseteq is a partial order, which means:

① Reflexivity: $S \sqsubseteq S$

② Antisymmetry: if $S \sqsubseteq t$ & $t \sqsubseteq S$
then $S = t$

③ Transitivity: if $S \sqsubseteq t$ & $t \sqsubseteq u$
then $S \sqsubseteq u$

⑥ Since selectivity of S_{i-1} is more than that of $S_{i-2} \dots$

(eg: $|S_3|=4$ $|S_2|=6$ $|S_1|=4$)

$$S_i = S_{i-1} \times S_{i-1}$$

algorithm for mining frequent patterns

Consider unioning S_{i-1} & S'_{i-1} s.t

$$|S_{i-1} \cap S'_{i-1}| = i-2$$

i.e. S_{i-1} & S'_{i-1} have exactly 1 element differing

$$\Rightarrow |S_{i-1} \cup S'_{i-1}| = i$$

eg: $S_{i-1} = \{N, K, E\}$ $S'_{i-1} = \{K, E, Y\}$
 $|S_{i-1} \cap S'_{i-1}| = |\{K, E\}| = 2 = i-2$
 $S_i = \{K, E, Y, N\}$

idea behind apriori

These 2 search paradigms span the entire spectrum of methods for constructing "interpretable" features

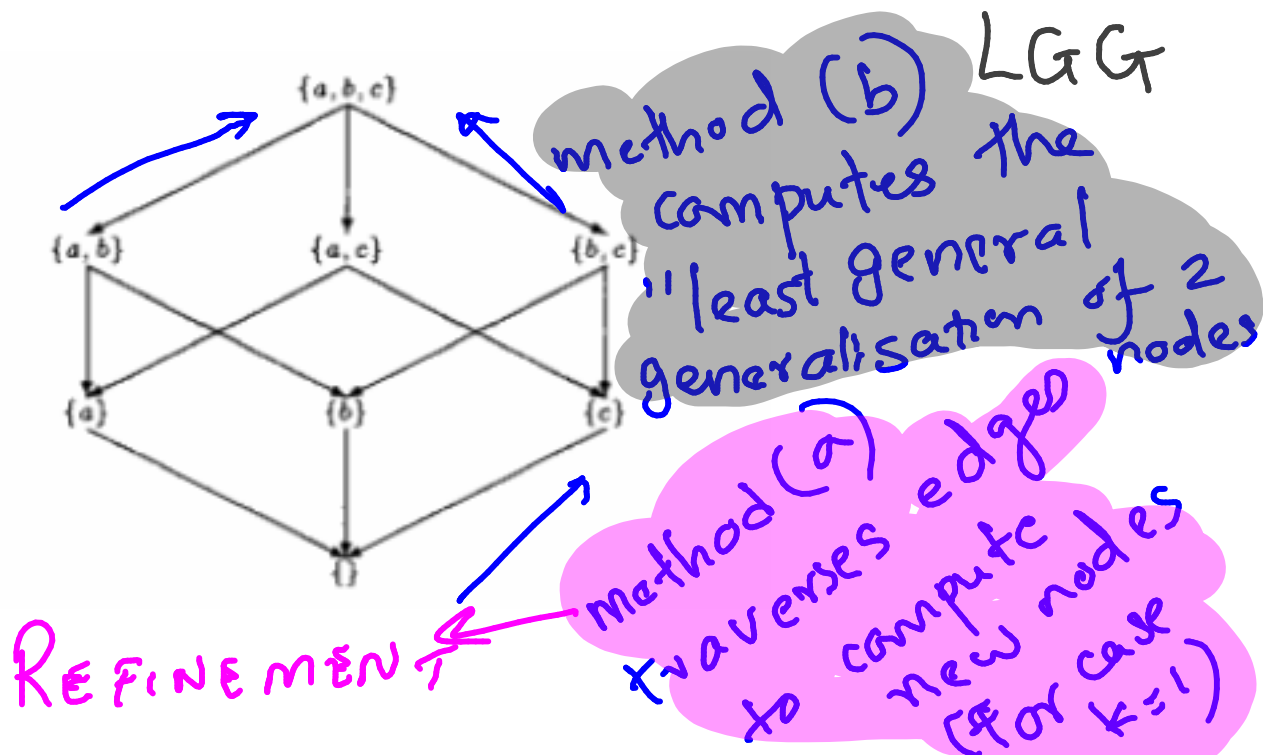


Figure 1.1: The lattice structure of $\langle S, \subseteq \rangle$, where S is the power set of $\{a, b, c\}$

