# Introduction to Convexity

Amitabh Basu

Compiled on Thursday 6[th] December, 2018 at 21:43

NOTES:1

# Contents

NOTES:          2

NOTES:                    3

# Why study convexity?

1. Convex optimization: least squares problem (linear regression), compressed sensing, classification.

2. Farkas' lemma : Fundamental theorem of asset pricing/No arbitrage theorem. Mention convexity assumption in finance for risk measures

3. Von Neumann's minimax theorem/existence of Nash equilibria

4. Statistical learning: nonnegative matrix factorization problem

5. Helly's Theorem: At least 1/3rd area cut off, Voting in agreeable societies

6. Hugely important tool in combinatorial optimization: canonical example – Transshipment problem

7. Radon's theorem: VC dimension of halfspaces from statistical learning theory

# 1   Definitions and Preliminaries

We will focus on $\mathbb{R}^d$ for arbitrary $d \in \mathbb{N}$: $\mathbf{x} = (x_1, \ldots, x_d) \in \mathbb{R}^d$. We will use the notation $\mathbb{R}^d_+$ to denote the set of all vectors with nonnegative coordinates. We will also use $\mathbf{e}^i$, $i = 1, \ldots, d$ to denote the $i$-th unit vector, i.e., the vector which has 1 in the $i$-th coordinate and 0 in every other coordinate.

**Definition 1.1.** A *norm on $\mathbb{R}^d$* is a function $N : \mathbb{R}^d \to \mathbb{R}_+$ satisfying:

1. $N(\mathbf{x}) = 0$ if and only if $\mathbf{x} = \mathbf{0}$,

2. $N(\alpha \mathbf{x}) = |\alpha| N(\mathbf{x})$ for all $\alpha \in \mathbb{R}$ and $\mathbf{x} \in \mathbb{R}^d$,

3. $N(\mathbf{x} + \mathbf{y}) \leq N(\mathbf{x}) + N(\mathbf{y})$ for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$. (Triangle inequality)

**Example 1.2.** For any $p \geq 1$, define the $\ell^p$ norm on $\mathbb{R}^d$: $\|\mathbf{x}\|_p = (|x_1|^p + |x_2|^p + \ldots + |x_d|^p)^{\frac{1}{p}}$. $p = 2$ is also called the *standard Euclidean norm*; we will drop the subscript 2 to denote the standard norm: $\|\mathbf{x}\| = \sqrt{x_1^2 + x_2^2 + \ldots + x_d^2}$. The $\ell^\infty$ norm is defined as $\|\mathbf{x}\|_\infty = \max_{i=1}^n |x_i|$.

**Definition 1.3.** Any norm on $\mathbb{R}^d$ defines a distance between points in $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ as $d_N(\mathbf{x}, \mathbf{y}) := N(\mathbf{x} - \mathbf{y})$. This is called the *metric or distance induced by the norm.* Such a metric satisfies three important properties:

1. $d_N(\mathbf{x}, \mathbf{y}) = \mathbf{0}$ if and only if $\mathbf{x} = \mathbf{y}$,

2. $d_N(\mathbf{x}, \mathbf{y}) = d_N(\mathbf{y}, \mathbf{x})$ for all $\mathbf{x} \in \mathbb{R}^d$,

3. $d_N(\mathbf{x}, \mathbf{z}) \leq d_N(\mathbf{x}, \mathbf{y}) + d_N(\mathbf{y}, \mathbf{z})$ for all $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbb{R}^d$. (Triangle inequality)

NOTES:                                4

**Definition 1.4.** We also utilize the (standard) inner product of $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ : $\langle \mathbf{x}, \mathbf{y} \rangle = x_1 y_1 + x_2 y_2 + \ldots + x_d y_d$. (Note that $\|\mathbf{x}\|_2^2 = \langle \mathbf{x}, \mathbf{x} \rangle$). We say $\mathbf{x}$ and $\mathbf{y}$ are orthogonal if $\langle \mathbf{x}, \mathbf{y} \rangle = 0$.

**Definition 1.5.** For any norm $N$ and $\mathbf{x} \in \mathbb{R}^d$, $r \in \mathbb{R}_+$, we will call the set $B_N(\mathbf{x}, r) := \{\mathbf{y} \in \mathbb{R}^d : N(\mathbf{y} - \mathbf{x}) \leq r\}$ as the *ball around* $\mathbf{x}$ *of radius* $r$. $B_N(\mathbf{0}, 1)$ will be called the *unit ball for the norm* $N$.

A subset $X \subseteq \mathbb{R}^d$ is said to be *bounded* if there exists $R \in \mathbb{R}$ such that $X \subseteq B_N(\mathbf{0}, R)$.

**Definition 1.6.** Given any set $X \subseteq \mathbb{R}^d$ and a scalar $\alpha \in \mathbb{R}$,

$$\alpha X := \{\alpha \mathbf{x} : \mathbf{x} \in X\}.$$

Given any two sets $X, Y \subseteq \mathbb{R}^d$, we define the *Minkowski sum* of $X, Y$ as

$$X + Y := \{\mathbf{x} + \mathbf{y} : \mathbf{x} \in X, \mathbf{y} \in Y\}.$$

**Basic real analysis and topology.** For any subset of real numbers $S \subseteq \mathbb{R}$, we recall the concepts of the *infimum* $\inf S$ and the *supremum* $\sup S$.

Fix a norm $N$ on $\mathbb{R}^d$. A set $X \subseteq \mathbb{R}^d$ is called *open* if for every $\mathbf{x} \in X$, there exists $r \in \mathbb{R}_+$ such that $B_N(\mathbf{x}, r) \subseteq X$. A set $X$ is *closed* if its complement $\mathbb{R}^d \setminus X$ is open.

**Theorem 1.7.**   1. $\emptyset, \mathbb{R}^d$ are both open and closed.

2. An arbitrary union of open sets is open. An arbitrary intersection of closed sets is closed.

3. A finite intersection of open sets is open. A finite union of closed sets is closed.

A *sequence* in $\mathbb{R}^d$ is a countable ordered set of points: $\mathbf{x}^1, \mathbf{x}^2, \mathbf{x}^3, \ldots$. We say that *the sequence converges* or that *the limit of the sequence exists* if there exists a point $\mathbf{x}$ such that for every $\epsilon > 0$, there exists $M \in \mathbb{N}$ such that $N(\mathbf{x} - \mathbf{x}^n) \leq \epsilon$ for all $n \geq M$. $\mathbf{x}$ is called the *limit point*, or simply the *limit*, of the sequence and will also sometimes be denoted by $\lim_{n \to \infty} \mathbf{x}^n$.

**Theorem 1.8.** A set $X$ is closed if and only if for every convergent sequence in $X$, the limit of the sequence is also in $X$.

We introduce three important notions:

1. For any set $X \subseteq \mathbb{R}^d$, the *closure of* $X$ is the smallest closed set containing $X$ and will be denoted by $\mathrm{cl}(X)$.

2. For any set $X \subseteq \mathbb{R}^d$, the *interior of* $X$ is the largest open set contained inside $X$ and will be denoted by $\mathrm{int}(X)$.

3. For any set $X \subseteq \mathbb{R}^d$, the *boundary of* $X$ is defined as $\mathrm{bd}(X) := \mathrm{cl}(X) \setminus \mathrm{int}(X)$.

NOTES:                                                            5

**Definition 1.9.** A set in $\mathbb{R}^d$ that is closed and bounded is called *compact.*

**Theorem 1.10.** Let $C \subseteq \mathbb{R}^d$ be a compact set. Then every sequence $\{\mathbf{x}^i\}_{i\in\mathbb{N}}$ contained in $C$ (not necessarily convergent) has a convergent subsequence.

A function $f : \mathbb{R}^d \to \mathbb{R}^n$ is *continuous* if for every convergent sequence $\{\mathbf{x}^n\}_{n=1}^\infty \subseteq \mathbb{R}^d$, the following holds: $\lim_{n\to\infty} f(\mathbf{x}^n) = f(\lim_{n\to\infty} \mathbf{x}^n)$.

**Theorem 1.11.** [Weierstrass' Theorem] Let $f : \mathbb{R}^d \to \mathbb{R}$ be a continuous function. Let $X \subseteq \mathbb{R}^d$ be a compact subset. Then $\inf\{f(\mathbf{x}) : \mathbf{x} \in X\}$ is attained, i.e., there exists $\mathbf{x}^{\min} \in X$ such that $f(\mathbf{x}^{\min}) = \inf\{f(\mathbf{x}) : \mathbf{x} \in X\}$. Similarly, there exists $\mathbf{x}^{\max} \in X$ such that $f(\mathbf{x}^{\max}) = \sup\{f(\mathbf{x}) : \mathbf{x} \in X\}$.

**Theorem 1.12.** Let $f : \mathbb{R}^d \to \mathbb{R}^n$ be a continuous function, and $C$ be a compact set. Then $f(C)$ is compact.

We will also need to speak of differentiability of functions $f : \mathbb{R}^d \to \mathbb{R}$.

**Definition 1.13.** We say that $f : \mathbb{R}^d \to \mathbb{R}$ is differentiable at $\mathbf{x} \in \mathbb{R}^d$, if there exists a linear transformation $A : \mathbb{R}^d \to \mathbb{R}$ such that
$$\lim_{\mathbf{h}\to\mathbf{0}} \frac{|f(\mathbf{x}+\mathbf{h}) - f(\mathbf{x}) - A\mathbf{h}|}{|\mathbf{h}|} = 0.$$

If $f$ is differentiable at $\mathbf{x}$, then the linear transformation is unique and is called the *gradient of $f$*. It is commonly denoted by $\nabla f(\mathbf{x})$.

**Definition 1.14.** The partial derivative of $f$ at $\mathbf{x}$ in the $i$-th direction is defined as the real number
$$f_i'(\mathbf{x}) := \lim_{h\to 0} \frac{f(\mathbf{x} + h\mathbf{e}^i) - f(\mathbf{x})}{h},$$

if the limit exists.

**Basic facts about matrices.** The set of $m \times n$ matrices will be denoted by $\mathbb{R}^{m\times n}$. The rank of a matrix $A$ will be denoted by $\mathrm{rk}(A)$ – it is the maximum number of linearly independent rows of $A$, which is equal to the maximum number of linearly independent columns of $A$. When $m = n$, we say that matrix is *square*.

**Definition 1.15.** A square matrix $A \in \mathbb{R}^{n\times n}$ is called symmetric if $A_{ij} = A_{ji}$ for all $i,j \in \{1,\dots,n\}$.

**Definition 1.16.** Let $A \in \mathbb{R}^{n\times n}$. A vector $\mathbf{v} \in \mathbb{R}^n$ is called an *eigenvector* of $A$, if there exists $\lambda \in \mathbb{R}$ such that $A\mathbf{v} = \lambda\mathbf{v}$. $\lambda$ is called the eigenvalue of $A$ associated with $\mathbf{v}$.

**Theorem 1.17.** If $A \in \mathbb{R}^{n\times n}$ is symmetric then it has $n$ orthogonal eigenvectors $\mathbf{v}^1,\dots,\mathbf{v}^n$ all of unit Euclidean norm, with associated eigenvalues $\lambda_1,\dots,\lambda_n \in \mathbb{R}$. Moreover, if $S$ is the matrix whose columns are $\mathbf{v}^1,\dots,\mathbf{v}^n$ and $\Lambda$ is the diagonal matrix with $\lambda_1,\dots,\lambda_n$ as the diagonal entries, then $A = S\Lambda S^T$. Moreover, $\mathrm{rk}(A)$ equals the number of nonzero eigenvalues.

NOTES: 6

**Theorem 1.18.** Let $A \in \mathbb{R}^{n \times n}$ be a symmetric matrix of rank $r$. The following are equivalent.

1. All eigenvalues of $A$ are nonnegative.

2. There exists a matrix $B \in \mathbb{R}^{r \times n}$ with linearly independent rows such that $A = B^T B$.

3. $\mathbf{u}^T A \mathbf{u} \geq 0$ for all $\mathbf{u} \in \mathbb{R}^n$.

**Definition 1.19.** A symmetric matrix $A \in \mathbb{R}^{n \times n}$ satisfying any of the three conditions in Theorem 1.18 is called a *positive semidefinite (PSD)* matrix. If $\mathrm{rk}(A) = n$, i.e., all its eigenvalues are strictly positive, then $A$ is called *positive definite*.

**Exercise 1.** Show that any positive definite matrix $A \in \mathbb{R}^{d \times d}$ defines a norm on $\mathbb{R}^d$ via $N_A(\mathbf{x}) = \sqrt{\mathbf{x}^T A \mathbf{x}}$. This norm is called the *norm induced by $A$*.

# 2 Convex Sets

## 2.1 Definitions and basic properties

A set $X \subseteq \mathbb{R}^d$ is called a *convex set* if for all $\mathbf{x}, \mathbf{y} \in X$, the line segment $[\mathbf{x}, \mathbf{y}]$ lies entirely in $X$. More precisely, for all $\mathbf{x}, \mathbf{y} \in X$ and every $\lambda \in [0, 1]$, $\lambda \mathbf{x} + (1 - \lambda) \mathbf{y} \in X$.

**Example 2.1.** Some examples of convex sets:

1. In $\mathbb{R}$, the only examples of convex sets are intervals (closed, open, half open): $(a, b), (a, b], [a, b], (-\infty, b]$ etc.

2. Let $\mathbf{a} \in \mathbb{R}^d$ and $\delta \in \mathbb{R}$. The sets $H(\mathbf{a}, \delta) = \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{a}, \mathbf{x} \rangle = \delta\}$, $H^+(\mathbf{a}, \delta) = \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{a}, \mathbf{x} \rangle \geq \delta\}$ and $H^-(\mathbf{a}, \delta) = \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{a}, \mathbf{x} \rangle \leq \delta\}$ are all convex sets. Sets of the form $H(\mathbf{a}, \delta)$ are called *hyperplanes* and sets of the form $H^+(\mathbf{a}, \delta), H^-(\mathbf{a}, \delta)$ are called *halfspaces*.

3. $\{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_\infty \leq 1\}$ is a convex set.

4. $\{\mathbf{x} = (x_1, \ldots, x_d) \in \mathbb{R}^d : x_1 + x_2 t + x_3 t^2 + \ldots + x_d t^{d-1} \geq 0 \text{ for all } t \geq 0\}$ is a convex set.

5. $\{(x, y) \in \mathbb{R}^2 : x^2 + y^2 \leq 5\}$ is convex. More generally, the ball $\{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\| \leq C\}$ for any $C \geq 0$ is convex.

**Exercise 2.** Show that if $N : \mathbb{R}^d \to \mathbb{R}$ is a norm, then every ball $B_N(\mathbf{x}, R)$ with respect to $N$ is convex.

**Definition 2.2.** Let $A \in \mathbb{R}^{d \times d}$ be a positive definite matrix . The set $\{x \in \mathbb{R}^d : \mathbf{x}^T A \mathbf{x} \leq 1\}$ is called an *ellipsoid*. In other words, an ellipsoid is the unit ball associated with the norm induced by $A$ – see Exercise 1. Exercise 2 shows that ellipsoids are convex.

NOTES:

7

**Theorem 2.3.** [Operations that preserve convexity] The following are all true.

1. Let $X_i$, $i \in I$ be an arbitrary family of convex sets. Then $\cap_{i \in I} X_i$ is a convex set.

2. Let $X$ be a convex set and $\alpha \in \mathbb{R}$, then $\alpha X$ is a convex set.

3. Let $X, Y$ be convex sets, then $X + Y$ is convex.

4. Let $T : \mathbb{R}^d \to \mathbb{R}^m$ be any linear transformation. If $X \subseteq \mathbb{R}^d$ is convex, then $T(X)$ is a convex set. If $Y \subseteq \mathbb{R}^m$ is convex, then $T^{-1}(Y)$ is convex.

*Proof.*     1. Let $\mathbf{x}, \mathbf{y} \in \cap_{i \in I} X_i$. This implies that $\mathbf{x}, \mathbf{y} \in X_i$ for every $i \in I$. Since each $X_i$ is convex, for every $\lambda \in [0, 1]$, $\lambda \mathbf{x} + (1 - \lambda)\mathbf{y} \in X_i$ for all $i \in I$. Therefore, $\lambda \mathbf{x} + (1 - \lambda)\mathbf{y} \in \cap_{i \in I} X_i$.

The proofs of 2., 3. and 4. are very similar, are left for the reader. $\qquad\square$

**Remark 2.4.** Observe that item 4. in Example 2.1 can be interpreted as an (uncountable) intersection of halfpsaces. Thus, item 2 from that example and Theorem 2.3 together give another proof that item 4. describes a convex set.

**Definition 2.5.** Let $Y = \mathbf{y}^1, \ldots, \mathbf{y}^n \in \mathbb{R}^d$ be a finite set of points. The *set of all convex combinations* of $Y$ is defined as
$$\{\lambda_1 \mathbf{y}^1 + \lambda_2 \mathbf{y}^2 + \ldots + \lambda_n \mathbf{y}^n : \lambda_i \geq 0, \ \ \lambda_1 + \lambda_2 + \ldots + \lambda_n = 1\}.$$

**Proposition 2.6.** If $X$ is convex and $\mathbf{y}^1, \ldots \mathbf{y}^n \in X$, then every convex combination of $\mathbf{y}^1, \ldots, \mathbf{y}^n$ is in $X$.

*Proof.* We prove it by induction on $n$. If $n = 1$, then the conclusion is trivial. Else consider any $\lambda_1, \ldots, \lambda_n \geq 0$ such that $\lambda_1 + \ldots \lambda_n = 1$. Then

$$
\begin{aligned}
&\lambda_1 \mathbf{y}^1 + \lambda_2 \mathbf{y}^2 + \ldots + \lambda_n \mathbf{y}^n \\
=\ & (\lambda_1 + \ldots + \lambda_{n-1})(\tfrac{\lambda_1}{\lambda_1 + \ldots + \lambda_{n-1}} \mathbf{y}^1 + \tfrac{\lambda_2}{\lambda_1 + \ldots + \lambda_{n-1}} \mathbf{y}^2 + \ldots + \tfrac{\lambda_{n-1}}{\lambda_1 + \ldots + \lambda_{n-1}} \mathbf{y}^{n-1}) + \lambda_n \mathbf{y}^n \\
=\ & (1 - \lambda_n)\tilde{\mathbf{y}} + \lambda_n \mathbf{y}^n
\end{aligned}
$$

where $\tilde{\mathbf{y}} := \frac{\lambda_1}{\lambda_1 + \ldots + \lambda_{n-1}} \mathbf{y}^1 + \frac{\lambda_2}{\lambda_1 + \ldots + \lambda_{n-1}} \mathbf{y}^2 + \ldots + \frac{\lambda_{n-1}}{\lambda_1 + \ldots + \lambda_{n-1}} \mathbf{y}^{n-1}$ belongs to $X$ by the induction hypothesis. The rest follows from definition of convexity. $\qquad\square$

**Definition 2.7.** Given any set $X \subseteq \mathbb{R}^d$ (not necessarily convex), the convex hull of $X$, denoted by $\mathrm{conv}(X)$, is a convex set $C$ such that $X \subseteq C$ and for any other convex set $C'$, $X \subseteq C' \Rightarrow C \subseteq C'$, i.e., the convex hull of $X$ is the smallest (with respect to set inclusion) convex set containing $X$.

**Theorem 2.8.** For any set $X \subseteq \mathbb{R}^d$ (not necessarily convex),

$$\mathrm{conv}(X) = \bigcap (C : X \subseteq C, C \text{ convex}) = \{\lambda_1 x_1 + \ldots + \lambda_t x_t : x_1, \ldots, x_t \in X, \lambda_1, \ldots, \lambda_t \geq 0, \sum_{i=1}^{t} \lambda_i = 1\}.$$

In other words, the convex hull of $X$ is the union of the set of convex combinations of all possible finite subsets of $X$.

NOTES:                                                   8

*Proof.* Let $\hat{C} = \bigcap (C : X \subseteq C, C \text{ convex})$, which is a convex set by Theorem 2.3 and by definition $X \subseteq \hat{C}$. Consider any other convex set $C'$ such that $X \subseteq C'$. Then $C'$ appears in the intersection, and thus $\hat{C} \subseteq C'$. Thus, $\hat{C} = \text{conv}(X)$.

Next, let $\tilde{C} = \{\lambda_1 x_1 + \ldots + \lambda_t x_t : x_1, \ldots, x_t \in X, \lambda_1, \ldots, \lambda_t \geq 0, \sum_{i=1}^{t} \lambda_i = 1\}$. Then,

1. $\tilde{C}$ is convex. Consider two points $z_1, z_2 \in \tilde{C}$. Thus there exist two finite index sets $I_1, I_2$, two finite subsets of $X$ given by $X_1 = \{x_i^1 \in X : i \in I_1\}$ and $X_2 = \{x_i^2 \in X : i \in I_2\}$, and two subsets of nonnegative real numbers $\{\lambda_i^1 \geq 0, i \in I_1\}$, $\{\lambda_i^2 \geq 0, i \in I_2\}$ such that $\sum_{i \in I_j} \lambda_i^j = 1$ for $j = 1, 2$, with the following property : $z_j = \sum_{i \in I_j} \lambda_i^j x_i^j$ for $j = 1, 2$. Then for any $\lambda \in [0, 1]$, $\lambda z_1 + (1 - \lambda) z_2 = \lambda(\sum_{i \in I_1} \lambda_i^1 x_i^1) + (1 - \lambda)(\sum_{i \in I_2} \lambda_i^2 x_i^2)$. Consider the finite set $\tilde{X} = X_1 \cup X_2$, and for each $x \in \tilde{X}$, if $x = x_i \in X_1$ with $i \in I_1$ let $\mu_x = \lambda \cdot \lambda_i^1$, and if $x = x_i \in X_2$ with $i \in I_2$, let $\mu_x = (1 - \lambda) \cdot \lambda_i^2$. It is easy to check that $\sum_{x \in \tilde{X}} \mu_x = 1$, and $\lambda z_1 + (1 - \lambda) z_2 = \sum_{x \in \tilde{X}} \mu_x x$. Thus, $\lambda z_1 + (1 - \lambda) z_2 \in \tilde{C}$.

2. $X \subseteq \tilde{C}$. We simply use $\lambda = 1$ as the multiplier for a point from $X$.

3. Let $C'$ be any convex set such that $X \subseteq C'$. Since $C'$ is convex, every point of the form $\lambda_1 x_1 + \ldots + \lambda_t x_t$ where $x_1, \ldots, x_t \in X$, $\lambda_i \geq 0$, $\sum_{i=1}^{t} \lambda_i = 1$ belongs to $C'$ by Proposition 2.6. Thus, $\tilde{C} \subseteq C'$.

From 1., 2. and 3., we get that $\tilde{C} = \text{conv}(X)$. $\qquad\square$

## 2.2  Convex cones, affine sets and dimension

We say $X$ is convex if for all $\mathbf{x}, \mathbf{y} \in X$ and $\lambda, \gamma \geq 0$ such that $\lambda + \gamma = 1$, $\lambda \mathbf{x} + \gamma \mathbf{y} \in X$. What happens if we relax the conditions on $\lambda, \gamma$?

**Definition 2.9.** We have three possibilities:

1. We say that $X \subseteq \mathbb{R}^d$ is a *convex cone* if for all $\mathbf{x}, \mathbf{y} \in X$ and $\lambda, \gamma \geq 0$, $\lambda \mathbf{x} + \gamma \mathbf{y} \in X$.

2. We say that $X \subseteq \mathbb{R}^d$ is an *affine set* or an *affine subspace*, if for all $\mathbf{x}, \mathbf{y} \in X$ and $\lambda, \gamma \in \mathbb{R}$ such that $\lambda + \gamma = 1$, $\lambda \mathbf{x} + \gamma \mathbf{y} \in X$.

3. We say $X \subseteq \mathbb{R}^d$ is a *linear set* or a *linear subspace* if for all $\mathbf{x}, \mathbf{y} \in X$ and $\lambda, \gamma \in \mathbb{R}$, $\lambda \mathbf{x} + \gamma \mathbf{y} \in X$.

**Remark 2.10.** Since we relaxed the conditions on $\lambda, \gamma$, convex cones, affine sets and linear sets are all special cases of convex sets.

Similar to the definition of the convex hull of an arbitrary subset $X$, one can define the *conical hull* of $X$ as the set inclusion wise smallest convex cone containing $X$ denoted by $\text{cone}(X)$. Similarly, the *affine (linear) hull* of $X$ as the set inclusion wise smallest affine (linear) set containing $X$. The affine hull will be be denoted by $\text{aff}(X)$, and linear hull will be denoted by $\text{span}(X)$. One can verify the following analog of Theorem 2.8.

NOTES:                                    9

**Theorem 2.11.** Let $X \subseteq \mathbb{R}^n$. The following are all true.

1. $\text{cone}(X) = \bigcap(C : X \subseteq C, C \text{ is a convex cone}) = \{\lambda_1 x_1 + \ldots + \lambda_t x_t : x_1, \ldots, x_t \in X, \lambda_1, \ldots, \lambda_t \geq 0\}$.

2. $\text{aff}(X) = \bigcap(C : X \subseteq C, C \text{ is an affine set}) = \{\lambda_1 x_1 + \ldots + \lambda_t x_t : x_1, \ldots, x_t \in X, \sum_{i=1}^{t} \lambda_i = 1\}$.

3. $\text{span}(X) = \bigcap(C : X \subseteq C, C \text{ is a linear subspace}) = \{\lambda_1 x_1 + \ldots + \lambda_t x_t : x_1, \ldots, x_t \in X, \lambda_1, \ldots, \lambda_t \in \mathbb{R}\}$.

The following is a fundamental theorem of linear algebra.

**Theorem 2.12.** Let $X \subseteq \mathbb{R}^d$. The following are equivalent.

1. $X$ is a linear subspace.

2. There exists $0 \leq m \leq d$ and linearly independent vectors $\mathbf{v}^1, \ldots, \mathbf{v}^m \in X$ such that every $\mathbf{x} \in X$ can be written as $\mathbf{x} = \lambda_1 \mathbf{v}^1 + \ldots + \lambda_m \mathbf{v}^m$ for some reals $\lambda_i$, $i = 1, \ldots, m$, i.e., $X = \text{span}(\{\mathbf{v}^1, \ldots, \mathbf{v}^m\})$.

3. There exists a matrix $A \in \mathbb{R}^{(d-m) \times d}$ with full row rank such that $X = \{\mathbf{x} \in \mathbb{R}^d : A\mathbf{x} = \mathbf{0}\}$.

*Proof sketch.* We take for granted the fact that we can have at most $d$ linearly independent vectors in $\mathbb{R}^d$. This is something one can show using Gaussian elimination.

It is easy to verify that 2. $\Rightarrow$ 1. (because linear combinations of linear combinations are linear combinations). To see that 1. $\Rightarrow$ 2., starting with a linear subspace $X$, we construct a finite set $\mathbf{v}^1, \ldots, \mathbf{v}^m \in X$ satisfying the conditions of 2. We do this in an iterative fashion. Start by picking any arbitrary $\mathbf{v}^1 \in X$. If $X = \text{span}(\mathbf{v}^1)$, then we are done. Else, choose $\mathbf{v}^2 \in X \setminus \text{span}(\mathbf{v}^1)$. Again, if $X = \text{span}(\mathbf{v}^1, \mathbf{v}^2)$ then we are done, else choose $\mathbf{v}^3 \in X \setminus \text{span}(\mathbf{v}^1, \mathbf{v}^2)$. This process has to end after at most $d$ steps, because we cannot have more than $d$ linearly independent vectors in $\mathbb{R}^d$.

It is easy to verify 3. $\Rightarrow$ 1. To see that 1. $\Rightarrow$ 3., define the set $X^\perp := \{\mathbf{y} \in \mathbb{R}^d : \langle \mathbf{y}, \mathbf{x} \rangle = 0 \ \forall \mathbf{x} \in X\}$ (this is known as the *orthogonal complement* of $X$). It can be verified that $X^\perp$ is a linear subspace. Moreover, by the equivalence 1. $\Leftrightarrow$ 2., we know that 2. holds for $X^\perp$. So there exist linearly independent vectors $\mathbf{a}^1, \ldots, \mathbf{a}^k$ for some $0 \leq k \leq d$ such that $X^\perp = \text{span}(\mathbf{a}^1, \ldots, \mathbf{a}^k)$. Let $A$ be the $k \times d$ matrix which has $\mathbf{a}^1, \ldots, \mathbf{a}^k$ as rows. One can now verify that $X = \{\mathbf{x} \in \mathbb{R}^d : A\mathbf{x} = \mathbf{0}\}$. The fact that one can take $k = d - m$ where $m$ is the number from condition 2. needs additional work, which we skip here. $\qquad\square$

**Definition 2.13.** The number $m$ showing up in item 2. in the above theorem is called the *dimension* of $X$. The set of vectors $\{\mathbf{v}^1, \ldots, \mathbf{v}^m\}$ are called a *basis* for the linear subspace.

There is an analogous theorem for affine sets. For this, we need the concept of *affine independence* that is analogous to the concept of linear independence.

**Definition 2.14.** We say a set $X$ is affinely independent if there does not exist $\mathbf{x} \in X$ such that $\mathbf{x} \in \text{aff}(X \setminus \{\mathbf{x}\})$.

We now give several characterizations of affine independence.

NOTES: 10

**Proposition 2.15.** Let $X \subseteq \mathbb{R}^d$. The following are equivalent.

1. $X$ is an affinely independent set.

2. For every $\mathbf{x} \in X$, the set $\{\mathbf{v} - \mathbf{x} : \mathbf{v} \in X \setminus \{\mathbf{x}\}\}$ is linearly independent.

3. There exists $\mathbf{x} \in X$ such that the set $\{\mathbf{v} - \mathbf{x} : \mathbf{v} \in X \setminus \{\mathbf{x}\}\}$ is linearly independent.

4. The set of vectors $\{(\mathbf{x}, 1) \in \mathbb{R}^{d+1} : \mathbf{x} \in X\}$ is linearly independent.

5. $X$ is a finite set with vectors $\mathbf{x}^1, \ldots, \mathbf{x}^m$ such that $\lambda_1 \mathbf{x}^1 + \ldots + \lambda_m \mathbf{x}^m = 0, \lambda_1 + \ldots + \lambda_m = 0$ implies $\lambda_1 = \lambda_2 = \ldots = \lambda_m = 0$.

*Proof.* 1. $\Rightarrow$ 2. Consider an arbitrary $\mathbf{x} \in X$. Suppose to the contrary that $\{\mathbf{v} - \mathbf{x} : \mathbf{v} \in X \setminus \{\mathbf{x}\}\}$ is not linearly independent, i.e., there exist multipliers $\lambda_{\mathbf{v}}$, not all zero, such that $\sum_{\mathbf{v} \in X \setminus \{\mathbf{x}\}} \lambda_{\mathbf{v}} (\mathbf{v} - \mathbf{x}) = 0$. Rearranging terms, we get $\sum_{\mathbf{v} \in X \setminus \{\mathbf{x}\}} \lambda_{\mathbf{v}} \mathbf{v} = (\sum_{\mathbf{v} \in X \setminus \{\mathbf{x}\}} \lambda_v) \mathbf{x}$. We now consider two cases:

*Case 1:* $\sum_{\mathbf{v} \in X \setminus \{\mathbf{x}\}} \lambda_{\mathbf{v}} = 0$. In this case, since not all the $\lambda_{\mathbf{v}}$ are zero, let $\bar{\mathbf{v}} \in X \setminus \{\mathbf{x}\}$ be such that $\lambda_{\bar{\mathbf{v}}} \neq 0$. Since $\sum_{\mathbf{v} \in X \setminus \{\mathbf{x}\}} \lambda_{\mathbf{v}} \mathbf{v} = (\sum_{\mathbf{v} \in X \setminus \{\mathbf{x}\}} \lambda_v) \mathbf{x} = 0$, we obtain that $\bar{\mathbf{v}} = \sum_{\mathbf{v} \in X \setminus \{\mathbf{x}, \bar{\mathbf{v}}\}} \frac{\lambda_{\mathbf{v}}}{-\lambda_{\bar{\mathbf{v}}}} \mathbf{v}$. Since $\sum_{\mathbf{v} \in X \setminus \{\mathbf{x}\}} \lambda_{\mathbf{v}} = 0$, this shows that $\bar{v} \in \mathrm{aff}(X \setminus \{\mathbf{x}, \mathbf{v}\})$, contradicting the assumption that $X$ is affinely independent.

*Case 2:* $\sum_{\mathbf{v} \in X \setminus \{\mathbf{x}\}} \lambda_{\mathbf{v}} \neq 0$. We can write $\mathbf{x} = \sum_{\mathbf{v} \in X \setminus \{\mathbf{x}\}} \frac{\lambda_{\mathbf{v}}}{\sum_{\mathbf{v} \in X \setminus \{\mathbf{x}\}} \lambda_v} \mathbf{v}$. This implies that $\mathbf{x} \in \mathrm{aff}(X \setminus \{\mathbf{x}\})$ contradicting the assumption that $X$ is affinely independent.

2. $\Rightarrow$ 3. Obvious.

3. $\Rightarrow$ 4. Let $\bar{\mathbf{x}}$ be such that $\{\mathbf{v} - \bar{\mathbf{x}} : \mathbf{v} \in X \setminus \{\bar{\mathbf{x}}\}\}$ is linearly independent. This means that the vectors $\{(\mathbf{v} - \bar{\mathbf{x}}, 0) : \mathbf{v} \in X \setminus \{\bar{\mathbf{x}}\}\} \cup \{(\bar{\mathbf{x}}, 1)\}$ are also linearly independent. Thus the matrix with these vectors as columns has full column rank. Now if we add the the column $(\bar{\mathbf{x}}, 1)$ to the rest of the columns, this does not change the column rank, and thus the columns remain linearly independent. But the new matrix has precisely $\{(\mathbf{x}, 1) \in \mathbb{R}^{d+1} : \mathbf{x} \in X\}$ as its columns.

4. $\Rightarrow$ 5. Follows from the fact that if $\{(\mathbf{x}, 1) \in \mathbb{R}^{d+1} : \mathbf{x} \in X\}$ is linearly independent, then the set $X$ must be finite. Moreover, if $\sum_{\mathbf{x} \in X} \lambda_{\mathbf{x}} (\mathbf{x}, 1) = 0$ for some real numbers $\{\lambda_{\mathbf{x}}\}_{\mathbf{x} \in X}$, then $\lambda_{\mathbf{x}} = 0$ for all $\mathbf{x} \in X$.

5. $\Rightarrow$ 1. Consider any $\mathbf{x}^i \in X$. If $\mathbf{x}^i \in \mathrm{aff}(X \setminus \{\mathbf{x}^i\})$, then there exist multipliers $\lambda_j \in \mathbb{R}$, $j \neq i$ such that $\mathbf{x}^i = \sum_{j \neq i} \lambda_j \mathbf{x}^j$ and $\sum_{j \neq i} \lambda_j = 1$. This implies that $\sum_{j=1}^m \lambda_j \mathbf{x}^j = 0$ where $\lambda_i = -1$, and therefore $\lambda_1 + \ldots + \lambda_m = 0$, contradicting the hypothesis of 5. $\square$

We are now ready to state the affine version of Theorem 2.12.

**Theorem 2.16.** Let $X \subseteq \mathbb{R}^d$. The following are equivalent.

1. $X$ is an affine subspace.

2. There exists a linear subspace $L$ of dimension $0 \le m \le d$, such that $X - \mathbf{x} = L$ for every $\mathbf{x} \in X$.

3. There exist affinely independent vectors $\mathbf{v}^1, \ldots, \mathbf{v}^{m+1} \in X$ for $0 \le m \le d$ such that every $\mathbf{x} \in X$ can be written as $\mathbf{x} = \lambda_1 \mathbf{v}^1 + \ldots + \lambda_{m+1} \mathbf{v}^{m+1}$ for some reals $\lambda_i$, $i = 1, \ldots, m+1$ such that $\lambda_1 + \ldots + \lambda_{m+1} = 1$, i.e., $X = \text{aff}(\{\mathbf{v}^1, \ldots, \mathbf{v}^{m+1}\})$.

4. There exists a matrix $A \in \mathbb{R}^{(d-m) \times d}$ with full row rank and a vector $\mathbf{b} \in \mathbb{R}^{d-m}$ for some $0 \le m \le d$ such that $X = \{\mathbf{x} \in \mathbb{R}^d : A\mathbf{x} = \mathbf{b}\}$.

*Proof.* 1. $\Rightarrow$ 2. Fix an arbitrary $\mathbf{x}^\star \in X$. Define $L = X - \mathbf{x}^\star$. We first show that $L$ is a linear subspace: for any $\mathbf{y}^1, \mathbf{y}^2 \in X$, $\lambda(\mathbf{y}^1 - \mathbf{x}^\star) + \gamma(\mathbf{y}^2 - \mathbf{x}^\star) \in X - \mathbf{x}^\star$ for any $\lambda, \gamma \in \mathbb{R}$. Since $\lambda(\mathbf{y}^1 - \mathbf{x}^\star) + \gamma(\mathbf{y}^2 - \mathbf{x}^\star) + \mathbf{x}^\star = \lambda \mathbf{y}^1 + \gamma \mathbf{y}^2 + (1 - \lambda - \gamma)\mathbf{x}^\star$ and $X$ is an affine subset, therefore, $\lambda(\mathbf{y}^1 - \mathbf{x}^\star) + \gamma(\mathbf{y}^2 - \mathbf{x}^\star) + \mathbf{x}^\star \in X$. So, $\lambda(\mathbf{y}^1 - \mathbf{x}^\star) + \gamma(\mathbf{y}^2 - \mathbf{x}^\star) \in L$. Now, for any other $\bar{\mathbf{x}} \in X$, we need to show that $L = X - \bar{\mathbf{x}}$. Consider any $\mathbf{y} \in L$, i.e., $\mathbf{y} = \mathbf{x} - \mathbf{x}^\star$ for some $\mathbf{x} \in X$. Observe that $\mathbf{y} = (\mathbf{x} + \bar{\mathbf{x}} - \mathbf{x}^\star) - \bar{\mathbf{x}}$ and $\mathbf{x} + \bar{\mathbf{x}} - \mathbf{x}^\star \in X$ (because the coefficients all sum to 1). Therefore, $\mathbf{y} \in X - \bar{\mathbf{x}}$ showing that $L = X - \mathbf{x}^\star \subseteq X - \bar{\mathbf{x}}$. Switching the roles of $\mathbf{x}^\star$ and $\bar{\mathbf{x}}$, one can similarly show that $X - \bar{\mathbf{x}} \subseteq X - \mathbf{x}^\star = L$.

2. $\Rightarrow$ 1. Consider any $\mathbf{y}^1, \mathbf{y}^2 \in X$ and let $\lambda, \gamma \in \mathbb{R}$ such that $\lambda + \gamma = 1$. We need to show that $\lambda \mathbf{y}^1 + \gamma \mathbf{y}^2 \in X$. Since $X - \mathbf{y}^1$ is a linear subspace, $\gamma(\mathbf{y}^2 - \mathbf{y}^1) \in X - \mathbf{y}^1$. Thus, $\gamma(\mathbf{y}^2 - \mathbf{y}^1) + \mathbf{y}^1 = \lambda \mathbf{y}^1 + \gamma \mathbf{y}^2 \in X$.

The equivalence of 2., 3. and 4. follows from Theorem 2.12. $\qquad\square$

**Definition 2.17** (Dimension of convex sets)**.** If $X$ is an affine subspace and $\mathbf{x} \in X$, the linear subspace $X - \mathbf{x}$ is called the *linear subspace parallel to $X$* and the dimension of $X$ is the dimension of the linear subspace $X - \mathbf{x}$. For any convex set $X$, the dimension of $X$ is the dimension of $\text{aff}(X)$ and will be denoted by $\dim(X)$.

**Lemma 2.18.** If $X$ is a set of affinely independent points, then $\dim(\text{aff}(X)) = |X| - 1$.

*Proof.* Fix any $\mathbf{x} \in X$. By Theorem 2.16, $L = \text{aff}(X) - \mathbf{x}$ is a linear subspace. We claim that $(X \setminus \{\mathbf{x}\}) - \mathbf{x}$ is a basis for $L$. The verification of this claim is left to the reader. $\qquad\square$

**Proposition 2.19.** Let $X$ be a convex set. $\dim(X)$ equals one less than the maximum number of affinely independent points in $X$.

*Proof.* Let $X_0 \subseteq X$ be a maximum sized set of affinely independent points in $X$. By Problem 5 in HW I, $\text{aff}(X_0) \subseteq \text{aff}(X)$. Since $X_0$ is a maximum sized set of affinely independent points in $X$, any $\mathbf{x} \in X$ must lie in $\text{aff}(X_0)$. Therefore, $X \subseteq \text{aff}(X_0)$. Since $\text{aff}(X_0)$ is an affine set, by definition of affine hull of $X$, we have $\text{aff}(X) \subseteq \text{aff}(X_0)$. Therefore, $\text{aff}(X) = \text{aff}(X_0)$, implying that $\dim(\text{aff}(X_0)) = \dim(\text{aff}(X))$. By Lemma 2.18, we thus obtain $|X_0| - 1 = \dim(\text{aff}(X))$. $\qquad\square$

## 2.3 Representations of convex sets

A large part of modern convex geometry is concerned with algorithms for computing with or optimizing over convex sets. For algorithmic purposes, we need ways to describe a convex set, so that it can be stored in a computer compactly and computations can be performed with it.

NOTES: 12

### 2.3.1 Extrinsic description: separating hyperplanes

Perhaps the most primitive convex set in $\mathbb{R}^d$ is the halfspace – see item 2. in Example 2.1. Moreover, a halfspace is a *closed* convex set. By Theorem 2.3, the intersection of an arbitrary family of halfspaces is a closed convex set. Perhaps the most fundamental theorem of convexity is that the converse is true.

**Theorem 2.20** (Separating Hyperplane Theorem)**.** Let $C \subseteq \mathbb{R}^d$ be a closed convex set and let $\mathbf{x} \notin C$. There exists a halfspace that contains $C$ and does not contain $\mathbf{x}$. More precisely, there exists $\mathbf{a} \in \mathbb{R}^d \setminus \{\mathbf{0}\}, \delta \in \mathbb{R}$ such that $\langle \mathbf{a}, \mathbf{y} \rangle \leq \delta$ for all $\mathbf{y} \in C$ and $\langle \mathbf{a}, \mathbf{x} \rangle > \delta$. The hyperplane $\{\mathbf{y} \in \mathbb{R}^d : \langle \mathbf{a}, \mathbf{y} \rangle = \delta\}$ is called a *separating hyperplane* for $C$ and $\mathbf{x}$.

*Proof.* If $C$ is empty, then any halfspace that does not contain $\mathbf{x}$ suffices. Otherwise, consider any $\bar{\mathbf{x}} \in C$ and let $r = \|\mathbf{x} - \bar{\mathbf{x}}\|$. Let $\bar{C} = C \cap B(\mathbf{x}, r)$. Since $C$ is closed and $B(\mathbf{x}, r)$ is compact, $\bar{C}$ is compact. One can also verify that the function $f(\mathbf{y}) = \|\mathbf{y} - \mathbf{x}\|$ is a continuous function on $\mathbb{R}^d$. Therefore, by Weierstrass' Theorem (Theorem 1.11), there exists $\mathbf{x}^\star \in \bar{C}$ such that $\|\mathbf{x} - \mathbf{x}^\star\| \leq \|\mathbf{x} - \mathbf{y}\|$ for all $\mathbf{y} \in \bar{C}$, and therefore in fact $\|\mathbf{x} - \mathbf{x}^\star\| \leq \|\mathbf{x} - \mathbf{y}\|$ for all $\mathbf{y} \in C$.

Let $\mathbf{a} = \mathbf{x} - \mathbf{x}^\star$ and let $\delta = \langle \mathbf{a}, \mathbf{x}^\star \rangle$. Note that $\mathbf{a} \neq \mathbf{0}$ because $\mathbf{x} \notin C$ and $\mathbf{x}^\star \in C$. Also note that $\langle \mathbf{a}, \mathbf{x} \rangle = \langle \mathbf{a}, \mathbf{a} + \mathbf{x}^\star \rangle = \|\mathbf{a}\|^2 + \delta > \delta$. Thus, it remains to check that $\langle \mathbf{a}, \mathbf{y} \rangle \leq \delta$ for all $\mathbf{y} \in C$. For any $\mathbf{y} \in C$, all the points $\alpha \mathbf{y} + (1 - \alpha)\mathbf{x}^\star$, $\alpha \in (0, 1)$ are in $C$ by convexity. Therefore, by the extremal property of $\mathbf{x}^\star$, we have

$$
\begin{array}{rrcll}
& \|\mathbf{x} - \mathbf{x}^\star\|^2 & \leq & \|\mathbf{x} - (\alpha \mathbf{y} + (1 - \alpha)\mathbf{x}^\star)\|^2 & \forall \alpha \in (0, 1) \\
\Rightarrow & 0 & \leq & \alpha^2 \|\mathbf{y} - \mathbf{x}^\star\|^2 - 2\alpha \langle \mathbf{x} - \mathbf{x}^\star, \mathbf{y} - \mathbf{x}^\star \rangle & \forall \alpha \in (0, 1) \\
\Rightarrow & 2\langle \mathbf{x} - \mathbf{x}^\star, \mathbf{y} - \mathbf{x}^\star \rangle & \leq & \alpha \|\mathbf{y} - \mathbf{x}^\star\|^2 & \forall \alpha \in (0, 1)
\end{array}
$$

Letting $\alpha \to 0$ in the last inequality yields that $0 \geq \langle \mathbf{x} - \mathbf{x}^\star, \mathbf{y} - \mathbf{x}^\star \rangle = \langle \mathbf{a}, \mathbf{y} - \mathbf{x}^\star \rangle$. Thus, $\langle \mathbf{a}, \mathbf{y} \rangle \leq \langle \mathbf{a}, \mathbf{x}^\star \rangle = \delta$ for all $\mathbf{y} \in C$. $\square$

**Corollary 2.21.** Every closed convex set can be written as the intersection of some family of halfpsaces. In other words, a subset $X \subseteq \mathbb{R}^d$ is a closed convex set if and only if there exists a family of tuples $(\mathbf{a}^i, \delta^i)$, $i \in I$ (where $I$ may be an uncountable index set) such that $X = \cap_{i \in I} H^-(\mathbf{a}^i, \delta^i)$.

**Definition 2.22.** A *finite* intersection of halfpsaces is called a *polyhedron*. In other words, $P \subseteq \mathbb{R}^d$ is a polyhedron if and only if there exist vectors $\mathbf{a}^1, \dots, \mathbf{a}^m \in \mathbb{R}^d$ and real numbers $b^1, \dots, b^m$ such that $P = \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{a}^i, \mathbf{x} \rangle \leq b^i \ \ i = 1, \dots, m\}$. The shorthand $P = \{\mathbf{x} \in \mathbb{R}^d : A\mathbf{x} \leq \mathbf{b}\}$ is often employed, where $A$ is the $m \times d$ matrix with $\mathbf{a}^1, \dots, \mathbf{a}^m$ as rows, and $\mathbf{b} = (b^1, \dots, b^m) \in \mathbb{R}^m$.

Thus, a polyhedron is completely described by specifying a matrix $A \in \mathbb{R}^{m \times d}$ and a vector $b \in \mathbb{R}^m$.

**Question 1.** How would one show that the unit ball for the standard Euclidean norm in $\mathbb{R}^d$ is **not** a polyhedron?

Another related, and very useful, result is the following.

NOTES: 13

**Theorem 2.23** (Supporting Hyperplane Theorem). Let $C \subseteq \mathbb{R}^d$ be a convex set and let $\mathbf{x} \in \mathrm{bd}(C)$. Then, there exists $\mathbf{a} \in \mathbb{R}^d \setminus \{\mathbf{0}\}, \delta \in \mathbb{R}$ such that $\langle \mathbf{a}, \mathbf{y} \rangle \leq \delta$ for all $\mathbf{y} \in C$ and $\langle \mathbf{a}, \mathbf{x} \rangle = \delta$. The hyperplane $\{\mathbf{y} \in \mathbb{R}^d : \langle \mathbf{a}, \mathbf{y} \rangle = \delta\}$ is called a *supporting hyperplane* for $C$ at $\mathbf{x}$.

*Proof.* Since $\mathrm{bd}(C) = \mathrm{bd}(\mathbb{R}^d \setminus \mathrm{cl}(C))$, $\mathbf{x} \in \mathrm{bd}(\mathbb{R}^d \setminus \mathrm{cl}(C))$. Since $\mathbb{R}^d \setminus \mathrm{cl}(C)$ is an open set, there exists a sequence $\{\mathbf{x}^i\}_{i \in \mathbb{N}}$ such that $\mathbf{x}^i \to \mathbf{x}$ and each $\mathbf{x}^i \notin \mathrm{cl}(C)$. By Theorem 2.20, for each $\mathbf{x}^i$, there exists $\mathbf{a}^i$ such that $\langle \mathbf{a}^i, \mathbf{y} \rangle < \langle \mathbf{a}^i, \mathbf{x}^i \rangle$ for all $\mathbf{y} \in C$. By scaling the vectors $\mathbf{a}^i$, we can assume that $\|\mathbf{a}^i\| = 1$ for all $i \in \mathbb{N}$.

Since the set of unit norm vectors is a compact set, by Theorem 1.10, one can pick a convergent subsequence $\mathbf{a}^{i_k} \to \mathbf{a}$ such that $\langle \mathbf{a}^{i_k}, \mathbf{y} \rangle < \langle \mathbf{a}^{i_k}, \mathbf{x}^{i_k} \rangle$ for all $\mathbf{y} \in C$. Taking the limit on both sides, we obtain $\langle \mathbf{a}, \mathbf{y} \rangle \leq \langle \mathbf{a}, \mathbf{x} \rangle$ for all $\mathbf{y} \in C$. We simply set $\delta = \langle \mathbf{a}, \mathbf{x} \rangle$. Note also that since $\|\mathbf{a}^i\| = 1$ for all $i \in \mathbb{N}$, we must have $\|\mathbf{a}\| = 1$, and so $\mathbf{a} \neq \mathbf{0}$. $\square$

**How to represent general convex sets: Separation oracles.** We have seen that polyhedra can be represented by a matrix $A$ and a right hand side $b$. Norm balls can be represented by the center $\mathbf{x}$ and the radius $R$. Ellipsoids can be represented by PD matrices $A$. What about general convex sets? This problem is gotten around by assuming that one has "black-box" access to the convex set via a *separation oracle*. More formally, we say that a convex set $C \subseteq \mathbb{R}^d$ is equipped with a separation oracle $O$ that takes as input any vector $\mathbf{x} \in \mathbb{R}^d$ and gives the following output: If $\mathbf{x} \in C$, the output is "YES", and if $\mathbf{x} \notin C$, then the output is a tuple $(\mathbf{a}, \delta) \in \mathbb{R}^d \times \mathbb{R}$ such that $\{\mathbf{y} \in \mathbb{R}^d : \langle \mathbf{a}, \mathbf{y} \rangle = \delta\}$ is a separating hyperplane for $\mathbf{x}$ and $C$.

**Farkas' lemma: A glimpse into polyhedral theory.** A nice characterization of solutions to systems of linear equations is given in linear algebra, which can be viewed as the most basic type of "theorem of the alternative".

**Theorem 2.24.** Let $A \in \mathbb{R}^{d \times n}$ and $\mathbf{b} \in \mathbb{R}^d$. Exactly one of the following is true.

1. $A\mathbf{x} = \mathbf{b}$ has a solution.

2. There exists $\mathbf{u} \in \mathbb{R}^d$ such that $\mathbf{u}^T A = \mathbf{0}$ and $\mathbf{u}^T \mathbf{b} \neq 0$.

What if we are interested in *nonnegative solutions* to linear equations? Farkas' lemma is a characterization of such solutions.

**Theorem 2.25.** [Farkas' Lemma] Let $A \in \mathbb{R}^{d \times n}$ and $\mathbf{b} \in \mathbb{R}^d$. Exactly one of the following is true.

1. $A\mathbf{x} = \mathbf{b}, \ \mathbf{x} \geq \mathbf{0}$ has a solution.

2. There exists $\mathbf{u} \in \mathbb{R}^d$ such that $\mathbf{u}^T A \leq \mathbf{0}$ and $\mathbf{u}^T \mathbf{b} > 0$.

Before we dive into the proof of Farkas' Lemma, we need a technical result.

**Lemma 2.26.** Let $\mathbf{a}^1, \ldots, \mathbf{a}^n \in \mathbb{R}^d$. Then $\mathrm{cone}(\{\mathbf{a}^1, \ldots, \mathbf{a}^n\})$ is closed.

NOTES: 14

*Proof.* We will complete the proof of this lemma when we do Caratheodory's theorem (See the end of Section 2.4). $\qquad\square$

*Proof of Theorem 2.25.* Let $\mathbf{a}^1, \dots, \mathbf{a}^n \in \mathbb{R}^d$ be the columns of the matrix $A$. By Lemma 2.26, the cone $C = \{A\mathbf{x} : \mathbf{x} \geq 0\}$ is closed. Then we have two cases, either $\mathbf{b} \in C$ or $\mathbf{b} \notin C$. In the first case, we end up in Case 1 of the statement of the theorem. In the second case, by Theorem 2.20, there exists $\mathbf{u} \in \mathbb{R}^d$ and $\delta \in \mathbb{R}$ such that $\langle \mathbf{u}, \mathbf{y} \rangle \leq \delta$ for all $\mathbf{y} \in C$ and $\langle \mathbf{u}, \mathbf{b} \rangle > \delta$. Since $\mathbf{0} \in C$, we must have $\delta \geq \langle \mathbf{u}, \mathbf{0} \rangle = 0$. This already shows that $\langle \mathbf{u}, \mathbf{b} \rangle > 0$.

