

Introduction to Convexity

Amitabh Basu

Compiled on Monday 9th December, 2019 at 16:35

Contents

5	1 Definitions and Preliminaries	4
6	Basic real analysis and topology.	4
7	Basic facts about matrices.	6
8	2 Convex Sets	6
9	2.1 Definitions and basic properties	6
10	2.2 Convex cones, affine sets and dimension	8
11	2.3 Representations of convex sets	11
12	2.3.1 Extrinsic description: separating hyperplanes	11
13	How to represent general convex sets: Separation oracles.	12
14	Farkas' lemma: A glimpse into polyhedral theory.	12
15	Duality/Polarity.	12
16	2.3.2 Intrinsic description: faces, extreme points, recession cone, lineality space	14
17	2.3.3 A remark about extrinsic and intrinsic descriptions	20
18	2.4 Combinatorial theorems: Helly-Radon-Carathéodory	20
19	An application to learning theory: VC-dimension of halfspaces.	21
20	Application to centerpoints.	23
21	2.5 Polyhedra	25
22	2.5.1 The Minkowski-Weyl Theorem	27
23	2.5.2 Valid inequalities and feasibility	28
24	2.5.3 Faces of polyhedra	30
25	2.5.4 Implicit equalities, dimension of polyhedra and facets	31
26	3 Convex Functions	33
27	3.1 General properties, epigraphs, subgradients	34
28	3.2 Continuity properties	38
29	3.3 First-order derivative properties	40
30	3.4 Second-order derivative properties	42
31	3.5 Sublinear functions, support functions and gauges	43
32	Gauges.	43
33	Support functions.	47
34	Generalized Cauchy-Schwarz/Holder's inequality.	48
35	One-to-one correspondence between closed, convex sets and closed, sublinear functions.	49
36	3.6 Directional derivatives, subgradients and subdifferential calculus	50
37	4 Optimization	54
38	Algorithmic setup: First-order oracles.	56
39	4.1 Subgradient algorithm	56
40	Subgradient Algorithm.	57
41	4.2 Generalized inequalities and convex mappings	59
42	4.3 Convex optimization with generalized inequalities	60
43	4.3.1 Lagrangian duality for convex optimization with generalized constraints	61
44	4.3.2 Solving the Lagrangian dual problem	63
45	4.3.3 Explicit examples of the Lagrangian dual	63
46	Conic optimization.	64
47	Convex optimization with explicit constraints and objective.	64
48	A closer look at linear programming duality.	65
49	4.3.4 Strong duality: sufficient conditions and complementary slackness	65
50	Slater's condition for strong duality.	65
51	Closed cone condition for strong duality in conic optimization.	66

52	Complementary slackness.	67
53	4.3.5 Saddle point interpretation of the Lagrangian dual	68
54	4.4 Cutting plane schemes	68
55	General cutting plane scheme	69
56	Center of Gravity Method.	71
57	Ellipsoid method.	72

1 Definitions and Preliminaries

We will focus on \mathbb{R}^d for arbitrary $d \in \mathbb{N}$: $\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$. 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 , $i = 1, \dots, d$ 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 $N : \mathbb{R}^d \rightarrow \mathbb{R}_+$ 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 ℓ^p norm on \mathbb{R}^d : $\|\mathbf{x}\|_p = (|x_1|^p + |x_2|^p + \dots + |x_d|^p)^{\frac{1}{p}}$. $p = 2$ is also called the *standard Euclidean norm*; we will drop the subscript 2 to denote the standard norm: $\|\mathbf{x}\| = \sqrt{x_1^2 + x_2^2 + \dots + x_d^2}$. The ℓ^∞ norm is defined as $\|\mathbf{x}\|_\infty = \max_{i=1}^d |x_i|$.

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. $d_N(\mathbf{x}, \mathbf{y}) = 0$ 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)

Definition 1.4. We also utilize the (standard) inner product of $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$: $\langle \mathbf{x}, \mathbf{y} \rangle = x_1y_1 + x_2y_2 + \dots + x_dy_d$. (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* . We will drop the subscript N when we speak of the standard Euclidean norm and there is no chance of confusion in the context.

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 denote the concept of the *infimum* by $\inf S$ and the *supremum* by $\sup S$.

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.

88 A *sequence* in \mathbb{R}^d is a countable ordered set of points: $\mathbf{x}^1, \mathbf{x}^2, \mathbf{x}^3, \dots$ and will often be denoted by $\{\mathbf{x}^i\}_{i \in \mathbb{N}}$.
 89 We say that *the sequence converges* or that *the limit of the sequence exists* if there exists a point \mathbf{x} such that
 90 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
 91 simply the *limit*, of the sequence and will also sometimes be denoted by $\lim_{n \rightarrow \infty} \mathbf{x}^n$. Although the definition
 92 of the limit is made here with respect to a particular norm, it is a well-known fact that the concept actually
 93 does not depend on the choice of the norm.

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

96 We introduce three important notions:

- 97 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
 98 $\text{cl}(X)$.
- 99 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
 100 by $\text{int}(X)$.
- 101 3. For any set $X \subseteq \mathbb{R}^d$, the *boundary* of X is defined as $\text{bd}(X) := \text{cl}(X) \setminus \text{int}(X)$.

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

103 **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
 104 convergent) has a convergent subsequence.

105 A function $f : \mathbb{R}^d \rightarrow \mathbb{R}^n$ is *continuous* if for every convergent sequence $\{\mathbf{x}^i\}_{i \in \mathbb{N}} \subseteq \mathbb{R}^d$, the following holds:
 106 $\lim_{i \rightarrow \infty} f(\mathbf{x}^i) = f(\lim_{n \rightarrow \infty} \mathbf{x}^i)$.

107 **Theorem 1.11.** [Weierstrass' Theorem] Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a continuous function. Let $X \subseteq \mathbb{R}^d$ be a
 108 nonempty, compact subset. Then $\inf\{f(\mathbf{x}) : \mathbf{x} \in X\}$ is attained, i.e., there exists $\mathbf{x}^{\min} \in X$ such that
 109 $f(\mathbf{x}^{\min}) = \inf\{f(\mathbf{x}) : \mathbf{x} \in X\}$. Similarly, there exists $\mathbf{x}^{\max} \in X$ such that $f(\mathbf{x}^{\max}) = \sup\{f(\mathbf{x}) : \mathbf{x} \in X\}$.

110 A generalization of the above theorem is the following.

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

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

Definition 1.13. We say that $f : \mathbb{R}^d \rightarrow \mathbb{R}^n$ is differentiable at $\mathbf{x} \in \mathbb{R}^d$, if there exists a linear transformation
 $A : \mathbb{R}^d \rightarrow \mathbb{R}^n$ such that

$$\lim_{\mathbf{h} \rightarrow \mathbf{0}} \frac{\|f(\mathbf{x} + \mathbf{h}) - f(\mathbf{x}) - A\mathbf{h}\|}{\|\mathbf{h}\|} = 0.$$

113 If f is differentiable at \mathbf{x} , then the linear transformation A is unique. It is commonly called the *differential*
 114 *or total derivative* of f and is denoted by $f'(\mathbf{x})$. When $n = 1$, it is commonly called *gradient* of f and is
 115 denoted by $\nabla f(\mathbf{x})$.

Definition 1.14. The partial derivative of $f : \mathbb{R}^d \rightarrow \mathbb{R}$ at \mathbf{x} in the i -th direction is defined as the real
 number

$$f'_i(\mathbf{x}) := \lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\mathbf{e}^i) - f(\mathbf{x})}{h},$$

116 if the limit exists.

117 **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
118 A will be denoted by $\text{rk}(A)$ – it is the maximum number of linearly independent rows of A , which is equal to
119 the maximum number of linearly independent columns of A . When $m = n$, we say that the matrix is *square*.

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

121 **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
122 that $A\mathbf{v} = \lambda\mathbf{v}$. λ is called the eigenvalue of A associated with \mathbf{v} .

123 **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
124 Euclidean norm, with associated eigenvalues $\lambda_1, \dots, \lambda_n \in \mathbb{R}$. Moreover, if S is the matrix whose columns
125 are $\mathbf{v}^1, \dots, \mathbf{v}^n$ and Λ is the diagonal matrix with $\lambda_1, \dots, \lambda_n$ as the diagonal entries, then $A = S\Lambda S^T$.

126 Moreover, $\text{rk}(A)$ equals the number of nonzero eigenvalues.

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

- 128 1. All eigenvalues of A are nonnegative.
- 129 2. There exists a matrix $B \in \mathbb{R}^{r \times n}$ with linearly independent rows such that $A = B^T B$.
- 130 3. $\mathbf{u}^T A \mathbf{u} \geq 0$ for all $\mathbf{u} \in \mathbb{R}^n$.

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

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

136 2 Convex Sets

137 2.1 Definitions and basic properties

138 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
139 precisely, for all $\mathbf{x}, \mathbf{y} \in X$ and every $\lambda \in [0, 1]$, $\lambda\mathbf{x} + (1 - \lambda)\mathbf{y} \in X$.

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

- 141 1. In \mathbb{R} , the only examples of convex sets are intervals (closed, open, half open): (a, b) , $(a, b]$, $[a, b]$, $(-\infty, b]$
142 etc.
- 143 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
144 $H^-(\mathbf{a}, \delta) = \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{a}, \mathbf{x} \rangle \leq \delta\}$ are all convex sets. Sets of the form $H(\mathbf{a}, \delta)$ are called *hyperplanes*
145 and sets of the form $H^+(\mathbf{a}, \delta), H^-(\mathbf{a}, \delta)$ are called *halfspaces*.
- 146 3. $\{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_\infty \leq 1\}$ is a convex set.
- 147 4. $\{\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d : x_1 + x_2 t + x_3 t^2 + \dots + x_d t^{d-1} \geq 0 \text{ for all } t \geq 0\}$ is a convex set.
- 148 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
149 convex.

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

151 **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
152 ellipsoid. In other words, an ellipsoid is the unit ball associated with the norm induced by A – see Exercise 1.

153 Exercise 2 shows that ellipsoids are convex.

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

- 155 1. Let $X_i, i \in I$ be an arbitrary family of convex sets. Then $\bigcap_{i \in I} X_i$ is a convex set.
 156 2. Let X be a convex set and $\alpha \in \mathbb{R}$, then αX is a convex set.
 157 3. Let X, Y be convex sets, then $X + Y$ is convex.
 158 4. Let $T : \mathbb{R}^d \rightarrow \mathbb{R}^m$ be any affine transformation, i.e., $T(\mathbf{x}) = A\mathbf{x} + \mathbf{b}$ for some matrix $A \in \mathbb{R}^{m \times d}$ and
 159 vector $\mathbf{b} \in \mathbb{R}^m$. 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
 160 convex.

161 *Proof.* 1. Let $\mathbf{x}, \mathbf{y} \in \bigcap_{i \in I} X_i$. This implies that $\mathbf{x}, \mathbf{y} \in X_i$ for every $i \in I$. Since each X_i is convex, for
 162 every $\lambda \in [0, 1]$, $\lambda\mathbf{x} + (1 - \lambda)\mathbf{y} \in X_i$ for all $i \in I$. Therefore, $\lambda\mathbf{x} + (1 - \lambda)\mathbf{y} \in \bigcap_{i \in I} X_i$.
 163 The proofs of 2., 3. and 4. are very similar, are left for the reader. \square

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

Definition 2.5. Let $Y = \{\mathbf{y}^1, \dots, \mathbf{y}^n\} \subset \mathbb{R}^d$ be a finite set of points. The *set of all convex combinations* of
 Y is defined as

$$\{\lambda_1 \mathbf{y}^1 + \lambda_2 \mathbf{y}^2 + \dots + \lambda_n \mathbf{y}^n : \lambda_i \geq 0, \lambda_1 + \lambda_2 + \dots + \lambda_n = 1\}.$$

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

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

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

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

170 **Definition 2.7.** Given any set $X \subseteq \mathbb{R}^d$ (not necessarily convex), the convex hull of X , denoted by $\text{conv}(X)$,
 171 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
 172 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),

$$\text{conv}(X) = \bigcap \{C : X \subseteq C, C \text{ convex}\} = \{\lambda_1 \mathbf{x}_1 + \dots + \lambda_t \mathbf{x}_t : \mathbf{x}_1, \dots, \mathbf{x}_t \in X, \lambda_1, \dots, \lambda_t \geq 0, \sum_{i=1}^t \lambda_i = 1\}.$$

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

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

178 Next, let $\tilde{C} = \{\lambda_1 \mathbf{x}_1 + \dots + \lambda_t \mathbf{x}_t : \mathbf{x}_1, \dots, \mathbf{x}_t \in X, \lambda_1, \dots, \lambda_t \geq 0, \sum_{i=1}^t \lambda_i = 1\}$. Then,

179 1. \tilde{C} is convex. Consider two points $\mathbf{z}_1, \mathbf{z}_2 \in \tilde{C}$. Thus there exist two finite index sets I_1, I_2 , two
180 finite subsets of X given by $X_1 = \{\mathbf{x}_i^1 \in X : i \in I_1\}$ and $X_2 = \{\mathbf{x}_i^2 \in X : i \in I_2\}$, and two
181 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 $\sum_{i \in I_j} \lambda_i^j = 1$
182 for $j = 1, 2$, with the following property: $\mathbf{z}_j = \sum_{i \in I_j} \lambda_i^j \mathbf{x}_i^j$ for $j = 1, 2$. Then for any $\lambda \in [0, 1]$,
183 $\lambda \mathbf{z}_1 + (1 - \lambda) \mathbf{z}_2 = \lambda (\sum_{i \in I_1} \lambda_i^1 \mathbf{x}_i^1) + (1 - \lambda) (\sum_{i \in I_2} \lambda_i^2 \mathbf{x}_i^2)$. Consider the finite set $\tilde{X} = X_1 \cup X_2$, and
184 for each $\mathbf{x} \in \tilde{X}$, if $\mathbf{x} = \mathbf{x}_i^1 \in X_1$ with $i \in I_1$ let $\mu_{\mathbf{x}} = \lambda \cdot \lambda_i^1$, and if $\mathbf{x} = \mathbf{x}_i^2 \in X_2$ with $i \in I_2$, let
185 $\mu_{\mathbf{x}} = (1 - \lambda) \cdot \lambda_i^2$. It is easy to check that $\sum_{\mathbf{x} \in \tilde{X}} \mu_{\mathbf{x}} = 1$, and $\lambda \mathbf{z}_1 + (1 - \lambda) \mathbf{z}_2 = \sum_{\mathbf{x} \in \tilde{X}} \mu_{\mathbf{x}} \mathbf{x}$. Thus,
186 $\lambda \mathbf{z}_1 + (1 - \lambda) \mathbf{z}_2 \in \tilde{C}$.

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

188 3. Let C' be any convex set such that $X \subseteq C'$. Since C' is convex, every point of the form $\lambda_1 \mathbf{x}_1 + \dots + \lambda_t \mathbf{x}_t$
189 where $\mathbf{x}_1, \dots, \mathbf{x}_t \in X$, $\lambda_i \geq 0$, $\sum_{i=1}^t \lambda_i = 1$ belongs to C' by Proposition 2.6. Thus, $\tilde{C} \subseteq C'$.

