

Convex Optimization (CS709)

Instructor: Saketh

Contents

Contents	i
1 Mathematical Program – its parts, Review of Vector spaces	3
2 Sub-spaces, Basis, Inner-product spaces	7
3 Inner product spaces, Induced Norms, Orthogonality, Orthogonal Basis	11
4 Limits, Hilbert Space, Direct Sum, Affine Sets	15
5 Affine Sets and Cones	19
6 Cones and Convex sets	21
7 Convex sets and their Polars	23
8 Separation theorem and Polyhedrons	25
9 Polyhedrons and Linear Functions	27
10 Affine and Conic Functions	31
11 Introduction to Convex Functions	33
12 Conjugate, Sub-gradient, second-order derivative conditions for convexity	35

13 Introduction to Gradients	39
14 First and Second Order Conditions for Convexity	41
15 Convex Programs and Four Fundamental Questions	45
16 Characterizing Optimal Set	49
17 KKT Conditions for Simple Cases	51
18 KKT Conditions	55
19 Introduction to Duality in MPs and PCCP Duality	57
20 Lagrange Duality	61
21 Duality Example and Conic Duality	63
22 Conic Quadratic Programs	65
23 Semi Definite Programs	67
24 SDP examples, Lagrange Relaxation and Geometric Programs	69

Lecture 1

- Closer look at an optimization problem
 - Provided an example of real-world machine learning problem: that of support estimation¹.
 - Discussed how the machine learning problem is typically posed as an optimization problem. Infact, we discussed three different ones².
 - Each English language description of the optimization was converted into one in Math language, which is called as a Mathematical Program.
- Formal definition of **Mathematical Program (MP)**: A symbol that is of one of the following (equivalent) forms:

$$(1.1) \quad \begin{array}{ll} \min_{x \in \mathcal{X}} & \mathcal{O}(x) \\ \text{s.t.} & x \in \mathcal{F} \end{array}$$

or

$$(1.2) \quad \begin{array}{ll} \min_{x \in \mathcal{X}} & f(x) \\ \text{s.t.} & g_i(x) \leq 0, \forall i = 1, \dots, m \end{array}$$

- Defined various components of an *MP*:

¹Interested students may look at research.microsoft.com/pubs/69731/tr-99-87.pdf

²Hence there might be multiple ways of posing a real-world problem as an optimization problem. Each may have its own merits and de-merits from the application's perspective.

1. **Domain of the MP (\mathcal{X})** — the domain in which the variable(s) lives. This is also the domain for the functions \mathcal{O}/f and all g_i i.e., $\mathcal{O} : \mathcal{X} \mapsto \mathbb{R}$, $f : \mathcal{X} \mapsto \mathbb{R}$ and $g_i : \mathcal{X} \mapsto \mathbb{R}$. In the examples, the domain was once set of Euclidean vectors, once a mix of Euclidean vectors and matrices, once a set of functions. In this course we will focus on domains that are sets of vectors in what are known as “Hilbert Spaces”. Hilbert spaces are straight-forward generalizations of Euclidean spaces. This will be our first topic of technical discussion.
 2. **Feasibility set ($\mathcal{F} = \{x \in \mathcal{X} \mid g_i(x) \leq 0 \forall i = 1, \dots, m\}$)** — this is a subset of the domain in which the variables are restricted to lie in. We will study special subsets³ of the domain, which have nice properties and are easy to deal with. The focus in this course is on *MPs* with Feasibility set as a “convex set”. Again, $g_i : \mathcal{X} \mapsto \mathbb{R}$.
 3. **Objective function (\mathcal{O}/f)** — the function of the variable which is to be minimized $\mathcal{O} : \mathcal{X} \mapsto \mathbb{R}$, $f : \mathcal{X} \mapsto \mathbb{R}$. We will study some special real-valued functions on \mathcal{X} which have some interesting properties. The focus of this course is on *MPs* with “convex” objective functions.
- Defined the value or **optimal value of the *MP***: as infimum (greatest lower bound)⁴ of the set of objective function values on the feasibility set i.e., $(15.1) = \inf(\{\mathcal{O}(x) \mid x \in \mathcal{F}\})$. Similarly, $(15.2) = \inf(\{f(x) \mid x \in \mathcal{X}, g_i(x) \leq 0 \forall i = 1, \dots, m\})$. By convention we define $-\infty$ as the value of the *MP* for which this set of function values is not bounded below. Again, by convention, we define value of an *MP* with $\mathcal{F} = \phi$ (feasibility set is empty) as ∞ . With this convention, note that all *MPs* have a well-defined (optimal) value.
 - Identified and defined the related problem (argmin/argmax):

$$(1.3) \quad \begin{aligned} & \arg \min_{x \in \mathcal{X}} \mathcal{O}(x), \\ & \text{s.t. } x \in \mathcal{F}, \end{aligned}$$

which is defined as that $x^* \in \mathcal{F}$ such that $\mathcal{O}(x^*)$ is equal to the optimal value of the corresponding *MP*. Note that such an x^* may not always exist.

- In course of the lectures, we will:
 1. analyze each component of a (convex) *MP* in detail (this subject goes with the name “convex analysis”)

³For us, subset means subset or equal to.

⁴Please revise notions of maximum, minimum, GLB(infimum), LUB(supremum) and their existence results, atleast for sets of real numbers. <http://en.wikipedia.org/wiki/Supremum> should be enough.

2. analyze *MPs* with convex objective functions and convex Feasibility sets in finite dimensional “Hilbert spaces” (Euclidean spaces for now) — which are called as **Convex Programs (CPs)**. Some of the key questions we will answer are: when is an *MP bounded, solvable?* Can we characterize an *optimal solution?* Is it unique? etc.
 3. understand the very important and useful notion of duality which gives ways of arriving at *equivalent* optimization problems for the given problem — this may lead to deep insights into the problem/solution-structure or may lead to efficient solving techniques.
 4. Study standard *CPs* for which off-the-shelf generic solvers are available.
 5. Study special (scalable?) optimization techniques which work on generic *CPs*.
- We started revising vector spaces⁵:
 - Given a non-empty set V endowed with two operations $+_V$ (vector addition: $+_V : V \times V \mapsto V$) and \cdot_V (scalar multiplication: $\cdot_V : \mathbb{R} \times V \mapsto V$), if $(V, +_V)$ form an Abelian group, and the operator \cdot_V is associative such that $v \in V \Rightarrow 1 \cdot_V v = v$ and the distributive laws governing the interaction of $+_V$ and \cdot_V hold, then the triplet $\mathcal{V} = (V, +_V, \cdot_V)$ is called a vector space and elements of V are called as vectors.
 - We gave a lot of examples of vector spaces — those with matrices, polynomials, random variables, functions etc. We identified the additive identity (0_V) in each case. We gave a couple of examples of spaces that are not vector spaces.

⁵Go through pages 1–12 in [Sheldon Axler, 1997]. Also go through related exercises.

Lecture 2

- Defined notion of **linear combination** of two vectors: given $\lambda_1, \lambda_2 \in \mathbb{R}$ and $v_1, v_2 \in V$, then, $\lambda_1 v_1 + \lambda_2 v_2$ is the linear combination of v_1 and v_2 with weights λ_1 and λ_2 . By induction, one can define linear combination of a finite number of vectors.

- Defined **linear span**¹ of set S :

$$LIN(S) = \left\{ \sum_{i=1}^m \lambda_i v_i \mid \lambda_i \in \mathbb{R} \forall i = 1, \dots, m, v_i \in S \forall i = 1, \dots, m, m \in \mathbb{N} \right\};$$

this is the set of all possible linear combinations with the elements of the set S .

- A set S is called a spanning set of vector space $\mathcal{V} = (V, +_V, \cdot_V)$ iff $LIN(S) = V$.
- Let $\mathcal{V} = (V, +_V, \cdot_V)$ be a vector space and $W \subset V$ such that $\mathcal{W} = (W, +_V, \cdot_V)$ is itself a vector space², then \mathcal{W} is said to be a **sub-space** of \mathcal{V} and the set W is known as a **linear set or linear variety**.
- We gave many examples of linear sets.
- We then talked about compact representations of vector spaces. The first answer was spanning set³.
- **A vector space is finite-dimensional if there exists a spanning set of finite size/cardinality**. A vector space is infinite-dimensional if there exists no spanning set of finite size.
- Obvious question was whether we can get some kind of minimal spanning set? We outlined a procedure:

¹ $|S|$ represents cardinality of the set S .

²this condition is equivalent to: W being closed under linear combinations.

³Atleast one spanning set always exists: the set itself.

- Start with a spanning set S .
 - If $0_V \in S$, then remove it from S .
 - Verify if v_1 can be written as lin. comb. of the others. If yes, then remove it; else let it remain.
 - Repeat this for all elements of S .
- Note that at the end of the procedure one is left with a spanning set. More importantly, the spanning set is special: it is a linearly independent set⁴. A set $S = \{v_1, \dots, v_m\}$ is said to be linearly independent iff $\lambda_1 v_1 + \dots + \lambda_m v_m = 0 \Rightarrow \lambda_1 = \lambda_2 = \dots = \lambda_m = 0$.
 - A spanning set that is linearly independent is called as a **basis**. Infact, we just showed that every finite-dimensional vector space has a basis.
 - Theorem 2.6 in Sheldon Axler [1997] says that cardinality of a linearly independent set is always lesser than that of a spanning set. From this it easily follows that cardinality of any basis of a vector space is the same. Hence basis is indeed the smallest spanning set.
 - The common cardinality of all bases is called the **dimensionality of the vector space**. In other words, to describe a n -dimensional vector space using linear combinations we will require n (linearly independent) vectors⁵.
 - Interestingly, basis also give a way to strike an equivalence between any n -dimensional vector space and the n -dimensional Euclidean space:
 - Let $B = \{v_1, \dots, v_n\}$ denote the basis of the n -dimensional vector space in question. Let $v = \lambda_1 v_1 + \dots + \lambda_n v_n$ and $w = \alpha_1 v_1 + \dots + \alpha_n v_n$. It is easy to show that $v = w \Leftrightarrow \lambda_i = \alpha_i \forall i = 1, \dots, n$. This shows that every vector in the vector space can be mapped to exactly one vector in \mathbb{R}^n i.e., there exists a bijection from V to \mathbb{R}^n .
 - What is more interesting is: linear combinations are preserved under this bijective map i.e., $\gamma \cdot_V v +_V \delta \cdot_V w$ is mapped to the Euclidean vector corresponding to the same linear combination of the corresponding images i.e., mapped to $\gamma \begin{bmatrix} \lambda_1 \\ \vdots \\ \lambda_n \end{bmatrix} + \delta \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix}$.

⁴Take this as an exercise.

⁵Further compression using linear combinations is not possible

- Hence a basis is like a pair of goggles, through which the **vector space** looks like a Euclidean space (of same dimension). This statement should help us in visualizing spaces of matrices, polynomials etc. and hopefully also in solving MPs involving them as variables.
- We noted basis for all examples of vector/sub-spaces we considered. In each case we noted the dimensionality too.
- A basis gives an inner/constitutional/compositional/primal description (a description of an object with help of parts in it) of the vector space it spans. Looking at the example of x-y plane in \mathbb{R}^3 , the question arose: can we give an alternate description of x-y plane that is not constitutional? and perhaps which is simpler? One way out is to define x-y plane as all those vectors “orthogonal” to the unit vector on the z-axis. This description uses a vector not in the subspace under question (and hence we get a outer/non-compositional/dual description). Since the notion of orthogonality in Euclidean spaces springs out from the notion of dot product, we went ahead generalizing this notion to arbitrary vector spaces. This generalization of dot-product is called an **inner product**: given a vector space $\hat{V} = (V, +_V, \cdot_V)$, a function $\langle \rangle_V : V \times V \mapsto \mathbb{R}$ is called an inner-product iff:

- $v \in V \Rightarrow \langle v, v \rangle_V \geq 0, \langle v, v \rangle_V = 0 \Leftrightarrow v = 0$ (positive-definiteness).
- $v, w \in V \Rightarrow \langle v, w \rangle_V = \langle w, v \rangle_V$ (symmetry).
- $u, v, w \in V$ and $\alpha, \beta \in \mathbb{R} \Rightarrow \langle \alpha \cdot_V u +_V \beta \cdot_V v, w \rangle_V = \alpha \langle u, w \rangle_V + \beta \langle v, w \rangle_V$ (distributive law).

The quadruple $\mathcal{V} = (V, +_V, \cdot_V, \langle \rangle_V)$ is known as an inner-product space.

- We gave many examples of inner-products⁶:
 - Euclidean Spaces $(\mathbb{R}^n, +, \cdot)$:
 - * Dot product: $\langle x, y \rangle = x^\top y$
 - * $\langle x, y \rangle_W = x^\top W y$, where W is a given diagonal matrix with positive entries
 - * $\langle x, y \rangle_W = x^\top W y$, where W is a given positive-definite matrix⁷.
 - Space of matrices $(\mathbb{R}^{n \times n}, +, \cdot)$:

⁶Student should prove correctness of each example. Some may require Cauchy-Schwartz inequality to be proved in next lecture..

⁷With this we said a sphere in this space will actually look like ellipse. In special case where W is also diagonal, the ellipse will actually be axis-parallel. **Nikunj conjectured that all inner-products on Euclidean space are of this form. Bonus marks will be awarded to the student who (correctly) proves or disproves it.**

- * Frobenius inner-product⁸: $\langle M, N \rangle_F = \sum_{i=1}^n \sum_{j=1}^n M_{ij}N_{ij}$.
 - * $\langle M, N \rangle_W = \sum_{i=1}^n \sum_{j=1}^n M_{ij}N_{ij}W_{ij}$, where $W_{ij} > 0 \forall i, j$.
 - * $\langle M, N \rangle_W = \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \sum_{l=1}^n M_{ij}N_{kl}W_{ijkl}$, where W is a $n \times n \times n \times n$ matrix/tensor that is *positive definite*⁹.
- Space of all L_2 functions¹⁰ $f : \mathbb{R} \mapsto \mathbb{R}$:
- * $\langle f, g \rangle = \int_{\mathbb{R}} f(x)g(x) dx$
 - * $\langle f, g \rangle_w = \int_{\mathbb{R}} f(x)g(x)w(x) dx$, where w is a function that takes on positive values only.
 - * $\langle f, g \rangle_w = \int_{\mathbb{R}} \int_{\mathbb{R}} f(x)g(y)w(x, y) dx dy$, where w is a *positive-definite function*¹¹.
- Space of all mean zero random variables with finite second moment ($\mathbb{E}[X^2] < \infty$):
- * $\langle X, Y \rangle = \mathbb{E}[XY]$
- The W matrix or the w function that plays a key role in defining the geometry of the space is called the kernel¹².

⁸resembles the dot product.

⁹The condition on W for this expression being inner-product is the natural definition of positive-definiteness for tensors.

¹⁰Refer en.wikipedia.org/wiki/Lp_space

¹¹The condition on w for the expression being an inner-product is the natural definition of positive definite function. Refer en.wikipedia.org/wiki/Mercer's_theorem for a key result about positive-definite functions. Similar to positive-definite matrices, for positive-definite tensors/functions one can talk about eigen-value-decomposition, with positive eigen-values!

¹²Observe that the kernel in case of the random variables example is the joint probability density function (assuming continuous random variables).

Lecture 3

- We began by proving an important result that follows from the definition of an inner-product: Cauchy-Schwartz inequality¹.
- We then generalized the notion of norm to abstract vector spaces: Given a vector space $\hat{\mathcal{V}} = (V, +_V, \cdot_V)$ and a function $\|\cdot\|_V : V \mapsto \mathbb{R}$, we say that the function $\|\cdot\|_V$ is a norm in the vector space $\hat{\mathcal{V}}$ iff:
 - $v \in V \Rightarrow \|v\|_V \geq 0, \|v\|_V = 0 \Leftrightarrow v = 0_V$ (Non-negativity)
 - $v \in V, \alpha \in \mathbb{R} \Rightarrow \|\alpha \cdot_V v\|_V = |\alpha| \|v\|_V$ (Distribution with \cdot_V)
 - $v, w \in V \Rightarrow \|v +_V w\|_V \leq \|v\|_V + \|w\|_V$ (Distribution with $+_V$ or triangle inequality)

The quadruple $\mathcal{V} = (V, +_V, \cdot_V, \|\cdot\|_V)$ is known as a normed vector space.

- We showed² that the function $\|v\|_V$ defined by $\|v\|_V = \sqrt{\langle v, v \rangle_V}$ is in fact a valid norm as per the above definition³. It is called the (inner-product) induced norm.
- We noted the expressions for induced norms in the various examples of inner-product spaces: e.g., Euclidean norm (induced by dot product), Frobenius norm (induced by Frobenius inner-product), standard deviation of random variable (induced by the inner-product in space of mean zero random variables). We noted expressions for spheres, ellipsoids in all these spaces.
- We gave examples of norms that are not induced norms⁴: e.g., for Euclidean vector x , $\|x\|_p = (\sum_{i=1}^n |x_i|^p)^{\frac{1}{p}}, p \geq 1, x_\infty = \max_{i \in \{1, \dots, n\}} |x_i|$ etc. From now

¹Refer pg. 104 in Sheldon Axler [1997]. Note that the proof in the book is different from the one that is done in the lecture.