Now suppose to the contrary that for some $\mathbf{a}^i$, $\langle \mathbf{u}, \mathbf{a}^i \rangle > 0$. Thus, there exists $\bar{\lambda} \geq 0$ such that $\bar{\lambda} \langle \mathbf{u}, \mathbf{a}^i \rangle > \delta$ (for example, take $\bar{\lambda} = \frac{|\delta| + 1}{\langle \mathbf{u}, \mathbf{a}^i \rangle}$). Since $\mathbf{y} := \bar{\lambda} \mathbf{a}^i \in C$, this implies that $\langle \mathbf{u}, \mathbf{y} \rangle > \delta$, contradicting that $\langle \mathbf{u}, \mathbf{y} \rangle \leq \delta$ for all $\mathbf{y} \in C$. $\qquad\square$

**Duality/Polarity.** With every linear space, one can associate a "dual" linear space which is its orthogonal complement.

**Definition 2.27.** Let $X \subseteq \mathbb{R}^d$ be a linear subspace. We define $X^\perp := \{\mathbf{y} \in \mathbb{R}^d : \langle \mathbf{y}, \mathbf{x} \rangle = 0 \ \forall \mathbf{x} \in X\}$ as the *orthogonal complement* of $X$.

The following is well-known from linear algebra.

**Proposition 2.28.** $X^\perp$ is a linear subspace. Moreover, $(X^\perp)^\perp = X$.

There is a way to generalize this idea of associating a dual object to convex sets.

**Definition 2.29.** Let $X \subseteq \mathbb{R}^d$ be any set. The set defined as

$$X^\circ := \{\mathbf{y} \in \mathbb{R}^d : \langle \mathbf{y}, \mathbf{x} \rangle \leq 1 \ \forall \mathbf{x} \in X\}$$

is called the *polar* of $X$.

**Proposition 2.30.** The following are all true.

1. $X^\circ$ is a closed, convex set for any $X \subseteq \mathbb{R}^d$ (not necessarily convex).

2. $(X^\circ)^\circ = \mathrm{cl}(\mathrm{conv}(X \cup \{\mathbf{0}\}))$. In particular, if $X$ is a closed convex set containing the origin, then $(X^\circ)^\circ = X$.

3. If $X$ is a convex cone, then $X^\circ = \{\mathbf{y} \in \mathbb{R}^d : \langle \mathbf{y}, \mathbf{x} \rangle \leq 0 \ \forall \mathbf{x} \in X\}$.

4. If $X$ is a linear subspace, then $X^\circ = X^\perp$.

*Proof.* 1. Follows from the fact that $X^\circ$ can be written as the intersection of closed halfspaces:

$$X^\circ = \bigcap_{\mathbf{x} \in X} \{\mathbf{y} \in \mathbb{R}^d : \langle \mathbf{y}, \mathbf{x} \rangle \leq 1\}.$$

NOTES: 15

2. Observe that $X \subseteq (X^\circ)^\circ$. Also, $\mathbf{0} \in (X^\circ)^\circ$, because $\mathbf{0}$ is always in the polar of any set. Since $(X^\circ)^\circ$ is a closed convex set by 1., we must have $\mathrm{cl}(\mathrm{conv}(X \cup \{\mathbf{0}\})) \subseteq (X^\circ)^\circ$.

To show the reverse inclusion, we show that if $\mathbf{y} \notin \mathrm{cl}(\mathrm{conv}(X \cup \{\mathbf{0}\}))$ then $\mathbf{y} \notin (X^\circ)^\circ$. Thus, we need to show that there exists $\mathbf{z} \in (X^\circ)$ such that $\langle \mathbf{y}, \mathbf{z} \rangle > 1$. Since $\mathbf{y} \notin \mathrm{cl}(\mathrm{conv}(X \cup \{\mathbf{0}\}))$, by Theorem 2.20, there exists $\mathbf{a} \in \mathbb{R}^d$, $\delta \in \mathbb{R}$ such that $\langle \mathbf{a}, \mathbf{y} \rangle > \delta$ and $\langle \mathbf{a}, \mathbf{x} \rangle \leq \delta$ for all $\mathbf{x} \in \mathrm{cl}(\mathrm{conv}(X \cup \{\mathbf{0}\}))$. Since $\mathbf{0} \in \mathrm{cl}(\mathrm{conv}(X \cup \{\mathbf{0}\}))$, we obtain that $0 \leq \delta$. We now consider two cases:

*Case 1: $\delta > 0$.* Set $\mathbf{z} = \frac{\mathbf{a}}{\delta}$. Now, $\langle \mathbf{z}, \mathbf{x} \rangle \leq 1$ for all $\mathbf{x} \in X$ because $\langle \mathbf{a}, \mathbf{x} \rangle \leq \delta$ for all $\mathbf{x} \in \mathrm{cl}(\mathrm{conv}(X \cup \{\mathbf{0}\})) \supseteq X$. Therefore, $\mathbf{z} \in X^\circ$. Moreover, $\langle \mathbf{z}, \mathbf{y} \rangle > 1$ because $\langle \mathbf{a}, \mathbf{y} \rangle > \delta$. So we are done.

*Case 2: $\delta = 0$.* Define $\epsilon := \langle \mathbf{a}, \mathbf{y} \rangle > \delta = 0$. Set $\mathbf{z} = \frac{2\mathbf{a}}{\epsilon}$. Then, $\langle \mathbf{z}, \mathbf{y} \rangle = 2 > 1$. Also, for every $\mathbf{x} \in X \subseteq \mathrm{cl}(\mathrm{conv}(X \cup \{\mathbf{0}\}))$, we obtain that $\langle \mathbf{z}, \mathbf{x} \rangle = \frac{2}{\epsilon}\langle \mathbf{a}, \mathbf{x} \rangle \leq \frac{2}{\epsilon}\delta = 0 \leq 1$. Thus, $\mathbf{z} \in X^\circ$. Thus, we are done.

3. and 4. are left to the reader. $\qquad\square$

---

**Example 2.31.** If $p, q \geq 1$ such that $\frac{1}{p} + \frac{1}{q} = 1$ (allowing for $p$ or $q$ to be $\infty$), then $B^\circ_{\ell^p}(\mathbf{0}, 1) = B_{\ell^q}(\mathbf{0}, 1)$. This example illustrates the use of the fundamental *Holder's inequality*.

**Proposition 2.32** (Holder's inequality)**.** If $p, q \geq 1$ such that $\frac{1}{p} + \frac{1}{q} = 1$ (allowing for $p$ or $q$ to be $\infty$), then
$$|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\|_p \|\mathbf{y}\|_q,$$
for every $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$. Moreover, if $p, q > 1$ then equality holds if and only if $|\mathbf{x}_i| = |\mathbf{y}_i|^{\frac{q}{p}}$.

The special case with $p = q = 2$ is known as the *Cauchy-Schwarz inequality*. We won't prove Holder's inequality here, but we will use it to derive the polarity relation between $\ell_p$ unit balls. We only show that $B_{\ell^q}(\mathbf{0}, 1) = B^\circ_{\ell^p}(\mathbf{0}, 1)$ for any $p, q > 1$ such that $\frac{1}{p} + \frac{1}{q} = 1$. The case $p = 1, q = \infty$ is considered in Problem 6 from "HW for Week III".

First, we show that $B_{\ell^q}(\mathbf{0}, 1) \subseteq B^\circ_{\ell^p}(\mathbf{0}, 1)$ Consider any $\mathbf{y} \in B_{\ell^q}(\mathbf{0}, 1)$ and consider any $\mathbf{x} \in B_{\ell^p}$. By Cauchy-Schwarz, we obtain that $\langle \mathbf{x}, \mathbf{y} \rangle \leq \|\mathbf{x}\|_p \|\mathbf{y}\|_q \leq 1$. Thus, $B_{\ell^q}(\mathbf{0}, 1) \subseteq B^\circ_{\ell^p}(\mathbf{0}, 1)$. To show the reverse inclusion $B^\circ_{\ell^p}(\mathbf{0}, 1) \subseteq B_{\ell^q}(\mathbf{0}, 1)$, consider any $\mathbf{y} \in B^\circ_{\ell^p}(\mathbf{0}, 1)$. We would like to show that $\mathbf{y} \in B_{\ell^q}(\mathbf{0}, 1)$, i.e., $\|\mathbf{y}\|_q \leq 1$. Suppose to the contrary that $\|\mathbf{y}\|_q > 1$. Consider $\mathbf{x}$ defined as follows: for each $i = 1, \ldots, d$, $\mathbf{x}_i$ has the same sign as $\mathbf{y}_i$, and $|\mathbf{x}_i| = |\mathbf{y}_i|^{\frac{q}{p}}$. Set $\tilde{\mathbf{x}} = \frac{\mathbf{x}}{\|\mathbf{x}\|_p}$. Now,

$$\langle \mathbf{y}, \tilde{\mathbf{x}} \rangle = \frac{1}{\|\mathbf{x}\|_p}\langle \mathbf{x}, \mathbf{y} \rangle = \frac{1}{\|\mathbf{x}\|_p}(\|\mathbf{x}\|_p \|\mathbf{y}\|_q) = \|\mathbf{y}\|_q > 1,$$

contradicting the fact that $\mathbf{y} \in B^\circ_{\ell^p}(\mathbf{0}, 1)$, because $\|\tilde{\mathbf{x}}\|_p = 1$. The second equality follows from Proposition 2.32 because of the special choice of $\mathbf{x}$.

NOTES: 16

### 2.3.2 Intrinsic description: faces, extreme points, recession cone, lineality space

We have seen that given any set $X$ of points in $\mathbb{R}^d$, the convex hull of $X$ – the smallest convex set containing $X$ – can be expressed as the set of all convex combinations of finite subsets of $X$ (Theorem 2.8). One possibility to represent a convex set $C$ *intrinsically* is to give a minimal subset $X \subseteq C$ such that all points in $C$ can be expressed as convex combinations of points in $X$, i.e., $C = \text{conv}(X)$. In particular, if $X$ is a finite set, then we can use $X$ to represent $C$ in a computer: implicitly, $C$ is the convex hull of the set $X$. We are going to get to such a "minimal" intrinsic description.

**Definition 2.33** (Faces and extreme points)**.** Let $C$ be a convex set. A convex subset $F \subseteq C$ is called an *extreme subset* or a *face* of $C$, if for any $\mathbf{x} \in F$ the following holds: $\mathbf{x}^1, \mathbf{x}^2 \in C, \frac{\mathbf{x}^1 + \mathbf{x}^2}{2} = \mathbf{x}$ implies that $\mathbf{x}^1, \mathbf{x}^2 \in F$. This is equivalent to saying that there is no point in $F$ that can be expressed as a convex combination of points in $C \setminus F$ – see Problem 10 from "HW for Week III".

A face of dimension 0 is called an *extreme point*. In other words, $\mathbf{x}$ is an extreme point of $C$ if the following holds: $\mathbf{x}^1, \mathbf{x}^2 \in C, \frac{\mathbf{x}^1 + \mathbf{x}^2}{2} = \mathbf{x}$ implies that $\mathbf{x}^1 = \mathbf{x}^2 = \mathbf{x}$. We denote the set of extreme points of $C$ by $\text{ext}(C)$.

The one-dimensional faces of a convex set are called its *edges*. If $k = \dim(C)$, then the $(k-1)$-dimensional faces are called *facets*. We will see below that the only $k$-dimensional face of $C$ is $C$ itself. Any face of $C$ that is not $C$ or $\emptyset$ is called a *proper* face of $C$.

**Definition 2.34.** Let $C$ be a convex set. We define the *relative interior* of $C$ as the set of all $\mathbf{x} \in C$ for which there exists $\epsilon > 0$ such that for all $\mathbf{y} \in \text{aff}(C)$, $\mathbf{x} + \epsilon\left(\frac{\mathbf{y} - \mathbf{x}}{\|\mathbf{y} - \mathbf{x}\|}\right) \in C$. We denote it by $\text{relint}(C)$.[1]

We define the *relative boundary* of $C$ to be $\text{relbd}(C) := \text{cl}(C) \setminus \text{relint}(C)$.

**Exercise 3.** Let $C$ be convex and $\mathbf{x} \in C$. Suppose that for all $\mathbf{y} \in \text{aff}(C)$, there exists $\epsilon_{\mathbf{y}}$ such that $\mathbf{x} + \epsilon_{\mathbf{y}}(\mathbf{y} - \mathbf{x}) \in C$. Show that $\mathbf{x} \in \text{relint}(C)$.

This exercise shows that it suffices to have a different $\epsilon$ for every direction; this implies a universal $\epsilon$ for every direction.

**Exercise 4.** Show that $\text{relint}(C)$ is nonempty for any nonempty convex set $C$.

**Lemma 2.35.** Let $C$ be a convex set of dimension $k$. The only $k$ dimensional face of $C$ is $C$ itself.

*Proof.* Let $F \subsetneq C$ be a proper face of $C$. Let $\mathbf{x} \in C \setminus F$. Let $X \subseteq F$ be a maximum set of affinely independent points in $F$. We claim that $X \cup \{\mathbf{x}\}$ is affinely independent. This immediately implies that $\dim(C) > \dim(F)$ and we will be done.

Suppose to the contrary that $\mathbf{x} \in \text{aff}(X)$. Then consider $\mathbf{x}^\star \in \text{relint}(F)$ (which is nonempty by Exercise 4). By definition, there exists $\epsilon > 0$ such that $\mathbf{y} = \mathbf{x}^\star + \epsilon(\mathbf{x} - \mathbf{x}^\star) \in F$. But this means that $\mathbf{y} = (1 - \epsilon)\mathbf{x}^\star + \epsilon\mathbf{x}$. Since $\mathbf{y} \in F$, and $\mathbf{x} \notin F$, this contradicts that $F$ is a face. $\qquad\square$

---

[1]For the reader familiar with the concept of a relative topology: the relative interior of $C$ is the interior of $C$ with respect to the relative topology of $\text{aff}(C)$.

NOTES:

**Lemma 2.36.** Let $C$ be a convex set and let $F \subseteq C$ be a face of $C$. If $\mathbf{x}$ is an extreme point of $F$, then $\mathbf{x}$ is an extreme point of $C$.

*Proof.* Left to the reader. $\qquad\qquad\square$

**Lemma 2.37.** Let $C \subseteq \mathbb{R}^d$ be convex. Let $\mathbf{a} \in \mathbb{R}^d$ and $\delta \in \mathbb{R}$ be such that $C \subseteq \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{a}, \mathbf{x} \rangle \leq \delta\}$. Then, the set $F = C \cap \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{a}, \mathbf{x} \rangle = \delta\}$ is a face of $C$.

*Proof.* Let $\bar{\mathbf{x}} \in F$ and $\mathbf{x}^1, \mathbf{x}^2 \in C$ such that $\frac{\mathbf{x}^1 + \mathbf{x}^2}{2} = \bar{\mathbf{x}}$. By the hypothesis, $\langle \mathbf{a}, \mathbf{x}^i \rangle \leq \delta$ for $i = 1, 2$. If for either $i = 1, 2$, $\langle \mathbf{a}, \mathbf{x}^i \rangle < \delta$, then

$$\langle \mathbf{a}, \bar{\mathbf{x}} \rangle = \left\langle \mathbf{a}, \frac{\mathbf{x}^1 + \mathbf{x}^2}{2} \right\rangle = \frac{\langle \mathbf{a}, \bar{\mathbf{x}}^1 \rangle + \langle \mathbf{a}, \bar{\mathbf{x}}^2 \rangle}{2} < \delta$$

contradicting that $\mathbf{x} \in F$. Therefore, we must have $\langle \mathbf{a}, \mathbf{x}^i \rangle = \delta$ for $i = 1, 2$ and thus, $\mathbf{x}^1, \mathbf{x}^2 \in F$. $\qquad\square$

**Definition 2.38.** A face $F$ of a convex set $C$ is called an *exposed face* if there exists $\mathbf{a} \in \mathbb{R}^d$ and $\delta \in \mathbb{R}$ be such that $C \subseteq \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{a}, \mathbf{x} \rangle \leq \delta\}$ and $F = C \cap \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{a}, \mathbf{x} \rangle = \delta\}$. We will sometimes make it explicit and say that $F$ is an *exposed face induced by* $(\mathbf{a}, \delta)$.

By working with the affine hull and the relative interior, and using Problem 3 from "HW for Week II", a stronger version of the supporting hyperplane theorem can be shown to be true.

**Theorem 2.39** (Supporting Hyperplane Theorem - II). Let $C \subseteq \mathbb{R}^d$ be convex and $\mathbf{x} \in \mathrm{relbd}(C)$. There exists $\mathbf{a} \in \mathbb{R}^d$ and $\delta \in \mathbb{R}$ such that all of the following hold:

(i) $\langle \mathbf{a}, \mathbf{y} \rangle \leq \delta$ for all $\mathbf{y} \in C$,

(ii) $\langle \mathbf{a}, \mathbf{x} \rangle = \delta$, and

(iii) there exists $\bar{\mathbf{y}} \in C$ such that $\langle \mathbf{a}, \bar{\mathbf{y}} \rangle < \delta$. This third condition says that $C$ is not completely contained in the hyperplane $\{\mathbf{y} \in \mathbb{R}^d : \langle \mathbf{a}, \mathbf{y} \rangle = \delta\}$.

An important consequence of the above discussion is the following theorem about the relative boundary of a closed, convex set $C$.

**Theorem 2.40.** Let $C \subseteq \mathbb{R}^d$ be a closed, convex set and $\mathbf{x} \in C$. $\mathbf{x}$ is contained in a proper face of $C$ if and only if $\mathbf{x} \in \mathrm{relbd}(C)$.

*Proof.* If $\mathbf{x} \in \mathrm{relbd}(C)$, then by Theorem 2.39 there exists $\mathbf{a} \in \mathbb{R}^d$ and $\delta \in \mathbb{R}$ such that the three conditions in Theorem 2.39 hold. By Lemma 2.37, $F = C \cap \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{a}, \mathbf{x} \rangle = \delta\}$ is a face of $C$, and it is proper face because of condition (iii) in Theorem 2.39.

Now let $\mathbf{x} \in F$ where $F$ is a proper face of $C$. Sicne $C$ is closed, it suffices to show that $\mathbf{x} \notin \mathrm{relint}(C)$. Suppose to the contrary that $\mathbf{x} \in \mathrm{relint}(C)$. Let $\bar{\mathbf{x}} \in C \setminus F$. Observe that $2\mathbf{x} - \bar{\mathbf{x}} \in \mathrm{aff}(C)$. Since $\mathbf{x}$ is assumed

NOTES: 18

to be in the relative interior of $C$, there exists $\epsilon > 0$ such that $\mathbf{y} = \epsilon((2\mathbf{x} - \bar{\mathbf{x}}) - \mathbf{x}) + \mathbf{x} \in C$. Rearranging terms, we obtain that

$$\mathbf{x} = \frac{\epsilon}{\epsilon + 1}\bar{\mathbf{x}} + \frac{1}{\epsilon + 1}\mathbf{y}.$$

Since $\mathbf{x} \in F$ and $\bar{\mathbf{x}} \notin F$, this contradicts the fact that $F$ is a face. Thus, $\mathbf{x} \notin \mathrm{relint}(C)$ and so $\mathbf{x} \in \mathrm{relbd}(C)$. $\qquad\square$

In our search for a subset $X \subseteq C$ such that $C = \mathrm{conv}(X)$, it is clear that $X$ must contain all extreme points. But is it sufficient to include all extreme points? In other words, is it true that $C = \mathrm{conv}(\mathrm{ext}(C))$? No! A simple counterexample is $\mathbb{R}^d_+$. Its only extreme point is 0. Another weird example is the set $\{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\| < 1\}$ – this set has NO extreme points! As you might suspect, the problem is that these sets are not compact, i.e., closed and bounded.

**Theorem 2.41** (Krein-Milman Theorem)**.** If $C$ a compact convex set, then $C = \mathrm{conv}(\mathrm{ext}(C))$.

*Proof.* The proof is going to use induction on the dimension of $C$. First, if $C$ is the empty set, then the statement is a triviality. So we assume $C$ is nonempty.

For the base case with $\dim(C) = 0$, i.e., $C = \{\mathbf{x}\}$ is a single point, the statement follows because $\{\mathbf{x}\}$ is an extreme point of $C$, and $C = \mathrm{conv}(\{\mathbf{x}\})$. For the induction step, consider any point $\mathbf{x} \in C$. We consider two cases:

*Case 1:* $\mathbf{x} \in \mathrm{relbd}(C)$. By Theorem 2.40, $\mathbf{x}$ is contained in a proper face $F$ of $C$. By Lemma 2.35, $\dim(F) < \dim(C)$. By the induction hypothesis applied to $F$ (note that $F$ is also compact using Problem 14 from "HW for Week III"), we can express $\mathbf{x}$ as a convex combination of extreme points of $F$, which by Lemma 2.36, shows that $\mathbf{x}$ is a convex combination of extreme points of $C$.

*Case 2:* $\mathbf{x} \in \mathrm{relint}(C)$. Let $\ell \subseteq \mathrm{aff}(C)$ be any affine set of dimension one (i.e., a line) going through $\mathbf{x}$. Since $C$ is compact, $\ell \cap C$ is a line segment. The end points $\mathbf{x}^1, \mathbf{x}^2$ of $\ell \cap C$ must be in the relative boundary of $C$. By the previous case, $\mathbf{x}^1, \mathbf{x}^2$ can be expressed as the convex combination of extreme points in $C$. Since $\mathbf{x}$ is a convex combination of $\mathbf{x}^1$ and $\mathbf{x}^2$, and a convex combination of convex combinations is a convex combination, we can express $\mathbf{x}$ as the convex combination of extreme points of $C$. $\qquad\square$

What about non-compact sets? Let us relax the condition of being bounded. So we want to describe closed, convex sets. It turns out that there is a nice way to deal with unboundedness. We introduce the necessary concepts next.

**Proposition 2.42.** Let $C$ be a closed, convex set, and $\mathbf{r} \in \mathbb{R}^d$. The following are equivalent:

1. There exists $\mathbf{x} \in C$ such that $\mathbf{x} + \lambda\mathbf{r} \in C$ for all $\lambda \geq 0$.

2. For every $\mathbf{x} \in C$, $\mathbf{x} + \lambda\mathbf{r} \in C$ for all $\lambda \geq 0$.

NOTES: 19

*Proof.* We only need to show 1. ⟹ 2.; the reverse implication is trivial. Let $\bar{\mathbf{x}}$ be such that $\bar{\mathbf{x}} + \lambda\mathbf{r} \in C$ for all $\lambda \geq 0$. Consider any arbitrary $\mathbf{x}^\star \in C$. Suppose to the contrary that there exists $\lambda' \geq 0$ such that $\mathbf{y} = \mathbf{x}^\star + \lambda'\mathbf{r} \notin C$. By Theorem 2.20, there exist $\mathbf{a} \in \mathbb{R}^d$, $\delta \in \mathbb{R}$ such that $\langle \mathbf{a}, \mathbf{y} \rangle > \delta$ and $\langle \mathbf{a}, \mathbf{x} \rangle \leq \delta$ for all $\mathbf{x} \in C$. This means that $\langle \mathbf{a}, \mathbf{r} \rangle > 0$ because otherwise, $\langle \mathbf{a}, \mathbf{y} \rangle = \langle \mathbf{a}, \mathbf{x}^\star \rangle + \lambda'\langle \mathbf{a}, \mathbf{r} \rangle = \delta + \lambda'\langle \mathbf{a}, \mathbf{r} \rangle \leq \delta$ causing a contradiction. But then, if we choose $\bar{\lambda} = \frac{|\delta - \langle \mathbf{a}, \bar{\mathbf{x}} \rangle| + 1}{\langle \mathbf{a}, \mathbf{r} \rangle}$, we would obtain that

$$\langle \mathbf{a}, \bar{\mathbf{x}} + \lambda\mathbf{r} \rangle = \langle \mathbf{a}, \bar{\mathbf{x}} \rangle + \bar{\lambda}\langle \mathbf{a}, \mathbf{r} \rangle = \langle \mathbf{a}, \bar{\mathbf{x}} \rangle + |\delta - \langle \mathbf{a}, \bar{\mathbf{x}} \rangle| + 1 \geq \delta + 1 > \delta,$$

contradicting the assumption that $\bar{\mathbf{x}} + \bar{\lambda}\mathbf{r} \in C$. □

**Definition 2.43.** Any $\mathbf{r} \in \mathbb{R}^d$ that satisfies the conditions in Proposition 2.42 is called a *recession direction* for $C$.

**Proposition 2.44.** The set of all recession directions of a closed, convex set is a closed, convex cone.

*Proof.* Fix any point $\mathbf{x}$ in the closed convex set $C$. Using condition 1. of Proposition 2.42, we see $\mathbf{r} \in \mathbb{R}^d$ is a recession direction if and only if for every $\lambda \geq 0$, $\mathbf{r} \in \frac{1}{\lambda}(C - \mathbf{x})$. Therefore,

$$\mathrm{rec}(C) = \bigcap_{\lambda \geq 0} \frac{1}{\lambda}(C - \mathbf{x}).$$

Each term in the intersection is a closed, convex set. Therefore, $\mathrm{rec}(C)$ is a closed, convex set. It is easy to see that for any $\mathbf{r} \in \mathrm{rec}(C)$, $\lambda\mathbf{r} \in \mathrm{rec}(C)$ also for every $\lambda \geq 0$. Thus, $\mathrm{rec}(C)$ is a closed, convex cone. □

**Definition 2.45.** We call the cone of recession directions the *recession cone* of $C$ and is denoted by $\mathrm{rec}(C)$. The set $\mathrm{rec}(C) \cap -\mathrm{rec}(C)$ is a linear subspace and is called the *lineality space* of $C$. It will be denoted by $\mathrm{lin}(C)$.

**Exercise 5.** Show that Proposition 2.42 remains true if $\lambda \geq 0$ is replaced by $\lambda \in \mathbb{R}$ in both conditions. Show that $\mathrm{lin}(C)$ is exactly the set of all $\mathbf{r} \in \mathbb{R}^d$ that satisfy these modified conditions.

Proposition 2.42 immediately gives the following corollary.

**Corollary 2.46.** Let $C$ be a closed convex set and let $F \subseteq C$ be a closed, convex subset. Then $\mathrm{rec}(F) \subseteq \mathrm{rec}(C)$.

*Proof.* Left as an exercise. □

Here is a characterization of compact convex sets.

**Theorem 2.47.** A closed convex set $C$ is compact if and only if $\mathrm{rec}(C) = \{\mathbf{0}\}$.

NOTES: 20

*Proof.* We leave it to the reader to check that if $C$ is compact, then $\mathrm{rec}(C) = \{\mathbf{0}\}$. For the other direction, assume that $\mathrm{rec}(C) = \{\mathbf{0}\}$. Suppose to the contrary that $C$ is not bounded, i.e., there exists a sequence of points $\mathbf{y}^i \in C$ such that $\|\mathbf{y}^i\| \to \infty$. Let $\mathbf{x} \in C$ be any point and consider the set of unit norm vectors $\mathbf{r}^i = \frac{\mathbf{y}^i - \mathbf{x}}{\|\mathbf{y}^i - \mathbf{x}\|}$. Since this is a sequence of unit norm vectors, by Theorem 1.10, there is a convergent subsequence $\{\mathbf{r}^{i_k}\}_{k=1}^\infty$ converging to $\mathbf{r}$ also with unit norm. We claim that $\mathbf{r}$ is a recession direction, giving a contradiction to $\mathrm{rec}(C) = \{\mathbf{0}\}$. To see this, for any $\lambda \geq 0$, let $N \in \mathbb{N}$ such that $\|\mathbf{y}^{i_k} - \mathbf{x}\| > \lambda$ for all $k \geq N$. We now observe that

$$\mathbf{x} + \lambda \mathbf{r}^{i_k} = \frac{(\|\mathbf{y}^{i_k} - \mathbf{x}\| - \lambda)}{\|\mathbf{y}^{i_k} - \mathbf{x}\|}\mathbf{x} + \frac{\lambda}{\|\mathbf{y}^{i_k} - \mathbf{x}\|}(\mathbf{x} + \mathbf{r}^i\|\mathbf{y}^{i_k} - \mathbf{x}\|) = \frac{(\|\mathbf{y}^{i_k} - \mathbf{x}\| - \lambda)}{\|\mathbf{y}^{i_k} - \mathbf{x}\|}\mathbf{x} + \frac{\lambda}{\|\mathbf{y}^{i_k} - \mathbf{x}\|}\mathbf{y}^{i_k} \in C$$

for all $k \geq N$. Letting $k \to \infty$, since $C$ is closed, we obtain that $\mathbf{x} + \lambda\mathbf{r} = \lim_{k\to\infty} \mathbf{x} + \lambda\mathbf{r}^{i_k} \in C$. $\qquad \square$

We next consider closed convex sets whose lineality space is $\{\mathbf{0}\}$.

**Definition 2.48.** If $\mathrm{lin}(C) = \{\mathbf{0}\}$ then $C$ is called *pointed*.

The main result about pointed closed convex sets says that you can decompose them into convex combinations of extreme points and recession directions.

**Theorem 2.49.** If $C$ is a closed, convex set that is pointed, then $C = \mathrm{conv}(\mathrm{ext}(C)) + \mathrm{rec}(C)$.

*Proof.* The proof follows the same lines as Theorem 2.41. We prove by induction on dimension of $C$. If $\dim(C) = 0$, then $C$ is a single point, and we are done.

We may assume $C$ is nonempty. Consider any $\mathbf{x} \in C$ and then two cases:

*Case 1:* $\mathbf{x} \in \mathrm{relbd}(C)$. By Theorem 2.40, $\mathbf{x}$ is contained in a proper face $F$ of $C$. By Lemma 2.35, $\dim(F) < \dim(C)$. By the induction hypothesis applied to $F$ (note that $F$ is also closed using Problem 14 from "HW for Week III"), we can express $\mathbf{x} = \mathbf{x}' + \mathbf{d}$, where $\mathbf{x}'$ is a convex combination of extreme points of $F$ and $\mathbf{d}$ is a recession direction for $F$. By Lemma 2.36, shows that $\mathbf{x}'$ is a convex combination of extreme points of $C$. By Corollary 2.46, $\mathbf{d} \in \mathrm{rec}(C)$.

*Case 2:* $\mathbf{x} \in \mathrm{relint}(C)$. Let $\ell$ be any affine set of dimension one (i.e., a line) going through $\mathbf{x}$. Since $C$ contains no lines ($C$ is pointed), $\ell \cap C$ is either a line segment, i.e., $\mathbf{x}$ is the convex combination of $\mathbf{x}^1, \mathbf{x}^2 \in \mathrm{relbd}(C)$, or $\ell \cap C$ is a half-line, i.e, $\mathbf{x} = \mathbf{x}' + \mathbf{d}$, where $\mathbf{x}' \in \mathrm{relbd}(C)$ and $\mathbf{d} \in \mathrm{rec}(C)$.

In the first case, using Case 1, for each $i = 1, 2$, $\mathbf{x}^i$ can be expressed as $\mathbf{x}^i = \mathbf{y}^i + \mathbf{d}^i$, where $\mathbf{y}^i$ is a convex combination of extreme points in $C$, and $\mathbf{d}^i \in \mathrm{rec}(C)$. Since $\mathbf{x}$ is a convex combination of $\mathbf{x}^1$ and $\mathbf{x}^2$, this shows that $\mathbf{x} \in \mathrm{conv}(\mathrm{ext}(C)) + \mathrm{rec}(C)$.

In the second case, applying Case 1 to $\mathbf{x}'$, we express $\mathbf{x}' = \mathbf{y}' + \mathbf{d}'$ where $\mathbf{y}'$ is a convex combination of extreme points in $C$, and $\mathbf{d}' \in \mathrm{rec}(C)$. Thus, $\mathbf{x} = \mathbf{y}' + \mathbf{d}' + \mathbf{d}$ and we have the desired representation. $\qquad \square$

Lets make this description even more "minimal". For this we will need to understand the structure of pointed cones.

NOTES:                                                      21

**Proposition 2.50.** Let $D$ be a closed, convex cone. The following are equivalent.

1. $D$ is pointed.

2. $D^\circ$ is full-dimensional, i.e., $\dim(D^\circ) = d$.

3. $\mathbf{0}$ is an exposed face of $D$.

4. There exists a compact, convex subset $B \subset D \setminus \{\mathbf{0}\}$ such that for every $\mathbf{d} \in D \setminus \{\mathbf{0}\}$, there exists a unique $\mathbf{b} \in B$ such that $\mathbf{d} = \lambda\mathbf{b}$ for some $\lambda > 0$. In particular, $D = \text{cone}(B)$.

*Proof.* 1. $\Rightarrow$ 2. If $D^\circ$ is not full-dimensional, then $\text{aff}(D^\circ)$ is a linear space of dimension strictly less than $d$, and so $\text{aff}(D^\circ)^\perp \neq \{\mathbf{0}\}$. Since $D^\circ \subseteq \text{aff}(D^\circ)$, using Problem 3 from "HW for Week III", and property 2. and 4. in Proposition 2.30, we obtain that $\text{aff}(D^\circ)^\perp = \text{aff}(D^\circ)^\circ \subseteq (D^\circ)^\circ = D$. Since $\text{aff}(D^\circ)^\perp$ is a linear space, this implies that $\text{aff}(D^\circ)^\perp \subseteq \text{lin}(D)$, contradicting the assumption that $D$ is pointed.

2. $\Rightarrow$ 3. By Problem 5 from "HW for Week II", $\text{int}(D^\circ) \neq \emptyset$. Choose any $\mathbf{y} \in \text{int}(D^\circ)$. Since $D^\circ = \{\mathbf{y} \in \mathbb{R}^d : \langle \mathbf{x}, \mathbf{y}\rangle \leq 0 \ \forall \mathbf{x} \in D\}$, using Problem 3 from "HW for Week II", we obtain that $\langle \mathbf{y}, \mathbf{x}\rangle < 0$ for every $\mathbf{x} \in D$. This shows that the exposed face induced by $(\mathbf{y}, 0)$ is exactly $\{\mathbf{0}\}$.

3. $\Rightarrow$ 4. Let $\mathbf{0}$ be an exposed face induced by $(\mathbf{y}, 0)$. Define $B := D \cap \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{y}, \mathbf{x}\rangle = -1\}$. It is clear from the definition that $\mathbf{0} \notin B$. We now show that $B$ is compact. It is the intersection of closed sets, so it is closed. By Theorem 2.47, it suffices to show that $\text{rec}(B) = \{\mathbf{0}\}$. Suppose to the contrary that there exists $\mathbf{r} \in \text{rec}(B) \setminus \{\mathbf{0}\}$. Consider any point $\bar{\mathbf{x}} \in B$. Since $\langle \mathbf{y}, \bar{\mathbf{x}}\rangle = -1$ and $\langle \mathbf{y}, \bar{\mathbf{x}} + \mathbf{r}\rangle = -1$, we obtain that $\langle \mathbf{y}, \mathbf{r}\rangle = 0$. Now, by Proposition 2.42, we obtain that $\mathbf{0} + \mathbf{r} \in D$, i.e., $\mathbf{r} \in D$. But then $\langle \mathbf{y}, \mathbf{r}\rangle = 0$ contradicting that $\mathbf{0}$ is an exposed face of $D$ induced by $(\mathbf{y}, 0)$.

We next consider any $\mathbf{d} \in D$. By our assumption, $\langle \mathbf{y}, \mathbf{d}\rangle < 0$. Thus, setting $\mathbf{b} = \frac{\mathbf{d}}{|\langle \mathbf{y}, \mathbf{d}\rangle|}$, we obtain that $\langle \mathbf{y}, \mathbf{b}\rangle = -1$ and thus, $\mathbf{b} \in B$. To show uniqueness, consider $\mathbf{b}^1, \mathbf{b}^2 \in B$ both satisfying the condition. This means, $\mathbf{b}^2 = \lambda\mathbf{b}^1$ for some $\lambda > 0$. Therefore,

$$\lambda\langle \mathbf{y}, \mathbf{b}^1\rangle = \langle \mathbf{y}, \mathbf{b}^2\rangle = -1 = \langle \mathbf{y}, \mathbf{b}^1\rangle$$

showing that $\lambda = 1$. This shows uniqueness of $\mathbf{b}$.

4. $\Rightarrow$ 1. If $D$ is not pointed, then there exists $\mathbf{x} \in D \setminus \{\mathbf{0}\}$ such that $-\mathbf{x} \in D$. Moreover, there exists $\lambda_1 > 0$ such that $\mathbf{x}^1 = \lambda_1\mathbf{x} \in B$ and $\lambda_2 > 0$ such $\mathbf{x}^2 = \lambda_2(-\mathbf{x}) \in B$. Since $B$ is convex, $\frac{\lambda_2}{\lambda_1 + \lambda_2}\mathbf{x}^1 + \frac{\lambda_1}{\lambda_1 + \lambda_2}\mathbf{x}^2 = \mathbf{0}$ is in $B$, contradicting the assumption. $\square$

**Definition 2.51.** For any closed convex cone $D$, any subset $B \subseteq D$ satisfying condition 4. of Proposition 2.50 is called a *base* of $D$.

The proof of Proposition 2.50 also shows the following.

**Corollary 2.52.** Let $D$ be a closed, convex cone. $D$ is pointed if and only if there exists a hyperplane $H$ such that $H \cap D$ is a base of $D$.

NOTES: 22

**Remark 2.53.** In fact, it can be shown that any base of a pointed cone $D$ must be of the form $H \cap D$ for some hyperplane $H$. We skip the proof of this fact from these notes.

**Definition 2.54.** Let $D$ be a closed, convex cone. An edge of $D$ is called an *extreme ray* of $D$. We say that $\mathbf{r} \in D$ *spans an extreme ray* if $\{\lambda \mathbf{r} : \lambda \geq 0\}$ is an extreme ray. The set of extreme rays of $D$ will be denoted by $\operatorname{extr}(D)$.

**Proposition 2.55.** Let $D$ be a closed, convex cone and $\mathbf{r} \in D \setminus \{\mathbf{0}\}$. $\mathbf{r}$ spans an extreme ray of $D$ if and only if for all $\mathbf{r}^1, \mathbf{r}^2 \in D$ such that $\mathbf{r} = \frac{\mathbf{r}^1 + \mathbf{r}^2}{2}$, there exist $\lambda_1, \lambda_2 \geq 0$ such that $\mathbf{r}^1 = \lambda_1 \mathbf{r}$ and $\mathbf{r}^2 = \lambda_2 \mathbf{r}$.

*Proof.* Left as an exercise. $\qquad\square$

Here is an analogue of the Krein-Milman Theorem (Theorem 2.41) for closed convex cones.

**Theorem 2.56.** If $D$ is a pointed, closed, convex cone, then $D = \operatorname{cone}(\operatorname{extr}(D))$.

*Proof.* By Proposition 2.50, there exists a base $B$ for $D$. Since $B$ is compact, $B = \operatorname{conv}(\operatorname{ext}(B))$ by Theorem 2.41. It is easy to verify that the ray spanned by each $\mathbf{r} \in \operatorname{ext}(B)$ is an extreme ray for $D$, and vice versa, any extreme ray of $D$ is spanned by some $\mathbf{r} \in \operatorname{ext}(B)$. Moreover, using the fact that $B = \operatorname{conv}(\operatorname{ext}(B))$, it immediately follows that $D = \operatorname{cone}(\operatorname{extr}(D))$. $\qquad\square$

**Slight abuse of notation.** For a closed convex set $C$, we will also use $\operatorname{extr}(C)$ to denote $\operatorname{extr}(\operatorname{rec}(C))$. We will also say these are the extreme rays of $C$.

Now we can write a sharper version of Theorem 2.49:

**Corollary 2.57.** If $C$ is a closed, convex set that is pointed, then $C = \operatorname{conv}(\operatorname{ext}(C)) + \operatorname{cone}(\operatorname{extr}(C))$.

Thus, to describe a pointed closed convex set, we just need to specify its extreme points and its extreme rays. We finally deal with general closed convex sets that are not necessarily pointed. The idea is that the lineality space can be "factored out".

**Lemma 2.58.** If $C$ is a closed convex set, then $C \cap \operatorname{lin}(C)^\perp$ is pointed.

*Proof.* Define $\hat{C} = C \cap \operatorname{lin}(C)^\perp$. $\hat{C}$ is closed because it is the intersection of two closed sets. By Corollary 2.46, $\operatorname{rec}(\hat{C}) \subseteq \operatorname{rec}(C)$. Therefore, $\operatorname{lin}(\hat{C}) = \operatorname{rec}(\hat{C}) \cap -\operatorname{rec}(\hat{C}) \subseteq \operatorname{rec}(C) \cap -\operatorname{rec}(C) = \operatorname{lin}(C)$. By the same reasoning, $\operatorname{lin}(\hat{C}) \subseteq \operatorname{lin}(\operatorname{lin}(C)^\perp) = \operatorname{lin}(C)^\perp$. Since $\operatorname{lin}(C) \cap \operatorname{lin}(C)^\perp = \{\mathbf{0}\}$, we obtain that $\operatorname{lin}(\hat{C}) = \{\mathbf{0}\}$. $\qquad\square$

**Theorem 2.59.** Let $C$ be a closed convex set and let $\hat{C} = C \cap \operatorname{lin}(C)^\perp$. Then

$$C = \operatorname{conv}(\operatorname{ext}(\hat{C})) + \operatorname{cone}(\operatorname{extr}(\hat{C})) + \operatorname{lin}(C).$$

*Proof.* We first observe that $C = \hat{C} + \operatorname{lin}(C)$. Indeed, for any $\mathbf{x} \in C$, we can express $\mathbf{x} = \mathbf{x}' + \mathbf{r}$ where $\mathbf{x}' \in \operatorname{lin}(C)^\perp$ and $\mathbf{r} \in \operatorname{lin}(C)$ (since $\operatorname{lin}(C) + \operatorname{lin}(C)^\perp = \mathbb{R}^n$). We also know that $\mathbf{x}' = \mathbf{x} - \mathbf{r} \in C$ because $\mathbf{r} \in \operatorname{lin}(C)$. Thus, $\mathbf{x}' \in \hat{C}$ and we are done. $\hat{C}$ is pointed by Lemma 2.58 and applying Corollary 2.57 gives the desired result. $\qquad\square$

NOTES: 23

Thus, a general closed convex set $C$ can be specified by giving a set of generators for its lineality space $\mathrm{lin}(C)$, and the extreme points and vectors spanning the extreme rays of the set $C \cap \mathrm{lin}(C)^\perp$. In Section 2.5, we will see that polyhedra are precisely those convex sets $C$ that have a finite number of extreme points and extreme rays for $C \cap \mathrm{lin}(C)^\perp$. So we see that polyhedra are especially easy to describe intrinsically: simply specify the finite list of extreme points, vectors spanning the extreme rays and a finite list of generators of $\mathrm{lin}(C)$.

### 2.3.3 A remark about extrinsic and intrinsic descriptions

You may have already observed that although a closed convex set can be represented as the intersection of halfspaces, such a representation is not unique. For example, consider the circle in $\mathbb{R}^2$. You can represent it by intersecting all its tangent halfspaces. On the other hand, if you throw away any finite subset of these halfspaces, you still get the same set. In fact, there is a representation which uses only countably many halfspaces. Thus, the same convex set can have many different representations as the intersection of halfspaces. Moreover, there is usually no way to choose a "canonical" representation, i.e., there is no set of representating halfspaces such that *any representation* will always include this "canonical" set of halfspaces (this situation will get a little better with polyhedra).

On the other hand, the intrinsic representation for a closed convex set is more "canonical". To begin with, consider the compact case. We express a compact $C$ as $\mathrm{conv}(\mathrm{ext}(C))$. We cannot remove any extreme point, because it cannot be represented as the convex combination of other points. Thus, this representation is unique/minimal/canonical in the sense that for any $X$ such that $C = \mathrm{conv}(X)$, we must have $\mathrm{ext}(C) \subseteq X$. With closed, convex sets that are pointed, we have a little more flexibility in choosing the representation because one can choose a different set of vectors to span the extreme rays. Even so, upto scaling, the representation is unique. More precisely, suppose $C$ is a closed, convex, pointed set that we express as

$$C = \mathrm{conv}(E) + \mathrm{cone}(R),$$

where $E = \mathrm{ext}(C)$ and $R$ is set of vectors each of which spans a different extreme ray of $\mathrm{rec}(C)$ and every extreme ray is spanned by some vector in $R$. Now, if we find another representation

$$C = \mathrm{conv}(E') + \mathrm{cone}(R'),$$

for some sets $E', R' \subseteq \mathbb{R}^d$, then we must have

(i) $E \subseteq E'$ and

(ii) for every $\mathbf{r} \in R$, there must be some nonnegative scaling of $\mathbf{r}$ present in $R'$.

Finally with closed, convex sets that are not pointed, we get an additional level of flexibility because of the non-trivial lineality space. Even so, there exists a canonical triple of sets $E, R, L \subseteq \mathbb{R}^d$ (see Theorem 2.59), such that

$$C = \mathrm{conv}(E) + \mathrm{cone}(R) + \mathrm{span}(L)$$

NOTES: 24

such that for any other triple, $E', R', L'$ satisfying

$$C = \text{conv}(E') + \text{cone}(R') + \text{span}(L'),$$

we must have

    (i) for every $\mathbf{v} \in E$, there exists $\mathbf{v}' \in E'$ such that $\mathbf{v} - \mathbf{v}' \in \text{lin}(C)$,

    (ii) for every $\mathbf{r} \in R$, there exists $\mathbf{r}' \in R'$ and $\lambda \geq 0$ such that $\mathbf{r} - \lambda \mathbf{r}' \in \text{lin}(C)$, and

    (iii) $\text{span}(L) = \text{span}(L') = \text{lin}(C)$.

    The same thing can be said about the extrinsic and intrinsic descriptions of an affine subspace: conditions 3. and 4. in Theorem 2.16, or a linear subspace: conditions 2. and 3. in Theorem 2.12.

## 2.4   Combinatorial theorems: Helly-Radon-Carathéodory

We will discuss three foundational results that expose combinatorial aspects of convexity. We begin with Radon's Theorem.

**Theorem 2.60** (Radon's Theorem)**.** Let $X \subseteq \mathbb{R}^d$ be a set of size at least $d + 2$. Then $X$ can be partitioned as $X = X_1 \uplus X_2$ into nonmpety sets $X_1, X_2$, such that $\text{conv}(X_1) \cap \text{conv}(X_2) \neq \emptyset$.

*Proof.* Since we can have at most $d+1$ affinely independent points in $\mathbb{R}^d$ (see condition 2. in Proposition 2.15), and $X$ has at least $d + 2$ points, there exists a subset $\{\mathbf{x}^1, \ldots, \mathbf{x}^k\} \subseteq X$ such that $\{\mathbf{x}^1, \ldots, \mathbf{x}^k\}$ is affinely dependent. By using characterization 5. in Proposition 2.15, there exist multipliers $\lambda_1, \ldots, \lambda_k \in \mathbb{R}$, not all zero, such that $\lambda_1 + \ldots + \lambda_k = 0$ and $\lambda_1 \mathbf{x}^1 + \ldots \lambda_k \mathbf{x}^k = 0$. Define $P := \{i : \lambda_i \geq 0\}$ and $N := \{i : \lambda_i < 0\}$. Since the $\lambda_i$'s are not all zero and $\lambda_1 + \ldots + \lambda_k = 0$, $P$ and $N$ both contain indices such that corresponding multiplier is non-zero. Moreover, $\sum_{i \in P} \lambda_i = \sum_{i \in N} (-\lambda_i)$ since $\lambda_1 + \ldots + \lambda_k = 0$, and $\sum_{j \in P} \lambda_j \mathbf{x}^j = \sum_{j \in N} (-\lambda_j) \mathbf{x}^j$ since $\lambda_1 \mathbf{x}^1 + \ldots + \lambda_k \mathbf{x}^k = 0$. Thus, we obtain that

$$\mathbf{y} = \sum_{j \in P} \frac{\lambda_j}{\sum_{i \in P} \lambda_i} \mathbf{x}^j = \sum_{j \in N} \frac{(-\lambda_j)}{\sum_{i \in N} (-\lambda_i)} \mathbf{x}^j,$$

showing that $\mathbf{y} \in \text{conv}(X_P) \cap \text{conv}(X_N)$ where $X_P = \{\mathbf{x}^i : i \in P\}$ and $X_N = \{\mathbf{x}^i : i \in N\}$. One can now simply define $X_1 = X_P$ and $X_2 = X \setminus X_P$. Note that $X_1, X_2$ are nonempty because $P$ and $N$ are nonempty sets. $\qquad\square$

**An application to learning theory: VC-dimension of halfspaces.** An important concept in learning theory is the *Vapnik-Çervonenkis (VC) dimension* of a family of subsets [5]. Let $\mathcal{F}$ be a family of subsets of $\mathbb{R}^d$ (possibly infinite).

**Definition 2.61.** A set $X \subseteq \mathbb{R}^q$ is said to be *shattered* by $\mathcal{F}$, if for every subset $X' \subseteq X$, there exists a set $F \in \mathcal{F}$ such that $X' = F \cap X$. The VC-dimension of $\mathcal{F}$ is defined as

$$\sup\{m \in \mathbb{N} : \text{ there exists a set } X \subseteq \mathbb{R}^d \text{ of size } m \text{ that can be shattered by } \mathcal{F}.\}$$

**Proposition 2.62.** Let $\mathcal{F}$ be the family of halfspaces in $\mathbb{R}^d$. The VC-dimension of $\mathcal{F}$ is $d + 1$.

*Proof.* For any $m \leq d + 1$, let $X$ be a set of $m$ affinely independent points. Now, for any subset $X' \subseteq X$, we claim that $\text{conv}(X') \cap \text{conv}(X \setminus X') = \emptyset$ (Verify!!). When we study polyhedra in Section 2.5, we will see that $\text{conv}(X')$ and $\text{conv}(X \setminus X')$ are compact convex sets. By Problem 7, there exists a separating hyperplane for these two sets, giving a halfspace $H$ such that $X' = H \cap X$.

Let $m \geq d + 2$. Consider any set $X$ with $m$ points. By Theorem 2.60, one can partition $X = X_1 \uplus X_2$ with $X_1, X_2$ nonempty such that there exists $\mathbf{y} \in \text{conv}(X_1) \cap \text{conv}(X_2)$. Let $X' = X_1$. Consider any halfspace $H$ such that $X' \subseteq H$. Since $H$ is convex, $\mathbf{y} \in H$. By Problem 11 in "HW for Week IV", we obtain that $H \cap X_2 \neq \emptyset$. Thus, $X$ cannot be shattered by the family of halfspaces in $\mathbb{R}^d$. $\qquad\square$

See Chapters 12 and 13 of [2] for more on VC dimension.

An extremely important corollary of Radon's Theorem is known as Helly's theorem concerning the intersection of a family of convex sets.

**Theorem 2.63** (Helly's Theorem)**.** Let $X_1, \ldots, X_k \subseteq \mathbb{R}^d$ be a family of convex sets. If $X_1 \cap \ldots, \cap X_k = \emptyset$, then there is a subfamily $X_{i_1}, \ldots, X_{i_m}$ for some $m \leq d + 1$, with $i_h \in \{1, \ldots, k\}$ for each $h = 1, \ldots, m$ such that $X_{i_1} \cap \ldots, \cap X_{i_m} = \emptyset$. Thus, there is a subfamily of size at most $d + 1$ that already certifies the empty intersection.

618
619
620
621
622
623
624

*Proof.* We prove by induction on $k$. The base case is if $k \leq d + 1$, then we are done. Assume we know the statement to be true for all families of convex sets with $\bar{k}$ elements for some $\bar{k} \geq d + 1$. Consider a family of $\bar{k} + 1$ convex sets $X_1, X_2, \ldots, X_{\bar{k}+1}$. Define a new family $C_1, \ldots, C_{\bar{k}}$, where $C_i = X_i$ if $i \leq \bar{k} - 1$ and $C_{\bar{k}} = X_{\bar{k}} \cap X_{\bar{k}+1}$. Since $\emptyset = X_1 \cap \ldots \cap X_{\bar{k}+1} = C_1 \cap \ldots \cap C_{\bar{k}}$, we can use the induction hypothesis on this new family and obtain a subfamily $C_{i_1}, \ldots, C_{i_m}$ such that $C_{i_1} \cap \ldots \cap C_{i_m} = \emptyset$ and $m \leq d + 1$. If $m \leq d$ or none of the $C_{i_h}$, $h = 1, \ldots, m$ equals $C_{\bar{k}}$, then we are done. So we assume that $m = d + 1$ and $C_{i_m} = C_{\bar{k}} = X_{\bar{k}} \cap X_{\bar{k}+1}$.

To simplify notation, let us relabel everything and define $D_h := C_{i_h} = X_{i_h}$, $h = 1, \ldots, d$ and $D_{d+1} = X_{\bar{k}}$ and $D_{d+2} = X_{\bar{k}+1}$. We thus know that $D_1 \cap \ldots \cap D_{d+2} = \emptyset$. We may assume that each subfamily of $d + 1$

sets from $D_1, \ldots, D_{d+2}$ has a nonempty intersection, because otherwise we will be done. Let these common intersection points be

$$\mathbf{x}^i \in \cap_{h \neq i} D_h, \ \ i = 1, \ldots, d+2.$$

By Theorem 2.60, there exists a partition $\{1, \ldots, d+2\} = L \uplus R$ where $L, R$ are nonempty sets, such that there exists $\mathbf{y} \in \operatorname{conv}(\{\mathbf{x}^i\}_{i \in L}) \cap \operatorname{conv}(\{\mathbf{x}^i\}_{i \in R})$. Now, we claim that $\mathbf{y} \in D_h$ for each $h \in \{1, \ldots, d+2\}$ arriving at a contradiction to $D_1 \cap \ldots \cap D_{d+2} = \emptyset$. Indeed, Consider any $h^\star \in \{1, \ldots, d+2\}$. Either $L$ or $R$ does not contain it. Suppose $L$ does not contain it. Then for each $i \in L$, $\mathbf{x}^i \in \cap_{h \neq i} D_h \subseteq D_{h^\star}$ because $i \neq h^\star$. Since $D_{h^\star}$ is convex, this shows that $\mathbf{y} \in \operatorname{conv}(\{\mathbf{x}^i\}_{i \in L}) \subseteq D_{h^\star}$. $\qquad \square$

A corollary for infinite families is often useful, as long as we assume compactness for the elements in the family.