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

191 2.2 Convex cones, affine sets and dimension

192 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
193 relax the conditions on λ, γ ?

194 **Definition 2.9.** Let $X \subseteq \mathbb{R}^d$ be a nonempty set. We have three possibilities:

- 195 1. We say that $X \subseteq \mathbb{R}^d$ is a *convex cone* if for all $\mathbf{x}, \mathbf{y} \in X$ and $\lambda, \gamma \geq 0$, $\lambda \mathbf{x} + \gamma \mathbf{y} \in X$.
- 196 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
197 $\lambda + \gamma = 1$, $\lambda \mathbf{x} + \gamma \mathbf{y} \in X$.
- 198 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 $\lambda, \gamma \in \mathbb{R}$, $\lambda \mathbf{x} + \gamma \mathbf{y} \in X$.

199 **Remark 2.10.** Since we relaxed the conditions on λ, γ , convex cones, affine sets and linear sets are all
200 special cases of convex sets.

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

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

- 207 1. $\text{cone}(X) = \bigcap \{C : X \subseteq C, C \text{ is a convex cone}\} = \{\lambda_1 \mathbf{x}_1 + \dots + \lambda_t \mathbf{x}_t : \mathbf{x}_1, \dots, \mathbf{x}_t \in X, \lambda_1, \dots, \lambda_t \geq 0\}$.
- 208 2. $\text{aff}(X) = \bigcap \{C : X \subseteq C, C \text{ is an affine set}\} = \{\lambda_1 \mathbf{x}_1 + \dots + \lambda_t \mathbf{x}_t : \mathbf{x}_1, \dots, \mathbf{x}_t \in X, \sum_{i=1}^t \lambda_i = 1\}$.
- 209 3. $\text{span}(X) = \bigcap \{C : X \subseteq C, C \text{ is a linear subspace}\} = \{\lambda_1 \mathbf{x}_1 + \dots + \lambda_t \mathbf{x}_t : \mathbf{x}_1, \dots, \mathbf{x}_t \in X, \lambda_1, \dots, \lambda_t \in \mathbb{R}\}$.

210 The following is a fundamental theorem of linear algebra.

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

- 212 1. X is a linear subspace.
- 213 2. There exists $0 \leq m \leq d$ and linearly independent vectors $\mathbf{v}^1, \dots, \mathbf{v}^m \in X$ such that every $\mathbf{x} \in X$ can
214 be written as $\mathbf{x} = \lambda_1 \mathbf{v}^1 + \dots + \lambda_m \mathbf{v}^m$ for some reals λ_i , $i = 1, \dots, m$, i.e., $X = \text{span}(\{\mathbf{v}^1, \dots, \mathbf{v}^m\})$.
- 215 3. There exists a matrix $A \in \mathbb{R}^{(d-m) \times d}$ with full row rank such that $X = \{\mathbf{x} \in \mathbb{R}^d : A\mathbf{x} = \mathbf{0}\}$.

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

218 It is easy to verify that 2. \Rightarrow 1. (because linear combinations of linear combinations are linear combi-
 219 nations). To see that 1. \Rightarrow 2., starting with a linear subspace X , we construct a finite set $\mathbf{v}^1, \dots, \mathbf{v}^m \in X$
 220 satisfying the conditions of 2. We do this in an iterative fashion. Start by picking any arbitrary $\mathbf{v}^1 \in X$. If
 221 $X = \text{span}(\{\mathbf{v}^1\})$, 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
 222 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
 223 cannot have more than d linearly independent vectors in \mathbb{R}^d .

224 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
 225 is known as the *orthogonal complement* of X). It can be verified that X^\perp is a linear subspace. Moreover, by
 226 the equivalence 1. \Leftrightarrow 2., we know that 2. holds for X^\perp . So there exist linearly independent vectors $\mathbf{a}^1, \dots, \mathbf{a}^k$
 227 for some $0 \leq k \leq d$ such that $X^\perp = \text{span}(\{\mathbf{a}^1, \dots, \mathbf{a}^k\})$. Let A be the $k \times d$ matrix which has $\mathbf{a}^1, \dots, \mathbf{a}^k$ as
 228 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
 229 the number from condition 2. needs additional work, which we skip here. \square

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

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

234 **Definition 2.14.** We say a set X is affinely independent if there does not exist $\mathbf{x} \in X$ such that $\mathbf{x} \in$
 235 $\text{aff}(X \setminus \{\mathbf{x}\})$.

236 We now give several characterizations of affine independence.

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

- 238 1. X is an affinely independent set.
- 239 2. For every $\mathbf{x} \in X$, the set $\{\mathbf{v} - \mathbf{x} : \mathbf{v} \in X \setminus \{\mathbf{x}\}\}$ is linearly independent.
- 240 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.
- 241 4. The set of vectors $\{(\mathbf{x}, 1) \in \mathbb{R}^{d+1} : \mathbf{x} \in X\}$ is linearly independent.
- 242 5. X is a finite set with vectors $\mathbf{x}^1, \dots, \mathbf{x}^m$ such that $\lambda_1 \mathbf{x}^1 + \dots + \lambda_m \mathbf{x}^m = \mathbf{0}, \lambda_1 + \dots + \lambda_m = 0$ implies
 243 $\lambda_1 = \lambda_2 = \dots = \lambda_m = 0$.

244 *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
 245 not linearly independent, i.e., there exist multipliers $\lambda_{\mathbf{v}}$, not all zero, such that $\sum_{\mathbf{v} \in X \setminus \{\mathbf{x}\}} \lambda_{\mathbf{v}}(\mathbf{v} - \mathbf{x}) = \mathbf{0}$.
 246 Rearranging terms, we get $\sum_{\mathbf{v} \in X \setminus \{\mathbf{x}\}} \lambda_{\mathbf{v}} \mathbf{v} = (\sum_{\mathbf{v} \in X \setminus \{\mathbf{x}\}} \lambda_{\mathbf{v}}) \mathbf{x}$. We now consider two cases:

247 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 $\lambda_{\bar{\mathbf{v}}} \neq 0$.
 248 Since $\sum_{\mathbf{v} \in X \setminus \{\mathbf{x}\}} \lambda_{\mathbf{v}} \mathbf{v} = (\sum_{\mathbf{v} \in X \setminus \{\mathbf{x}\}} \lambda_{\mathbf{v}}) \mathbf{x} = \mathbf{0}$, 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 $\sum_{\mathbf{v} \in X \setminus \{\mathbf{x}\}} \lambda_{\mathbf{v}} =$
 249 0 , we obtain that $\sum_{\mathbf{v} \in X \setminus \{\mathbf{x}, \bar{\mathbf{v}}\}} \frac{\lambda_{\mathbf{v}}}{-\lambda_{\bar{\mathbf{v}}}} = 1$ and thus $\bar{\mathbf{v}} \in \text{aff}(X \setminus \{\mathbf{x}, \bar{\mathbf{v}}\})$, contradicting the assumption that X
 250 is affinely independent.

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

253 2. \Rightarrow 3. Obvious.

254 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
 255 $\{(\mathbf{v} - \bar{\mathbf{x}}, 0) : \mathbf{v} \in X \setminus \{\bar{\mathbf{x}}\}\} \cup \{(\bar{\mathbf{x}}, 1)\}$ are also linearly independent. Thus the matrix with these vectors as
 256 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. If $\{(\mathbf{x}, 1) \in \mathbb{R}^{d+1} : \mathbf{x} \in X\}$ is linearly independent, then the set X must be finite with elements $\mathbf{x}^1, \dots, \mathbf{x}^m$. Moreover, for any $\lambda_1, \dots, \lambda_m$ such that $\lambda_1 \mathbf{x}^1 + \dots + \lambda_m \mathbf{x}^m = \mathbf{0}$, $\lambda_1 + \dots + \lambda_m = 0$ we have $\sum_{\mathbf{x} \in X} \lambda_{\mathbf{x}} (\mathbf{x}, 1) = \mathbf{0}$. By linear independence of the set $\{(\mathbf{x}, 1) \in \mathbb{R}^{d+1} : \mathbf{x} \in X\}$, $\lambda_1 = \dots = \lambda_m = 0$.

5. \Rightarrow 1. Consider any $\mathbf{x}^i \in X$. If $\mathbf{x}^i \in \text{aff}(X \setminus \{\mathbf{x}^i\})$, 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 $\sum_{j=1}^m \lambda_j \mathbf{x}^j = \mathbf{0}$ where $\lambda_i = -1$, and therefore $\lambda_1 + \dots + \lambda_m = 0$, 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 $0 \leq m \leq d$, such that $X - \mathbf{x} = L$ for every $\mathbf{x} \in X$.
3. There exist affinely independent vectors $\mathbf{v}^1, \dots, \mathbf{v}^{m+1} \in X$ for $0 \leq m \leq d$ such that every $\mathbf{x} \in X$ can be written as $\mathbf{x} = \lambda_1 \mathbf{v}^1 + \dots + \lambda_{m+1} \mathbf{v}^{m+1}$ for some reals λ_i , $i = 1, \dots, m+1$ such that $\lambda_1 + \dots + \lambda_{m+1} = 1$, i.e., $X = \text{aff}(\{\mathbf{v}^1, \dots, \mathbf{v}^{m+1}\})$.
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 $0 \leq m \leq d$ such that $X = \{\mathbf{x} \in \mathbb{R}^d : A\mathbf{x} = \mathbf{b}\}$.

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

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. \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 nonempty convex set X , the dimension of X is the dimension of $\text{aff}(X)$ and will be denoted by $\dim(X)$. As a matter of convention, we take the dimension of the empty set to be -1 .

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. \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 for Week 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 $|X_0| - 1 = \dim(\text{aff}(X))$. \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.

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 $\bar{C} = C \cap B(\bar{\mathbf{x}}, r)$. Since C is closed and $B(\bar{\mathbf{x}}, r)$ 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 $\mathbf{x}^* \in \bar{C}$ such that $\|\mathbf{x} - \mathbf{x}^*\| \leq \|\mathbf{x} - \mathbf{y}\|$ for all $\mathbf{y} \in \bar{C}$, and therefore in fact $\|\mathbf{x} - \mathbf{x}^*\| \leq \|\mathbf{x} - \mathbf{y}\|$ for all $\mathbf{y} \in C$.

Let $\mathbf{a} = \mathbf{x} - \mathbf{x}^*$ and let $\delta = \langle \mathbf{a}, \mathbf{x}^* \rangle$. Note that $\mathbf{a} \neq \mathbf{0}$ because $\mathbf{x} \notin C$ and $\mathbf{x}^* \in C$. Also note that $\langle \mathbf{a}, \mathbf{x} \rangle = \langle \mathbf{a}, \mathbf{a} + \mathbf{x}^* \rangle = \|\mathbf{a}\|^2 + \delta > \delta$. 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 $\alpha \mathbf{y} + (1 - \alpha)\mathbf{x}^*$, $\alpha \in (0, 1)$ are in C by convexity. Therefore, by the extremal property of \mathbf{x}^* , we have

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

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

Corollary 2.21. Every closed convex set can be written as the intersection of some family of halfspaces. 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, δ^i) , $i \in I$ (where I may be an uncountable index set) such that $X = \bigcap_{i \in I} H^-(\mathbf{a}^i, \delta^i)$.

Definition 2.22. A *finite* intersection of halfspaces 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 $\mathbf{b} \in \mathbb{R}^m$.

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.

Theorem 2.23 (Supporting Hyperplane Theorem). Let $C \subseteq \mathbb{R}^d$ be a convex set and let $\mathbf{x} \in \text{bd}(C)$. 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 $\text{bd}(C) = \text{bd}(\mathbb{R}^d \setminus \text{cl}(C))$, $\mathbf{x} \in \text{bd}(\mathbb{R}^d \setminus \text{cl}(C))$. Since $\mathbb{R}^d \setminus \text{cl}(C)$ is an open set, there exists a sequence $\{\mathbf{x}^i\}_{i \in \mathbb{N}}$ such that $\mathbf{x}^i \rightarrow \mathbf{x}$ and each $\mathbf{x}^i \notin \text{cl}(C)$. 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 $i \in \mathbb{N}$.

334 Since the set of unit norm vectors is a compact set, by Theorem 1.10, one can pick a convergent sub-
 335 sequence $\mathbf{a}^{i_k} \rightarrow \mathbf{a}$ 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
 336 $\langle \mathbf{a}, \mathbf{y} \rangle \leq \langle \mathbf{a}, \mathbf{x} \rangle$ 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
 337 have $\|\mathbf{a}\| = 1$, and so $\mathbf{a} \neq \mathbf{0}$. \square

338 **How to represent general convex sets: Separation oracles.** We have seen that polyhedra can be
 339 represented by a matrix A and a right hand side b . Norm balls can be represented by the center \mathbf{x} and
 340 the radius R . Ellipsoids can be represented by positive definite matrices A . What about general convex
 341 sets? This problem is gotten around by assuming that one has “black-box” access to the convex set via a
 342 *separation oracle*. More formally, we say that a convex set $C \subseteq \mathbb{R}^d$ is equipped with a separation oracle O
 343 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
 344 $\mathbf{x} \notin C$, 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
 345 for \mathbf{x} and C .

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

349 **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.

- 350 1. $A\mathbf{x} = \mathbf{b}$ has a solution.
- 351 2. There exists $\mathbf{u} \in \mathbb{R}^d$ such that $\mathbf{u}^T A = \mathbf{0}$ and $\mathbf{u}^T \mathbf{b} \neq 0$.

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

354 **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.

- 355 1. $A\mathbf{x} = \mathbf{b}$, $\mathbf{x} \geq \mathbf{0}$ has a solution.
- 356 2. There exists $\mathbf{u} \in \mathbb{R}^d$ such that $\mathbf{u}^T A \leq \mathbf{0}$ and $\mathbf{u}^T \mathbf{b} > 0$.

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

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

359 *Proof.* We will complete the proof of this lemma when we do Caratheodory’s theorem (see the end of
 360 Section 2.4). \square

361 *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
 362 $C = \{A\mathbf{x} : \mathbf{x} \geq \mathbf{0}\}$ is closed. We now have two cases, either $\mathbf{b} \in C$ or $\mathbf{b} \notin C$. In the first case, we end up in
 363 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 $\delta \in \mathbb{R}$
 364 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
 365 shows that $\langle \mathbf{u}, \mathbf{b} \rangle > 0$.

366 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 $\bar{\lambda} \langle \mathbf{u}, \mathbf{a}^i \rangle > \delta$
 367 (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 $\langle \mathbf{u}, \mathbf{y} \rangle \leq \delta$
 368 for all $\mathbf{y} \in C$. \square

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

371 **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
 372 *orthogonal complement* of X .

373 The following is well-known from linear algebra.

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

375 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

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

376 is called the *polar* of X .