²Refer pg.105 in Sheldon Axler [1997] for a proof.

³It now follows that every inner-product space has an associated normed vector space, formed with the induced norm.

⁴Interested students may attempt proving this.

onwards, wherever we say norm or write norm it means the induced norm unless specifically mentioned otherwise.

- It is natural to now define distance $dist(u, v) = \|u - v\|_V$.
- One can define **projection of a vector $v \in V$ onto a set $S \subset V$** : $P_S(v) = \arg \min_{x \in S} \|v - x\|_V$, i.e., the vector in S that is closest to v .
- We define cosine of angle between two vectors⁵: $\cos(\angle u, v) \equiv \frac{\langle u, v \rangle_V}{\|u\|_V \|v\|_V}$. According to this, angle $\angle u, v = 0^\circ \Leftrightarrow u = \alpha v, \alpha > 0$, $\angle u, v = 180^\circ \Leftrightarrow u = -\alpha v, \alpha > 0$ and $\angle u, v = 90^\circ \Leftrightarrow \langle u, v \rangle_V = 0$. This defines **orthogonality** for two vectors (vectors with angle between them as 90°).
- We gave examples of orthogonal vectors: e.g., any pair of symmetric matrix and skew-symmetric matrix are orthogonal; any function having a root at c is orthogonal to the function taking zero value everywhere except at c ; any pair of un-correlated random variables are orthogonal⁶ etc.
- We then realized that when provided with a basis where each element is of unit length and every pair of elements are orthogonal to each other, any finite dimensional inner-product space is equivalent to the Euclidean space with dot product as the inner-product. Such a special basis is called **orthogonal or orthonormal basis**.
- We outlined the Gram-Schmidt procedure⁷ for inductively constructing an orthogonal basis from a basis. Hence any finite dimensional inner-product space (IPS) has an orthogonal basis and is equivalent to the Euclidean one. As a corollary, all geometric results true in Euclidean IPS hold in all IPSs⁸.
- We then went ahead to answer the question of compact representation of vector space (now using notion of orthogonality rather than that of linear combinations): let $S \subset V$ be a linear set and B be a basis of it. Let $dim(V) = n$ and $dim(S) = m \leq n$. Consider the set $S^\perp \equiv \{v \in V \mid \langle v, u \rangle_V = 0 \forall u \in S\}$, called as the **orthogonal complement (or Dual space/set) of S**. We then noted the following results⁹:
 - S^\perp is a linear set and hence forms a subspace of \mathcal{V} . Let B^\perp be a basis for S^\perp .

⁵Cauchy-Schwartz inequality guarantees that this is a valid definition of cosine.

⁶We also noted that the cosine of angle between random vectors is commonly known as Pearson's correlation coefficient

⁷Refer section 6.20 in Sheldon Axler [1997].

⁸Take proving parallelogram law, Pythagoras theorem as exercises.

⁹Students should take proving each result as an exercise. The proof of the complementarity in terms of dimension follows from the Rank-Nullity theorem.

- $\{0_V\} = S \cap S^\perp$, which forms the smallest subspace¹⁰.
 - $\dim(S) + \dim(S^\perp) = n$ (hence the name includes the word complement).
 - $B \cup B^\perp$ is a basis for V .
 - $(S^\perp)^\perp = S$.
 - $S = \{v \in V \mid \langle v, u \rangle_V = 0 \ \forall u \in B^\perp\}$ (we will refer to this representation of S in terms of B^\perp as the **dual/outer/sculptor's representation of a linear set**) and $S^\perp = \{v \in V \mid \langle v, u \rangle_V = 0 \ \forall u \in B\}$.
- Hence, the basis of S^\perp is called the **dual basis** of S and vice-versa.
 - Note that when $m \leq \lfloor \frac{n}{2} \rfloor$ the basis is the most compact representation for S ; whereas in the other case the dual basis is the most compact representation.
 - We will call a linear set of dimension one less than that of the entire vector space as a **hyperplane** (through origin)¹¹.
 - **Mandatory reading:**
 - Sections A.1.1-A.1.4, A.2, A.4.1, A.7 in Nemirovski [2005].

¹⁰In fact, if $\{S_\lambda \mid \lambda \in \Lambda\}$ is a collection (possibly uncountable) of linear sets indexed by elements in the index set Λ , then $\bigcap_{\lambda \in \Lambda} S_\lambda$ is itself a linear set.

¹¹This is generalization of concept of line (through origin) in a plane and that of plane through origin in 3-d space etc. Needless to say, the dual representation is the most efficient representation (unless vector space is trivial) for a hyperplane. This is the reason from school days we are always taught to describe planes using equations (like $w^\top x = 0$) and rarely using the vectors in the plane!

Lecture 4

- We re-emphasized on the usefulness of primal and dual representations. We illustrated the example of symmetric matrices where the dual representation is definitely more compact than the primal representation.
- Once the notion of norm exists one can define limits/convergence: Let $\{x_n\}$ be a sequence of vectors. If for any given $\epsilon > 0, \exists N \ni \forall n \geq N$, we have: $\|x - x_n\| < \epsilon$, then the sequence is said to converge to x i.e., $\{x_n\} \rightarrow x$. x is called the limit of the sequence $\{x_n\}$ i.e., $\lim_{n \rightarrow \infty} x_n = x$.
- We know that Euclidean spaces have no gaps (they are complete spaces) i.e., every Cauchy sequence converges. Because of our equivalence, all finite dim. inner-product spaces are also complete and hence qualify to be called as **Hilbert spaces**¹ (complete normed vector spaces are called as Banach spaces).
- We talked about an operation called direct summing that will enable us to “join” a Hilbert space of say Euclidean vectors with that of say matrices: Given two inner-product/Hilbert spaces $\mathcal{V}_1 = (V_1, +_1, \cdot_1, \langle \rangle_1)$ and $\mathcal{V}_2 = (V_2, +_2, \cdot_2, \langle \rangle_2)$, we defined the direct sum of those, $\mathcal{V} = \mathcal{V}_1 \oplus \mathcal{V}_2$, which is another inner-product space defined as $\mathcal{V} = (V, +, \cdot, \langle \rangle)$, where $V \equiv V_1 \times V_2 = \{(v_1, v_2) \mid v_1 \in V_1, v_2 \in V_2\}$ i.e., V is the Cartesian product (or sometimes called as direct product) of the sets V_1 and V_2 . Given two vectors $v = (v_1, v_2), w = (w_1, w_2) \in V$, we have: $v + w \equiv (v_1 +_1 w_1, v_2 +_2 w_2)$, $\alpha \cdot v \equiv (\alpha \cdot_1 v_1, \alpha \cdot_2 v_2)$ and $\langle v, w \rangle \equiv \langle v_1, w_1 \rangle_1 + \langle v_2, w_2 \rangle_2$. This is the natural way of stacking up arbitrary spaces to form big space. Note that with such a direct sum, the following two sub-spaces are orthogonal complements of each other: $\hat{V}_1 = \{(v_1, 0_2) \mid v_1 \in V_1\}$ and $\hat{V}_2 = \{(0_1, v_2) \mid v_2 \in V_2\}$ (here, $0_1, 0_2$ denote the identity elements in $\mathcal{V}_1, \mathcal{V}_2$ respectively).
- This completed our study/review of domain of an MP. We then began with study of the next component of an MP, which is the feasibility set (as noted

¹Refer http://en.wikipedia.org/wiki/Hilbert_space.

earlier, is a subset of the domain). In other words, we began a study of special subsets of Hilbert spaces: linear sets, Affine sets, conic sets, convex sets. In case of each of these category of subsets, we will study the definition, primal and dual representations, some algebra and topology results.

- We will be concerned with the following algebraic operations over sets²:

Union: $\cup_{\lambda \in \Lambda} S_\lambda \equiv \{v \mid v \in S_\lambda \text{ for some } \lambda \in \Lambda\}$.

Intersection: $\cap_{\lambda \in \Lambda} S_\lambda \equiv \{v \mid v \in S_\lambda \forall \lambda \in \Lambda\}$. In the following we assume $\Lambda = \{1, \dots, n\}$.

Direct/Cartesian Product: $S_1 \times \dots \times S_n \equiv \{(v_1, \dots, v_n) \mid v_i \in S_i \forall i = 1, \dots, n\}$.

Linear Combination: $\alpha_1 S_1 + \dots + \alpha_n S_n \equiv \{\alpha_1 \cdot v_1 + \dots + \alpha_n \cdot v_n \mid v_i \in S_i \forall i = 1, \dots, n\}$. Here $\alpha_i \in \mathbb{R}$.

Complement: $S_1^c = \{v \in V \mid v \notin S_1\}$.

SetDifference: The set difference of S_1 and S_2 (denoted by $S_1 \setminus S_2$) is defined as $S_1 \cap S_2^c$.

- We will be concerned with the following topological concepts:

Closed Set: A set $S \subset V$ is said to be closed iff the limit of every convergent sequence in S belongs to S .

Open Set: A set $S \subset V$ is said to be open iff for every $s \in S$, there exists a $\epsilon > 0$ such that $N_\epsilon(s) \subset S$. Here, $N_\epsilon(s) \equiv \{v \in V \mid \|s - v\| \leq \epsilon\}$ is the ϵ neighborhood of s or equivalently a ball of radius ϵ centered at s .

Bounded Set: A set $S \subset V$ is said to be bounded iff there exists a finite radius $r > 0$ such that a ball of that radius centered at origin contains the set i.e., $S \subset N_r(0_V)$.

Interior: The set of all interior points of S is called the interior: $int(S)$. A vector/point $s \in S$ is called an interior point of S iff there exists a $\epsilon > 0$ such that $N_\epsilon(s) \subset S$. A set S is said to have interior or is said to have volume iff the interior is non-empty i.e., $int(S) \neq \phi$.

Boundary: Boundary of S is those vectors in S that are not in the interior of S : $\delta(S) \equiv S \setminus int(S)$.

Compact Set: A set that is closed and bounded is called a compact set.

²We will assume $\{S_\lambda \mid \lambda \in \Lambda\}$ is a collection of subsets of V , indexed by the set Λ . This index set could be uncountable.

- We also noted some standard results³ regarding these topological concepts:
 1. Complementarity of open and closedness: S is closed if and only if S^c is open.
 2. Intersection of (possibly uncountable no.) closed sets is closed; union of (possibly uncountable no.) open sets is open.
 3. Union of finite number of closed sets is closed and intersection of finite number of open sets is open.
 4. Heine-Borel theorem⁴.
 5. Bolzano-Weierstrass theorem⁵.

- We took up Linear sets first as we already know its definition and primal, dual representations⁶:
 - We noted that intersection of (possibly uncountable no.) linear sets is linear.
 - Union of linear sets need not be linear.
 - Linear combinations of linear sets are linear.
 - Complement of linear set is not linear.
 - Cartesian product of linear sets is linear.
 - Linear sets are always closed (since they form spaces equivalent to Euclidean ones)
 - All linear sets except that with the entire set of vectors are not open (and don't have volume).
 - All linear sets except the trivial one are not bounded.

- We modeled linear sets looking at planes/lines through origin. We then defined Affine sets to model those that may not contain the origin: A set $A \subset V$ is said to be an Affine set iff it can be written as $A = \{a_0\} + L$, where $a_0 \in V$ and $L \subset V$ is a linear set.

- We noted the following results for Affine sets:

³Again, since we deal with domains equivalent to Euclidean ones, we will take these as standard Analysis results and not prove them here.

⁴Refer en.wikipedia.org/wiki/Heine-Borel_theorem

⁵Refer en.wikipedia.org/wiki/Bolzano-Weierstrass_theorem. Also section A.4.4 in Nemirovski [2005]

⁶Students should prove these claims.

- Given A , the associated linear set L is fixed. Infact, $L = A - A$. However a_0 can be replaced by any $a \in A$.
- Suppose L is the linear set associated with A and its basis is $\{l_1, \dots, l_m\}$, then $a \in A \Rightarrow a = a_0 + \sum_{i=1}^m \lambda_i l_i$. Further, it is easy to see any l_i can be written as $a_i - a_0$ (for some $a_i \in A$). So, $a = (1 - \sum_{i=1}^m \lambda_i) a_0 + \sum_{i=1}^m \lambda_i a_i$. In other words, any vector in an Affine set in a vector space can be uniquely written as $[\rho_0 \ \rho_1 \ \dots \ \rho_m]^\top$ where $\sum_{i=1}^m \rho_i = 1$ (and hence equivalent to a hyperplane that does not pass through origin in \mathbb{R}^{m+1}).
- **Mandatory reading:**
 - Sections A.4 and A.3 in Nemirovski [2005].
- **Optional reading:** section 1 in Rockafellar [1996].

Lecture 5

- We summarized the results with Affine sets detailed in [section A.3 in Nemirovski \[2005\]](#)
 - Definition is shifted linear set. In particular, all linear sets are Affine.
 - Primal view is set closed under Affine combination. The notion of Affine basis is immediate (Infact, we write this using the basis of the linear set associated with the Affine set).
 - Dual view is (finite) intersection of hyperplanes that need not pass through origin.
 - Dimension of an Affine set is that of the associated linear set and it requires $n + 1$ elements to form a n -dimensional Affine set.
 - All Affine sets are closed (follows from the result with linear sets) and none except the entire set of vectors is an open set. Affine sets except the trivial ones (i.e., singleton vectors) are un-bounded.
 - Sum and intersection of Affine sets is Affine; whereas union need not be. Complement will not be an Affine set.
- We next looked at two special sets associated with a hyperplane (through origin): If the hyperplane is given by $H = \{u \in V \mid \langle v_H, u \rangle = 0\}$, then we call the set $H_+ = \{u \in V \mid \langle v_H, u \rangle \geq 0\}$ as its positive half-space and the set $H_- = \{u \in V \mid \langle v_H, u \rangle \leq 0\}$ as its negative half-space. We then wished to study sets which are intersections of such half-spaces:
- We defined **cone** as a set that is closed under **conic combinations** of its elements. $\sum_{i=1}^n \lambda_i v_i$, where each $\lambda_i \geq 0$, is called as the conic combination of the vectors v_1, \dots, v_n with weights $\lambda_1, \dots, \lambda_n$.
- **Conic hull** of a set S is defined as: $CONIC(S) = \{\sum_{i=1}^n \lambda_i v_i \mid v_i \in S, \lambda_i \geq 0 \forall i, n \in \mathbb{N}\}$.

- Hence K is a cone iff $K = \text{CONIC}(K)$.
- S is called a **conicly spanning set** of K iff $K = \text{CONIC}(S)$. We then took many examples of cones: i) cones in Euclidean spaces: we constructed them by taking some set S and looking at $\text{CONIC}(S)$ e.g., was the half-space (through origin), ice-cream cone, infinite wedge etc. ii) set of all psd matrices¹ iii) set of all kernels over a given L_2 space².
- We noted that in some examples the conicly spanned set was finite and in some it was not possible to get a finite-sized conicly spanning set. The cones that have finite-sized conicly spanning sets are called as **polyhedral cones**³.
- The obvious question now is can we get compact representations of a cone K : one way is to look for the smallest⁴ conicly spanning set of K . The other way perhaps is to describe cones using inner-products (i.e., dual description)⁵.
- Drawing an analogy to the notion of orthogonal complement (dual space) of a linear set, we defined **dual cone of a cone K** : $K^* = \{v \in V \mid \langle v, u \rangle \geq 0 \ \forall u \in K\}$. It is an easy exercise to show that K^* is indeed a cone.
- Meghshyam noted that the notion of dual cone is consistent with that of orthogonal complement in the sense that in the special case the cone K is a linear set, then dual cone is nothing but the orthogonal complement of K .
- For all the examples we noted the dual cones.
- We realized interesting cases where dual cone is same as the original cone. Such cones are called as self-dual cones. e.g., ice-cream cone, cone of psd/kernels (in the vector space of all symmetric matrices/functions).
- **Mandatory reading:**
 - Section B.1.4, B.2.6.B in Nemirovski [2005], section 2.6.1 in Boyd and Vandenberghe [2004].
- Optional reading: relevant parts on cones in section 2,14 in Rockafellar [1996].

¹The conicly spanning set was the set of all symmetric rank one matrices of the corresponding dimension.

²The conicly spanning set was the set of all $k(x, y) \equiv \bar{k}(x)\bar{k}(y)$

³Wedge is polyhedral; whereas ice-cream is not.

⁴It turns out that the notion of basis or Affine basis cannot be simply carried over to cones. Hence we will postpone the answer until we discuss convex sets.