**Corollary 2.64.** Let $\mathcal{X}$ be a (possibly infinite) family of compact, convex sets. If $\cap_{X \in \mathcal{X}} X = \emptyset$, then there is a subfamily $X_{i_1}, \ldots, X_{i_m}$ for some $m \leq d+1$, with $i_h \in \{1, \ldots, k\}$ for each $h = 1, \ldots, m$ such that $X_{i_1} \cap \ldots, \cap X_{i_m} = \emptyset$. Thus, there is a subfamily of size at most $d+1$ that already certifies the empty intersection.

*Proof.* By a standard result in topology, if the intersection of an infinite family of compact sets is empty, then there is a finite subfamily whose intersection is also empty. One can now apply Theorem 2.63 to this finite subfamily and obtain a subfamily of is at most $d+1$. $\qquad \square$

**Application to centerpoints.** Helly's theorem can be used to extend the notion of median to distributions on $\mathbb{R}^d$ with $d \geq 2$. Let $\mu$ be any probability distribution on $\mathbb{R}^d$. For any point $\mathbf{x} \in \mathbb{R}^d$, define

$$f_\mu(\mathbf{x}) := \inf\{\mu(H) : H \text{ halfspace such that } \mathbf{x} \in H\}.$$

Define the *centerpoint or median* with respect $\mu$ as any $\mathbf{x}$ in the set $\mathrm{C}_\mu := \arg\max_{\mathbf{x} \in \mathbb{R}^d} f_\mu(\mathbf{x})$. It can be shown that this set is nonempty for all probability distributions $\mu$. For $d = 1$, this gives the standard notion of a median, and one can show that for any probability distribution $\mu$ on $\mathbb{R}$, $f_\mu(\mathbf{x}) = \frac{1}{2}$ for any centerpoint/median $\mathbf{x}$. In higher dimensions, unfortunately, one cannot guarantee a value of $\frac{1}{2}$. In fact, given the uniform distribution on a triangle in $\mathbb{R}^2$, one can show that the centroid $\mathbf{x}$ of the triangle is the unique centerpoint, and has value $f_\mu(\mathbf{x}) = \frac{4}{9} < \frac{1}{2}$. So can one guarantee any lower bound? Or can we find distributions whose centerpoint values are arbitrarily low? Grünbraum [4] proved a lower bound for the value of a centerpoint, irrespective of the distribution. The only assumption is a mild regularity condition on the distribution: for any halfspace $H$ and any $\delta > 0$, there exists a closed halfspace $H' \subseteq \mathbb{R}^d \setminus H$ such that $\mu(H') \geq \mu(\mathbb{R}^d \setminus H) - \delta$.

**Theorem 2.65.** Let $\mu$ be any probability distribution on $\mathbb{R}^d$ satisfying the above assumption. There exists a point $\mathbf{x} \in \mathbb{R}^d$ such that $f_\mu(\mathbf{x}) \geq \frac{1}{d+1}$.

*Proof.* Given any $\alpha \in \mathbb{R}$, let $\mathcal{H}_\alpha$ be the set of all halfspaces $H$ such that $\mu(H) \geq \alpha$. It is not hard to check that if $\alpha > 0$, then $D_\alpha := \cap_{H \in \mathcal{H}_\alpha} H$ is a compact, convex set. Indeed, for any coordinate indexed by $i = 1, \ldots, d$, there must exist some $\delta_1^i, \delta_2^i$ such that the halfspaces $H_1^i := \{\mathbf{x} \in \mathbb{R}^d : \mathbf{x}_i \leq \delta_1^i\}$ and $H_2^i := \{\mathbf{x} \in \mathbb{R}^d : \mathbf{x}_i \geq \delta_2^i\}$ satisfy $\mu(H_1^i) \geq \alpha$ and $\mu(H_2^i) \geq \alpha$. Thus, $D_\alpha$ is contained in the box $\{\mathbf{x} \in \mathbb{R}^d : \delta_2^i \leq \mathbf{x}_i \leq \delta_1^i, \ i = 1, \ldots, d\}$.

We now claim that for any $\mathbf{x} \in D_\alpha$, we have $f_\mu(\mathbf{x}) \geq 1 - \alpha$. To see this, consider any halfspace $H = \{\mathbf{y} \in \mathbb{R}^d : \langle \mathbf{a}, \mathbf{y} \rangle \leq \delta\}$ that contains $\mathbf{x} \in D_\alpha$. We will show that $\mu(\mathbb{R}^d \setminus H) \leq \alpha$. Indeed, if $\mu(\mathbb{R}^d \setminus H) > \alpha$, then some halfspace $H'$ contained in $\mathbb{R}^d \setminus H$ also has mass at least $\alpha$. This would imply that $H'$ contains all of $D_\alpha$ and, therefore, $\mathbf{x} \in H'$. But since $H' \subseteq \mathbb{R}^d \setminus H$, this contradicts the fact that $\mathbf{x} \in H$.

Therefore, it suffices to show that $D_{\frac{d}{d+1}+\epsilon}$ is nonempty for every $\epsilon > 0$, because using compactness $\cap_{\epsilon > 0} D_{\frac{d}{d+1}+\epsilon}$ is nonempty, and any point $\mathbf{x}$ in this set will satisfy $f_\mu(\mathbf{x}) \geq \frac{1}{d+1}$.

Now let's fix an $\epsilon > 0$. We want to show that $D_{\frac{d}{d+1}+\epsilon}$ is nonempty. By standard measure-theoretic arguments, there exists a ball $B$ centered at the origin such that $\mu(B) \geq 1 - \frac{\epsilon}{2}$ and $D_{\frac{d}{d+1}+\epsilon} \subseteq B$, because $D_\alpha := \cap_{H \in \mathcal{H}_\alpha} H$ is a compact.

Define $\mathcal{C} = \{B \cap H : H \text{ is a closed half space with } \mu(H) \geq \frac{d}{d+1} + \epsilon\}$. Thus, $\mathcal{C}$ is a family of compact sets such that $D_{\frac{d}{d+1}+\epsilon} = \bigcap\{C : C \in \mathcal{C}\}$. For any subset $\{C_1, \ldots, C_{h(S)}\} \subseteq \mathcal{C}$ of size $d+1$, we claim

$$\mu(C_1^c \cup \ldots \cup C_{d+1}^c) \leq 1 - (d+1)\frac{\epsilon}{2}.$$

This is because each $C_i^c = B^c \cup H_i^c$ for some half space $H_i$ satisfying $\mu(H_i^c) \leq \frac{1}{d+1} - \epsilon$. Since $\mu(B^c) \leq \frac{\epsilon}{2}$,
we obtain that $\mu(C_i^c) \leq \frac{1}{d+1} - \frac{\epsilon}{2}$. Therefore,

$$\mu(C_1 \cap \ldots \cap C_{h(S)}) = 1 - (\mu(C_1^c \cup \ldots C_{h(S)}^c)) \geq 1 - (1 - (d+1)\frac{\epsilon}{2}) = (d+1)\frac{\epsilon}{2} > 0.$$

This implies that $C_1 \cap \ldots \cap C_{h(S)} \neq \emptyset$. By Corollary 2.64, $\bigcap\{C : C \in \mathcal{C}\}$ is nonempty and so $D_{\frac{d}{d+1}+\epsilon}$ is nonempty. $\square$

Another useful theorem is Carathéodory's theorem which says that if a point $\mathbf{x}$ can be expressed as the convex combination of some other set $X \subseteq \mathbb{R}^d$ of points, then there is a subset of $X' \subseteq X$ of size at most $d + 1$ such that $\mathbf{x} \in \text{conv}(X')$. We state the conical version first, and then the convex version.

**Theorem 2.66** (Carathéodory's Theorem – cone version)**.** Let $X \subseteq \mathbb{R}^d$ (not necessarily convex) and let $\mathbf{x} \in \text{cone}(X)$. There exists a subset $X' \subseteq X$ such that $X'$ is linearly independent (and thus, $|X'| \leq d$), and $\mathbf{x} \in \text{cone}(X')$.

*Proof.* Since $\mathbf{x} \in \text{cone}(X)$, by Theorem 2.11, we can find a finite set $\{\mathbf{x}^1, \ldots, \mathbf{x}^k\} \subseteq X$ such that $\mathbf{x} \in \text{cone}(\{\mathbf{x}^1, \ldots, \mathbf{x}^k\})$. Choose a minimal such set, i.e., there is not strict subset of $\{\mathbf{x}^1, \ldots, \mathbf{x}^k\}$ whose conical hull contains $\mathbf{x}$. This implies that $\mathbf{x} = \lambda_1 \mathbf{x}^1 + \ldots + \lambda_k \mathbf{x}^k$ for some $\lambda_i > 0$ for each $i = 1, \ldots k$. We claim that $\mathbf{x}^1, \ldots, \mathbf{x}^k$ are linearly independent. Suppose to the contrary that there exist multipliers $\gamma_1, \ldots, \gamma_k \in \mathbb{R}$, not all zero, such that $\gamma_1 \mathbf{x}^1 + \ldots + \gamma_k \mathbf{x}^k = \mathbf{0}$. By changing the signs of the $\gamma_i$'s if necessary, we may assume that there exists $j \in \{1, \ldots, k\}$ such that $\gamma_j > 0$. Define

$$\theta = \min_{j: \gamma_j > 0} \frac{\lambda_j}{\gamma_j}, \qquad \lambda_i' = \lambda_i - \theta \gamma_i \quad \forall i = 1, \ldots, k.$$

Observe that $\lambda_i' \geq 0$ for all $i = 1, \ldots, k$ and

$$\lambda_1' \mathbf{x}^1 + \ldots, \lambda_k' \mathbf{x}^k = \lambda_1 \mathbf{x}^1 + \ldots \lambda_k \mathbf{x}^k - \theta(\gamma_1 \mathbf{x}^1 + \ldots + \gamma_k \mathbf{x}^k) = \lambda_1 \mathbf{x}^1 + \ldots \lambda_k \mathbf{x}^k = \mathbf{x}.$$

However, at least one of the $\lambda_i'$'s is zero (corresponding to an index in $\arg\min_{j: \gamma_j > 0} \frac{\lambda_j}{\gamma_j}$), contradicting the minimal choice of $\{\mathbf{x}^1, \ldots, \mathbf{x}^k\}$. $\qquad\square$

**Theorem 2.67** (Carathéodory's Theorem – convex version)**.** Let $X \subseteq \mathbb{R}^d$ (not necessarily convex) and let $\mathbf{x} \in \text{conv}(X)$. There exists a subset $X' \subseteq X$ such that $X'$ is affinely independent (and thus, $|X'| \leq d + 1$), and $\mathbf{x} \in \text{conv}(X')$.

*Proof.* Consider the set $Y \subseteq \mathbb{R}^{d+1}$ defined by $Y := \{(\mathbf{y}, 1) : \mathbf{y} \in X\}$. Now, $\mathbf{x} \in \text{conv}(X)$ is equivalent to saying that $(\mathbf{x}, 1) \in \text{cone}(Y)$. We get the desired result by applying Theorem 2.66 and condition 4. of Proposition 2.15. $\qquad\square$

We can finally furnish the proof of Lemma 2.26.

*Proof of Lemma 2.26.* Consider a convergent sequence $\{\mathbf{x}^i\}_{i \in \mathbb{N}} \subseteq \text{cone}(\{\mathbf{a}^1, \ldots, \mathbf{a}^n\})$ converging to $\mathbf{x} \in \mathbb{R}^d$. By Theorem 2.66, every $\mathbf{x}^i$ is in the conical hull of some linearly independent subset of $\{\mathbf{a}^1, \ldots, \mathbf{a}^n\}$. Since there are only finitely many linearly independent subsets of $\{\mathbf{a}^1, \ldots, \mathbf{a}^n\}$, one of these subsets contains infinitely many elements of the sequence $\{x^i\}_{i \in \mathbb{N}}$. Thus, after passing to that subsequence, we may assume that $\{\mathbf{x}^i\}_{i \in \mathbb{N}} \subseteq \text{cone}(\{\bar{\mathbf{a}}^1, \ldots, \bar{\mathbf{a}}^k\})$ where $\{\bar{\mathbf{a}}^1, \ldots, \bar{\mathbf{a}}^k\}$ are linearly independent. For each $\mathbf{x}^i$, there exists $\boldsymbol{\lambda}^i \in \mathbb{R}_+^k$ such that $\mathbf{x}^i = \boldsymbol{\lambda}_1^i \bar{\mathbf{a}}^1 + \ldots + \boldsymbol{\lambda}_k^i \bar{\mathbf{a}}^k$. Since $\{\mathbf{x}^i\}_{i \in \mathbb{N}}$ is a convergent sequence, it is also a bounded set. This implies that $\{\boldsymbol{\lambda}^i\}_{i \in \mathbb{N}}$ is a bounded set in $\mathbb{R}_+^k$ because $\bar{\mathbf{a}}^1, \ldots, \bar{\mathbf{a}}^k$ are all linearly independent. Thus, by

NOTES: 29

Theorem 1.10 there is a convergent subsequence $\boldsymbol{\lambda}^{i_k} \to \boldsymbol{\lambda} \in \mathbb{R}_+^k$. Note that $\mathbf{x}^{i_k} = A\boldsymbol{\lambda}^{i_k}$, where $A \in \mathbb{R}^{d \times k}$ is the matrix with $\bar{\mathbf{a}}^1, \ldots, \bar{\mathbf{a}}^k$ as columns. Taking limits,

$$\mathbf{x} = \lim_{k \to \infty} \mathbf{x}^{i_k} = \lim_{k \to \infty} A\boldsymbol{\lambda}^{i_k} = A\boldsymbol{\lambda}.$$

Since $\boldsymbol{\lambda} \in \mathbb{R}_+^k$, we find that $\mathbf{x} \in \mathrm{cone}(\{\bar{\mathbf{a}}^1, \ldots, \bar{\mathbf{a}}^k\}) \subseteq \mathrm{cone}(\{\mathbf{a}^1, \ldots, \mathbf{a}^n\})$. $\qquad\square$

Here is another result that proves handy in many situations.

**Theorem 2.68.** Let $X \subseteq \mathbb{R}^d$ be a compact set (not necessarily convex). Then $\mathrm{conv}(X)$ is compact.

*Proof.* By Theorem 2.67, every $\mathbf{x} \in \mathrm{conv}(X)$ is the convex combination of some $d+1$ points in $X$. Define the following function $f : \underbrace{\mathbb{R}^d \times \ldots \times \mathbb{R}^d}_{d+1 \text{ times}} \times \mathbb{R}^{d+1} \to \mathbb{R}^d$ as follows:

$$f(\mathbf{y}^1, \ldots, \mathbf{y}^{d+1}, \boldsymbol{\lambda}) = \lambda_1 \mathbf{y}^1 + \ldots + \lambda_{d+1} \mathbf{y}^{d+1}.$$

It is easily verified that $f$ is a continuous function (each coordinate of $f(\cdot)$ is a bilinear quadratic function of the input). We now observe that $\mathrm{conv}(X)$ is the image of $\underbrace{X \times \ldots \times X}_{d+1 \text{ times}} \times \Delta^{d+1}$ under $f$, where

$$\Delta^{d+1} := \{\boldsymbol{\lambda} \in \mathbb{R}_+^{d+1} : \lambda_1 + \ldots + \lambda_{d+1} = 1\}.$$

Since $X$ and $\Delta^{d+1}$ are compact sets, we obtain the result by applying Theorem 1.12. $\qquad\square$

## 2.5 Polyhedra

Recall that a polyhedron is any convex set that can be obtained by intersecting a finite number of halfspaces (Definition 2.22). Polyhedra, in a sense, are the nicest convex sets to work with because of this finiteness property. For example, our first result will be that a polyhedron can have only finitely many extreme points.

Even so, one thing to keep in mind is that the same polyhedron can be described as the intersection of two completely different finite families of halfspaces. This brings into sharp focus the non-uniqueness of extrinsic descriptions discussed in Section 2.3.3. Consider the following systems of halfspace/inequalities.

$$
\begin{array}{rcl}
-x_1 & \leq & 0 \\
x_1 + x_2 & \leq & 0 \\
x_1 - x_2 & \leq & 0 \\
-x_1 - x_2 - x_3 & \leq & 0 \\
x_2 + x_3 & \leq & 5
\end{array}
\qquad\qquad
\begin{array}{rcl}
2x_1 + x_2 & \leq & 0 \\
-x_1 + x_2 & \leq & 0 \\
x_1 - 2x_2 & \leq & 0 \\
x_1 - 2x_3 & \leq & 0 \\
2x_1 + x_2 + 2x_3 & \leq & 10
\end{array}
$$

Both these systems describe the same polyhedron $P = \mathrm{conv}\{(0,0,0), (0,0,5)\}$ in $\mathbb{R}^3$. However, if a polyhedron is given by its list of extreme points and extreme rays, this ambiguity disappears. Moreover, having

NOTES:                                         30

⁶⁶⁹ these two alternate extrinsic/intrinsic descriptions is very useful as many properties become easier to see
⁶⁷⁰ in one description, compared to the other description. Let us, therefore, start by making some important
⁶⁷¹ observations about extreme points and extreme rays of a polyhedron.

⁶⁷² **Definition 2.69.** Let $P$ be a polyhedron. Let $A \in \mathbb{R}^{m \times d}$ with rows $\mathbf{a}^1, \ldots, \mathbf{a}^m$ and $\mathbf{b} \in \mathbb{R}^m$ such that
⁶⁷³ $P = \{\mathbf{x} \in \mathbb{R}^d : A\mathbf{x} \leq \mathbf{b}\}$. Given any $\mathbf{x} \in P$, define $\mathrm{tight}(\mathbf{x}, A, \mathbf{b}) := \{i : \langle \mathbf{a}^i, \mathbf{x} \rangle = \mathbf{b}_i\}$. For brevity, when
⁶⁷⁴ $A$ and $\mathbf{b}$ are clear from the context, we will shorten this to $\mathrm{tight}(\mathbf{x})$. We also use the notation $A_{\mathrm{tight}(\mathbf{x})}$ to
⁶⁷⁵ denote the submatrix formed by taking the rows of $A$ indexed by $\mathrm{tight}(\mathbf{x})$. Similarly, $\mathbf{b}_{\mathrm{tight}(\mathbf{x})}$ will denote
⁶⁷⁶ the subvector of $\mathbf{b}$ indexed by $\mathrm{tight}(\mathbf{x})$.

⁶⁷⁷ **Theorem 2.70.** Let $P = \{\mathbf{x} \in \mathbb{R}^d : A\mathbf{x} \leq \mathbf{b}\}$ be a polyhedron given by $A \in \mathbb{R}^{m \times d}$ and $\mathbf{b} \in \mathbb{R}^m$. Let $\mathbf{x} \in P$.
⁶⁷⁸ Then, $\mathbf{x}$ is an extreme point of $P$ if and only if $A_{\mathrm{tight}(\mathbf{x})}$ has rank equal to $d$, i.e., the rows of $A$ indexed by
⁶⁷⁹ $\mathrm{tight}(\mathbf{x})$ span $\mathbb{R}^d$.

*Proof.* ($\Leftarrow$) Suppose $A_{\mathrm{tight}(\mathbf{x})}$ has rank equal to $d$; we want to establish that $\mathbf{x}$ is an extreme point. Consider
any $\mathbf{x}^1, \mathbf{x}^2 \in P$ such that $\mathbf{x} = \frac{\mathbf{x}^1 + \mathbf{x}^2}{2}$. For each $i \in \mathrm{tight}(\mathbf{x})$, $\langle \mathbf{a}^i, \mathbf{x}^1 \rangle \leq \mathbf{b}_i$ and similarly, $\langle \mathbf{a}^i, \mathbf{x}^2 \rangle \leq \mathbf{b}_i$. Now,
we observe that

$$\mathbf{b}_i = \langle \mathbf{a}^1, \mathbf{x} \rangle = \frac{\langle \mathbf{a}^i, \mathbf{x}^1 \rangle}{2} + \frac{\langle \mathbf{a}^i, \mathbf{x}^2 \rangle}{2} \leq \mathbf{b}_i.$$

⁶⁸⁰ Thus, the inequality must be an equality. Therefore, for each $i \in \mathrm{tight}(\mathbf{x})$, $\langle \mathbf{a}^i, \mathbf{x}^1 \rangle = \mathbf{b}_i$ and similarly,
⁶⁸¹ $\langle \mathbf{a}^i, \mathbf{x}^2 \rangle = \mathbf{b}_i$. In other words, we have that $A_{\mathrm{tight}(\mathbf{x})}\mathbf{x} = \mathbf{b}_{\mathrm{tight}(\mathbf{x})}$, and $A_{\mathrm{tight}(\mathbf{x})}\mathbf{x}^j = \mathbf{b}_{\mathrm{tight}(\mathbf{x})}$ for $j = 1, 2$.
⁶⁸² Since rank of $A_{\mathrm{tight}(\mathbf{x})} = d$, the system of equations must have a unique solution. This means $\mathbf{x} = \mathbf{x}^1 = \mathbf{x}^2$.
⁶⁸³ This shows that $\mathbf{x}$ is extreme.

($\Rightarrow$) Suppose to the contrary that $\mathbf{x}$ is extreme and $A_{\mathrm{tight}(\mathbf{x})}$ has rank strictly less than $d$ (note that its
rank is less than or equal to $d$ because it has $d$ columns). Thus, there exists a non-zero $\mathbf{r} \in \mathbb{R}^d$ such that
$A_{\mathrm{tight}(\mathbf{x})}\mathbf{r} = 0$. Define

$$\epsilon := \min\{\min\{\frac{\mathbf{b}_j - \langle \mathbf{a}^j, \mathbf{x} \rangle}{\langle \mathbf{a}^j, \mathbf{r} \rangle} : \langle \mathbf{a}^j, \mathbf{r} \rangle > 0\}, \min\{\frac{\mathbf{b}_j - \langle \mathbf{a}^j, \mathbf{x} \rangle}{-\langle \mathbf{a}^j, \mathbf{r} \rangle} : \langle \mathbf{a}^j, \mathbf{r} \rangle < 0\}\}$$

⁶⁸⁴ Note that $\epsilon > 0$. We now claim that $\mathbf{x}^1 := \mathbf{x} + \epsilon\mathbf{r} \in P$ and $\mathbf{x}^2 := \mathbf{x} - \epsilon\mathbf{r} \in P$. This would show that
⁶⁸⁵ $\mathbf{x} = \frac{\mathbf{x}^1 + \mathbf{x}^2}{2}$ with $\mathbf{x}^1 \neq \mathbf{x}^2$ (because $\mathbf{r} \neq 0$ and $\epsilon > 0$), contradicting extremality.
⁶⁸⁶ To finish the proof, we need to check that $A\mathbf{x}^1 \leq \mathbf{b}$ and $A\mathbf{x}^2 \leq \mathbf{b}$. We will do the calculations for $\mathbf{x}^1$ –
⁶⁸⁷ the calculations for $\mathbf{x}^2$ are similar. Consider any $j \in \{1, \ldots, m\}$. If $j \in \mathrm{tight}(\mathbf{x})$, the since $A_{\mathrm{tight}}\mathbf{r} = 0$, we
⁶⁸⁸ obtain that $\langle \mathbf{a}^j, \mathbf{x}^1 \rangle = \langle \mathbf{a}^j, \mathbf{x} \rangle + \epsilon\langle \mathbf{a}^j, \mathbf{r} \rangle = \langle \mathbf{a}^j, \mathbf{x} \rangle = \mathbf{b}_j$. If $j \notin \mathrm{tight}(\mathbf{x})$, then we consider two cases:

⁶⁸⁹ *Case 1:* $\langle \mathbf{a}^j, \mathbf{r} \rangle > 0$. Since $\epsilon \leq \frac{\mathbf{b}_j - \langle \mathbf{a}^j, \mathbf{x} \rangle}{\langle \mathbf{a}^j, \mathbf{r} \rangle}$, we obtain that $\langle \mathbf{a}^j, \mathbf{x}^1 \rangle = \langle \mathbf{a}^j, \mathbf{x} \rangle + \epsilon\langle \mathbf{a}^j, \mathbf{r} \rangle \leq \mathbf{b}_j$.

⁶⁹⁰ *Case 2:* $\langle \mathbf{a}^j, \mathbf{r} \rangle < 0$. In this case, $\langle \mathbf{a}^j, \mathbf{x}^1 \rangle = \langle \mathbf{a}^j, \mathbf{x} \rangle + \epsilon\langle \mathbf{a}^j, \mathbf{r} \rangle < \mathbf{b}_j$, simply because $\epsilon > 0$ and $\langle \mathbf{a}^j, \mathbf{r} \rangle < 0$. □

⁶⁹¹ This immediately gives the following.

NOTES: 31

**Corollary 2.71.** Any polyhedron $P \subseteq \mathbb{R}^d$ has a finite number of extreme points.

*Proof.* Let $A \in \mathbb{R}^{m \times d}$ and $\mathbf{b} \in \mathbb{R}^m$ be such that $P = \{\mathbf{x} \in \mathbb{R}^d : A\mathbf{x} \le \mathbf{b}\}$. From Theorem 2.70, for any extreme point, $A_{\text{tight}(\mathbf{x})}$ has rank $d$. There are only finitely many subsets $I \subseteq \{1, \ldots, m\}$ such that the submatrix $A_I$ is of rank $d$. Moreover, for any $I \subseteq \{1, \ldots, m\}$ such that $A_I$ is rank $d$ such that $A_I \mathbf{x} = \mathbf{b}_I$ has a solution, the set of solutions to $A_I \mathbf{x} = \mathbf{b}_I$ is unique. This shows that there are only finitely many extreme points. □

What about the extreme rays? First we define *polyhedral cones*.

**Definition 2.72.** A convex cone that is also a polyhedron is called a polyhedral cone.

**Proposition 2.73.** Let $D$ be a convex cone. $D \subseteq \mathbb{R}^d$ is a polyhedral cone if and only if there exists a matrix $A \in \mathbb{R}^{m \times d}$ for some $m \in N$ such that $D = \{\mathbf{x} \in \mathbb{R}^d : A\mathbf{x} \le \mathbf{0}\}$.

*Proof.* We simply have to show the forward direction, the reverse is easy. Assume $D$ is a polyhedral cone. Thus, it is polyhedron and so there exists a matrix $A \in \mathbb{R}^{m \times d}$ and $\mathbf{b} \in \mathbb{R}^m$ for some $m \in \mathbb{N}$ such that $D = \{\mathbf{x} : A\mathbf{x} \le \mathbf{b}\}$. Since $D$ is a closed, convex cone (closed because all polyhedra are closed), $\text{rec}(D) = D$. By Problem 1 in "HW for Week IV", we obtain that $D = \text{rec}(D) = \{\mathbf{x} : A\mathbf{x} \le \mathbf{0}\}$. □

Problem 1 in "HW for Week IV" also immediately implies the following.

**Proposition 2.74.** If $P$ is a polyhedron, then $\text{rec}(P)$ is a polyhedral cone.

**Theorem 2.75.** Let $D = \{\mathbf{x} : A\mathbf{x} \le \mathbf{0}\}$ be a polyhedral cone and let $\mathbf{r} \in D$. $\mathbf{r}$ spans an extreme ray if and only if $A_{\text{tight}(\mathbf{r})}$ has rank $d - 1$.

*Proof.* ($\Leftarrow$) Let $A_{\text{tight}(\mathbf{r})}$ have rows $\bar{\mathbf{a}}^1, \ldots, \bar{\mathbf{a}}^k$. Each $F_i := D \cap \{\mathbf{y} : \langle \bar{\mathbf{a}}^i, \mathbf{y} \rangle = 0\}$ for each $i = 1, \ldots, k$ is an exposed face of $D$. By Problem 13 in "HW for Week III", $F := \cap_{i=1}^k F_i$ is a face of $D$. Since $A_{\text{tight}(\mathbf{r})}$ has rank $d - 1$, the set $\{\mathbf{x} : A_{\text{tight}(\mathbf{x})}\mathbf{x} = \mathbf{0}\}$ is a 1-dimensional linear subspace. Since $F \subseteq \{\mathbf{x} : A_{\text{tight}(\mathbf{x})}\mathbf{x} = \mathbf{0}\}$, $F$ is a 1-dimensional face of $D$ and hence an extreme ray. Note that $\mathbf{r} \in F$ and thus $\mathbf{r}$ spans $F$.

($\Rightarrow$) Suppose $\mathbf{r}$ spans the 1-dimensional face $F$. Recall that this means that any $\mathbf{x} \in F$ is a scaling of $\mathbf{r}$. Rank of $A_{\text{tight}(\mathbf{r})}$ cannot be $d$ since then $\mathbf{r}$ is an extreme point of $D$ and $\mathbf{r} = \mathbf{0}$ by Problem 3 in "HW for Week IV". This would contradict that $\mathbf{r}$ spans an extreme ray of $D$. Thus, rank of $A_{\text{tight}(\mathbf{r})} \le d - 1$. If it is strictly less, then consider any $\mathbf{r}' \in \{\mathbf{x} : A_{\text{tight}(\mathbf{x})}\mathbf{x} = \mathbf{0}\}$ that is linearly independent to $\mathbf{r}$ – such an $\mathbf{r}'$ exists if rank of $A_{\text{tight}(\mathbf{r})} \le d - 2$. Define

$$\epsilon := \min\{\min\{\frac{-\langle \mathbf{a}^j, \mathbf{r} \rangle}{\langle \mathbf{a}^j, \mathbf{r}' \rangle} : \langle \mathbf{a}^j, \mathbf{r}' \rangle > 0\}, \min\{\frac{-\langle \mathbf{a}^j, \mathbf{r} \rangle}{-\langle \mathbf{a}^j, \mathbf{r}' \rangle} : \langle \mathbf{a}^j, \mathbf{r}' \rangle < 0\}\}$$

Note that $\epsilon > 0$. We now claim that $\mathbf{r}^1 := \mathbf{r} + \epsilon \mathbf{r}' \in D$ and $\mathbf{r}^2 := \mathbf{r} - \epsilon \mathbf{r}' \in D$. This would show that $\mathbf{r} = \frac{\mathbf{r}^1 + \mathbf{r}^2}{2}$. Moreover, since $\mathbf{r}'$ and $\mathbf{r}$ are linearly independent, $\mathbf{r}^1, \mathbf{r}^2$ are not scalings of $\mathbf{r}$. This contradicts Proposition 2.55.

NOTES:                                                          32

717      To finish the proof, we need to check that $A\mathbf{r}^1 \leq \mathbf{0}$ and $A\mathbf{r}^2 \leq \mathbf{0}$. This is the same set of calculations as
718 in the proof of Theorem 2.70.           $\square$

719      Analogous to Corollary 2.71, we have:

720 **Corollary 2.76.** Any polyhedral cone $D$ has finitely many extreme rays.

721 ### 2.5.1    The Minkowski-Weyl Theorem

722 We can now state the first part of the famous Minkowski-Weyl theorem.

723 **Theorem 2.77** (Minkowski-Weyl Theorem – Part I)**.** Let $P \subseteq \mathbb{R}^d$ be a polyhedron. Then there exist finite
724 sets $V, R \subseteq \mathbb{R}^d$ such that $P = \text{conv}(V) + \text{cone}(R)$.

725 *Proof.* Let $L$ be a finite set of vectors spanning $\text{lin}(P)$ ($L$ is taken as the empty set if $\text{lin}(P) = \{\mathbf{0}\}$). Note
726 that $\text{lin}(P) = \text{cone}(L \cup -L)$. Define $\hat{P} = P \cap \text{lin}(P)^{\perp}$. By Problem 1 (iii) in "HW for Week V", $\hat{P}$ is also a
727 polyhedron. By Corollary 2.71, we obtain that $V := \text{ext}(\hat{P})$ is a finite set. Moreover, by Proposition 2.74,
728 $\text{rec}(\hat{P})$ is a polyhedral cone. By Corollary 2.76, $\text{extr}(\text{rec}(\hat{P}))$ is a finite set. Define $R = \text{extr}(\text{rec}(\hat{P})) \cup L \cup -L$.
729 By Theorem 2.59, $P = \text{conv}(\text{ext}(\hat{P})) + \text{cone}(\text{rec}(\hat{P})) + \text{lin}(P) = \text{conv}(V) + \text{cone}(R)$.     $\square$

730      We now make an observation about polars.

731 **Lemma 2.78.** Let $V, R \subseteq \mathbb{R}^d$ be finite sets and let $X = \text{conv}(V) + \text{cone}(R)$. Then $X$ is a closed, convex set.

732 *Proof.* $\text{conv}(V)$ is compact, by Theorem 2.68, and $\text{cone}(R)$ is closed by Lemma 2.26. By Problem 6 in "HW
733 for Week I/II" we obtain that $X = \text{conv}(V) + \text{conv}(R)$ is closed. Since the Minkowski sum of convex sets is
734 convex (property 3. in Theorem 2.3), $X$ is also convex.     $\square$

**Theorem 2.79.** Let $V = \{\mathbf{v}^1, \ldots, \mathbf{v}^k\} \subseteq \mathbb{R}^d$, and $R = \{\mathbf{r}^1, \ldots, \mathbf{r}^n\} \subseteq \mathbb{R}^d$ with $k \geq 1$ and $n \geq 0$. Let
$X = \text{conv}(V) + \text{cone}(R)$. Then

$$X^{\circ} = \left\{ \mathbf{y} \in \mathbb{R}^d : \begin{array}{ll} \langle \mathbf{v}^i, \mathbf{y} \rangle \leq 1 & i = 1, \ldots, k \\ \langle \mathbf{r}^j, \mathbf{y} \rangle \leq 0 & j = 1, \ldots, n \end{array} \right\}.$$

*Proof.* Define $\tilde{X} := \left\{ \mathbf{y} \in \mathbb{R}^d : \begin{array}{ll} \langle \mathbf{v}^i, \mathbf{y} \rangle \leq 1 & i = 1, \ldots, k \\ \langle \mathbf{r}^i, \mathbf{y} \rangle \leq 0 & i = 1, \ldots, n \end{array} \right\}$. We first verify that $\tilde{X} \subseteq X^{\circ}$, i.e., $\langle \mathbf{y}, \mathbf{x} \rangle \leq 1$ for

all $\mathbf{y} \in \tilde{X}$ and $\mathbf{x} \in X$. By definition of $X$, we can write $\mathbf{x} = \sum_{i=1}^{k} \lambda_i \mathbf{v}^i + \sum_{j=1}^{n} \mu_j \mathbf{r}^j$ for some $\lambda_i, \mu_j \geq 0$ such
that $\sum_{i=1}^{k} \lambda_i = 1$. Thus,

$$\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{i=1}^{k} \lambda_i \langle \mathbf{v}^i, \mathbf{y} \rangle + \sum_{j=1}^{n} \mu_j \langle \mathbf{r}^j, \mathbf{y} \rangle \leq 1,$$

735 since $\langle \mathbf{v}^i, \mathbf{y} \rangle \leq 1$ for $i = 1, \ldots, k$, and $\langle \mathbf{r}^j, \mathbf{y} \rangle \leq 0$ for $j = 1, \ldots, n$.

NOTES:                  33

To see that $X^\circ \subseteq \tilde{X}$, consider any $\mathbf{y} \in X^\circ$. Since $\langle \mathbf{x}, \mathbf{y} \rangle \leq 1$ for all $\mathbf{x} \in X$, we must have $\langle \mathbf{v}^i, \mathbf{y} \rangle \leq 1$ for $i = 1, \ldots, k$ since $\mathbf{v}^i \in X$. Suppose to the contrary that $\langle \mathbf{r}^j, \mathbf{y} \rangle > 0$ for some $j \in \{1, \ldots, n\}$. Then there exists $\lambda > 0$ such that $\langle \mathbf{v}^1 + \lambda \mathbf{r}^j, \mathbf{y} \rangle > 1$. But this contradicts the fact that $\langle \mathbf{x}, \mathbf{y} \rangle \leq 1$ for all $\mathbf{x} \in X$ because $\mathbf{v}^1 + \lambda \mathbf{r}^j \in X$, by definition of $X$. Therefore, $\langle \mathbf{r}^j, \mathbf{y} \rangle \leq 0$ for $j = 1, \ldots, n$ and thus, $\mathbf{y} \in \tilde{X}$. $\qquad\square$

This has the following corollary.

**Corollary 2.80.** Let $P$ be a polyhedron. Then $P^\circ$ is a polyhedron.

*Proof.* If $P = \emptyset$, then $P^\circ = \mathbb{R}^d$, which is a polyhedron. Else, by Theorem 2.77, there exist finite sets $V, R \subseteq \mathbb{R}^d$ such that $P = \text{conv}(V) + \text{cone}(R)$, with $V \neq \emptyset$. By Theorem 2.79, $P^\circ$ is the intersection of finitely many halfspaces, and is thus a polyhedron. $\qquad\square$

We now prove the converse of Theorem 2.77.

**Theorem 2.81** (Minkowski-Weyl Theorem – Part II)**.** Let $V, R \subseteq \mathbb{R}^d$ be finite sets and let $X = \text{conv}(V) + \text{cone}(R)$. Then $X \subseteq \mathbb{R}^d$ is a polyhedron.

*Proof.* The case when $X$ is empty is trivial. So we consider $X$ is nonempty. Take any $\mathbf{t} \in X$ and define $X' = X - \mathbf{t}$. Now, it is easy to see $X$ is polyhedron if and only if $X'$ is a polyhedron (Verify!!). So it suffices to show that $X'$ is a polyhedron. Note that $X' = \text{conv}(V') + \text{cone}(R)$ where $V' = V - \mathbf{t}$, which is a nonempty set because $V$ is nonempty (since $X$ is assumed to be nonempty). By Theorem 2.79, $(X')^\circ$ is a polyhedron. By Lemma 2.78, $X'$ is a closed, convex set, and also $\mathbf{0} \in X'$. Therefore, $X' = ((X')^\circ)^\circ$ by condition 2. in Theorem 2.30. Applying Corollary 2.80 with $P = (X')^\circ$, we obtain that $((X')^\circ)^\circ = X'$ is a polyhedron. $\qquad\square$

Collecting Theorems 2.77 and 2.81 together, we have the full-blown Minkowski-Weyl Theorem.

**Theorem 2.82** (Minkowski-Weyl Theorem – full version)**.** Let $X \subseteq \mathbb{R}^d$. Then the following are equivalent.

(i) ($\mathcal{H}$-description) There exists $m \in \mathbb{N}$, a matrix $A \in \mathbb{R}^{m \times d}$ and a vector $\mathbf{b} \in \mathbb{R}^m$ such that $X = \{\mathbf{x} \in \mathbb{R}^d : A\mathbf{x} \leq \mathbf{b}\}$.

(ii) ($\mathcal{V}$-description) There exist finite sets $V, R \subseteq \mathbb{R}^d$ such that $X = \text{conv}(V) + \text{cone}(R)$.

A compact version is often useful.

**Theorem 2.83** (Minkowski-Weyl Theorem – compact version)**.** Let $X \subseteq \mathbb{R}^d$. Then $X$ is a bounded polyhedron if and only if $X$ is the convex hull of a finite set of points.

*Proof.* Left as an exercise. $\qquad\square$

### 2.5.2   Valid inequalities and feasibility

**Definition 2.84.** Let $X \subseteq \mathbb{R}^d$ (not necessarily convex) and let $\mathbf{a} \in \mathbb{R}^d, \delta \in \mathbb{R}$. We say that $\langle \mathbf{a}, \mathbf{x} \rangle \leq \delta$ is a *valid inequality/halfspace* for $X$ if $X \subseteq H^-(\mathbf{a}, \delta)$.

Consider a polyhedron $P = \{\mathbf{x} \in \mathbb{R}^d : A\mathbf{x} \leq \mathbf{b}\}$ with $A \in \mathbb{R}^{m \times d}, \mathbf{b} \in \mathbb{R}^m$. For any vector $\mathbf{y} \in \mathbb{R}_+^m$, the inequality $\langle \mathbf{y}^T A, \mathbf{x} \leq \mathbf{y}^T \mathbf{b}$ is clearly a valid inequality for $P$. The next theorem says that all valid inequalities are of this form, upto a translation.

**Theorem 2.85.** Let $P = \{\mathbf{x} \in \mathbb{R}^d : A\mathbf{x} \leq \mathbf{b}\}$ with $A \in \mathbb{R}^{m \times d}, \mathbf{b} \in \mathbb{R}^m$ be a nonempty polyhedron. Let $\mathbf{c} \in \mathbb{R}^d, \delta \in \mathbb{R}$. Then $\langle \mathbf{c}, \mathbf{x} \rangle \leq \delta$ is a valid inequality for $P$ if and only if there exists $\mathbf{y} \in \mathbb{R}_+^m$ such that $\mathbf{c}^T = \mathbf{y}^T A$ and $\mathbf{y}^T \mathbf{b} \leq \delta$.

*Proof.* ($\Leftarrow$) Suppose there exists $\mathbf{y} \in \mathbb{R}_+^m$ such that $\mathbf{c}^T = \mathbf{y}^T A$ and $\mathbf{y}^T \mathbf{b} \leq \delta$. The validity of $\langle \mathbf{c}, \mathbf{x} \rangle \leq \delta$ is clear from the following relations for any $\mathbf{x} \in P$:

$$\langle \mathbf{c}, \mathbf{x} \rangle = \langle \mathbf{y}^T A, \mathbf{x} \rangle = \mathbf{y}^T (A\mathbf{x}) \leq \mathbf{y}^T \mathbf{b} \leq \delta,$$

where the first inequality follows from the fact that $\mathbf{x} \in P$ implies $A\mathbf{x} \leq \mathbf{b}$ and $\mathbf{y}$ is nonnegative.

($\Rightarrow$) Let $\langle \mathbf{c}, \mathbf{x} \rangle \leq \delta$ be a valid inequality for $P$. Suppose to the contrary that there is no nonnegative solution to $\mathbf{c}^T = \mathbf{y}^T A$ and $\mathbf{y}^T \mathbf{b} \leq \delta$. This is equivalent to saying that the following system has no solution in $\mathbf{y}, \lambda$:

$$A^T \mathbf{y} = \mathbf{c}, \quad \mathbf{b}^T \mathbf{y} + \lambda = \delta, \quad \mathbf{y} \geq 0, \lambda \geq 0.$$

Setting this up in matrix notation, we have no nonnegative solutions to

$$\left[ \begin{array}{cc} A^T & 0 \\ \mathbf{b}^T & 1 \end{array} \right] \left[ \begin{array}{c} \mathbf{y} \\ \lambda \end{array} \right] = \left[ \begin{array}{c} \mathbf{c} \\ \delta \end{array} \right].$$

By Farkas' Lemma (Theorem 2.25), there exists $\mathbf{u} = (\bar{\mathbf{u}}, \mathbf{u}_{d+1}) \in \mathbb{R}^{d+1}$ such that

$$\bar{\mathbf{u}}^T A^T + \mathbf{u}_{d+1} \mathbf{b}^T \leq \mathbf{0}, \quad \mathbf{u}_{d+1} \leq 0, \quad \text{and} \quad \bar{\mathbf{u}}^T \mathbf{c} + \mathbf{u}_{d+1} \delta > 0. \tag{2.1}$$

We now consider two cases:

*Case 1:* $\mathbf{u}_{d+1} = 0$. Plugging into (2.1), we obtain $\bar{\mathbf{u}}^T A^T \leq \mathbf{0}$, i.e. $A\bar{\mathbf{u}} \leq \mathbf{0}$, and $\langle \mathbf{c}, \bar{\mathbf{u}} \rangle > 0$. By Problem 1 in "HW for Week IV", $\bar{\mathbf{u}} \in \mathrm{rec}(P)$. Consider any $\mathbf{x} \in P$ (we assume $P$ is nonempty). Let $\mu = \frac{1 + (\delta - \langle \mathbf{c}, \mathbf{x} \rangle)}{\langle \mathbf{c}, \bar{\mathbf{u}} \rangle} > 0$. Now $\mathbf{x} + \mu \bar{\mathbf{u}} \in P$ since $\bar{\mathbf{u}} \in \mathrm{rec}(P)$. However, $\langle \mathbf{c}, \mathbf{x} + \mu \bar{\mathbf{u}} \rangle = \delta + 1 > \delta$, contradicting that $\langle \mathbf{c}, \mathbf{x} \rangle \leq \delta$ is a valid inequality for $P$.

*Case 2:* $\mathbf{u}_{d+1} < 0$. By rearranging (2.1), we have $A\bar{\mathbf{u}} \leq (-\mathbf{u}_{d+1})\mathbf{b}$ and $\langle \mathbf{c}, \bar{\mathbf{u}} \rangle > (-\mathbf{u}_{d+1})\delta$. By setting $\mathbf{x} = \frac{\bar{\mathbf{u}}}{-\mathbf{u}_{d+1}}$, obtain that $\mathbf{x} \in P$ and $\langle \mathbf{c}, \mathbf{x} \rangle > \delta$, contradicting that $\langle \mathbf{c}, \mathbf{x} \rangle \leq \delta$ is a valid inequality for $P$. $\square$

NOTES:                                                    35

**Definition 2.86.** Let $\mathbf{c} \in \mathbb{R}^d$ and $\delta_1, \delta_2$. If $\delta_1 \leq \delta_2$, then the inequality/halfspace $\langle \mathbf{c}, \mathbf{x} \rangle \leq \delta_1$ is said to *dominate* the inequality/halfspace $\langle \mathbf{c}, \mathbf{x} \rangle \leq \delta_2$.

**Remark 2.87.** Let $P = \{\mathbf{x} \in \mathbb{R}^d : A\mathbf{x} \leq \mathbf{b}\}$ with $A \in \mathbb{R}^{m \times d}, \mathbf{b} \in \mathbb{R}^m$ be a polyhedron. Then $\langle \mathbf{c}, \mathbf{x} \rangle \leq \delta$ is called a *consequence of* $A\mathbf{x} \leq \mathbf{b}$ if there exists $\mathbf{y} \in \mathbb{R}^m_+$ such that $\mathbf{c}^T = \mathbf{y}^T A$ and $\delta = \mathbf{y}^T \mathbf{b}$. Another way to think of Theorem 2.85 is that it says the geometric property of being a valid inequality is the same as the algebraic property of being a consequence:

> [Alternate version of Theorem 2.85] Let $P = \{\mathbf{x} \in \mathbb{R}^d : A\mathbf{x} \leq \mathbf{b}\}$ be a nonempty polyhedron. Then $\langle \mathbf{c}, \mathbf{x} \rangle \leq \delta$ is a valid inequality for $P$ if and only if $\langle \mathbf{c}, \mathbf{x} \rangle \leq \delta$ is dominated by a consequence of $A\mathbf{x} \leq \mathbf{b}$.

A version of Theorem 2.85 for empty polyhedra is also useful. It can be interpreted as the existence of a short certificate of infeasibility of polyhedra.

**Theorem 2.88.** Let $P = \{\mathbf{x} \in \mathbb{R}^d : A\mathbf{x} \leq \mathbf{b}\}$ with $A \in \mathbb{R}^{m \times d}, \mathbf{b} \in \mathbb{R}^m$ be a polyhedron. Then $P = \emptyset$ if and only if $\langle \mathbf{0}, \mathbf{x} \rangle \leq -1$ is a consequence of $A\mathbf{x} \leq \mathbf{b}$.

*Proof.* It is easy to see that if $\langle \mathbf{0}, \mathbf{x} \rangle \leq -1$ is a consequence of $A\mathbf{x} \leq \mathbf{b}$ then $P = \emptyset$, because any point that satisfies $A\mathbf{x} \leq \mathbf{b}$ must satisfy every consequence of it, and no point satisfies $\langle \mathbf{0}, \mathbf{x} \rangle \leq -1$.

So now assume $P = \emptyset$. This means that there is no solution to $A\mathbf{x} \leq \mathbf{b}$. This is equivalent to saying that there is no solution to $A\mathbf{x}^1 - A\mathbf{x}^2 + \mathbf{s} = \mathbf{b}$ with $\mathbf{x}^1, \mathbf{x}^2, \mathbf{s} \geq 0$.[2] In matrix notation, this means there are no nonnegative solutions to

$$\begin{bmatrix} A & -A & I \end{bmatrix} \begin{bmatrix} \mathbf{x}^1 \\ \mathbf{x}^2 \\ \mathbf{s} \end{bmatrix} = \mathbf{b}.$$

By Farkas' Lemma (Theorem 2.25), there exists $\mathbf{u} \in \mathbb{R}^m$ such that

$$\mathbf{u}^T A \leq \mathbf{0}, \ \ \mathbf{u}^T(-A) \leq \mathbf{0}, \ \ \mathbf{u} \leq \mathbf{0}, \quad \text{and} \quad \mathbf{u}^T \mathbf{b} > 0.$$

Define $\mathbf{y} = \frac{-\mathbf{u}}{\mathbf{u}^T \mathbf{b}} \geq \mathbf{0}$. Then $\mathbf{y}^T A = \mathbf{0}$ and $\mathbf{y}^T \mathbf{b} = -1$, showing that $\langle \mathbf{0}, \mathbf{x} \rangle \leq -1$ is a consequence of $A\mathbf{x} \leq \mathbf{b}$. $\qquad \square$

### 2.5.3 Faces of polyhedra

Faces for polyhedra are very structured. Firstly, every face is an exposed face – something that is not true for general closed, convex sets. Secondly, there is an algebraic characterization of faces in terms of the describing inequalities of a polyhedron.

**Theorem 2.89.** Let $P = \{\mathbf{x} \in \mathbb{R}^d : A\mathbf{x} \leq \mathbf{b}\}$ with $A \in \mathbb{R}^{m \times d}, \mathbf{b} \in \mathbb{R}^m$. Let $F \subseteq P$ such that $F \neq \emptyset, P$. The following are equivalent.

---

[2]This is easily seen by the the transformation $\mathbf{x} = \mathbf{x}^1 - \mathbf{x}^2$.

(i) $F$ is a face of $P$.

(ii) $F$ is an exposed face of $P$.

(iii) There exists a subset $I \subseteq \{1, \ldots, m\}$ such that $F = \{\mathbf{x} \in P : A_I \mathbf{x} = \mathbf{b}_I\}$.

*Proof.* $(i) \Rightarrow (ii)$. Consider $\bar{\mathbf{x}} \in \mathrm{relint}(F)$ (which exists by Exercise 4). Since $F$ is a proper face, by Theorem 2.40, $\bar{\mathbf{x}} \in \mathrm{relbd}(P)$. By Theorem 2.39, there exists a supporting hyperplane at $\bar{\mathbf{x}}$ given by $\langle \mathbf{a}, \mathbf{x} \rangle \leq \delta$. Let $\{\mathbf{y} \in P : \langle \mathbf{a}, \mathbf{y} \rangle = \delta\}$ be the corresponding exposed face. Since $\mathbf{x} \in \mathrm{relint}(F)$, one can show that $F \subseteq \{\mathbf{y} \in P : \langle \mathbf{a}, \mathbf{y} \rangle = \delta\}$ (Verify!!). Thus, there exists an exposed face containing $F$. Let $F'$ be the minimal (with respect to set inclusion) exposed face of $P$ that contains $F$, i.e., for any other exposed face $F'' \supseteq F$, we have $F' \subseteq F''$. Let this exposed face $F'$ by defined by the valid inequality $\langle \mathbf{c}^1, \mathbf{x} \rangle \leq \delta_1$ for $P$.

If $F = F'$, then we are done because $F'$ is an exposed face. Otherwise, $F \subsetneq F'$, and so $F$ is a face of $F'$. Therefore, $\bar{\mathbf{x}} \in \mathrm{relbd}(F')$. Applying Theorem 2.39 to $F'$ and $\bar{\mathbf{x}}$, we obtain $\mathbf{c}^2 \in \mathbb{R}^d, \delta_2 \in \mathbb{R}$ such that $F \subseteq F' \cap \{\mathbf{y} \in \mathbb{R}^d : \langle \mathbf{c}^2, \mathbf{y} \rangle = \delta_2\}$, and there exists $\bar{\mathbf{y}} \in F'$ such that $\langle \mathbf{c}^2, \bar{\mathbf{y}} \rangle < \delta_2$. Using Theorem 2.83, we find finite sets $V, R$ such that $P = \mathrm{conv}(V) + \mathrm{cone}(R)$. Notice that since $P \subseteq H^-(\mathbf{c}^1, \delta_1)$, we must have $\langle \mathbf{c}^1, \mathbf{v} \rangle \leq \delta_1$ for all $\mathbf{v} \in V$ and $\langle \mathbf{c}^1, \mathbf{r} \rangle \leq 0$ for all $\mathbf{r} \in R$.