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

- 378 1. X° is a closed, convex set for any $X \subseteq \mathbb{R}^d$ (not necessarily convex).
 379 2. $(X^\circ)^\circ = \text{cl}(\text{conv}(X \cup \{\mathbf{0}\}))$. In particular, if X is a closed convex set containing the origin, then
 380 $(X^\circ)^\circ = X$.
 381 3. If X is a convex cone, then $X^\circ = \{\mathbf{y} \in \mathbb{R}^d : \langle \mathbf{y}, \mathbf{x} \rangle \leq 0 \quad \forall \mathbf{x} \in X\}$.
 382 4. If X is a linear subspace, then $X^\circ = X^\perp$.

Proof. 1. Follows from the fact that X° can be written as the intersection of closed halfspaces:

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

- 383 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
 384 a closed convex set by 1., we must have $\text{cl}(\text{conv}(X \cup \{\mathbf{0}\})) \subseteq (X^\circ)^\circ$.

385 To show the reverse inclusion, we show that if $\mathbf{y} \notin \text{cl}(\text{conv}(X \cup \{\mathbf{0}\}))$ then $\mathbf{y} \notin (X^\circ)^\circ$. Thus, we need
 386 to show that there exists $\mathbf{z} \in X^\circ$ such that $\langle \mathbf{y}, \mathbf{z} \rangle > 1$. Since $\mathbf{y} \notin \text{cl}(\text{conv}(X \cup \{\mathbf{0}\}))$, by Theorem 2.20,
 387 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 $\mathbf{x} \in \text{cl}(\text{conv}(X \cup \{\mathbf{0}\}))$. Since
 388 $\mathbf{0} \in \text{cl}(\text{conv}(X \cup \{\mathbf{0}\}))$, we obtain that $0 \leq \delta$. We now consider two cases:

389 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 $\mathbf{x} \in \text{cl}(\text{conv}(X \cup$
 390 $\{\mathbf{0}\})) \supseteq X$. 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.

391 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
 392 $\mathbf{x} \in X \subseteq \text{cl}(\text{conv}(X \cup \{\mathbf{0}\}))$, 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
 393 are done.

- 394 3. and 4. are left to the reader. □

395

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

Proposition 2.32 (Hölder's inequality). If $p, q \geq 1$ such that $\frac{1}{p} + \frac{1}{q} = 1$ (allowing for p or q to be ∞), 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 Hölder's inequality here, but we will use it to derive the polarity relation between ℓ_p unit balls. We only show that $B_{\ell^q}(\mathbf{0}, 1) = (B_{\ell^p}(\mathbf{0}, 1))^\circ$ 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 $B_{\ell^q}(\mathbf{0}, 1) \subseteq (B_{\ell^p}(\mathbf{0}, 1))^\circ$. Consider any $\mathbf{y} \in B_{\ell^q}(\mathbf{0}, 1)$ and consider any $\mathbf{x} \in B_{\ell^p}$. By Hölder's inequality, we obtain that $\langle \mathbf{x}, \mathbf{y} \rangle \leq \|\mathbf{x}\|_p \|\mathbf{y}\|_q \leq 1$. Thus, $B_{\ell^q}(\mathbf{0}, 1) \subseteq (B_{\ell^p}(\mathbf{0}, 1))^\circ$. To show the reverse inclusion $(B_{\ell^p}(\mathbf{0}, 1))^\circ \subseteq B_{\ell^q}(\mathbf{0}, 1)$, consider any $\mathbf{y} \in (B_{\ell^p}(\mathbf{0}, 1))^\circ$. 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 $i = 1, \dots, d$, \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 $\mathbf{y} \in (B_{\ell^p}(\mathbf{0}, 1))^\circ$, because $\|\tilde{\mathbf{x}}\|_p = 1$. The second equality follows from Proposition 2.32 because of the special choice of \mathbf{x} .

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 $\text{relbd}(C) := \text{cl}(C) \setminus \text{relint}(C)$.

¹For 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 $\text{aff}(C)$.

417 **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
 418 $\mathbf{x} + \epsilon_{\mathbf{y}}(\mathbf{y} - \mathbf{x}) \in C$. Show that $\mathbf{x} \in \text{relint}(C)$.

419 This exercise shows that it suffices to have a different ϵ for every direction; this implies a universal ϵ for
 420 every direction.

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

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

423 *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
 424 points in F . We claim that $X \cup \{\mathbf{x}\}$ is affinely independent. This immediately implies that $\dim(C) > \dim(F)$
 425 and we will be done.

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

429 **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}
 430 is an extreme point of C .

431 *Proof.* Left to the reader. \square

432 **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,
 433 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

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

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

435 **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 $\delta \in \mathbb{R}$ such
 436 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
 437 say that F is an *exposed face induced by* (\mathbf{a}, δ) .

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

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

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

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

444 (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
 445 in the hyperplane $\{\mathbf{y} \in \mathbb{R}^d : \langle \mathbf{a}, \mathbf{y} \rangle = \delta\}$.

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

448 **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
 449 only if $\mathbf{x} \in \text{relbd}(C)$.

450 *Proof.* If $\mathbf{x} \in \text{relbd}(C)$, then by Theorem 2.39 there exists $\mathbf{a} \in \mathbb{R}^d$ and $\delta \in \mathbb{R}$ such that the three conditions
 451 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 a proper face
 452 because of condition (iii) in Theorem 2.39.

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

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

453 Since $\mathbf{x} \in F$ and $\bar{\mathbf{x}} \notin F$, this contradicts the fact that F is a face. Thus, $\mathbf{x} \notin \text{relint}(C)$ and so $\mathbf{x} \in$
 454 $\text{relbd}(C)$. \square

455 In our search for a subset $X \subseteq C$ such that $C = \text{conv}(X)$, it is clear that X must contain all extreme
 456 points. But is it sufficient to include all extreme points? In other words, is it true that $C = \text{conv}(\text{ext}(C))$?
 457 No! A simple counterexample is \mathbb{R}_+^d . Its only extreme point is $\mathbf{0}$. Another weird example is the set
 458 $\{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\| < 1\}$ – this set has NO extreme points! As you might suspect, the problem is that these sets
 459 are not compact, i.e., closed and bounded.

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

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

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

466 *Case 1: $\mathbf{x} \in \text{relbd}(C)$.* By Theorem 2.40, \mathbf{x} is contained in a proper face F of C . By Lemma 2.35, $\dim(F) <$
 467 $\dim(C)$. By the induction hypothesis applied to F (note that F is also compact using Problem 14 from “HW
 468 for Week III”), we can express \mathbf{x} as a convex combination of extreme points of F , which by Lemma 2.36,
 469 shows that \mathbf{x} is a convex combination of extreme points of C .

470 *Case 2: $\mathbf{x} \in \text{relint}(C)$.* Let $\ell \subseteq \text{aff}(C)$ be any affine set of dimension one (i.e., a line) going through \mathbf{x} . Since
 471 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
 472 of C . By the previous case, $\mathbf{x}^1, \mathbf{x}^2$ can be expressed as the convex combination of extreme points in C .
 473 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
 474 combination, we can express \mathbf{x} as the convex combination of extreme points of C . \square

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

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

- 479 1. There exists $\mathbf{x} \in C$ such that $\mathbf{x} + \lambda \mathbf{r} \in C$ for all $\lambda \geq 0$.
- 480 2. For every $\mathbf{x} \in C$, $\mathbf{x} + \lambda \mathbf{r} \in C$ for all $\lambda \geq 0$.

Proof. Since C is nonempty, we only need to show 1. \Rightarrow 2.; the reverse implication is trivial. Let $\bar{\mathbf{x}} \in C$ be such
 that $\bar{\mathbf{x}} + \lambda \mathbf{r} \in C$ for all $\lambda \geq 0$. Consider any arbitrary $\mathbf{x}^* \in C$. Suppose to the contrary that there exists $\lambda' \geq 0$
 such that $\mathbf{y} = \mathbf{x}^* + \lambda' \mathbf{r} \notin C$. 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, $\langle \mathbf{a}, \mathbf{y} \rangle = \langle \mathbf{a}, \mathbf{x}^* \rangle + \lambda' \langle \mathbf{a}, \mathbf{r} \rangle \leq \delta + \lambda' \langle \mathbf{a}, \mathbf{r} \rangle \leq \delta$
 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

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

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

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

484 **Proposition 2.44.** The set of all recession directions of a nonempty, closed, convex set is a closed, convex
 485 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,

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

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

488 **Definition 2.45.** Let C be any nonempty, closed, convex set. We call the cone of recession directions
 489 the *recession cone* of C and it is denoted by $\text{rec}(C)$. The set $\text{rec}(C) \cap -\text{rec}(C)$ is a linear subspace and
 490 is called the *lineality space* of C . It will be denoted by $\text{lin}(C)$. As a matter of convention, we say that
 491 $\text{rec}(C) = \text{lin}(C) = \{\mathbf{0}\}$ when C is empty.

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

494 Proposition 2.42 immediately gives the following corollary.

495 **Corollary 2.46.** Let C be a closed convex set and let $F \subseteq C$ be a closed, convex subset. Then $\text{rec}(F) \subseteq$
 496 $\text{rec}(C)$.

497 *Proof.* Left as an exercise. \square

498 Here is a characterization of compact convex sets.

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

Proof. We leave it to the reader to check that if C is compact, then $\text{rec}(C) = \{\mathbf{0}\}$. For the other direction, assume that $\text{rec}(C) = \{\mathbf{0}\}$. Suppose to the contrary that C is not bounded, i.e., there exists a sequence of points $\mathbf{y}^i \in C$ such that $\|\mathbf{y}^i\| \rightarrow \infty$. 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 $\text{rec}(C) = \{\mathbf{0}\}$. 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

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

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

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

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

503 The main result about pointed closed convex sets says that you can decompose them into convex combi-
 504 nations of extreme points and recession directions.

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

506 *Proof.* The proof follows the same lines as Theorem 2.41. We may assume C is nonempty since otherwise
 507 $\text{ext}(C) = \{\}$. We prove by induction on dimension of C . If $\dim(C) = 0$, then C is a single point, and we are
 508 done. Consider any $\mathbf{x} \in C$ and then two cases:

509 *Case 1: $\mathbf{x} \in \text{relbd}(C)$.* By Theorem 2.40, \mathbf{x} is contained in a proper face F of C . By Lemma 2.35, $\dim(F) <$
 510 $\dim(C)$. By the induction hypothesis applied to F (note that F is also closed using Problem 14 from “HW
 511 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
 512 \mathbf{d} is a recession direction for F . By Lemma 2.36, \mathbf{x}' is a convex combination of extreme points of C . By
 513 Corollary 2.46, $\mathbf{d} \in \text{rec}(C)$.

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

517 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
 518 combination of extreme points in C , and $\mathbf{d}^i \in \text{rec}(C)$. Since \mathbf{x} is a convex combination of \mathbf{x}^1 and \mathbf{x}^2 , this
 519 shows that $\mathbf{x} \in \text{conv}(\text{ext}(C)) + \text{rec}(C)$.

520 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
 521 extreme points in C , and $\mathbf{d}' \in \text{rec}(C)$. Thus, $\mathbf{x} = \mathbf{y}' + \mathbf{d}' + \mathbf{d}$ and we have the desired representation. \square

522 Lets make this description even more “minimal”. For this we will need to understand the structure of
 523 pointed cones.

524 **Proposition 2.50.** Let $D \subseteq \mathbb{R}^d$ be a closed, convex cone. The following are equivalent.

- 525 1. D is pointed.
- 526 2. D° is full-dimensional, i.e., $\dim(D^\circ) = d$.
- 527 3. $\mathbf{0}$ is an exposed face of D .
- 528 4. There exists a compact, convex subset $B \subset D \setminus \{\mathbf{0}\}$ such that every $\mathbf{d} \in D \setminus \{\mathbf{0}\}$ can be uniquely
 529 written in the form $\mathbf{d} = \lambda \mathbf{b}$, where $\mathbf{b} \in B$ and $\lambda > 0$. In particular, $D = \text{cone}(B)$.

530 *Proof.* 1. \Rightarrow 2. If D° is not full-dimensional, then $\text{aff}(D^\circ)$ is a linear space of dimension strictly less than d ,
 531 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.
 532 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
 533 space, this implies that $\text{aff}(D^\circ)^\perp \subseteq \text{lin}(D)$, contradicting the assumption that D is pointed.

534 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 $D^\circ = \{\mathbf{y} \in$
 535 $\mathbb{R}^d : \langle \mathbf{x}, \mathbf{y} \rangle \leq 0 \ \forall \mathbf{x} \in D\}$, using Problem 3 from “HW for Week II”, we obtain that $\langle \mathbf{y}, \mathbf{x} \rangle < 0$ for every
 536 $\mathbf{x} \in D \setminus \{\mathbf{0}\}$. This shows that the exposed face induced by $(\mathbf{y}, 0)$ is exactly $\{\mathbf{0}\}$.

537 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
 538 from the definition that $\mathbf{0} \notin B$. Since it is the intersection of a convex cone and an affine set, B is also convex.
 539 We now show that B is compact. It is the intersection of closed sets, so it is closed. By Theorem 2.47, it
 540 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
 541 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,
 542 we obtain that $\mathbf{0} + \mathbf{r} \in D$, i.e., $\mathbf{r} \in D$. But then $\langle \mathbf{y}, \mathbf{r} \rangle = 0$ contradicts the fact that $\mathbf{0}$ is an exposed face of
 543 D induced by $(\mathbf{y}, 0)$.

We next consider any $\mathbf{d} \in D \setminus \{\mathbf{0}\}$. 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,

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

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

545 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 $\lambda_1 > 0$
 546 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
 547 in B , contradicting the assumption. \square

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

550 The proof of Proposition 2.50 also shows the following.

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

553 **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
 554 some hyperplane H . We skip the proof of this fact from these notes.

555 **Definition 2.54.** Let D be a closed, convex cone. An edge of D is called an *extreme ray* of D . We say that
 556 $\mathbf{r} \in D$ 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
 557 by $\text{extr}(D)$.

558 **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
 559 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 $\mathbf{r}^2 = \lambda_2 \mathbf{r}$.

560 *Proof.* Left as an exercise. □

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

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

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

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

569 Now we can write a sharper version of Theorem 2.49:

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

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

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

575 *Proof.* Define $\hat{C} = C \cap \text{lin}(C)^\perp$. \hat{C} is closed because it is the intersection of two closed sets. By Corollary 2.46,
 576 $\text{rec}(\hat{C}) \subseteq \text{rec}(C)$. Therefore, $\text{lin}(\hat{C}) = \text{rec}(\hat{C}) \cap -\text{rec}(\hat{C}) \subseteq \text{rec}(C) \cap -\text{rec}(C) = \text{lin}(C)$. By the same reasoning,
 577 $\text{lin}(\hat{C}) \subseteq \text{lin}(\text{lin}(C)^\perp) = \text{lin}(C)^\perp$. Since $\text{lin}(C) \cap \text{lin}(C)^\perp = \{\mathbf{0}\}$, we obtain that $\text{lin}(\hat{C}) = \{\mathbf{0}\}$. □

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

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

578 *Proof.* We first observe that $C = \hat{C} + \text{lin}(C)$. Indeed, for any $\mathbf{x} \in C$, we can express $\mathbf{x} = \mathbf{x}' + \mathbf{r}$ where
 579 $\mathbf{x}' \in \text{lin}(C)^\perp$ and $\mathbf{r} \in \text{lin}(C)$ (since $\text{lin}(C) + \text{lin}(C)^\perp = \mathbb{R}^n$). We also know that $\mathbf{x}' = \mathbf{x} - \mathbf{r} \in C$ because
 580 $\mathbf{r} \in \text{lin}(C)$. Thus, $\mathbf{x}' \in \hat{C}$ and we are done. \hat{C} is pointed by Lemma 2.58 and applying Corollary 2.57 gives
 581 the desired result. □

582 Thus, a general closed convex set C can be specified by giving a set of generators for its lineality space
 583 $\text{lin}(C)$, and the extreme points and vectors spanning the extreme rays of the set $C \cap \text{lin}(C)^\perp$. In Section 2.5,
 584 we will see that polyhedra are precisely those convex sets C that have a finite number of extreme points and
 585 extreme rays for $C \cap \text{lin}(C)^\perp$. So we see that polyhedra are especially easy to describe intrinsically: simply
 586 specify the finite list of extreme points, vectors spanning the extreme rays and a finite list of generators of
 587 $\text{lin}(C)$.