⁵It is easy to see that the description of half-space using inner-products is simply its definition and requires less number of vectors than its primal description through conic combinations. This motivates our definition of Dual cone.

Lecture 6

- We attempted proving an interesting result: for a closed cone¹ K , we have $(K^*)^* = K$. While it was easy to see that $K \subset (K^*)^*$, we said it is not straightforward to show the converse. We noted that a *separation theorem*, which we will state and prove in coming lectures on convex sets, will help proving it. Infact we mentioned all duality concepts including that of notion of subgradients for convex functions follow from this basic fundamental theorem².
- For now, we assumed that the above conjecture is true and hence dual description of a closed cone is immediate: closed cone is always intersection of some half-spaces. Infact, all sets that are intersections of half-spaces are cones. Hence, this could have been taken as the definition of closed cones.
- The following results about algebra with cones are true³:
 - $\cap_{i \in I} K_i$ is a cone; whereas $K_1 \cup K_2$ need not be a cone. Complement of cone is not a cone.
 - $K_1 + \dots + K_n = CONIC(K_1 \cup \dots \cup K_n)$.
 - $(\cap_{i=1}^n K_i)^* = \sum_{i=1}^n K_i^*$. (Dubovitski-Milutin lemma)
- We then defined a **convex set** C : C is convex iff $u, v \in C \Rightarrow \lambda_1 v + \lambda_2 w \in C \forall \lambda_1, \lambda_2 \geq 0, \lambda_1 + \lambda_2 = 1$.
- $\sum_{i=1}^n \lambda_i v_i$, where each $\lambda_i \geq 0$ and $\sum_{i=1}^n \lambda_i = 1$, is called as the **convex combination** of the vectors v_1, \dots, v_n with weights $\lambda_1, \dots, \lambda_n$. By induction, C is convex iff $C = CONV(C)$, where $CONV(C)$ is the **convex-hull** of C , which is the set of all convex combinations with elements of C .

¹ $\mathbb{R}_{++}^n \cup \{0\}$ is an example of a cone that is not closed.

²We may call it the **Fundamental theorem of convex analysis**.

³Students should take them as exercise. You may need separation theorem. Here, every $K_i, i \in I$ represents a cone (need not be polyhedral unless mentioned otherwise).

- We then gave many examples of convex sets in Euclidean and matrix spaces: polygons, spheres, Birkhoff polytope, set of all probability density functions.
- We say C is a polytope iff $\exists S \subset C \ni C = \text{CONV}(S), |S| < \infty$.
- We then defined a n -dimensional simplex⁴ as $\text{CONV}(S)$, where S is an affinely independent set of size $n + 1$.
- In general, for a set S , dimensionality of S is defined as that of the affine-hull of it i.e., $\text{dim}(S) \equiv \text{dim}(\text{AFF}(S))$.
- **Mandatory reading:**
 - Sections B.1.1-B.1.5 in Nemirovski [2005], sections 2.1-2.3 in Boyd and Vandenberghe [2004].
- **Optional reading:** sections 2,3 in Rockafellar [1996].

⁴We will later on note that all convex sets have simplices in them (triangulation of sets) and hence they have volume (when restricted to the affine hull of the set).

Lecture 7

- We wished to come-up with a dual description of convex sets. Intuitively we argued that convex sets are intersections of (arbitrary no.) half-spaces that need not pass through origin. This motivated the definition of **polar of a set C** as $C' \equiv \text{polar}(C) \equiv \{x \mid \langle x, c \rangle \leq 1 \ \forall c \in C\}$.
- It is easy to show that C' is a convex set for any arbitrary C .
- We noted that in case C is a linear set, then polar is same as orthogonal complement and in case C is a cone, then polar is same as dual cone. Hence the definition of polar is a natural extension of those in case of linear and conic sets.
- Infact, the following interesting duality result can be proved:

Theorem 7.0.1. *A set C is closed convex containing origin if and only if $C = \text{polar}(\text{polar}(C))$.*

Refer proposition B.2.2 in Nemirovski [2005] for proof..

- Again, like in the case of cones, this duality result also is proved using the (yet to be proved) separation theorem!
- From this theorem it is clear that C is closed convex if and only if C can be written as intersection of (perhaps arbitrary no.) half-spaces (that need not pass through origin)¹. In other words, this could have been the alternate definition of closed convex sets.
- We then illustrated how the polars of various closed convex sets (containing origin), presented earlier as examples, looked like.

¹This can be proved by simply shifting the origin to lie inside the set C and then applying the duality theorem above and then shifting back the origin.

- Infact, more duality results as given in Exercise B.14 and B.15 in Nemirovski [2005] can be proved².
- **Mandatory reading:**
 - Section B.2.6 in Nemirovski [2005].
- Optional reading: section 14 in Rockafellar [1996].

²This is a student exercise.

Lecture 8

- We began by giving a brief summary of the primal/dual representations and results regarding algebra and topology of all special sets we studied till now (please refer to first page in Appendix).
- We further stressed on an important topological result that every convex set has volume/interior (when restricted to its dimensionality i.e., when restricted to its affine-hull; this is the notion of [relative interior/volume](#)). [Please refer theorem B.1.1 in Nemirovski \[2005\]](#). An easy proof of this is: prove that every convex set contains a simplex of the same dimension in it. Then showing a simplex has volume is easy (inscribed circle/hypersphere exists).
- To illustrate the point that sometimes in proving a set is (closed) convex, neither of the primal or dual definitions are easy to verify: we looked at (non-empty) $C = \{x \in \mathbb{R}^n \mid x^\top Ax + b^\top x + c \leq 0\}$. The claim was this (non-empty) set is convex whenever $A \succeq 0$. We illustrated a very easy proof of this using the following useful 1-dimensional characterization of convex sets: A set C is convex if and only if its (non-empty) intersection with any line is convex.
- We then went ahead and gave a simple proof of the separation theorem. Please refer to the appendix (pages 2-3) for the definition of separation and the proof. Though we may not use explicitly, the general separation theorem [Theorem B.2.5 in Nemirovski \[2005\]](#) is a good thing to know.
- We then wrote down the result of separation theorem when applied to a polyhedral cone and a vector, slightly differently, leading to the Farkas lemma (refer sec.B.2.4 in Nemirovski [2005]), and saw that duality sometimes helps us answer difficult questions by posing the difficult question as an easy question on a dual. Here is one way of writing Farkas lemma:

Lemma 8.0.2. *Consider two sets of linear inequalities (S_1) given by: $Ax = b, x \geq 0$ (here, x is the dummy variable) and (S_2) given by $A^\top y \geq 0, b^\top y < 0$*

(here, y is the dummy variable). Separation theorem gives that (S_1) is solvable/consistent/feasible if and only if (S_2) is not-solvable/in-consistent/infeasible.

There are many ways of writing down such results and in general are called as “Theorems on Alternative”. Some of them appear in theorem 1.2.1 and exercises 1.2-1.4 in Nemirovski [2005]. We will realize later that one way of deriving duals of optimization problems in infact by using such theorems on alternative.

- We then asked the question are there convex sets where dual description is more efficient than primal and vice-versa. For a unit circle at origin: $C = \{x \mid \|x\| \leq 1\}$, the (“smallest”) convexly spanning set is $\{x \mid \|x\| = 1\}$ and (efficient) dual description is $\{x \mid x^\top y \leq 1 \ \forall \|y\| = 1\}$. So, this is the case of self-duality and hence both primal and dual descriptions are equally efficient. So is the case with polytopes (atleast intuitively), provided the dimensionality is same as that of the space. In case of polytopes of dimension less than that of the space, the primal is more efficient than dual.
- However there are striking examples of convex sets where the dual description is “infinitely better” than the primal: this is the case of half-spaces, cones, shifted cones etc. We then defined [convex sets that are intersections of finite number of halfspaces as polyhedron or polyhedral set](#).
- With some examples we conjectured the Minkowski-Weyl theorem¹: A set C is polyhedron if and only if there exist finite sets S, T such that $C = CONIC(S) + CONV(T)$.
- While we postponed the proof of this to next lecture, it is easy to see the following consequence of above result: every polytope is a polyhedron.
- **Mandatory reading:**
 - [Section B.1.6, B.2.5 in Nemirovski \[2005\]. Section 2.5 in Boyd and Vandenberghe \[2004\]](#)
- **Optional reading:** section 11 in Rockafellar [1996].

¹This result is analogous to the fact that affine sets are shifted linear sets: polyhedrons are polyhedral cones shifted with a polytope.

Lecture 9

- We proved the Minkowski-Weyl theorem using the proof methodology in Lauritzen [2010] (refer theorem 4.5.1 in Lauritzen [2010]). Appendix pages 3-5 provide a proof for the same. Here are the key steps:
 1. Consider a cone whose projection is the polyhedron.
 2. Assuming a cone with finite dual description is polyhedral, write down the primal description of the cone.
 3. The required projection (that is polyhedral) is already in form of polytope + polyhedral cone.
 4. Then indeed prove that polyhedral cone has finite dual description. Since dual of dual cone is the original cone, this statement gives that cones with finite dual description are polyhedral:
 - (a) The trick is to write down the polyhedral cone as projection of some cone with finite dual description.
 - (b) Since projections of cones with finite dual description have finite dual description, we have the required result. In our case, the projection required was onto the last few dimensions. This we obtained by systematic elimination of variables. Please refer theorem 1.2.2 in Lauritzen [2010].
- We then went ahead to study the final ingredient of mathematical programs, which is real-valued functions defined on (subsets of finite-dimensional) Hilbert spaces.
- We began with few examples of such functions $f : S \mapsto \mathbb{R}$ (here, $S \subset V$). Then talked about some important sets associated with functions:
 - Graph of f : $graph(f) \equiv \{(x, y) \mid x \in dom(f), f(x) = y\} \subset V \times \mathbb{R}$. This lies in the space that is direct sum of \mathcal{V} and \mathbb{R} .

- Epigraph of f : $epi(f) \equiv \{(x, y) \mid x \in dom(f), f(x) \leq y\} \subset V \times \mathbb{R}$. This lies in the space that is direct sum of \mathcal{V} and \mathbb{R} .
 - Level-set of f at α : $L_\alpha(f) \equiv \{x \in dom(f) \mid f(x) \leq \alpha\} \subset V$. This lies in the space of \mathcal{V} itself.
- We said that we will study special functions whose associated sets (graph/epigraph/level-sets) are special (i.e., convex/conic/affine/linear).
 - We began our study with linear functions: $f : L \mapsto \mathbb{R}$, where $L \subset V$ is a linear set, is a linear function if and only if linear combinations are preserved under the function i.e., $\lambda_1, \lambda_2 \in \mathbb{R}, x_1, x_2 \in L \Rightarrow f(\lambda_1 x_1 + \lambda_2 x_2) = \lambda_1 f(x_1) + \lambda_2 f(x_2)$.
 - After giving some examples we noted the following important result that was very easy to prove: f is linear if and only if $graph(f)$ is a hyperplane (in direct sum space of vectors $L \times \mathbb{R}$)¹ that passes through origin:
 1. We first showed f is linear if and only if $graph(f)$ is a linear set. This was straight-forward to prove.
 2. We then argued that $dim(graph(f))$ is atleast $dim(L)$ (as f must be defined at atleast $dim(L)$ no. linearly-independent points; infact defining f at all points/vectors in any basis will completely define the function). Also, since (x, y) whenever $y \neq f(x)$ is not a point in the graph, the dimensionality of the linear set is not $dim(L) + 1$. Hence $dim(graph(f))$ must be $dim(L)$.
 - We then proved the Riesz representation theorem (the set of linear functions on L is equivalent to L itself)²:
 1. For linear f (since graph is a hyperplane), we have that there exists $(u, v) \in L \times \mathbb{R}$ such that, $graph(f) = \{(x, y) \mid \langle (u, v), (x, y) \rangle = 0\}$. It is easy to see that $v \neq 0$ (otherwise graph is parallel to \mathbb{R} axis). Dividing the entire hyperplane equation by v , we obtain $f(x) = \langle -\frac{u}{v}, x \rangle$.
 2. Also any function of the form $f(x) = \langle l, x \rangle$ where $l \in L$ is a linear function on L .
 3. Moreover, if $f_1 = \langle l_1, x \rangle$ and $f_2 = \langle l_2, x \rangle$, then $f_1 = f_2$ if and only if $l_1 = l_2$.

¹In the following, whenever we talk about graphs, epi-graphs we will always assume the space is affine-hull of the domain, which makes arguments easy.

²We commented that this self-duality of set of linear functions with a linear set is special for finite-dim spaces and not true for the infinite dimensional ones.

- More importantly this theorem is giving us a dual definition (defn. in terms of inner-product) of linear functions: linear functions are exactly inner-product forms with one argument fixed. More specifically, $f : L \mapsto \mathbb{R}$ if and only if $\exists l \in L \ni f(x) = \langle l, x \rangle$. With this, one can give many many examples of linear functions in various spaces.
- **Mandatory reading:**
 - Section B.2.8 in Nemirovski [2005];
- Optional reading: sections B.2.1-B.2.3 in Nemirovski [2005] (these were not covered in lectures but very useful to know); relevant parts of sections 17,19,21 in Rockafellar [1996].

Lecture 10

- Once linear functions are studied, affine functions (and results with them) are immediate:
 - $f : A \mapsto \mathbb{R}$, where A is an affine set, is affine function if and only if affine combinations are preserved under f i.e., $\lambda_1, \lambda_2 \in \mathbb{R}, \lambda_1 + \lambda_2 = 1, x_1, x_2 \in L \Rightarrow f(\lambda_1 x_1 + \lambda_2 x_2) = \lambda_1 f(x_1) + \lambda_2 f(x_2)$.
 - Needless to say, all linear functions are affine.
 - Again, f is affine if and only if $graph(f)$ is an affine set of dimensionality same as A .
 - Let L_A be the linear set associated with A . Because of the above result, it turns out that f is affine if and only if there exists a $u \in L_A, b \in \mathbb{R}$ such that $f(x) = \langle u, x \rangle + b$.
- We then noted that the epigraphs of linear and affine functions are halfspaces that pass through and need not pass through origin¹. With this motivation we defined conic functions:
 - $f : K \mapsto \mathbb{R}$, where K is a cone, is called a conic function² if and only if conic combinations are under-estimated under f i.e., $\lambda_1, \lambda_2 \geq 0, x_1, x_2 \in K \Rightarrow f(\lambda_1 x_1 + \lambda_2 x_2) \leq \lambda_1 f(x_1) + \lambda_2 f(x_2)$.
 - We gave many examples: all norms are conic, all semi-norms are conic. We gave examples of conic functions that are not defined on entire V , those that are not even, that whose value can be negative.
 - It was easy to show that f is conic if and only if $epi(f)$ is conic. We say that f is **closed conic** if and only if $epi(f)$ is closed conic set.

¹Also, the level-sets of linear and affine functions are affine sets.

²Rockafellar uses the term support functions instead of conic functions. We will define support functions later (and realize to be same as conic functions).

- Using this, appendix (page 6-7) provides the dual description of closed conic functions: $f : K \mapsto \mathbb{R}$, where K is closed cone, is closed conic if and only if³ $\exists S \subset V$ such that f is maximum⁴ over set of all linear functions defined by S i.e., $f(x) = \max_{y \in S} \langle x, y \rangle \forall x \in K$. Some books use the term “support function of S ” to describe such functions⁵.
- The proof in appendix infact gives the form of S in terms of dual cone of $\text{epi}(f)$. In cases where this dual cone is itself an epigraph⁶ (of some f^*), we call f^* as dual of f . It follows from the derivation in appendix that:

$$f^*(x) = \max_{u \in V} \langle -u, x \rangle, \\ \text{s.t. } f(u) \leq 1,$$

where $x \in \text{dom}(f^*) = \{y \in V \mid \langle y, z \rangle \geq 0 \forall z \ni f(z) = 0\}$

- It is easy to see that dual function of $f(x) = \|x\|_M$ ($M \succ 0$) is $f^*(x) = \|x\|_{M^{-1}}$. Infact, there is a special name for dual of a function that is a norm: it is called [dual norm](#)⁷. Hence, $\|\cdot\|_{M^{-1}}$ is the dual of norm of $\|\cdot\|_M$ and vice-versa.
- Needlessly to say, by the very defn., we have: $f^{**} = f$ (whenever f is closed conic).
- Optional reading: section 13 in Rockafellar [1996].

³Proof in appendix is for the “only if” part. The “if” part is easy to prove and left as exercise.

⁴maximum over functions is defined as point-wise maximum i.e., function value of maximum is equal to maximum of function values.

⁵So our result is that support function concept is same as closed conic function.

⁶We gave examples of cones where this does not happen.

⁷Note that our defn. matches the defn. of dual norm in Boyd and Vandenberghe [2004].