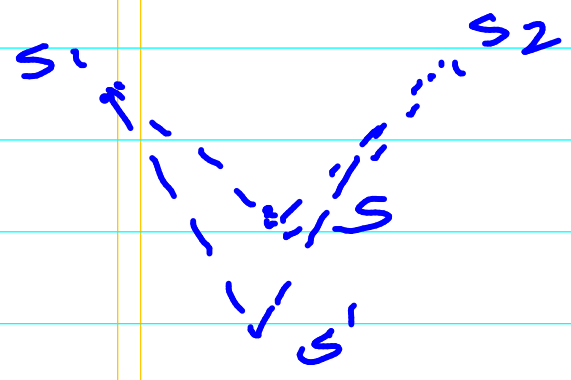
Given a p.o $\langle S, \subseteq \rangle$ or $\langle S, \preceq \rangle$
 the lgg of s_1 & s_2 is

$s_1 \sqcup s_2 = s$ iff $\textcircled{1} s_1 \preceq s$
 $s_2 \preceq s$
 $\textcircled{2}$ if $\exists s'$ s.t.
 $s_1 \preceq s'$ $s_2 \preceq s'$ then $s \preceq s'$

lgg = least upper bound (lub)
 $s \preceq s'$

Given a p.o $\langle S, \sqsubseteq \rangle$ or $\langle S, \leq \rangle$
 the glb (greatest lower bound) of s_1 & s_2 is

$$s_1 \sqcap s_2 = s \quad \text{iff} \quad \textcircled{1} \begin{matrix} s \leq s_1 \\ s \leq s_2 \end{matrix}$$



$\textcircled{2}$ If $\exists s'$ st $s' \leq s_1$ & $s' \leq s_2$ then $s' \leq s$

A Partial Order $\langle S, \sqsubseteq \rangle$ or $\langle S, \leq \rangle$ for which a lub and glb exist in S for any $s_1 \in S$ & $s_2 \in S$ is called a LATTICE

NOTE: In the table that follows, definitions of \leq could be reversed to make them definitions for \geq . Eg. if $s_1 \leq s_2$ then there exists \sqsupseteq st $s_1 \leq s_2$ iff $s_2 \sqsupseteq s_1$, st $s_1 \sqcap s_2 = s_1 \sqcup s_2$ & $s_1 \sqcup s_2 = s_1 \sqcap s_2$

Kernel

1) Set kernels
(each eg is bag of features)
say $\Sigma = \text{vocab of features}$

2) String kernels
(eg eg is a sequence of features)
say $\Sigma = \text{vocab of features}$

glb is longest common subsequence (LCS)

Corresponding S ,
a possible \leq and
whether a lattice

[$S = \text{implicit feature space}$]

1) $S = 2^\Sigma$

\leq is \subseteq

• Yes, as seen above, we have a lattice
• " \sqcup " ie lub is " \cup " (union)

• " \sqcap " ie glb is " \cap "

• Downward refinement = deleting element from set

• Upward refinement = adding element to set

2) $S = \bigcup_{k=1}^{\infty} \underbrace{\Sigma \times \Sigma \times \dots \times \Sigma}_{k \text{ times}}$
set of all subseq uences

$S_1 \leq S_2$ if characters in S_1 occur in S_2 in same order

• $S_1 \sqcup S_2 = S \rightarrow$ (verify this exists)

• $S_1 \sqcap S_2 = S$

• Downward refinement = deleting element from sequence

• Upward refinement = adding

Algos for computing lcs

3) Parse tree kernel
 (each eg is a subtree of the parse tree such that all productions are complete)

3) • S = subset of first order logic language (subset language restricted to that representing parse trees)
 • $\leq \equiv \theta$ subsumption in definite first order logic (FOL) (sections 1.4.6 & 2.1.1 of ↓)