588 **2.3.3 A remark about extrinsic and intrinsic descriptions**

589 You may have already observed that although a closed convex set can be represented as the intersection of
 590 halfspaces, such a representation is not unique. For example, consider the circle in \mathbb{R}^2 . You can represent
 591 it by intersecting all its tangent halfspaces. On the other hand, if you throw away any finite subset of
 592 these halfspaces, you still get the same set. In fact, there is a representation which uses only countably
 593 many halfspaces. Thus, the same convex set can have many different representations as the intersection of
 594 halfspaces. Moreover, there is usually no way to choose a “canonical” representation, i.e., there is no set of
 595 representing halfspaces such that *any representation* will always include this “canonical” set of halfspaces
 596 (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 $\text{conv}(\text{ext}(C))$. 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 $C = \text{conv}(X)$, we must have $\text{ext}(C) \subseteq X$.
 With closed, convex sets that have a nontrivial recession cone, the situation is a bit more subtle. First, there
 is more flexibility in choosing the representation because one can choose a different set of vectors to span
 the extreme rays. One might think that this is just a scaling issue and the following result holds: if C is a
 pointed, closed, convex set, and we consider any “intrinsic” representation

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

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

- 598 (i) $\text{ext}(C) \subseteq E$ and
- 599 (ii) for every \mathbf{r} that spans an extreme ray of $\text{rec}(C)$, there must be some nonnegative scaling of \mathbf{r} present
 600 in R .

601 While the above holds for polyhedra and many other closed, convex sets, it is not true in general. We
 602 leave it as an exercise to find a closed, convex set that violates the above claim.

603 **2.4 Combinatorial theorems: Helly-Radon-Carathéodory**

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

606 **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
 607 as $X = X_1 \uplus X_2$ into 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 $\{\mathbf{x}^1, \dots, \mathbf{x}^k\} \subseteq X$ such that $\{\mathbf{x}^1, \dots, \mathbf{x}^k\}$ is affinely
 dependent. By using characterization 5. in Proposition 2.15, there exist multipliers $\lambda_1, \dots, \lambda_k \in \mathbb{R}$, not all
 zero, such that $\lambda_1 + \dots + \lambda_k = 0$ and $\lambda_1 \mathbf{x}^1 + \dots + \lambda_k \mathbf{x}^k = \mathbf{0}$. Define $P := \{i : \lambda_i \geq 0\}$ and $N := \{i : \lambda_i < 0\}$. Since
 the λ_i ’s are not all zero and $\lambda_1 + \dots + \lambda_k = 0$, 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 $\lambda_1 + \dots + \lambda_k = 0$, and $\sum_{j \in P} \lambda_j \mathbf{x}^j = \sum_{j \in N} (-\lambda_j) \mathbf{x}^j$
 since $\lambda_1 \mathbf{x}^1 + \dots + \lambda_k \mathbf{x}^k = \mathbf{0}$. Thus, we obtain that

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

608 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
 609 simply define $X_1 = X_P$ and $X_2 = X \setminus X_P$. □

610

An application to learning theory: VC-dimension of halfspaces. An important concept in learning theory is the *Vapnik-Červonenkis (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 $X \subseteq \mathbb{R}^d$ 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

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

611

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 from “HW for Week II”, 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 $X = X_1 \uplus X_2$ 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 . \square

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

612

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

613

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

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 $X_1, X_2, \dots, X_{\bar{k}+1}$. Define a new family $C_1, \dots, C_{\bar{k}}$, where $C_i = X_i$ if $i \leq \bar{k} - 1$ and $C_{\bar{k}} = X_{\bar{k}} \cap X_{\bar{k}+1}$. Since $\emptyset = X_1 \cap \dots \cap X_{\bar{k}+1} = C_1 \cap \dots \cap C_{\bar{k}}$, we can use the induction hypothesis on this new family and obtain a subfamily C_{i_1}, \dots, C_{i_m} such that $C_{i_1} \cap \dots \cap C_{i_m} = \emptyset$ and $m \leq d + 1$. If $m \leq d$ or none of the C_{i_h} , $h = 1, \dots, m$ equals $C_{\bar{k}}$, then we are done. So we assume that $m = d + 1$ and $C_{i_m} = C_{\bar{k}} = X_{\bar{k}} \cap X_{\bar{k}+1}$.

618

To simplify notation, let us relabel everything and define $D_h := C_{i_h} = X_{i_h}$, $h = 1, \dots, d$ and $D_{d+1} = X_{\bar{k}}$ and $D_{d+2} = X_{\bar{k}+1}$. We thus know that $D_1 \cap \dots \cap D_{d+2} = \emptyset$. We may assume that each subfamily of $d + 1$ sets from D_1, \dots, D_{d+2} has a nonempty intersection, because otherwise we will be done. Let these common intersection points be

$$\mathbf{x}^i \in \bigcap_{h \neq i} D_h, \quad i = 1, \dots, d + 2.$$

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

625

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

626

632 **Corollary 2.64.** Let \mathcal{X} be a (possibly infinite) family of compact, convex sets. If $\bigcap_{X \in \mathcal{X}} X = \emptyset$, then there
633 is a subfamily X_{i_1}, \dots, X_{i_m} for some $m \leq d + 1$, with $i_h \in \{1, \dots, k\}$ for each $h = 1, \dots, m$ such that
634 $X_{i_1} \cap \dots \cap X_{i_m} = \emptyset$. Thus, there is a subfamily of size at most $d + 1$ that already certifies the empty
635 intersection.

636 *Proof.* By a standard result in topology, if the intersection of an infinite family of compact sets is empty,
637 then there is a finite subfamily whose intersection is also empty. One can now apply Theorem 2.63 to this
638 finite subfamily and obtain a subfamily of size at most $d + 1$. \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 μ 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 to μ as any \mathbf{x} in the set $C_\mu := \arg \max_{\mathbf{x} \in \mathbb{R}^d} f_\mu(\mathbf{x})$. It can be shown that this set is nonempty for all probability distributions μ . For $d = 1$, this gives the standard notion of a median, and one can show that for any probability distribution μ 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 μ 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}_α be the set of all halfspaces H such that $\mu(H) \geq \alpha$. It is not hard to check that if $\alpha < 1$, then $D_\alpha := \bigcap_{H \in \mathcal{H}_\alpha} H$ is a compact, convex set. Indeed, for any coordinate indexed by $i = 1, \dots, d$, there must exist some δ_1^i, δ_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_α is contained in the box $\{\mathbf{x} \in \mathbb{R}^d : \delta_2^i \leq \mathbf{x}_i \leq \delta_1^i, i = 1, \dots, d\}$.

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 α . This would imply that H' contains all of D_α 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 and the fact that $D_\alpha \subseteq D_\beta$ when $\alpha \leq \beta$, we would have $\bigcap_{\epsilon > 0} D_{\frac{d}{d+1} + \epsilon}$ 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 $D_\alpha := \bigcap_{H \in \mathcal{H}_\alpha} H$ is compact, as observed earlier.

Define $\mathcal{C} = \{B \cap H : H \text{ is a closed halfspace with } \mu(H) \geq \frac{d}{d+1} + \epsilon\}$. 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 $\{C_1, \dots, C_{d+1}\} \subseteq \mathcal{C}$ of size $d + 1$, we claim

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

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,

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

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

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

644 **Theorem 2.66** (Carathéodory's Theorem – cone version). Let $X \subseteq \mathbb{R}^d$ (not necessarily convex) and let
 645 $\mathbf{x} \in \text{cone}(X)$. There exists a subset $X' \subseteq X$ such that X' is linearly independent (and thus, $|X'| \leq d$), and
 646 $\mathbf{x} \in \text{cone}(X')$.

Proof. Since $\mathbf{x} \in \text{cone}(X)$, by Theorem 2.11, we can find a finite set $\{\mathbf{x}^1, \dots, \mathbf{x}^k\} \subseteq X$ such that $\mathbf{x} \in \text{cone}(\{\mathbf{x}^1, \dots, \mathbf{x}^k\})$. Choose a minimal such set, i.e., there is no strict subset of $\{\mathbf{x}^1, \dots, \mathbf{x}^k\}$ whose conical hull contains \mathbf{x} . This implies that $\mathbf{x} = \lambda_1 \mathbf{x}^1 + \dots + \lambda_k \mathbf{x}^k$ for some $\lambda_i > 0$ for each $i = 1, \dots, k$. We claim that $\mathbf{x}^1, \dots, \mathbf{x}^k$ are linearly independent. Suppose to the contrary that there exist multipliers $\gamma_1, \dots, \gamma_k \in \mathbb{R}$, not all zero, such that $\gamma_1 \mathbf{x}^1 + \dots + \gamma_k \mathbf{x}^k = \mathbf{0}$. By changing the signs of the γ_i 's if necessary, we may assume that there exists $j \in \{1, \dots, k\}$ such that $\gamma_j > 0$. Define

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

Observe that $\lambda'_i \geq 0$ for all $i = 1, \dots, k$ and

$$\lambda'_1 \mathbf{x}^1 + \dots + \lambda'_k \mathbf{x}^k = \lambda_1 \mathbf{x}^1 + \dots + \lambda_k \mathbf{x}^k - \theta(\gamma_1 \mathbf{x}^1 + \dots + \gamma_k \mathbf{x}^k) = \lambda_1 \mathbf{x}^1 + \dots + \lambda_k \mathbf{x}^k = \mathbf{x}.$$

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

649 **Theorem 2.67** (Carathéodory's Theorem – convex version). Let $X \subseteq \mathbb{R}^d$ (not necessarily convex) and let
 650 $\mathbf{x} \in \text{conv}(X)$. There exists a subset $X' \subseteq X$ such that X' is affinely independent (and thus, $|X'| \leq d + 1$),
 651 and $\mathbf{x} \in \text{conv}(X')$.

652 *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
 653 to saying that $(\mathbf{x}, 1) \in \text{cone}(Y)$. We get the desired result by applying Theorem 2.66 and condition 4. of
 654 Proposition 2.15. □

655 We can finally furnish the proof of Lemma 2.26.

Proof of Lemma 2.26. Consider a convergent sequence $\{\mathbf{x}^i\}_{i \in \mathbb{N}} \subseteq \text{cone}(\{\mathbf{a}^1, \dots, \mathbf{a}^n\})$ 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 $\{\mathbf{a}^1, \dots, \mathbf{a}^n\}$. Since there are only finitely many linearly independent subsets of $\{\mathbf{a}^1, \dots, \mathbf{a}^n\}$, one of these subsets contains infinitely many elements of the sequence $\{\mathbf{x}^i\}_{i \in \mathbb{N}}$. Thus, after passing to that subsequence, we may assume that $\{\mathbf{x}^i\}_{i \in \mathbb{N}} \subseteq \text{cone}(\{\bar{\mathbf{a}}^1, \dots, \bar{\mathbf{a}}^k\})$ where $\bar{\mathbf{a}}^1, \dots, \bar{\mathbf{a}}^k$ are linearly independent. For each \mathbf{x}^i , there exists $\boldsymbol{\lambda}^i \in \mathbb{R}_+^k$ such that $\mathbf{x}^i = \boldsymbol{\lambda}^i \bar{\mathbf{a}}^1 + \dots + \boldsymbol{\lambda}^i_k \bar{\mathbf{a}}^k$. If we denote by $A \in \mathbb{R}^{d \times k}$ the matrix whose columns are $\bar{\mathbf{a}}^1, \dots, \bar{\mathbf{a}}^k$, then $\mathbf{x}^i = A \boldsymbol{\lambda}^i$ and $\boldsymbol{\lambda}^i = A^{-1} \mathbf{x}^i$ for every $i \in \mathbb{N}$. 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 it is the image of a bounded set under the linear (and therefore continuous) map A^{-1} . Thus, by Theorem 1.10 there is a convergent subsequence $\boldsymbol{\lambda}^{i_k} \rightarrow \boldsymbol{\lambda} \in \mathbb{R}_+^k$. Taking limits,

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

656 Since $\boldsymbol{\lambda} \in \mathbb{R}_+^k$, we find that $\mathbf{x} \in \text{cone}(\{\bar{\mathbf{a}}^1, \dots, \bar{\mathbf{a}}^k\}) \subseteq \text{cone}(\{\mathbf{a}^1, \dots, \mathbf{a}^n\})$. □

657 Here is another result that proves handy in many situations.

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

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

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

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 $\text{conv}(X)$ is the image of $\underbrace{X \times \dots \times X}_{d+1 \text{ times}} \times \Delta^{d+1}$ under f , where

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

659 Since X and Δ^{d+1} are compact sets, we obtain the result by applying Theorem 1.12. □

660 2.5 Polyhedra

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

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

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

667 Both these systems describe the same polyhedron $P = \text{conv}\{(0, 0, 0), (0, 0, 5)\}$ in \mathbb{R}^3 . However, if a polyhedron
 668 is given by its list of extreme points and extreme rays, this ambiguity disappears. Moreover, having
 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 $\mathbf{a}^1, \dots, \mathbf{a}^m$ 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 $\text{tight}(\mathbf{x}, A, \mathbf{b}) := \{i : \langle \mathbf{a}^i, \mathbf{x} \rangle = \mathbf{b}_i\}$. For brevity, when
 674 A and \mathbf{b} are clear from the context, we will shorten this to $\text{tight}(\mathbf{x})$. We also use the notation $A_{\text{tight}(\mathbf{x})}$ to
 675 denote the submatrix formed by taking the rows of A indexed by $\text{tight}(\mathbf{x})$. Similarly, $\mathbf{b}_{\text{tight}(\mathbf{x})}$ will denote
 676 the subvector of \mathbf{b} indexed by $\text{tight}(\mathbf{x})$.

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 $A_{\text{tight}(\mathbf{x})}$ has rank equal to d , i.e., the rows of A indexed by
 679 $\text{tight}(\mathbf{x})$ span \mathbb{R}^d .