Lecture 11

- We began by trying to show that the function $f : S^n \mapsto \mathbb{R}$ given by $f(M) =$ maximum eigen-value of M is a conic function. We realized that the easiest way of showing this is through the dual defn. of conics i.e., showing that f is infact a support function (of some set S):
 - We first wrote down the following fact (proved in lecture): $f(M) = \max_{\|x\| \leq 1} x^\top Mx$.
 - With the above we have that: f is support function of the set $S = \{X \in S^n \mid X \succeq 0, \text{trace}(X) \leq 1, \text{rank}(X) = 1\}$.
- We then defined the next obvious, the convex functions: $f : C \mapsto \mathbb{R}$, where C is a convex set, is called a convex function if and only if convex combinations are under-estimated under f i.e., $\lambda_1, \lambda_2 \geq 0, \lambda_1 + \lambda_2 = 1, x_1, x_2 \in C \Rightarrow f(\lambda_1 x_1 + \lambda_2 x_2) \leq \lambda_1 f(x_1) + \lambda_2 f(x_2)$.
- It was not difficult to verify that the functions: x^2 , $\|x\|^2$ and $-\log(x)$ are convex. Infact in proving each one of them, we used the AM-GM inequality in some form.
- We then went ahead and said using (non-trivial) induction, one can show with a convex f : $f(\sum_{i=1}^n \lambda_i x_i) \leq \sum_{i=1}^n \lambda_i f(x_i)$, whenever $x_i \in C \forall i$ and $\lambda_i \geq 0 \forall i, \sum_{i=1}^n \lambda_i = 1$.
- Defining a random variable X that takes the value x_i with probability λ_i , the above result says: $f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)]$. A more deeper result is that this inequality known as the **Jensen's inequality** holds for any (perhaps non-discrete) random variable X . We commented that many fundamental inequalities may be derived from this, including the Holder's inequality (refer section 3.1.9 in Boyd and Vandenberghe [2004]). From Holder's inequality it follows that dual of $f(x) = \|x\|_p (p \geq 1)$ is $f^*(x) = \|x\|_q$ where $\frac{1}{p} + \frac{1}{q} = 1$.

- It was again easy to show $f : C \mapsto \mathbb{R}$ is convex function if and only if $\text{epi}(f)$ is convex.
- We noted examples of convex functions whose epigraphs are not closed: $f(x) = 1, x \in (-1, 1)$ (domain is open convex set); $g(x) = 1$ if $x \in (-1, 1)$ and $g(x) = 2$ if $x \in \{-1, 1\}$ (domain is closed convex set). Since we give dual defns. for closed convex sets, we do so for functions too... We define closed convex functions: f is closed convex iff $\text{epi}(f)$ is closed convex. Needless to say, the domain of closed convex functions must be themselves closed.
- Sections 3.1.1, 3.1.7,3.1.8,3.1.9 in Boyd and Vandenberghe [2004]; C.1 in Nemirovski [2005].
- Optional reading: relevant parts in section 4 in Rockafellar [1996].

Lecture 12

- We began by realizing the dual defn. for closed convex functions. The derivation is on page 8 of Appendix.
- Motivated by the dual form, we defined the notion of conjugate: given any function $f : V \mapsto \mathbb{R}$, the conjugate of f is defined as $f'(x) = \max_{y \in V} \langle x, y \rangle - f(y)$. Other names for conjugate are: Fenchel dual, Fenchel conjugate, Legendre transform¹.
- Pages 9,10 of appendix provide a proof for the following result: Closed convex functions are exactly those where $f'' = f$.
- We showed that for $f(x) = \frac{1}{2}x^\top Qx$ (where $Q \succ 0$), $f'(x) = \frac{1}{2}x^\top Q^{-1}x$.
- We noted that global properties of f are local properties of f' and vice-versa. As an example we said, $f'(0) = -\min_{x \in V} f(x)$.
- We then noted the fact that for a closed convex function² $f : C \mapsto \mathbb{R}$, at any point $x_0 \in C$, there exists a supporting hyperplane i.e., $\exists(u_{x_0}, v_{x_0}) \in V \times \mathbb{R} \ni \langle (u_{x_0}, v_{x_0}), (x - x_0, y - f(x_0)) \rangle \leq 0 \forall (x, y) \in \text{epi}(f)$. One way of proving this is to simply use the generic separation theorem (theorem B.2.5 in Nemirovski [2005]), that gives that the $\text{epi}(f)$ and the point $(x_0, f(x_0))$ are non-strictly separable. The other way, which I prefer, is to recall from the previous conjugacy proof that $f(x) = \max_{(y, \alpha) \in \text{epi}(f^*)} \langle x, y \rangle - \alpha$. The fact that there is a supporting hyperplane at $(x_0, f(x_0))$ is proved if we prove that this supremum is attained at some point (y^*, α^*) in the (closed set) $\text{epi}(f^*)$ for $x = x_0$. This is easy to prove: now there must exist a sequence $(y_n, \alpha_n) \subset \text{epi}(f^*)$

¹One can extend the definition to functions that are limited to some domain that is not the entire set of vectors by following the trick suggested in section 3.1.2 in Boyd and Vandenberghe [2004]. In following many times we simplify proofs etc. by assuming functions are always defined over entire set of vectors.

²To simplify arguments lets assume $\dim(C) = \dim(V)$ unless mentioned otherwise.

such that $\langle (x_0, -1), (y_n, \alpha_n) \rangle \rightarrow f(x_0)$ (by the very definition of supremum). Also, it is easy to see that this sequence itself is bounded and hence by Heine Borel theorem there must exist a sub-sequence that converges and since $\text{epi}(f^*)$ is closed, it must converge in the set³. In other words, the supremum is achieved⁴.

- It was easy to see that $v_{x_0} \leq 0$ (simply use the supporting hyperplane's inequality at points (x_0, y) where $y \geq f(x_0)$).
- We then noted that if $x_0 \in \text{int}(C)$ ⁵, then $v_{x_0} \neq 0$. This is easy to prove: since $x_0 \in \text{int}(C)$, $x_0 + \rho u_{x_0} \in C$ for some small enough $\rho > 0$. In case $v = 0$, the supporting hyperplane inequality says that $\rho \langle u_{x_0}, u_{x_0} \rangle \leq 0$, which is impossible. Hence dividing the whole inequality by v_{x_0} and re-writing it at the point $(x, y) = (x, f(x))$ gives: $f(x) \geq f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle$, where $\nabla f(x_0) \equiv -\frac{u_{x_0}}{v_{x_0}} \in V$.
- The above inequality says that at every $x_0 \in \text{int}(C)$, there is an affine function that is always below the given function and is equal to it at $(x_0, f(x_0))$.
- This motivates the following definition: for any function $f : S \subset V \mapsto \mathbb{R}$ and $x_0 \in S$, if there exists a $\nabla f(x_0) \in V$ such that $f(x) \geq f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle$, then f is said to be **sub-differentiable**⁶ at x_0 . The vector $\nabla f(x_0) \in V$ is called **the sub-gradient**. The inequality above is called the **sub-gradient inequality of f at x_0** .
- Needless to say, there might be many vectors satisfying the sub-gradient inequality⁷ of f at x_0 . The set of all such sub-gradients at a point is called the **sub-differential of f at x_0 : $\partial f(x_0)$** . It is easy to see that the sub-differential set of any function at a point is either empty or a convex, bounded set.
- The discussion above basically shows that a closed convex function is sub-differentiable at all (relatively) interior points of the domain⁸.

³This limit provides the required (u_{x_0}, v_{x_0})

⁴This argument is analogous to the one we made to prove that projection of a point onto a closed set always exists.

⁵remember we are always concerned with the “relative” interior i.e., the space is restricted to affine-hull of C , which with our assumption is V

⁶Since the RHS expression is analogous to Taylors series expanded to first term, this name involves the term “differentiable”.

⁷We gave examples of such cases in lecture. Atleast from the examples, we were hinting at “non-differentiability” being the cause for this multiplicity.

⁸In fact, this is true for any convex function and is easy to prove based on above discussion. Please refer prop. C.6.5 in Nemirovski [2005] for details.

- The converse, if a function defined on a convex set is sub-differentiable everywhere, then it is convex, is also true and is easy to prove.
- Hence for everywhere sub-differentiable functions, satisfying sub-gradient inequality at all points is the definition of a convex function.
- We showed that for the convex (infact, conic) function $f(x) = \|x\|$, $\nabla f(x_0) = \frac{x_0}{\|x_0\|}$ for $x_0 \neq 0$ and any vector in unit sphere around origin is a sub-gradient of f at 0.
- Analogous to the case of convex sets, we said the following result is true: a function $f : C \mapsto \mathbb{R}$ is convex if and only if its restriction to any line in C is convex i.e., for any $x_1, x_2 \in C$, the function $g : [0, 1] \mapsto \mathbb{R}$ given by $g(t) = f(tx_1 + (1 - t)x_2)$ is convex.
- We then recalled the college-day result that a function double-differentiable function $f : (a, b) \mapsto \mathbb{R}$ is convex if and only if $\frac{d^2f(x)}{dx^2} \geq 0 \forall x \in (a, b)$ (**Prop. C.2.1 in Nemirovski [2005]**). Also, one can show that if $f : [a, b] \mapsto \mathbb{R}$ is continuous everywhere and convex on (a, b) , then f is convex on $[a, b]$.
- We then argued that the above two results i) convexity is essentially a 1-d concept ii) non-negative double derivative defines convexity in 1-d, when used in a combination can turn out to be a powerful tool in proving convexity of functions. As an example we showed that $f(x) = x^\top Ax$ (where A is symmetric) is convex if and only if $A \succeq 0$.
- Also note that, level-sets of convex function are convex. Hence the above double differential criteria may sometimes help proving convexity of some sets.
- **Mandatory reading: sections C.6.3, C.6.2, C.2.2 in Nemirovski [2005]; section 3.3 in Boyd and Vandenberghe [2004]**
- **Optional reading: section 12, 26, 23 in Rockafellar [1996].**

Lecture 13

- In order to provide a first order (and later second order) characterization of convexity, we began by extending the notion of differentiability for real-valued functions defined over Hilbert spaces.
- Since the notion of continuity is more basic than differentiability, we spent few minutes on it: a function $f : S \subset V \mapsto \mathbb{R}$ is said to be **continuous**¹ at $x_0 \in S$ iff for every sequence $\{x_n\} \subset S \rightarrow x_0$ we have $\{f(x_n)\} \rightarrow f(x_0)$. We commented that convex functions have elegant continuity² properties. **Please study section C.4 in Nemirovski [2005].**
- We then recalled the notion of differentiability for $f : \mathbb{R} \mapsto \mathbb{R}$ as: f is said to be **differentiable** at a point $x_0 \in \mathbb{R}$ iff the instantaneous slope given by $\lim_{h \rightarrow 0} \frac{f(x_0+h)-f(x_0)}{h}$ exists. In such a case the limit (instantaneous slope) is called as the derivative of f at x_0 .
- We also recalled that the above limit exists iff the left limit, $\lim_{h \uparrow 0} \frac{f(x_0+h)-f(x_0)}{h}$, known as the left derivative is equal to the right limit, $\lim_{h \downarrow 0} \frac{f(x_0+h)-f(x_0)}{h}$, known as the right derivative.
- We then looked at $f : V \mapsto \mathbb{R}$ and wanted to extend the notion of differentiability. We said one way is to look at instantaneous slope, but now we can have one such slope in each direction $u \in V$! We gave the following definition: the **directional derivative** of f at $x_0 \in V$ along the vector (direction) $u \in V$ is given by $D_f(x_0; u) \equiv \lim_{h \rightarrow 0} \frac{f(x_0+hu)-f(x_0)}{h}$, provided the limit exists. One can also define the corresponding left, right limits and we call them as the left directional derivative $D_f^-(x_0; u)$ and the right directional derivative $D_f^+(x_0; u)$.
- We then said, atleast in case of functions on \mathbb{R} , the most useful result³ involv-

¹This is a simple extension of definition of continuity to functions on Hilbert spaces.

²All convex functions are continuous in the (relative) interior and infact, they are locally Lipschitz continuous.

³Taylor series result.

ing derivative/differentiability is that there exists an affine function that well approximates the function f in a neighborhood of x_0 . This was immediate by re-writing $\frac{df(x_0)}{dx} = \lim_{h \rightarrow 0} \frac{f(x_0+h)-f(x_0)}{h}$ as $\lim_{x \rightarrow x_0} \frac{f(x)-f(x_0)-\frac{df(x_0)}{dx}(x-x_0)}{|x-x_0|} = 0$, which is hence an alternate equivalent definition of differentiability in case of functions on \mathbb{R} .

- We then tried to extend this notion of differentiability to functions on Hilbert spaces: let $S \subset V$ and $f : S \mapsto \mathbb{R}$ be given. Then f is said to be **differentiable** at $x_0 \in S$ iff there exists a affine function $A(x)$ such that $A(x_0) = f(x_0)$ and $\lim_{x \rightarrow x_0} \frac{f(x)-A(x)}{\|x-x_0\|} = 0$ ⁴. Since we know that general form of affine functions is linear function + constant number, we know that existence of such $A(x)$ is same⁵ as existence of $g_f(x_0) \in V$, called the **gradient** of f at x_0 , such that $\lim_{x \rightarrow x_0} \frac{f(x)-(f(x_0)+\langle g_f(x_0), x-x_0 \rangle)}{\|x-x_0\|} = 0$.
- It is easy to see that in case $V = \mathbb{R}$, the gradient notion is exactly same as that of derivative. We next wanted to see if some relation exists, in general (i.e., if $V \neq \mathbb{R}$), between directional derivative and gradient.
- Interestingly, it turns out that $\langle g_f(x_0), u \rangle = D_f(x_0; u)$ for any $u \in V$ and any $x_0 \in S$, where the function is differentiable. In particular, if the function is differentiable at a point, then all directional derivatives exist at that point. This relation between gradient and directional derivative was simple to prove, but nevertheless as we will see in subsequent lectures, an extremely useful result for convex optimization.
- The next immediate question was: is it true that if a function has all directional derivative at a point, then it is differentiable⁶? Unfortunately, the answer is NO and http://people.whitman.edu/~hundredr/courses/M225/Ch14/Example_DirectionalDeriv.pdf provides a counter example.
- In case $V = \mathbb{R}^n$, taking u as unit vectors with all entries zero except one entry, we realized that the gradient vector is simply the Euclidean vector with entries as the partial derivatives. Infact, this trick can be repeated for any function on $V(\neq \mathbb{R}^n)$ by considering an appropriate orthonormal basis and is most of the time, the easiest way to compute gradients⁷. Using this we computed gradient of $f(X) = \|X\|_F^2$ at X_0 as $2X_0$.

⁴The conditions basically say that the $f(x)$ is well approximated by $A(x)$ in a neighborhood of x_0 .

⁵General form of affine $A(x)$ such that $A(x_0) = f(x_0)$ is $A(x) = f(x_0) + \langle g_f(x_0), x - x_0 \rangle$.

⁶Note that this was true for all functions on \mathbb{R} .

⁷In next lecture we will hint of other ways for computing gradients.

Lecture 14

- We began by noting that the equality $\langle g_f(x_0), u \rangle = D_f(x_0; u)$ implies that the instantaneous direction of maximum increase of the function at x_0 is $g_f(x_0)$ and moreover, no other direction can match the increase. We said that this motivates greedy algorithms like gradient descent (which we will study later) for optimization.
- By motivating examples we conjunctured the following two important results:
 1. Let $f : S \subset V \mapsto \mathbb{R}$ be a convex function that is sub-differentiable at x_0 . Then f is differentiable at x_0 if and only if the sub-differential set at x_0 has a single element. Infact in this case, $g_f(x_0)$ is equal to this unique sub-gradient. The proof for “only if” part is easy¹: by the subgradient inequality, for any sub-gradient $\nabla f(x_0)$, any $u \in V$ and any $h > 0$, we must have $\langle \nabla f(x_0), u \rangle \leq \frac{f(x_0+hu)-f(x_0)}{h}$. Taking limits, this implies that $\langle \nabla f(x_0), u \rangle \leq D_f^+(x_0; u) = D_f(x_0; u) = \langle g_f(x_0), u \rangle$. Since this is true for all $u \in V$, we must have that $\nabla f(x_0) = g_f(x_0)$. For proof of the “if” part, please refer theorem 25.1 in Rockafellar [1996].
 2. This result talks about whether the direction of sub-gradient instantaneously increases the function or not: for this, we define the notion of **tangent cone** of a set S at a point $s \in S$, denoted by $\mathcal{T}_S(s)$, as $\{u \mid \exists h > 0 \ni s + hu \in S\}$. It is easy to see that this set is indeed a cone, provided the set S is convex. We call the dual cone of the tangent cone as the **normal cone**, denoted by $\mathcal{N}_S(s)$. Let $f : C \mapsto \mathbb{R}$ be a convex function and $x_0 \in \text{int}(C)$ (relative interior). Assume that $\partial f(x_0) \neq \{0\}$ (i.e., sub-differential set is not uniquely zero). Then², $\mathcal{N}_{L_f(f(x_0))}(x_0) = -\text{CONIC}(\partial f(x_0))$. In lecture, we proved the easy part: $\mathcal{T}_{L_f(f(x_0))}(x_0) \subset -(\partial f(x_0))^*$.