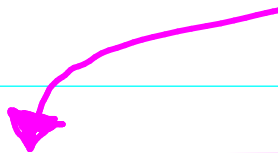
• \sqcap & \sqcup are defined in terms of glb & lub in definite FOL (sections 2.1.1 onwards of ↓)

4) Rational kernel

4) • S = All weighted finite state transducers

5) Graph kernel

4) (same as above, ie defined in terms of subset of FOL just as in the case of parse trees)
 Additionally, " \leq " could correspond to **SUBGRAPH ISOMORPHISM** discussed in earlier lectures



6) KLOG

5) S = definite clause FOL. Entire chapter 2 of ↓

For all above cases, existence of partial order implies the existence of a simple "antimonotonic quality criterion" (ie $Q(s) = \text{true}$ iff $S \subseteq \text{example}$)

Q is called an anti-monotonic quality criterion iff

$$Q(s) = \text{true} \Rightarrow Q(s') = \text{true}$$

for all $s' \leq s$

Consider the following Q 's

(A) $D =$

e_1^+	e_1^-
e_2^+	e_2^-
\vdots	\vdots
e_n^+	e_n^-

n_1 true & n_2 -ve eggs.
Let S be the space of "features" (like one of the above) & let

S could be any
- lattice with
some " \leq "

each $e_i^+ \in S$

& each $e_j^- \in S$

$$Q(s) = \text{true} \quad \text{iff} \quad |\{e_i^+ \mid s \leq e_i^+\}| \geq k$$

Is $Q(s)$ antimonotonic? & = false o/w

$$\{e_i^+ \mid s' \leq e_i^+\} \supseteq \{e_i^+ \mid s \leq e_i^+\}$$

$$\text{if } s' \leq s \text{ \& } s \leq e_i^+ \Rightarrow s' \leq e_i^+$$

(since \leq is
transitive)

\therefore if $s' \leq s$

$$|\{e_i^+ \mid s' \leq e_i^+\}| \geq |\{e_i^+ \mid s \leq e_i^+\}|$$

$\Rightarrow Q$ will be antimonotonic

(B) D is same as before

$$Q(s) = \text{true} \quad \text{iff}$$

$$= \text{false o/w}$$

$$|\{e_i^+ \mid s \leq e_i^+\}| - |\{e_j^- \mid s \leq e_j^-\}| \geq k$$

i.e. (pos-neg) loss should $\geq k$

Ans: Q is NOT antimonotonic

Proof: By counter example [H/w]

(C) D is same as before

$Q(S) = \text{True}$ iff $\text{max margin loss} \geq k$

Ans: Q is NOT antimonotonic

Proof: By _____ [H/w]

Q: For general quality criteria $Q(S)$ that are antimonotonic, you could use the frequent pattern mining algo (using subroutines (a) or (b) discussed earlier) to search for all the "good" patterns in $\langle S, \leq \rangle$

● But if $Q(s)$ is not antimonotonic, what can you do? → Data Mining problem

● Instead of asking for → Machine Learning problem

$$S^* = \{s \in S \mid Q(s) = \text{True}\}$$

if I asked for

$$S^* = \operatorname{argmax}_{S' \subseteq S} \text{objective}(S')$$

$$S' \subseteq S$$

pos-neg

eg:
likelihood
of data D
based on S'

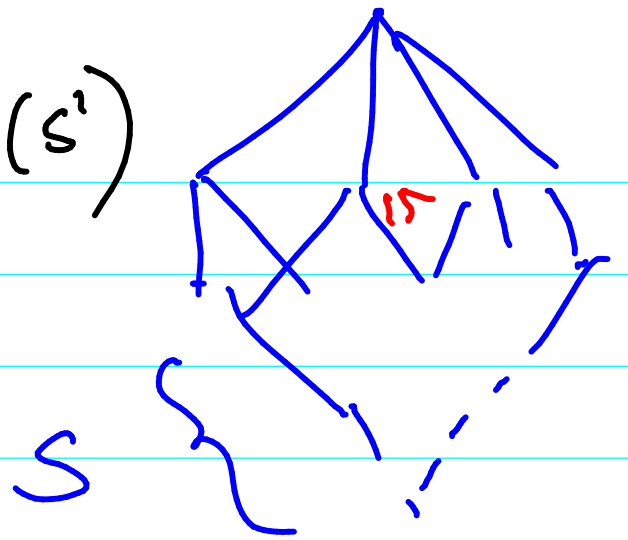
hinge loss
with regulariser
based on D

how do you search for S^* effectively?