Proof. (\Leftarrow) Suppose $A_{\text{tight}(\mathbf{x})}$ 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 $i \in \text{tight}(\mathbf{x})$, $\langle \mathbf{a}^i, \mathbf{x}^1 \rangle \leq \mathbf{b}_i$ and similarly, $\langle \mathbf{a}^i, \mathbf{x}^2 \rangle \leq \mathbf{b}_i$. Now,
 we observe that

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

680 Thus, the inequality must be an equality. Therefore, for each $i \in \text{tight}(\mathbf{x})$, $\langle \mathbf{a}^i, \mathbf{x}^1 \rangle = \mathbf{b}_i$ and similarly,
 681 $\langle \mathbf{a}^i, \mathbf{x}^2 \rangle = \mathbf{b}_i$. In other words, we have that $A_{\text{tight}(\mathbf{x})}\mathbf{x} = \mathbf{b}_{\text{tight}(\mathbf{x})}$, and $A_{\text{tight}(\mathbf{x})}\mathbf{x}^j = \mathbf{b}_{\text{tight}(\mathbf{x})}$ for $j = 1, 2$.
 682 Since the rank of $A_{\text{tight}(\mathbf{x})}$ is 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 $A_{\text{tight}(\mathbf{x})}$ 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 $A_{\text{tight}(\mathbf{x})}\mathbf{r} = \mathbf{0}$. Define

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

684 Note that $\epsilon > 0$ because whenever $\langle \mathbf{a}^j, \mathbf{r} \rangle \neq 0$ we have that $j \notin \text{tight}(\mathbf{x})$ and thus all the numerators are
 685 strictly positive. 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 $\mathbf{x} = \frac{\mathbf{x}^1 + \mathbf{x}^2}{2}$
 686 with $\mathbf{x}^1 \neq \mathbf{x}^2$ (because $\mathbf{r} \neq \mathbf{0}$ and $\epsilon > 0$), contradicting extremality.

687 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 –
 688 the calculations for \mathbf{x}^2 are similar. Consider any $j \in \{1, \dots, m\}$. If $j \in \text{tight}(\mathbf{x})$, then since $A_{\text{tight}(\mathbf{x})}\mathbf{r} = \mathbf{0}$,
 689 we obtain that $\langle \mathbf{a}^j, \mathbf{x}^1 \rangle = \langle \mathbf{a}^j, \mathbf{x} \rangle + \epsilon \langle \mathbf{a}^j, \mathbf{r} \rangle = \langle \mathbf{a}^j, \mathbf{x} \rangle = \mathbf{b}_j$. If $j \notin \text{tight}(\mathbf{x})$, then we consider two cases:

690 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 $\langle \mathbf{a}^j, \mathbf{x}^1 \rangle = \langle \mathbf{a}^j, \mathbf{x} \rangle + \epsilon \langle \mathbf{a}^j, \mathbf{r} \rangle \leq \mathbf{b}_j$.

691 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 $\langle \mathbf{a}^j, \mathbf{r} \rangle < 0$. \square

692 This immediately gives the following.

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

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

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

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

701 **Proposition 2.73.** Let $D \subseteq \mathbb{R}^d$ be a convex cone. D is a polyhedral cone if and only if there exists a matrix
 702 $A \in \mathbb{R}^{m \times d}$ for some $m \in \mathbb{N}$ such that $D = \{\mathbf{x} \in \mathbb{R}^d : A\mathbf{x} \leq \mathbf{0}\}$.

703 *Proof.* We simply have to show the forward direction, the reverse is easy. Assume D is a polyhedral cone.
 704 Thus, it is a 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
 705 $D = \{\mathbf{x} : A\mathbf{x} \leq \mathbf{b}\}$. Since D is a closed, convex cone (closed because all polyhedra are closed), $\text{rec}(D) = D$.
 706 By Problem 1 in “HW for Week IV”, we obtain that $D = \text{rec}(D) = \{\mathbf{x} : A\mathbf{x} \leq \mathbf{0}\}$. \square

707 Problem 1 in “HW for Week IV” also immediately implies the following.

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

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

711 *Proof.* (\Leftarrow) Let $A_{\text{tight}(\mathbf{r})}$ have rows $\bar{\mathbf{a}}^1, \dots, \bar{\mathbf{a}}^k$. Each $F_i := D \cap \{\mathbf{x} : \langle \bar{\mathbf{a}}^i, \mathbf{x} \rangle = 0\}$ for each $i = 1, \dots, k$ is an
 712 exposed face of D . By Problem 13 in “HW for Week III”, $F := \bigcap_{i=1}^k F_i$ is a face of D . Since $A_{\text{tight}(\mathbf{r})}$ has
 713 rank $d - 1$, the set $\{\mathbf{x} : A_{\text{tight}(\mathbf{r})}\mathbf{x} = \mathbf{0}\}$ is a 1-dimensional linear subspace. Since $F \subseteq \{\mathbf{x} : A_{\text{tight}(\mathbf{r})}\mathbf{x} = \mathbf{0}\}$,
 714 F is a 1-dimensional face of D (it cannot be 0 dimensional because it contains $\{\mathbf{0}\}$ and $\mathbf{r} \neq \mathbf{0}$) and hence an
 715 extreme ray. Since $\mathbf{r} \in F$, we have that \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 $A_{\text{tight}(\mathbf{r})} \leq d - 1$. If it is

strictly less, then consider any $\mathbf{r}' \in \{\mathbf{x} : A_{\text{tight}(\mathbf{r})}\mathbf{x} = \mathbf{0}\}$ that is linearly independent to \mathbf{r} – such an \mathbf{r}' exists if rank of $A_{\text{tight}(\mathbf{r})} \leq d - 2$. Define

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

716 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
 717 $\mathbf{r} = \frac{\mathbf{r}^1 + \mathbf{r}^2}{2}$. Moreover, since \mathbf{r}' and \mathbf{r} are linearly independent, $\mathbf{r}^1, \mathbf{r}^2$ are not scalings of \mathbf{r} . This contradicts
 718 Proposition 2.55.

719 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
 720 in the proof of Theorem 2.70. \square

721 Analogous to Corollary 2.71, we have:

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

723 2.5.1 The Minkowski-Weyl Theorem

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

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

727 *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
 728 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 VI”, \hat{P} is also
 729 a polyhedron. By Corollary 2.71, we obtain that $V := \text{ext}(\hat{P})$ is a finite set. Moreover, by Proposition 2.74,
 730 $\text{rec}(\hat{P})$ is a polyhedral cone. By Corollary 2.76, $\text{extr}(\text{rec}(\hat{P}))$ is a finite set. Define $R = \text{extr}(\text{rec}(\hat{P})) \cup L \cup -L$.
 731 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

732 We now make an observation about polars.

733 **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.

734 *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
 735 for Week II” we obtain that $X = \text{conv}(V) + \text{conv}(R)$ is closed. Since the Minkowski sum of convex sets is
 736 convex (property 3. in Theorem 2.3), X is also convex. \square

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

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

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

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

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

738 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 $\langle \mathbf{v}^i, \mathbf{y} \rangle \leq 1$
 739 for $i = 1, \dots, k$ since $\mathbf{v}^i \in X$. Suppose to the contrary that $\langle \mathbf{r}^j, \mathbf{y} \rangle > 0$ for some $j \in \{1, \dots, n\}$. Then there
 740 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
 741 $\mathbf{v}^1 + \lambda \mathbf{r}^j \in X$, by definition of X . Therefore, $\langle \mathbf{r}^j, \mathbf{y} \rangle \leq 0$ for $j = 1, \dots, n$ and thus, $\mathbf{y} \in \tilde{X}$. \square

742 This has the following corollary.

743 **Corollary 2.80.** Let P be a polyhedron. Then P° is a polyhedron.

744 *Proof.* If $P = \emptyset$, then $P^\circ = \mathbb{R}^d$, which is a polyhedron. Else, by Theorem 2.77, there exist finite sets
 745 $V, R \subseteq \mathbb{R}^d$ such that $P = \text{conv}(V) + \text{cone}(R)$, with $V \neq \emptyset$. By Theorem 2.79, P° is the intersection of finitely
 746 many halfspaces, and is thus a polyhedron. \square

747 We now prove the converse of Theorem 2.77.

748 **Theorem 2.81** (Minkowski-Weyl Theorem – Part II). Let $V, R \subseteq \mathbb{R}^d$ be finite sets and let $X = \text{conv}(V) +$
 749 $\text{cone}(R)$. Then $X \subseteq \mathbb{R}^d$ is a polyhedron.

750 *Proof.* The case when X is empty is trivial. So we consider X is nonempty. Take any $\mathbf{t} \in X$ and define
 751 $X' = X - \mathbf{t}$. Now, it is easy to see X is a polyhedron if and only if X' is a polyhedron (Verify!!). So it
 752 suffices to show that X' is a polyhedron. Note that $X' = \text{conv}(V') + \text{cone}(R)$ where $V' = V - \mathbf{t}$, which is
 753 a nonempty set because V is nonempty (since X is assumed to be nonempty). By Theorem 2.79, $(X')^\circ$ is
 754 a polyhedron. By Lemma 2.78, X' is a closed, convex set, and also $\mathbf{0} \in X'$. Therefore, $X' = ((X')^\circ)^\circ$ by
 755 condition 2. in Theorem 2.30. Applying Corollary 2.80 with $P = (X')^\circ$, we obtain that $((X')^\circ)^\circ = X'$ is a
 756 polyhedron. \square

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

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

759 (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 $X = \{\mathbf{x} \in$
 760 $\mathbb{R}^d : \mathbf{A}\mathbf{x} \leq \mathbf{b}\}$.

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

762 A compact version is often useful.

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

765 *Proof.* Left as an exercise. \square

766 2.5.2 Valid inequalities and feasibility

767 **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
 768 *valid inequality/halfspace* for X if $X \subseteq H^-(\mathbf{a}, \delta)$.

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

772 **Theorem 2.85.** Let $P = \{\mathbf{x} \in \mathbb{R}^d : \mathbf{A}\mathbf{x} \leq \mathbf{b}\}$ with $A \in \mathbb{R}^{m \times d}, \mathbf{b} \in \mathbb{R}^m$ be a nonempty polyhedron. Let
 773 $\mathbf{c} \in \mathbb{R}^d, \delta \in \mathbb{R}$. 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
 774 $\mathbf{c}^T = \mathbf{y}^T A$ 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 $\mathbf{x} \in P$:

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

775 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}, λ :

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

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

$$\begin{bmatrix} A^T & \mathbf{0} \\ \mathbf{b}^T & 1 \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ \lambda \end{bmatrix} = \begin{bmatrix} \mathbf{c} \\ \delta \end{bmatrix}.$$

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

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

777 We now consider two cases:

778 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
779 "HW for Week IV", $\bar{\mathbf{u}} \in \text{rec}(P)$. Consider any $\mathbf{x} \in P$ (we assume P is nonempty). Let $\mu = \frac{1 + \langle \mathbf{c}, \bar{\mathbf{u}} \rangle}{\langle \mathbf{c}, \bar{\mathbf{u}} \rangle} > 0$.
780 Now $\mathbf{x} + \mu \bar{\mathbf{u}} \in P$ since $\bar{\mathbf{u}} \in \text{rec}(P)$. 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
781 inequality for P .

782 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
783 $\mathbf{x} = \frac{\bar{\mathbf{u}}}{-\mathbf{u}_{d+1}}$, we obtain that $A\mathbf{x} \leq \mathbf{b}$ and $\langle \mathbf{c}, \mathbf{x} \rangle > \delta$, contradicting that $\langle \mathbf{c}, \mathbf{x} \rangle \leq \delta$ is a valid inequality for
784 P . \square

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

787 **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
788 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
789 think of Theorem 2.85 is that it says the geometric property of being a valid inequality is the same as the
790 algebraic property of being a consequence:

791 [Alternate version of Theorem 2.85] Let $P = \{\mathbf{x} \in \mathbb{R}^d : A\mathbf{x} \leq \mathbf{b}\}$ be a nonempty polyhedron.
792 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
793 of $A\mathbf{x} \leq \mathbf{b}$.

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

796 **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
797 only if $\langle \mathbf{0}, \mathbf{x} \rangle \leq -1$ is a consequence of $A\mathbf{x} \leq \mathbf{b}$.

798 *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
799 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

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

800 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
801 $A\mathbf{x} \leq \mathbf{b}$. \square

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

802 **2.5.3 Faces of polyhedra**

803 Faces for polyhedra are very structured. Firstly, every face is an exposed face – something that is not true
 804 for general closed, convex sets. Secondly, there is an algebraic characterization of faces in terms of the
 805 describing inequalities of a polyhedron. This is the content of Theorem 2.90 below. First, we make some
 806 simpler observations.

807 **Theorem 2.89.** The following are both true.

- 808 1. Every face of a polyhedron is a polyhedron.
 809 2. Every polyhedron has finitely many faces.

810 *Proof.* We prove the theorem for the case of bounded polyhedra, i.e., polytopes. The extension of this proof
 811 idea to the unbounded case is left as an exercise.

812 Let P be any polytope and let $F \subseteq P$ be a face. If $\mathbf{x} \in F$ is an extreme point of F , then by Problem 15
 813 in “HW for Week IV”, $\{\mathbf{x}\}$ is a face of P and therefore \mathbf{x} is an extreme point of P . Since P has finitely
 814 many extreme points by Corollary 2.71, F has only finitely many extreme points. By Problem 14 in “HW
 815 for Week III”, F is closed, and since P is compact, so is F . By the Krein-Milman theorem (Theorem 2.41),
 816 F is the convex hull of its finitely many extreme points. By the Minkowski-Weyl theorem (Theorem 2.81),
 817 F is a polyhedron. Moreover, since we showed that any face is the convex hull of some subset of extreme
 818 points of P , there can only be finitely many faces since P has finitely many extreme points. \square

819 **Theorem 2.90.** 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
 820 following are equivalent.

- 821 (i) F is a face of P .
 822 (ii) F is an exposed face of P .
 823 (iii) There exists a subset $I \subseteq \{1, \dots, m\}$ such that $F = \{\mathbf{x} \in P : A_I \mathbf{x} = \mathbf{b}_I\}$.

824 *Proof.* (i) \Rightarrow (ii). Consider $\bar{\mathbf{x}} \in \text{relint}(F)$ (which exists by Exercise 4). Since F is a proper face, by
 825 Theorem 2.40, $\bar{\mathbf{x}} \in \text{relbd}(P)$. By Theorem 2.39, there exists a supporting hyperplane at $\bar{\mathbf{x}}$ given by $\langle \mathbf{a}, \mathbf{x} \rangle \leq \delta$.
 826 Let $\{\mathbf{y} \in P : \langle \mathbf{a}, \mathbf{y} \rangle = \delta\}$ be the corresponding exposed face. Since $\mathbf{x} \in \text{relint}(F)$, one can show that
 827 $F \subseteq \{\mathbf{y} \in P : \langle \mathbf{a}, \mathbf{y} \rangle = \delta\}$ (Verify!!). Thus, there exists an exposed face containing F . Let F' be the minimal
 828 (with respect to set inclusion) exposed face of P that contains F , i.e., for any other exposed face $F'' \supseteq F$,
 829 we have $F' \subseteq F''$. Note that such a minimal exposed face exists because we have only finitely many faces
 830 by Theorem 2.89. Let this exposed face F' be defined by the valid inequality $\langle \mathbf{c}^1, \mathbf{x} \rangle \leq \delta_1$ for P .