**Claim 1.** One can always choose $\lambda \geq 0$ such that $\lambda \mathbf{c}^1 + \mathbf{c}^2, \lambda \delta_1 + \delta_2$ satisfy

$$\langle \lambda \mathbf{c}^1 + \mathbf{c}^2, \mathbf{v} \rangle \leq \lambda \delta_1 + \delta_2 \text{ for all } \mathbf{v} \in V, \quad \langle \lambda \mathbf{c}^1 + \mathbf{c}^2, \mathbf{r} \rangle \leq 0, \text{ for all } \mathbf{r} \in R.$$

*Proof of Claim.* The relations can be rearranged to say

$$\langle \mathbf{c}^2, \mathbf{v} \rangle - \delta_2 \leq \lambda(\delta_1 - \langle \mathbf{c}^1, \mathbf{v} \rangle) \text{ for all } \mathbf{v} \in V \quad \langle \mathbf{c}^2, \mathbf{r} \rangle \leq \lambda(-\langle \mathbf{c}^1, \mathbf{r} \rangle), \text{ for all } \mathbf{r} \in R. \qquad (2.2)$$

First, recall that $0 \leq \delta_1 - \langle \mathbf{c}^1, \mathbf{v} \rangle$ for all $\mathbf{v} \in V$ and $0 \leq -\langle \mathbf{c}^1, \mathbf{r} \rangle$ for all $\mathbf{r} \in R$. Notice that since $F' \subseteq H^-(\mathbf{c}^2, \delta_2)$, if $\langle \mathbf{c}^1, \mathbf{v} \rangle = \delta_1$ for some $\mathbf{v} \in V$, this means that $\mathbf{v} \in F'$ and therefore $\langle \mathbf{c}^2, \mathbf{v} \rangle \leq \delta_2$. Similarly, if $\langle \mathbf{c}^1, \mathbf{r} \rangle = 0$ for some $\mathbf{r} \in R$, this means that $\mathbf{r} \in \mathrm{rec}(F')$ and therefore $\langle \mathbf{c}^2, \mathbf{r} \rangle \leq 0$. Thus, the following choice of

$$\lambda := \max \left\{ 0, \max_{\mathbf{v} \in V : \langle \delta_1 - \langle \mathbf{c}^1, \mathbf{v} \rangle > 0} \frac{\langle \mathbf{c}^2, \mathbf{v} \rangle - \delta_2}{\delta_1 - \langle \mathbf{c}^1, \mathbf{v} \rangle}, \max_{\mathbf{r} \in R : -\langle \mathbf{c}^1, \mathbf{r} \rangle > 0} \frac{\langle \mathbf{c}^2, \mathbf{r} \rangle}{-\langle \mathbf{c}^1, \mathbf{r} \rangle} \right\}$$

satisfies (2.2). $\qquad \square$

Using the $\lambda$ from the above claim, $X = P \cap \{\mathbf{y} \in \mathbb{R}^d : \langle \lambda \mathbf{c}^1 + \mathbf{c}^2, \mathbf{y} \rangle = \lambda \delta_1 + \delta_2\}$ is an exposed face of $P$ containing $F$. Moreover, $\langle \lambda \mathbf{c}^1 + \mathbf{c}^2, \mathbf{y} \rangle \leq \lambda \delta_1 + \delta_2$ is valid for $F'$ because the inequality is a nonnegative combination of the two valid inequalities $\langle \mathbf{c}^1, \bar{\mathbf{y}} \rangle \leq \delta_1, \langle \mathbf{c}^2, \bar{\mathbf{y}} \rangle \leq \delta_2$ for $F'$. Therefore, $X \subseteq F'$. But $\bar{\mathbf{y}}$ satisfies this inequality strictly, because it satisfies $\langle \mathbf{c}^2, \bar{\mathbf{y}} \rangle < \delta_2$, so $X \subsetneq F$. This contradicts the minimality of $F'$.

NOTES:                                                37

$(ii) \Rightarrow (iii)$. Let $\mathbf{c} \in \mathbb{R}^d, \delta \in \mathbb{R}$ be such that $F = P \cap \{\mathbf{x} : \langle \mathbf{c}, \mathbf{x} \rangle = \delta\}$. By Theorem 2.85, there exists $\mathbf{y} \in \mathbb{R}^m_+$ such that $\mathbf{c}^T = \mathbf{y}^T A$ and $\delta \geq \mathbf{y}^T \mathbf{b}$. Consider any $\mathbf{x} \in F$ (recall that $F$ is assumed to be nonempty). Then

$$\delta = \langle \mathbf{c}, \mathbf{x} \rangle = \langle \mathbf{y}^T A, \mathbf{x} \rangle = \mathbf{y}^T A\mathbf{x} \leq \mathbf{y}^T \mathbf{b} \leq \delta. \tag{2.3}$$

Thus, equality must hold everywhere and $\mathbf{y}^T \mathbf{b} = \delta$. Moreover, $\mathbf{y}^T A\mathbf{x} = \mathbf{y}^T \mathbf{b}$ for all $\mathbf{x} \in F$, which implies that $\mathbf{y}^T (A\mathbf{x} - \mathbf{b}) = 0$ for all $\mathbf{x} \in F$. This last relation says that for any $i \in \{1, \ldots, m\}$, if $\mathbf{y}_i > 0$ then $\langle \mathbf{a}^i, \mathbf{x} \rangle = \mathbf{b}_i$ for every $\mathbf{x} \in F$. Thus, setting $I = \{i : \mathbf{y}_i > 0\}$, we immediately obtain that $A_I \mathbf{x} = \mathbf{b}_I$ for all $\mathbf{x} \in F$. Consider any $\bar{\mathbf{x}} \in P$ satisfying $A_I \bar{\mathbf{x}} = \mathbf{b}_I$. Therefore, $\mathbf{y}^T A\bar{\mathbf{x}} = \mathbf{y}^T \mathbf{b}$ since $\mathbf{y}_i = 0$ for $i \notin I$. Therefore, $\mathbf{c}^T \mathbf{x} = \mathbf{y}^T A\bar{\mathbf{x}} = \mathbf{y}^T \mathbf{b} = \delta$, and thus, $\mathbf{x} \in P \cap \{\mathbf{x} : \langle \mathbf{c}, \mathbf{x} \rangle = \delta\} = F$.

$(iii) \Rightarrow (i)$. By definition, $F = \cap_{i \in I} F_i$, where $F_i = \{\mathbf{x} \in P : \langle \mathbf{a}^i, \mathbf{x} \rangle = \mathbf{b}_i\}$. By definition, each $F_i$ is an exposed face, and thus a face. By Problem 13 in "HW for Week III", the intersection of faces is a face and thus, $F$ is a face. $\qquad \square$

Here are some nice consequences of Theorem 2.89.

**Theorem 2.90.** The following are both true.

1. Every polyhedron has finitely many faces.

2. Every face of a polyhedron is a polyhedron.

### 2.5.4   Implicit equalities, dimension of polyhedra and facets

Given a polyhedron $P = \{\mathbf{x} : A\mathbf{x} \leq \mathbf{b}\}$ how can we decide the dimension of $P$? The concept of implicit equalities is important for this.

**Definition 2.91.** Let $A \in \mathbb{R}^{m \times d}$ and $\mathbf{b} \in \mathbb{R}^m$. We say that the inequality $\langle \mathbf{a}^i, \mathbf{x} \rangle \leq \mathbf{b}_i$ for some $i \in \{1, \ldots, m\}$ is an *implicit equality for the polyhedron* $P = \{\mathbf{x} : A\mathbf{x} \leq \mathbf{b}\}$ if $P \subseteq \{\mathbf{x} : \langle \mathbf{a}^i, \mathbf{x} \rangle = \mathbf{b}_i\}$, i.e., $P \subseteq H(\mathbf{a}^i, \mathbf{b}_i)$. We denote the subsystem of implicit equalities of $A\mathbf{x} \leq \mathbf{b}$ by $A_= \mathbf{x} \leq \mathbf{b}_=$. We will also use $A_+ \mathbf{x} \leq \mathbf{b}_+$ to denote the inequalities in $A\mathbf{x} \leq \mathbf{b}$ that are NOT implicit equalities.

Note that for each $i$ such that $\langle \mathbf{a}^i, \mathbf{x} \rangle \leq \mathbf{b}$ is not an implicit equality, there exists $\mathbf{x} \in P$ such that $\langle \mathbf{a}^i, \mathbf{x} \rangle < \mathbf{b}_i$.

**Exercise 6.** Let $P = \{\mathbf{x} : A\mathbf{x} \leq \mathbf{b}\}$. Show that there exists $\bar{\mathbf{x}} \in P$ such that $A_= \bar{\mathbf{x}} = \mathbf{b}_=$ and $A_+ \bar{\mathbf{x}} < \mathbf{b}_+$. Show the stronger statement that $\text{relint}(P) = \{\mathbf{x} \in \mathbb{R}^d : A_= \mathbf{x} = \mathbf{b}_=, \ A_+ \mathbf{x} < \mathbf{b}_+\}$.

We can completely characterize the affine hull of a polyhedron, and consequently its dimension, in terms of the implicit equalities.

**Proposition 2.92.** Let $A \in \mathbb{R}^{m \times d}$ and $\mathbf{b} \in \mathbb{R}^m$ and $P = \{\mathbf{x} : A\mathbf{x} \leq \mathbf{b}\}$. Then

$$\text{aff}(P) = \{\mathbf{x} \in \mathbb{R}^d : A_= \mathbf{x} = \mathbf{b}_=\} = \{\mathbf{x} \in \mathbb{R}^d : A_= \mathbf{x} \leq \mathbf{b}_=\}.$$

NOTES:                                                38

*Proof.* It is easy to verify that $\mathrm{aff}(P) \subseteq \{\mathbf{x} \in \mathbb{R}^d : A_=\mathbf{x} = \mathbf{b}_=\} \subseteq \{\mathbf{x} \in \mathbb{R}^d : A_=\mathbf{x} \leq \mathbf{b}_=\}$. We show that $\{\mathbf{x} \in \mathbb{R}^d : A_=\mathbf{x} \leq \mathbf{b}_=\} \subseteq \mathrm{aff}(P)$. Consider any $\mathbf{y}$ satisfying $A_=\mathbf{y} \leq \mathbf{b}_=$. Using Exercise 6, choose any $\bar{\mathbf{x}} \in P$ such that $A_=\bar{\mathbf{x}} = \mathbf{b}_=$ and $A_+\bar{\mathbf{x}} < \mathbf{b}_+$. If $A_+\mathbf{y} \leq \mathbf{b}_+$, then $\mathbf{y} \in P \subseteq \mathrm{aff}(P)$ and we are done. Otherwise, set

$$\mu := \min_{i : \langle \mathbf{a}^i, \mathbf{y} \rangle > \mathbf{b}_i} \left\{ \frac{\mathbf{b}_i - \langle \mathbf{a}^i, \bar{\mathbf{x}} \rangle}{\langle \mathbf{a}^i, \mathbf{y} \rangle - \langle \mathbf{a}^i, \bar{\mathbf{x}} \rangle} \right\}.$$

Observe that since $\langle \mathbf{a}^i, \mathbf{y} \rangle > \mathbf{b}_i > \langle \mathbf{a}^i, \bar{\mathbf{x}} \rangle$ for each $i$ considered in the minimum, we have $0 < \mu < 1$. One can check that $(1 - \mu)\bar{\mathbf{x}} + \mu\mathbf{y} \in P$. This shows that $\mathbf{y} \in \mathrm{aff}(P)$, because $\mathbf{y}$ is on the line joining two points in $P$, namely $\bar{\mathbf{x}}$ and $(1 - \mu)\bar{\mathbf{x}} + \mu\mathbf{y}$. $\square$

Combined with part 4. of Theorem 2.16, this gives the following corollary.

**Corollary 2.93.** Let $A \in \mathbb{R}^{m \times d}$ and $\mathbf{b} \in \mathbb{R}^m$ and $P = \{\mathbf{x} : A\mathbf{x} \leq \mathbf{b}\}$. Then

$$\dim(P) = d - \mathrm{rank}(A_=).$$

As we have seen before, a given description $P = \{\mathbf{x} : A\mathbf{x} \leq \mathbf{b}\}$ for a polyhedron may be redundant, in the sense, that we can remove some of the inequalities, and still have the same set $P$. This motivates the following definition.

**Definition 2.94.** Let $A \in \mathbb{R}^{m \times d}$ and $\mathbf{b} \in \mathbb{R}^m$. We say that the inequality $\langle \mathbf{a}^i, \mathbf{x} \rangle \leq \mathbf{b}_i$ for some $i \in \{1, \ldots, m\}$ is *redundant for the polyhedron* $P = \{\mathbf{x} : A\mathbf{x} \leq \mathbf{b}\}$ if $P = \{\mathbf{x} : A_{-i}\mathbf{x} \leq \mathbf{b}_{-i}\}$, where $A_{-i}$ denotes the matrix $A$ without row $i$ and $\mathbf{b}_{-i}$ is the vector $\mathbf{b}$ with the $i$-th coordinate removed. Otherwise, if $P \subsetneq \{\mathbf{x} : A_{-i}\mathbf{x} \leq \mathbf{b}_{-i}\}$, then $\langle \mathbf{a}^i, \mathbf{x} \rangle \leq \mathbf{b}_i$ is said to be *irredundant for* $P$. The system $A\mathbf{x} \leq \mathbf{b}$ is said to be an irredundant system if every inequality is irredundant for $P = \{\mathbf{x} : A\mathbf{x} \leq \mathbf{b}\}$.

The following characterization of facets of a polyhedron is quite useful, specially in combinatorial optimization and polyhedral combinatorics.

**Theorem 2.95.** Let $P = \{\mathbf{x} \in \mathbb{R}^d : A\mathbf{x} \leq \mathbf{b}\}$ be nonempty with $A \in \mathbb{R}^{m \times d}, \mathbf{b} \in \mathbb{R}^m$ giving an irredundant system. Let $F \subseteq P$. The following are equivalent.

(i) $F$ is a facet of $P$, i.e., $F$ is a face with $\dim(F) = \dim(P) - 1$.

(ii) $F$ is maximal, proper face of $P$, i.e., for any proper face $F' \supseteq F$, we must have $F' = F$.

(iii) There exists a unique $i \in \{1, \ldots, m\}$ such that $F = \{\mathbf{x} \in P : \langle \mathbf{a}^i, \mathbf{x} \rangle = \mathbf{b}_i\}$ and $\langle \mathbf{a}^i, \mathbf{x} \rangle \leq \mathbf{b}_i$ is not an implicit equality.

*Proof.* $(i) \Rightarrow (ii)$. Suppose to the contrary that there exists a proper face $F' \supsetneq F$. Observe that $F$ is a face of $F'$ by Problem 15 in "HW for Week IV", and so $F$ is a proper face of $F'$. By Lemma 2.35,

dim($F'$) > dim($F$) = dim($P$) − 1. So, dim($F'$) = dim($P$). This contradicts the fact that $F'$ is proper face, by Lemma 2.35.

$(ii) \Rightarrow (iii)$. By Theorem 2.89, there exists a subset of indices $I \subseteq \{1, \ldots, m\}$ such that $F = \{\mathbf{x} \in \mathbb{R}^d : A\mathbf{x} \le \mathbf{b}, \; A_I\mathbf{x} = \mathbf{b}_I\}$. If all the inequalities indexed by $I$ are implicit equalities for $P$, then $F = P$, contradicting the assumption that $F$ is a proper face. So there exists $i \in I$ such that $\langle \mathbf{a}^i, \mathbf{x} \rangle \le \mathbf{b}_i$ is not an implicit equality. Let $F' = \{\mathbf{x} \in P : \langle \mathbf{a}^i, \mathbf{x} \rangle = \mathbf{b}_i\}$ be the face defined by this inequality; since $\langle \mathbf{a}^i, \mathbf{x} \rangle \le \mathbf{b}_i$ is not an implicit equality, $F'$ is a proper face of $P$. Also observe that $F \subseteq F'$. Hence $F = F' = \{x \in P : \langle \mathbf{a}^i, \mathbf{x} \rangle = \mathbf{b}_i\}$ by maximality of $F$. To show uniqueness of $i$, we would like to show that $I = \{i\}$. We show this by exhibiting $\mathbf{x}^0 \in F$ with the following property: for any $j \ne i$ such that $\langle \mathbf{a}^j, \mathbf{x} \rangle \le \mathbf{b}_j$ is not an implicit equality, we have $\langle \mathbf{a}^j, \mathbf{x}^0 \rangle < \mathbf{b}_j$. To see this, let $\mathbf{x}^1 \in P$ such that $A_=\mathbf{x}^1 = \mathbf{b}_=$ and $A_+\mathbf{x}^1 < \mathbf{b}_+$ (such an $\mathbf{x}^1$ exists by Exercise 6). Since $A\mathbf{x} \le \mathbf{b}$ is an irredundant system, if we remove the inequality indexed by $i$, then we get some new points that satisfy the rest of the inequalities, but which violate $\langle \mathbf{a}^i, \mathbf{x} \rangle \le \mathbf{b}_i$. More precisely, there exists $\mathbf{x}^2 \in \mathbb{R}^d$ such that $A_=\mathbf{x}^2 = \mathbf{b}_=$, $A_+^{-i}\mathbf{x}^2 \le \mathbf{b}_+^{-i}$ and $\langle \mathbf{a}^i, \mathbf{x}^2 \rangle > \mathbf{b}_i$, where $A_+^{-i}\mathbf{x} \le \mathbf{b}_+^{-i}$ denotes the system $A_+\mathbf{x} \le \mathbf{b}_+$ without the inequality indexed by $i$. Since $\langle \mathbf{a}^i, \mathbf{x}^1 \rangle < \mathbf{b}_i$ and $\langle \mathbf{a}^i, \mathbf{x}^2 \rangle > \mathbf{b}_i$, there exists a convex combination of $\mathbf{x}^1, \mathbf{x}^2$ such that this convex combination $\mathbf{x}^0$ satisfies $\langle \mathbf{a}^i, \mathbf{x}^0 \rangle = \mathbf{b}_i$. Since $A_=\mathbf{x}^1 = \mathbf{b}_=$ and $A_=\mathbf{x}^2 = \mathbf{b}_=$, we must have $A_=\mathbf{x}^0 = \mathbf{b}_=$. Moreover, since $A_+\mathbf{x}^1 < \mathbf{b}_+$ and $A_+^{-i}\mathbf{x}^2 \le \mathbf{b}_+^{-i}$, we must have that for any $j \ne i$ indexing an inequality in $A_+\mathbf{x} \le \mathbf{b}_+$, it must satisfy $\langle \mathbf{a}^j, \mathbf{x}^0 \rangle < \mathbf{b}_j$. Thus, we are done.

$(iii) \Rightarrow (i)$. By Theorem 2.89, $F$ is a face. We now establish that dim($F$) = dim($P$) − 1. Let $\mathcal{J}$ denote the set of indices that index inequalities in $A\mathbf{x} \le \mathbf{b}$ that are not implicit equalities. Since there exists a unique $i \in \mathcal{J}$ such that $F = \{x \in P : \langle \mathbf{a}^i, \mathbf{x} \rangle = \mathbf{b}_i\}$, this means that for any $j \in \mathcal{J} \setminus i$, there exists $\mathbf{x}^j \in F$ such that $\langle \mathbf{a}^j, \mathbf{x}^j \rangle < \mathbf{b}_j$. Now let $\mathbf{x}^0 = \frac{1}{|\mathcal{J}|} \sum_{j \in \mathcal{J} \setminus \{i\}} \mathbf{x}^j$, and observe that $\mathbf{x}^0 \in F$ and for any $j \in \mathcal{J} \setminus i$, we have $\langle \mathbf{a}^j, \mathbf{x}^0 \rangle < \mathbf{b}_j$. Let us describe the polyhedron $F$ by the system $\tilde{A}\mathbf{x} \le \tilde{\mathbf{b}}$ that appends the inequality $\langle -\mathbf{a}^i, \mathbf{x} \rangle \le -\mathbf{b}_i$ to the system $A\mathbf{x} \le \mathbf{b}$.

**Claim 2.** rank($\tilde{A}_=$) = rank($A_=$) + 1.

*Proof.* The properties of $\mathbf{x}^0$ show that the matrix $\tilde{A}_=$ is simply the matrix $A_=$ appended with $\mathbf{a}^i$. So it suffices to show that $\mathbf{a}^i$ is not a linear combination of the rows of $A_=$. Suppose to the contrary that $\mathbf{a}^i = \mathbf{y}^T A_=$ for some $\mathbf{y} \in \mathbb{R}^k$ where $k$ is the number of rows of $A_=$. If $\mathbf{b}_i < \mathbf{y}^T \mathbf{b}_=$, then $P$ is empty because any $\mathbf{x} \in P$ satisfies $A_=\mathbf{x} = \mathbf{b}_=$, and therefore must satisfy $\mathbf{y}^T A_=\mathbf{x} = \mathbf{y}^T \mathbf{b}_=$ and this contradicts $\mathbf{y}^T A_=\mathbf{x} = \langle \mathbf{a}^i, \mathbf{x} \rangle \le \mathbf{b}_i$. If $\mathbf{b}_i \ge \mathbf{y}^T \mathbf{b}_=$, then $\langle \mathbf{a}^i, \mathbf{x} \rangle \le \mathbf{b}_i$ is redundant for $P$, as every $\mathbf{x}$ satisfying $A_=\mathbf{x} = \mathbf{b}_=$ satisfies $\langle \mathbf{a}^i, \mathbf{x} \rangle \le \mathbf{b}_i$. $\square$

Using Corollary 2.93, we obtain that dim($F$) = $d$ − rank($\tilde{A}_=$) = $d$ − rank($A_=$) − 1 = dim($P$) − 1. $\square$

A consequence of this characterization of facets is that full-dimensional polyhedra have a unique system describing them, upto scaling.

**Definition 2.96.** We say that the inequality $\langle \mathbf{a}, \mathbf{x} \rangle \leq \delta$ is *equivalent* to the inequality $\langle \mathbf{a}', \mathbf{x} \rangle \leq \delta'$ if there exists $\lambda \geq 0$ such that $\mathbf{a}' = \lambda \mathbf{a}$ and $\delta' = \lambda \delta$. Equivalent inequalities define the same halfspace, i.e., $H^-(\mathbf{a}, \delta) = H^-(\mathbf{a}', \delta')$.

**Theorem 2.97.** Let $P$ be a full-dimensional polyehdron. Let $A \in \mathbb{R}^{m \times d}$ matrix, $A' \in \mathbb{R}^{p \times d}$, $\mathbf{b} \in \mathbb{R}^m$ and $\mathbf{b}' \in \mathbb{R}^p$ be such that $A\mathbf{x} \leq \mathbf{b}$ and $A'\mathbf{x} \leq \mathbf{b}'$ are both irredundant systems describing $P$, i.e.,

$$\{\mathbf{x} \in \mathbb{R}^d : A\mathbf{x} \leq \mathbf{b}\} = \{\mathbf{x} \in \mathbb{R}^d : A'\mathbf{x} \leq \mathbf{b}'\} = P.$$

Then both systems are the same upto permutation and scaling. More precisely, the following holds:

1. $m = p$.

2. There exists permutation $\sigma : \{1, \ldots, m\} \to \{1, \ldots, m\}$ such that for each $i \in \{1, \ldots, m\}$, $\langle \mathbf{a}^i, \mathbf{x} \rangle \leq \mathbf{b}_i$ is equivalent to $\langle \mathbf{a}'^{\sigma(i)}, \mathbf{x} \rangle \leq \mathbf{b}'_{\sigma(i)}$.

*Proof.* Left as an exercise. $\qquad\square$

# 3   Convex Functions

We now turn our attention to convex functions, as a step towards optimization. In this context, we will need to sometimes talk about the extended real numbers $\mathbb{R} \cup \{-\infty, +\infty\}$. One reason is that in optimization problems, many times a supremum may be $\infty$ or an infimum may be $-\infty$, and using them on the same footing as the reals makes certain statements nicer, without having to exclude annoying special cases. For this, one needs to set up some convenient rules for arithmetic over $\mathbb{R} \cup \{-\infty, +\infty\}$:

- $x + \infty = \infty$ for any $x \in \mathbb{R} \cup \{+\infty\}$.

- $x(+\infty) = +\infty$ for all $x > 0$. We will avoid situations where we need to consider $0 \cdot +\infty$.

- $x < \infty$ for all $x \in \mathbb{R}$.

## 3.1   General properties, epigraphs, subgradients

**Definition 3.1.** A function $f : \mathbb{R}^d \to \mathbb{R} \cup \{\infty\}$ is called *convex* if

$$f(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y}),$$

for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ and $\lambda \in (0, 1)$. If the inequality is strict for all $\mathbf{x} \neq \mathbf{y}$, then the function is called *strictly convex*. The *domain* (sometimes also called *effective domain*) of $f$ is defined as

$$\mathrm{dom}(f) := \{\mathbf{x} \in \mathbb{R}^d : f(\mathbf{x}) < +\infty\}.$$

A function $g$ is said to be *(strictly) concave* if $-g$ is (strictly) convex.

NOTES:                                                    41

926    The domain of a convex function is easily seen to be convex.

927    **Proposition 3.2.** Let $f : \mathbb{R}^d \to \mathbb{R} \cup \{\infty\}$ be a convex function. Then $\mathrm{dom}(f)$ is a convex set.

928    *Proof.* Left as an exercise.                                                                                                                    □

929    The following subfamily of convex functions is nicer to deal with from an algorithmic perspective.

**Definition 3.3.** A function $f : \mathbb{R}^d \to \mathbb{R} \cup \{\infty\}$ is called *strongly convex* with *modulus of strong convexity* $c > 0$ if
$$f(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) + \leq \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y}) - \frac{1}{2}c\lambda(1 - \lambda)\|\mathbf{x} - \mathbf{y}\|^2,$$

930    for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ and $\lambda \in (0, 1)$.

931    The above definition will become particularly intuitive when we speak of differentiable convex functions
932    in Section 3.3. Even so, the following proposition sheds some light on strongly convex functions.

933    **Proposition 3.4.** A function $f : \mathbb{R}^d \to \mathbb{R} \cup \{\infty\}$ is strongly convex modulus of strong convexity $c > 0$ if
934    and only if the function $g(\mathbf{x}) := f(\mathbf{x}) - \frac{1}{2}c\|\mathbf{x}\|^2$ is convex.

935    Convex functions have a natural convex set associated with them, called the *epigraph*. Many properties of
936    convex functions can be obtained by just analyzing the corresponding epigraph and using all the technology
937    built in Section 2. We give the formal definition for general functions below; very informally, it is "the region
938    above the graph of a function".

**Definition 3.5.** Let $f : \mathbb{R}^d \to \mathbb{R} \cup \{\infty\}$ be any function (not necessarily convex). The *epigraph* of $f$ is
defined as
$$\mathrm{epi}(f) := \{(\mathbf{x}, t) \in \mathbb{R}^n \times \mathbb{R} : f(\mathbf{x}) \leq t\}.$$

939    Note that $\mathrm{epi}(f) \subseteq \mathbb{R}^d \times \mathbb{R}$, so it lives in a space whose dimension is one more than the space over which
940    the function is defined, just like the graph of the function. **Note also that the epigraph is nonempty**
941    **if and only if the function is not identically equal to** $+\infty$. Convex functions are precisely those
942    functions whose epigraphs are convex.

943    **Proposition 3.6.** Let $f : \mathbb{R}^d \to \mathbb{R} \cup \{\infty\}$ be any function. $f$ is convex if and only if $\mathrm{epi}(f)$ is a convex set.

944    *Proof.* ($\Rightarrow$) Consider any $(\mathbf{x}^1, t_1), (\mathbf{x}^2, t_2) \in \mathrm{epi}(f)$, and any $\lambda \in (0, 1)$.
945    The result is a consequence of the following sequence of implications:

$$
\begin{aligned}
&(\mathbf{x}^1, t_1) \in \mathrm{epi}(f), \ \ (\mathbf{x}^2, t_2) \in \mathrm{epi}(f), \ \ f \text{ is convex} \\
\Rightarrow \ \ & f(\mathbf{x}^1) \leq t_1, \ \ f(\mathbf{x}^2) \leq t_2, \ \ f(\lambda \mathbf{x}^1 + (1 - \lambda)\mathbf{x}^2) \leq \lambda f(\mathbf{x}^1) + (1 - \lambda)f(\mathbf{x}^2) \\
\Rightarrow \ \ & f(\lambda \mathbf{x}^1 + (1 - \lambda)\mathbf{x}^2) \leq \lambda t_1 + (1 - \lambda)t_2 \\
\Rightarrow \ \ & (\lambda \mathbf{x}^1 + (1 - \lambda)\mathbf{x}^2, \lambda t_1 + (1 - \lambda)t_2) \in \mathrm{epi}(f)
\end{aligned}
$$

NOTES:                                          42

( ⟸ ) Consider the any $\mathbf{x}^1, \mathbf{x}^2 \in \mathbb{R}^d$ and $\lambda \in (0,1)$. The points $(\mathbf{x}^1, f(\mathbf{x}^1)), (\mathbf{x}^2, f(\mathbf{x}^2))$ both lie in epi($f$). By convexity of epi($f$), we have that $(\lambda \mathbf{x}^1 + (1-\lambda)\mathbf{x}^2, \lambda f(\mathbf{x}^1) + (1-\lambda)f(\mathbf{x}^2)) \in$ epi($f$). This implies that $f(\lambda \mathbf{x}^1 + (1-\lambda)\mathbf{x}^2) \leq \lambda f(\mathbf{x}^1) + (1-\lambda)f(\mathbf{x}^2)$, showing that $f$ is convex. □

Just like the class of *closed*, convex sets are nicer to deal with compared sets that simply convex but not closed (mainly because of the separating/supporting hyperplane theorem), it will be convenient to isolate a similar class of "nicer" convex functions.

**Definition 3.7.** A function is said to be a *closed, convex function* if its epigraph is a closed, convex set.

One can associate another family of convex sets with a convex function.

**Definition 3.8.** Let $f : \mathbb{R}^d \to \mathbb{R} \cup \{\infty\}$ be any function. Given $\alpha \in \mathbb{R}$, the *$\alpha$-sublevel set* of $f$ is the set

$$f_\alpha := \{\mathbf{x} \in \mathbb{R}^d : f(\mathbf{x}) \leq \alpha\}.$$

The following can be verified by the reader.

**Proposition 3.9.** All sublevel sets of a convex function are convex sets.

The converse of Proposition 3.9 is *not true*. Functions whose sublevel sets are all convex are called *quasi-convex*.

**Example 3.10.**     1. *Indicator function.* For any subset $X \subseteq \mathbb{R}^d$, define

$$I_X(\mathbf{x}) := \left\{ \begin{array}{cl} 0 & \text{if } \mathbf{x} \in X \\ +\infty & \text{if } \mathbf{x} \notin X \end{array} \right.$$

Then $I_X$ is convex if and only if $X$ is convex.

2. *Linear/Affine function.* Let $\mathbf{a} \in \mathbb{R}^d$ and $\delta \in \mathbb{R}$. Then the function $\mathbf{x} \mapsto \langle \mathbf{a}, \mathbf{x} \rangle + \delta$ is called an *affine function* (if $\delta = 0$, this is a *linear function*). It is easily verified that affine functions are convex.

3. *Norms and Distances.* Let $N : \mathbb{R}^d \to \mathbb{R}$ be a norm (see Definition 1.1). Then $N$ is convex (Verify !!). Let $C$ be a nonempty convex set. Then the distance function associated with the norm $N$, defined as

$$d_C^N(\mathbf{x}) := \inf_{\mathbf{y} \in C} N(\mathbf{y} - \mathbf{x})$$

is a convex function.

4. *Maximum of affine functions/Piecewise linear/Polyhedral function.* Let $\mathbf{a}^1, \ldots, \mathbf{a}^m \in \mathbb{R}^d$ and $\delta_1, \ldots, \delta_m \in \mathbb{R}$. The function

$$f(\mathbf{x}) := \max_{i=1,\ldots,m} (\langle \mathbf{a}^i, \mathbf{x} \rangle + \delta_i)$$

NOTES:                                    43

is a convex function. Let us verify this. Consider any $\mathbf{x}^1, \mathbf{x}^2 \in \mathbb{R}^d$ and $\lambda \in (0,1)$. Then,

$$
\begin{aligned}
f(\lambda \mathbf{x}^1 + (1-\lambda)\mathbf{x}^2) &= \max_{i=1,\ldots,m}(\langle \mathbf{a}^i, \lambda \mathbf{x}^1 + (1-\lambda)\mathbf{x}^2 \rangle + \delta_i) \\
&= \max_{i=1,\ldots,m}\left(\lambda(\langle \mathbf{a}^i, \mathbf{x}^1 \rangle + \delta_i) + (1-\lambda)(\langle \mathbf{a}^i, \mathbf{x}^2 \rangle + \delta_i)\right) \\
&\leq \max_{i=1,\ldots,m}\left(\lambda(\langle \mathbf{a}^i, \mathbf{x}^1 \rangle + \delta_i)\right) + \max_{i=1,\ldots,m}\left((1-\lambda)(\langle \mathbf{a}^i, \mathbf{x}^2 \rangle + \delta_i)\right) \\
&= \lambda \max_{i=1,\ldots,m}(\langle \mathbf{a}^i, \mathbf{x}^1 \rangle + \delta_i) + (1-\lambda) \max_{i=1,\ldots,m}(\langle \mathbf{a}^i, \mathbf{x}^2 \rangle + \delta_i) \\
&= \lambda f(\mathbf{x}^1) + (1-\lambda)f(\mathbf{x}^2)
\end{aligned}
$$

The inequality follows from the fact that if $\ell_1,\ldots,\ell_m$ and $u_1,\ldots,u_m$ be two sets of $m$ real numbers for some $m \in \mathbb{N}$, then $\max_{i=1,\ldots,m}(\ell_i + u_i) \leq \max_{i=1,\ldots,m}\ell_i + \max_{i=1,\ldots,m}u_i$.

An important consequence of the definition of convexity for functions is Jensen's inequality which sees its uses in diverse areas of science and engineering.

**Theorem 3.11.** [Jensen's Inequality] Let $f : \mathbb{R}^d \to \mathbb{R} \cup \{\infty\}$ be any function. Then $f$ is convex if and only if for any finite set of points $\mathbf{x}^1,\ldots,\mathbf{x}^n \in \mathbb{R}^d$ and $\lambda_1,\ldots,\lambda_n \geq 0$ such that $\lambda_1 + \ldots, \lambda_n = 1$, the following holds:

$$
f(\lambda_1 \mathbf{x}^1 + \ldots + \lambda_n \mathbf{x}^n) \leq \lambda_1 f(\mathbf{x}^1) + \ldots \lambda_n f(\mathbf{x}^n).
$$

*Proof.* ($\Leftarrow$) Just use the hypothesis with $n = 2$.

($\Rightarrow$) It suffices to show the inequality when all $\lambda_i > 0$. If any $f(\mathbf{x}^i)$ is $+\infty$, then the inequality holds trivially. So we assume that each $f(\mathbf{x}^i) < +\infty$. By Proposition 3.6, $\text{epi}(f)$ is a convex set. For each $i = 1,\ldots,m$, the point $(\mathbf{x}^i, f(\mathbf{x}^i) \in \text{epi}(f)$ by definition of $\text{epi}(f)$. Since $\text{epi}(f)$ is convex, $\sum_{i=1}^m \lambda_i(\mathbf{x}^i, f(\mathbf{x}^i) \in \text{epi}(f)$, i.e., $(\lambda_1 \mathbf{x}^1 + \ldots + \lambda_n \mathbf{x}^n, \lambda_1 f(\mathbf{x}^1) + \ldots \lambda_n f(\mathbf{x}^n)) \in \text{epi}(f)$. Therefore, $f(\lambda_1 \mathbf{x}^1 + \ldots + \lambda_n \mathbf{x}^n) \leq \lambda_1 f(\mathbf{x}^1) + \ldots \lambda_n f(\mathbf{x}^n)$. $\square$

Recall Theorem 2.3 that showed convexity of a set is preserved under certain operations. We would like to develop a similar result for convex functions.

**Theorem 3.12.** [Operations that preserve the property of being a (closed) convex function] Let $f_i : \mathbb{R}^d \to \mathbb{R} \cup \{+\infty\}$, $i \in I$ be a family of (closed) convex functions where the index set $I$ is potentially infinite. The following are all true.

1. (Nonnegative combinations). If $I$ is a finite set, and $\alpha_i \geq 0$, $i \in I$ be a corresponding set of nonnegative reals, then $\sum_{i \in I} \alpha_i f_i$ is a (closed) convex function.

2. (Taking supremums). The function defined as $g(\mathbf{x}) := \sup_{i \in I} f_i(\mathbf{x})$ is a (closed) convex function (even when $I$ is uncountable infinite).

3. (Pre-Composition with an affine function). Let $A \in \mathbb{R}^{m \times d}$ and $\mathbf{b} \in \mathbb{R}^m$ and let $f : \mathbb{R}^m \to \mathbb{R}$ be any (closed) convex function on $\mathbb{R}^m$. Then $g(\mathbf{x}) := f(A\mathbf{x} + \mathbf{b})$ as a function from $\mathbb{R}^d \to \mathbb{R}$ is a (closed) convex function.

NOTES: 44

4. (Post-Composition with an increasing convex function). Let $h : \mathbb{R} \to \mathbb{R}$ be a (closed) convex function that is also increasing, i.e., $h(x) \geq h(y)$ when $x \geq y$. Let $f : \mathbb{R}^d \to \mathbb{R}$ be a (closed) convex function such that for some $\mathbf{x} \in \mathbb{R}^d$, $f(\mathbf{x}) \in \text{dom}(h)$. We adopt the convention that $h(+\infty) = +\infty$. Then $h(f(\mathbf{x}))$ as a function from $\mathbb{R}^d \to \mathbb{R}$ is a (closed) convex function.

*Proof.*     1. Let $F = \sum_{i \in I} \alpha_i f_i$. Consider any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ and $\lambda \in (0, 1)$. Then

$$
\begin{aligned}
F(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) &= \sum_{i \in I} \alpha_i f_i(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) \\
&\leq \sum_{i \in I} \alpha_i (\lambda f_i(\mathbf{x}) + (1 - \lambda)f_i(\mathbf{y})) \\
&= \lambda \sum_{i \in I} \alpha_i f_i(\mathbf{x}) + (1 - \lambda) \sum_{i \in I} \alpha_i f_i(\mathbf{y}) \\
&= \lambda F(\mathbf{x}) + (1 - \lambda)F(\mathbf{y})
\end{aligned}
$$

We use the non negativity of $\alpha_i$ in the inequality on the second displayed line above. We omit the proof of closedness of the function.

2. The main observation is that $\text{epi}(g) = \cap_{i \in I} \text{epi}(f_i)$ because $g(\mathbf{x}) \leq t$ if and only if $f_i(\mathbf{x}) \leq t$ for all $i \in I$. Since the intersection of (closed) convex sets is a (closed) convex set (part 1. of Theorem 2.3), we have the result.

3. The main observation is that for any $\mathbf{x} \in \mathbb{R}^d$ and $t \in \mathbb{R}$, $(\mathbf{x}, t) \in \text{epi}(g)$ if and only if $(A\mathbf{x}+\mathbf{b}, t) \in \text{epi}(f)$. Define the affine map $T : \mathbb{R}^d \times \mathbb{R} \to \mathbb{R}^m \times \mathbb{R}$ as follows $T(\mathbf{x}, t) = (A\mathbf{x}+b, t)$. Then $\text{epi}(g) = T^{-1}(\text{epi}(f)$. Since the pre-image of a (closed) convex set with respect to an affine transformation is (closed) convex (part 4. of Theorem 2.3), we obtain that $\text{epi}(g)$ is (closed) convex.

4. Left as an exercise.

$\square$

We can now see some more interesting examples of convex functions.

**Example 3.13.**     1. Let $\mathbf{a}^i \in \mathbb{R}^d$ and $\delta_i \in \mathbb{R}$ for some index set $i \in I$. Then the function

$$
f(\mathbf{x}) := \sup_{i \in I}(\langle \mathbf{a}^i, \mathbf{x} \rangle + \delta_i)
$$

is closed convex. This is an alternate proof of the convexity of the maximum of finitely many affine functions – part 4. of Example 3.10.

2. Consider the vector space $V$ of symmetric $n \times n$ matrices. One can view $V$ as $\mathbb{R}^{\frac{n(n+1)}{2}}$. Let $k \leq n$. Consider the function $f_k : V \to \mathbb{R}$ which takes a matrix $X$ and maps it to $f(X)$ which is the sum of the $k$ largest eigenvalues of $X$. Then $f_k$ is a convex function. This is seen by the following argument. Given any $Y \in V$ define the linear function $A_Y$ on $V$ as follows: $A_Y(X) = \sum_{i,j} X_{ij}Y_{ij}$. Then

$$
f_k(X) = \sup_{Y \in \Omega} A_{YY^T}(X),
$$

where $\Omega$ is the set of $n \times k$ matrices with $k$ orthonormal columns in $\mathbb{R}^n$. This shows that $f_k$ is the supremum of linear functions, and by Theorem 3.12, it is closed convex.

We see in part 1. of Example 3.13 that the supremum of affine functions is convex. We will show below that, in fact, every convex function is the supremum of some family of affine functions. This is analogous to the fact that all closed convex sets are the intersection of some family of halfspaces. We build up to this with an important definition.

**Definition 3.14.** Let $f : \mathbb{R}^d \to \mathbb{R} \cup \{\infty\}$ be any function. Let $\mathbf{x} \in \mathrm{dom}(f)$. Then $\mathbf{a} \in \mathbb{R}^d$ is said to define an *affine support of $f$ at $\mathbf{x}$* if $f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \mathbf{a}, \mathbf{y} - \mathbf{x} \rangle$ for all $\mathbf{y} \in \mathbb{R}^d$.

**Theorem 3.15.** Let $f : \mathbb{R}^d \to \mathbb{R}$ be any function. Then $f$ is closed convex if and only if there exists an affine support of $f$ at every $\mathbf{x} \in \mathbb{R}^d$.

*Proof.* ($\Rightarrow$) Consider any $\mathbf{x} \in \mathbb{R}^d$. By definition of closed convex, $\mathrm{epi}(f)$ is a closed convex set. Moreover, $(\mathbf{x}, f(\mathbf{x})) \in \mathrm{bd}(\mathrm{epi}(f))$. By Theorem 2.23, there exists $(\bar{\mathbf{a}}, r) \in \mathbb{R}^d \times \mathbb{R}$ and $\delta \in \mathbb{R}$ such that $\bar{\mathbf{a}}$ and $r$ are not both 0, and $\langle \bar{\mathbf{a}}, \mathbf{y} \rangle + rt \leq \delta$ for all $(\mathbf{y}, t) \in \mathrm{epi}(f)$, and $\langle \bar{\mathbf{a}}, \mathbf{x} \rangle + r f(\mathbf{x}) = \delta$.

We claim that $r < 0$. Suppose to the contrary that $r \geq 0$. First consider the case that $\bar{\mathbf{a}} = \mathbf{0}$, then $r > 0$. $(\mathbf{x}, t) \in \mathrm{epi}(f)$ for all $t \geq f(\mathbf{x})$. But this contradicts that $rt = \langle \bar{\mathbf{a}}, \mathbf{y} \rangle + rt \leq \delta$ for all $t \geq f(\mathbf{x})$ and $r f(\mathbf{x}) = \langle \bar{\mathbf{a}}, \mathbf{x} \rangle + r f(\mathbf{x}) = \delta$. Next consider the case that $\bar{\mathbf{a}} \neq \mathbf{0}$. Consider any $\mathbf{y} \in \mathbb{R}^d$ satisfying $\langle \bar{\mathbf{a}}, \mathbf{y} \rangle > \delta$. Since $f$ is real valued, there exists $(\mathbf{y}, t) \in \mathrm{epi}(f)$ for some $t \geq 0$. Since $r \geq 0$, this contradicts that $\langle \bar{\mathbf{a}}, \mathbf{y} \rangle + rt \leq \delta$.

Now set $\mathbf{a} = \frac{\bar{\mathbf{a}}}{-r}$. $\langle \bar{\mathbf{a}}, \mathbf{x} \rangle + r f(\mathbf{x}) = \delta$ and $\langle \bar{\mathbf{a}}, \mathbf{y} \rangle + r f(\mathbf{y}) \leq \delta$ for all $\mathbf{y} \in \mathbb{R}^d$ together imply that $\langle \bar{\mathbf{a}}, \mathbf{y} \rangle \leq (-r) f(\mathbf{y}) + \langle \bar{\mathbf{a}}, \mathbf{x} \rangle + r f(\mathbf{x})$. Rearranging, we obtain that $f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \mathbf{a}, \mathbf{y} - \mathbf{x} \rangle$ for all $\mathbf{y} \in \mathbb{R}^d$.

($\Leftarrow$) By definition of affine support, for every $\mathbf{x} \in \mathbb{R}^d$, there exists $\mathbf{a}_{\mathbf{x}} \in \mathbb{R}^d$ such that $f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \mathbf{a}_{\mathbf{x}}, \mathbf{y} - \mathbf{x} \rangle$ for all $\mathbf{y} \in \mathbb{R}^d$. This implies that, in fact,

$$f(\mathbf{y}) = \sup_{\mathbf{x} \in \mathbb{R}^d} (f(\mathbf{x}) + \langle \mathbf{a}_{\mathbf{x}}, \mathbf{y} - \mathbf{x} \rangle),$$

because setting $\mathbf{x} = \mathbf{y}$ on the right hand side gives $f(\mathbf{y})$. Thus, $f$ is the supremum of a family of affine functions, which by Example 3.13, shows that $f$ is closed convex. $\square$

**Remark 3.16.**    1. Any convex function that is finite valued everywhere is closed convex. This follows from a continuity result we will prove later. We skip the details in these notes. Thus, in the forward direction of Theorem 3.15, one may weaken the hypothesis to just convex, as opposed to closed convex.

2. In the reverse direction of Theorem 3.15, one may weaken the hypothesis to having *local* affine support everywhere. A function $f : \mathbb{R}^d \to \mathbb{R}$ is said to have local affine support at $\mathbf{x}$ if there exists $\epsilon > 0$ (depending on $\mathbf{x}$) such that $f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \mathbf{a}, \mathbf{y} - \mathbf{x} \rangle$ for all $\mathbf{y} \in B(\mathbf{x}, \epsilon)$. We will omit the proof of this extension of Theorem 3.15 here. See Chapter on "Convex Functions" in [3].

Affine supports for convex functions have been given a special name.

NOTES:                                                                            46

**Definition 3.17.** Let $f : \mathbb{R}^d \to \mathbb{R} \cup \{+\infty\}$ be a convex function. For any $\mathbf{x} \in \text{dom}(f)$, an affine support at $x$ is called a *subgradient* of $f$ at $\mathbf{x}$. The set of all subgradients at $\mathbf{x}$ is denoted by $\partial f(\mathbf{x})$ and is called the *subdifferential* of $f$ at $\mathbf{x}$.

**Theorem 3.18.** Let $f : \mathbb{R}^d \to \mathbb{R} \cup \{+\infty\}$ be a convex function. For any $\mathbf{x} \in \text{dom}(f)$, the subdifferential $\partial f(\mathbf{x})$ at $\mathbf{x}$ is a closed, convex set.

*Proof.* Note that

$$\partial f(\mathbf{x}) := \{\mathbf{a} \in \mathbb{R}^d : \langle \mathbf{y} - \mathbf{x}, \mathbf{a} \rangle \leq f(\mathbf{y}) - f(\mathbf{x}) \ \ \forall \mathbf{y} \in \mathbb{R}^d\}.$$

Since the right hand side of the above equation is the intersection of a family of halfspaces, this shows that $\partial f(\mathbf{x})$ is a closed, convex set. $\qquad\square$

## 3.2  Continuity properties

Convex functions enjoy strong continuity properties in the relative interior of their domains[3]. This fact is very useful in many contexts, especially in optimization, because this is useful in showing that minimizers and maximizers exist when optimizing convex functions that show up in practice, via Weierstrass' theorem (Theorem 1.11).

**Proposition 3.19.** Let $f : \mathbb{R}^d \to \mathbb{R} \cup \{+\infty\}$ be a convex function. Take $\mathbf{x}^\star \in \mathbb{R}^d$ and suppose that for some $\epsilon > 0$ and $m, M \in \mathbb{R}$, the inequalities

$$m \leq f(\mathbf{x}) \leq M$$

hold for all $\mathbf{x}$ in the ball $B(\mathbf{x}^\star, 2\epsilon)$. Then for all $\mathbf{x}, \mathbf{y} \in B(\mathbf{x}^\star, \epsilon)$, it holds that

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq \left(\frac{M - m}{\epsilon}\right) \|\mathbf{x} - \mathbf{y}\|. \tag{3.1}$$

In particular, $f$ is locally Lipschitz about $\mathbf{x}^\star$.

*Proof.* Take $\mathbf{x}, \mathbf{y} \in B(\mathbf{x}^\star, \epsilon)$. Define $\mathbf{z} = \mathbf{y} + \epsilon\left(\frac{\mathbf{y} - \mathbf{x}}{\|\mathbf{y} - \mathbf{x}\|}\right)$. Note that

$$\|\mathbf{z} - \mathbf{x}^\star\| = \left\|\mathbf{y} + \epsilon\left(\frac{\mathbf{y} - \mathbf{x}}{\|\mathbf{y} - \mathbf{x}\|}\right) - \mathbf{x}^\star\right\| \leq \|\mathbf{y} - \mathbf{x}^\star\| + \left\|\epsilon\left(\frac{\mathbf{y} - \mathbf{x}}{\|\mathbf{y} - \mathbf{x}\|}\right)\right\| \leq \epsilon + \epsilon = 2\epsilon.$$

Thus $\mathbf{z} \in B(\mathbf{x}^\star, 2\epsilon)$. Also,

$$\mathbf{y} = \left(\frac{\|\mathbf{y} - \mathbf{x}\|}{\epsilon + \|\mathbf{y} - \mathbf{x}\|}\right)\mathbf{z} + \left(1 - \frac{\|\mathbf{y} - \mathbf{x}\|}{\epsilon + \|\mathbf{y} - \mathbf{x}\|}\right)\mathbf{x},$$

---

[3]This section was written by Joseph Paat.

NOTES:                                                      47

showing that $\mathbf{y}$ is a convex combination of $\mathbf{x}$ and $\mathbf{z}$. Therefore we may apply the convexity of $f$ to see

$$
\begin{aligned}
f(\mathbf{y}) &\leq \left( \frac{\|\mathbf{y} - \mathbf{x}\|}{\epsilon + \|\mathbf{y} - \mathbf{x}\|} \right) f(\mathbf{z}) + \left( 1 - \frac{\|\mathbf{y} - \mathbf{x}\|}{\epsilon + \|\mathbf{y} - \mathbf{x}\|} \right) f(\mathbf{x}) \\
&= f(\mathbf{x}) + \left( \frac{\|\mathbf{y} - \mathbf{x}\|}{\epsilon + \|\mathbf{y} - \mathbf{x}\|} \right) (f(\mathbf{z}) - f(\mathbf{x})) \\
&\leq f(\mathbf{x}) + \left( \frac{\|\mathbf{y} - \mathbf{x}\|}{\epsilon} \right) (M - m) \qquad\qquad \text{using the bounds on } f \text{ in } B(\mathbf{x}^\star, 2\epsilon).
\end{aligned}
$$

Hence $f(\mathbf{y}) - f(\mathbf{x}) \leq \left( \frac{\|\mathbf{y} - \mathbf{x}\|}{\epsilon} \right) (M - m)$.

Repeating this argument by swapping the roles of $\mathbf{x}$ and $\mathbf{y}$, we get $f(\mathbf{x}) - f(\mathbf{y}) \leq \left( \frac{\|\mathbf{y} - \mathbf{x}\|}{\epsilon} \right) (M - m)$.
Therefore (3.2) holds. $\qquad\square$

**Proposition 3.20.** Let $f : \mathbb{R}^d \to \mathbb{R} \cup \{+\infty\}$ be a convex function. Consider any compact, convex subset $S \subseteq \mathrm{dom}(f)$ and let $\mathbf{x}^\star \in \mathrm{relint}(S)$. Then there is a $\epsilon_{\mathbf{x}^\star} > 0$ and values $m_{\mathbf{x}^\star}, M_{\mathbf{x}^\star} \in \mathbb{R}$ so that

$$
m_{\mathbf{x}^\star} \leq f(\mathbf{x}) \leq M_{\mathbf{x}^\star} \tag{3.2}
$$

for all $\mathbf{x} \in B(\mathbf{x}^\star, 2\epsilon_{\mathbf{x}^\star}) \cap S$.

*Proof.* Let $\mathbf{v}^1, \ldots, \mathbf{v}^\ell$ be vectors that span the linear space parallel to $\mathrm{aff}(S)$ (see Theorem 2.16). By definition of relative interior, since $\mathbf{x} \in \mathrm{aff}(S)$, there exists $\epsilon > 0$ such that $\mathbf{x}^\star + \epsilon \mathbf{v}^j$ and $\mathbf{x}^\star - \epsilon \mathbf{v}^j$ are both in $S$ for $j = 1, \ldots, \ell$. Denote the set of points $\mathbf{x}^\star \pm \epsilon \mathbf{v}^j$ as $\mathbf{x}_1, \ldots, \mathbf{x}_k \in S$ ($k = 2\ell$), and define $S' := \mathrm{conv}\{\mathbf{x}_1, \ldots, \mathbf{x}_k\}$. Observe that $\mathbf{x}^\star \in \mathrm{relint}(S')$ and $\mathrm{aff}(S') = \mathrm{aff}(S)$. Set $M_{\mathbf{x}^\star} = \max\{f(\mathbf{x}_i) : i = 1, \ldots, k\}$. Using Problem 3 from "HW for Week VII", it follows that $f(\mathbf{x}) \leq M_{\mathbf{x}^\star}$ for all $\mathbf{x} \in S'$.