¹Refer prop. C.6.5 in Nemirovski [2005] for an alternate proof.

²Bonus 5marks to the student who proves this conjuncture.

- One can write down the following from result 1 above:
 - For convex functions, we have a new definition of differentiability: sub-gradient being unique. SO one way to compute gradient is, guess a sub-gradient and prove it is the only vector that enables satisfaction of sub-gradient inequality. Then this sub-gradient must be the gradient. In particular, we now know that $f(x) = \|x\|$ is not differentiable at 0.
 - Another definition that follows (from midsem question) is: a closed convex function f is differentiable at x_0 if and only if the set $\arg \max_{y \in V} \langle y, x_0 \rangle - f'(y)$ has a single element and in this case, that element is the gradient.
 - In particular, it is clear that gradient is a sub-gradient and satisfies the sub-gradient inequality. This gives a new definition of convexity: for everywhere differentiable functions, convexity is same as satisfaction of sub-gradient inequality by the gradient.
- One can write down the following from result 2 above:
 - A convex function is differentiable with non-zero gradient at x_0 if and only if, the normal cone $\mathcal{N}_{L_f(f(x_0))}(x_0)$ has a single element.
 - A sub-gradient direction need not be a direction of instantaneous increase of a function. But there will exist some sub-gradient that is a direction of instantaneous increase of the function.
- Owing to the above, from now onwards we will denote a gradient also by $\nabla f(x_0)$. From the situation it must be understood whether ∇f represents gradient or sub-gradient.
- We then defined the notion of twice-differentiability³: A function $f : S \subset \mathbb{R}^n \mapsto \mathbb{R}$ is said to be **twice-differentiable** at $x_0 \in S$ iff there exists a $H(x_0) \in S^n$, called the **Hessian**, such that $\lim_{x \rightarrow x_0} \frac{f(x) - (f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle + \frac{1}{2}(x - x_0)^\top H(x_0)(x - x_0))}{\|x - x_0\|^2} = 0$.
- Taking limit in the above such that $x = x_0 + hu, h \rightarrow 0$ gives: $\frac{1}{2}u^\top H(x_0)u = \lim_{h \rightarrow 0} \frac{f(x_0 + hu) - f(x_0) - h\langle \nabla f(x_0), u \rangle}{h^2}$. This gives: i) The diagonal entries of Hessian are simply double partial derivatives wrt. each co-ordinate. One can prove that the Hessian is infact the matrix all possible double partial derivatives and this proof is beyond the scope of this course. ii) a twice-differentiable everywhere f is convex if and only if the Hessian is psd at all interior points.
- **Mandatory reading: sections C.3, C.2.2 in Nemirovski [2005]; all sections and especially 3.1.3 and 3.1.4 in Boyd and Vandenberghe [2004].**

³In lecture we noted how to extend defn. to $V \neq \mathbb{R}^n$.

- Optional reading: relevant parts in sections 23-25 in Rockafellar [1996].

Lecture 15

- After completing the topic of Convex Analysis, now that we understand each component of an MP, we wish to study a special class of MPs known as Convex Programs. This marks the beginning of Convex Optimization theory. We recalled some definitions:
- **Mathematical Program (MP)**: A symbol that is of the following form¹:

$$(15.1) \quad \begin{array}{ll} \min_{x \in \mathcal{X}} & \mathcal{O}(x) \\ \text{s.t.} & x \in \mathcal{F} \end{array}$$

- More specifically, if \mathcal{F} is written as $\{x \mid g_i(x) \leq 0, \forall i = 1, \dots, m\}$, then the MP is called a **Ordinary Mathematical Program (OMP)**:

$$(15.2) \quad \begin{array}{ll} \min_{x \in \mathcal{X}} & \mathcal{O}(x) \\ \text{s.t.} & g_i(x) \leq 0, \forall i = 1, \dots, m \end{array}$$

- A MP/OMP is said to be **un-constrained** iff $\mathcal{F} = V$.
- We then define the following for a given MP/OMP:
 1. **Domain** — \mathcal{X} . This is the domain for the functions \mathcal{O} and all g_i i.e., $\mathcal{O} : \mathcal{X} \mapsto \mathbb{R}$ and $g_i : \mathcal{X} \mapsto \mathbb{R}$.
 2. **Feasibility set** — $\mathcal{F} = \{x \in \mathcal{X} \mid g_i(x) \leq 0 \forall i = 1, \dots, m\}$. An element of this set is called as a **Feasible solution**. If this set is non-empty, then we say that the MP/OMP is **feasible**.

¹We assume $\mathcal{F} \subset \mathcal{X}$.

3. **Objective function** — \mathcal{O} .
4. **Optimal value or (simply) value** of the MP/OMP is defined as $\inf (\{f(x) \mid x \in \mathcal{F}\})$ whenever the program is feasible. The value is defined as ∞ in case the program is infeasible. The program is said to be **bounded** iff its value is not $-\infty$.
5. **Optimal solution or (simply) solution** of the MP/OMP is defined as that $x^* \in \mathcal{F}$ such that $f(x^*) \leq f(x) \forall x \in \mathcal{F}$. The set of all optimal solutions is called as the **optimal set or solution set** and is denoted by:

$$(15.3) \quad \begin{array}{ll} \arg \min_{x \in \mathcal{X}} & \mathcal{O}(x) \\ \text{s.t.} & x \in \mathcal{F} \end{array}$$

or

$$(15.4) \quad \begin{array}{ll} \arg \min_{x \in \mathcal{X}} & \mathcal{O}(x) \\ \text{s.t.} & g_i(x) \leq 0, \forall i = 1, \dots, m \end{array}$$

The program is said to be **solvable** iff the solution set is non-empty.

- A MP is said to be a **Convex Program (CP)** iff \mathcal{O}, \mathcal{F} are convex.
- An OMP is said to be an **Ordinary Convex Program (OCP)** iff $\mathcal{O}, g_i \forall i = 1, \dots, m$ are convex.
- Analogous defns. can be given for maximization problems.
- We then went through the five example MPs shown in the first lecture and determined whether they are convex or not. This provided us with examples and non-examples of CPs².
- The most fundamental questions regarding an MP ofcourse are “when is an MP bounded/feasible³/solvable/uniquely-solvable and can we characterize the solution set?”. It turns out that the answers (atleast partially) to these questions are neat and elegant for CPs.
- In particular, we answer the following four fundamental questions:

²Bonus marks to that student who shows that the first example can be written as a CP. You just need to complete the proceeding in the lecture.

³We will later see that feasibility and boundedness turn out to be some kind of duals and hence dealing with one of them is enough.

1. Arrive at suff. cond. for boundedness of CP. We conjectured that a CP with bounded feasibility set is itself bounded. The proof follows from i) the sub-gradient inequality⁴, which shows that a convex function is always greater than or equal to an affine function, and ii) linear functions are bounded on bounded sets.
2. Arrive at suff. cond. for solvability of CP. We conjectured that a CP is solvable if the Feasibility set is compact and the objective is closed convex (or continuous). Refer pg 11 of Appendix for a sketch of proof of this.
3. Arrive at suff.cond. for unique solvability.
4. Arrive at nec. and suff. cond. for an optimal solution. This will help in characterizing the optimal set. This, we commented, is the most useful result of this course.

⁴Recall that every convex function is sub-differentiable at a (relatively) interior point and such a point in turn exists for any convex set. (prop. C.6.5 and theorem B.1.1 in Nemirovski [2005]).

Lecture 16

- We began by noting that a convex program with closed strictly convex objective and a compact feasibility set is uniquely solvable. Please refer pg 12 of appendix for a proof.
- We then wanted to characterize the optimal set of a CP using something that is “easy to compute/available” while solving programs. We recalled one characterization already proved in practice problems: $\partial f'(0) = \arg \min_{x \in V} f(x)$. In particular, this equality shows that the optimal set for any¹ CP is either empty or closed convex. This is because the sub-differential set is closed convex (or empty). However, we commented that this is NOT a very useful characterization in practice as conjugate may not be available and perhaps an optimization problem is need to be solved to obtain it.
- We then note the following extremely useful (and trivial to prove) theorem:

Theorem 16.0.3. *Given an unconstrained CP, x^* is an optimal solution if and only if $0 \in \partial \mathcal{O}(x^*)$. In other words, the optimal set is $\{v \in V \mid 0 \in \partial \mathcal{O}(v)\}$. In particular, this gives a characterization for solvability of unconstrained convex programs: an unconstrained CP is solvable if and only if there exists atleast one $v \in V$ such that 0 is a sub-gradient of the objective at that v .*

- We then considered applications of this theorem:
 1. Consider $\min_{x \in \mathbb{R}^n} \frac{1}{2}x^\top Px + q^\top x + r$, where $P \succeq 0$. This is an example of an unconstrained convex program and hence the above theorem applies: x^* is optimal if and only if $Px^* + q = 0$. From linear algebra, we have that such an x^* exists if and only if $q \in \mathcal{C}(P)$, the column space of P . And in such a case, one can perform Gaussian elimination to arrive at x^* . In particular, if $P \succ 0$, we have that the CP is uniquely solvable and

¹One can include programs that are constrained by defining the function to be ∞ outside the feasibility set.

$x^* = -P^{-1}q$. Also, in this case, it is easy to see that the optimal value of the program is $-\frac{1}{2}q^\top P^{-1}q + r$.

2. We proved the Schur complement lemma using this theorem.

- We then conjectured and proved the following theorem:

Theorem 16.0.4. *Given a [differential CP](#)², x^* is optimal if and only if $x^* \in \mathcal{F}$ and $\nabla f(x^*) \in \mathcal{N}_{\mathcal{F}}(x^*)$.*

Please refer page 14 in the appendix for a proof.

- We commented that this is a fundamental theorem and extremely useful in convex optimization theory. Infact we said that we will derive many (perhaps) familiar optimality conditions using this theorem in the subsequent lectures.
- For the sake of completeness, we wish to write an analogous theorem for convex programs that need not be differentiable (needlessly to say, sub-differentiability at all feasible points is assumed): given a feasible solution x^* of CP, x^* is an optimal solution of (P) if and only if for every $u \in T_{\mathcal{F}}(x^*)$ we can identify a $\nabla f(x^*) \in \partial f(x^*)$ such that $\langle \nabla f(x^*), u \rangle \geq 0$. The proof follows from a characterization of the support function of the sub-differential set in terms of one-sided directional derivative. Interested students are requested to read section 23 from Rockafellar [1996] and specifically theorem 23.4 in it. Though more general, we won't further use this result as it is not elegant. Henceforth we will focus on optimality conditions for various differentiable CPs.
- **Mandatory reading:** [Section C.5 in Nemirovski \[2005\]](#).
- **Optional reading:** section 27 in Rockafellar [1996]

²A CP with a differentiable (everywhere in the domain) objective.

Lecture 17

- We began by noting that in the theorem 16.0.4, the gradient is more familiar to us, even in the sense of computing it¹, than the other ingredient namely, normal cone. Hence we began studying normal cones for some special sets.
- We began with the case of an open feasibility set. Here, it is easy to see that the tangent cone at any point (since every point is in the interior) contains all directions and hence the normal cone has only the 0 vector. Hence in this case, again the gradient being zero characterizes optimality. Infact, we wrote down the following:

Corollary 17.0.5. *Let program (CP) be differential and $x^* \in \text{int}(\mathcal{F})$. Then, x^* is optimal if and only if $\nabla f(x^*) = 0$.*

- We then looked at cases where the constraints are linear inequalities (i.e., define a polyhedral set), called **Polyhedrally Constrained Convex Program (PCCP)**:

$$(17.1) \quad \begin{array}{ll} \min_{x \in \mathcal{X}} & \mathcal{O}(x) \\ \text{s.t.} & \langle a_i, x \rangle \leq b_i, \forall i = 1, \dots, m \end{array}$$

and conjunctured the following result:

Corollary 17.0.6. *Let (CP) be a differential PCCP with domain \mathcal{X} being an open set. Then x^* is optimal if and only if:*

1. $x^* \in \mathcal{F} = \{x \in \mathcal{X} \mid \langle a_i, x \rangle \leq b_i, \forall i = 1, \dots, m\}$ (*Feasibility cond.*).
2. \exists a $\lambda^* \in \mathbb{R}^m$ such that:
 - (a) $\lambda^* \geq 0$ (*Non-negativity cond.*).
 - (b) $\nabla f(x^*) + \sum_{i=1}^m \lambda_i^* a_i = 0$ (*Gradient cond.*).

¹For e.g. we know derivatives of various standard functions etc. Moreover, are familiar with finite difference methods for numerical differentiation etc.

$$(c) \lambda_i^* (\langle a_i, x^* \rangle - b_i) = 0 \quad \forall i = 1, \dots, m \quad (\text{Complementary-slackness cond.}).$$

Another interesting way of writing the same result is: (CP) is solvable if and only if there exists a point $(x^*, \lambda^*) \in \mathcal{X} \times \mathbb{R}^m$ such that the above 4 conditions, called hence-forth as the **KKT conditions** hold. For convenience, we will call a point satisfying these 4 conditions as a **KKT point**.

- While the proof of this theorem was done in lecture, here we will use the proof from the subsequent lecture that generalizes this corollary (and not repeated here). What is more important is the following comments:
 - The KKT conditions allow us to pose the problem of verifying whether a x^* is optimal, in terms of verifying solvability/feasibility/consistency of a set of linear inequalities (in terms of λ). This is elegant, as linear inequalities is a well-understood subject.
 - Through the example of maximizing entropy of a distribution with given moments (and an eg. of a simple MP involving matrices), we made the point that KKT conditions help in realizing profound facts about the optimal solution (and sometimes even the analytical expression for optimal set), without the need to have absolutely any knowledge about the domain from which the MP originated! This was cool. For e.g., we obtained the famous thermodynamics relation that $p(E) \propto \exp \frac{-E}{kT}$ by simply writing down the KKT conditions.
 - We argued that in case there are linear equalities instead of linear inequalities, everything in the corollary will remain same except that the non-negativity conditions on the corresponding λ_i s should not be insisted upon.
 - We defined a PCCP with a linear objective as a **Linear Program (LP)**. So the above corollary applies to any LP with an open domain.
 - We defined a PCCP with a convex quadratic objective as a **Quadratic Program (QP)**. So the above corollary applies to any QP with an open domain.
 - Suppose it is known that a diff. PCCP (with open domain) is solvable. Since it is then guaranteed that a KKT point (x^*, λ^*) exists, the first part of it, which is x^* , we know gives an optimal solution for the PCCP. The other part, which is λ^* , infact gives an idea about activity of the constraints: we say a constraint $g(x) \leq 0$ is **active** at x^* iff $g(x^*) = 0$. It is easy to see that if $\lambda_i^* > 0$ (for some i), then by the complementary slackness cond., we have that the constraint $\langle a_i, x \rangle \leq b_i$ is active at x^* .

And conversely, if the i^{th} constraint is not active at x^* , then $\lambda_i^* = 0$. In fact, because of this complementarity relation, the name complementary slackness.

Lecture 18

- We began proving the following corollary that characterizes optimality in case of differentiable OCPs:

Corollary 18.0.7. *Let (15.2) be a regular convex program: an ordinary convex program with differentiable \mathcal{O} , g_i defined over an open domain satisfying the Slater's condition — there exists a $x_0 \in \mathcal{X} \ni g_i(x_0) < 0$ for all i such that g_i is a non-affine function. Then the following statements are true:*

– x^* is optimal iff:

1. $x^* \in \mathcal{X}, g_i(x^*) \leq 0 \forall i = 1, \dots, m.$ (feasibility cond.)

2. $\exists \lambda^* \in \mathbb{R}^m \ni$

(a) $\lambda^* \geq 0$ (non-negativity cond.)

(b) $\nabla f(x^*) + \sum_{i=1}^m \lambda_i^* \nabla g_i(x^*) = 0$ (gradient cond.)

(c) $\lambda_i^* g_i(x^*) = 0 \forall i = 1, \dots, m$ (complementary slackness cond.)

(These conditions are hence-forth referred to as the KKT conditions and they generalize the previous cases).

– The regular CP is solvable iff there exists a KKT point i.e., a point $(x^*, \lambda^*) \in \mathcal{X} \times \mathbb{R}^m$ satisfying the above KKT conditions.

Please refer Appendix I in the previous year's notes for a proof.