The Machine Learning problem can be more interesting & generic:-

Approaches to solve the search problem (say machine learning setting)

Find $S^* = \underset{S \subseteq S}{\text{argmin}} \text{obj}(S')$



option ①

DFS (exhaustive search) ~~X~~

Traverse subsets of all nodes using DFS to find best

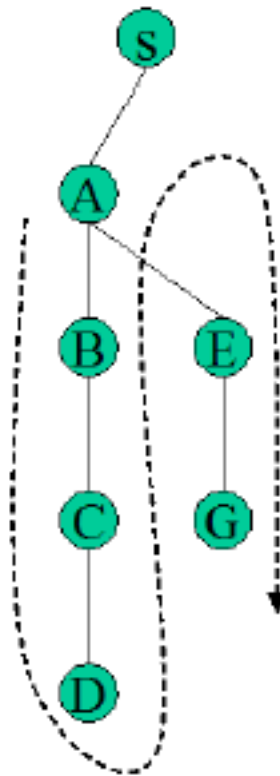
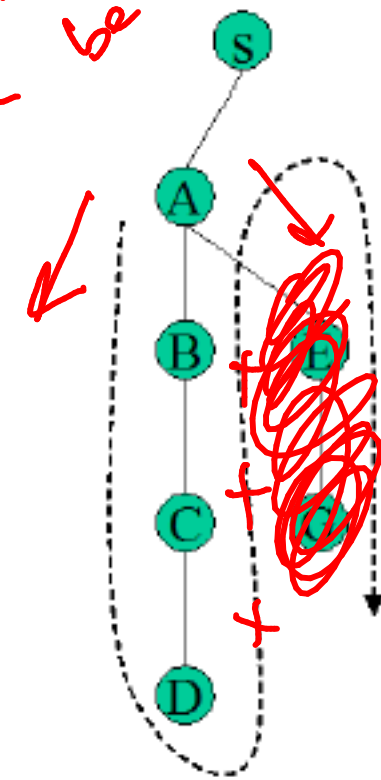


Figure 6.5: The DFS tree from S to G through the map in Figure 6.3.

Read about all these search algs in chapter 6 of

Option 2
 Greedy hill climbing = Greedy DFS
 (Greedy for DFS)

At A, find most "promising path" and choose between B & E



To choose next node at A

$$N^* = \operatorname{argmax}_{N} \operatorname{obj}(\{S, A, N\})$$

$N = B$ or
 $N = E$

i.e

Figure 6.5: The DFS tree from S to G through the map

$$N^* = \operatorname{argmax}_{N} \operatorname{obj} \left[\{S, A, B\}, \{S, A, E\} \right]$$

Examples of algos/approaches that use this idea

① ML Rules: Greedy search over simple conjunctions of features
 objective used = Likelihood of

naive bayes
model over the
set S'

(b) kFOLL: Greedy search over
space of first order
features

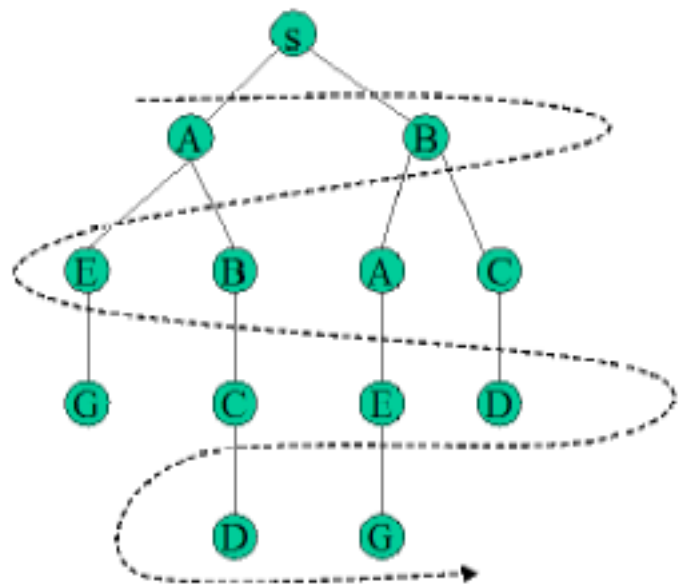
objective used = SVM objective /
SVM accuracy