Now since $f$ is convex, by Theorem 3.15, there is some affine support function $L(\mathbf{x}) = \langle \mathbf{a}, (\mathbf{x} - \mathbf{x}^\star) \rangle + f(\mathbf{x}^\star)$ for $f$ at $\mathbf{x}^\star$. Define $m_{\mathbf{x}^\star} = \min\{L(\mathbf{x}_i) : i = 1, \ldots, k\}$. Consider any point $\mathbf{x} = \sum_{i=1}^k \lambda_i \mathbf{x}_i \in S'$, where $\lambda_1, \ldots, \lambda_k$ are convex coefficients, and observe that

$$
L(\mathbf{x}) = \left\langle \mathbf{a}, \left( \sum_{i=1}^k \lambda_i \mathbf{x}_i \right) - \mathbf{x}^\star \right\rangle + f(\mathbf{x}^\star) = \sum_{i=1}^k \lambda_i \left( \langle \mathbf{a}, \mathbf{x}_i - \mathbf{x}^\star \rangle + f(\mathbf{x}^\star) \right) = \sum_{i=1}^{d+1} \lambda_i L(\mathbf{x}_i) \geq m_{\mathbf{x}^\star}.
$$

Since $L$ is an affine support, it follows that $f(\mathbf{x}) \geq L(\mathbf{x}) \geq m_{\mathbf{x}^\star}$ for all $\mathbf{x} \in S'$. Finally, as $\mathbf{x}^\star \in \mathrm{relint}(S')$ and $\mathrm{aff}(S') = \mathrm{aff}(S)$, there is some $\epsilon > 0$ so that $B(\mathbf{x}^\star, 2\epsilon) \cap S \subseteq S'$.
$\qquad\square$

**Theorem 3.21.** Let $f : \mathbb{R}^d \to \mathbb{R} \cup \{+\infty\}$ be a convex function. Let $D \subseteq \mathrm{relint}(\mathrm{dom}(f))$ be a convex, compact subset. Then there is a constant $L = L(D) \geq 0$ so that

$$
|f(\mathbf{x}) - f(\mathbf{y})| \leq L \|\mathbf{x} - \mathbf{y}\| \tag{3.3}
$$

for all $\mathbf{x}, \mathbf{y} \in D$. In particular, $f$ is locally Lipschitz continuous over the relative interior of its domain.

NOTES: 48

*Proof.* Let $S$ be a compact set such that $D \subseteq \mathrm{relint}(S) \subseteq \mathrm{relint}(\mathrm{dom}(f))$. From Proposition 3.20, for every $\mathbf{x} \in \mathrm{relint}(S)$, there is a tuple $(\epsilon_{\mathbf{x}}, m_{\mathbf{x}}, M_{\mathbf{x}})$ so that $m_{\mathbf{x}} \leq f(\mathbf{y}) \leq M_{\mathbf{x}}$ for all $\mathbf{y} \in B(\mathbf{x}, 2\epsilon_{\mathbf{x}}) \cap S$. Proposition 3.19 then implies that there is some $L_{\mathbf{x}} \geq 0$ so that $|f(\mathbf{y}) - f(\mathbf{z})| \leq L_{\mathbf{x}} \|\mathbf{z} - \mathbf{y}\|$ for all $\mathbf{z}, \mathbf{y} \in B(\mathbf{x}, \epsilon_{\mathbf{x}})$. Note that the collection $\{B(\mathbf{x}, \epsilon_{\mathbf{x}}) \cap S : \mathbf{x} \in D\}$ forms an open cover of $S$ (in the relative topology of $\mathrm{aff}(S)$). Therefore, as $S$ is compact, there exists a finite set $\{x_1, ..., x_k\} \subset S$ so that $S \subseteq \bigcup_{i=1}^{k} B(\mathbf{x}_i, \epsilon_{\mathbf{x}_i})$. Set $L = \max\{L_{\mathbf{x}_i} : i \in [k]\}$.

Now take $\mathbf{y}, \mathbf{z} \in S$. The line segment $[\mathbf{y}, \mathbf{z}]$ can be divided into finitely many segments $[\mathbf{y}, \mathbf{z}] = [\mathbf{y}_1, \mathbf{y}_2] \cup [\mathbf{y}_2, \mathbf{y}_3] \cup ... \cup [\mathbf{y}_{q-1}, \mathbf{y}_q]$, where $\mathbf{y}_1 = \mathbf{y}$, $\mathbf{y}_q = \mathbf{z}$, and each interval $[\mathbf{y}_i, \mathbf{y}_{i+1}]$ is contained in some ball $B(\mathbf{x}_j, \epsilon_{\mathbf{x}_j})$ for $j \in [k]$. Without loss of generality, we may assume that $q - 1 \leq k$ and $[\mathbf{y}_i, \mathbf{y}_{i+1}] \subseteq B(\mathbf{x}_i, \epsilon_{\mathbf{x}_i})$ for each $i \in [q-1]$. It follows that

$$
\begin{aligned}
|f(\mathbf{y}) - f(\mathbf{z})| &= \left| f(\mathbf{y}_1) + \left( \sum_{i=2}^{q-1} f(y_i) \right) - \left( \sum_{i=2}^{q-1} f(y_i) \right) - f(y_q) \right| \\
&= \left| \sum_{i=1}^{q-1} f(\mathbf{y}_i) - f(\mathbf{y}_{i+1}) \right| \\
&\leq \sum_{i=1}^{q-1} |f(\mathbf{y}_i) - f(\mathbf{y}_{i+1})| \\
&\leq \sum_{i=1}^{q-1} L_{\mathbf{x}_i} \|\mathbf{y}_i - \mathbf{y}_{i+1}\| \\
&\leq \sum_{i=1}^{q-1} L \|\mathbf{y}_i - \mathbf{y}_{i+1}\| \\
&= L \|\mathbf{y}_1 - \mathbf{y}_q\| = L \|\mathbf{y} - \mathbf{z}\|.
\end{aligned}
$$

Hence $f$ is Lipschitz over $S$ with constant $L$. $\qquad\square$

## 3.3   First-order derivative properties

A convex function enjoys very strong differentiability properties. We will first state some useful results without proof. See the Chapter on "Convex Functions" in Gruber [3] for full proofs.

**Theorem 3.22.** Let $f : \mathbb{R}^d \to \mathbb{R} \cup \{+\infty\}$ be a convex function and let $\mathbf{x} \in \mathrm{int}(\mathrm{dom}(f))$. Then $f$ is differentiable at $\mathbf{x}$ if and only if the partial derivative $f_i'(\mathbf{x})$ exists for all $i = 1, \ldots, d$.

**Theorem 3.23.** [Reidemeister's Theorem] Let $f : \mathbb{R}^d \to \mathbb{R} \cup \{+\infty\}$ be a convex function. Then $f$ is differentiable almost everywhere in $\mathrm{int}(\mathrm{dom}(f))$, i.e., the subset of $\mathrm{int}(\mathrm{dom}(f))$ where $f$ is not differentiable has Lebesgue measure 0.

NOTES: 49

We now prove the central relationships between the gradient $\nabla f$ and convexity. We first observe some facts about convex functions on the real line.

**Proposition 3.24.** Let $f : \mathbb{R} \to \mathbb{R}$ be a convex function. Then for any real numbers $x < y < z$, we must have

$$\frac{f(y) - f(x)}{y - x} \le \frac{f(z) - f(x)}{z - x} \le \frac{f(z) - f(y)}{z - y}.$$

Moreover, if $f$ is strictly convex, then these inequalities are strict.

*Proof.* Since $y \in (x, z)$, there exists $\alpha \in (0, 1)$ such that $y = \alpha x + (1 - \alpha)z$. Now we follow the inequalities:

$$\begin{array}{rl}
\frac{f(y)-f(x)}{y-x} & = \frac{f(\alpha x + (1-\alpha)z) - f(x)}{\alpha x + (1-\alpha)z - x} \\
& \le \frac{\alpha f(x) + (1-\alpha)f(z) - f(x)}{\alpha x + (1-\alpha)z - x} \\
& = \frac{f(z) - f(x)}{z - x}
\end{array}.$$

Similarly,

$$\begin{array}{rl}
\frac{f(z)-f(y)}{z-y} & = \frac{f(z) - f(\alpha x + (1-\alpha)z)}{z - \alpha x - (1-\alpha)z} \\
& \ge \frac{f(z) - \alpha f(x) - (1-\alpha)f(z)}{z - \alpha x - (1-\alpha)z} \\
& = \frac{f(z) - f(x)}{z - x}
\end{array}.$$

The strict convexity implication is clear from the above. $\qquad\square$

An immediate corollary is the following relationship between the derivative of a function on the real line and convexity.

**Proposition 3.25.** Let $f : \mathbb{R} \to \mathbb{R}$ be a differentiable function. Then $f$ is convex if and only if $f'$ is an increasing function, i.e., $f'(x) \le f'(y)$ for all $x \le y \in \mathbb{R}$. Moreover, $f$ is strictly convex if and only if $f'$ is strictly increasing. $f$ is strongly convex with strong convexity modulus $c > 0$ if and only if $f'(x) \ge f'(y) + c(x - y)$ for all $x \ge y \in \mathbb{R}$.

*Proof.* ($\Rightarrow$) Recall that $f'(x) = \lim_{t \to 0_+} \frac{f(x+t) - f(x)}{t}$. But for every $0 < t < y - x$, we have $\frac{f(x+t) - f(x)}{t} \le \frac{f(y) - f(x)}{y - x}$ by Proposition 3.24. Thus, $f'(x) \le \frac{f(y) - f(x)}{y - x}$. By a similar argument, we obtain $f'(y) \ge \frac{f(y) - f(x)}{y - x}$. This gives the relation.

($\Leftarrow$) Consider any $x, z \in R$ and $\alpha \in (0, 1)$. Let $y = \alpha x + (1 - \alpha)z$. By the mean value theorem, there exists $t_1 \in [x, y]$ such that $\frac{f(y) - f(x)}{y - x} = f'(t_1)$ and $t_2 \in [y, z]$ such that $\frac{f(z) - f(y)}{z - y} = f'(t_2)$. Since $t_2 \ge t_1$ and we assume $f'$ is increasing, then $f'(t_2) \ge f'(t_1)$. This implies that

$$\frac{f(z) - f(y)}{z - y} \ge \frac{f(y) - f(x)}{y - x}.$$

Substituting $y = \alpha x + (1 - \alpha)z$ and rearranging, we obtain that $f(\alpha x + (1 - \alpha)z) \le \alpha f(x) + (1 - \alpha)f(z)$.

The argument for strict convexity follows by replacing all inequalities by strict inequalities. $\qquad\square$

NOTES: 50

We can now prove the main result of this subsection. A key idea behind the results below is that one can reduce testing convexity of a function on $\mathbb{R}^d$ to testing convexity of any one-dimensional "slice" of it. More precisely,

**Proposition 3.26.** Let $f : \mathbb{R}^d \to \mathbb{R}$ be a function. Then $f$ is convex if and only if for every $\mathbf{x}, \mathbf{r} \in \mathbb{R}^d$, the function $\phi : \mathbb{R} \to \mathbb{R}$ defined by $\phi(t) = f(\mathbf{x} + t\mathbf{r})$ is convex.

*Proof.* Left as an exercise. □

**Theorem 3.27.** Let $f : \mathbb{R}^d \to \mathbb{R}$ be differentiable everywhere. Then the following are all equivalent.

1. $f$ is convex.

2. $f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$ for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$.

3. $\langle \nabla f(\mathbf{y}) - \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \geq 0$ for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$.

A characterization of strict convexity is obtained if all the above inequalities are considered strict for all $\mathbf{x} \neq \mathbf{y} \in \mathbb{R}^d$. A characterization of strong convexity with modulus $c > 0$ is obtained if 2. is replaced with $f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{1}{2}c\|\mathbf{y} - \mathbf{x}\|^2$ for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, and 3. is replaced with $\langle \nabla f(\mathbf{y}) - \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \geq c\|\mathbf{y} - \mathbf{x}\|^2$ for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$.

*Proof.* 1. $\Rightarrow$ 2. Consider any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$. For every $\alpha > 0$, convexity of $f$ implies that $f((1-\alpha)\mathbf{x} + \alpha\mathbf{y}) \leq (1-\alpha)f(\mathbf{x}) + \alpha f(\mathbf{y})$. Rearranging, we obtain

$$\frac{f((1-\alpha)\mathbf{x}+\alpha\mathbf{y})-f(\mathbf{x})}{\alpha} \leq f(\mathbf{y}) - f(\mathbf{x})$$
$$\Rightarrow \quad \frac{f(\mathbf{x}+\alpha(\mathbf{y}-\mathbf{x}))-f(\mathbf{x})}{\alpha} \leq f(\mathbf{y}) - f(\mathbf{x})$$

Letting $\alpha \to 0$ on the left hand side, we obtain the directional derivative $\langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$ and 2. is established.

2. $\Rightarrow$ 3. By switching the roles of $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, we obtain the following

$$\begin{array}{l} f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \\ f(\mathbf{x}) \geq f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \end{array} \ .$$

Adding these inequalities together we obtain 3.

3. $\Rightarrow$ 1. Consider any $\bar{\mathbf{x}}, \bar{\mathbf{y}} \in \mathbb{R}^d$ and define the function $\phi(t) := f(\bar{\mathbf{x}} + t(\bar{\mathbf{y}} - \bar{\mathbf{x}}))$. Observe that $\phi'(t) = \langle \nabla f(\bar{\mathbf{x}} + t(\bar{\mathbf{y}} - \bar{\mathbf{x}})), \bar{\mathbf{y}} - \bar{\mathbf{x}} \rangle$ for any $t \in \mathbb{R}$. For $t_2 > t_1$, we have that

$$\begin{aligned} \phi'(t_2) - \phi'(t_1) &= \langle \nabla f(\bar{\mathbf{x}} + t_2(\bar{\mathbf{y}} - \bar{\mathbf{x}})), \bar{\mathbf{y}} - \bar{\mathbf{x}} \rangle - \langle \nabla f(\bar{\mathbf{x}} + t_1(\bar{\mathbf{y}} - \bar{\mathbf{x}})), \bar{\mathbf{y}} - \bar{\mathbf{x}} \rangle \\ &= \langle \nabla f(\bar{\mathbf{x}} + t_2(\bar{\mathbf{y}} - \bar{\mathbf{x}})) - \nabla f(\bar{\mathbf{x}} + t_1(\bar{\mathbf{y}} - \bar{\mathbf{x}})), \bar{\mathbf{y}} - \bar{\mathbf{x}} \rangle \\ &= \frac{1}{t_2 - t_1} \langle \nabla f(\bar{\mathbf{x}} + t_2(\bar{\mathbf{y}} - \bar{\mathbf{x}})) - \nabla f(\bar{\mathbf{x}} + t_1(\bar{\mathbf{y}} - \bar{\mathbf{x}})), (t_2 - t_1)(\bar{\mathbf{y}} - \bar{\mathbf{x}}) \rangle \\ &= \frac{1}{t_2 - t_1} \langle \nabla f(\bar{\mathbf{x}} + t_2(\bar{\mathbf{y}} - \bar{\mathbf{x}})) - \nabla f(\bar{\mathbf{x}} + t_1(\bar{\mathbf{y}} - \bar{\mathbf{x}})), (t_2(\bar{\mathbf{y}} - \bar{\mathbf{x}}) - \bar{\mathbf{x}}) - (t_1(\bar{\mathbf{y}} - \bar{\mathbf{x}}) - \bar{\mathbf{x}}) \rangle \\ &\geq 0 \end{aligned}$$

where the last inequality follows from the fact that $\langle \nabla f(\mathbf{y}) - \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \geq 0$ for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, and $t_2 > t_1$. Therefore, by Proposition 3.25, we obtain that $\phi(t)$ is a convex function in $t$. By Proposition 3.26, $f$ is convex. □

NOTES: 51

## 3.4 Second-order derivative properties

A simple consequence of Proposition 3.25 for twice differentiable functions on the real line is the following.

**Corollary 3.28.** Let $f : \mathbb{R} \to \mathbb{R}$ be a twice differentiable function. Then $f$ is convex if and only if $f''(x) \geq 0$ for all $x \in \mathbb{R}$. If $f''(x) > 0$, then $f$ is strictly convex.

**Remark 3.29.** From Proposition 3.25, we know strict convexity of $f$ is equivalent to the condition that $f'$ is strictly increasing. However, this is not equivalent to $f''(x) > 0$, the implication only goes in one direction. This is why we lose the other direction when discussing strict convexity in Corollary 3.28. As a concrete example, consider $f(x) = x^4$ which is strictly convex, but the second derivative is 0 at $x = 0$.

This enables one to characterize convexity of $f : \mathbb{R}^d \to \mathbb{R}$ in terms of its Hessian, which will be denoted by $\nabla^2 f$.

**Theorem 3.30.** Let $f : \mathbb{R}^d \to \mathbb{R}$ be a twice differentiable function. Then the following are all true.

1. $f$ is convex if and only if $\nabla^2 f(\mathbf{x})$ is positive semidefinite (PSD) for all $\mathbf{x} \in \mathbb{R}^d$.

2. If $\nabla^2 f(\mathbf{x})$ is positive definite (PD) for all $\mathbf{x} \in \mathbb{R}^d$, then $f$ is strictly convex.

3. $f$ is strongly convex with modulus $c > 0$ if and only if $\nabla^2 f(\mathbf{x}) - cI$ is positive semidefinite (PSD) for all $\mathbf{x} \in \mathbb{R}^d$.

*Proof.* 1. ($\Rightarrow$) Let $\mathbf{x} \in \mathbb{R}^d$ and we would like to show that $\nabla^2 f(\mathbf{x})$ is positive semidefinite. Consider any $\mathbf{r} \in \mathbb{R}^d$. Define the function $\phi(t) = f(\mathbf{x} + t\mathbf{r})$. By Proposition 3.26, $\phi$ is convex. By Corollary 3.28, $0 \leq \phi''(0) = \langle \nabla^2 f(\mathbf{x})\mathbf{r}, \mathbf{r} \rangle$. Since the choice of $\mathbf{r}$ was arbitrary, this shows that $\nabla^2 f(\mathbf{x})$ is positive semidefinite.

($\Leftarrow$) Assume $\nabla^2 f(\mathbf{x})$ is positive semidefinite fo all $\mathbf{x} \in \mathbb{R}^d$, and consider $\bar{\mathbf{x}}, \mathbf{r} \in \mathbb{R}^d$. Define the function $\phi(t) = f(\bar{\mathbf{x}} + t\mathbf{r})$. Now $\phi''(t) = \langle \nabla^2 f(\bar{\mathbf{x}} + t\mathbf{r})\mathbf{r}, \mathbf{r} \rangle \geq 0$, since $\nabla^2 f(\bar{\mathbf{x}} + t\mathbf{r})$ is positive semidefinite. By Corollary 3.28, $\phi$ is convex. By Proposition 3.26, $f$ is convex.

2. This follows from the same construction as in 1. above, and the sufficient condition that if the second derivative of one-dimensional function is strictly positive, then the function is strictly convex.

3. We omit the proof of the characterization of strong convexity. $\qquad \square$

## 3.5 Sublinear functions, support functions and gauges

We will now introduce a more structured subfamily of convex functions which is easier to deal with analytically, and yet has very important uses in diverse areas.

**Definition 3.31.** A function $f : \mathbb{R}^d \to \mathbb{R} \cup \{+\infty\}$ is called *sublinear* if it satisfies the following two properties:

(i) $f$ is *positively homogeneous*, i.e., $f(\lambda\mathbf{r}) = \lambda f(\mathbf{r})$ for all $\mathbf{r} \in \mathbb{R}^d$ and $\lambda > 0$.

(ii) $f$ is *subadditive*, i.e., $f(\mathbf{x} + \mathbf{y}) \leq f(\mathbf{x}) + f(\mathbf{y})$ for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$.

Here is the connection with convexity.

**Proposition 3.32.** Let $f : \mathbb{R}^d \to \mathbb{R} \cup \{+\infty\}$. Then the following are equivalent:

1. $f$ is sublinear.

2. $f$ is convex and positively homogeneous.

3. $f(\lambda_1\mathbf{x}^1 + \lambda_2\mathbf{x}^2) \leq \lambda_1 f(\mathbf{x}^1) + \lambda_2 f(\mathbf{x}^2)$ for all $\mathbf{x}^1, \mathbf{x}^2 \in \mathbb{R}^d$ and $\lambda_1, \lambda_2 > 0$.

*Proof.* Left as an exercise. $\qquad\square$

A characterization via epigraphs is also possible.

**Proposition 3.33.** Let $f : \mathbb{R}^d \to \mathbb{R} \cup \{+\infty\}$ such that $f(\mathbf{0}) = 0$. Then $f$ is sublinear if and only if $\mathrm{epi}(f)$ is a convex cone in $\mathbb{R}^d \times \mathbb{R}$.

*Proof.* ($\Rightarrow$) From Proposition 3.32, we know that $f$ is convex and positively homogeneous. From Proposition 3.6, this implies that $\mathrm{epi}(f)$ is convex. So we only need to verify that if $(\mathbf{x}, t) \in \mathrm{epi}(f)$ then $\lambda(\mathbf{x}, t) = (\lambda\mathbf{x}, \lambda t) \in \mathrm{epi}(f)$ for all $\lambda \geq 0$. If $\lambda = 0$, then the result follows from the assumption that $f(\mathbf{0}) = 0$. Now consider $\lambda > 0$. Since $(\mathbf{x}, t) \in \mathrm{epi}(f)$, we have $f(\mathbf{x}) \leq t$ and by positive homogeneity of $f$, $f(\lambda\mathbf{x}) = \lambda f(\mathbf{x}) \leq \lambda t$, and so $(\lambda\mathbf{x}, \lambda t) \in \mathrm{epi}(f)$.

($\Leftarrow$) From Proposition 3.6 and the assumption that $\mathrm{epi}(f)$ is a convex cone, we get that $f$ is convex. We now verify that $f$ is positively homogeneous; by Proposition 3.32, we will be done. We first verify that for all $\lambda > 0$ and $\mathbf{x} \in \mathbb{R}^d$, $f(\lambda\mathbf{x}) \leq \lambda f(\mathbf{x})$. Since $\mathrm{epi}(f)$ is a convex cone and $(\mathbf{x}, f(\mathbf{x})) \in \mathrm{epi}(f)$, we have that $\lambda(\mathbf{x}, f(\mathbf{x})) = (\lambda\mathbf{x}, \lambda f(\mathbf{x})) \in \mathrm{epi}(f)$. This implies that $f(\lambda\mathbf{x}) \leq \lambda f(\mathbf{x})$.

Now, for any particular $\bar{\lambda} > 0$ and $\bar{\mathbf{x}} \in \mathbb{R}^d$, we have that $f(\bar{\lambda}\bar{\mathbf{x}}) \leq \bar{\lambda} f(\bar{\mathbf{x}})$. But using the above observation with $\lambda = \frac{1}{\bar{\lambda}}$ and $\mathbf{x} = \bar{\lambda}\mathbf{x}$, we obtain that $f(\frac{1}{\bar{\lambda}}\bar{\lambda}\bar{\mathbf{x}}) \leq \frac{1}{\bar{\lambda}} f(\bar{\lambda}\bar{\mathbf{x}})$, i.e., $\bar{\lambda} f(\bar{\mathbf{x}}) \leq f(\bar{\lambda}\bar{\mathbf{x}})$. Hence, we must have $f(\bar{\lambda}\bar{\mathbf{x}}) = \bar{\lambda} f(\bar{\mathbf{x}})$. $\qquad\square$

**Gauges.** One easily observes that any norm $N : \mathbb{R}^d \to \mathbb{R}$ is a sublinear function – recall Definition 1.1. In fact, a norm has the additional "symmetry" property that $N(\mathbf{x}) = N(-\mathbf{x})$. Since a sublinear function is convex (Proposition 3.32), and sublevel sets of convex sets are convex, we immediately know that the unit norm balls $B_N(\mathbf{0}, 1) = \{\mathbf{x} \in \mathbb{R}^d : N(\mathbf{x}) \leq 1\}$ are convex sets. Because of the "symmetry property" of norms, these unit norm balls are also "symmetric" about the origin. This merits a definition.

**Definition 3.34.** A convex set $C \subseteq \mathbb{R}^d$ is said to be *centrally symmetric about the origin*, if $\mathbf{x} \in C$ implies that $-\mathbf{x} \in C$. Sometimes we will abbreviate this to say $C$ is centrally symmetric.

NOTES: 53

1179    We now summarize the above discussion in the following observation.

**Proposition 3.35.** Let $N : \mathbb{R}^d \to \mathbb{R}$ be a norm. Then the unit norm ball $B_N(\mathbf{0}, 1) = \{\mathbf{x} \in \mathbb{R}^d : N(\mathbf{x}) \leq 1\}$ is a centrally symmetric, closed convex set.

One can actually prove a converse to the above statement, which will establish a nice one-to-one correspondence between norms and centrally symmetric convex sets. We first generalize the notion of a norm to a family of sublinear functions called "gauge functions".

**Definition 3.36.** Let $C \subseteq \mathbb{R}^d$ be a closed, convex set such that $\mathbf{0} \in C$. Define the following function $\gamma_C : \mathbb{R}^d \to \mathbb{R} \cup \{+\infty\}$ as

$$\gamma_C(\mathbf{r}) = \inf\{\lambda > 0 : \mathbf{r} \in \lambda C\}.$$

$\gamma_C$ is called the *gauge* or the *Minkowski functional* of $C$.

**Exercise 7.** Show that $\gamma_C$ is finite valued everywhere if and only if $\mathbf{0} \in \text{int}(C)$.

The following is a useful observation for the analysis of gauge functions.

**Lemma 3.37.** Let $C \subseteq \mathbb{R}^d$ be a closed convex set such that $\mathbf{0} \in C$, and let $\mathbf{r} \in \mathbb{R}^d$ be any vector. Then the set $\{\lambda > 0 : \mathbf{r} \in \lambda C\}$ is either empty or a convex interval of the real line of the form $(a, +\infty)$ or $[a, +\infty)$.

*Proof.* Define $I := \{\lambda > 0 : \mathbf{r} \in \lambda C\}$ and suppose it is nonempty. It suffices to show that if $\bar{\lambda} \in I$ then for all $\lambda \geq \bar{\lambda}$, $\lambda \in I$. This follows from the fact that $\bar{\lambda} \in I$ implies that $\frac{1}{\bar{\lambda}}\mathbf{r} \in C$. For any $\lambda \geq \bar{\lambda}$, we have $\frac{1}{\lambda}\mathbf{r} = \frac{\bar{\lambda}}{\lambda}(\frac{1}{\bar{\lambda}}\mathbf{r}) + (\frac{\lambda - \bar{\lambda}}{\lambda})\mathbf{0}$ which is in $C$ because $C$ is convex and $\mathbf{0} \in C$. $\qquad\square$

A useful intuition to keep in mind is that for any $\mathbf{r}$ the gauge function value $\gamma_C(\mathbf{r})$ gives you a factor to scale $\mathbf{r}$ with so that you end up on the boundary of $C$. More precisely,

**Proposition 3.38.** Let $C \subseteq \mathbb{R}^d$ be a closed, convex set such that $\mathbf{0} \in C$. Suppose $\mathbf{r} \in \mathbb{R}^d$ such that $0 < \gamma_C(\mathbf{r}) < \infty$. Then $\frac{1}{\gamma_C(\mathbf{r})}\mathbf{r} \in \text{relbd}(C)$.

*Proof.* From Lemma 3.37, we have that for all $\lambda > \gamma_C(\mathbf{r})$, we have that $\mathbf{r} \in \lambda C$, i.e., $\frac{1}{\lambda}\mathbf{r} \in C$. Taking the limit $\lambda \downarrow \gamma_C(\mathbf{r})$ and using the fact that $C$ is closed, we obtain that $\frac{1}{\gamma_C(\mathbf{r})}\mathbf{r} \in C$. If $\frac{1}{\gamma_C(\mathbf{r})}\mathbf{r} \in \text{relint}(C)$, then we can scale $\frac{1}{\gamma_C(\mathbf{r})}\mathbf{r}$ by $\alpha > 1$ and obtain that $\frac{\alpha}{\gamma_C(\mathbf{r})}\mathbf{r} \in C$, which would imply that $\mathbf{r} \in \frac{\gamma_C(\mathbf{r})}{\alpha}C\}$, contradicting the fact that $\gamma_C(\mathbf{r}) = \inf\{\lambda > 0 : \mathbf{r} \in \lambda C\}$, since $\frac{\gamma_C(\mathbf{r})}{\alpha} < \gamma_C(\mathbf{r})$. $\qquad\square$

The following theorem relates geometric properties of $C$ with analytical properties of the gauge function. These relations are extremely handy to keep in mind.

**Theorem 3.39.** Let $C \subseteq \mathbb{R}^d$ be a closed, convex set such that $\mathbf{0} \in C$. Then the following are all true.

1. $\gamma_C$ is a nonnegative, sublinear function.

NOTES:

2. $C = \{\mathbf{x} \in \mathbb{R}^d : \gamma_C(\mathbf{x}) \leq 1\}$.

3. $\mathrm{rec}(C) = \{\mathbf{r} \in \mathbb{R}^d : \gamma_C(\mathbf{r}) = 0\}$.

4. If $\mathbf{0} \in \mathrm{relint}(C)$, then $\mathrm{relint}(C) = \{\mathbf{x} \in \mathbb{R}^d : \gamma_C(\mathbf{x}) < 1\}$.

*Proof.*    1. Although 1. can be proved directly from the definition of the gauge, we postpone its proof until we speak of *support functions* below.

2. We now first show that $C \subseteq \{\mathbf{x} \in \mathbb{R}^d : \gamma_C(\mathbf{x}) \leq 1\}$. This is because $\mathbf{x} \in C$ implies that $1 \in \{\lambda > 0 : \mathbf{x} \in \lambda C\}$ and therefore, $\inf\{\lambda > 0 : \mathbf{x} \in \lambda C\} \leq 1$.

Now, we verify that $\{\mathbf{x} \in \mathbb{R}^d : \gamma_C(\mathbf{x}) \leq 1\} \subseteq C$. $\gamma_C(\mathbf{x}) \leq 1$ implies that $\inf\{\lambda > 0 : \mathbf{x} \in \lambda C\} \leq 1$ and since $\{\lambda > 0 : \mathbf{x} \in \lambda C\}$ is convex by Lemma 3.37, this means that either $1 \in \{\lambda > 0 : \mathbf{r} \in \lambda C\}$, and thus $\mathbf{x} \in C$ or $1 = \inf\{\lambda > 0 : \mathbf{x} \in \lambda C\} = \gamma_C(\mathbf{x})$. By Proposition 3.38, we have that $1 \cdot \mathbf{x} \in C$.

3. Since $\{\lambda > 0 : \mathbf{r} \in \lambda C\}$ is convex, as proved in part 2., we observe that $\gamma_C(\mathbf{r}) = 0$ if and only if $\frac{1}{\lambda}\mathbf{r} \in C\}$ for all $\lambda > 0$. Since $\mathbf{0} \in C$, this is equivalent to saying that $t\mathbf{r} \in C$ for all $t \geq 0$; more explicitly, $\mathbf{0} + t\mathbf{r} \in C$ for all $t \geq 0$. This is equivalent to saying that $\mathbf{r}$ satisfies Definition 2.43 of $\mathrm{rec}(C)$.

4. Consider any $\mathbf{x} \in \mathrm{relint}(C)$. By definition of relative interior, there exists $\lambda > 1$ such that $\lambda\mathbf{x} \in C$. By part 2. above, $\gamma_C(\lambda\mathbf{x}) \leq 1$ and by part 1. above, $\gamma_C$ is positively homogeneous, and thus, $\gamma_C(\mathbf{x}) \leq \frac{1}{\lambda} < 1$.

Now suppose $\mathbf{x} \in \mathbb{R}^d$ such that $\gamma_C(\mathbf{x}) < 1$. If $\gamma_C(\mathbf{x}) = 0$, then $\mathbf{x} \in \mathrm{rec}(C)$ by part 3. above. Since $\mathbf{0} \in \mathrm{relint}(C)$, we also have $\mathbf{x} = \mathbf{0} + \mathbf{x} \in \mathrm{relint}(C)$. Now suppose, $0 < \gamma_C(\mathbf{x}) < 1$. By part 2. above, $\mathbf{x} \in C$. Suppose to the contrary that $\mathbf{x} \notin \mathrm{relint}(C)$. By Theorem 2.40, $\mathbf{x}$ is contained in a proper face $F$ of $C$. Since $\mathbf{0} \in \mathrm{relint}(C)$, $\mathbf{0}$ is not contained in $F$. Also, $\gamma_C(\frac{\mathbf{x}}{\gamma_C(\mathbf{x})}) = 1$ by positive homogeneity of $\gamma_C$, from part 1. above. Therefore, $\frac{\mathbf{x}}{\gamma_C(\mathbf{x})} \in C$. However, $\mathbf{x} = (1 - \gamma_C(\mathbf{x}))\mathbf{0} + \gamma_C(\mathbf{x})(\frac{\mathbf{x}}{\gamma_C(\mathbf{x})})$. Since $\gamma_C(\mathbf{x}) < 1$ and $\mathbf{0} \notin F$, this would contradict the fact that $F$ is a face. $\qquad\square$

We derive some immediate consequences.

**Corollary 3.40.** Let $C \subseteq \mathbb{R}^d$ be a closed, convex set containing the origin. Then $C$ is compact if and only if $\gamma(\mathbf{r}) > 0$ for all $\mathbf{r} \in \mathbb{R}^d \setminus \{\mathbf{0}\}$.

**Corollary 3.41.** [Uniqueness of the gauge] Let $C$ be a compact convex set containing the origin in its interior, i.e., $\mathbf{0} \in \mathrm{int}(C)$. Let $f : \mathbb{R}^d \to \mathbb{R}$ be any sublinear function. Then $C = \{\mathbf{x} \in \mathbb{R}^d : f(\mathbf{x}) \leq 1\}$ if and only if $f = \gamma_C$.

*Proof.* The sufficiency follows from Theorem 3.39, part 2. For the necessity, suppose to the contrary that $f(\mathbf{x}) \neq \gamma_C(\mathbf{x})$ for some $\mathbf{x} \in \mathbb{R}^d$. We first observe that $\mathbf{x} \neq \mathbf{0}$ because $f(\mathbf{0}) = 0 = \gamma_C(\mathbf{0})$ by positive homogeneity and the fact that $f$ is continuous (Theorem 3.21) because $f$ is convex (Proposition 3.32).

NOTES:                                    55

First suppose $f(\mathbf{x}) > \gamma_C(\mathbf{x})$. Since $C$ is compact, we know that $\gamma_C(\mathbf{x}) > 0$. Consider that point $\frac{1}{\gamma_C(\mathbf{x})}\mathbf{x}$. By Proposition 3.38, $\mathbf{x} \in \mathrm{relbd}(C)$. However, since $f$ is positively homogeneous, $f(\frac{1}{\gamma_C(\mathbf{x})}\mathbf{x}) = \frac{1}{\gamma_C(\mathbf{x})}f(\mathbf{x}) > 1$ because $f(\mathbf{x}) > \gamma_C(\mathbf{x})$. This contradicts that $C = \{\mathbf{x} \in \mathbb{R}^d : f(\mathbf{x}) \leq 1\}$.

Next suppose $f(\mathbf{x}) < \gamma_C(\mathbf{x})$. If $f(\mathbf{x}) \leq 0$, then by positive homogeneity, $f(\lambda\mathbf{x}) \leq 0$ for all $\lambda \geq 0$. Thus, $\lambda\mathbf{x} \in C$ for all $\lambda \geq 0$ by the assumption that $C = \{\mathbf{x} \in \mathbb{R}^d : f(\mathbf{x}) \leq 1\}$. This means that $\mathbf{x} \in \mathrm{rec}(C)$ which contradicts the fact that $C$ is compact (see Theorem 2.47). Thus, we may assume that $f(\mathbf{x}) > 0$.

Now let $\mathbf{y} = \frac{1}{f(\mathbf{x})}\mathbf{x}$. By positive homogeneity of $\gamma_C$, we obtain that $\gamma_C(\mathbf{y}) = \gamma_C(\frac{1}{f(\mathbf{x})}\mathbf{x}) = \frac{\gamma_C(\mathbf{x})}{f(\mathbf{x})} > 1$. Therefore, $\mathbf{y} \notin C$ by Theorem 3.39, part 2. However, $f(\mathbf{y}) = 1$, which contradicts the assumption that $C = \{\mathbf{x} \in \mathbb{R}^d : f(\mathbf{x}) \leq 1\}$. $\qquad\square$

The proof of Corollary 3.41 also implies the following.

**Corollary 3.42.** [Uniqueness of the gauge-II] Let $C$ be a closed, convex set (not necessarily compact) containing the origin in its interior, i.e., $\mathbf{0} \in \mathrm{int}(C)$. Let $f : \mathbb{R}^d \to \mathbb{R}$ be any *nonnegative*, sublinear function. Then $C = \{\mathbf{x} \in \mathbb{R}^d : f(\mathbf{x}) \leq 1\}$ if and only if $f = \gamma_C$.

Consequently, for every nonnegative, sublinear function $f$, there exists a closed, convex set $C$ such that $f = \gamma_C$.

We also make the following observation on when the gauge function can take $+\infty$ as a value.

**Lemma 3.43.** Let $C$ be a closed, convex set with $\mathbf{0} \in C$. Then the gauge $\gamma_C$ is finite valued everywhere (i.e., $\gamma_C(\mathbf{x}) < \infty$ for all $\mathbf{x} \in \mathbb{R}^d$) if and only if $\mathbf{0} \in \mathrm{int}(C)$.

*Proof.* ($\Longrightarrow$) Suppose $\mathbf{0}$ is not in the interior, i.e., $\mathbf{0}$ is on the boundary of $C$. By the Supporting Hyperplane Theorem 2.23, there exist $\mathbf{a} \in \mathbb{R}^d \setminus \{\mathbf{0}\}$ and $\delta \in \mathbb{R}$ such that $C \subseteq H^-(\mathbf{a}, \delta)$ and $\langle \mathbf{a}, \mathbf{0}\rangle = \delta$. Thus, $\delta = 0$. Now consider any $\mathbf{r} \in \mathbb{R}^d$ such that $\langle \mathbf{a}, \mathbf{r}\rangle > 0$. However, since $C \subseteq H^-(\mathbf{a}, 0)$, it follows that $\lambda C \subseteq H^-(\mathbf{a}, 0)$ for all $\lambda > 0$. Therefore, the set $\{\lambda > 0 : \mathbf{r} \in \lambda C\}$ is empty, and we conclude that $\gamma_C(\mathbf{r}) = \infty$. In fact, this shows that $\gamma_C$ takes value $\infty$ on the entire "open" halfspace $\{\mathbf{r} \in \mathbb{R}^d : \langle \mathbf{a}, \mathbf{r}\rangle > 0\}$.

($\Longleftarrow$) Assume $\mathbf{0} \in \mathrm{int}(C)$ and consider any $\mathbf{x} \in \mathbb{R}^d$. Since $\mathbf{0} \in \mathrm{int}(C)$, there exists $\epsilon > 0$ such that $\epsilon\mathbf{x} \in C$. Thus, $\frac{1}{\epsilon}$ is in the set $\{\lambda > 0 : \mathbf{x} \in \lambda C\}$, and so the infimum over this set is finite valued. Thus, $\gamma_C(\mathbf{x}) < \infty$ for all $\mathbf{x} \in \mathbb{R}^d$. $\qquad\square$

We can now finally settle the correspondence between norms and centrally symmetric, compact convex sets.

**Theorem 3.44.** Let $N : \mathbb{R}^d \to \mathbb{R}$ be a norm. Then $B_N(\mathbf{0}, 1) = \{\mathbf{x} \in \mathbb{R}^d : N(\mathbf{x}) \leq 1\}$ is a centrally symmetric, compact convex set with $\mathbf{0}$ in its interior. Moreover, $\gamma_{B_N(\mathbf{0},1)} = N$.

Conversely, let $B$ be a centrally symmetric, compact convex set containing $\mathbf{0}$ in its interior. Then $\gamma_B$ is a norm on $\mathbb{R}^d$ and $B = B_{\gamma_B}(\mathbf{0}, 1)$.

*Proof.* For the first part, since $N$ is sublinear, it is convex (by Proposition 3.32). By definition, $B_N(\mathbf{0}, 1) = \{\mathbf{x} \in \mathbb{R}^d : N(\mathbf{x}) \leq 1\}$ is a sublevel set for $N$, and is thus a convex set. It is closed, since $N$ is continuous by Theorem 3.21. Since $N(\mathbf{x}) = N(-\mathbf{x})$, this also shows that $B_N(\mathbf{0}, 1)$ is centrally symmetric. We now show that $\mathrm{rec}(B_N(\mathbf{0}, 1)) = \{\mathbf{0}\}$; this will imply that it is compact by Theorem 2.47. Consider any nonzero vector $\mathbf{r}$, and let $N(\mathbf{r}) = M > 0$. Then, $\frac{2}{M}\mathbf{r} = \mathbf{0} + \frac{2}{M}\mathbf{r}$, but $N(\frac{2}{M}\mathbf{r}) = 2$. Thus, $\frac{2}{M}\mathbf{r} \notin B_N(\mathbf{0}, 1)$, and so $\mathbf{r}$ cannot be a recession direction for $B_N(\mathbf{0}, 1)$.

We verify that $\mathbf{0} \in \mathrm{int}(B_N(\mathbf{0}, 1))$. If not, then by the Supporting Hyperplane Theorem 2.23, there exists $\mathbf{a} \in \mathbb{R}^d \setminus \{\mathbf{0}\}$ and $\delta \in \mathbb{R}$ such that $B_N(\mathbf{0}, 1) \subseteq H^-(\mathbf{a}, \delta)$ and $\langle \mathbf{a}, \mathbf{0} \rangle = \delta$. Thus, $\delta = 0$. Now, since $\mathbf{a} \neq 0$, $N(\mathbf{a}) > 0$. Thus, $N(\frac{\mathbf{a}}{N(\mathbf{a})}) = 1$ and by definition, $\frac{\mathbf{a}}{N(\mathbf{a})} \in B_N(\mathbf{0}, 1)$. However, $\langle \mathbf{a}, \frac{\mathbf{a}}{N(\mathbf{a})} \rangle = \frac{\|\mathbf{a}\|^2}{N(\mathbf{a})} > 0$ which contradicts the fact that $B_N(\mathbf{0}, 1) \subseteq H^-(\mathbf{a}, 0)$. Therefore, from Corollary 3.41, we obtain that $N = \gamma_{B_N(\mathbf{0}, 1)}$.

For the second part, we know that $\gamma_B$ is sublinear, and since $B$ is compact, $\gamma_B(\mathbf{r}) > 0$ for all $\mathbf{r} \neq 0$ by Corollary 3.40. Since $0 \in \mathrm{int}(B)$, Lemma 3.43 implies that $\gamma_C$ is finite valued everywhere. To confirm that $\gamma_B$ is a norm, all that remains to be checked is that $\gamma_B(\mathbf{x}) = \gamma_B(-\mathbf{x})$ for all $\mathbf{x} \neq \mathbf{0}$. Suppose to the contrary that $\gamma_B(\mathbf{x}) > \gamma_B(-\mathbf{x})$ (note that this is without loss of generality). This implies that $\gamma_B(\frac{1}{\gamma_B(-\mathbf{x})}\mathbf{x}) > 1$. Therefore, $\frac{1}{\gamma_B(-\mathbf{x})}\mathbf{x} \notin B$ by Theorem 3.39, part 2. However, $\gamma_B(-\frac{1}{\gamma_B(-\mathbf{x})}\mathbf{x}) = \frac{1}{\gamma_B(-\mathbf{x})}\gamma_B(-\mathbf{x}) = 1$ showing that $-\frac{1}{\gamma_B(-\mathbf{x})}\mathbf{x} \in B$ by Theorem 3.39, part 2. This contradicts the fact that $B$ is centrally symmetric. Thus, $\gamma_B$ is a norm on $\mathbb{R}^d$. Moreover, by Theorem 3.39, part 2., $B = \{\mathbf{x} \in \mathbb{R}^d : \gamma_B(\mathbf{x}) \leq 1\} = B_{\gamma_B}(\mathbf{0}, 1)$. $\square$

Let us build towards a more computational approach to the gauge. First, lets give an explicit formula for the gauge of a halfspace containing the origin.

**Example 3.45.** Let $H := H^-(\mathbf{a}, \delta)$ be a halfspace defined by some $\mathbf{a} \in \mathbb{R}^d$ and $\delta \in \mathbb{R}$ such that $\mathbf{0} \in H^-(\mathbf{a}, \delta)$. We assume that we have normalized $\delta$ to be 0 or 1. If $\delta = 0$, then

$$\gamma_H(\mathbf{r}) = \begin{cases} 0 & \text{if } \langle \mathbf{a}, \mathbf{r} \rangle \leq 0 \\ +\infty & \text{if } \langle \mathbf{a}, \mathbf{r} \rangle > 0 \end{cases}$$

If $\delta = 1$, then

$$\gamma_H(\mathbf{r}) = \max\{0, \langle \mathbf{a}, \mathbf{r} \rangle\}.$$

The above calculation, along with the next theorem, gives powerful computational tools for gauge functions.

**Theorem 3.46.** Let $C_i$, $i \in I$ be a (not necessarily finite) family of closed, convex sets, and let $C = \cap_{i \in I} C_i$. Then

$$\gamma_C = \sup_{i \in I} \gamma_{C_i}.$$

*Proof.* Consider any $\mathbf{r} \in \mathbb{R}^d$. Let us define $A_i = \{\lambda > 0 : \mathbf{r} \in \lambda C_i\}$ for each $i \in I$, and define $A = \{\lambda > 0 : \mathbf{r} \in \lambda C\}$. Observe that $A = \cap A_i$. If any $A_i$ is empty, then $\gamma_{C_i} = \infty$, and $A$ is empty and therefore $\gamma_C = \infty$, and the equality holds. Now suppose all $A_i$'s are nonempty, and so by Lemma 3.37, each $A_i$ is of the form

NOTES: 57

$(a_i, \infty)$ or $[a_i, \infty)$. If $A = \emptyset$, then it must mean that $a_i \to \infty$. Since $\gamma_{C_i}(\mathbf{r}) = \inf A_i = a_i$, this shows that $\sup_{i \in I} \gamma_{C_i}(\mathbf{r}) = \infty$. Moreover, $A = \emptyset$ implies that $\gamma_C(\mathbf{r}) = \inf A = \infty$. Finally, consider the case that $A$ is nonempty. Then since $A = \cap A_i$, $\gamma_C(\mathbf{r}) = a = \sup_{i \in I} a_i = \sup_{i \in I} \gamma_{C_i}(\mathbf{r})$. $\qquad\square$

This shows that gauge functions for polyhedra can be computed very easily.

**Corollary 3.47.** Let $P$ be a polyhedron containing the origin in its interior. Thus, there exist $\mathbf{a}^1, \dots, \mathbf{a}^m \in \mathbb{R}^d$ such that
$$P = \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{a}^i, \mathbf{x} \rangle \leq 1 \ \ i = 1, \dots, m\}.$$

Then
$$\gamma_P(\mathbf{r}) = \max\{0, \langle \mathbf{a}^1, \mathbf{r} \rangle, \dots, \langle \mathbf{a}^m, \mathbf{r} \rangle\}.$$

*Proof.* Use the formula from 3.45 and Theorem 3.46. $\qquad\square$

**Support functions.** While gauges are good in the sense that they are a nice generalization of norms from centrally symmetric convex bodies to asymmetric convex bodies, there is a drawback. Gauges are a strict subset of sublinear functions because they are always nonnegative, while there are many sublinear functions that take negative values. We would like to establish a one-to-one correspondence between sublinear functions and all closed, convex sets. Note that the correspondence via the epigraph only establishes a correspondence with closed, convex cones, and that too not all closed, convex cones are covered. The right definition, it turns out, is inspired by optimization of linear functions over closed, convex sets.

**Definition 3.48.** Let $S \subseteq \mathbb{R}^d$ be any set. The *support function* for $S$ is a function on $\mathbb{R}^d$ defined as
$$\sigma_S(\mathbf{r}) = \sup_{\mathbf{x} \in S} \langle \mathbf{r}, \mathbf{x} \rangle.$$

The following is easy to verify, and aspects of it were already explored in the midterm and HWs.

**Proposition 3.49.** Let $S \subseteq \mathbb{R}^d$. Then
$$\sigma_S = \sigma_{\mathrm{cl}(S)} = \sigma_{\mathrm{conv}(S)} = \sigma_{\mathrm{cl}(\mathrm{conv}(S))}.$$

**Proposition 3.50.** Let $S \subseteq \mathbb{R}^d$. Then $\sigma_S$ is a closed, sublinear function, i.e., its epigraph is a closed, convex cone.

*Proof.* We first check that $\sigma_S$ is sublinear. We check positive homogeneity. For any $\mathbf{r} \in \mathbb{R}^d$ and $\lambda > 0$,
$$\sigma_S(\lambda \mathbf{r}) = \sup_{\mathbf{x} \in S} \langle \lambda \mathbf{r}, \mathbf{x} \rangle = \sup_{\mathbf{x} \in S} \lambda \langle \mathbf{r}, \mathbf{x} \rangle = \lambda \sup_{\mathbf{x} \in S} \langle \mathbf{r}, \mathbf{x} \rangle = \lambda \sigma_S(\mathbf{r}).$$

NOTES: 58

We check subadditivity. Let $\mathbf{r}^1, \mathbf{r}^2 \in \mathbb{R}^d$. Then,

$$
\begin{aligned}
\sigma_S(\mathbf{r}^1 + \mathbf{r}^2) &= \sup_{\mathbf{x} \in S} \langle \mathbf{r}^1 + \mathbf{r}^2, \mathbf{x} \rangle \\
&= \sup_{\mathbf{x} \in S} (\langle \mathbf{r}^1, \mathbf{x} \rangle + \langle \mathbf{r}^2, \mathbf{x} \rangle) \\
&\leq \sup_{\mathbf{x} \in S} \langle \mathbf{r}^1, \mathbf{x} \rangle + \sup_{\mathbf{x} \in S} \langle \mathbf{r}^2, \mathbf{x} \rangle \\
&= \sigma_S(\mathbf{r}^1) + \sigma_S(\mathbf{r}^2).
\end{aligned}
$$

Since $\sigma_S$ is the supremum of linear functions $\langle \mathbf{x}, \mathbf{r} \rangle$, $\mathbf{x} \in S$, epi$(f)$ is the intersection of closed halfspaces, which shows that it is closed. The fact that it is a convex cone follows from Proposition 3.33. $\qquad\square$

We now establish a fundamental correspondence between gauges and support functions via polarity.

**Theorem 3.51.** Let $C$ be a closed convex set containing the origin. Then

$$
\gamma_C = \sigma_{C^\circ}.
$$

*Proof.* Recall that $C = (C^\circ)^\circ$ by Proposition 2.30 part 2. Unwrapping the definitions, this says that

$$
C = \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{a}, \mathbf{x} \rangle \leq 1 \ \forall \mathbf{a} \in C^\circ\} = \cap_{\mathbf{a} \in C^\circ} H^-(\mathbf{a}, 1).
$$

By Theorem 3.46 and Example 3.45, we obtain that

$$
\gamma_C(\mathbf{r}) = \sup_{\mathbf{a} \in C^\circ} \gamma_{H^-(\mathbf{a}, 1)}(\mathbf{r}) = \sup_{\mathbf{a} \in C^\circ} \max\{0, \langle \mathbf{a}, \mathbf{r} \rangle\}.
$$

Since $\mathbf{0} \in C^\circ$, the last term above can be written as $\sup_{\mathbf{a} \in C^\circ} \langle \mathbf{a}, \mathbf{r} \rangle = \sigma_{C^\circ}(\mathbf{r})$. $\qquad\square$

**Example 3.52.** Consider the polyhedron

$$
P = \{\mathbf{x} \in \mathbb{R}^2 : -\mathbf{x}_1 - \mathbf{x}_2 \leq 1, \ \tfrac{1}{2}\mathbf{x}_1 - \mathbf{x}_2 \leq 1, \ -\mathbf{x}_1 + \tfrac{1}{2}\mathbf{x}_2 \leq 1\}.
$$

From Corollary 3.47, we obtain that

$$
\gamma_P(\mathbf{r}) = \max\{0, -\mathbf{r}_1 - \mathbf{r}_2, \tfrac{1}{2}\mathbf{r}_1 - \mathbf{r}_2, -\mathbf{r}_1 + \tfrac{1}{2}\mathbf{r}_2\},
$$

and by Theorem 3.39 part 2., we obtain that $P = \{\mathbf{x} \in \mathbb{R}^2 : \gamma_P(\mathbf{x}) \leq 1\}$. Now consider the function

$$
f(\mathbf{r}) = \max\{-\mathbf{r}_1 - \mathbf{r}_2, \tfrac{1}{2}\mathbf{r}_1 - \mathbf{r}_2, -\mathbf{r}_1 + \tfrac{1}{2}\mathbf{r}_2\}.
$$

It turns out that $P = \{\mathbf{x} \in \mathbb{R}^2 : f(\mathbf{x}) \leq 1\}$ because

$$
\begin{aligned}
\mathbf{x} \in P \quad &\Leftrightarrow \quad -\mathbf{x}_1 - \mathbf{x}_2 \leq 1, \ \tfrac{1}{2}\mathbf{x}_1 - \mathbf{x}_2 \leq 1, \ -\mathbf{x}_1 + \tfrac{1}{2}\mathbf{x}_2 \leq 1 \\
&\Leftrightarrow \quad \max\{-\mathbf{x}_1 - \mathbf{x}_2, \tfrac{1}{2}\mathbf{x}_1 - \mathbf{x}_2, -\mathbf{x}_1 + \tfrac{1}{2}\mathbf{x}_2\} \leq 1 \\
&\Leftrightarrow \quad f(\mathbf{x}) \leq 1.
\end{aligned}
$$

NOTES:

Notice that $f((1,1)) = -\frac{1}{2} \neq 0 = \gamma_P((1,1))$. Also, $f$ is sublinear because $f$ is the support function of the set $S = \{(-1,-1),(\frac{1}{2},-1),(-1,\frac{1}{2})\}$. This shows that Corollary 3.41 really breaks down if the assumption of compactness is removed. Even so, given a closed, convex set $C$, any sublinear function that has a set $C$ as its 1-sublevel set must match the gauge on $\mathbb{R}^d \setminus \mathrm{int}(\mathrm{rec}(C))$ (see Problem 7 from "HW for Week IX"). If you are interested in learning more about representing closed, convex sets as the sublevel sets of sublinear functions, please see [1] on exciting new results.