831 If $F = F'$, then we are done because F' is an exposed face. Otherwise, $F \subsetneq F'$, and so F is a face of
 832 F' . Therefore, $\bar{\mathbf{x}} \in \text{relbd}(F')$. 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
 833 $F \subseteq F' \cap \{\mathbf{y} \in \mathbb{R}^d : \langle \mathbf{c}^2, \mathbf{y} \rangle = \delta_2\}$, 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
 834 find finite sets V, R such that $P = \text{conv}(V) + \text{cone}(R)$. Notice that since $P \subseteq H^-(\mathbf{c}^1, \delta_1)$, we must have
 835 $\langle \mathbf{c}^1, \mathbf{v} \rangle \leq \delta_1$ 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

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

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

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

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 $\mathbf{r} \in \text{rec}(F')$ and therefore $\langle \mathbf{c}^2, \mathbf{r} \rangle \leq 0$. Thus, the following choice of

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

837 satisfies (2.2). □

838 Using the λ 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
 839 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
 840 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
 841 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' .

842 (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
 843 $\mathbf{y} \in \mathbb{R}_+^m$ 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).
 844 Then

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

845 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
 846 that $\mathbf{y}^T (A \mathbf{x} - \mathbf{b}) = \mathbf{0}$ for all $\mathbf{x} \in F$. This last relation says that for any $i \in \{1, \dots, m\}$, if $y_i > 0$ then
 847 $\langle \mathbf{a}^i, \mathbf{x} \rangle = \mathbf{b}_i$ for every $\mathbf{x} \in F$. Thus, setting $I = \{i : y_i > 0\}$, we immediately obtain that $A_I \mathbf{x} = \mathbf{b}_I$ for all
 848 $\mathbf{x} \in F$. 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 $y_i = 0$ for $i \notin I$. Therefore,
 849 $\mathbf{c}^T \bar{\mathbf{x}} = \mathbf{y}^T A \bar{\mathbf{x}} = \mathbf{y}^T \mathbf{b} = \delta$, and thus, $\bar{\mathbf{x}} \in P \cap \{\mathbf{x} : \langle \mathbf{c}, \mathbf{x} \rangle = \delta\} = F$.

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

853 2.5.4 Implicit equalities, dimension of polyhedra and facets

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

856 **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 $i \in \{1, \dots, m\}$
 857 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
 858 denote the subsystem of implicit equalities of $A \mathbf{x} \leq \mathbf{b}$ by $A^\# \mathbf{x} = \mathbf{b}^\#$. We will also use $A^+ \mathbf{x} \leq \mathbf{b}^+$ to denote
 859 the inequalities in $A \mathbf{x} \leq \mathbf{b}$ that are NOT implicit equalities.

860 Note that for each i such that $\langle \mathbf{a}^i, \mathbf{x} \rangle \leq \mathbf{b}_i$ is not an implicit equality, there exists $\mathbf{x} \in P$ such that
 861 $\langle \mathbf{a}^i, \mathbf{x} \rangle < \mathbf{b}_i$.

862 **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 $A^+ \bar{\mathbf{x}} < \mathbf{b}^+$.
 863 Show the stronger statement that $\text{relint}(P) = \{\mathbf{x} \in \mathbb{R}^d : A^\# \mathbf{x} = \mathbf{b}^\#, A^+ \mathbf{x} < \mathbf{b}^+\}$.

864 We can completely characterize the affine hull of a polyhedron, and consequently its dimension, in terms
 865 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

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

Proof. It is easy to verify that $\text{aff}(P) \subseteq \{\mathbf{x} \in \mathbb{R}^d : A^\# \mathbf{x} = \mathbf{b}^\#\} \subseteq \{\mathbf{x} \in \mathbb{R}^d : A^\# \mathbf{x} \leq \mathbf{b}^\#\}$. We show that
 $\{\mathbf{x} \in \mathbb{R}^d : A^\# \mathbf{x} \leq \mathbf{b}^\#\} \subseteq \text{aff}(P)$. 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 $\mathbf{y} \in P \subseteq \text{aff}(P)$ and we are done. Otherwise, set

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

866 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
 867 check that $(1 - \mu)\bar{\mathbf{x}} + \mu\mathbf{y} \in P$. This shows that $\mathbf{y} \in \text{aff}(P)$, because \mathbf{y} is on the line joining two points in P ,
 868 namely $\bar{\mathbf{x}}$ and $(1 - \mu)\bar{\mathbf{x}} + \mu\mathbf{y}$. \square

869 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

$$\dim(P) = d - \text{rank}(A^\#).$$

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

873 **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 $i \in \{1, \dots, m\}$
 874 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
 875 without row i and \mathbf{b}_{-i} is the vector \mathbf{b} with the i -th coordinate removed. Otherwise, if $P \subsetneq \{\mathbf{x} : A_{-i}\mathbf{x} \leq \mathbf{b}_{-i}\}$,
 876 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
 877 every inequality is irredundant for $P = \{\mathbf{x} : A\mathbf{x} \leq \mathbf{b}\}$.

878 The following characterization of facets of a polyhedron is quite useful, specially in combinatorial opti-
 879 mization and polyhedral combinatorics.

880 **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
 881 system. Let $F \subseteq P$. The following are equivalent.

- 882 (i) F is a facet of P , i.e., F is a face with $\dim(F) = \dim(P) - 1$.
- 883 (ii) F is a maximal, proper face of P , i.e., for any proper face $F' \supseteq F$, we must have $F' = F$.
- 884 (iii) There exists a unique $i \in \{1, \dots, m\}$ 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
 885 implicit equality.

886 *Proof.* (i) \Rightarrow (ii). Suppose to the contrary that there exists a proper face $F' \supseteq F$. Observe that F is
 887 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,
 888 $\dim(F') > \dim(F) = \dim(P) - 1$. So, $\dim(F') = \dim(P)$. This contradicts the fact that F' is proper face,
 889 by Lemma 2.35.

890 (ii) \Rightarrow (iii). By Theorem 2.90, there exists a subset of indices $I \subseteq \{1, \dots, m\}$ such that $F = \{\mathbf{x} \in$
 891 $\mathbb{R}^d : A\mathbf{x} \leq \mathbf{b}, A_I\mathbf{x} = \mathbf{b}_I\}$. If all the inequalities indexed by I are implicit equalities for P , then $F = P$,
 892 contradicting the assumption that F is a proper face. So there exists $i \in I$ such that $\langle \mathbf{a}^i, \mathbf{x} \rangle \leq \mathbf{b}_i$ is
 893 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
 894 $\langle \mathbf{a}^i, \mathbf{x} \rangle \leq \mathbf{b}_i$ is not an implicit equality, F' is a proper face of P . Also observe that $F \subseteq F'$. Hence
 895 $F = F' = \{\mathbf{x} \in P : \langle \mathbf{a}^i, \mathbf{x} \rangle = \mathbf{b}_i\}$. To show uniqueness of i , we would like to show that $I = \{i\}$. We show
 896 this by exhibiting $\mathbf{x}^0 \in F$ with the following property: for any $j \neq i$ such that $\langle \mathbf{a}^j, \mathbf{x} \rangle \leq \mathbf{b}_j$ is not an implicit
 897 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
 898 \mathbf{x}^1 exists by Exercise 6). Since $A\mathbf{x} \leq \mathbf{b}$ is an irredundant system, if we remove the inequality indexed by i ,
 899 then we get some new points that satisfy the rest of the inequalities, but which violate $\langle \mathbf{a}^i, \mathbf{x} \rangle \leq \mathbf{b}_i$. More
 900 precisely, there exists $\mathbf{x}^2 \in \mathbb{R}^d$ such that $A^\# \mathbf{x}^2 = \mathbf{b}^\#$, $A_{-i}^+ \mathbf{x}^2 \leq \mathbf{b}_{-i}^+$ and $\langle \mathbf{a}^i, \mathbf{x}^2 \rangle > \mathbf{b}_i$, where $A_{-i}^+ \mathbf{x} \leq \mathbf{b}_{-i}^+$
 901 denotes the system $A^+ \mathbf{x} \leq \mathbf{b}^+$ without the inequality indexed by i . Since $\langle \mathbf{a}^i, \mathbf{x}^1 \rangle < \mathbf{b}_i$ and $\langle \mathbf{a}^i, \mathbf{x}^2 \rangle > \mathbf{b}_i$,
 902 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
 903 $A^\# \mathbf{x}^1 = \mathbf{b}^\#$ 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 $A_{-i}^+ \mathbf{x}^2 \leq \mathbf{b}_{-i}^+$, we
 904 must have that for any $j \neq i$ indexing an inequality in $A^+ \mathbf{x} \leq \mathbf{b}^+$, \mathbf{x}^0 must satisfy $\langle \mathbf{a}^j, \mathbf{x}^0 \rangle < \mathbf{b}_j$. Thus, we
 905 are done.

906 (iii) \Rightarrow (i). By Theorem 2.90, F is a face. We now establish that $\dim(F) = \dim(P) - 1$. Let \mathcal{J} denote
 907 the set of indices that index inequalities in $A\mathbf{x} \leq \mathbf{b}$ that are not implicit equalities. Since there exists a

908 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 $\mathbf{x}^j \in F$
 909 such that $\langle \mathbf{a}^j, \mathbf{x}^j \rangle < \mathbf{b}_j$. Now let $\mathbf{x}^0 = \frac{1}{|\mathcal{J}|-1} \sum_{j \in \mathcal{J} \setminus \{i\}} \mathbf{x}^j$, and observe that $\mathbf{x}^0 \in F$ (since F is convex) and
 910 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 $\tilde{A}\mathbf{x} \leq \tilde{\mathbf{b}}$ that
 911 appends the inequality $\langle -\mathbf{a}^i, \mathbf{x} \rangle \leq -\mathbf{b}_i$ to the system $A\mathbf{x} \leq \mathbf{b}$.

912 **Claim 2.** $\text{rank}(\tilde{A}^\ominus) = \text{rank}(A^\ominus) + 1$.

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

918 Using Corollary 2.93, we obtain that $\dim(F) = d - \text{rank}(\tilde{A}^\ominus) = d - \text{rank}(A^\ominus) - 1 = \dim(P) - 1$. \square

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

921 **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
 922 exists $\lambda \geq 0$ such that $\mathbf{a}' = \lambda \mathbf{a}$ and $\delta' = \lambda \delta$. Equivalent inequalities define the same halfspace, i.e.,
 923 $H^-(\mathbf{a}, \delta) = H^-(\mathbf{a}', \delta')$.

Theorem 2.97. Let P be a full-dimensional polyhedron. Let $A \in \mathbb{R}^{m \times d}$, $A' \in \mathbb{R}^{p \times d}$, $\mathbf{b} \in \mathbb{R}^m$ 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.,

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

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

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

928 *Proof.* Left as an exercise. \square

929 3 Convex Functions

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

- 935 • $x + \infty = \infty$ for any $x \in \mathbb{R} \cup \{+\infty\}$.
- 936 • $x(+\infty) = +\infty$ for all $x > 0$. We will avoid situations where we need to consider $0 \cdot (+\infty)$.
- 937 • $x < \infty$ for all $x \in \mathbb{R}$.

938 3.1 General properties, epigraphs, subgradients

Definition 3.1. A function $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ 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

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

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

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

941 **Proposition 3.2.** Let $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ be a convex function. Then $\text{dom}(f)$ is a convex set.

942 *Proof.* Left as an exercise. □

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

Definition 3.3. A function $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ is called *strongly convex* with *modulus of strong convexity* $c > 0$ if

$$f(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y}) - \frac{1}{2}c\lambda(1 - \lambda)\|\mathbf{x} - \mathbf{y}\|^2,$$

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

945 The following proposition sheds some light on strongly convex functions.

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

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

Definition 3.5. Let $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ be any function (not necessarily convex). The *epigraph* of f is defined as

$$\text{epi}(f) := \{(\mathbf{x}, t) \in \mathbb{R}^n \times \mathbb{R} : f(\mathbf{x}) \leq t\}.$$

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

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

957 *Proof.* (\Rightarrow) Consider any $(\mathbf{x}^1, t_1), (\mathbf{x}^2, t_2) \in \text{epi}(f)$, and any $\lambda \in (0, 1)$.

958 The result is a consequence of the following sequence of implications:

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

959 (\Leftarrow) Consider the any $\mathbf{x}^1, \mathbf{x}^2 \in \mathbb{R}^d$ and $\lambda \in (0, 1)$. We wish to show that $f(\lambda \mathbf{x}^1 + (1 - \lambda)\mathbf{x}^2) \leq \lambda f(\mathbf{x}^1) + (1 - \lambda)f(\mathbf{x}^2)$.
 960 If $f(\mathbf{x}^1) = +\infty$ or $f(\mathbf{x}^2) = +\infty$, then relation holds trivially. So we assume $f(\mathbf{x}^1), f(\mathbf{x}^2) < +\infty$.
 961 The points $(\mathbf{x}^1, f(\mathbf{x}^1)), (\mathbf{x}^2, f(\mathbf{x}^2))$ both lie in $\text{epi}(f)$. By convexity of $\text{epi}(f)$, we have that $(\lambda \mathbf{x}^1 + (1 - \lambda)\mathbf{x}^2, \lambda f(\mathbf{x}^1) + (1 - \lambda)f(\mathbf{x}^2)) \in \text{epi}(f)$. This implies that $f(\lambda \mathbf{x}^1 + (1 - \lambda)\mathbf{x}^2) \leq \lambda f(\mathbf{x}^1) + (1 - \lambda)f(\mathbf{x}^2)$,
 962 showing that f is convex. □

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

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

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

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

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

969 The following can be verified by the reader.

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

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

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

$$I_X(\mathbf{x}) := \begin{cases} 0 & \text{if } \mathbf{x} \in X \\ +\infty & \text{if } \mathbf{x} \notin X \end{cases}$$

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

974 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*
 975 *function* (if $\delta = 0$, this is a *linear function*). It is easily verified that affine functions are convex.

3. *Norms and Distances.* Let $N : \mathbb{R}^d \rightarrow \mathbb{R}$ 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

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

976 is a convex function.

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

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

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,

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

977 The inequality follows from the fact that if ℓ_1, \dots, ℓ_m and u_1, \dots, u_m are two sets of m real numbers
 978 for some $m \in \mathbb{N}$, then $\max_{i=1, \dots, m} (\ell_i + u_i) \leq \max_{i=1, \dots, m} \ell_i + \max_{i=1, \dots, m} u_i$.