(c) nFOLL: Like kFOLL, searches
over space of first
order features

objective used = likelihood
of Naive Bayes (MLRules)

(d) Breadth
first search

Exhaustive
search (X not feasible)



④ BEAM SEARCH = greedy BFS

Breadth first search but retain & expand only top few promising nodes at each level

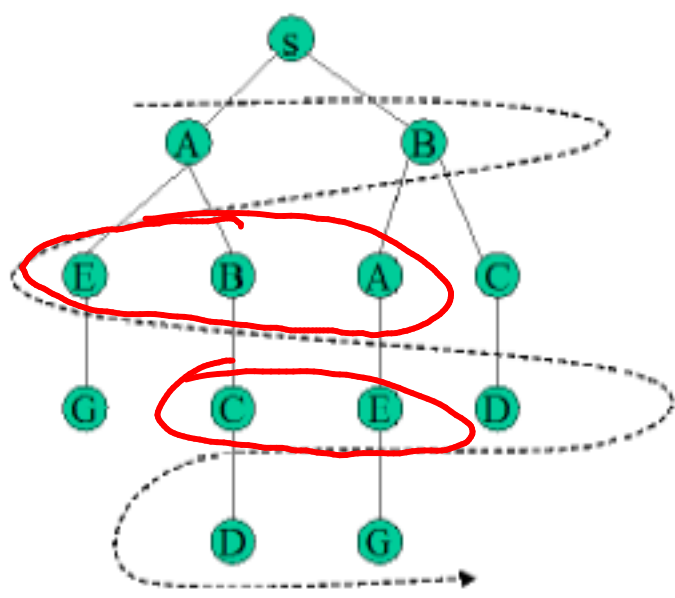
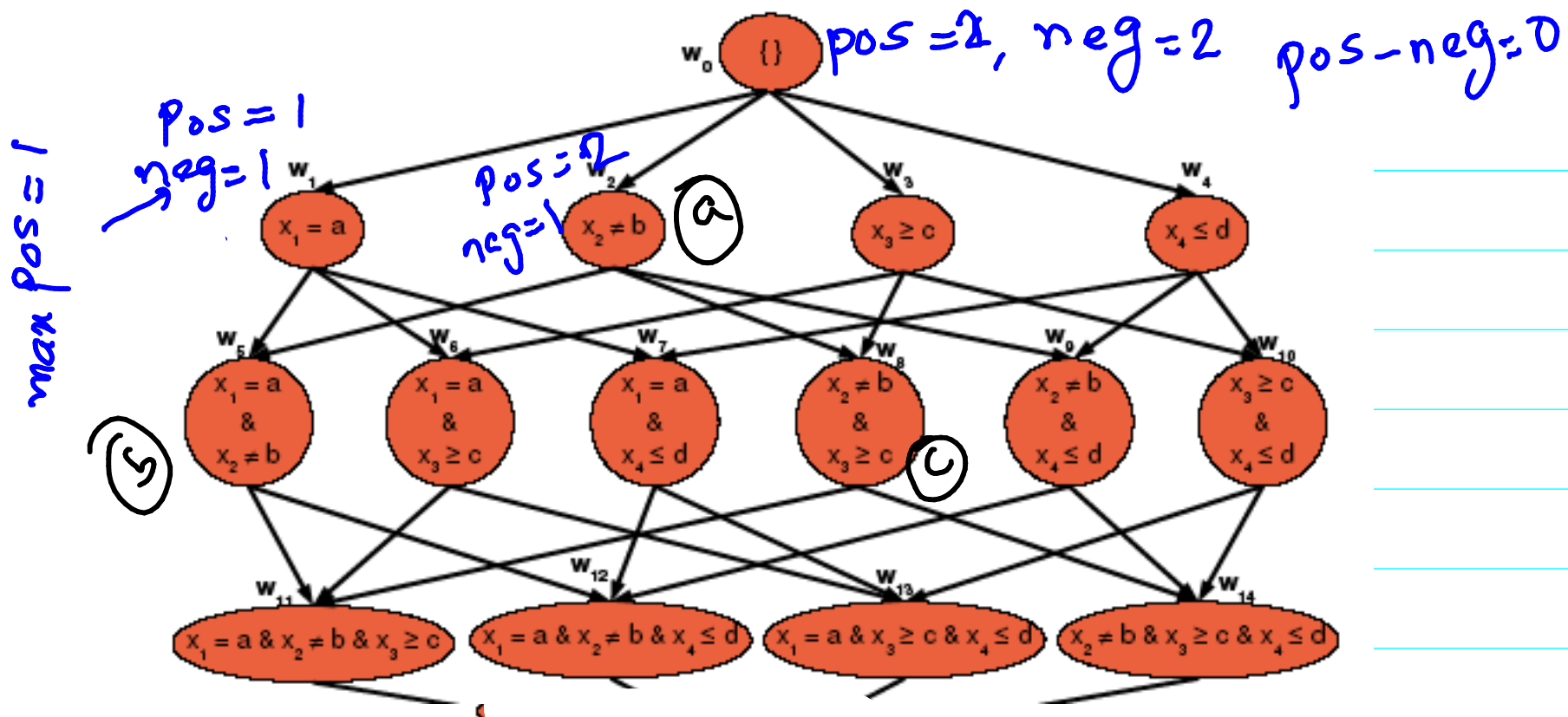