---

**Generalized Cauchy-Schwarz/Holder's inequality.** Using our relationship between norms and gauges and support functions, we can write an inequality which vastly generalizes Holder's inequality (and consequently, Cauchy-Schwarz' inequality) – see Proposition 2.32.

**Theorem 3.53.** Let $C \subseteq \mathbb{R}^d$ be a compact, convex set containing the origin in its interior. Then

$$\langle \mathbf{x}, \mathbf{y} \rangle \leq \gamma_C(\mathbf{x})\sigma_C(\mathbf{y}) \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d.$$

*Proof.* Consider any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$. Since $C$ is compact, $\gamma_C(\mathbf{x}) > 0$ by Corollary 3.40, and $\sigma_C(\mathbf{y}) < \infty$. By Proposition 3.38, $\frac{\mathbf{x}}{\gamma_C(\mathbf{x})} \in C$, and therefore,

$$\langle \frac{\mathbf{x}}{\gamma_C(\mathbf{x})}, \mathbf{y} \rangle \leq \sup_{\mathbf{z} \in C}\langle \mathbf{z}, \mathbf{y} \rangle = \sigma_C(\mathbf{y}).$$

This immediately implies $\langle \mathbf{x}, \mathbf{y} \rangle \leq \gamma_C(\mathbf{x})\sigma_C(\mathbf{y})$. $\qquad\square$

**Corollary 3.54.** Let $C \subseteq \mathbb{R}^d$ be a compact, convex set containing the origin in its interior. Then

$$\langle \mathbf{x}, \mathbf{y} \rangle \leq \gamma_C(\mathbf{x})\gamma_{C^\circ}(\mathbf{y}) \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d.$$

*Proof.* Follows from Theorems 3.53 and 3.51. $\qquad\square$

The above corollary generalizes Holder's inequality by recalling that when $\frac{1}{p} + \frac{1}{q} = 1$, then the $\ell^p$ and $\ell^q$ unit balls are polars of each other. Note that Theorem 3.53 and Corollary 3.54 have no assumption of centrally symmetric sets, so they strictly generalize the norm inequalities of Holder and Cauchy-Schwarz.

---

**One-to-one correspondence between closed, convex sets and closed, sublinear functions.** Proposition 3.50 shows that support functions are closed, sublinear functions. Proposition 3.49 shows that two different sets, e.g., $S$ and $\mathrm{conv}(S)$, may give rise to the same sublinear function $\sigma_S = \sigma_{\mathrm{conv}(S)}$ via the support function construction. In other words, if we consider the mapping $S \to \sigma_S$ as a mapping from the family of subsets of $\mathbb{R}^d$ to the family of closed, sublinear functions, this mapping is not injective. But if we restrict to closed, convex sets, it can shown that this mapping is injective.

NOTES: 60

**Exercise 8.** Let $C_1, C_2$ be closed, convex sets. Then $\sigma_{C_1} = \sigma_{C_2}$ if and only if $C_1 = C_2$.

A natural question now is whether the mapping $C \to \sigma_C$ from the family of closed, convex sets to the family of closed, sublinear functions is *onto*. The answer is yes! **Thus, all closed, sublinear functions are support functions and vice versa**.

**Theorem 3.55.** Let $f : \mathbb{R}^d \to \mathbb{R} \cup \{+\infty\}$ be a sublinear function that is also closed. Then the set

$$C_f := \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{r}, \mathbf{x} \rangle \leq f(\mathbf{r}) \ \forall \mathbf{r} \in \mathbb{R}^d\} = \cap_{\mathbf{r} \in \mathbb{R}^d} H^-(\mathbf{r}, f(\mathbf{r})) \tag{3.4}$$

is a closed, convex set. Moreover, $\sigma_{C_f} = f$.

Conversely, if $C$ is a closed, convex set, then $C_{\sigma_C} = C$.

*Proof.* We will prove the assertion when $f$ is finite valued everywhere; the proof for general $f$ is more tedious and does not provide any additional insight, in our opinion, and will be skipped here.

Since $C_f$ is defined as the intersection of a family of halfspace (indexed by $\mathbb{R}^d$), $C_f$ is a closed, convex set. We now establish that $\sigma_{C_f} = f$. For any $\mathbf{r} \in \mathbb{R}^d$, since $C_f \subseteq H^-(\mathbf{r}, f(\mathbf{r}))$, we must have that $\sigma_{C_f}(\mathbf{r}) = \langle \mathbf{r}, \mathbf{x} \in C \rangle \leq f(\mathbf{r})$. To show that $\sigma_{C_f}(\mathbf{r}) \geq f(\mathbf{r})$, it suffices to exhibit $\mathbf{y} \in C_f$ such that $\langle \mathbf{r}, \mathbf{y} \rangle = f(\mathbf{r})$. Consider epi$(f)$, which by Proposition 3.33, is a closed convex cone (since $f$ is assumed to be closed). By Theorem 2.23, there exists a supporting hyperplane for epi$(f)$ at $(\mathbf{r}, f(\mathbf{r}))$. Let this hyperplane by defined by $(\mathbf{y}, \eta) \in \mathbb{R}^d \times \mathbb{R}$ and $\alpha \in \mathbb{R}$ such that epi$(f) \subseteq H^-((\mathbf{y}, \eta), \alpha)$. Using Problems 8 and 9 from "HW for Week IX", one can assume that $\alpha = 0$ and $\eta < 0$. After normalizing, this means that epi$(f) \subseteq H^-((\mathbf{y}/-\eta, -1), 0)$. This implies that for every $\mathbf{r}' \in \mathbb{R}^d$, $(\mathbf{r}', f(\mathbf{r}')) \in H^-((\mathbf{y}/-\eta, -1), 0)$, which implies that $\langle \mathbf{r}', \frac{\mathbf{y}}{-\eta} \rangle \leq f(\mathbf{r}')$ for all $\mathbf{r}' \in \mathbb{R}^d$. So, $\frac{\mathbf{y}}{-\eta} \in C_f$. Moreover, since $H((\mathbf{y}/-\eta, -1), 0)$ is a supporting hyperplane at $(\mathbf{r}, f(\mathbf{r}))$, we must have $\langle \mathbf{r}, \frac{\mathbf{y}}{-\eta} \rangle - f(\mathbf{r}) = 0$. So, we are done.

We now show that $C_{\sigma_C} = C$ for any closed, convex set $C$. Consider any $\mathbf{x} \in C$. Then $\langle \mathbf{r}, \mathbf{x} \rangle \leq \sup_{\mathbf{y} \in C} \langle \mathbf{r}, \mathbf{y} \rangle = \sigma_C(\mathbf{r})$. Therefore, $\mathbf{x} \in H^-(\mathbf{r}, \sigma(\mathbf{r}))$ for all $\mathbf{r} \in \mathbb{R}^d$. This shows that $\mathbf{x} \in C_{\sigma_C}$, and therefore, $C \subseteq C_{\sigma_C}$. To show the reverse inclusion, consider any $\mathbf{y} \notin C$. Since $C$ is a closed, convex set, there exists a separating hyperplane $H(\mathbf{a}, \delta)$ such that $C \subseteq H^-(\mathbf{a}, \delta)$ and $\langle \mathbf{a}, \mathbf{y} \rangle > \delta$. $C \subseteq H^-(\mathbf{a}, \delta)$ implies that $\sigma_C(\mathbf{a}) = \sup_{\mathbf{x} \in C} \langle \mathbf{a}, \mathbf{x} \rangle \leq \delta$. Since $C_{\sigma_C}$ has $\langle \mathbf{a}, \mathbf{x} \rangle \leq \sigma_C(\mathbf{a})$ as a defining halfspace, and $\langle \mathbf{a}, \mathbf{y} \rangle > \delta \geq \sigma_C(\mathbf{a})$, we observe that $\mathbf{y} \notin C_{\sigma_C}$. $\qquad\square$

One can associate a nice picture with the above construction of $C_f$ associated with the sublinear function $f$, which corresponds to the following proposition.

**Proposition 3.56.** Let $f : \mathbb{R}^d \to \mathbb{R}$ be a sublinear function, and let $C_f$ be defined as in Theorem 3.55. Then $\mathbf{y} \in C_f$ if and only if $(\mathbf{y}, -1) \in$ epi$(f)^\circ$. In other words, $C_f = \{\mathbf{y} \in \mathbb{R}^d : (\mathbf{y}, -1) \in$ epi$(f)^\circ\}$.
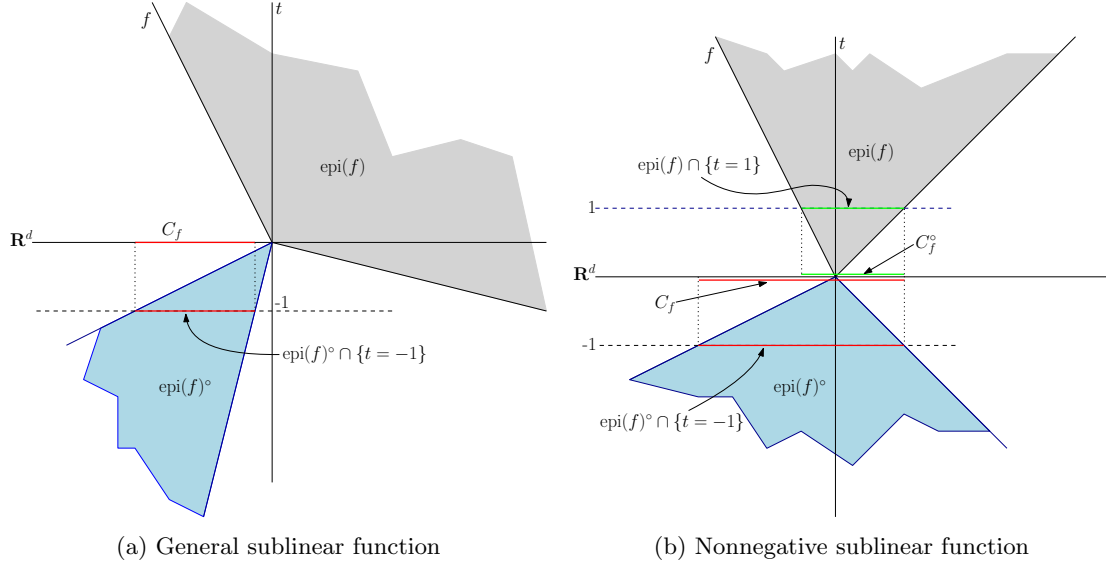
NOTES: 61

(a) General sublinear function      (b) Nonnegative sublinear function

Figure 1: Illustration of Propositions 3.56 and 3.57

*Proof.* We simply observe the following equivalences.

$$
\begin{array}{rll}
\mathbf{y} \in C_f & \Leftrightarrow \ \langle \mathbf{r}, \mathbf{y} \rangle \leq f(\mathbf{r}) & \forall \mathbf{r} \in \mathbb{R}^d \\
& \Leftrightarrow \ \langle \mathbf{r}, \mathbf{y} \rangle \leq t & \forall \mathbf{r} \in \mathbb{R}^d, t \in \mathbb{R} \text{ such that } f(\mathbf{r}) \leq t \\
& \Leftrightarrow \ \langle \mathbf{r}, \mathbf{y} \rangle - t \leq 0 & \forall \mathbf{r} \in \mathbb{R}^d, t \in \mathbb{R} \text{ such that } f(\mathbf{r}) \leq t \\
& \Leftrightarrow \ \langle (\mathbf{r}, t), (\mathbf{y}, -1) \rangle \leq 0 & \forall \mathbf{r} \in \mathbb{R}^d, t \in \mathbb{R} \text{ such that } f(\mathbf{r}) \leq t \\
& \Leftrightarrow \ \langle (\mathbf{y}, -1), (\mathbf{r}, t) \rangle \leq 0 & \forall (\mathbf{r}, t) \in \mathrm{epi}(f) \\
& \Leftrightarrow \ (\mathbf{y}, -1) \in \mathrm{epi}(f)^\circ &
\end{array}
$$

1355                                                              □

1356      When $f$ is a nonnegative sublinear function, even more can be said.

1357 **Proposition 3.57.** Let $f : \mathbb{R}^d \to \mathbb{R}$ be a sublinear function that nonnegative everywhere, and let $C_f$
1358 be defined as in Theorem 3.55. Then $f = \gamma_{(C_f)^\circ}$, i.e., $f$ is the gauge function for $(C_f)^\circ$. Consequently,
1359 $(C_f)^\circ = \{\mathbf{y} \in \mathbb{R}^d : (\mathbf{y}, 1) \in \mathrm{epi}(f)\} = \{\mathbf{y} \in \mathbb{R}^d : f(\mathbf{y}) \leq 1\}$.

1360 *Proof.* Since $f \geq 0$, $\mathrm{epi}(f) \subseteq \{(\mathbf{r}, t) : t \geq 0\}$. Therefore, $(\mathbf{0}, -1) \in \mathrm{epi}(f)^\circ$. By Proposition 3.56, $\mathbf{0} \in C_f$.
1361 Moreover, by Theorems 3.55 and 3.51, $f = \sigma_{C_f} = \gamma_{(C_f)^\circ}$. By Theorem 3.39 part 2., this shows that

NOTES:                                            62

$(C_f)^\circ = \{\mathbf{y} \in \mathbb{R}^d : f(\mathbf{y}) \leq 1\}$. By Problem 10 from the "HW for Week IX", we have that $(C_f)^\circ = \{\mathbf{y} \in \mathbb{R}^d : (\mathbf{y}, 1) \in \mathrm{epi}(f)\} = \{\mathbf{y} \in \mathbb{R}^d : f(\mathbf{y}) \leq 1\}$. $\square$

## 3.6 Directional derivatives, subgradients and subdifferential calculus

Let us look at directional derivatives of convex functions more closely. Let $f : \mathbb{R}^d \to \mathbb{R}$ be any function and let $\mathbf{x} \in \mathbb{R}^d$, and $\mathbf{r} \in \mathbb{R}^d$. We define the *directional derivative of $f$ at $\mathbf{x}$ in the direction $\mathbf{r}$* as:

$$f'(\mathbf{x}; \mathbf{r}) := \lim_{t \downarrow 0} \frac{f(\mathbf{x} + t\mathbf{r}) - f(\mathbf{x})}{t}, \tag{3.5}$$

if that limit exists. We will be speaking of $f'(\mathbf{x}; \cdot)$ as a function from $\mathbb{R}^d \to \mathbb{R}$. When the function $f$ is convex, this function has very nice properties.

**Lemma 3.58.** If $f : \mathbb{R}^d \to \mathbb{R}$ is convex, the expression $\frac{f(\mathbf{x}+t\mathbf{r})-f(\mathbf{x})}{t}$ is a non-decreasing function of $t$.

*Proof.* By Proposition 3.26, the function $\phi(t) = f(\mathbf{x} + t\mathbf{r})$ is a convex function. By Proposition 3.24, we observe that $\frac{\phi(t)-\phi(0)}{t}$ is a non-decreasing function of $t$. $\square$

**Proposition 3.59.** Let $f : \mathbb{R}^d \to \mathbb{R}$ be a convex function, and let $\mathbf{x} \in \mathbb{R}^d$. Then the limit in (3.5) exists for all $\mathbf{r} \in \mathbb{R}^d$ and the function $f'(\mathbf{x}; \cdot) : \mathbb{R}^d \to \mathbb{R}$ is sublinear.

*Proof.* By Proposition 3.26, the function $\phi(t) = f(\mathbf{x}+t\mathbf{r})$ is a convex function, and $f'(\mathbf{x}; \mathbf{r}) = \lim_{t\downarrow 0} \frac{\phi(t)-\phi(0)}{t}$. By Lemma 3.58, we observe that $\frac{\phi(t)-\phi(0)}{t}$ is a non-decreasing function of $t$, and restricting to $t > 0$, $\frac{\phi(t)-\phi(0)}{t}$ is lower bounded by the value at $t = -1$, i.e., $\frac{\phi(-1)-\phi(0)}{-1}$. Therefore, $\lim_{t\downarrow 0} \frac{\phi(t)-\phi(0)}{t}$ exists and is in fact equal to $\inf_{t>0} \frac{\phi(t)-\phi(0)}{t}$.

We now prove positive homogeneity of $f'(\mathbf{x}; \cdot)$. For any $\mathbf{r} \in \mathbb{R}^d$ and $\lambda > 0$, we obtain that

$$
\begin{aligned}
f'(\mathbf{x}; \lambda\mathbf{r}) &= \lim_{t\downarrow 0} \frac{f(\mathbf{x}+t\lambda\mathbf{r})-f(\mathbf{x})}{t} \\
&= \lim_{t\downarrow 0} \lambda \frac{f(\mathbf{x}+t\lambda\mathbf{r})-f(\mathbf{x})}{\lambda t} \\
&= \lambda \lim_{t\downarrow 0} \frac{f(\mathbf{x}+t\lambda\mathbf{r})-f(\mathbf{x})}{\lambda t} \\
&= \lambda \lim_{t'\downarrow 0} \frac{f(\mathbf{x}+t'\mathbf{r})-f(\mathbf{x})}{t'} \\
&= \lambda f'(\mathbf{x}; \mathbf{r}).
\end{aligned}
$$

We next establish that $f'(\mathbf{x}; \cdot)$ is convex. Consider any $\mathbf{r}^1, \mathbf{r}^2 \in \mathbb{R}^d$ and $\lambda \in (0, 1)$.

$$
\begin{aligned}
f'(\mathbf{x}; \lambda\mathbf{r}^1 + (1-\lambda)\mathbf{r}^2) &= \lim_{t\downarrow 0} \frac{f(\mathbf{x}+t(\lambda\mathbf{r}^1+(1-\lambda)\mathbf{r}^2))-f(\mathbf{x})}{t} \\
&= \lim_{t\downarrow 0} \lambda \frac{f(\lambda\mathbf{x}+(1-\lambda)\mathbf{x}+t(\lambda\mathbf{r}^1+(1-\lambda)\mathbf{r}^2))-\lambda f(\mathbf{x})-(1-\lambda)f(\mathbf{x})}{t} \\
&= \lim_{t\downarrow 0} \frac{f(\lambda(\mathbf{x}+t\mathbf{r}^1)+(1-\lambda)(\mathbf{x}+t\mathbf{r}^2))-\lambda f(\mathbf{x})-(1-\lambda)f(\mathbf{x})}{t} \\
&\leq \lim_{t\downarrow 0} \frac{\lambda f(\mathbf{x}+t\mathbf{r}^1)+(1-\lambda)f(\mathbf{x}+t\mathbf{r}^2)-\lambda f(\mathbf{x})-(1-\lambda)f(\mathbf{x})}{t} \\
&= \lambda \lim_{t\downarrow 0} \frac{f(\mathbf{x}+t\mathbf{r}^1)-f(\mathbf{x})}{t} + (1-\lambda)\lim_{t\downarrow 0} \frac{f(\mathbf{x}+t\mathbf{r}^2)-f(\mathbf{x})}{t} \\
&= \lambda f'(\mathbf{x}; \mathbf{r}^1) + (1-\lambda)f'(\mathbf{x}; \mathbf{r}^2),
\end{aligned}
$$

NOTES: 63

where the inequality follows from convexity of $f$. By Proposition 3.32, the function $f$ is sublinear. □

There is a nice connection with subgradients and subdifferentials – recall Definition 3.17. Also recall the construction of the closed, convex set $C_f$ from a sublinear function $f$ from Theorem 3.55.

**Theorem 3.60.** Let $f : \mathbb{R}^d \to \mathbb{R}$ be a convex function, and let $\mathbf{x} \in \mathbb{R}^d$. Then

$$\partial f(\mathbf{x}) = C_{f'(\mathbf{x};\cdot)}.$$

In other words, $f'(\mathbf{x};\cdot)$ is the support function for the subdifferential $\partial f(\mathbf{x})$.

*Proof.* Recall from Definitions 3.14 and 3.17 that

$$\begin{aligned} \partial f(\mathbf{x}) &= \{\mathbf{s} \in \mathbb{R}^d : \langle \mathbf{s}, \mathbf{y} - \mathbf{x} \rangle \leq f(\mathbf{y}) - f(\mathbf{x}) \ \forall \mathbf{y} \in \mathbb{R}^d\} \\ &= \{\mathbf{s} \in \mathbb{R}^d : \langle \mathbf{s}, \mathbf{r} \rangle \leq f(\mathbf{x} + \mathbf{r}) - f(\mathbf{x}) \ \forall \mathbf{r} \in \mathbb{R}^d\}. \end{aligned}$$

Thus, we have the following equivalences.

$$\begin{aligned} \mathbf{s} \in \partial f(\mathbf{x}) \quad &\Leftrightarrow \quad \langle \mathbf{s}, \mathbf{r} \rangle \leq f(\mathbf{x} + \mathbf{r}) - f(\mathbf{x}) && \forall \mathbf{r} \in \mathbb{R}^d \\ &\Leftrightarrow \quad \langle \mathbf{s}, t\mathbf{r} \rangle \leq f(\mathbf{x} + t\mathbf{r}) - f(\mathbf{x}) && \forall \mathbf{r} \in \mathbb{R}^d, t > 0 \\ &\Leftrightarrow \quad \langle \mathbf{s}, \mathbf{r} \rangle \leq \frac{f(\mathbf{x} + t\mathbf{r}) - f(\mathbf{x})}{t} && \forall \mathbf{r} \in \mathbb{R}^d, t > 0 \\ &\Leftrightarrow \quad \langle \mathbf{s}, \mathbf{r} \rangle \leq f'(\mathbf{x}; \mathbf{r}) && \forall \mathbf{r} \in \mathbb{R}^d \\ &\Leftrightarrow \quad \mathbf{s} \in C_{f'(\mathbf{x};\mathbf{r})} && \forall \mathbf{r} \in \mathbb{R}^d, \end{aligned}$$

where the second-to-last equivalence follows the fact that $\frac{f(\mathbf{x}+t\mathbf{r})-f(\mathbf{x})}{t}$ is a decreasing function of $t$ by Lemma 3.58, and the last equivalence follows from the definition of $C_{f'(\mathbf{x};\mathbf{r})}$ in (3.4). □

A characterization of differentiability for convex functions can be obtained using these concepts.

**Theorem 3.61.** Let $f : \mathbb{R}^d \to \mathbb{R}$ be a convex function, and let $\mathbf{x} \in \mathbb{R}^d$. Then the following are equivalent.

(i) $f$ is differentiable at $\mathbf{x}$.

(ii) $f'(\mathbf{x};\cdot)$ is a linear function given by $f'(\mathbf{x},\mathbf{r}) = \langle \mathbf{a_x}, \mathbf{r} \rangle$ for some $\mathbf{a_x} \in \mathbb{R}^d$.

(iii) $\partial f(\mathbf{x})$ is a singleton, i.e., there is a unique subgradient for $f$ at $\mathbf{x}$.

Moreover, if any of the above conditions hold then $\nabla f(\mathbf{x}) = \mathbf{a_x} = \mathbf{s}$, where $\mathbf{s}$ is the unique subgradient in $\partial f(\mathbf{x})$.

*Proof.* $(i) \implies (ii)$. If $f$ is differentiable, then it is well-known from calculus that $f'(\mathbf{x};\mathbf{r}) = \langle \nabla f(\mathbf{x}), \mathbf{r} \rangle$; thus, setting $\mathbf{a_x} = \nabla f(\mathbf{x})$ suffices.
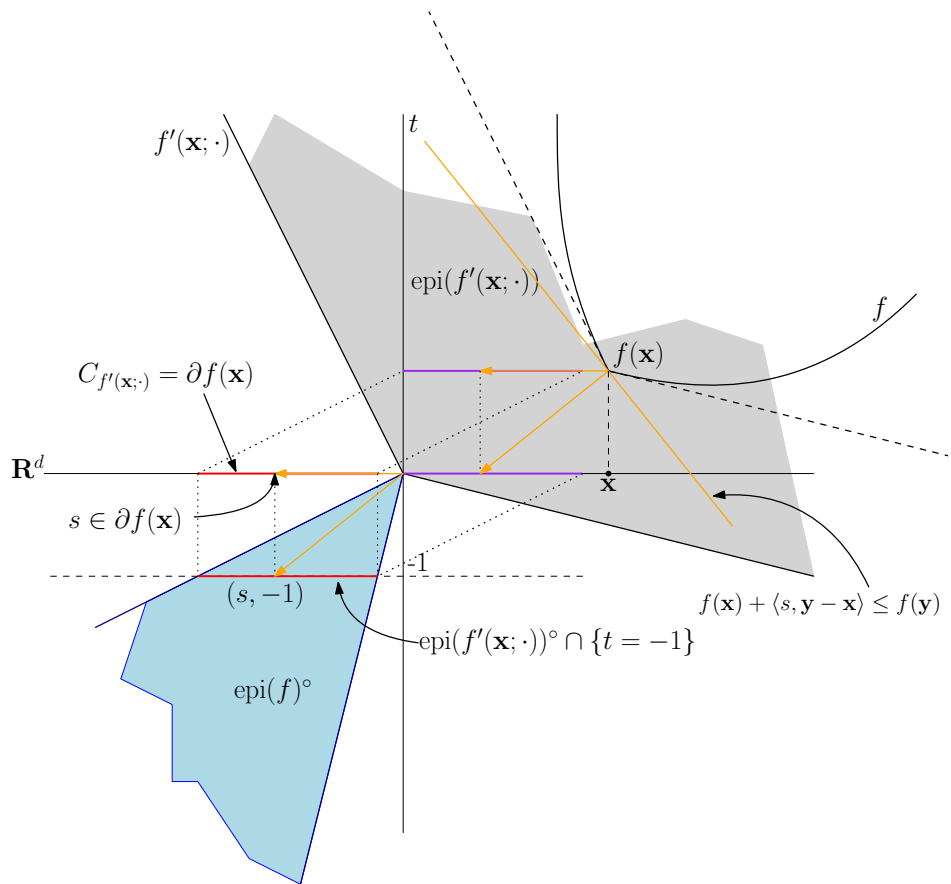
NOTES: 64

Figure 2: A picture illustrating the relationship between the sublinear function $f'(\mathbf{x}; \cdot)$, the set $C_{f'(\mathbf{x}; \cdot)}$, the subgradient $\partial f(\mathbf{x})$, and an affine support hyperplane given by an element $s \in \partial f(\mathbf{x})$. Recall the relationships from Figure 1.

$(ii) \implies (iii)$. By Theorem 3.60 and (3.4), we obtain that

$$
\begin{aligned}
\partial f(\mathbf{x}) &= C_{f'(\mathbf{x};\cdot)} \\
&= \{\mathbf{s} \in \mathbb{R}^d : \langle \mathbf{s}, \mathbf{r} \rangle \le f'(\mathbf{x}; \mathbf{r}) \;\; \forall \mathbf{r} \in \mathbb{R}^d\} \\
&= \{\mathbf{s} \in \mathbb{R}^d : \langle \mathbf{s}, \mathbf{r} \rangle \le \langle \mathbf{a_x}, \mathbf{r} \rangle \;\; \forall \mathbf{r} \in \mathbb{R}^d\}.
\end{aligned}
$$

We now observe that if $\langle \mathbf{s}, \mathbf{r} \rangle \le \langle \mathbf{a_x}, \mathbf{r} \rangle$ for all $\mathbf{r} \in \mathbb{R}^d$, then we must have $\mathbf{s} = \mathbf{a_x}$. Therefore, $\partial f(\mathbf{x}) = \{\mathbf{a_x}\}$.

$(iii) \implies (i)$. Let $\mathbf{s}$ be the unique subgradient at $\mathbf{x}$. We will establish that

$$
\lim_{\mathbf{h} \to \mathbf{0}} \frac{|f(\mathbf{x} + \mathbf{h}) - f(\mathbf{x}) - \langle \mathbf{s}, \mathbf{h} \rangle|}{\|\mathbf{h}\|} = 0,
$$

thus showing that $f$ is differentiable at $\mathbf{x}$ with gradient $\mathbf{s}$. In other words, given any $\delta > 0$, we must find $\epsilon > 0$ such that $\mathbf{h} \in B(\mathbf{0}, \epsilon)$ implies that $\frac{|f(\mathbf{x}+\mathbf{h}) - f(\mathbf{x}) - \langle \mathbf{s}, \mathbf{h} \rangle|}{\|\mathbf{h}\|} < \delta$.

Suppose to the contrary that for some $\delta > 0$, for every $k \ge 1$ there exists $\mathbf{h}_k$ such that $\|\mathbf{h}_k\| =: t_k \le \frac{1}{k}$ and $\frac{|f(\mathbf{x}+\mathbf{h}_k) - f(\mathbf{x}) - \langle \mathbf{s}, \mathbf{h}_k \rangle|}{t_k} \ge \delta$. Since $\frac{\mathbf{h}_k}{t_k}$ is a sequence of unit norm vectors, by Theorem 1.10, there is a convergent subsequence which converges to $\mathbf{r}$ with unit norm. To keep the notation easy, we relabel indices so that $\{\mathbf{h}_k\}_{k=1}^\infty$ is the convergent sequence. Using Theorem 3.21, there exists a constrant $L := L(B(\mathbf{0}, 1))$ such that $|f(\mathbf{y}) - f(\mathbf{z})| \le L\|\mathbf{y} - \mathbf{z}\|$ for all $\mathbf{y}, \mathbf{z} \in B(\mathbf{0}, 1)$. Noting that $\mathbf{h}_k$ and $t_k \mathbf{r}$ for all $k \ge 1$ are in the unit ball $B(\mathbf{0}, 1)$ (since $t_k \le \frac{1}{k}$),

$$
\begin{aligned}
\delta &\le \frac{|f(\mathbf{x}+\mathbf{h}_k) - f(\mathbf{x}) - \langle \mathbf{s}, \mathbf{h}_k \rangle|}{t_k} \\
&\le \frac{|f(\mathbf{x}+\mathbf{h}_k) - f(\mathbf{x}+t_k\mathbf{r})| + |f(\mathbf{x}+t_k\mathbf{r}) - f(\mathbf{x}) - \langle \mathbf{s}, t_k\mathbf{r} \rangle| + |\langle \mathbf{s}, t_k\mathbf{r} \rangle - \langle \mathbf{s}, \mathbf{h}_k \rangle|}{t_k} \\
&\le \frac{L\|t_k\mathbf{r} - \mathbf{h}_k\|}{t_k} + \frac{|f(\mathbf{x}+t_k\mathbf{r}) - f(\mathbf{x}) - \langle \mathbf{s}, t_k\mathbf{r} \rangle|}{t_k} + \frac{|\langle \mathbf{s}, t_k\mathbf{r} \rangle - \langle \mathbf{s}, \mathbf{h}_k \rangle|}{t_k} \\
&\le L\|\mathbf{r} - \tfrac{\mathbf{h}_k}{t_k}\| + |\tfrac{f(\mathbf{x}+t_k\mathbf{r}) - f(\mathbf{x})}{t_k} - \langle \mathbf{s}, \mathbf{r} \rangle| + \|\mathbf{s}\|\|\mathbf{r} - \tfrac{\mathbf{h}_k}{t_k}\| \\
&= (L + \|\mathbf{s}\|)\|\mathbf{r} - \tfrac{\mathbf{h}_k}{t_k}\| + |\tfrac{f(\mathbf{x}+t_k\mathbf{r}) - f(\mathbf{x})}{t_k} - \langle \mathbf{s}, \mathbf{r} \rangle|
\end{aligned}
$$

By letting $k \to \infty$, the last expression in the above goes to 0, contradicting that $\delta > 0$. $\qquad\square$

The following rules for manipulating subgradients and subdifferentials will be useful from an algorithmic perspective when we discuss optimization in the next section.

**Theorem 3.62. Subdifferential calculus.** The following are all true.

1. Let $f_1, f_2 : \mathbb{R}^d \to \mathbb{R}$ be convex functions and let $t_1, t_2 \ge 0$. Then

$$
\partial(t_1 f_1 + t_2 f_2)(\mathbf{x}) = t_1 \partial f_1(\mathbf{x}) + t_2 \partial f_2(\mathbf{x}) \text{ for all } \mathbf{x} \in \mathbb{R}^d.
$$

2. Let $A \in \mathbb{R}^{m \times d}$ and $b \in \mathbb{R}^m$ and let $T(\mathbf{x}) = A\mathbf{x} + \mathbf{b}$ be the corresponding affine map from $\mathbb{R}^d \to \mathbb{R}^m$ and let $g : \mathbb{R}^m \to \mathbb{R}$ be a convex function. Then

$$
\partial(g \circ T)(\mathbf{x}) = A^T \partial g(A\mathbf{x} + \mathbf{b}) \text{ for all } \mathbf{x} \in \mathbb{R}^d.
$$

NOTES: 66

3. Let $f_j : \mathbb{R}^d \to \mathbb{R}$, $j \in J$ be convex functions for some (possibly infinite) index set $J$, and let $f = \sup_{j \in J} f_j$. Then

$$\mathrm{cl}(\mathrm{conv}(\cup_{j \in J(\mathbf{x})} \partial f_j(\mathbf{x}))) \subseteq \partial f(\mathbf{x}),$$

where $J(\mathbf{x})$ is the set of indices $j$ such that $f_j(\mathbf{x}) = f(\mathbf{x})$. Moreover, equality holds in the above relation, if one can impose a topology on $J$ such that $J(\mathbf{x})$ is a compact set.

# 4 Optimization

We now being our study of the general convex optimization problem

$$\inf_{\mathbf{x} \in C} f(\mathbf{x}), \tag{4.1}$$

where $f : \mathbb{R}^d \to \mathbb{R}$ is a convex function, and $C$ is a closed, convex set. We first observe that local minimizers are global minimizers for convex optimization problems.

**Definition 4.1.** Let $g : \mathbb{R}^d \to \mathbb{R}$ be any function (not necessarily convex) and let $X \subseteq \mathbb{R}^d$ be any set (not necessarily convex). Then $\mathbf{x}^\star \in X$ is said to be a *local minimizer* for the problem $\inf_{\mathbf{x} \in X} g(\mathbf{x})$ is there exists $\epsilon > 0$ such that $g(\mathbf{y}) \geq g(\mathbf{x}^\star)$ for all $\mathbf{y} \in B(\mathbf{x}^\star, \epsilon) \cap X$.

$\mathbf{x}^\star \in X$ is said to be a *global minimizer* if $g(\mathbf{y}) \geq g(\mathbf{x}^\star)$ for all $\mathbf{y} \in X$.

Note that if $C$ is a compact, convex set, then (4.1) has a global minimizer by Weierstrass' Theorem (Theorem 1.11), because convex functions are continuous over the relative interior of their domain (Theorem 3.21).

**Theorem 4.2.** Any local minimizer for (4.1) is a global minimizer.

*Proof.* Let $\mathbf{x}^\star$ be a local minimizer, i.e., there exists $\epsilon > 0$ such that $f(\mathbf{y}) \geq f(\mathbf{x}^\star)$ for all $\mathbf{y} \in B(\mathbf{x}^\star, \epsilon) \cap C$. Suppose to the contrary that there exists $\bar{\mathbf{y}} \in C$ such that $f(\bar{\mathbf{y}}) < f(\mathbf{x}^\star)$. Then $\bar{\mathbf{y}} \notin B(\mathbf{x}^\star, \epsilon)$; otherwise, it would contradict $f(\mathbf{y}) \geq f(\mathbf{x}^\star)$ for all $\mathbf{y} \in B(\mathbf{x}^\star, \epsilon) \cap C$. Consider the line segment $[\mathbf{x}^\star, \bar{\mathbf{y}}]$. It must intersect $B(\mathbf{x}^\star), \epsilon)$ in a point other than $\mathbf{x}^\star$. Therefore, there exists $1 > \lambda > 0$ such that $\bar{\mathbf{x}} = \lambda \mathbf{x}^\star + (1 - \lambda)\bar{\mathbf{y}}$ is in $B(\mathbf{x}^\star, \epsilon)$. By convexity of $f$, $f(\bar{\mathbf{x}}) \leq \lambda f(\mathbf{x}^\star) + (1 - \lambda)f(\bar{\mathbf{y}})$. Since $\lambda \in (0, 1)$ and $f(\bar{\mathbf{y}}) < f(\mathbf{x}^\star)$, this implies that $f(\bar{x}) < f(\mathbf{x}^\star)$. Moreover, since $C$ is convex, $\bar{\mathbf{x}} \in C$, and so $\bar{\mathbf{x}} \in B(\mathbf{x}^\star, \epsilon) \cap C$. This contradicts that $f(\mathbf{y}) \geq f(\mathbf{x}^\star)$ for all $\mathbf{y} \in B(\mathbf{x}^\star, \epsilon) \cap C$ $\hfill\square$

We now give a characterization of global minimizers of (4.1) in terms of the local geometry of $C$ and the first order properties of $f$, i.e., its subdifferential $\partial f$. We first need some concepts related to the local geometry of a convex set.

**Definition 4.3.** Let $C \subseteq \mathbb{R}^d$ be a convex set, and let $\mathbf{x} \in C$. Define the *cone of feasible directions* as

$$F_C(\mathbf{x}) = \{\mathbf{r} \in \mathbb{R}^d : \exists \epsilon > 0 \text{ such that } \mathbf{x} + \epsilon \mathbf{r} \in C\}.$$

NOTES: 67

$F_C(\mathbf{x})$ may not be a closed cone – consider $C$ as the unit circle in $\mathbb{R}^2$ and $\mathbf{x} = (-1, 0)$; then $F_C(\mathbf{x}) = \{\mathbf{r} \in \mathbb{R}^2 : \mathbf{r}_1 > 0\} \cup \{\mathbf{0}\}$. It is much nicer to work with its closure.

**Definition 4.4.** Let $C \subseteq \mathbb{R}^d$ be a convex set, and let $\mathbf{x} \in C$. The *tangent cone of $C$ at $\mathbf{x}$* is $T_C(\mathbf{x}) := \mathrm{cl}(F_C(\mathbf{x}))$.

The final concept related to the local geometry of closed, convex sets will be the *normal cone.*

**Definition 4.5.** Let $C \subseteq \mathbb{R}^d$ be a convex set, and let $\mathbf{x} \in C$. The *normal cone of $C$ at $\mathbf{x}$* is $N_C(\mathbf{x}) := \{\mathbf{r} \in \mathbb{R}^d : \langle \mathbf{r}, \mathbf{x} \rangle \geq \langle \mathbf{r}, \mathbf{y} \rangle \ \forall \mathbf{y} \in C\}$.

The normal cone $N_C(\mathbf{x})$ is the set of vectors $\mathbf{r} \in \mathbb{R}^d$ such that $\mathbf{x}$ is the maximizer over $C$ for the corresponding linear functional $\langle \mathbf{r}, \cdot \rangle$, i.e., $\langle \mathbf{r}, \mathbf{x} \rangle = \sup_{\mathbf{y} \in C} \langle \mathbf{r}, \mathbf{y} \rangle$. Moreover, since $N_C(\mathbf{x}) = \{\mathbf{r} \in \mathbb{R}^d : \langle \mathbf{r}, \mathbf{y} - \mathbf{x} \rangle \leq 0 \ \forall \mathbf{y} \in C\}$ which is an intersection of halfspaces with the origin on the boundary, it is immediate that $N_C$ is a closed, convex cone.

**Proposition 4.6.** Let $C \subseteq \mathbb{R}^d$ be a convex set, and let $\mathbf{x} \in C$. Then $F_C(\mathbf{x}), T_C(\mathbf{x})$ and $N_C(\mathbf{x})$ are all convex cones, with $T_C(\mathbf{x}), N_C(\mathbf{x})$ being closed, convex cones. Moreover, $N_C(\mathbf{x}) = T_C(\mathbf{x})^\circ$, i.e., the tangent cone and the normal cone are polars of each other.

*Proof.* See Problem 4 in "HW for Week X". $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $\square$

We are now ready to state the characterization of a global minimizer of (4.1), in terms of the local geometry of $C$ and the first-order information of $f$.

**Theorem 4.7.** Let $f : \mathbb{R}^d \to \mathbb{R}$ be a convex function, and $C$ be a closed, convex set. Then the following are all equivalent.

1. $\mathbf{x}^\star$ is a global minimizer of (4.1).

2. $f'(\mathbf{x}^\star; \mathbf{y} - \mathbf{x}^\star) \geq 0$ for all $\mathbf{y} \in C$.

3. $f'(\mathbf{x}^\star; \mathbf{r}) \geq 0$ for all $\mathbf{r} \in T_C(\mathbf{x}^\star)$.

4. $\mathbf{0} \in \partial f(\mathbf{x}^\star) + N_C(\mathbf{x}^\star)$.

*Proof.* 1. $\implies$ 2. Since $f(\mathbf{z}) \geq f(\mathbf{x}^\star)$ for all $\mathbf{z} \in C$, in particular this holds for $\mathbf{z} = \mathbf{x} + t(\mathbf{y} - \mathbf{x})$ for all $0 \leq t \leq 1$. Therefore, $\frac{f(\mathbf{x}^\star + t(\mathbf{y} - \mathbf{x}^\star)) - f(\mathbf{x}^\star)}{t} \geq 0$ for all $t \in (0, 1)$. Taking the limit as $t \to 0$, we obtain that $f'(\mathbf{x}^\star; \mathbf{y} - \mathbf{x}^\star) \geq 0$.

2. $\implies$ 3. We first show that $f'(x^\star; \mathbf{r}) \geq 0$ for all $\mathbf{x} \in F_C(\mathbf{x})$. Let $\epsilon > 0$ such that $\mathbf{y} = \mathbf{x}^\star + \epsilon \mathbf{r} \in C$. By assumption, $0 \leq f'(\mathbf{x}^\star; \mathbf{y} - \mathbf{x}^\star) = f'(\mathbf{x}^\star; \epsilon \mathbf{r}) = \epsilon f'(\mathbf{x}; \mathbf{r})$, using the positive homogeneity of $f'(\mathbf{x}^\star; \cdot)$, since $f'(\mathbf{x}^\star; \cdot)$ is sublinear by Proposition 3.59. Diving by $\epsilon$, we obtain that $f'(\mathbf{x}^\star; \mathbf{r}) \geq 0$ for all $\mathbf{r} \in F_C(\mathbf{x}^\star)$.

Since $f'(\mathbf{x}^\star; \cdot)$ is sublinear, it is convex by Proposition 3.32, and thus, it is continuous by Theorem 3.21. Consequently, it must be nonnegative on $T_C(\mathbf{x}) = \text{cl}(F_C(\mathbf{x}))$, because it is nonnegative on $F_C(\mathbf{x})$.

3. $\implies$ 4. Suppose to the contrary that $\mathbf{0} \notin \partial f(\mathbf{x}^\star) + N_C(\mathbf{x}^\star)$. Since $f$ is assumed to be finite-valued everywhere, $\text{dom}(f) = \mathbb{R}^d$. Thus, by Problem 15 in "HW for Week X", $\partial f(\mathbf{x}^\star)$ is a compact, convex set. Moreover, $N_C(\mathbf{x}^\star)$ is a closed, convex cone by Proposition 4.6. Therefore, by Problem 6 in "HW for Week I/II", $\partial f(\mathbf{x}^\star) + N_C(\mathbf{x}^\star)$ is a closed, convex set. By the separating hyperplane theorem (Theorem 2.20), there exist $\mathbf{a} \in \mathbb{R}^d, \delta \in \mathbb{R}$ such that $0 = \langle \mathbf{a}, \mathbf{0} \rangle > \delta \geq \langle \mathbf{a}, \mathbf{v} \rangle$ for all $\mathbf{v} \in \partial f(\mathbf{x}^\star) + N_C(\mathbf{x}^\star)$.

First, we claim that $\langle \mathbf{a}, \mathbf{n} \rangle \leq 0$ for all $\mathbf{n} \in N_C(\mathbf{x}^\star)$. Otherwise, consider $\bar{\mathbf{n}} \in N_C(\mathbf{x}^\star)$ such that $\langle \mathbf{a}, \bar{\mathbf{n}} \rangle > 0$. Since $N_C(\mathbf{x}^\star)$ is a closed, convex cone, $\lambda \bar{\mathbf{n}} \in N_C(\mathbf{x}^\star)$ for all $\lambda \geq 0$. But then consider any $\mathbf{s} \in \partial f(\mathbf{x})$ (which is nonempty by Problem 15 in "HW for Week X") and the set of points $\mathbf{s} + \lambda \bar{\mathbf{n}}$. Since $\langle \mathbf{a}, \bar{\mathbf{n}} \rangle > 0$, we can find $\lambda \geq 0$ large enough such that $\langle \mathbf{a}, \mathbf{s} + \lambda \bar{\mathbf{n}} \rangle > \delta$, contradicting that $\delta \geq \langle \mathbf{a}, \mathbf{v} \rangle$ for all $\mathbf{v} \in \partial f(\mathbf{x}^\star) + N_C(\mathbf{x}^\star)$.

Since $\langle \mathbf{a}, \mathbf{n} \rangle \leq 0$ for all $\mathbf{n} \in N_C(\mathbf{x}^\star)$, we obtain that $\mathbf{a} \in N_C(\mathbf{x}^\star)^\circ = T_C(\mathbf{x}^\star)$, by Proposition 4.6. Now we use the fact that $\partial f(\mathbf{x}^\star) \subseteq \partial f(\mathbf{x}^\star) + N_C(\mathbf{x}^\star)$, since $\mathbf{0} \in N_C(\mathbf{x}^\star)$. This implies that $\langle \mathbf{a}, \mathbf{s} \rangle \leq \delta < 0$ for all $\mathbf{s} \in \partial f(\mathbf{x}^\star)$. Since $\partial f(\mathbf{x}^\star)$ is a compact, convex set, this implies that $\sup_{\mathbf{s} \in \partial f(\mathbf{x}^\star)} \langle \mathbf{a}, \mathbf{s} \rangle < 0$. From Theorem 3.60, $f'(\mathbf{x}^\star; \mathbf{a}) = \sigma_{\partial f(\mathbf{x}^\star)}(\mathbf{a}) = \sup_{\mathbf{s} \in \partial f(\mathbf{x}^\star)} \langle \mathbf{a}, \mathbf{s} \rangle < 0$. This contradicts the assumption of 3., because we showed above that $\mathbf{a} \in T_C(\mathbf{x}^\star)$.

4. $\implies$ 1. Consider any $\mathbf{y} \in C$. Since $\mathbf{0} \in \partial f(\mathbf{x}^\star) + N_C(\mathbf{x}^\star)$, there exist $\mathbf{s} \in \partial f(\mathbf{x}^\star)$ and $\mathbf{n} \in N_C(\mathbf{x}^\star)$ such that $\mathbf{0} = \mathbf{s} + \mathbf{n}$. Now, $\mathbf{y} - \mathbf{x}^\star \in T_C(\mathbf{x}^\star)$ and so $\langle \mathbf{y} - \mathbf{x}^\star, \mathbf{n} \rangle \leq 0$ by Proposition 4.6. Since we have

$$0 = \langle \mathbf{y} - \mathbf{x}^\star, \mathbf{0} \rangle = \langle \mathbf{y} - \mathbf{x}^\star, \mathbf{s} \rangle + \langle \mathbf{y} - \mathbf{x}^\star, \mathbf{n} \rangle,$$

this implies that $\langle \mathbf{y} - \mathbf{x}^\star, \mathbf{s} \rangle \geq 0$. By definition of subgradient, $f(\mathbf{y}) \geq f(\mathbf{x}^\star) + \langle \mathbf{s}, \mathbf{y} - \mathbf{x}^\star \rangle \geq f(\mathbf{x}^\star)$. Since the choice of $\mathbf{y} \in C$ was arbitrary, this shows that $\mathbf{x}^\star$ is a global minimizer. $\square$

**Algorithmic setup: First-order oracles.** To tackle the problem (4.1) computationally, we have to set up a precise way to access the values/subgradients of the function $f$ and test if given points belong to the set $C$ or not. To make this algorithmically clean, we define *first-order oracles*.

**Definition 4.8.** A *first order oracle* for a convex function $f : \mathbb{R}^d \to \mathbb{R}$ is an oracle/algorithm/black-box that takes as input any $\mathbf{x} \in \mathbb{R}^d$ and returns $f(\mathbf{x})$ and some $\mathbf{s} \in \partial f(\mathbf{x})$. A *first order oracle* for a closed, convex set $C \subseteq \mathbb{R}^d$ is an oracle/algorithm/black-box that takes as input any $\mathbf{x} \in \mathbb{R}^d$ and either correctly reports that $\mathbf{x} \in C$ or correctly reports a separating hyperplane separating $\mathbf{x}$ from $C$, i.e., it returns $\mathbf{a} \in \mathbb{R}^d, \delta \in \mathbb{R}$ such that $C \subseteq H^-(\mathbf{a}, \delta)$ and $\langle \mathbf{a}, \mathbf{x} \rangle > \delta$. Such an oracle is also known as a *separation oracle*.

## 4.1 Subgradient algorithm

To build up towards an algorithm that assumes only first-order oracles for $f$ and $C$, we will first look at the situation where we have a first order oracle for $f$, and a *stronger* oracle for $C$ which, given any $\mathbf{x} \in \mathbb{R}^d$, can report the closest point in $C$ to $\mathbf{x}$. Recall that in the proof of Theorem 2.20, we had shown that such a closest point always exists as long as $C$ is a closed, convex set.

NOTES: 69

**Definition 4.9.** $\operatorname{Proj}_C(\mathbf{x})$ will denote the closest point (under the standard Euclidean norm) in $C$ to $\mathbf{x}$.

Note that an oracle that reports $\operatorname{Proj}_C(\mathbf{x})$ for any $\mathbf{x} \in \mathbb{R}^d$ is stronger than a separation oracle for $C$, because $\operatorname{Proj}_C(\mathbf{x}) = \mathbf{x}$ if and only if $\mathbf{x} \in C$, and when $\operatorname{Proj}_C(\mathbf{x}) \neq \mathbf{x}$, then one can use $\mathbf{a} = \mathbf{x} - \operatorname{Proj}_C(\mathbf{x})$ and $\delta = \langle \mathbf{a}, \operatorname{Proj}_C(\mathbf{x}) \rangle$ as a separating hyperplane; see the proof of Theorem 2.20. Even so, for "simple" sets $C$, computing $\operatorname{Proj}_C(\mathbf{x})$ is not a difficult task. For example, when $C = \mathbb{R}_+^d$, then $\operatorname{Proj}_C(\mathbf{x}) = \mathbf{y}$, where $\mathbf{y}_i = \max\{0, \mathbf{x}_i\}$ for all $i = 1, \ldots, d$.

We now give a simple and elegant algorithm to solve the problem 4.1 when one has access to an oracle that can output $\operatorname{Proj}_C(\mathbf{x})$ for any $\mathbf{x} \in \mathbb{R}^d$, and a first-order oracle for $f$. The algorithm does not assume any properties beyond convexity for the function $f$ (e.g., differentiability). Note that, in particular, when we have no constraints, i.e., $C = \mathbb{R}^n$, then $\operatorname{Proj}_C(\mathbf{x}) = \mathbf{x}$ for all $\mathbf{x} \in \mathbb{R}^n$. Therefore, this algorithm can be used for *unconstrained optimization of general convex functions* with only a first-order oracle for $f$.

### Subgradient Algorithm.

1. Choose any sequence $h_0, h_1, \ldots,$ of strictly positive numbers. Let $\mathbf{x}_0 \in \mathbb{R}^d$.

2. For $i = 1, 2, \ldots,$ do

    (a) Use the first-order oracle for $f$ to get some $\mathbf{s}^i \in \partial f(\mathbf{x}^i)$.

    (b) Set $\mathbf{x}^{i+1} = \operatorname{Proj}_C\big(\mathbf{x}^i - h_i \frac{\mathbf{s}^i}{\|\mathbf{s}^i\|}\big)$.

The points $\mathbf{x}^0, \mathbf{x}^1, \ldots$ will be called the *iterates* of the Subgradient Algorithm. We now do a simple convergence analysis for the algorithm. First, a simple observation about the point $\operatorname{Proj}_C(\mathbf{x})$.