- We illustrated the power of the KKT conditions using two examples:
 - Problem 1 in Quiz2 of the previous year. This was an elegant case KKT conditions infact gave an analytical expression for the optimal solution.
 - Problem 5(c) in section 2.1 in previous year problem-set. Here KKT conditions were used to compare and related the optimal solutions and hence the optimal values of two different programs.
- The following comments are noteworthy:

- Apart from differentiability assumptions, a regular convex program is as general as any convex program and hence this elegant result (KKT conditions) is also generic and hence extremely useful. Infact, these provided necessary (but not sufficient) conditions for optimality for regular programs (refer theorem 9.1.1 in Fletcher [2000]).
 - Historically, Kuhn-Tucker published this result in 1950s; later it was discovered that Karush wrote the same in his MSc thesis in 1930s, which was unpublished. So now the conditions are named after all the three: Karush-Kuhn-Tucker (KKT) conditions.
 - The books (see below) provided alternate proofs of the same corollary. It is worth going through them.
- There is one more elegant case where the optimality conditions are simple and infact look very similar to those in case of Linear Programs. This is the case of [conic program](#)¹:

$$(18.1) \quad \begin{array}{ll} \min_{x \in \mathbb{R}^n} & c^\top x, \\ \text{s.t.} & b - A^\top x \in K. \end{array}$$

Here $K \subset \mathbb{R}^m$ is a cone. Please refer to the problem in Quiz-2 and its solution for the derivation of KKT in the case of conic programs:

Corollary 18.0.8. *Let (18.1) be a conic program. Then,*

– x^* is optimal iff:

1. $x^* \in \mathcal{F} = \{x \in \mathbb{R}^n \mid b - A^\top x \in K\}$. (feasibility cond.)

2. $\exists \lambda^* \in \mathbb{R}^m \ni$

(a) $\lambda^* \in K^*$ (non-negativity cond.)

(b) $c + A\lambda^* = 0$ (gradient cond.)

(c) $(b - A^\top x^*)^\top \lambda^* = 0$ (complementary slackness cond.)

(These conditions are hence-forth referred to as the KKT conditions and they generalize the Linear Program case.).

- **Mandatory Reading:** Section D.2.3.B in Nemirovski [2005]; Sections 5.5.3 and 5.5.4 (more importantly, the solved examples) in Boyd and Vandenberghe [2004].
- **Optional Reading:** Section 28 in Rockafellar [1996].

¹We will first define conic programs with Euclidean space as the domain and later on talk about other extensions.

Lecture 19

- Analogous to duality in sets, functions, we wanted a notion of duality for MPs.
- Recall that the key thought process in defining duals previously was:
 1. Is there an outer/sculptural description of sets? In case of MPs, since one of the most important aspects is the optimal value a natural inner and outer description of this number is “it is the infimum of all its upper bounds” (this is like “inner” description since every feasible point of the given MP will provide an upper bound; moreover, the given minimization is simply finding this least upper bound) and “it is the supremum of all its lower bounds” (this is like “outer” description).
 2. Procedurally, recall that in order to answer above question, we defined dual quantities like: orthogonal complement, dual cone, polar set, dual and conjugate functions. Hence in case of MPs we are looking for a maximization problem, which at every feasible point gives a lower bound and in fact reaches the optimal value of the given (primal) MP at optimality.
 3. Additionally, recall that while defining dual cone/polar-set etc. we were careful so that these quantities were nice (i.e., conic/convex even if the original sets/functions were not conic/convex). Similarly for MPs, we would like the dual program to be a convex one (so that the dual at least does not pose computational challenges).
- In summary, we basically motivated the following definition of dual for any given MP (we hence-forth call the given MP as the **primal**): given a primal program, dual of it is any program that satisfies the following properties:
 1. At feasible point of the dual, its objective value is a lower bound to the optimal value of the primal. This is known as the principle of **weak duality**.
 2. The optimal values of the primal and dual are the same. This is called as the principle of **strong duality**.

3. The dual is itself a convex program.

- Given this definition we were in search of schemes/procedures for systematically writing down duals. We started with the case of DPCCPs (with open domain), where we already characterized optimality:

Theorem 19.0.9. *Let (P), the following program, be a differentiable PCCP with an open domain:*

$$(19.1) \quad \begin{aligned} \min_{x \in \mathcal{X}} \quad & f(x), \\ \text{s.t.} \quad & \langle a_i, x \rangle \leq b_i \quad \forall i = 1, \dots, m. \end{aligned}$$

Assume that (P) is solvable. Then the following is a dual of (P):

$$(19.2) \quad \max_{\lambda \in \mathbb{R}_+^m} \quad - \sum_{i=1}^m \lambda_i b_i - f' \left(- \sum_{i=1}^m \lambda_i a_i \right),$$

where f' is the conjugate of f with domain \mathcal{X} .

The proof of this is given in appendix pages 15-17.

- Infact, even if (P) were not convex (i.e., f is arbitrary), one can still write down the program (19.2) and it is easy to see that in this case weak duality holds (i.e., (P) \leq (19.2)) and (19.2) is convex. Hence the above DCCP dual is useful even for polyhedrally constrained non-convex programs for obtaining a lower bound.
- We took the example of a strictly convex QP and wrote its dual using the above technique. Interestingly, the dual of this was again a strictly convex QP. We call such programs as self-dual.
- From the above we have the following LP duality corollary¹:

Corollary 19.0.10. *Let (P), the following program, be an LP:*

$$(19.3) \quad \begin{aligned} \min_{x \in V} \quad & \langle c, x \rangle, \\ \text{s.t.} \quad & \langle a_i, x \rangle \leq b_i \quad \forall i = 1, \dots, m. \end{aligned}$$

Assume that (P) is solvable. Then the following is a dual of (P):

$$(19.4) \quad \begin{aligned} \max_{\lambda \in \mathbb{R}^m} \quad & - \sum_{i=1}^m \lambda_i b_i, \\ \text{s.t.} \quad & \lambda \geq 0, \quad \sum_{i=1}^m \lambda_i a_i + c = 0, \end{aligned}$$

¹proof follows from computing the conjugate of the linear function in the objective.

- Infact, theorem 1.2.2 in Nemirovski [2005] provides a stronger version of this theorem. Please refer to it and its proof.
- Looking at LP dual we made few comments about duals (which are true for any dual):
 1. The space of primal variables need not be the same of dual variables. For e.g. in the case of LP, the primal variable lives in some vector space \mathcal{V} , whereas the dual variables live in \mathbb{R}^m . We commented that this can sometimes be put to use for computational/programming ease.
 2. The dimensionality of primal variable determines the number of constraints in the dual and vice-versa. Hence there is always a trade-off in using the forms. Some algorithms might suit primal form and some dual form. Infact, some of the state-of-the-art algorithms maintain both primal and dual solutions in order to make the best use of this trade-off.
 3. Needless to say, dual gives insight into the optimal solution. Next lecture we will illustrate this using an example.
 4. Most importantly, dual often provides a different view of the problem. And such an insight could motivate further efficient algorithms. This we said is the most important use of writing duals — to view the same optimization problem in different ways. We will illustrate this also with an example.
- **Mandatory Reading:** sections 1.1-1.2 in Nemirovski [2005].

Lecture 20

- We took the program that defines support function of 1-norm ball (i.e., ∞ -norm's dual form). Wrote down its dual using LP duality result. Two observations were note worthy: i) writing an optimization problem in the standard form itself could be non-trivial (and many times it is infact cumbersome) ii) from the dual perhaps one could easily read-off the optimal value (as it perhaps is a simple problem).
- We then took the case of DOCPs and proved the following theorem:

Theorem 20.0.11. *Given the following ordinary convex program (P):*

$$(20.1) \quad \begin{aligned} & \min_{x \in \mathcal{X}} && f(x), \\ & \text{s.t.} && g_i(x) \leq 0 \quad \forall i = 1, \dots, m, \end{aligned}$$

where all the functions involved are differentiable, the domain is open and the program is solvable. Then the following is a dual of (P):

$$(20.2) \quad \max_{\lambda \in \mathbb{R}_+^m} \underline{L}(\lambda),$$

where $\underline{L}(\lambda) \equiv \min_{x \in \mathcal{X}} L(x, \lambda)$ and $L(x, \lambda) \equiv f(x) + \sum_{i=1}^m \lambda_i g_i(x)$ (the former function is called as the Lagrange dual function and the latter function is called as the Lagrange function or simply the Lagrangian).

Please refer appendix pages 18-20 for a proof.

- Infact, the above theorem can be strengthened and for details please refer theorem D.2.2 in Nemirovski [2005]. This version is very useful as it does not assume differentiability!
- Again, it is easy to see that the dual above when written for a (P) that is not necessarily convex, gives weak duality but perhaps not strong duality.

However, the dual will remain convex. Hence this trick of writing the Lagrange dual (sometimes this trick is referred to as Lagrange relaxation) is useful in obtaining easy-to-compute lower bounds for non-convex programs.

- We also then realized that the primal can also be written in terms of the Lagrangian function: $(P) = \min_{x \in \mathcal{X}} \bar{L}(x)$, where $\bar{L}(x) \equiv \max_{\lambda \in \mathbb{R}_+^m} L(x, \lambda)$ (this function can be called as the Lagrangian primal function). In other words the following min-max interchange corollary is immediate:

Corollary 20.0.12. *In the context of theorem 20.0.11, the following min-max interchange is allowed:*

$$\min_{x \in \mathcal{X}} \max_{\lambda \in \mathbb{R}_+^m} L(x, \lambda) = \max_{\lambda \in \mathbb{R}_+^m} \min_{x \in \mathcal{X}} L(x, \lambda).$$

Please refer appendix D.3 in Nemirovski [2005] for more on such min-max theorems.

- **Mandatory Reading:** appendix section D in Nemirovski [2005], chapter 5 in Boyd and Vandenberghe [2004].

Lecture 21

- We began with an example illustrating Lagrange duality scheme. We took the problem of finding the maximally separating hyperplane between two given (finite) sets of points. The following points are note-worthy: i) the journey from the English sentence description of the problem to the DOCP form was indeed long. We introduced dummy variables etc. ii) thankfully, there was an analytical expression for Lagrangian dual function¹ iii) since we invariably introduce dummy variables in primal to bring it in DOCP form, the dual also contains variables/inequalities that perhaps can be eliminated. After this filtering of the dual, one may find an elegant interpretation for the dual — in this case the dual is the problem of finding the minimum distance between the convex hulls of the given sets. We commented that – Lagrange duality gave this elegant geometrical result and hence must indeed be a useful notion of duality in practice.
- We then derived the following dual for the conic program case:

Theorem 21.0.13. *Given the following conic program (P), which is solvable:*

$$(21.1) \quad \begin{array}{ll} \min_{x \in \mathbb{R}^n} & c^\top x, \\ \text{s.t.} & b - A^\top x \in K \subset \mathbb{R}^m. \end{array}$$

A dual of the above is:

$$(21.2) \quad \begin{array}{ll} \max_{y \in \mathbb{R}^m} & -b^\top y, \\ \text{s.t.} & Ay + c = 0, y \in K^* \end{array}$$

Please refer appendix pages 21-22 for a proof of this.

¹However we had to be smart in writing down terms in compact/vector form etc.; else the process is cumbersome. Also in many cases, analytical expression for the Lagrange dual may not exist.

- We noted that whenever we start with a self-dual cone K . the conic dual is as elegant and simple as the case of LPs (LP duality is infact a special case of this); however the conic program is far more generic.
- Needless to say, the above dual can also be written when the K in (P) is not a cone. In this case, still weak duality is guaranteed and the dual remains convex; however strong duality many not hold. Hence this conic duality scheme is also useful in obtaining lower bounds for non-convex programs.
- The Nemirovski [2005] book presents a slightly different conic duality theorem, which is worth studying (refer theorem 1.7.1).
- **Mandatory Reading:** chapter 1 in Nemirovski [2005].

Lecture 22

- For few lectures our objective is to study special sub-classes of convex programs — their form, dual, examples etc.
- Linear Programs and Quadratic Programs, defined earlier, are the most familiar to all of us and we must be familiar with some examples also. Previously we wrote down their duals also¹. Hence we wont further focus on them. There are a number of free toolboxes for solving these.
- The next bigger class is convex Quadratically Constrained Quadratic Programs (convex QCQPs):

$$(22.1) \quad \begin{aligned} \min_{x \in \mathbb{R}^n} \quad & \frac{1}{2}x^\top P x + q^\top x + r, \\ \text{s.t.} \quad & \frac{1}{2}x^\top P_i x + q_i^\top x + r_i \leq 0 \quad \forall i = 1, \dots, m, \end{aligned}$$

where all $P, P_i \succeq 0$. We will later (while dealing with SDPs) write their dual. CPLEX is an efficient solver for QCQPs and there are many others.

- In this lecture, we will focus on a super class of all of these, called as conic quadratic programs (CQs) or Second Order Cone Programs (SOCPs):

$$(22.2) \quad \begin{aligned} \min_{x \in \mathbb{R}^n} \quad & c^\top x, \\ \text{s.t.} \quad & \|A_i^\top x + b_i\| \leq c_i^\top x + d_i \quad \forall i = 1, \dots, m. \end{aligned}$$

The constraints of the form as those in the above program as called as conic quadratic constraints or second order cone constraints.

- The first key observation was that the above program can be written as a conic program with the cone as Cartesian/direct product of ice-cream cones (which

¹We leave the case of convex QPs that are non-strictly convex for later analysis. While dealing with SDPs we will write down dual for (non-convex) QCQPs, which will cover this case.

is self-dual). Hence by applying conic duality we obtained a dual of it. Please refer to section 2.1 in Nemirovski [2005] for details. Interestingly, the dual is also a conic-quadratic and hence self-dual.

- We then gave many examples of sets and functions that can be expressed as conic quadratics. Please refer section 2.3 in Nemirovski [2005].
- Importantly there are very efficient solvers for CQs: Mosek, SeDuMi and SDPT3.
- **Mandatory Reading:** chapter 2 in Nemirovski [2005] and Lobo et al. [1998].

Lecture 23

- Given that we studied conic programs with ice-cream cones, we next wished to study those with the cone as the psd cone (which is also self-dual in the space of symm. matrices). For that we had to first study linear transformations from one Hilbert space \mathcal{V} to another \mathcal{W} . The following comments were immediate:

1. If $\dim(V) = n$ and $\dim(W) = m$ and given a set of bases for each space, then there is a bijection between the set of all linear functions from V to W and the set of all matrices of size $m \times n$.
2. If M_l is the matrix associated with $l : V \mapsto W$, then we call the linear function associated with M_l^\top as the **adjoint** of l and denote it by $l^\top : W \mapsto V$.
3. Moreover, if the bases were orthogonal (for both spaces), then we have the following: $\langle l(v), w \rangle_W = \langle v, l^\top(w) \rangle_V$.

- We then wrote down the expression for conic program in arbitrary spaces:

$$(23.1) \quad \begin{aligned} \min_{x \in V} \quad & \langle c, x \rangle_V, \\ \text{s.t.} \quad & b -_W l(x) \in K \subset W \end{aligned}$$

- Using the fact that all (finite dim) Hilbert spaces are essentially Euclidean, the following dual of the conic program was immediate:

$$(23.2) \quad \begin{aligned} \max_{y \in W} \quad & -\langle b, y \rangle_W, \\ \text{s.t.} \quad & l^\top(y) + c = 0, \quad y \in K^*. \end{aligned}$$

- We then took the special case of (23.1) with $V = \mathbb{R}^n$, $l(x) = \sum_{i=1}^n x_i A_i$ (where all $A_i \in S^m$) and K as the cone of all psd matrices of size m . This we defined as **Semi Definite Program (SDP)**:

$$(23.3) \quad \begin{aligned} \min_{x \in \mathbb{R}^n} \quad & c^\top x, \\ \text{s.t.} \quad & B - \sum_{i=1}^n x_i A_i \succeq 0 \end{aligned}$$

- The constraints in the form as those in SDP are called as [Linear Matrix Inequalities \(LMIs\)](#).
- We then computed the adjoint $l^\top(Z)$ as Euclidean vector with entries as $\langle Z, A_i \rangle_F$. Using this the following dual for SDP is immediate:

$$(23.4) \quad \begin{array}{ll} \max_{Y \in S^m} & -\langle B, Y \rangle_F, \\ \text{s.t.} & \langle A_i, Y \rangle_F + c_i = 0 \quad \forall i = 1, \dots, m, \quad Y \succeq 0. \end{array}$$

Hence SDPs are (in dual form) simply programs with linear objective and linear constraints with an additional psd constraint.

- We then showed that every CQ constraint can be written as an LMI (refer pg 109 in Nemirovski [2005]). Moreover, multiple LMIs can be written as a single LMI by bundling the matrices in a block diagonal fashion. Thus every CQ program (and hence every LP, QP, QCQP) can be written as an SDP. Ofcourse the converse is not true.
- We then wrote down the problem of finding the most spherical ellipsoid that tightly encloses a finite set of given points as an MP. In the subsequent lecture we will write it as an SDP.
- Finally, Mosek, SeDuMi and SDPT3 all are efficient SDP solvers. Thus SDP is a very generic class of convex programs, with elegant theoretical results and freely available efficient toolboxes.
- **Mandatory Reading:** chapter 3.1 in Nemirovski [2005]; chapters 3,6 in Sheldon Axler [1997] for a review of linear maps.