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

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

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

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

982 (\Rightarrow) If any $f(\mathbf{x}^i)$ is $+\infty$, then the inequality holds trivially. So we assume that each $f(\mathbf{x}^i) < +\infty$. By
 983 Proposition 3.6, $\text{epi}(f)$ is a convex set. For each $i = 1, \dots, m$, the point $(\mathbf{x}^i, f(\mathbf{x}^i)) \in \text{epi}(f)$ by definition of
 984 $\text{epi}(f)$. 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., $(\lambda_1 \mathbf{x}^1 + \dots + \lambda_n \mathbf{x}^n, \lambda_1 f(\mathbf{x}^1) + \dots + \lambda_n f(\mathbf{x}^n)) \in$
 985 $\text{epi}(f)$. Therefore, $f(\lambda_1 \mathbf{x}^1 + \dots + \lambda_n \mathbf{x}^n) \leq \lambda_1 f(\mathbf{x}^1) + \dots + \lambda_n f(\mathbf{x}^n)$. \square

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

988 **Theorem 3.12.** [Operations that preserve the property of being a (closed) convex function] Let $f_i : \mathbb{R}^d \rightarrow$
 989 $\mathbb{R} \cup \{+\infty\}$, $i \in I$ be a family of (closed) convex functions where the index set I is potentially infinite. The
 990 following are all true.

- 991 1. (Nonnegative combinations). If I is a finite set, and $\alpha_i \geq 0$, $i \in I$ is a corresponding set of nonnegative
 992 reals, then $\sum_{i \in I} \alpha_i f_i$ is a (closed) convex function.
- 993 2. (Taking supremums). The function defined as $g(\mathbf{x}) := \sup_{i \in I} f_i(\mathbf{x})$ is a (closed) convex function (even
 994 when I is uncountable infinite).
- 995 3. (Pre-Composition with an affine function). Let $A \in \mathbb{R}^{m \times d}$ and $\mathbf{b} \in \mathbb{R}^m$ and let $f : \mathbb{R}^m \rightarrow \mathbb{R}$ be any
 996 (closed) convex function on \mathbb{R}^m . Then $g(\mathbf{x}) := f(A\mathbf{x} + \mathbf{b})$ as a function from $\mathbb{R}^d \rightarrow \mathbb{R}$ is a (closed)
 997 convex function.
- 998 4. (Post-Composition with an increasing convex function). Let $h : \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$ be a (closed) convex
 999 function that is also increasing, i.e., $h(x) \geq h(y)$ when $x \geq y$. Let $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ be a (closed)
 1000 convex function. We adopt the convention that $h(+\infty) = +\infty$. Then $h(f(\mathbf{x}))$ as a function from
 1001 $\mathbb{R}^d \rightarrow \mathbb{R}$ 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

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

1002 We use the nonnegativity of α_i in the inequality on the second displayed line above. We omit the proof
 1003 of closedness of the function.

- 1004 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
 1005 $i \in I$. Since the intersection of (closed) convex sets is a (closed) convex set (part 1. of Theorem 2.3),
 1006 we have the result.
- 1007 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 $(A\mathbf{x} + \mathbf{b}, t) \in \text{epi}(f)$.
 1008 Define the affine map $T : \mathbb{R}^d \times \mathbb{R} \rightarrow \mathbb{R}^m \times \mathbb{R}$ as follows: $T(\mathbf{x}, t) = (A\mathbf{x} + \mathbf{b}, t)$. Then $\text{epi}(g) = T^{-1}(\text{epi}(f))$.
 1009 Since the pre-image of a (closed) convex set with respect to an affine transformation is (closed) convex
 1010 (part 4. of Theorem 2.3), we obtain that $\text{epi}(g)$ is (closed) convex.
- 1011 4. Left as an exercise.

1012 \square

1013 We can now see some more interesting examples of convex functions.

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

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

1014 is closed convex. This is an alternate proof of the convexity of the maximum of finitely many affine
 1015 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 $f_k : V \rightarrow \mathbb{R}$ 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

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

1016 where Ω is the set of $n \times k$ matrices with k orthonormal columns in \mathbb{R}^n . This shows that f_k is the
1017 supremum of linear functions, and by Theorem 3.12, it is closed convex.

1018 We see in part 1. of Example 3.13 that the supremum of affine functions is convex. We will show below
1019 that, in fact, every convex function is the supremum of some family of affine functions. This is analogous
1020 to the fact that all closed convex sets are the intersection of some family of halfspaces. We build up to this
1021 with an important definition.

1022 **Definition 3.14.** Let $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ be any function. Let $\mathbf{x} \in \text{dom}(f)$. Then $\mathbf{a} \in \mathbb{R}^d$ is said to define
1023 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 $\mathbf{y} \in \mathbb{R}^d$.

1024 A useful picture to keep in mind is the following fact: $\mathbf{a} \in \mathbb{R}^d$ is an affine support of f at \mathbf{x} if and only if
1025 the hyperplane $\langle \mathbf{a}, \mathbf{y} \rangle - t \leq \langle \mathbf{a}, \mathbf{x} \rangle - f(\mathbf{x})$ is a supporting hyperplane for the epigraph of f at $(\mathbf{x}, f(\mathbf{x}))$.

1026 **Theorem 3.15.** Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be any function. Then f is closed convex if and only if there exists an
1027 affine support of f at every $\mathbf{x} \in \mathbb{R}^d$.

1028 *Proof.* (\Rightarrow) Consider any $\mathbf{x} \in \mathbb{R}^d$. By definition of closed convex, $\text{epi}(f)$ is a closed convex set. Moreover,
1029 $(\mathbf{x}, f(\mathbf{x})) \in \text{bd}(\text{epi}(f))$. 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
1030 both zero, and $\langle \bar{\mathbf{a}}, \mathbf{y} \rangle + rt \leq \delta$ for all $(\mathbf{y}, t) \in \text{epi}(f)$, and $\langle \bar{\mathbf{a}}, \mathbf{x} \rangle + rf(\mathbf{x}) = \delta$.

1031 We claim that $r < 0$. Suppose to the contrary that $r \geq 0$. First consider the case that $\bar{\mathbf{a}} = \mathbf{0}$, then
1032 $r > 0$. $(\mathbf{x}, t) \in \text{epi}(f)$ for all $t \geq f(\mathbf{x})$. But this contradicts that $rt = \langle \bar{\mathbf{a}}, \mathbf{y} \rangle + rt \leq \delta$ for all $t \geq f(\mathbf{x})$
1033 and $rf(\mathbf{x}) = \langle \bar{\mathbf{a}}, \mathbf{x} \rangle + rf(\mathbf{x}) = \delta$. Next consider the case that $\bar{\mathbf{a}} \neq \mathbf{0}$. Consider any $\mathbf{y} \in \mathbb{R}^d$ satisfying
1034 $\langle \bar{\mathbf{a}}, \mathbf{y} \rangle > \delta$. Since f is real valued, there exists $(\mathbf{y}, t) \in \text{epi}(f)$ for some $t \geq 0$. Since $r \geq 0$, this contradicts
1035 that $\langle \bar{\mathbf{a}}, \mathbf{y} \rangle + rt \leq \delta$.

1036 Now set $\mathbf{a} = \frac{\bar{\mathbf{a}}}{-r}$. $\langle \bar{\mathbf{a}}, \mathbf{x} \rangle + rf(\mathbf{x}) = \delta$ and $\langle \bar{\mathbf{a}}, \mathbf{y} \rangle + rf(\mathbf{y}) \leq \delta$ for all $\mathbf{y} \in \mathbb{R}^d$ together imply that
1037 $\langle \bar{\mathbf{a}}, \mathbf{y} \rangle \leq (-r)f(\mathbf{y}) + \langle \bar{\mathbf{a}}, \mathbf{x} \rangle + rf(\mathbf{x})$. 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,

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

1038 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
1039 functions, which by Example 3.13, shows that f is closed convex. \square

1040 **Remark 3.16.** 1. Any convex function that is finite valued everywhere is closed convex. This follows
1041 from a continuity result we will prove later. We skip the details in these notes. Thus, in the forward
1042 direction of Theorem 3.15, one may weaken the hypothesis to just convex, as opposed to closed convex.

1043 2. In the reverse direction of Theorem 3.15, one may weaken the hypothesis to having *local* affine support
1044 everywhere. A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is said to have local affine support at \mathbf{x} if there exists $\epsilon > 0$
1045 (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
1046 extension of Theorem 3.15 here. See Chapter on “Convex Functions” in [3].

1047 3. The proof of Theorem 3.15 also shows that if $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ and $\mathbf{x} \in \text{dom}(f)$, then there exists
1048 an affine support of f at \mathbf{x} .

1049 Affine supports for convex functions have been given a special name.

1050 **Definition 3.17.** Let $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ be a convex function. For any $\mathbf{x} \in \text{dom}(f)$, an affine support at
 1051 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
 1052 *subdifferential* of f at \mathbf{x} .

1053 **Theorem 3.18.** Let $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ be a convex function. For any $\mathbf{x} \in \text{dom}(f)$, the subdifferential
 1054 $\partial f(\mathbf{x})$ at \mathbf{x} is a closed, convex set.

Proof. Note that

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

1055 Since the above set is the intersection of a family of halfspaces, this shows that $\partial f(\mathbf{x})$ is a closed, convex
 1056 set. \square

1057 3.2 Continuity properties

1058 Convex functions enjoy strong continuity properties in the relative interior of their domains³. This fact is
 1059 very useful in many contexts, especially in optimization, because this is useful in showing that minimizers
 1060 and maximizers exist when optimizing convex functions that show up in practice, via Weierstrass' theorem
 1061 (Theorem 1.11).

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

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

1062 hold for all \mathbf{x} in the ball $B(\mathbf{x}^*, 2\epsilon)$. Then for all $\mathbf{x}, \mathbf{y} \in B(\mathbf{x}^*, \epsilon)$, it holds that

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

1063 In particular, f is locally Lipschitz continuous about \mathbf{x}^* .

1064 *Proof.* Take $\mathbf{x}, \mathbf{y} \in B(\mathbf{x}^*, \epsilon)$ with $\mathbf{x} \neq \mathbf{y}$. Define $\mathbf{z} = \mathbf{y} + \epsilon \left(\frac{\mathbf{y} - \mathbf{x}}{\|\mathbf{y} - \mathbf{x}\|} \right)$. Note that

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

1065 Thus $\mathbf{z} \in B(\mathbf{x}^*, 2\epsilon)$. Also,

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

1066 showing that \mathbf{y} is a convex combination of \mathbf{x} and \mathbf{z} . Therefore we may apply the convexity of f to see

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

1067 Hence $f(\mathbf{y}) - f(\mathbf{x}) \leq \left(\frac{\|\mathbf{y} - \mathbf{x}\|}{\epsilon} \right) (M - m)$.

1068 Repeating this argument by swapping the roles of \mathbf{x} and \mathbf{y} , we get $f(\mathbf{x}) - f(\mathbf{y}) \leq \left(\frac{\|\mathbf{y} - \mathbf{x}\|}{\epsilon} \right) (M - m)$.

1069 Therefore (3.2) holds. \square

³This section was written by Joseph Paat.

1070 **Proposition 3.20.** Let $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ be a convex function. Consider any compact, convex subset
 1071 $S \subseteq \text{dom}(f)$ and let $\mathbf{x}^* \in \text{relint}(S)$. Then there is an $\epsilon_{\mathbf{x}^*} > 0$ and values $m_{\mathbf{x}^*}, M_{\mathbf{x}^*} \in \mathbb{R}$ so that

$$m_{\mathbf{x}^*} \leq f(\mathbf{x}) \leq M_{\mathbf{x}^*} \quad (3.2)$$

1072 for all $\mathbf{x} \in B(\mathbf{x}^*, 2\epsilon_{\mathbf{x}^*}) \cap S$.

1073 *Proof.* Let $\mathbf{v}^1, \dots, \mathbf{v}^\ell$ be vectors that span the linear space parallel to $\text{aff}(S)$ (see Theorem 2.16). By definition
 1074 of relative interior, since $\mathbf{x}^* \in \text{aff}(S)$, there exists $\epsilon > 0$ such that $\mathbf{x}^* + \epsilon\mathbf{v}^j$ and $\mathbf{x}^* - \epsilon\mathbf{v}^j$ are both in S for
 1075 $j = 1, \dots, \ell$. Denote the set of points $\mathbf{x}^* \pm \epsilon\mathbf{v}^j$ as $\mathbf{x}_1, \dots, \mathbf{x}_k \in S$ ($k = 2\ell$), and define $S' := \text{conv}\{\mathbf{x}_1, \dots, \mathbf{x}_k\}$.
 1076 Observe that $\mathbf{x}^* \in \text{relint}(S')$ and $\text{aff}(S') = \text{aff}(S)$. Set $M_{\mathbf{x}^*} = \max\{f(\mathbf{x}_i) : i = 1, \dots, k\}$. Using Problem 3
 1077 from “HW for Week VIII”, it follows that $f(\mathbf{x}) \leq M_{\mathbf{x}^*}$ for all $\mathbf{x} \in S'$.

Now since f is convex, by Theorem 3.15 (see Remark 3.16 part 3.), there is some affine support function
 $L(\mathbf{x}) = \langle \mathbf{a}, (\mathbf{x} - \mathbf{x}^*) \rangle + f(\mathbf{x}^*)$ for f at \mathbf{x}^* . Define $m_{\mathbf{x}^*} = \min\{L(\mathbf{x}_i) : i = 1, \dots, k\}$. Consider any point
 $\mathbf{x} = \sum_{i=1}^k \lambda_i \mathbf{x}_i \in S'$, where $\lambda_1, \dots, \lambda_k$ are convex coefficients, and observe that

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

1078 Since L is an affine support, it follows that $f(\mathbf{x}) \geq L(\mathbf{x}) \geq m_{\mathbf{x}^*}$ for all $\mathbf{x} \in S'$. Finally, as $\mathbf{x}^* \in \text{relint}(S')$
 1079 and $\text{aff}(S') = \text{aff}(S)$, there is some $\epsilon > 0$ so that $B(\mathbf{x}^*, 2\epsilon) \cap S \subseteq S'$.

1080 □

1081 **Theorem 3.21.** Let $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ be a convex function. Let $D \subseteq \text{relint}(\text{dom}(f))$ be a convex,
 1082 compact subset. Then there is a constant $L = L(D) \geq 0$ so that

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

1083 for all $\mathbf{x}, \mathbf{y} \in D$. In particular, f is locally Lipschitz continuous over the relative interior of its domain.

1084 *Proof.* Let S be a compact set such that $D \subseteq \text{relint}(S) \subseteq \text{relint}(\text{dom}(f))$. From Proposition 3.20, for
 1085 every $\mathbf{x} \in \text{relint}(S)$, 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 $\mathbf{y} \in B(\mathbf{x}, 2\epsilon_{\mathbf{x}}) \cap S$.
 1086 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 $\mathbf{z}, \mathbf{y} \in$
 1087 $B(\mathbf{x}, \epsilon_{\mathbf{x}})$. 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
 1088 of $\text{aff}(S)$). Therefore, as S is compact, there exists a finite set $\{x_1, \dots, x_k\} \subset S$ so that $S \subseteq \bigcup_{i=1}^k B(\mathbf{x}_i, \epsilon_{\mathbf{x}_i})$.
 1089 Set $L = \max\{L_{\mathbf{x}_i} : i \in \{1, \dots, k\}\}$.