**Lemma 4.10.** Let $C \subseteq \mathbb{R}^d$ be a closed, convex set, let $\mathbf{x}^\star \in C$ and $\mathbf{x} \in \mathbb{R}^d$ (not necessarily in $C$). Then

$$\|\operatorname{Proj}_C(\mathbf{x}) - \mathbf{x}^\star\| \leq \|\mathbf{x} - \mathbf{x}^\star\|.$$

*Proof.* The proof of Theorem 2.20 shows that if we set $\mathbf{a} = \mathbf{x} - \operatorname{Proj}_C(\mathbf{x})$, then $\langle \mathbf{a}, \operatorname{Proj}_C(\mathbf{x}) - \mathbf{y} \rangle \geq 0$ for all $\mathbf{y} \in C$; in particular, $\langle \mathbf{a}, \operatorname{Proj}_C(\mathbf{x}) - \mathbf{x}^\star \rangle \geq 0$. We now observe that

$$
\begin{aligned}
\|\mathbf{x} - \mathbf{x}^\star\|^2 &= \|\mathbf{x} - \operatorname{Proj}_C(\mathbf{x}) + \operatorname{Proj}_C(\mathbf{x}) - \mathbf{x}^\star\|^2 \\
&= \|\mathbf{a} + \operatorname{Proj}_C(\mathbf{x}) - \mathbf{x}^\star\|^2 \\
&= \|\mathbf{a}\|^2 + \|\operatorname{Proj}_C(\mathbf{x}) - \mathbf{x}^\star\|^2 + 2\langle \mathbf{a}, \operatorname{Proj}_C(\mathbf{x}) - \mathbf{x}^\star \rangle \\
&\geq \|\operatorname{Proj}_C(\mathbf{x}) - \mathbf{x}^\star\|^2,
\end{aligned}
$$

since $\langle \mathbf{a}, \operatorname{Proj}_C(\mathbf{x}) - \mathbf{x}^\star \rangle \geq 0$. $\qquad \square$

**Theorem 4.11.** Let $f : \mathbb{R}^d \to \mathbb{R}$ be a convex function, and let $\mathbf{x}^\star \in \arg\min_{\mathbf{x} \in C} f(\mathbf{x})$. Suppose $\mathbf{x}_0 \in B(\mathbf{x}^\star, R)$ for some real number $R \geq 0$. Let $M := M(B(\mathbf{x}^\star, R))$ be a Lipschitz constant for $f$, guaranteed to exist by

NOTES: 70

Theorem 3.21, i.e., $|f(\mathbf{x}) - f(\mathbf{y})| \le M\|\mathbf{x} - \mathbf{y}\|$ for all $\mathbf{x}, \mathbf{y} \in B(\mathbf{x}^\star, R)$. Let $\mathbf{x}^0, \mathbf{x}^1, \ldots$ be the sequence of iterates obtained by the Subgradient Algorithm above. Then,

$$\min_{i=0,\ldots,k} f(\mathbf{x}^i) \le f(\mathbf{x}^\star) + M\left(\frac{R^2 + \sum_{i=0}^{k} h_i^2}{2\sum_{i=0}^{k} h_i}\right).$$

*Proof.* Define $r_i = \|\mathbf{x}^i - \mathbf{x}^\star\|$ and $v_i = \frac{\langle \mathbf{s}^i, \mathbf{x}^i - \mathbf{x}^\star\rangle}{\|\mathbf{s}^i\|}$ for $i = 0, 1, 2\ldots,$. We next observe that

$$
\begin{aligned}
r_{i+1}^2 &= \|\text{Proj}_C\big(\mathbf{x}^i - h_i\tfrac{\mathbf{s}^i}{\|\mathbf{s}^i\|}\big) - \mathbf{x}^\star\|^2 \\
&\le \|\mathbf{x}^i - h_i\tfrac{\mathbf{s}^i}{\|\mathbf{s}^i\|} - \mathbf{x}^\star\|^2 \qquad \text{by Lemma 4.10} \\
&= \|\mathbf{x}^i - \mathbf{x}^\star\|^2 + h_i^2 - 2h_i v_i \\
&= r_i^2 + h_i^2 - 2h_i v_i
\end{aligned}
$$

Adding these inequalities for $i = 0, 1, \ldots, k$, we obtain that

$$r_{k+1}^2 \le r_0^2 + \sum_{i=0}^{k} h_i^2 - 2\sum_{i=0}^{k} h_i v_i. \tag{4.2}$$

Let $v_{\min} = \min_{i=0,\ldots,k} v_i$ and let $i^{\min}$ be such that $v_{\min} = v_{i^{\min}}$. Using the fact that $r_0^2 = \|\mathbf{x}^0 - \mathbf{x}^\star\|^2 \le R^2$, and that $r_{k+1}^2 \ge 0$, we obtain from (4.2) that

$$v_{\min}\Big(2\sum_{i=0}^{k} h_i\Big) \le 2\sum_{i=0}^{k} h_i v_i \le R^2 + \sum_{i=0}^{k} h_i^2.$$

Consequently,

$$v_{\min} \le \frac{R^2 + \sum_{i=0}^{k} h_i^2}{2\sum_{i=0}^{k} h_i}. \tag{4.3}$$

Consider the hyperplane $H := H(\mathbf{s}^{i^{\min}}, \langle \mathbf{s}^{i^{\min}}, \mathbf{x}^{i^{\min}}\rangle)$ passing through $\mathbf{x}^{i^{\min}}$, orthogonal to $\mathbf{s}^{i^{\min}}$. Let $\bar{\mathbf{x}}$ be the point on $H$ closest to $\mathbf{x}^\star$. By Problem 12 in "HW for Week XI", $v_{\min} = \|\bar{\mathbf{x}} - \mathbf{x}^\star\|$. Moreover, $v_{\min} \le v_0 \le \|\mathbf{x}^0 - \mathbf{x}^\star\| \le R$. Therefore, $\bar{\mathbf{x}} \in B(\mathbf{x}^\star, R)$. Using the Lipschitz constant $M$, we obtain that $f(\bar{\mathbf{x}}) \le f(\mathbf{x}^\star) + Mv_{\min}$. Finally, since $\mathbf{s}^{i^{\min}} \in \partial f(\mathbf{x}^{i^{\min}})$, we must have that $f(\bar{\mathbf{x}}) \ge f(\mathbf{x}^{i^{\min}}) + \langle \mathbf{s}^{i^{\min}}, \bar{\mathbf{x}} - \mathbf{x}^{i^{\min}}\rangle = f(\mathbf{x}^{i^{\min}})$, since $\bar{\mathbf{x}} \in H$. Therefore, we obtain

$$\min_{i=0,\ldots,k} f(\mathbf{x}^i) \le f(\mathbf{x}^{i^{\min}}) \le f(\bar{\mathbf{x}}) \le f(\mathbf{x}^\star) + Mv_{\min} \le f(\mathbf{x}^\star) + M\left(\frac{R^2 + \sum_{i=0}^{k} h_i^2}{2\sum_{i=0}^{k} h_i}\right),$$

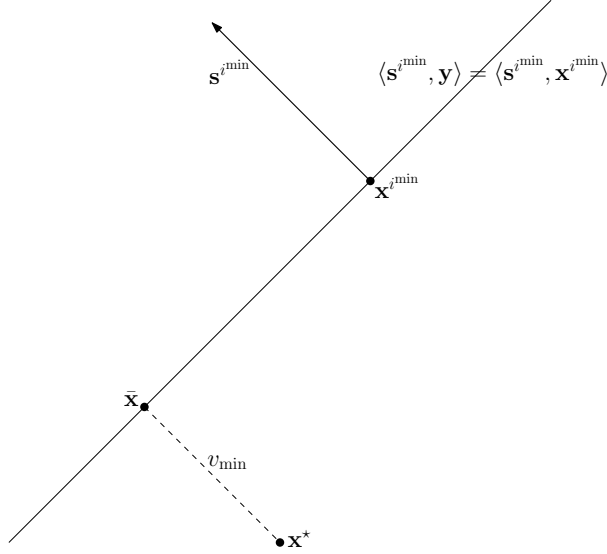where the last inequality follows from (4.3); see Figure 3. $\qquad\square$

NOTES: 71

Figure 3: Using $v_{\min}$ to bound the function value.

If we fix the number of steps of the algorithm to be $N \in \mathbb{N}$, then the choice of $h_0, \ldots, h_N$ that minimizes $\frac{R^2 + \sum_{i=0}^{k} h_i^2}{2 \sum_{i=0}^{k} h_i}$ is where $h_i = \frac{R}{\sqrt{N+1}}$ for all $i = 0, \ldots, N$, which yields the following corollary.

**Corollary 4.12.** Let $f : \mathbb{R}^d \to \mathbb{R}$ be a convex function, and let $\mathbf{x}^\star \in \arg\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$. Suppose $\mathbf{x}_0 \in B(\mathbf{x}^\star, R)$ for some real number $R \geq 0$. Let $M := M(B(\mathbf{x}^\star, R))$ be a Lipschitz constant for $f$. Let $N \in \mathbb{N}$ be any natural number, and set $h_i = \frac{R}{\sqrt{N+1}}$ for all $i = 0, \ldots, N$. Then the iterates of the Subgradient Algorithm, with this choice of $h_i$, satisfy

$$\min_{i=0,\ldots,N} f(\mathbf{x}^i) \leq f(\mathbf{x}^\star) + \frac{MR}{\sqrt{N+1}}.$$

Turning this around, if we want to be within $\epsilon$ of the optimal value $f(\mathbf{x}^\star)$ for some $\epsilon > 0$, we should run the Subgradient Algorithm for $\frac{M^2 R^2}{\epsilon^2}$ iterates, with $h_i = \frac{\epsilon}{M}$.

If we theoretically let the algorithm run for infinitely many steps, we would hope to make the difference between $\min_i f(\mathbf{x}^i)$ and $f(\mathbf{x}^\star)$ go to 0 in the limit. This, of course, depends on the choice of the sequence $h_0, h_1, \ldots$ so that the expression $\frac{R^2 + \sum_{i=0}^{k} h_i^2}{2 \sum_{i=0}^{k} h_i} \to 0$ as $k \to \infty$. There is a general sufficient condition that guarantees this.

NOTES: 72

**Corollary 4.13.** Let $\{h_i\}_{i=0}^{\infty}$ be a sequence of strictly positive real numbers such that $\lim_{i \to \infty} h_i = 0$ and $\sum_{i=1}^{\infty} h_i = \infty$. Then, for any real number $R$,

$$\lim_{k \to \infty} \frac{R^2 + \sum_{i=0}^{k} h_i^2}{2 \sum_{i=0}^{k} h_i} = 0.$$

**Remark 4.14.** Corollary 4.12 shows that the subgradient algorithm has a convergence that is *independent* of the dimension! Now matter how large $d$ is, as long as one can access subgradients for $f$ and project to $C$, the number of iterations needed to converge to within $\epsilon$ is $O(\frac{1}{\epsilon^2})$. This is important to keep in mind for applications where the dimension is extremely large.

## 4.2 Generalized inequalities and convex mappings

We first review the notion of a partially ordered set.

**Definition 4.15.** Let $X$ be any set. A *partial order* on $X$ is binary relation on $X$, i.e., a subset $\mathcal{R} \subseteq X \times X$ that satisfies certain conditions. We will denote $x \preccurlyeq y$ for $x, y \in X$ if $(x, y) \in \mathcal{R}$. The conditions are as follows:

1. $x \preccurlyeq x$ for all $x \in X$.

2. $x \preccurlyeq y$ and $y \preccurlyeq z$ implies $x \preccurlyeq z$.

3. $x \preccurlyeq y$ and $y \preccurlyeq x$ if and only if $x = y$.

We would like to be able to define partial orders on $\mathbb{R}^m$ for any $m \geq 1$. In doing so, we want to be mindful of the vector space structure of $\mathbb{R}^m$.

**Definition 4.16.** We will say that a binary relation on $\mathbb{R}^m$ is a *generalized inequality*, if it satisfies the following conditions.

1. $\mathbf{x} \preccurlyeq \mathbf{x}$ for all $\mathbf{x} \in \mathbb{R}^m$.

2. $\mathbf{x} \preccurlyeq \mathbf{y}$ and $\mathbf{y} \preccurlyeq \mathbf{z}$ implies $\mathbf{x} \preccurlyeq \mathbf{z}$.

3. $\mathbf{x} \preccurlyeq \mathbf{y}$ and $\mathbf{y} \preccurlyeq \mathbf{x}$ if and only if $\mathbf{x} = \mathbf{y}$.

4. $\mathbf{x} \preccurlyeq \mathbf{y}$ implies $\mathbf{x} + \mathbf{z} \preccurlyeq \mathbf{y} + \mathbf{z}$ for all $\mathbf{z} \in \mathbb{R}^m$.

5. $\mathbf{x} \preccurlyeq \mathbf{y}$ implies $\lambda \mathbf{x} \preccurlyeq \lambda \mathbf{y}$ for all $\lambda \geq 0$.

Generalized inequalities have an elegant geometric characterization.

NOTES: 73

**Proposition 4.17.** Let $K \subseteq \mathbb{R}^m$ be a closed, convex, pointed cone. Then, the relation on $\mathbb{R}^m$ defined by $\mathbf{x} \preccurlyeq_K \mathbf{y}$ if and only if $\mathbf{y} - \mathbf{x} \in K$, is a generalized inequality. In this case, we say that $\preccurlyeq_K$ is the generalized inequality *induced by $K$*.

Conversely, any generalized inequality $\preccurlyeq$ is induced by a unique cone given by $K_{\preccurlyeq} = \{\mathbf{x} \in \mathbb{R}^d : \mathbf{0} \preccurlyeq \mathbf{x}\}$. In other words, $\preccurlyeq$ is the same relation as $\preccurlyeq_{K_{\preccurlyeq}}$.

*Proof.* Left as an exercise. $\square$

**Example 4.18.** Here are some examples of generalized inequalities.

1. $K = \mathbb{R}^m_+$ induces the generalized inequality $\mathbf{x} \preccurlyeq_K \mathbf{y}$ if and only if $\mathbf{x}_i \leq \mathbf{y}_i$ for all $i = 1 \ldots, m$. This is often abbreviated to $\mathbf{x} \leq \mathbf{y}$, and is sometimes called the "canonical" generalized inequality on $\mathbb{R}^m$.

2. $K = \{\mathbf{x} \in \mathbb{R}^d : \sqrt{\mathbf{x}_1^2 + \ldots + \mathbf{x}_{d-1}^2} \leq \mathbf{x}_d\}$. This cone is called the *Lorentz cone*, and the corresponding generalized inequality is called a *second order cone constraints (SOCC)*.

3. Let $m = n^2$ for some $n \in \mathbb{N}$, i.e., consider the space $\mathbb{R}^{n^2}$. Identifying $\mathbb{R}^{n^2} = \mathbb{R}^{n \times n}$ with some ordering of the coordinates, we think of $\mathbb{R}^{n^2}$ as the space of all $n \times n$ matrices. Let $K$ be the cone of all symmetric matrices that are positive semidefinite; see Definition 1.19. The corresponding generalized inequality on $\mathbb{R}^{n^2}$ is called the *positive semidefinite cone constraint*.

We would like to extend the notion of convex functions to vector valued maps, for which we will use the notion of generalized inequalities.

**Definition 4.19.** Let $\preccurlyeq_K$ be a generalized inequality on $\mathbb{R}^m$ induced by the cone $K$. We say that $G : \mathbb{R}^d \to \mathbb{R}^m$ is a *K-convex mapping* if

$$G(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) \preccurlyeq_K \lambda G(\mathbf{x}) + (1 - \lambda)G(\mathbf{y})$$

for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ and $\lambda \in (0, 1)$.

**Example 4.20.** Here are some examples of $K$-convex mappings.

1. Let $K \subseteq \mathbb{R}^m$ be any closed, convex, pointed cone. If $G : \mathbb{R}^d \to \mathbb{R}^m$ is an affine map, i.e., there exist a matrix $A \in \mathbb{R}^{m \times d}$ and a vector $\mathbf{b} \in \mathbb{R}^m$ such that $G(\mathbf{x}) = A\mathbf{x} + \mathbf{b}$, then $G$ is a $K$-convex mapping.

2. Let $m = n^2$ for some $n \in \mathbb{N}$, i.e., consider the space $\mathbb{R}^{n^2}$ and let $\preccurlyeq$ be the *positive semidefinite cone constraint* from part 3. of Example 4.18, i.e., induced by the cone $K$ of positive semidefinite matrices. Let $A_0, A_1, \ldots, A_d$ be fixed $p \times n$ matrices, for some $p \in \mathbb{N}$ (not necessarily equal to $n$). Define $G : \mathbb{R}^d \times \mathbb{R} \to \mathbb{R}^{n^2}$ to be the mapping

$$G(\mathbf{x}, s) = (A_0 + \mathbf{x}_1 A_1 + \ldots + \mathbf{x}_d A_d)^T (A_0 + \mathbf{x}_1 A_1 + \ldots + \mathbf{x}_d A_d) - s I_n,$$

where $I_n$ is the $n \times n$ identity matrix. Then $G$ is a $K$-convex mapping.

NOTES: 74

3. Let $K = \mathbb{R}^m_+$, and let $g_1, \ldots, g_m : \mathbb{R}^d \to \mathbb{R}$ be convex functions. Let $G : \mathbb{R}^d \to \mathbb{R}^m$ be defined as $G(\mathbf{x}) = (g_1(\mathbf{x}), \ldots, g_m(\mathbf{x}))$, then $G$ is a $K$-convex mapping.

## 4.3 Convex optimization with generalized inequalities

We can now define a very general framework for convex optimization problems, which is more concrete than the abstraction level of black-box first-order oracles, but is still flexible enough to incorporate the majority of convex optimization problems that show up in practice.

**Definition 4.21.** Let $f : \mathbb{R}^d \to \mathbb{R}$ be a convex function, let $K \subseteq \mathbb{R}^m$ be a closed, convex cone, and let $G : \mathbb{R}^d \to \mathbb{R}^m$ be a $K$-convex mapping. Then $f, K, G$ define a *convex optimization problem with generalized constraints* given as follows

$$\inf\{f(\mathbf{x}) \ : \ G(\mathbf{x}) \preccurlyeq_K \mathbf{0}\}. \tag{4.4}$$

Problem 3 in "HW for Week XI" shows that the set $C = \{\mathbf{x} \in \mathbb{R}^d : G(\mathbf{x}) \preccurlyeq_K \mathbf{0}\}$ is a convex set, when $G$ is a $K$-convex mapping. Thus, (4.4) is a special case of (4.1).

**Example 4.22.** Let us look at some concrete examples of (4.4).

1. **Linear/Quadratic Programming.** Let $f(\mathbf{x}) = \langle \mathbf{c}, \mathbf{x} \rangle$ for some $\mathbf{c} \in \mathbb{R}^d$, let $K = \mathbb{R}^m_+$ and let $G : \mathbb{R}^d \to \mathbb{R}^m$ be an affine map, i.e., $G(\mathbf{x}) = A\mathbf{x} - \mathbf{b}$ for some matrix $A \in \mathbb{R}^{m \times d}$ and a vector $\mathbf{b} \in \mathbb{R}^m$. Then (4.4) becomes

$$\inf\{\langle \mathbf{c}, \mathbf{x} \rangle \ : \ A\mathbf{x} \leq \mathbf{b}\}$$

   which is the problem of minimizing a linear function over a polyhedron. This is more commonly known as a *linear program*, in accordance with the fact that the objective and the constraints are all linear.

   If $f(\mathbf{x}) = \mathbf{x}^T Q \mathbf{x} + \langle \mathbf{c}, \mathbf{x} \rangle$ where $Q$ is a given $d \times d$ positive semidefinite matrix, and $\mathbf{c} \in \mathbb{R}^d$, then $f$ is a convex function (see Problem 14 from "HW for Week XI"). With $K$ and $G$ as above, (4.4) is called a *convex quadratic program*.

2. **Semidefinite Programming.** Let $m = n^2$ for some $n \in \mathbb{N}$ and consider the space $\mathbb{R}^{n^2}$. Let $f(\mathbf{x}) = \langle \mathbf{c}, \mathbf{x} \rangle$ for some $\mathbf{c} \in \mathbb{R}^d$, let $K \subseteq \mathbb{R}^{n^2}$ be the positive semidefinite cone, including the positive semidefinite cone constraint, and let $G : \mathbb{R}^d \to \mathbb{R}^{n^2}$ be an affine map, i.e., there exist $n \times n$ matrices $F_0, F_1, \ldots, F_d$ such that $G(\mathbf{x}) = F_0 + \mathbf{x}_1 F_1 + \ldots + \mathbf{x}_d F_d$. Then (4.4) becomes

$$\inf\{\langle \mathbf{c}, \mathbf{x} \rangle \ : \ -F_0 - \mathbf{x}_1 F_1 - \ldots - \mathbf{x}_d F_d \text{ is a PSD matrix}\}.$$

   This is known as a *semidefinite program*.

3. **Convex optimization with explicit constraints.** Let $f, g_1, \ldots, g_m : \mathbb{R}^d \to \mathbb{R}$ be convex functions. Define $K = \mathbb{R}^m_+$ and define $G : \mathbb{R}^d \to \mathbb{R}^m$ as $G(\mathbf{x}) = (g_1(\mathbf{x}), \ldots, g_m(\mathbf{x}))$, which is the $K$-convex mapping from Example 4.20. Then (4.4) becomes

$$\inf\{f(\mathbf{x}) \ : \ g_1(\mathbf{x}) \leq 0, \ldots, g_m(\mathbf{x}) \leq 0\}.$$

NOTES:

### 4.3.1  Lagrangian duality for convex optimization with generalized constraints

Given that the Subgradient Algorithm is a simple and elegant method for solving unconstrained problems, or problems with "simple" constraint sets $C$ (i.e., when one can compute $\mathrm{Proj}_C(\mathbf{x})$ efficiently), we will try to transform convex optimization problems with more complicated constraints into ones with simple constraints. This is the motivation for what is known as *Lagrangian duality*.

Note that problem (4.4) is equivalent to the problem

$$\inf_{\mathbf{x}\in\mathbb{R}^d} f(\mathbf{x}) + I_{-K}(G(\mathbf{x})), \tag{4.5}$$

where $I_{-K}$ is the indicator function for the cone $-K$. It can be shown that the function $I_{-K} \circ G$ is a convex function – see Problem 4 from "HW for Week XI". Thus, problem (4.5) is an unconstrained convex optimization problem. However, indicator functions are nasty to deal with because they are not finite valued, and thus, obtaining subgradient at all points becomes impossible. Thus, we try to replace $I_{-K}$ with a "nicer" penalty function $p : \mathbb{R}^m \to \mathbb{R}$, which is not that wildly discontinuous, and is finite-valued everywhere. So we would be looking at the problem

$$\inf_{\mathbf{x}\in\mathbb{R}^d} f(\mathbf{x}) + p(G(\mathbf{x})), \tag{4.6}$$

What properties should we require from our penalty function? First we would like problem (4.6) to be a convex problem, thus, we impose that

$$p \circ G : \mathbb{R}^d \to \mathbb{R} \text{ is a convex function.} \tag{4.7}$$

Next, from an optimization perspective, we would like to have guaranteed relationship between the function $f(\mathbf{x}) + I_{-K}(G(\mathbf{x}))$ and the function $f(\mathbf{x}) + p(G(\mathbf{x}))$. It turns out that a nice property to have is the guarantee that $f(\mathbf{x}) + p(G(\mathbf{x})) \le f(\mathbf{x}) + I_{-K}(G(\mathbf{x}))$ for all $\mathbf{x} \in \mathbb{R}^d$. This can be achieved by imposing that

$$p \text{ is an underestimator of } I_{-K}, \text{ i.e., } p \le I_{-K}. \tag{4.8}$$

Lagrangian duality theory is the study of penalty functions $p$ that are *linear* on $\mathbb{R}^m$, and satisfy the two conditions highlighted above. Now a function $p : \mathbb{R}^m \to \mathbb{R}$ is linear if and only if there exists $\mathbf{c} \in \mathbb{R}^m$ such that $p(\mathbf{z}) = \langle \mathbf{c}, \mathbf{z} \rangle$. The following proposition characterizes linear functions that satisfy the two conditions above.

**Proposition 4.23.** Let $p : \mathbb{R}^m \to \mathbb{R}$ be a linear function given by $p(\mathbf{z}) = \langle \mathbf{c}, \mathbf{z} \rangle$ for some $\mathbf{c} \in \mathbb{R}^m$. Then the following are equivalent:

1. $p$ satisfies condition (4.8).

2. $\mathbf{c} \in -K^\circ$, i.e., $-\mathbf{c}$ is in the polar of $K$.

NOTES: 76

3. $p$ satisfies conditions (4.7) and (4.8).

*Proof.* (1. $\implies$ 2.) Condition (4.8) is equivalent to saying that $p(\mathbf{z}) \leq 0$ for all $\mathbf{z} \in -K$, i.e.,

$$
\begin{aligned}
& \langle \mathbf{c}, \mathbf{z} \rangle \leq 0 && \text{for all } \mathbf{z} \in -K \\
\Leftrightarrow\ & \langle \mathbf{c}, -\mathbf{z} \rangle \leq 0 && \text{for all } \mathbf{z} \in K \\
\Leftrightarrow\ & \langle -\mathbf{c}, \mathbf{z} \rangle \leq 0 && \text{for all } \mathbf{z} \in K \\
\Leftrightarrow\ & -\mathbf{c} \in K^\circ \\
\Leftrightarrow\ & \mathbf{c} \in -K^\circ
\end{aligned}
$$

(2. $\implies$ 3.) We showed above that assuming $\mathbf{c} \in -K^\circ$ is equivalent to condition (4.8). We now check that $\mathbf{c} \in -K^\circ$ implies (4.7). Since $G$ is a $K$-convex mapping, we have that for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ and $\lambda \in (0,1)$,

$$
\begin{aligned}
& \langle \mathbf{c}, \lambda G(\mathbf{x}) + (1-\lambda)G(\mathbf{y}) - G(\lambda \mathbf{x} + (1-\lambda)\mathbf{y}) \rangle \geq 0 \\
\implies\ & \langle \mathbf{c}, \lambda G(\mathbf{x}) \rangle + \langle \mathbf{c}, (1-\lambda)G(\mathbf{y}) \rangle \geq \langle \mathbf{c}, G(\lambda \mathbf{x} + (1-\lambda)\mathbf{y}) \rangle \\
\implies\ & \lambda \langle \mathbf{c}, G(\mathbf{x}) \rangle + (1-\lambda)\langle \mathbf{c}, G(\mathbf{y}) \rangle \geq \langle \mathbf{c}, G(\lambda \mathbf{x} + (1-\lambda)\mathbf{y}) \rangle \\
\implies\ & \lambda p(G(\mathbf{x})) + (1-\lambda)p(G(\mathbf{y})) \geq p(G(\lambda \mathbf{x} + (1-\lambda)\mathbf{y}))
\end{aligned}
$$

Hence, condition (4.7) is satisfied.

(3. $\implies$ 1.) Trivial. $\qquad\square$

**Definition 4.24.** The set $-K^\circ$ is important in Lagrangian duality, and a separate notation and name has been invented: $-K^\circ$ is called the *dual cone* of $K$ and is denoted by $K^\star$.

The above discussions show that for any $\mathbf{y} \in K^\star$, the optimal value of the (4.6), with $p$ given by $p(\mathbf{z}) = \langle \mathbf{y}, \mathbf{z} \rangle$, is a lower bound on the optimal value of (4.4). This motivates definition of the so-called *dual function* $\mathcal{L} : \mathbb{R}^m \to \mathbb{R}$ associated with (4.4) as follows:

$$
\mathcal{L}(\mathbf{y}) := \inf_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + \langle \mathbf{y}, G(\mathbf{x}) \rangle \tag{4.9}
$$

We state the lower bound property formally.

**Proposition 4.25.** [Weak Duality] Let $f : \mathbb{R}^d \to \mathbb{R}$ be convex, let $K \subseteq \mathbb{R}^m$ be a closed, convex cone, and let $G : \mathbb{R}^d \to \mathbb{R}^m$ be a $K$-convex mapping. Let $\mathcal{L} : \mathbb{R}^m \to \mathbb{R}$ be as defined in (4.9). Then, for all $\bar{\mathbf{x}} \in \mathbb{R}^d$ such that $G(\bar{\mathbf{x}}) \preccurlyeq_K$ and all $\bar{\mathbf{y}} \in K^\star$, we must have $\mathcal{L}(\bar{\mathbf{y}}) \leq f(\bar{\mathbf{x}})$. Consequently, $\mathcal{L}(\bar{\mathbf{y}}) \leq \inf\{f(\mathbf{x}) \ : \ G(\mathbf{x}) \preccurlyeq_K \mathbf{0}\}$.

*Proof.* We simply follow the inequalities

$$
\begin{aligned}
\mathcal{L}(\bar{\mathbf{y}}) &= \inf_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + \langle \mathbf{y}, G(\mathbf{x}) \rangle \\
&\leq f(\bar{\mathbf{x}}) + \langle \mathbf{y}, G(\bar{\mathbf{x}}) \rangle \\
&\leq f(\bar{\mathbf{x}}),
\end{aligned}
$$

where the last inequality holds because $G(\mathbf{x}) \preccurlyeq_K$ and $\bar{\mathbf{y}} \in K^\star$, and so $\langle \mathbf{y}, G(\bar{\mathbf{x}}) \rangle \leq 0$. $\qquad\square$

NOTES: 77

Proposition 4.25 shows that any $\mathbf{y} \in K^\star$ provides the lower bound $\mathcal{L}(\mathbf{y})$ on the optimal value of the optimization problem (4.4). The *Lagrangian dual optimization problem* is the problem of finding the $\mathbf{y} \in K^\star$ that provides the *best/largest* lower bound. In other words, the Lagrangian dual problem is defined as

$$\sup_{\mathbf{y} \in K^\star} \mathcal{L}(\mathbf{y}), \tag{4.10}$$

and Proposition 4.25 can be restated as

$$\sup\{\mathcal{L}(\mathbf{y}) : \mathbf{y} \in K^\star\} \;\; \leq \;\; \inf\{f(\mathbf{x}) \;:\; G(\mathbf{x}) \preccurlyeq_K \mathbf{0}\}. \tag{4.11}$$

If we have equality in (4.11), then to solve (4.4), one can instead solve (4.10). This merits a definition.

**Definition 4.26** (Strong Duality)**.** We say that we have a *zero duality gap* if equality holds in (4.11). In addition, if the supremum in (4.10) is attained for some $\mathbf{y} \in K^\star$, then we say that *strong duality* holds.

### 4.3.2 Solving the Lagrangian dual problem

Before we investigate conditions under which we have zero duality gap or strong duality, let us try to see how one use the subgradient algorithm to solve (4.10).

**Proposition 4.27.** $\mathcal{L}(\mathbf{y})$ is a concave function of $\mathbf{y}$.

*Proof.* We have to show that $-\mathcal{L}(\mathbf{y})$ is a convex function of $\mathbf{y}$. This follows from the fact that

$$
\begin{aligned}
-\mathcal{L}(\mathbf{y}) &= -\inf_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + \langle \mathbf{y}, G(\mathbf{x}) \rangle \\
&= \sup_{\mathbf{x} \in \mathbb{R}^d} -f(\mathbf{x}) + \langle \mathbf{y}, -G(\mathbf{x}) \rangle,
\end{aligned}
$$

i.e., $-\mathcal{L}(\mathbf{y})$ is the supremum of affine functions of $\mathbf{y}$ of the form $-f(\mathbf{x}) + \langle \mathbf{y}, -G(\mathbf{x}) \rangle$. By part 2. of Theorem 3.12, $-\mathcal{L}(\mathbf{y})$ is convex in $\mathbf{y}$. $\qquad\square$

We could now use the subgradient algorithm to solve (4.10), if we had a first order oracle for $\mathcal{L}(\mathbf{y})$ and an algorithm to project to $K^\star$. We show that a subgradient for $-\mathcal{L}(\mathbf{y})$ can be found by solving an unconstrained convex optimization problem.

**Proposition 4.28.** Let $\bar{\mathbf{y}} \in \mathbb{R}^m$ and let $\bar{\mathbf{x}} \in \arg\inf_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + \langle \bar{\mathbf{y}}, G(\mathbf{x}) \rangle$. Then $-G(\bar{\mathbf{x}}) \in \partial(-\mathcal{L})(\bar{\mathbf{y}})$.

*Proof.* We express $-\mathcal{L}(\mathbf{y}) = \sup_{\mathbf{x} \in \mathbb{R}^d} -f(\mathbf{x}) + \langle \mathbf{y}, -G(\mathbf{x}) \rangle$ as the supremum of affine functions, and use part 3. of Theorem 3.62, and the fact that the subdifferential of the affine function $-f(\bar{\mathbf{x}}) + \langle \mathbf{y}, -G(\bar{\mathbf{x}}) \rangle$, at $\bar{\mathbf{y}}$ is simply $-G(\bar{\mathbf{x}})$. $\qquad\square$

Now if we have an algorithm that can compute $\mathrm{Proj}_{K^\star}(\mathbf{y})$ for all $\mathbf{y} \in \mathbb{R}^m$, then using Propositions 4.27 and 4.28, one can solve the Lagrangian dual problem (4.10), where in each iteration of the algorithm, one solves the unconstrained problem $\inf_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + \langle \bar{\mathbf{y}}, G(\mathbf{x}) \rangle$ for a given $\bar{\mathbf{y}} \in K^\star$. This can, in turn, be solved by the subgradient algorithm if one has the appropriate first order oracles for $f(\mathbf{x})$ and $\langle \bar{\mathbf{y}}, G(\mathbf{x}) \rangle$.

NOTES: 78

### 4.3.3 Explicit examples of the Lagrangian dual

We will now explore some special settings of convex optimization problems with generalized inequalities, and see that the Lagrangian dual has a particularly nice form.

**Conic optimization.** Let $K \subseteq \mathbb{R}^m$ be a closed, convex, pointed cone. Let $G : \mathbb{R}^d \to \mathbb{R}^m$ be an affine map given by $G(\mathbf{x}) = A\mathbf{x} - \mathbf{b}$, where $A \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$. Let $f : \mathbb{R}^d \to \mathbb{R}$ be a linear function given by $f(\mathbf{x}) = \langle \mathbf{c}, \mathbf{x} \rangle$ for some $\mathbf{c} \in \mathbb{R}^d$. Then Problem 4.4 becomes

$$\inf\{\langle \mathbf{c}, \mathbf{x} \rangle : A\mathbf{x} \preccurlyeq_K \mathbf{b}\}. \tag{4.12}$$

For a fixed cone $K$, problems of the form (4.12) with are called *conic optimization problems over the cone $K$*. As we pick different data $A, \mathbf{b}, \mathbf{c}$, we get different instances of a conic optimization problem over the cone $K$. A special case is when $K = \mathbb{R}^m_+$, which is known as *linear programming or linear optimization* – see Example 4.22 – which is the problem of optimizing a linear function over a polyhedron.

Let us investigate the dual function of (4.12). Recall that $\mathcal{L}(\mathbf{y}) = \inf_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + \langle \mathbf{y}, G(\mathbf{x}) \rangle$, which in this case becomes

$$
\begin{aligned}
\inf_{\mathbf{x} \in \mathbb{R}^d} \langle \mathbf{c}, \mathbf{x} \rangle + \langle \mathbf{y}, A\mathbf{x} - \mathbf{b} \rangle &= \inf_{\mathbf{x} \in \mathbb{R}^d} \langle \mathbf{c}, \mathbf{x} \rangle + \langle \mathbf{y}, A\mathbf{x} \rangle - \langle \mathbf{y}, \mathbf{b} \rangle \\
&= \inf_{\mathbf{x} \in \mathbb{R}^d} \langle \mathbf{c}, \mathbf{x} \rangle + \langle A^T\mathbf{y}, \mathbf{x} \rangle - \langle \mathbf{y}, \mathbf{b} \rangle \\
&= \inf_{\mathbf{x} \in \mathbb{R}^d} \langle \mathbf{c} + A^T\mathbf{y}, \mathbf{x} \rangle - \langle \mathbf{y}, \mathbf{b} \rangle.
\end{aligned}
$$

Now, if $\mathbf{c} + A^T\mathbf{y} \neq \mathbf{0}$, then the infimum above is clearly $-\infty$. And if $\mathbf{c} + A^T\mathbf{y} = \mathbf{0}$, then the infimum is $-\langle \mathbf{b}, \mathbf{y} \rangle$. Therefore, for (4.12), the dual function is given by

$$
\mathcal{L}(\mathbf{y}) = \begin{cases} -\infty & \mathbf{c} + A^T\mathbf{y} \neq \mathbf{0} \\ -\langle \mathbf{b}, \mathbf{y} \rangle & \mathbf{c} + A^T\mathbf{y} = \mathbf{0} \end{cases} \tag{4.13}
$$

Therefore,

$$\sup_{\mathbf{y} \in K^\star} \mathcal{L}(\mathbf{y}) = \sup\{-\langle \mathbf{b}, \mathbf{y} \rangle : A^T\mathbf{y} = -\mathbf{c}, \ \mathbf{y} \in K^\star\} = -\inf\{\langle \mathbf{b}, \mathbf{y} \rangle : A^T\mathbf{y} = -\mathbf{c}, \ \mathbf{y} \in K^\star\}.$$

To remove the slightly annoying minus sign in front of $\mathbf{c}$ above, it is more standard to write (4.12) as $-\sup\{\langle -\mathbf{c}, \mathbf{x} \rangle : A\mathbf{x} \preccurlyeq_K \mathbf{b}\}$, and then replace $-\mathbf{c}$ with $\mathbf{c}$ throughout the above derivation. Thus, the standard primal dual pairs for conic optimization problems are

$$\sup\{\langle \mathbf{c}, \mathbf{x} \rangle : A\mathbf{x} \preccurlyeq_K \mathbf{b}\} \ \leq \ \inf\{\langle \mathbf{b}, \mathbf{y} \rangle : A^T\mathbf{y} = \mathbf{c}, \ \mathbf{y} \in K^\star\}. \tag{4.14}$$

*Linear Programming/Optimization.* Specializing to the linear programming case with $K = \mathbb{R}^m_+$ and observing that $K^\star = K = \mathbb{R}^m_+$ (see Problem 2 from "HW for Week III"), we obtain the primal dual pair

$$\sup\{\langle \mathbf{c}, \mathbf{x} \rangle : A\mathbf{x} \leq \mathbf{b}\} \ \leq \ \inf\{\langle \mathbf{b}, \mathbf{y} \rangle : A^T\mathbf{y} = \mathbf{c}, \ \mathbf{y} \geq \mathbf{0}\}. \tag{4.15}$$

NOTES:                                    79

*Semidefinite Programming/Optimization.* Another special case is that of semidefinite optimization. This is the situation when $m = n^2$ and $K$ is the cone of positive semidefinite matrices. $G : \mathbb{R}^d \to \mathbb{R}^{n^2}$ is an affine map from $\mathbb{R}^d$ to the space of $n \times n$ matrices. To avoid dealing with asymmetric matrices, $G$ is always assumed to be of the form $G(\mathbf{x}) = \mathbf{x}_1 A_1 + \ldots + \mathbf{x}_d A_d - A_0$, where $A_0, A_1, \ldots, A_d$ are $n \times n$ symmetric matrices[4]. If one works though the algebra in this case and uses the fact that the positive semidefinite cone is *self-dual,* *i.e.,* $K = K^\star$, (4.14) becomes

$$\sup\{\langle \mathbf{c}, \mathbf{x} \rangle : \mathbf{x}_1 A_1 + \ldots + \mathbf{x}_d A_d - A_0 \text{ is a PSD matrix}\} \quad \leq \quad \inf\{\langle A_0, Y \rangle : \langle A_i, Y \rangle = \mathbf{c}_i, \ Y \text{ is a PSD matrix}\},$$

where $\langle X, Z \rangle = \sum_{i,j} X_{ij} Z_{ij}$ for any pair $X, Z$ of $n \times n$ symmetric matrices.

**Convex optimization with explicit constraints and objective.** Recall part 3. if Example 4.22, where $K = \mathbb{R}_+^m$, $f, g_1, \ldots, g_m : \mathbb{R}^d \to \mathbb{R}$ are convex functions, and $G : \mathbb{R}^d \to \mathbb{R}^m$ as $G(\mathbf{x}) = (g_1(\mathbf{x}), \ldots, g_m(\mathbf{x}))$, giving the explicit problem

$$\inf\{f(\mathbf{x}) \ : \ g_1(\mathbf{x}) \leq 0, \ldots, g_m(\mathbf{x}) \leq 0\}.$$

In this case, since $K^\star = K = \mathbb{R}_+^m$ (see Problem 2 from "HW for Week III"), the dual problem is

$$\sup_{\mathbf{y} \in K^\star} \mathcal{L}(\mathbf{y}) = \sup_{\mathbf{y} \geq \mathbf{0}} \inf_{\mathbf{x} \in \mathbb{R}^d} \{f(\mathbf{x}) + \mathbf{y}_1 g_1(\mathbf{x}) + \ldots, \mathbf{y}_m g_m(\mathbf{x})\}.$$

**A closer look at linear programming duality.** Consider the following linear program:

$$\begin{array}{rrcl} \sup & 2x_1 - 1.5x_2 & & \\ & x_1 + x_2 & \leq & 1 \\ & x_1 - x_2 & \leq & 1 \\ & -x_1 + x_2 & \leq & 1 \\ & -x_1 - x_2 & \leq & 1 \end{array} \tag{4.16}$$

To solve this problem, let us make some simple observations. If we multiply the first inequality by 0.5, the second inequality by 3.5, the third by 1.75 and the fourth by 0.25 and add all these scaled inequalities, then we obtain the inequality $2x_1 - 1.5x_2 \leq 6$. Now any $\mathbf{x} \in \mathbb{R}^2$ satisfying the constraints of the above linear program must also satisfy this new inequality. This shows that our supremum is *at most* 6. Now if we choose another set of multipliers : $0.25, 1.75, 0, 0$ (in order), then we obtain the inequality $2x_1 - 1.5x_2 \leq 2$, which gives a better bound of $2 \leq 6$ on the optimal solution value. Now, consider the point $x_1 = 1, x_2 = 0$: this have value $2 \cdot 1 - 1.5 \cdot 0 = 2$. Since we have an upper bound of 2 from the above arguments, we know that $x_1 = 1, x_2 = 0$ is actually the optimal solution to the above linear program! Thus, we have provided the optimal solution, and a quick certificate of its optimality. If you think about how we were deriving the upper bounds of 6 and 2, we were looking for nonnegative multipliers $y_1, y_2, y_3, y_4$ such that the corresponding

---

[4]Dealing with asymmetric matrices is not hard, but involves little details that can be overlooked for this exposition, and don't provide any great insight.

NOTES: 80

combination of the inequalities gives us $2x_1 - 1.5x_2$ on the left hand side, and the upper bound was simply the right hand side of the combined inequality, which is, $y_1 + y_2 + y_3 + y_4$. If the right hand side is to end up as $2x_1 - 1.5x_2$, then we must have $y_1 + y_2 - y_3 - y_4 = 2$ and $y_1 - y_2 + y_3 - y_4 = -1.5$. To get the best upper bound, we want to find the minimum value of $y_1 + y_2 + y_3 + y_4$ such that $y_1 + y_2 - y_3 - y_4 = 2$ and $y_1 - y_2 + y_3 - y_4 = -1.5$, and all $y_i$'s are nonnegative. But this is exactly the dual problem in (4.15). We hope this gives the reader a more "hands-on" perspective on the Lagrangian dual of a linear program.

### 4.3.4  Strong duality: sufficient conditions and complementary slackness

In the above example of the linear program in (4.16), it turned out that we could find a primal feasible solution and a dual feasible solution that have the same value, which shows that we have strong duality, and certifies the optimality of the two solutions. We will see below that this always happens for linear programs. For general conic optimization problems, or a convex optimization problem with generalized inequalities, this does not always hold and one may not even have zero duality gap. We now supply two conditions under which strong duality is obtained. Linear programming strong duality will be a special case of the second condition.

**Slater's condition for strong duality.**   The following is perhaps the most well-known sufficient condition in convex optimization that guarantees strong duality.

**Theorem 4.29.** [Slater's condition] Let $f : \mathbb{R}^d \to \mathbb{R}$ be convex, let $K \subseteq \mathbb{R}^m$ be a closed, convex cone, and let $G : \mathbb{R}^d \to \mathbb{R}^m$ be a $K$-convex mapping. Let $\mathcal{L} : \mathbb{R}^m \to \mathbb{R}$ be as defined in (4.9). If there exists $\bar{\mathbf{x}}$ such that $-G(\bar{\mathbf{x}}) \in \text{int}(K)$ and $\inf\{f(\mathbf{x}) : G(\mathbf{x}) \preccurlyeq_K \mathbf{0}\}$ is a finite value, then there exists $\mathbf{y}^\star \in K^\star$ such that $\sup_{\mathbf{y} \in K^\star} \mathcal{L}(\mathbf{y}) = \mathcal{L}(\mathbf{y}^\star) = \inf\{f(\mathbf{x}) : G(\mathbf{x}) \preccurlyeq_K \mathbf{0}\}$, i.e., strong duality holds.

Before we begin the proof, we need to establish a slight variant of the separating hyperplane theorem, that does not make any closedness or compactness assumptions.

**Proposition 4.30.** Let $A, B \subseteq \mathbb{R}^d$ be convex sets (not necessarily closed) such that $A \cap B = \emptyset$. Then there exist $\mathbf{a} \in \mathbb{R}^d$ and $\delta \in \mathbb{R}$ such that $\langle \mathbf{a}, \mathbf{x} \rangle \geq \langle \mathbf{a}, \mathbf{y} \rangle$ for all $\mathbf{x} \in A, \mathbf{y} \in B$.

*Proof.* Left as an exercise. $\qquad\square$

*Proof of Theorem 4.29.* Let $\mu_0 = \inf\{f(\mathbf{x}) : G(\mathbf{x}) \preccurlyeq_K \mathbf{0}\} < \infty$. Define the sets

$$A = \{(\mathbf{z}, r) \in \mathbb{R}^m \times \mathbb{R} : \exists \mathbf{x} \in \mathbb{R}^d \text{ such that } f(\mathbf{x}) \leq r, \ G(\mathbf{x}) \preccurlyeq_K \mathbf{z}\},$$
$$B = \{(\mathbf{z}, r) \in \mathbb{R}^m \times \mathbb{R} : r < \mu_0, \ \mathbf{z} \preccurlyeq_K \mathbf{0}\}.$$

It is not hard to verify that $A, B$ are convex. Moreover, since $\mu_0 = \inf\{f(\mathbf{x}) : G(\mathbf{x}) \preccurlyeq_K \mathbf{0}\} < \infty$, it is also not hard to verify that $A \cap B = \emptyset$. By Proposition 4.30, there exists $\mathbf{a} \in \mathbb{R}^m, \gamma \in \mathbb{R}$ such that

$$\langle \mathbf{a}, \mathbf{z}_1 \rangle + \gamma r_1 \geq \langle \mathbf{a}, \mathbf{z}_2 \rangle + \gamma r_2 \tag{4.17}$$

for all $(\mathbf{z}_1, r_1) \in A$ and $(\mathbf{z}_2, r_2) \in B$.

NOTES:                                                  81

**Claim 3.** $\mathbf{a} \in K^\star$ and $\gamma \geq 0$.

*Proof of Claim.* Suppose the contrary that $\mathbf{a} \notin K^\star$. Then $\mathbf{a} \notin -K^\circ = (-K)^\circ$. Thus, there exists $\bar{\mathbf{z}} \in -K$, i.e., $\bar{\mathbf{z}} \preceq_K \mathbf{0}$, such that $\langle \mathbf{a}, \bar{\mathbf{z}} \rangle > 0$. Now (4.17) holds with $\mathbf{z}_1 = G(\bar{\mathbf{x}})$ ($\bar{\mathbf{x}}$ is the point in the hypothesis of the theorem), $r_1 = f(\bar{\mathbf{x}})$, $r_2 = \mu_0 - 1$ and $z_2 = \lambda \bar{\mathbf{z}}$ for all $\lambda \geq 0$. But since $\langle \mathbf{a}, \bar{\mathbf{z}} \rangle > 0$, the inequality (4.17) would be violated for large enough $\lambda$. Thus, we must have $\mathbf{a} \in K^\star$.

Similarly, (4.17) holds with $\mathbf{z}_1 = G(\bar{\mathbf{x}})$, $r_1 = f(\bar{\mathbf{x}})$, $z_2 = \bar{\mathbf{z}}$ and $\underline{all}$ $r_2 < \mu_0$. If $\gamma < 0$, then letting $r_2 \to -\infty$ would violate (4.17). $\qquad\square$

We now show that, in fact, $\gamma > 0$. Substitute $\mathbf{z}_1 = G(\bar{\mathbf{x}})$ ($\bar{\mathbf{x}}$ is the point in the hypothesis of the theorem), $r_1 = f(\bar{\mathbf{x}})$, $r_2 = \mu_0 - 1$ and $z_2 = \mathbf{0}$ in (4.17). If $\gamma = 0$, then this relation becomes

$$\langle \mathbf{a}, G(\bar{\mathbf{x}}) \rangle \geq \mathbf{0}.$$

However, since $-G(\bar{\mathbf{x}}) \in \text{int}(K)$ and therefore, $\langle \mathbf{a}, G(\bar{\mathbf{x}}) \rangle < 0$ (see Problem 3 from "HW for Week II"). By Claim 3, $\gamma > 0$.

Let $\mathbf{y}^\star := \frac{\mathbf{a}}{\gamma}$; by Claim 3, $\mathbf{y}^\star \in K^\star$. We will now show that for every $\epsilon > 0$, $\mathcal{L}(\mathbf{y}^\star) \geq \mu_0 - \epsilon$. This will establish the result because this means $\mathcal{L}(\mathbf{y}^\star) \geq \mu_0$ and since $\mathcal{L}(\mathbf{y}) \leq \mu_0$ for all $\mathbf{y} \in K^\star$ by Proposition 4.25, we must have $\sup_{\mathbf{y} \in K^\star} \mathcal{L}(\mathbf{y}) = \mathcal{L}(\mathbf{y}^\star) = \mu_0$. Consider any $\mathbf{x} \in \mathbb{R}^d$. $\mathbf{z}_1 = G(\mathbf{x})$ and $r_1 = f(\mathbf{x})$ gives a point in $A$. Substituting into (4.17) with $\mathbf{z}_2 = 0$ and $r_2 = \mu_0 - \epsilon$, we obtain that $\langle \mathbf{a}, G(\mathbf{x}) \rangle + \gamma f(\mathbf{x}) \geq \gamma(\mu_0 - \epsilon)$. Dividing through by $\gamma$, we obtain

$$\langle \mathbf{y}^\star, G(\mathbf{x}) \rangle + f(\mathbf{x}) \geq \mu_0 - \epsilon.$$

This implies that $\mathcal{L}(\mathbf{y}^\star) = \inf_{\mathbf{x} \in \mathbb{R}^d} \langle \mathbf{y}^\star, G(\mathbf{x}) \rangle + f(\mathbf{x}) \geq \mu_0 - \epsilon$. $\qquad\square$

**Closed cone condition for strong duality in conic optimization.** Slater's condition applied to conic optimization problems translates into requiring that there is some $\bar{\mathbf{x}}$ such that $\mathbf{b} - A\bar{\mathbf{x}} \in \text{int}(K)$. Another very useful strong duality condition uses topological properties of the dual cone $K^*$.

**Theorem 4.31.** [Closed cone condition] Consider the conic optimization primal dual pair (4.14). Suppose the set $\{(A^T \mathbf{y}, \langle \mathbf{b}, \mathbf{y} \rangle) \in \mathbb{R}^d \times \mathbb{R} : \mathbf{y} \in K^*\}$ is closed and the dual is feasible, i.e., there exists $\mathbf{y} \in K^\star$ such that $A^T \mathbf{y} = \mathbf{c}$. Then we have zero duality gap. If the optimal dual value is finite, then strong duality holds in (4.14).

*Proof.* Since the dual is feasible, its optimal value is either $-\infty$ or finite. By weak duality (Proposition 4.25), in the first case we must have zero duality gap and the primal is infeasible. So we consider the case when the optimal value of the dual is finite, say $\mu_0 \in \mathbb{R}$. Let us label the set $S := \{(A^T \mathbf{y}, \langle \mathbf{b}, \mathbf{y} \rangle) : \mathbf{y} \in K^*\} \subseteq \mathbb{R}^d \times \mathbb{R}$. Notice that the optimal value of the dual is $\mu_0 = \inf\{r \in \mathbb{R} : (\mathbf{c}, r) \in S\}$. Since $S$ is closed, the set $\{r \in \mathbb{R} : (\mathbf{c}, r) \in S\}$ is closed because it is topologically the same as $S \cap (\mathbf{c} \times \mathbb{R})$. Therefore the infimum in $\inf\{r \in \mathbb{R} : (\mathbf{c}, r) \in S\}$ is over a closed subset of the real line. Hence, $(\mathbf{c}, \mu_0) \in S$ and so there exists $\mathbf{y}^\star \in K^\star$ such that $A^T \mathbf{y}^\star = \mathbf{c}$ and $\langle \mathbf{b}, \mathbf{y}^\star \rangle = \mu_0$.