Lecture 24

- Using the examples 18, 20b from chapter 3 in Nemirovski [2005] we have that the program from the previous lecture can be written as an SDP.
- Students are urged to study all the examples for SDP-representable functions/sets in chapter 3.
- We then set out to write the dual of a (non-convex) QCQP and interestingly the Lagrange dual turned out to be an SDP! Hence SDPs play an important role in bounding/approximating non-convex programs. Please refer section 3.4 in Nemirovski [2005] for details. As mentioned earlier, the process of writing the Lagrange dual for a non-convex program is called as Lagrange relaxation.
- We then went on to study another important class of CPs called Geometric Programs (GPs):

$$(24.1) \quad \begin{aligned} \min_{x \in \mathbb{R}_{++}^n} \quad & f(x), \\ \text{s.t.} \quad & g_i(x) \leq 1 \quad \forall i = 1, \dots, m_1, \\ & h_i(x) = 1 \quad \forall i = 1, \dots, m_2, \end{aligned}$$

where f, g_i are all posynomials and h_i are all monomials. Monomial is of the form $c \prod_{i=1}^n x_i^{a_i}$ (all x_i, c are positive and all $a_i \in \mathbb{R}$). Posynomial is simply a sum of finite number of monomials.

- We gave an example of problems that can be posed as a GP and realized that perhaps GPs can indeed model many problems.
- GPs need not be convex, but they can be written as a convex program by i) replacing x_i by e^{y_i} and then ii) taking log for all functions. With this the posynomials turn out as functions of the following form: $\log \left(\sum_i e^{b_i^\top y + d_i} \right)$, which is convex and the equality constraints become linear equalities.

- We then argued that the dual of this convexified GP is essentially maximizing entropy problem that is familiar to us.
- Most importantly, GPs can be efficiently solved and free solvers are available: e.g., GGPLAB¹.
- **Mandatory Reading:** chapter 3 in Nemirovski [2005] and Vandenberghe and Boyd [1996]; for GPs section 4.5.2 in Boyd and Vandenberghe [2004] and Boyd et al. [2007].

¹Download from <http://www.stanford.edu/~boyd/ggplab/>

Bibliography

- S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- S. Boyd, S.-J. Kim, L. Vandenberghe, and A. Hassibi. A Tutorial on Geometric Programming. *Optimization and Engineering*, 8(1):67–127, 2007.
- R. Fletcher. *Practical Methods of Optimization*. Wiley, 2000.
- Niels Lauritzen. Lectures on Convex Sets. Available at home.imf.au.dk/niels/leconset.pdf, 2010.
- M. Lobo, L. Vandenberghe, S. Boyd, and H. Lebret. Applications of Second-Order Cone Programming. *Linear Algebra and its Applications*, 284:193–228, 1998.
- A. Nemirovski. Lectures On Modern Convex Optimization. Available freely at www2.isye.gatech.edu/~nemirovs/Lect_ModConvOpt.pdf, 2005.
- R. T. Rockafellar. *Convex Analysis*. Princeton University Press, 1996.
- Sheldon Axler. *Linear Algebra Done Right*. Springer-Verlag, 1997.
- L. Vandenberghe and S. Boyd. Semidefinite Programming. *SIAM Review*, 38(1):49–95, 1996.

Subset (S)	Definition	Primal rep.	Dual rep.	Algebra	Topology
Linear	Forms vector space	Basis (B), $S = LIN(B)$	B^\perp for orth.comp. $S = \{v \in V \mid \langle v, b \rangle_V = 0 \forall b \in B^\perp\}$	Int., sum Union, comp.	Closed Open,bounded
Affine	Shifted linear set	AffineBasis (B), $S = AFF(B)$	B^\perp for orth.comp. (shift) $S = \{v \in V \mid \langle v, b \rangle_V = r_b \forall b \in B^\perp\}$	Int., sum Union, comp.	Closed Open,bounded
Cone	closed conic comb.	Conicly-spanning (B), $S = CONIC(B)$	B^* for dual cone $S = \{v \in V \mid \langle v, b \rangle_V \geq 0 \forall b \in B^*\}$	Int., sum Union, comp.	Closed Open,bounded
Convex	closed convex comb.	Convexly-spanning (B), $S = CONV(B)$	B' for polar cone ¹ $S = \{v \in V \mid \langle v, b \rangle_V \leq 1 \forall b \in B'\}$	Int., sum Union, comp.	Closed,Open,Bounded

29/Aug/2012

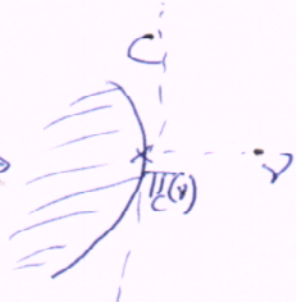
Result 1: A set $C \subset V$ is convex \Leftrightarrow ~~any~~ (non-empty) intersections of C with any line in V is convex.

Definition: Two non-empty sets $S, T \subset V$ are said to be ^(strictly) separable if

$$\exists a \in V \Rightarrow \sup_{s \in S} \langle a, s \rangle < \inf_{t \in T} \langle a, t \rangle$$

Separation Theorem: A closed convex set C and the singleton set $\{v\}$, where $v \notin C$ are (strictly) separable.

Proof: We will show this using notion of projection \rightarrow



Proof is in two steps:

(i) Existence of $\Pi_C(v)$: $\Pi_C(v) = \operatorname{argmin}_{x \in C} \|x - v\|$ always exists.

(ii) $\langle x - \Pi_C(v), v - \Pi_C(v) \rangle \leq 0 \quad \forall x \in C$.

Given these it is easy to see that $v - \Pi_C(v)$ is the vector ^{'a' in separation} ~~in the~~ defn.

(i) is a standard result in analysis and follows from the HEINE BOREL theorem (rough sketch given during lecture).

Proof of (ii): $\underline{TST} \quad \langle x - \pi_C(v), y - \pi_C(v) \rangle \leq 0 \quad \forall x \in C.$

We know $x \in C$ & $\pi_C(v) \in C \Rightarrow \lambda x + (1-\lambda)\pi_C(v) \in C \quad \forall \lambda \in (0,1)$

$$\text{Also, } \|\lambda x + (1-\lambda)\pi_C(v) - v\| \geq \|\pi_C(v) - v\| \quad \forall \lambda \in (0,1)$$

$$\Leftrightarrow \|\lambda(x - \pi_C(v)) + \pi_C(v) - v\|^2 \geq \|\pi_C(v) - v\|^2 \quad \forall \lambda \in (0,1)$$

$$\Leftrightarrow \lambda^2 \|x - \pi_C(v)\|^2 + 2\lambda \langle x - \pi_C(v), \pi_C(v) - v \rangle \geq 0 \quad \forall \lambda \in (0,1)$$

$$\Leftrightarrow \langle x - \pi_C(v), v - \pi_C(v) \rangle \leq \lambda \|x - \pi_C(v)\|^2 \quad \forall \lambda \in (0,1)$$

$$\Leftrightarrow \langle x - \pi_C(v), v - \pi_C(v) \rangle \leq 0. \quad \text{Here Proved.}$$

Theorem: A set $C \subset V$ is

~~polyhedral~~ polyhedron $\Leftrightarrow \exists S, T \subset V \ni$

(defn: convex set with finite dual descrip.) $C = \text{CONIC}(S) + \text{CONV}(T).$

Proof of \Rightarrow (we have proof of \Leftarrow as exercise)

TST: Given C is polyhedron i.e. $C = \left\{ x \mid \begin{array}{l} \langle x, v_1 \rangle \leq \alpha_1, \\ \langle x, v_m \rangle \leq \alpha_m \end{array} \right\}$

$\exists S, T \subset V \ni C = \text{CONIC}(S) + \text{CONV}(T).$

The proof is by construction and uses an ~~analogous~~ ^{analogous} result for cones.

Consider ~~Prove~~ $K = \left\{ (x, y) \mid \begin{aligned} &\langle (x, y), (y_1, -x_1) \rangle \leq 0, \\ &\vdots \\ &\langle (x, y), (y_m, -x_m) \rangle \leq 0, \\ &\langle (x, y), (0, -1) \rangle \leq 0 \end{aligned} \right\}$ (which is a cone)

$C \in V \times \mathbb{R}$ and is in the usual dual norm space.

It is easy to see that $x \in C \Leftrightarrow (x, 1) \in K$.

But K is polyhedral cone $\Rightarrow K = \text{CONIC}(\{(u_1, 0), \dots, (u_p, 0), (w_1, 1), \dots, (w_q, 1)\})$

$= \left\{ \sum_i \lambda_i (u_i, 0) + \sum_i \mu_i (w_i, 1) \mid \lambda_i \geq 0, \mu_i \geq 0 \right\}$

~~But $C = \{x \mid (x, 1) \in K\} = \left\{ \left(\sum_i \lambda_i u_i + \sum_i \mu_i w_i, \sum_i \mu_i \right) \mid \lambda_i \geq 0, \mu_i \geq 0 \right\}$~~

~~$= \left\{ \sum_i \lambda_i (u_i, 0) + \sum_i \mu_i (w_i, 1) \mid \lambda_i \geq 0, \mu_i \geq 0 \right\}$~~

But $C = \{x \mid (x, 1) \in K\} = \left\{ \sum_i \lambda_i u_i + \sum_i \mu_i w_i \mid \lambda_i \geq 0, \mu_i \geq 0, \sum_i \mu_i = 1 \right\}$

Hence, with, $S = \{u_1, \dots, u_p\}$, $T = \{w_1, \dots, w_q\}$, we have:

$$C = \text{CONIC}(S) + \text{CONV}(T)$$

Hence Proved.

Result:

~~K is a cone that is convex~~

K is polyhedral cone $\Rightarrow K$ is intersection of finite no. half spaces (through origin)

(Corollary: Since $K^{**} = K$, we have the other result: ~~that~~ a cone that is the intersection of finite no. half spaces is polyhedral)

\rightarrow To simplify notation proof is given for Euclidean spaces \leftarrow

Proof:

K is polyhedral cone $\Rightarrow K = \left\{ \sum_{i=1}^m \lambda_i v_i \mid \lambda_i \geq 0 \right\}$

for some $v_1, \dots, v_m \in \mathbb{R}^n$ ($m \in \mathbb{N}$)

$\Rightarrow K = \left\{ V\lambda \mid \lambda \geq 0 \right\}$

(V is matrix with cols. as v_i & λ is vector with entries as λ_i)

$= \left\{ y \mid y = V\lambda, \lambda \geq 0 \right\}$

Consider the cone: $\bar{K} = \left\{ \begin{bmatrix} x \\ y \end{bmatrix} \in \mathbb{R}^{n+m} \mid \begin{array}{l} y - Vx \leq 0 \\ Vx - y \leq 0 \\ -x \leq 0 \end{array} \right\}$

Again, $x \in K \Leftrightarrow x \in \mathbb{R}^n$ projection of \bar{K} onto the last n dimension.

It is an easy exercise to show that ~~a~~ ^{such} projection of a cone with finite dual description is another cone with finite dual description. Hence Proved.

Dual representation of conic functions

Given,

$f: K \rightarrow \mathbb{R}$ is a conic function.

Since $\text{epi}(f)$ is a conic set, we have,

$$\begin{aligned} \text{epi}(f) &= \left\{ (x, y) \mid \langle (u, v), (x, y) \rangle \geq 0 \quad \forall (u, v) \in (\text{epi}(f))^* \right\} \\ &= \left\{ (x, y) \mid \langle u, x \rangle + vy \geq 0 \quad \forall (u, v) \in (\text{epi}(f))^* \right\} \end{aligned}$$

If $v < 0$, then $y \leq \langle \frac{u}{v}, x \rangle$, which is not possible.

Hence $v \geq 0$. Also, not all 'v' are zero as its an epigraph.

However, some of 'v' may be zero. These determine the domain of f.

$$\begin{aligned} \text{epi}(f) &= \left\{ (x, y) \mid \begin{array}{l} \langle u, x \rangle \geq 0 \quad \forall (u, 0) \in (\text{epi}(f))^* \\ y \geq \max_{(u, v) \in V \times \mathbb{R}} \langle -\frac{u}{v}, x \rangle \quad \forall (u, v) \in (\text{epi}(f))^*, v \neq 0 \end{array} \right\} \\ &= \left\{ (x, y) \mid \begin{array}{l} \max_{(u, v) \in V \times \mathbb{R}} \langle -\frac{u}{v}, x \rangle \leq y, \\ \text{s.t. } (u, v) \in (\text{epi}(f))^*, v \neq 0 \end{array} \right\} \\ &= \left\{ (x, y) \mid \langle u, x \rangle \geq 0 \quad \forall (u, 0) \in (\text{epi}(f))^* \right\} \\ &= \left\{ (x, y) \mid f(x) \leq y, \quad x \in K \right\} \end{aligned}$$

Here,

$$f(x) = \max_{(u,v) \in V \times R} \left\langle \frac{-u}{v}, x \right\rangle$$

s.t. $(u,v) \in (\text{epi}(f))^*, v \neq 0$

$$\text{dom}(f) = K = \left\{ x \mid \begin{array}{c} \exists \\ \langle u, x \rangle \geq 0 \wedge (u,0) \in (\text{epi}(f))^* \end{array} \right\}$$

Towards dual defn. of convex functions

7/9/12

With the experience of convex sets and conic sets and the results of conic functions, we guess that convex functions are exactly those that are maximum of a set of affine functions.
(maximum, pointwise)

Claim 1: Let $f(x) = \max_{(\beta, \alpha) \in S} \langle \beta, x \rangle + \alpha$. Then f is ^{closed} convex, whenever $\text{dom}(f)$ is closed convex.

Proof: Consider $\text{epi}(f) = \left\{ (x, y) \mid \max_{(\beta, \alpha) \in S} \langle \beta, x \rangle + \alpha \leq y \right\}$

$= \left\{ (x, y) \mid x \in \text{dom}(f), \langle \beta, x \rangle + \alpha \leq y \ \forall (\beta, \alpha) \in S \right\}$

$= \left\{ (x, y) \mid x \in \text{dom}(f), \langle (\beta, -1), (x, y) \rangle \leq -\alpha \ \forall (\beta, \alpha) \in S \right\}$

\downarrow dual given convex \uparrow int. of half spaces is convex \downarrow (this form because epigraph)

$\text{epi}(f)$ is convex $\Rightarrow f$ is ^{closed} convex.

Claim 2: All closed convex functions ^{can} be written as maximum of a set of affine functions. ~~every real closed~~

Proof: See red marks in above proof.

Given $f: \mathcal{C}V \rightarrow \mathbb{R}$ that is closed convex, we have:

$$f'' = f$$

(conjugate of conjugate)

(We assume both f' and f'' are defined (real-valued) on V i.e. the max in defn. is never infinity)
 It is easy to extend this to case of functions with a domain $\subset V$, by considering the notion of extended-valued functions

Proof: $f''(x) = \max_{y \in V} \langle x, y \rangle - f'(y)$

$$= \max_{y \in V, \alpha \in \mathbb{R}} \langle x, y \rangle - \alpha$$

s.t. $f'(y) \leq \alpha$

$$= \max_{y \in V, \alpha \in \mathbb{R}} \langle x, y \rangle - \alpha$$

s.t. $\max_{z \in V} \langle y, z \rangle - f(z) \leq \alpha$

$$= \max_{y \in V, \alpha \in \mathbb{R}} \langle x, y \rangle - \alpha$$

s.t. $\langle y, z \rangle - \alpha \leq f(z) \quad \forall z \in V$

Constraints \Rightarrow any (y, α) define all affine functions less than f .
 (i.e. each $g(\cdot) = \langle y, \cdot \rangle - \alpha$ defines an affine minorant of f .)

Now let the feasibility set of above optimization prob. be \mathcal{F} .

$$\text{i.e. } (y, \alpha) \in \mathcal{F} \Leftrightarrow \langle y, z \rangle - \alpha \leq f(z) \quad \forall z \in V$$

$$\Leftrightarrow \text{epi}(f) \subseteq \left\{ (z, \beta) \mid \langle (y, -1), (z, \beta) \rangle \leq \alpha \right\}$$

Since \mathcal{F} is closed convex, we have $\text{epi}(f)$ is closed convex and thus

$$\text{epi}(f) = \left\{ (z, \beta) \mid \langle (y, -1), (z, \beta) \rangle \leq \alpha \quad \forall (y, \alpha) \in \mathcal{F} \right\} \quad \text{--- (I)}$$

However we also have:

$$f''(x) = \max_{(y, \alpha) \in \mathcal{F}} \langle x, y \rangle - \alpha$$

$$\begin{aligned} \text{epi}(f'') &= \left\{ (x, \nu) \mid \langle x, y \rangle - \alpha \leq \nu \quad \forall (y, \alpha) \in \mathcal{F} \right\} \\ &= \left\{ (x, \nu) \mid \langle (y, -1), (x, \nu) \rangle \leq \alpha \quad \forall (y, \alpha) \in \mathcal{F} \right\} \quad \text{--- (II)} \end{aligned}$$

(I) & (II) give that $\text{epi}(f) = \text{epi}(f'') \Rightarrow f = f''$. Hence Proved.