Figure 6.6: The BFS tree from S to G through the map in Figure 6.3.

Non-exhaustive

eg: Decision tree construction in propositional as well as first order logic [TILDE]

⑤ Branch & Bound search



$$\mathcal{D} = \left[\begin{array}{l|l} e_1^+ = \begin{cases} x_1 = a \\ x_2 \neq b \end{cases} & e_1^- = \begin{cases} x_1 \neq a \\ x_2 \neq b \end{cases} \\ e_2^+ = \begin{cases} x_2 \neq b \\ x_3 \geq c \end{cases} & e_2^- = \begin{cases} x_1 = a \\ x_3 \geq c \end{cases} \end{array} \right]$$

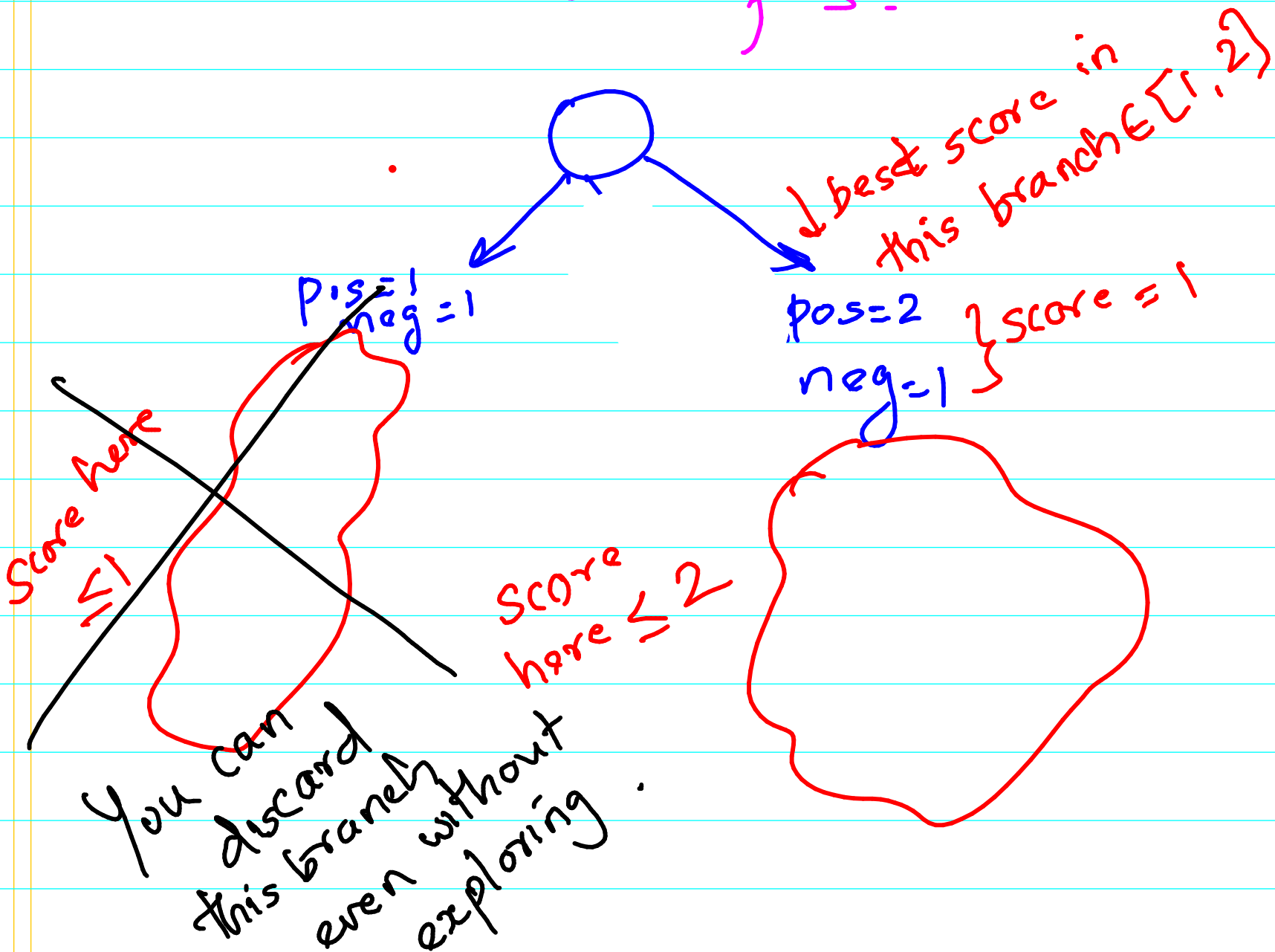
find s^* s.t. $s^* = \underset{s \in S}{\operatorname{argmax}} \operatorname{pos}(s) - \operatorname{neg}(s)$

if $\operatorname{pos}(s)$ & $\operatorname{neg}(s)$ are # of positive & negative examples covered by s ,

then $\forall s' \preceq s, \text{pos}(s') \leq \text{pos}(s)$
 $\preceq \implies 0 \leq \text{neg}(s') \leq \text{neg}(s)$

$$\implies \text{pos}(s') - \text{neg}(s') \leq \text{pos}(s)$$

$\text{pos}(s)$ is an upper bound on score of all children of s .



In fact, we can show that