NOTES: 82

Since $\mu_0 = \inf\{r \in \mathbb{R} : (\mathbf{c}, r) \in S\}$, for every $\epsilon > 0$, $(\mathbf{c}, \mu_0 - \epsilon) \notin S$. Therefore, there exists a separating hyperplane $(\mathbf{a}, \gamma) \in \mathbb{R}^d \times \mathbb{R}$ and $\delta \in \mathbb{R}$ such that $\langle \mathbf{a}, A^T\mathbf{y}\rangle + \gamma \cdot \langle \mathbf{b}, \mathbf{y}\rangle \leq \delta$ for all $\mathbf{y} \in K^\star$, and $\langle \mathbf{a}, \mathbf{c}\rangle + \gamma(\mu_0 - \epsilon) > \delta$. By Problem 8 from "HW for Week IX", we may assume $\delta = 0$. Therefore, we have

$$\langle \mathbf{a}, A^T\mathbf{y}\rangle + \gamma \cdot \langle \mathbf{b}, \mathbf{y}\rangle \leq 0 \text{ for all } \mathbf{y} \in K^\star, \tag{4.18}$$

$$\langle \mathbf{a}, \mathbf{c}\rangle + \gamma(\mu_0 - \epsilon) > 0 \tag{4.19}$$

Substituting $\mathbf{y}^\star$ in (4.18), we obtain that $\langle \mathbf{a}, \mathbf{c}\rangle + \gamma\mu_0 \leq 0$, and (4.19) tells us that $\langle \mathbf{a}, \mathbf{c}\rangle + \gamma\mu_0 > \gamma\epsilon$. This implies that $\gamma < 0$ since $\epsilon > 0$. Now (4.18) can be rewritten as $\langle A\mathbf{a} + \gamma\mathbf{b}, \mathbf{y}\rangle \leq 0$ for all $\mathbf{y} \in K^\star$ and (4.19) can be rewritten as $\langle \mathbf{a}, \mathbf{c}\rangle > -\gamma(\mu_0 - \epsilon)$. Dividing through both these relations by $-\gamma > 0$, and setting $\mathbf{x} = \frac{\mathbf{a}}{-\gamma}$, we obtain that $\langle A\mathbf{x} - \mathbf{b}, \mathbf{y}\rangle \leq 0$ for all $\mathbf{y} \in K^\star$ implying that $A\mathbf{x} \preccurlyeq_K \mathbf{b}$, and $\langle \mathbf{x}, \mathbf{c}\rangle > \mu_0 - \epsilon$. Thus, we have a feasible solution $\mathbf{x}$ for the primal with value at least $\mu_0 - \epsilon$. Since $\epsilon > 0$ was chosen arbitrarily, this shows that for every $\epsilon > 0$, the primal has optimal value better than $\mu_0 - \epsilon$. Therefore, the primal value must be $\mu_0$ and we have zero duality gap. The existence of $y^\star$ shows that we have strong duality. □

*Linear Programming strong duality.* The closed cone condition for strong duality implies that linear programs always enjoy strong duality when either the primal or the dual (or both) are feasible. This is because the cone $K = \mathbb{R}^m_+$ is a polyhedral cone and also self-dual, i.e., $K^\star = K = \mathbb{R}^M_+$. Since linear transformations of polyhedral cones are polyhedral (see part 5. of Problem 1 in "HW for Week V"), and hence closed, linear programs always satisfy the condition in Theorem 4.31. One therefore has the following table for the possible outcomes in the primal-dual linear programming pair.

| Primal \ Dual | Infeasible | Finite | Unbounded |
|---|---|---|---|
| Infeasible | Possible | Impossible | Possible |
| Finite | Impossible | Possible, Zero duality gap | Impossible |
| Unbounded | Possible | Impossible | Impossible |

An alternate proof of zero duality gap for linear programming follows from our results on polyhedral theory. We outline it here to illustrate that linear programming duality can be approached in different ways (although ultimately both proofs go back to the separating hyperplane theorem – Theorem 2.20). We consider two cases:

*Primal is infeasible.* In this case, we will show that is the dual is feasible, then the dual must be unbounded. Since the primal is infeasible, the polyhedron $A\mathbf{x} \leq \mathbf{b}$ is empty. By Theorem 2.88, there exists $\hat{\mathbf{y}} \geq \mathbf{0}$ such that $A^T\hat{\mathbf{y}} = \mathbf{0}$ and $\langle \mathbf{b}, \hat{\mathbf{y}}\rangle = -1$. Since the dual is feasible, consider any $\bar{\mathbf{y}} \geq 0$ such that $A^T\bar{\mathbf{y}} = \mathbf{c}$. Now, all points of the form $\bar{\mathbf{y}} + \lambda\hat{\mathbf{y}}$ are also feasible to the dual, and the corresponding value $\langle \mathbf{b}, \bar{\mathbf{y}} + \lambda\hat{\mathbf{y}}\rangle$ can be made to go to $-\infty$ because $\langle \mathbf{b}, \hat{\mathbf{y}}\rangle = -1$.

NOTES: 83

$1753$ *Primal is feasible.* If the primal is unbounded, then by weak duality, the dual must be infeasible. So let
$1754$ us consider the case that the primal has a finite value $\mu_0$. This means that the inequality $\langle \mathbf{c}, \mathbf{x} \rangle \leq \mu_0$ is a
$1755$ valid inequality for the polyhedron $A\mathbf{x} \leq \mathbf{b}$. By Theorem 2.85, there exists $\hat{\mathbf{y}} \geq 0$ such that $A^T \hat{\mathbf{y}} = \mathbf{c}$ and
$1756$ $\langle \mathbf{b}, \hat{\mathbf{y}} \rangle \leq \mu_0$. Therefore the dual has a solution $\hat{\mathbf{y}}$ whose objective value is equal to the primal value $\mu_0$. This
$1757$ guarantees strong duality.

$1758$ **Complementary slackness.** Complementary slackness is a useful necessary condition when we have
$1759$ primal and dual optimal solutions with zero duality gap.

$1760$ **Theorem 4.32.** Let $f : \mathbb{R}^d \to \mathbb{R}$ be convex, let $K \subseteq \mathbb{R}^m$ be a closed, convex cone, and let $G : \mathbb{R}^d \to \mathbb{R}^m$ be
$1761$ a $K$-convex mapping. Let $\mathcal{L} : \mathbb{R}^m \to \mathbb{R}$ be as defined in (4.9). Let $\mathbf{x}^\star$ be such that $G(\mathbf{x}^\star) \preccurlyeq_K \mathbf{0}$ and $\mathbf{y}^\star \in K^\star$
$1762$ such that $f(\mathbf{x}^\star) = \mathcal{L}(\mathbf{y}^\star)$. Then $\langle \mathbf{y}^\star, G(\mathbf{x}^\star) \rangle = 0$.

*Proof.* We simply observe that since $G(\mathbf{x}^\star) \preccurlyeq_K \mathbf{0}$ and $\mathbf{y}^\star \in K^\star$, we must have $\langle \mathbf{y}^\star, G(\mathbf{x}^\star) \rangle \leq 0$. Therefore,

$$f(\mathbf{x}^\star) \geq f(\mathbf{x}^\star) + \langle \mathbf{y}^\star, G(\mathbf{x}^\star) \rangle \geq \inf_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + \langle \mathbf{y}^\star, G(\mathbf{x}) \rangle = \mathcal{L}(\mathbf{y}^\star).$$

$1763$ Since $f(\mathbf{x}^\star) = \mathcal{L}(\mathbf{y}^\star)$ by assumption, equality must hold throughout above giving us $\langle \mathbf{y}^\star, G(\mathbf{x}^\star) \rangle = 0$.  □

$1764$ ### 4.3.5 Saddle point interpretation of the Lagrangian dual

$1765$ Let us go back to the original problem (4.4) and revisit the dual function $\mathcal{L}(\mathbf{y})$. Define the function

$$\hat{\mathcal{L}}(\mathbf{x}, \mathbf{y}) := f(\mathbf{x}) + \langle \mathbf{y}, G(\mathbf{x}) \rangle \tag{4.20}$$

$1766$ which is often called the *Lagrangian function* associated with (4.4). A characterization of a pair of optimal
$1767$ solutions to (4.4) and (4.10) can be obtained using saddle points of the Lagrangian function.

$1768$ **Theorem 4.33.** Let $f : \mathbb{R}^d \to \mathbb{R}$ be convex, let $K \subseteq \mathbb{R}^m$ be a closed, convex cone, and let $G : \mathbb{R}^d \to \mathbb{R}^m$ be
$1769$ a $K$-convex mapping. Let $\mathcal{L} : \mathbb{R}^m \to \mathbb{R}$ be as defined in (4.9) and $\hat{\mathcal{L}} : \mathbb{R}^d \times \mathbb{R}^m \to \mathbb{R}$ be as defined in (4.20).
$1770$ Let $\mathbf{x}^\star$ be such that $G(\mathbf{x}^\star) \preccurlyeq_K \mathbf{0}$ and $\mathbf{y}^\star \in K^\star$. then the following are equivalent.

$1771$     1. $\mathcal{L}(\mathbf{y}^\star) = f(\mathbf{x}^\star)$.

$1772$     2. $\hat{\mathcal{L}}(\mathbf{x}^\star, \hat{\mathbf{y}}) \leq \hat{\mathcal{L}}(\mathbf{x}^\star, \mathbf{y}^\star) \leq \hat{\mathcal{L}}(\hat{\mathbf{x}}, \mathbf{y}^\star)$, for all $\hat{\mathbf{x}} \in \mathbb{R}^d$ and $\hat{\mathbf{y}} \in K^\star$.

*Proof.* 1. $\implies$ 2. Consider any $\hat{\mathbf{x}} \in \mathbb{R}^d$ and $\hat{\mathbf{y}} \in K^\star$. We now derive the following chain of inequalities:

$$
\begin{aligned}
\hat{\mathcal{L}}(\mathbf{x}^\star, \hat{\mathbf{y}}) &= f(\mathbf{x}^\star) + \langle \hat{\mathbf{y}}, G(\mathbf{x}^\star) \rangle && \\
&\leq f(\mathbf{x}^\star) && && \text{since } \langle \hat{\mathbf{y}}, G(\mathbf{x}^\star) \rangle \leq 0 \text{ because } \hat{\mathbf{y}} \in K^\star, G(\mathbf{x}^\star) \preccurlyeq_K \mathbf{0} \\
&= f(\mathbf{x}^\star) + \langle \mathbf{y}^\star, G(\mathbf{x}^\star) \rangle && = \hat{\mathcal{L}}(\mathbf{x}^\star, \mathbf{y}^\star) && \text{since } \langle \mathbf{y}^\star, G(\mathbf{x}^\star) \rangle = 0 \text{ by Theorem 4.32} \\
&= \mathcal{L}(\mathbf{y}^\star) && && \text{since } \mathcal{L}(\mathbf{y}^\star) = f(\mathbf{x}^\star) \\
&= \inf_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + \langle \mathbf{y}^\star, G(\mathbf{x}) \rangle && \\
&\leq f(\hat{\mathbf{x}}) + \langle \mathbf{y}^\star, G(\hat{\mathbf{x}}) \rangle && \\
&= \hat{\mathcal{L}}(\hat{\mathbf{x}}, \mathbf{y}^\star) &&
\end{aligned}
$$

NOTES: 84

2. $\implies$ 1. Since $\hat{\mathcal{L}}(\mathbf{x}^\star, \hat{\mathbf{y}}) \leq \hat{\mathcal{L}}(\mathbf{x}^\star, \mathbf{y}^\star)$ for all $\hat{\mathbf{y}} \in K^\star$, we have that

$$\hat{\mathcal{L}}(\mathbf{x}^\star, \mathbf{y}^\star) = \sup_{\mathbf{y} \in K^\star} \hat{\mathcal{L}}(\mathbf{x}^\star, \hat{\mathbf{y}}) = \sup_{\mathbf{y} \in K^\star} f(\mathbf{x}^\star) + \langle \mathbf{y}, G(\mathbf{x}^\star) \rangle = f(\mathbf{x}^\star),$$

where the last equality follows from the fact that $\langle \mathbf{y}, G(\mathbf{x}^\star) \rangle \leq 0$ for all $\mathbf{y} \in K^\star$ and so the supremum is achieved for $\mathbf{y} = \mathbf{0}$. On the other hand, since $\hat{\mathcal{L}}(\mathbf{x}^\star, \mathbf{y}^\star) \leq \hat{\mathcal{L}}(\hat{\mathbf{x}}, \mathbf{y}^\star)$, for all $\hat{\mathbf{x}} \in \mathbb{R}^d$, we have that

$$\hat{\mathcal{L}}(\mathbf{x}^\star, \mathbf{y}^\star) = \inf_{\mathbf{x} \in \mathbb{R}^d} \hat{\mathcal{L}}(\mathbf{x}, \mathbf{y}^\star) = \inf_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + \langle \mathbf{y}^\star, G(\mathbf{x}) \rangle = \mathcal{L}(\mathbf{y}^\star).$$

Thus, we obtain that $f(\mathbf{x}^\star) = \hat{\mathcal{L}}(\mathbf{x}^\star, \mathbf{y}^\star) = \mathcal{L}(\mathbf{y}^\star)$. $\qquad\square$

Theorem 4.33 says that $\mathbf{x}^\star$ and $\mathbf{y}^\star$ are solutions for the primal problem (4.4) and dual problem (4.10) respectively, if and only if $(\mathbf{x}^\star, \mathbf{y}^\star)$ is a saddle point for the function $\hat{\mathcal{L}}(\mathbf{x}, \mathbf{y})$. This can be used to directly solve (4.4) and (4.10) simultaneously by searching for saddle-points of the function $\hat{\mathcal{L}}(\mathbf{x}, \mathbf{y})$. This approach can be useful, if one has analytical forms for $f$ and $G$ (with sufficient differentiable properties) so that finding saddle-points is a reasonable option.

## 4.4   Cutting plane schemes

We now go back to the most general convex optimization (4.1). As before, we make no assumptions on $f$ and $C$ except that we have access to first-order oracles for $f$ and $C$, i.e., for any $\mathbf{x} \in \mathbb{R}^d$, the oracle returns an element from the subdifferential $\partial f(\mathbf{x})$, and if $\mathbf{x} \notin C$ then it returns a separating hyperplane.

The subgradient algorithm from Section 4.1 can be used to solve (4.1) if one has access to the projection operator $\text{Proj}_C(\mathbf{x})$, which is stronger than a separation oracle. *Cutting plane schemes* are a class of algorithms that work with just a separation oracle. Moreover, the number of oracle calls is is quite different from the number of oracle calls made by the subgradient algorithm: on the one hand, they typically exhibit a logarithmic dependence of $\ln(\frac{MR}{\epsilon})$ on the initial data $M, R$ and error guarantee $\epsilon$ as opposed to the quadratic dependence $\frac{M^2 R^2}{\epsilon^2}$ of the subgradient algorithm; on the other other, cutting plane schemes have a polynomial dependence on the dimension $d$ of the problem (typically of the order of $d^2$), and such a dependence does not exist for the subgradient algorithm – see Remark 4.14.

We will present the algorithm and the analysis for the situation when $C$ is compact and full-dimensional. Hence the minimizer $\mathbf{x}^\star$ exists for (4.1) since $f$ is convex, and therefore, continuous by Theorem 3.21. There are ways to get around this assumption, but we will ignore this complication in this write-up.

**General cutting plane scheme**

1. Choose any $E_0 \supseteq C$.

2. For $i = 0, 1, 2, \ldots$, do

NOTES: 85

(a) Choose $\mathbf{x}^i \in E_i$.

(b) Call the separation oracle for $C$ with $\mathbf{x}^i$ as input.

    *Case 1:* $\mathbf{x}^i \in C$. Call the first order oracle for $f$ to get some $\mathbf{s}^i \in \partial f(\mathbf{x}^i)$.

    *Case 2:* $\mathbf{x}^i \notin C$. Set $\mathbf{s}^i$ to be the normal vector of some separating hyperplane for $\mathbf{x}^i$ from $C$.

(c) Set $E_{i+1} \supseteq E_i \cap \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{s}^i, \mathbf{x} \rangle \leq \langle \mathbf{s}^i, \mathbf{x}^i \rangle\}$.

The points $\mathbf{x}^0, \mathbf{x}^1, \ldots$ will be called the *iterates* of the Cutting Plane scheme.

**Remark 4.34.** The above general scheme actually defines a family of algorithms. We have two choices to make to get a particular algorithm out of this scheme. First, there must be a strategy/procedure to choose $\mathbf{x}^i \in E_i$ in step 2(a) in every iteration. Second, there should be a strategy to define $E_{i+1}$ as a superset of $E_i \cap \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{s}^i, \mathbf{x} \rangle \leq \langle \mathbf{s}^i, \mathbf{x}^i \rangle\}$ in step 2(c) of the scheme. Depending on what these two strategies are, we get different variants of the general cutting plane scheme. We will look at two variants below: the *center of gravity method* and the *ellipsoid method.*

Technically, we also have to make a choice for $E_0$ in Step 1, but this is usually given as part of the input to the problem: $E_0$ is usually large ball or polytope containing $C$ that is provided or known at the start.

We now start our analysis of cutting plane schemes. We introduce a useful notation to denote the polyhedron defined by the halfspaces obtained during the iterations of the cutting plane scheme.

**Definition 4.35.** Let $\mathbf{z}^1, \ldots, \mathbf{z}^k \subseteq \mathbb{R}^d$ and let $\mathbf{s}^1, \ldots, \mathbf{s}^k$ be the corresponding outputs of the first-order oracle, i.e., $\mathbf{s}^i \in \partial f(\mathbf{z}^i)$ if $\mathbf{z}^i \in C$, and $\mathbf{s}^i$ is the normal vector of a separating hyperplane if $\mathbf{z}^i \notin C$. Define

$$G(\mathbf{z}^1, \ldots, \mathbf{z}^k) := \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{s}^i, \mathbf{x} \rangle \leq \langle \mathbf{s}^i, \mathbf{z}^i \rangle \ \ i = 1, \ldots, k\}.$$

This polyhedron will be referred to as the *gradient polyhedron* of $\mathbf{z}^1, \ldots, \mathbf{z}^k$. The name is a bit of a misnomer, because we are considering general $f$, so we may have no gradients, and also some of the halfspaces could correspond to separating hyperplanes which have nothing to do with gradients. Even so we stick with this terminology.

**Definition 4.36.** Let $\mathbf{x}^0, \mathbf{x}^1, \ldots$ be the iterates of a cutting plane scheme. For any iteration $t \geq 0$, we define $h(t) := |C \cap \{\mathbf{x}^0, \ldots, \mathbf{x}^t\}|$, i.e., $h(t)$ is the number of feasible iterates until iteration $t$. We also define

$$S_t = C \cap G(\mathbf{x}^0, \ldots, \mathbf{x}^t).$$

As we shall see below, the volume of $S_t$ will be central in measuring our progress towards the optimal solution. We first observe in the next lemma that $S_t$ can be describe as the intersection of $C$ and the gradient polyhedron of only the feasible iterates.

**Lemma 4.37.** Let $\mathbf{x}^0, \mathbf{x}^1, \ldots$ be the iterates of a cutting plane scheme. Let the feasible iterates be denoted by $\{\mathbf{x}^{i_1}, \ldots, \mathbf{x}^{i_{h(t)}}\} = C \cap \{\mathbf{x}^0, \ldots, \mathbf{x}^t\}$, with $0 \leq i_1 \leq i_2 \leq \ldots \leq i_{h(t)}$. Then $S_t = C \cap G(\mathbf{x}^{i_1}, \ldots, \mathbf{x}^{i_{h(t)}})$.

NOTES:                                  86

*Proof.* Let $X_t = \{\mathbf{x}^0, \ldots, \mathbf{x}^t\}$. We derive the following relations.

$$
\begin{aligned}
S_t &= C \cap G(\mathbf{x}^0, \ldots, \mathbf{x}^t) \\
&= C \cap G(X_t \setminus \{\mathbf{x}^{i_1}, \ldots, \mathbf{x}^{i_{h(t)}}\}) \cap G(\mathbf{x}^{i_1}, \ldots, \mathbf{x}^{i_{h(t)}}) \\
&= C \cap G(\mathbf{x}^{i_1}, \ldots, \mathbf{x}^{i_{h(t)}}),
\end{aligned}
$$

where the last inequality follows since $C \subseteq G(X_t \setminus \{\mathbf{x}^{i_1}, \ldots, \mathbf{x}^{i_{h(t)}}\})$ because each $\mathbf{z} \in X_t \setminus \{\mathbf{x}^{i_1}, \ldots, \mathbf{x}^{i_{h(t)}}\}$ is infeasible, i.e., $\mathbf{z} \notin C$, and therefore, the corresponding vector $\mathbf{s}$ is a separating hyperplane for $\mathbf{z}$ and $C$. Thus, $C \subseteq \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{s}, \mathbf{x} \rangle \leq \langle \mathbf{s}, \mathbf{z} \rangle\}$. $\square$

Since our analysis will involve the volume of $S_t$, while our algorithm only works with the sets $E_t$, we need to establish a definite relationship between these two sets.

**Lemma 4.38.** Let $\mathbf{x}^0, \mathbf{x}^1, \ldots$ be the iterates of a cutting plane scheme. Then $E_{t+1} \supseteq S_t$ for all $t \geq 0$.

*Proof.* By definition $E_{i+1} \supseteq E_i \cap \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{s}^i, \mathbf{x} \rangle \leq \langle \mathbf{s}^i, \mathbf{x}^i \rangle\}$ for all $i = 0, \ldots, t$. By putting all these relationships together, we obtain that

$$
E_{t+1} \supseteq E_0 \cap G(\mathbf{x}^0, \ldots, \mathbf{x}^t) \supseteq C \cap G(\mathbf{x}^0, \ldots, \mathbf{x}^t) = S_t, \tag{4.21}
$$

where the second containment follows from the assumption that $E_0 \supseteq C$. $\square$

We now state our main structural result for the analysis of cutting plane schemes. We use $\mathrm{dist}(\mathbf{x}, X)$ to denote the distance of $\mathbf{x} \in \mathbb{R}^d$ from any subset $X \subseteq \mathbb{R}^d$, i.e., $\mathrm{dist}(\mathbf{x}, X) := \inf_{\mathbf{y} \in X} \|\mathbf{x} - \mathbf{y}\|$.

**Theorem 4.39.** Let $f : \mathbb{R}^d \to \mathbb{R}$ be a convex function and let $C$ be a compact, convex set. Let $\mathbf{x}^\star$ be the minimizer for (4.1). Let $\mathbf{x}^0, \mathbf{x}^1, \ldots$ be the iterates of any cutting plane scheme. Let the feasible iterates be denoted by $\{\mathbf{x}^{i_1}, \ldots, \mathbf{x}^{i_{h(t)}}\} = C \cap \{\mathbf{x}^0, \ldots, \mathbf{x}^t\}$, with $0 \leq i_1 \leq i_2 \leq \ldots \leq i_{h(t)}$. Define

$$
v_{\min}(t) := \min_{j = i_1, \ldots, i_{h(t)}} \mathrm{dist}(\mathbf{x}^\star, H(\mathbf{s}^j, \langle \mathbf{s}^j, \mathbf{x}^j \rangle)),
$$

i.e., $v_{\min}(t)$ is the minimum distance of $\mathbf{x}^\star$ from the hyperplanes $\{\mathbf{x} : \langle \mathbf{s}^j, \mathbf{x} \rangle = \langle \mathbf{s}^j, \mathbf{x}^j \rangle\}$, $j = i_1, \ldots i_{h(t)}$. Let $D$ be diameter of $C$, i.e., $D = \max_{\mathbf{x}, \mathbf{y} \in C} \|\mathbf{x} - \mathbf{y}\|$. Then the following are all true.

1. For any $t \geq 0$, if $\mathrm{vol}(E_{t+1}) < \mathrm{vol}(C)$ then $h(t) > 0$, i.e., there is at least one feasible iterate.

2. For any $t \geq 0$ such that $h(t) > 0$, $v_{\min}(t) \leq D\left(\frac{\mathrm{vol}(S_t)}{\mathrm{vol}(C)}\right)^{\frac{1}{d}} \leq D\left(\frac{\mathrm{vol}(E_{t+1})}{\mathrm{vol}(C)}\right)^{\frac{1}{d}}$ for all $t \geq 0$.

3. For any $t \geq 0$ such that $h(t) > 0$, $\min_{j = i_1, \ldots i_{h(t)}} f(\mathbf{x}^j) \leq f(\mathbf{x}^\star) + M v_{\min}(t) \leq f(\mathbf{x}^\star) + M D\left(\frac{\mathrm{vol}(E_{t+1})}{\mathrm{vol}(C)}\right)^{\frac{1}{d}}$, where $M = L(B_2(\mathbf{x}^\star, v_{\min}))$ is a Lipschitz constant for $f$ over $B_2(\mathbf{x}^\star, v_{\min})$ (see Theorem 3.21). This provides a bound on the value of the best feasible point seen upto iteration $t$, in comparison to the optimal value $f(\mathbf{x}^\star)$.

NOTES: 87

1841    Theorem 4.39 shows that if we can ensure $\mathrm{vol}(E_t) \to 0$ as $t \to \infty$, then we have a convergent algorithm.

1842    *Proof of Theorem 4.39.*     1. We prove the contrapositive. If $h(t) = 0$, then all iterates upto iteration $t$
1843    are infeasible, i.e., $\mathbf{x}^i \notin C$ for all $i = 1, \ldots, t$. This implies that all the vector $\mathbf{s}^i$ are normal vectors
1844    for separating hyperplanes. So $C \subseteq G(\mathbf{x}^0, \ldots, \mathbf{x}^t)$. Since $C \subseteq E_0$, this implies that $C = E_0 \cap C \subseteq$
1845    $E_0 \cap G(\mathbf{x}^0, \ldots, \mathbf{x}^t) \subseteq E_{t+1}$, where the last containment follows from the first containment in (4.21).
1846    Therefore, $\mathrm{vol}(C) \leq \mathrm{vol}(E_{t+1})$.

2. Let $\alpha = \frac{v_{\min}(t)}{D}$. Since $D$ is the diameter of $C$, we must have $C \subseteq B_2(\mathbf{x}^\star, D)$. Thus,

$$\alpha(C - \mathbf{x}^\star) + \mathbf{x}^\star \subseteq B_2(\mathbf{x}^\star, \alpha D) = B_2(\mathbf{x}^\star, v_{\min}(t)) \subseteq G(\mathbf{x}^{i_1}, \ldots, \mathbf{x}^{i_{h(t)}}),$$

1847    where the first equality follows from the definition of $\alpha$ and the final containment follows the definition
1848    of $v_{\min}(t)$. Since $\mathbf{x}^\star \in C$ and $C$ is convex, we know that $\alpha(C - \mathbf{x}^\star) + \mathbf{x}^\star = \alpha C + (1-\alpha)\mathbf{x}^\star \subseteq C$. Therefore,
1849    $\alpha(C - \mathbf{x}^\star) + \mathbf{x}^\star = C \cap (\alpha(C - \mathbf{x}^\star) + \mathbf{x}^\star) \subseteq C \cap G(\mathbf{x}^{i_1}, \ldots, \mathbf{x}^{i_{h(t)}}) = S_t$, where the last equality follows
1850    from Lemma 4.37. This implies that $\alpha^d \mathrm{vol}(C) = \mathrm{vol}(\alpha(C - \mathbf{x}^\star)) \leq \mathrm{vol}(S_t)$. Rearranging and using the
1851    definition of $\alpha$, we obtain that $v_{\min}(t) \leq D\left(\frac{\mathrm{vol}(S_t)}{\mathrm{vol}(C)}\right)^{\frac{1}{d}}$. By Lemma 4.38, $D\left(\frac{\mathrm{vol}(S_t)}{\mathrm{vol}(C)}\right)^{\frac{1}{d}} \leq D\left(\frac{\mathrm{vol}(E_{t+1})}{\mathrm{vol}(C)}\right)^{\frac{1}{d}}$.

3. It suffices to prove the first inequality; the second inequality follows from part 1. above. Let $i^{\min} \in$
$\{i_1, i_2, \ldots, i_{h(t)}\}$ be such that $v_{\min}(t) = \mathrm{dist}(\mathbf{x}^\star, H(\mathbf{s}^{i^{\min}}, \langle \mathbf{s}^{i^{\min}}, \mathbf{x}^{i^{\min}}\rangle))$. Denote $H := H(\mathbf{s}^{i^{\min}}, \langle \mathbf{s}^{i^{\min}}, \mathbf{x}^{i^{\min}}\rangle)$
passing through $\mathbf{x}^{i^{\min}}$, orthogonal to $\mathbf{s}^{i^{\min}}$. Let $\bar{\mathbf{x}}$ be the point on $H$ closest to $\mathbf{x}^\star$. Using the Lipschitz
constant $M$, we obtain that $f(\bar{\mathbf{x}}) \leq f(\mathbf{x}^\star) + M v_{\min}(t)$; see Figure 3.. Finally, since $\mathbf{s}^{i^{\min}} \in \partial f(\mathbf{x}^{i^{\min}})$,
we must have that $f(\bar{\mathbf{x}}) \geq f(\mathbf{x}^{i^{\min}}) + \langle \mathbf{s}^{i^{\min}}, \bar{\mathbf{x}} - \mathbf{x}^{i^{\min}}\rangle = f(\mathbf{x}^{i^{\min}})$, since $\bar{\mathbf{x}} \in H$ implying that
$\langle \mathbf{s}^{i^{\min}}, \bar{\mathbf{x}} - \mathbf{x}^{i^{\min}}\rangle = 0$. Therefore, we obtain

$$\min_{j = i_1, \ldots i_{h(t)}} f(\mathbf{x}^j) \leq f(\mathbf{x}^{i^{\min}}) \leq f(\bar{\mathbf{x}}) \leq f(\mathbf{x}^\star) + M v_{\min}(t).$$

1852    $\square$

1853    We now analyze two instantiations of the cutting plane scheme with concrete strategies to choose $\mathbf{x}^i$ and
1854    $E_{i+1}$ in each iteration $i$.

1855    **Center of Gravity Method.**    The first one is called the *center of gravity method*.

**Definition 4.40.** The *center of gravity* for any compact set $X \subseteq \mathbb{R}^d$ with non-zero volume is defined as

$$\frac{\int_X \mathbf{x}\, d\mathbf{x}}{\mathrm{vol}(X)}.$$

NOTES:                                    88

An important property of the center gravity of compact, convex sets was established by Grünbaum [4].

**Theorem 4.41.** Let $C \subseteq \mathbb{R}^d$ be a compact, convex set with center of gravity $\bar{\mathbf{x}}$. Then for every hyperplane $H$ such that $\mathbf{x} \in H$,
$$\frac{1}{e} \leq \left(\frac{d}{d+1}\right)^d \leq \frac{\mathrm{vol}(H^+ \cap C)}{\mathrm{vol}(C)} \leq 1 - \left(\frac{d}{d+1}\right)^d \leq 1 - \frac{1}{e},$$
where $H^+$ is a halfspace with boundary $H$.

Theorem 4.41 follows from the proof of Theorem 2 in [4] and will not be repeated here.

In the *center of gravity method*, $\mathbf{x}_i$ is chosen as the center of gravity of $E_i$ in Step 2(a) of the General cutting plane scheme, and $E_{i+1}$ is set to be *equal* to $E_i \cap \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{s}^i, \mathbf{x} \rangle \leq \langle \mathbf{s}^i, \mathbf{x}^i \rangle\}$ in Step 2(c) in the General cutting plane scheme. Theorem 4.41 then implies the following. Sometimes, the center of gravity method assumes that $E_0 = C$, where the central assumption is that one can compute the center of gravity of $C$ and any subset of it.

**Theorem 4.42.** In the center of gravity method, if $h(t) > 0$ for some iteration $t \geq 0$, then

$$\min_{j=i_1,\ldots i_{h(t)}} f(\mathbf{x}^j) \leq f(\mathbf{x}^\star) + MD(1 - \frac{1}{e})^{t/d}(\frac{\mathrm{vol}(E_0)}{\mathrm{vol}(C)})^{1/d},$$

where $D$ is the diameter of $C$ and $M$ is a Lipschitz constant for $f$ over $B_2(\mathbf{x}^\star, D)$.
    In particular, if $E_0 = C$, then $\min_{j=i_1,\ldots i_{h(t)}} f(\mathbf{x}^j) \leq f(\mathbf{x}^\star) + MD(1 - \frac{1}{e})^{t/d}$.

*Proof.* Follows from Theorem 4.39 part 3., and the fact that $B(\mathbf{x}^\star, v_{\min}) \subseteq B(\mathbf{x}^\star, D)$ implying that $M$ is a Lipschitz constant for $f$ over $B(\mathbf{x}^\star, v_{\min})$, and $\mathrm{vol}(E_{t+1}) \leq (1 - \frac{1}{e})^t \mathrm{vol}(E_0)$ by Theorem 4.41. □

By setting the error term $MD(1 - \frac{1}{e})^{t/d}(\frac{\mathrm{vol}(E_0)}{\mathrm{vol}(C)})^{1/d}$ less than equal to $\epsilon$ in Theorem 4.42, the following is an immediate consequence.

**Corollary 4.43.** For any $\epsilon > 0$, after $O(d\ln(\frac{MD}{\epsilon}) + \ln\left(\frac{\mathrm{vol}(E_0)}{\mathrm{vol}(C)}\right))$ iterations of the center of gravity method,

$$\min_{j=i_1,\ldots i_{h(t)}} f(\mathbf{x}^j) \leq f(\mathbf{x}^\star) + \epsilon.$$

In particular, if $E_0 = C$, then one needs $O(d\ln(\frac{MD}{\epsilon}))$ iterations.

NOTES:                                                      89

**Ellipsoid method.** The ellipsoid method is a cutting plane scheme where $E_0$ is assumed to be a large ball with radius $R$ around a known point $\mathbf{x}_0$ (typically $\mathbf{x}_0 = \mathbf{0}$) that is guaranteed to contain $C$. At every iteration $i$, $E_i$ is maintained to be an ellipsoid and in Step 2(a), $\mathbf{x}^i$ is chosen to be the center of $E_i$. In Step 2(c), $E_{i+1}$ is set to be an ellipsoid that contains $E_i \cap \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{s}^i, \mathbf{x} \rangle \leq \langle \mathbf{s}^i, \mathbf{x}^i \rangle\}$, such that $\mathrm{vol}(E_{i+1}) \leq (1 - \frac{1}{d^2+1})^{d/2} \mathrm{vol}(E_i)$. The technical bulk of the analysis goes into showing that such an ellipsoid $E_{i+1}$ *always* exists.

**Definition 4.44.** Recall from Definition 2.2 that an ellipsoid is the unit ball associated with the norm induced by a positive definite matrix, i.e., $E = \{\mathbf{x} \in \mathbb{R}^d : \mathbf{x}^T A \mathbf{x} \leq 1\}$ for some positive definite matrix $A$. First, we need to also consider translated ellipsoids so that the center is not $\mathbf{0}$ anymore. Secondly, for computational reasons involving inverses of matrices, we will actually define the following family of objects, which are just translated ellipsoids, just written in a different way. Given a positive definite matrix $H \in \mathbb{R}^{d \times d}$ and a point $\mathbf{y} \in \mathbb{R}^d$, we define

$$E(H, \mathbf{y}) := \{\mathbf{x} \in \mathbb{R}^d : (\mathbf{x} - \mathbf{y})^T H^{-1} (\mathbf{x} - \mathbf{y}) \leq 1\}.$$

The next proposition follows from unwrapping the definition. It shows that ellipsoids are simply the image of the Euclidean unit norm ball under an invertible linear transformation.

**Proposition 4.45.** Let $H \in \mathbb{R}^{d \times d}$ be a positive definite matrix and let $H^{-1} = B^T B$ for some invertible matrix $B \in \mathbb{R}^{d \times d}$. Then $E(H, \mathbf{y}) = \mathbf{y} + B^{-1}(B_2(\mathbf{0}, 1))$. Thus, $\mathrm{vol}(E(H, \mathbf{y})) = \det(B^{-1}) \mathrm{vol}(B_2(\mathbf{0}, 1)) = \sqrt{\det(H)} \mathrm{vol}(B_2(\mathbf{0}, 1))$.

In the following, we will utilize the following relation for any $\mathbf{w}, \mathbf{z} \in \mathbb{R}^d$ and $A \in \mathbb{R}^{d \times d}$

$$(\mathbf{w} + \mathbf{z})^T A (\mathbf{w} + \mathbf{z}) = \mathbf{w}^T A \mathbf{w} + 2 \mathbf{w}^T A \mathbf{z} + \mathbf{z}^T A \mathbf{z}. \tag{4.22}$$

**Theorem 4.46.** Let $H \in \mathbb{R}^{d \times d}$ and $\mathbf{y} \in \mathbb{R}^d$. Let $\mathbf{s} \in \mathbb{R}^d$ and let $E_+ = E(H, \mathbf{y}) \cap H^-(\mathbf{s}, \langle \mathbf{s}, \mathbf{y} \rangle)$. Define

$$\mathbf{y}_+ = \mathbf{y} - \frac{1}{d+1} \cdot \frac{H\mathbf{s}}{\sqrt{\mathbf{s}^T H \mathbf{s}}}$$

$$H_+ = \frac{d^2}{d^2-1}\left(H - \frac{2}{d+1} \cdot \frac{H\mathbf{s}\mathbf{s}^T H}{\mathbf{s}^T H \mathbf{s}}\right).$$

Then $E_+ \subseteq E(H_+, \mathbf{y}_+)$ and $\mathrm{vol}(E(H_+, \mathbf{y}_+)) \leq (1 - \frac{1}{(d+1)^2})^{d/2} \mathrm{vol}(E(H, \mathbf{y}))$.

*Proof.* We first prove $E_+ \subseteq E(H_+, \mathbf{y}_+)$. Consider any $\mathbf{x} \in E_+ = E(H, \mathbf{y}) \cap H^-(\mathbf{s}, \langle \mathbf{s}, \mathbf{y} \rangle)$. To ease notational burden, we denote $G = H^{-1}$ and $G_+ = H_+^{-1}$. A direct calculation shows that $G_+ = \frac{d^2-1}{d^2}(G + \frac{2}{d-1} \cdot \frac{\mathbf{s}\mathbf{s}^T}{\mathbf{s}^T H \mathbf{s}})$. Thus, $\mathbf{x}$ satisfies

$$(\mathbf{x} - \mathbf{y})^T G (\mathbf{x} - \mathbf{y}) \leq 1 \tag{4.23}$$

$$\langle \mathbf{s}, \mathbf{x} - \mathbf{y} \rangle \leq 0 \tag{4.24}$$

NOTES: 90

We now verify that

$$(\mathbf{x}-\mathbf{y}_+)^T G_+(\mathbf{x}-\mathbf{y}_+) = (\mathbf{x}-\mathbf{y}+\tfrac{1}{d+1}\cdot\tfrac{H\mathbf{s}}{\sqrt{\mathbf{s}^T H\mathbf{s}}})^T G_+(\mathbf{x}-\mathbf{y}+\tfrac{1}{d+1}\cdot\tfrac{H\mathbf{s}}{\sqrt{\mathbf{s}^T H\mathbf{s}}})$$
$$= (\mathbf{x}-\mathbf{y})^T G_+(\mathbf{x}-\mathbf{y})+\tfrac{2}{d+1}(\mathbf{x}-\mathbf{y})^T G_+(\tfrac{H\mathbf{s}}{\sqrt{\mathbf{s}^T H\mathbf{s}}})+(\tfrac{1}{d+1})^2\tfrac{\mathbf{s}^T H^T G_+ H\mathbf{s}}{\mathbf{s}^T H\mathbf{s}},$$

<sub></sub>where we use (4.22). Let us analyze the three terms separately. The first term simplifies to

$$(\mathbf{x}-\mathbf{y})^T G_+(\mathbf{x}-\mathbf{y}) = (\mathbf{x}-\mathbf{y})^T(\tfrac{d^2-1}{d^2}(G+\tfrac{2}{d-1}\cdot\tfrac{\mathbf{s}\mathbf{s}^T}{\mathbf{s}^T H\mathbf{s}}))(\mathbf{x}-\mathbf{y})$$
$$= \tfrac{d^2-1}{d^2}\left((\mathbf{x}-\mathbf{y})^T G(\mathbf{x}-\mathbf{y})+\tfrac{2}{d-1}\tfrac{(\mathbf{s}^T(\mathbf{x}-\mathbf{y}))^2}{\mathbf{s}^T H\mathbf{s}}\right)$$

The second term simplifies to

$$\tfrac{2}{d+1}(\mathbf{x}-\mathbf{y})^T G_+(\tfrac{H\mathbf{s}}{\sqrt{\mathbf{s}^T H\mathbf{s}}}) = \tfrac{2}{d+1}(\mathbf{x}-\mathbf{y})^T(\tfrac{d^2-1}{d^2}(G+\tfrac{2}{d-1}\cdot\tfrac{\mathbf{s}\mathbf{s}^T}{\mathbf{s}^T H\mathbf{s}}))(\tfrac{H\mathbf{s}}{\sqrt{\mathbf{s}^T H\mathbf{s}}})$$
$$= \tfrac{d^2-1}{d^2}\cdot\tfrac{2}{d+1}\left(\tfrac{\mathbf{s}^T(\mathbf{x}-\mathbf{y})}{\sqrt{\mathbf{s}^T H\mathbf{s}}}+\tfrac{2}{d-1}\cdot\tfrac{(\mathbf{x}-\mathbf{y})^T\mathbf{s}\mathbf{s}^T H\mathbf{s}}{\mathbf{s}^T H\mathbf{s}\cdot\sqrt{\mathbf{s}^T H\mathbf{s}}}\right)$$
$$= \tfrac{d^2-1}{d^2}\cdot\tfrac{2}{d+1}\left(\tfrac{\mathbf{s}^T(\mathbf{x}-\mathbf{y})}{\sqrt{\mathbf{s}^T H\mathbf{s}}}+\tfrac{2}{d-1}\cdot\tfrac{(\mathbf{x}-\mathbf{y})^T\mathbf{s}}{\sqrt{\mathbf{s}^T H\mathbf{s}}}\right)$$
$$= \tfrac{d^2-1}{d^2}\cdot\tfrac{2}{d-1}\left(\tfrac{\mathbf{s}^T(\mathbf{x}-\mathbf{y})}{\sqrt{\mathbf{s}^T H\mathbf{s}}}\right)$$

The third term simplifies to

$$(\tfrac{1}{d+1})^2\tfrac{\mathbf{s}^T H^T G_+ H\mathbf{s}}{\mathbf{s}^T H\mathbf{s}} = (\tfrac{1}{d+1})^2\tfrac{\mathbf{s}^T H(\tfrac{d^2-1}{d^2}(G+\tfrac{2}{d-1}\cdot\tfrac{\mathbf{s}\mathbf{s}^T}{\mathbf{s}^T H\mathbf{s}}))H\mathbf{s}}{\mathbf{s}^T H\mathbf{s}}$$
$$= \tfrac{d^2-1}{d^2}\cdot(\tfrac{1}{d+1})^2\left(\tfrac{\mathbf{s}^T H\mathbf{s}+\tfrac{2}{d-1}(\mathbf{s}^T H\mathbf{s})}{\mathbf{s}^T H\mathbf{s}}\right)$$
$$= \tfrac{d^2-1}{d^2}(\tfrac{1}{d^2-1}),$$

Putting all of it together, we obtain that

$$(\mathbf{x}-\mathbf{y}_+)^T G_+(\mathbf{x}-\mathbf{y}_+)=\frac{d^2-1}{d^2}\left((\mathbf{x}-\mathbf{y})^T G(\mathbf{x}-\mathbf{y})+\frac{2}{d-1}\frac{(\mathbf{s}^T(\mathbf{x}-\mathbf{y}))^2}{\mathbf{s}^T H\mathbf{s}}+\frac{2}{d-1}\left(\frac{\mathbf{s}^T(\mathbf{x}-\mathbf{y})}{\sqrt{\mathbf{s}^T H\mathbf{s}}}\right)+\frac{1}{d^2-1}\right) \quad (4.25)$$

<sub></sub>We now argue that $\frac{(\mathbf{s}^T(\mathbf{x}-\mathbf{y}))^2}{\mathbf{s}^T H\mathbf{s}}+\frac{\mathbf{s}^T(\mathbf{x}-\mathbf{y})}{\sqrt{\mathbf{s}^T H\mathbf{s}}}=\frac{\mathbf{s}^T(\mathbf{x}-\mathbf{y})}{\mathbf{s}^T H\mathbf{s}}(\sqrt{\mathbf{s}^T H\mathbf{s}}+\mathbf{s}^T(\mathbf{x}-\mathbf{y}))\le 0$. Since $\mathbf{s}^T(\mathbf{x}-\mathbf{y})\le 0$ by
<sub></sub>(4.24), it suffices to show that $\sqrt{\mathbf{s}^T H\mathbf{s}}+\mathbf{s}^T(\mathbf{x}-\mathbf{y})\ge 0$, or equivalently, that $|\mathbf{s}^T(\mathbf{x}-\mathbf{y})|\le\sqrt{\mathbf{s}^T H\mathbf{s}}$.

<sub></sub>**Claim 4.** $|\mathbf{s}^T(\mathbf{x}-\mathbf{y})|\le\sqrt{\mathbf{s}^T H\mathbf{s}}$.

*Proof of Claim.* Let the eigendecomposition of $H$ be given as $H=S\Lambda S^T$, where $S$ is the orthonormal matrix which has the eigenvectors of $H$ as columns, and $\Lambda$ is a diagonal matrix with the corresponding eigenvalues.

NOTES: 91

Then $H^{-1} = S\Lambda^{-1}S^T = G$. Now,

$$
\begin{aligned}
|\mathbf{s}^T(\mathbf{x}-\mathbf{y})| &= |\mathbf{s}^T S\Lambda^{\frac{1}{2}}\Lambda^{-\frac{1}{2}}S^T(\mathbf{x}-\mathbf{y})| \\
&= |\langle \Lambda^{\frac{1}{2}}S^T\mathbf{s}, \Lambda^{-\frac{1}{2}}S^T(\mathbf{x}-\mathbf{y})\rangle| \\
&\leq \|\Lambda^{\frac{1}{2}}S^T\mathbf{s}\|_2\|\Lambda^{-\frac{1}{2}}S^T(\mathbf{x}-\mathbf{y})\|_2 \\
&= \sqrt{(\Lambda^{\frac{1}{2}}S^T\mathbf{s})^T(\Lambda^{\frac{1}{2}}S^T\mathbf{s})}\sqrt{(\Lambda^{-\frac{1}{2}}S^T(\mathbf{x}-\mathbf{y}))^T(\Lambda^{-\frac{1}{2}}S^T(\mathbf{x}-\mathbf{y}))} \\
&= \sqrt{\mathbf{s}^T S\Lambda^{\frac{1}{2}}\Lambda^{\frac{1}{2}}S^T\mathbf{s}}\sqrt{(\mathbf{x}-\mathbf{y})^T S\Lambda^{-\frac{1}{2}}\Lambda^{-\frac{1}{2}}S^T(\mathbf{x}-\mathbf{y})} \\
&= \sqrt{\mathbf{s}^T H\mathbf{s}}\sqrt{(\mathbf{x}-\mathbf{y})^T G(\mathbf{x}-\mathbf{y})} \\
&\leq \sqrt{\mathbf{s}^T H\mathbf{s}},
\end{aligned}
$$

where there first inequality is the Cauchy-Schwarz inequality, and the last inequality follows from (4.23). □

This claim, together with (4.25), implies that

$$
\begin{aligned}
(\mathbf{x}-\mathbf{y}_+)^T G_+(\mathbf{x}-\mathbf{y}_+) &\leq \tfrac{d^2-1}{d^2}\left((\mathbf{x}-\mathbf{y})^T G(\mathbf{x}-\mathbf{y}) + \tfrac{1}{d^2-1}\right) \\
&\leq \tfrac{d^2-1}{d^2}\left(1+\tfrac{1}{d^2-1}\right) \\
&= 1,
\end{aligned}
$$

where the second inequality follows from (4.23).

We now prove the volume claim. Let $H = B^T B$ for some invertible matrix $B$. We use $I_d$ to denote the $d\times d$ identity matrix. By Proposition 4.45,

$$
\begin{aligned}
\frac{\mathrm{vol}(E(H_+,\mathbf{y}_+))}{\mathrm{vol}(E(H,\mathbf{y}))} &= \sqrt{\frac{\det(H_+)}{\det(H)}} \\
&= \sqrt{\frac{\det(\frac{d^2}{d^2-1}(H-\frac{2}{d+1}\cdot\frac{H\mathbf{s}\mathbf{s}^T H}{\mathbf{s}^T H\mathbf{s}}))}{\det(H)}} \\
&= \left(\tfrac{d^2}{d^2-1}\right)^{\frac{d}{2}}\sqrt{\frac{\det(H-\frac{2}{d+1}\cdot\frac{H\mathbf{s}\mathbf{s}^T H}{\mathbf{s}^T H\mathbf{s}})}{\det(H)}} \\
&= \left(\tfrac{d^2}{d^2-1}\right)^{\frac{d}{2}}\sqrt{\frac{\det(B^T B-\frac{2}{d+1}\cdot\frac{B^T B\mathbf{s}\mathbf{s}^T B^T B}{\mathbf{s}^T B^T B\mathbf{s}})}{\det(B^T B)}} \\
&= \left(\tfrac{d^2}{d^2-1}\right)^{\frac{d}{2}}\sqrt{\frac{\det(B^T(I_d-\frac{2}{d+1}\cdot\frac{B\mathbf{s}\mathbf{s}^T B^T}{\mathbf{s}^T B^T B\mathbf{s}})B)}{\det(B^T)\det(B)}} \\
&= \left(\tfrac{d^2}{d^2-1}\right)^{\frac{d}{2}}\sqrt{\frac{\det(B^T)\det(I_d-\frac{2}{d+1}\cdot\frac{B\mathbf{s}\mathbf{s}^T B^T}{\mathbf{s}^T B^T B\mathbf{s}})\det(B)}{\det(B^T)\det(B)}} \\
&= \left(\tfrac{d^2}{d^2-1}\right)^{\frac{d}{2}}\sqrt{\det(I_d-\tfrac{2}{d+1}\cdot\tfrac{B\mathbf{s}\mathbf{s}^T B^T}{\mathbf{s}^T B^T B\mathbf{s}})} \\
&= \left(\tfrac{d^2}{d^2-1}\right)^{\frac{d}{2}}\cdot\left(1-\tfrac{2}{d+1}\right)^{\frac{1}{2}},
\end{aligned}
$$

where the last equality follows from the fact that the matrix $\frac{B\mathbf{s}\mathbf{s}^T B^T}{\mathbf{s}^T B^T B\mathbf{s}} = \frac{\mathbf{a}\mathbf{a}^T}{\|\mathbf{a}\|^2}$ with $\mathbf{a} = B\mathbf{s}$, is a rank one positive semidefinite matrix with eigenvalue 1 with multiplicity 1, and eigenvalue 0 with multiplicity $d-1$.

NOTES:                                                      92

Now finally we observe that

$$
\begin{aligned}
\left(\tfrac{d^2}{d^2-1}\right)^{\frac{d}{2}} \cdot (1 - \tfrac{2}{d+1})^{\frac{1}{2}} &= \left(\tfrac{d^2}{d^2-1} \cdot (1 - \tfrac{2}{d+1})^{\frac{1}{d}}\right)^{\frac{d}{2}} \\
&\leq \left(\tfrac{d^2}{d^2-1} \cdot (1 - \tfrac{2}{d(d+1)})\right)^{\frac{d}{2}} \\
&= \left(\tfrac{d^2(d^2+d-2)}{d(d+1)(d^2-1)}\right)^{\frac{d}{2}} \\
&= (1 - \tfrac{1}{(d+1)^2})^{d/2}
\end{aligned}
$$

This completes the proof. $\qquad\square$

This can be used to give the guarantee of the ellipsoid method as follows.

**Theorem 4.47.** Using the ellipsoid method with $E_0 = B(\mathbf{x}_0, R)$, if $h(t) > 0$ for some iteration $t \geq 0$, then

$$
\min_{j=i_1,\dots i_{h(t)}} f(\mathbf{x}^j) \leq f(\mathbf{x}^\star) + MR\left(1 - \frac{1}{(d+1)^2}\right)^{t/2} \cdot \left(\frac{\mathrm{vol}(E_0)}{\mathrm{vol}(C)}\right)^{1/d} \leq MRe^{-\frac{t}{2(d+1)^2}} \cdot \left(\frac{\mathrm{vol}(E_0)}{\mathrm{vol}(C)}\right)^{1/d},
$$

where $M$ is a Lipschitz constant for $f$ over $B_2(\mathbf{x}_0, 2R)$.

*Proof.* The first inequality follows from Theorem 4.39 part 3., and the fact that $B(\mathbf{x}^\star, v_{\min}) \subseteq B(\mathbf{x}_0, 2R)$ implying that $M$ is a Lipschitz constant for $f$ over $B(\mathbf{x}^\star, v_{\min})$, and $\mathrm{vol}(E_{t+1}) \leq (1 - \frac{1}{e})^t \mathrm{vol}(E_0)$ by Theorem 4.46. The second inequality follows from the general inequality that $(1 + x) \leq e^x$ for all $x \in \mathbb{R}$. $\qquad\square$

By setting the error term $MRe^{-\frac{t}{2(d+1)^2}} \cdot \left(\frac{\mathrm{vol}(E_0)}{\mathrm{vol}(C)}\right)^{1/d}$ less than equal to $\epsilon$ in Theorem 4.47, the following is an immediate consequence.

**Corollary 4.48.** For any $\epsilon > 0$, after $2((d+1)^2 \ln(\frac{MR}{\epsilon}) + \frac{(d+1)^2}{d} \ln\left(\frac{\mathrm{vol}(E_0)}{\mathrm{vol}(C)}\right))$ iterations of the ellipsoid method,

$$
\min_{j=i_1,\dots i_{h(t)}} f(\mathbf{x}^j) \leq f(\mathbf{x}^\star) + \epsilon.
$$

In particular, if there exists $\rho > 0$ such that $B_2(\mathbf{z}, \rho) \subseteq C$ for some $\mathbf{z} \in C$, then after $2(d+1)^2 \ln(\frac{MR^2}{\epsilon\rho})$ iterations of the ellipsoid method, $\min_{j=i_1,\dots i_{h(t)}} f(\mathbf{x}^j) \leq f(\mathbf{x}^\star) + \epsilon$.

*Proof.* We simply use the fact that $\mathrm{vol}(B_2(\mathbf{z}, \lambda)) = \lambda^d \mathrm{vol}(B_2(\mathbf{0}, 1))$ for any $\mathbf{z} \in \mathbb{R}^d$ and $\lambda \geq 0$. $\qquad\square$

Because of the logarithmic dependence on the data $(M, R, \rho)$ and the error guarantee $\epsilon$, and the quadratic dependence on the dimension $d$, the ellipsoid method is said to have *polynomial* running time for convex optimization.

# References

[1] Michele Conforti, Gérard Cornuéjols, Aris Daniilidis, Claude Lemaréchal, and Jérôme Malick. Cut-generating functions and s-free sets. *Mathematics of Operations Research*, 40(2):276–391, 2014.

[2] Luc Devroye, László Györfi, and Gábor Lugosi. *A probabilistic theory of pattern recognition*, volume 31. Springer Science & Business Media, 2013.

[3] P.M. Gruber. *Convex and Discrete Geometry*, volume 336 of *Grundlehren der Mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences]*. Springer-Verlag, Berlin, 2007.

[4] B. Grünbaum. Partitions of mass-distributions and of convex bodies by hyperplanes. *Pacific J. Math.*, 10:1257–1261, 1960.

[5] Vladimir N Vapnik and A Ya Chervonenkis. On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability & Its Applications*, 16(2):264–280, 1971.

NOTES: 94