Result: The CP (15.1) is solvable whenever F is compact and either (i) δ is continuous convex &

(ii) δ is closed convex

(In fact, conditions (i) and (ii) are same as F is closed)

Proof: Since F is bounded, we have that CP is bounded

\Downarrow
 $\exists \delta^* \in \mathbb{R}$ such that

$$\min_{x \in X} \delta(x) = \delta^*$$

\Leftarrow s.t. $x \in F$

\Rightarrow \exists a sequence $\{\delta(x_n)\} \rightarrow \delta^*$

Since F is compact, by Heine Bolz theorem, we know that there exists a sub-sequence of $\{x_n\}$, say $\{x_{n_k}\}$ such that:

$$\{x_{n_k}\} \rightarrow x^* \in F$$

$$\Rightarrow \{\delta(x_{n_k})\} \rightarrow \delta(x^*) = \delta^*$$

($\because \delta$ is continuous)

\Downarrow
 x^* is ~~the~~ ^{an} optimal solution of CP

\Downarrow
CP is solvable.

~~Result~~ A (CP) with strictly ~~convex and~~ closed

Result: A (CP) with a strictly closed convex objective δ and a compact feasibility set (F) is uniquely solvable.

Proof: First of all we know such a (CP) is solvable.

Now suppose it is not uniquely solvable i.e. $\exists x_1^* \in F, x_2^* \in F \Rightarrow$

$$x_1^* \neq x_2^* \quad \& \quad \delta(x_1^*) \leq \delta(x) \quad \forall x \in F \Rightarrow \delta(x_1^*) \leq \delta(x_2^*)$$

$$\delta(x_2^*) \leq \delta(x) \quad \forall x \in F \Rightarrow \delta(x_2^*) \leq \delta(x_1^*) \quad \rightarrow \delta(x_1^*) = \delta(x_2^*)$$

Now consider $\lambda x_1^* + (1-\lambda)x_2^* \in F$ (where $\lambda \in (0,1)$)

we have by strict concavity, $\delta(\lambda x_1^* + (1-\lambda)x_2^*) < \lambda \delta(x_1^*) + (1-\lambda) \delta(x_2^*)$

$$= \delta(x_1^*)$$

which is not possible.

Result: J is convex $\Rightarrow T_J(x^*)$ is a cone. for any $x^* \in J$.

Proof: Recall, $T_J(x^*) = \{u \mid \exists h > 0 \ni x^* + hu \in J\}$

Now, $u \in T_J(x^*) \Rightarrow \exists h > 0 \ni x^* + hu \in J$

$$\Rightarrow \exists h > 0 \ni x^* + \left(\frac{h}{\alpha}\right)(\alpha u) \in J \quad \forall \alpha > 0$$

$$\Rightarrow \exists h > 0 \ni x^* + h(\alpha u) \in J \quad \forall \alpha > 0$$

$$\Rightarrow \alpha u \in T_J(x^*) \quad \forall \alpha > 0$$

Also $0 \in T_J(x^*)$ as $x^* \in J$

$$\therefore u \in T_J(x^*) \Rightarrow \alpha u \in T_J(x^*) \quad \forall \alpha \geq 0. \quad \text{--- (I)}$$

Now let's show $u_1 \in T_J(x^*) \Rightarrow u_1 + u_2 \in T_J(x^*)$ --- (II)

$u_2 \in T_J(x^*)$

(I) & (II) give the required result.

Proof of (II) $\exists h_1 > 0 \ni x^* + h_1 u_1 \in J$ & $\exists h_2 > 0 \ni x^* + h_2 u_2 \in J$

$$\frac{h_2}{h_1 + h_2} (x^* + h_1 u_1) + \frac{h_1}{h_1 + h_2} (x^* + h_2 u_2) \in J \quad (\because J \text{ is convex})$$

$$\Downarrow \\ x^* + \frac{h_1 h_2}{h_1 + h_2} (u_1 + u_2) \in J \Rightarrow u_1 + u_2 \in T_J(x^*)$$

Here proved.

Freedom: Given a differentiable (CP),

$$x^* \text{ is optimal} \Leftrightarrow x^* \in \mathcal{F}, \nabla f(x^*) \in \mathcal{N}_{\mathcal{F}}(x^*)$$

Proof: $x^* \in \mathcal{F}, f(x^*) \leq f(x) \forall x \in \mathcal{F}$

$$x^* \in \mathcal{F}, f(x^*) \leq f(x^* + hu) \forall u \in \mathcal{T}_{\mathcal{F}}(x^*), \forall h > 0, x^* + hu \in \mathcal{F}$$

$$x^* \in \mathcal{F}, \frac{f(x^* + hu) - f(x^*)}{h} \geq 0 \forall u \in \mathcal{T}_{\mathcal{F}}(x^*), \forall h > 0, x^* + hu \in \mathcal{F}$$

$$x^* \in \mathcal{F}, \min_{h > 0} \frac{f(x^* + hu) - f(x^*)}{h} \geq 0 \forall u \in \mathcal{T}_{\mathcal{F}}(x^*)$$

s.t. $x^* + hu \in \mathcal{F}$

$$x^* \in \mathcal{F}, \langle \nabla f(x^*), u \rangle \geq 0 \forall u \in \mathcal{T}_{\mathcal{F}}(x^*)$$

$$x^* \in \mathcal{F}, \nabla f(x^*) \in \mathcal{N}_{\mathcal{F}}(x^*)$$

Hence Proved.

but, from sub-gradient inequality we have:

$$\min_{h > 0} \frac{f(x^* + hu) - f(x^*)}{h} \geq \langle \nabla f(x^*), u \rangle$$

s.t. $x^* + hu \in \mathcal{F}$

$$\lim_{h \rightarrow 0} \frac{f(x^* + hu) - f(x^*)}{h} = \langle \nabla f(x^*), u \rangle$$

$$\therefore \min_{h > 0} \frac{f(x^* + hu) - f(x^*)}{h} = \langle \nabla f(x^*), u \rangle$$

s.t. $x^* + hu \in \mathcal{F}$

DPCCP. Duality

(P) $\min_{x \in X} f(x)$ $\xrightarrow{\text{differentiable}}$ Let f be feasible
s.t. $\langle a_i, x \rangle \leq b_i \quad \forall i \in J_m$ primal, given program.
ret of (P)
open

~~One way of writing dual is to write down expression for optimal value involving x^* (optimal value) but it is~~

Our goal is to write down a dual. We will do this by following steps:

(i) Think of an expression for optimal value without involving x^* (optimal value).
This can perhaps be done from KKT conditions.

(ii) See under what conditions the optimal value expression, as a function of the involved (dual) variables, forms a lower bound for the optimal value. From this the dual will follow.

We will follow the above strategy for all cases in this course.

Towards (i): Recall that

x^* is optimal iff (i) $\langle a_i, x^* \rangle \leq b_i \quad \forall i$

(ii) $\exists \lambda^* \in \mathbb{R}^m \ni$

(a) $\lambda^* \geq 0$

(b) $\sum \lambda_i^* a_i + \nabla f(x^*) = 0$

(c) $\lambda_i^* (\langle a_i, x^* \rangle - b_i) = 0 \quad \forall i$

This gives an analytical expression for $\nabla f(x^*) = -\sum_i \lambda_i^* a_i$.
Whereas we need analytical expression for $f(x^*)$.

It is clear that we will use the following result:

$$f(x^*) + f'(\nabla f(x^*)) = \langle \nabla f(x^*), x^* \rangle$$

↳ conjugate of f over domain X

$$\begin{aligned} \text{Hence } f(x^*) &= \langle -\sum_i \lambda_i^* a_i, x^* \rangle - f'(-\sum_i \lambda_i^* a_i) \\ &= -\sum_i \lambda_i^* \langle a_i, x^* \rangle - f'(-\sum_i \lambda_i^* a_i) \\ &= -\sum_i \lambda_i^* b_i - f'(-\sum_i \lambda_i^* a_i) \quad (\because \text{KKT}) \end{aligned}$$

$$\therefore \min_{x \in X} f(x) = -\sum_i \lambda_i^* b_i - f'(-\sum_i \lambda_i^* a_i)$$

s.t. $\langle a_i, x \rangle \leq b_i$

→ (I)

where λ_i^* satisfies KKT condition.

This motivates us to consider the following function:

$$g(\lambda) = -\sum_i \lambda_i b_i - f'(-\sum_i \lambda_i a_i)$$

We will now show that $f(x) \geq g(\lambda) \quad \forall$

~~$x \in X, \lambda \geq 0$~~
 $x \in \mathbb{R}^n, \lambda \geq 0$

Proof: Consider the Fenchel inequality:

$$f(x) + f'(y) \geq \langle x, y \rangle \quad \forall x \in X, y \in \text{dom}(f')$$

$$\begin{aligned} \Rightarrow f(x) + f'(-\sum_i \lambda_i a_i) &\geq \langle x, -\sum_i \lambda_i a_i \rangle \\ &\geq -\sum_i \lambda_i \langle a_i, x \rangle \\ &\geq -\sum_i \lambda_i b_i \quad \forall x \in X, \lambda \geq 0. \end{aligned}$$

~~$\forall x \in X, \lambda \geq 0$~~
 $\forall x \in \mathbb{R}^n, \lambda \geq 0$
 ~~$\forall x \in X, \lambda \geq 0$~~

Hence,

$$f(x) \geq -\sum_i \lambda_i b_i - f'(-\sum_i \lambda_i a_i)$$

over $x \in F$ over $\lambda \geq 0$

$$\min_{x \in X} f(x)$$

$$\geq \max_{\lambda \geq 0} -\sum \lambda_i b_i - f'(-\sum \lambda_i a_i)$$

s.t. $\langle a_i, x \rangle \leq b_i, \forall i$

(Weak duality)

From (I) we get strong duality (i.e. the RHS equals LHS at $\lambda = \lambda^*$)

$$\min_{x \in X} f(x) \quad \text{s.t.} \quad \langle a_i, x \rangle \leq b_i \quad = \quad \max_{\lambda \geq 0} -\sum \lambda_i b_i - f'(-\sum \lambda_i a_i) \quad \text{(D)}$$

Now, this can be called a dual provided we prove that (D) this is a convex program:

$$-(D) = \min_{\lambda \geq 0} \underbrace{f'(-\sum \lambda_i a_i)}_{\text{Convex on } \lambda} + \underbrace{\sum \lambda_i b_i}_{\text{Linear } \lambda}$$

∴ (D) is a convex program.

Lagrange duality for DCP

OCP

min
 $x \in X$

s.t.

$f(x)$ → differentiable.

$g_i(x) \leq 0 \quad \forall i=1, \dots, m$
→ differentiable.

Let \mathcal{F} be the feasibility set of this DCP

Again, we will start by writing analytical expression of $f(x^*)$:

x^* is optimal \Leftrightarrow (i) $g_i(x^*) \leq 0 \quad \forall i$

(ii) $\exists \lambda^* \in \mathbb{R}^m \exists$

(a) $\lambda^* \geq 0$

(b) $\nabla f(x^*) + \sum_i \lambda_i^* \nabla g_i(x^*) = 0$

(c) $\lambda_i^* (g_i(x^*)) = 0 \quad \forall i$

This motivates the definition of the following function:

$$L_\lambda(x) = L(x, \lambda) = f(x) + \sum_i \lambda_i g_i(x).$$

(abuse of notation)

→ This is called the Lagrange function or simply Lagrangian.

(b) gives $\nabla L_{\lambda^*}(x^*) = 0$

From the conjugate ^{result} defn, we have:

$$L_{\lambda^*}(x^*) + L'_{\lambda^*}(\nabla L_{\lambda^*}(x^*)) = \langle \nabla L_{\lambda^*}(x^*), x^* \rangle = 0$$

$$\Rightarrow \mathcal{L}_{\lambda^*}(x^*) = -\mathcal{L}'_{\lambda^*}(0)$$

$$\text{but } \mathcal{L}_{\lambda^*}(x^*) = f(x^*)$$

$$\Rightarrow f(x^*) = -\mathcal{L}'_{\lambda^*}(0) = -\max_{x \in X} (f(x) + \sum_i \lambda_i^* g_i(x))$$

$$= \min_{x \in X} f(x) + \sum_i \lambda_i^* g_i(x)$$

$$\equiv \underline{\mathcal{L}}(\lambda^*)$$

definition (notation)

$$\therefore \min_{x \in X} f(x) = \underline{\mathcal{L}}(\lambda^*), \text{ where}$$

$$\text{s.t. } g_i(x) \leq 0 \quad \forall i$$

(I)

$$\text{where, } \lambda^* \text{ satisfies KKT cond. \& } \underline{\mathcal{L}}(\lambda^*) = \min_{x \in X} f(x) + \sum_i \lambda_i^* g_i(x)$$

~~Therefore, it is clear that~~ Now, we claim:

$$f(x) \geq \underline{\mathcal{L}}(\lambda) \quad \forall x \in \mathcal{F}, \lambda \geq 0.$$

Proof:

$$\begin{aligned} f(x) &\geq f(x) + \sum_i \lambda_i g_i(x) \quad \forall \lambda \geq 0, x \in \mathcal{F} \\ &\geq \min_{x \in \mathcal{F}} f(x) + \sum_i \lambda_i g_i(x) \\ &\geq \min_{x \in X} f(x) + \sum_i \lambda_i g_i(x) = \underline{\mathcal{L}}(\lambda) \quad \forall \lambda \geq 0 \end{aligned}$$

Therefore,

$$\begin{array}{l} \min_{x \in X} f(x) \\ \text{s.t. } g_i(x) \leq 0 \quad \forall i \end{array} \geq \max_{\lambda \geq 0} \underline{L}(\lambda) \quad (\text{weak duality})$$

by \textcircled{I} we have:

$$\begin{array}{l} \min_{x \in X} f(x) \\ \text{s.t. } g_i(x) \leq 0 \quad \forall i \end{array} = \max_{\lambda \geq 0} \underline{L}(\lambda) \quad \textcircled{D} \quad (\text{strong duality})$$

It is easy to see that $\underline{L}(\lambda)$ is convex & hence \textcircled{D} is a convex program & qualifies to be called a dual.

Conic Duality (Euclidean Case)

$$\textcircled{P} \quad \min_{x \in \mathbb{R}^n} c^T x$$

$$\text{s.t.} \quad b - A^T x \in K \quad \text{core.}$$

Let \textcircled{F} be feasibility set of \textcircled{P}

From KKT: x^* is optimal \Leftrightarrow

- (i) $b - A^T x^* \in K$
- (ii) $\exists \lambda^* \in \mathbb{R}^m \Rightarrow$
 - (a) $\lambda^* \in K^*$
 - (b) $A \lambda^* + c = 0$
 - (c) $\lambda^{*T} (b - A^T x^*) = 0$

Now,
$$c^T x^* = (-A \lambda^*)^T x^* = -\lambda^{*T} A^T x^* + \lambda^{*T} b - \lambda^{*T} b$$

$$= \lambda^{*T} (b - A^T x^*) - \lambda^{*T} b = -\lambda^{*T} b$$

$\therefore \min_{x \in \mathbb{R}^n} c^T x = -\lambda^{*T} b \quad \textcircled{I}$

s.t. $b - A^T x \in K$

This motivates us to look at $g(\lambda) = -b^T \lambda$

Claim:

$c^T x \geq -b^T \lambda \quad \forall x \in F, \lambda \in K^*, A \lambda + c = 0$

Proof:

$$c^T x = -\lambda^T A^T x = \lambda^T (b - A^T x) - b^T \lambda \geq -b^T \lambda$$
(dual core defn.)

Hence,

$$\begin{array}{ll} \min_{x \in \mathbb{R}^n} & c^T x \\ \text{s.t.} & b - A^T x \in K \end{array} \geq \begin{array}{ll} \max_{\lambda \in \mathbb{R}^m} & -b^T \lambda \\ \text{s.t.} & \lambda \in K^*, A\lambda + c = 0 \end{array}$$

(weak duality)

In fact, by \textcircled{I} :

$$\begin{array}{ll} \min_{x \in \mathbb{R}^n} & c^T x \\ \text{s.t.} & b - A^T x \in K \end{array} = \begin{array}{ll} \max_{\lambda \in \mathbb{R}^m} & -b^T \lambda \\ \text{s.t.} & \lambda \in K^*, A\lambda + c = 0 \end{array}$$

(strong duality) \textcircled{D}

\textcircled{D} is a convex program & hence qualifies to be a dual.