node (a) = $\operatorname{argmax}_{s \in S} \text{pos}(s) - \text{neg}(s)$ → out branch & bound works for this

(*)

{ (b), (c) } = $\operatorname{argmax}_{S' \subseteq S} \sum_{s \in S'} \text{pos}(s) - \text{neg}(s)$

↓

(*)

While for pos-neg, a bound for (*) is not known, we have bounds found for

flow abt branch & bound for this?

Generalized HKL

$$\min_{w, b} \frac{1}{2} \left(\sum_{v \in \mathcal{V}} d_v \|w_{D(v)}\|_\rho \right)^2 + C \sum_{i=1}^m L \left(y^i, \sum_{v \in \mathcal{V}} w_v R_v(x^i) - b \right) = \text{obj}(\mathcal{V}' \subseteq \mathcal{V})$$

where $1 < \rho \leq 2$.

this obj ↓

Hierarchical norm regulariser

$S = \mathcal{V}$

We solve

$\operatorname{argmax}_{\mathcal{V}' \subseteq \mathcal{V}} \text{obj}(\mathcal{V}')$

using branch & bound

bound is just a sufficiency condition

We have also extended this branch and bound approach to a modified StructSVM formulation (Regulariser modified as above)

All this and more at

<http://www.cse.iitb.ac.in/~cs717/notes/classNotes/InterpretableFeatures/featureInductionTalk.pdf>