1090 Now take $\mathbf{y}, \mathbf{z} \in S$. The line segment $[\mathbf{y}, \mathbf{z}]$ can be divided into finitely many segments $[\mathbf{y}, \mathbf{z}] = [\mathbf{y}_1, \mathbf{y}_2] \cup$
 1091 $[\mathbf{y}_2, \mathbf{y}_3] \cup \dots \cup [\mathbf{y}_{q-1}, \mathbf{y}_q]$, 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 $B(\mathbf{x}_j, \epsilon_{\mathbf{x}_j})$
 1092 for $j \in \{1, \dots, k\}$. 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

1093 each $i \in \{1, \dots, q-1\}$. It follows that

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

1094 Hence f is Lipschitz continuous over S with constant L . □

1095 3.3 First-order derivative properties

1096 A convex function enjoys very strong differentiability properties. We will first state some useful results
1097 without proof. See the Chapter on “Convex Functions” in Gruber [3] for full proofs.

1098 **Theorem 3.22.** Let $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ be a convex function and let $\mathbf{x} \in \text{int}(\text{dom}(f))$. Then f is
1099 differentiable at \mathbf{x} if and only if the partial derivative $f'_i(\mathbf{x})$ exists for all $i = 1, \dots, d$.

1100 **Theorem 3.23.** [Reidemeister’s Theorem] Let $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ be a convex function. Then f is
1101 differentiable almost everywhere in $\text{int}(\text{dom}(f))$, i.e., the subset of $\text{int}(\text{dom}(f))$ where f is not differentiable
1102 has Lebesgue measure 0.

1103 We now prove the central relationships between the gradient ∇f and convexity. We first observe some
1104 facts about convex functions on the real line.

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

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

1105 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:

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

Similarly,

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

1106 The strict convexity implication is clear from the above. □

1107 An immediate corollary is the following relationship between the derivative of a function on the real line
 1108 and convexity.

1109 **Proposition 3.25.** Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be a differentiable function. Then f is convex if and only if f' is an
 1110 increasing function, i.e., $f'(x) \geq f'(y)$ for all $x \geq y \in \mathbb{R}$. Moreover, f is strictly convex if and only if
 1111 f' is strictly increasing. f is strongly convex with strong convexity modulus $c > 0$ if and only if $f'(x) \geq$
 1112 $f'(y) + c(x - y)$ for all $x \geq y \in \mathbb{R}$.

1113 *Proof.* (\Rightarrow) Recall that $f'(x) = \lim_{t \rightarrow 0^+} \frac{f(x+t) - f(x)}{t}$. But for every $0 < t < y - x$, we have $\frac{f(x+t) - f(x)}{t} \leq$
 1114 $\frac{f(y) - f(x)}{y - x}$ by Proposition 3.24. Thus, $f'(x) \leq \frac{f(y) - f(x)}{y - x}$. By a similar argument, we obtain $f'(y) \geq \frac{f(y) - f(x)}{y - x}$.
 1115 This gives the relation.

(\Leftarrow) Consider any $x, z \in \mathbb{R}$ 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 $t_2 \geq t_1$ and
 we assume f' is increasing, then $f'(t_2) \geq f'(t_1)$. This implies that

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

1116 Substituting $y = \alpha x + (1 - \alpha)z$ and rearranging, we obtain that $f(\alpha x + (1 - \alpha)z) \leq \alpha f(x) + (1 - \alpha)f(z)$.
 1117 □

1118 We can now prove the main result of this subsection. A key idea behind the results below is that one can
 1119 reduce testing convexity of a function on \mathbb{R}^d to testing convexity of any one-dimensional “slice” of it. More
 1120 precisely,

1121 **Proposition 3.26.** Let $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ be a function. Then f is convex if and only if for every
 1122 $\mathbf{x}, \mathbf{r} \in \mathbb{R}^d$, the function $\phi : \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$ defined by $\phi(t) = f(\mathbf{x} + t\mathbf{r})$ is convex.

1123 *Proof.* Left as an exercise. □

1124 **Theorem 3.27.** Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be differentiable everywhere. Then the following are all equivalent.

- 1125 1. f is convex.
- 1126 2. $f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$ for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$.
- 1127 3. $\langle \nabla f(\mathbf{y}) - \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \geq 0$ for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$.

1128 A characterization of strict convexity is obtained if all the above inequalities are considered strict for all
 1129 $\mathbf{x} \neq \mathbf{y} \in \mathbb{R}^d$. A characterization of strong convexity with modulus $c > 0$ is obtained if 2. is replaced with
 1130 $f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{1}{2}c\|\mathbf{y} - \mathbf{x}\|^2$ for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, and 3. is replaced with $\langle \nabla f(\mathbf{y}) - \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \geq$
 1131 $c\|\mathbf{y} - \mathbf{x}\|^2$ 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

$$\begin{aligned} & \frac{f((1-\alpha)\mathbf{x} + \alpha\mathbf{y}) - f(\mathbf{x})}{\alpha} \leq f(\mathbf{y}) - f(\mathbf{x}) \\ \Rightarrow & \frac{f(\mathbf{x} + \alpha(\mathbf{y} - \mathbf{x})) - f(\mathbf{x})}{\alpha} \leq f(\mathbf{y}) - f(\mathbf{x}) \end{aligned}$$

1132 Letting $\alpha \rightarrow 0$ 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

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

1133 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

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

1134 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 $t_2 > t_1$.
 1135 Therefore, by Proposition 3.25, we obtain that $\phi(t)$ is a convex function in t . By Proposition 3.26, f is
 1136 convex. \square

1137 3.4 Second-order derivative properties

1138 A simple consequence of Proposition 3.25 for twice differentiable functions on the real line is the following.

1139 **Corollary 3.28.** Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be a twice differentiable function. Then f is convex if and only if $f''(x) \geq 0$
 1140 for all $x \in \mathbb{R}$. If $f''(x) > 0$, then f is strictly convex.

1141 **Remark 3.29.** From Proposition 3.25, we know strict convexity of f is equivalent to the condition that f' is
 1142 strictly increasing. However, this is not equivalent to $f''(x) > 0$, the implication only goes in one direction.
 1143 This is why we lose the other direction when discussing strict convexity in Corollary 3.28. As a concrete
 1144 example, consider $f(x) = x^4$ which is strictly convex, but the second derivative is 0 at $x = 0$.

1145 This enables one to characterize convexity of $f : \mathbb{R}^d \rightarrow \mathbb{R}$ in terms of its Hessian, which will be denoted
 1146 by $\nabla^2 f$.

1147 **Theorem 3.30.** Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a twice differentiable function. Then the following are all true.

- 1148 1. f is convex if and only if $\nabla^2 f(\mathbf{x})$ is positive semidefinite (PSD) for all $\mathbf{x} \in \mathbb{R}^d$.
- 1149 2. If $\nabla^2 f(\mathbf{x})$ is positive definite (PD) for all $\mathbf{x} \in \mathbb{R}^d$, then f is strictly convex.
- 1150 3. f is strongly convex with modulus $c > 0$ if and only if $\nabla^2 f(\mathbf{x}) - cI$ is positive semidefinite (PSD) for
 1151 all $\mathbf{x} \in \mathbb{R}^d$.

1152 *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
 1153 $\mathbf{r} \in \mathbb{R}^d$. Define the function $\phi(t) = f(\mathbf{x} + t\mathbf{r})$. By Proposition 3.26, ϕ is convex. By Corollary 3.28,
 1154 $0 \leq \phi''(0) = \langle \nabla^2 f(\mathbf{x})\mathbf{r}, \mathbf{r} \rangle$. Since the choice of \mathbf{r} was arbitrary, this shows that $\nabla^2 f(\mathbf{x})$ is positive
 1155 semidefinite.

1156 (\Leftarrow) Assume $\nabla^2 f(\mathbf{x})$ is positive semidefinite for all $\mathbf{x} \in \mathbb{R}^d$, and consider $\bar{\mathbf{x}}, \mathbf{r} \in \mathbb{R}^d$. Define the function
 1157 $\phi(t) = f(\bar{\mathbf{x}} + t\mathbf{r})$. 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
 1158 Corollary 3.28, ϕ is convex. By Proposition 3.26, f is convex.

1159 2. This follows from the same construction as in 1. above, and the sufficient condition that if the second
 1160 derivative of a one-dimensional function is strictly positive, then the function is strictly convex.

1161 3. We omit the proof of the characterization of strong convexity.
 1162 \square

1163 3.5 Sublinear functions, support functions and gauges

1164 We will now introduce a more structured subfamily of convex functions which is easier to deal with analyti-
1165 cally, and yet has very important uses in diverse areas.

1166 **Definition 3.31.** A function $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ is called *sublinear* if it satisfies the following two properties:

1167 (i) f is *positively homogeneous*, i.e., $f(\lambda \mathbf{r}) = \lambda f(\mathbf{r})$ for all $\mathbf{r} \in \mathbb{R}^d$ and $\lambda > 0$.

1168 (ii) f is *subadditive*, i.e., $f(\mathbf{x} + \mathbf{y}) \leq f(\mathbf{x}) + f(\mathbf{y})$ for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$.

1169 Here is the connection with convexity.

1170 **Proposition 3.32.** Let $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$. Then the following are equivalent:

1171 1. f is sublinear.

1172 2. f is convex and positively homogeneous.

1173 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 $\lambda_1, \lambda_2 > 0$.

1174 *Proof.* Left as an exercise. □

1175 A useful property of sublinear functions is the following.

1176 **Proposition 3.33.** Let $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ be a sublinear function. Then either $f(\mathbf{0}) = 0$ or $f(\mathbf{0}) = +\infty$.

1177 A characterization of sublinear functions via epigraphs is also possible.

1178 **Proposition 3.34.** Let $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ such that $f(\mathbf{0}) = 0$. Then f is sublinear if and only if $\text{epi}(f)$
1179 is a convex cone in $\mathbb{R}^d \times \mathbb{R}$.

1180 *Proof.* (\Rightarrow) From Proposition 3.32, we know that f is convex and positively homogeneous. From Propo-
1181 sition 3.6, this implies that $\text{epi}(f)$ is convex. So we only need to verify that if $(\mathbf{x}, t) \in \text{epi}(f)$ then
1182 $\lambda(\mathbf{x}, t) = (\lambda \mathbf{x}, \lambda t) \in \text{epi}(f)$ for all $\lambda \geq 0$. If $\lambda = 0$, then the result follows from the assumption that
1183 $f(\mathbf{0}) = 0$. Now consider $\lambda > 0$. Since $(\mathbf{x}, t) \in \text{epi}(f)$, we have $f(\mathbf{x}) \leq t$ and by positive homogeneity of f ,
1184 $f(\lambda \mathbf{x}) = \lambda f(\mathbf{x}) \leq \lambda t$, and so $(\lambda \mathbf{x}, \lambda t) \in \text{epi}(f)$.

1185 (\Leftarrow) From Proposition 3.6 and the assumption that $\text{epi}(f)$ is a convex cone, we get that f is convex. We
1186 now verify that f is positively homogeneous; by Proposition 3.32, we will be done. We first verify that for
1187 all $\lambda > 0$ and $\mathbf{x} \in \mathbb{R}^d$, $f(\lambda \mathbf{x}) \leq \lambda f(\mathbf{x})$. Since $\text{epi}(f)$ is a convex cone and $(\mathbf{x}, f(\mathbf{x})) \in \text{epi}(f)$, we have that
1188 $\lambda(\mathbf{x}, f(\mathbf{x})) = (\lambda \mathbf{x}, \lambda f(\mathbf{x})) \in \text{epi}(f)$. This implies that $f(\lambda \mathbf{x}) \leq \lambda f(\mathbf{x})$.

1189 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
1190 with $\lambda = \frac{1}{\bar{\lambda}}$ and $\mathbf{x} = \bar{\lambda} \bar{\mathbf{x}}$, 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
1191 $f(\bar{\lambda} \bar{\mathbf{x}}) = \bar{\lambda} f(\bar{\mathbf{x}})$. □

1192 **Gauges.** One easily observes that any norm $N : \mathbb{R}^d \rightarrow \mathbb{R}$ is a sublinear function – recall Definition 1.1.
1193 In fact, a norm has the additional “symmetry” property that $N(\mathbf{x}) = N(-\mathbf{x})$. Since a sublinear function is
1194 convex (Proposition 3.32), and sublevel sets of convex sets are convex, we immediately know that the unit
1195 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,
1196 these unit norm balls are also “symmetric” about the origin. This merits a definition.

1197 **Definition 3.35.** A convex set $C \subseteq \mathbb{R}^d$ is said to be *centrally symmetric about the origin*, if $\mathbf{x} \in C$ implies
1198 that $-\mathbf{x} \in C$. Sometimes we will abbreviate this to say C is centrally symmetric.

1199 We now summarize the above discussion in the following observation.

1200 **Proposition 3.36.** Let $N : \mathbb{R}^d \rightarrow \mathbb{R}$ be a norm. Then the unit norm ball $B_N(\mathbf{0}, 1) = \{\mathbf{x} \in \mathbb{R}^d : N(\mathbf{x}) \leq 1\}$
1201 is a centrally symmetric, closed convex set.

1202 One can actually prove a converse to the above statement, which will establish a nice one-to-one corre-
 1203 spondence between norms and centrally symmetric convex sets. We first generalize the notion of a norm to
 1204 a family of sublinear functions called “gauge functions”.

Definition 3.37. Let $C \subseteq \mathbb{R}^d$ be a closed, convex set such that $\mathbf{0} \in C$. Define the following function
 $\gamma_C : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ as

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

1205 γ_C is called the *gauge* or the *Minkowski functional* of C .

1206 **Exercise 7.** Show that γ_C is finite valued everywhere if and only if $\mathbf{0} \in \text{int}(C)$.

1207 The following is a useful observation for the analysis of gauge functions.

1208 **Lemma 3.38.** 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
 1209 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 $[a, +\infty)$.

1210 *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
 1211 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
 1212 $\frac{1}{\lambda}\mathbf{r} = \frac{\bar{\lambda}}{\lambda}(\frac{1}{\bar{\lambda}}\mathbf{r}) + (\frac{\lambda-\bar{\lambda}}{\lambda})\mathbf{0}$ which is in C because C is convex and $\mathbf{0} \in C$. \square

1213 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
 1214 scale \mathbf{r} with so that you end up on the boundary of C . More precisely,

1215 **Proposition 3.39.** 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
 1216 $0 < \gamma_C(\mathbf{r}) < \infty$. Then $\frac{1}{\gamma_C(\mathbf{r})}\mathbf{r} \in \text{relbd}(C)$.

1217 *Proof.* From Lemma 3.38, 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
 1218 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
 1219 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
 1220 the fact that $\gamma_C(\mathbf{r}) = \inf\{\lambda > 0 : \mathbf{r} \in \lambda C\}$, since $\frac{\gamma_C(\mathbf{r})}{\alpha} < \gamma_C(\mathbf{r})$. \square

1221 The following theorem relates geometric properties of C with analytical properties of the gauge function.
 1222 These relations are extremely handy to keep in mind.

1223 **Theorem 3.40.** Let $C \subseteq \mathbb{R}^d$ be a closed, convex set such that $\mathbf{0} \in C$. Then the following are all true.

- 1224 1. γ_C is a nonnegative, sublinear function.
- 1225 2. $C = \{\mathbf{x} \in \mathbb{R}^d : \gamma_C(\mathbf{x}) \leq 1\}$.
- 1226 3. $\text{rec}(C) = \{\mathbf{r} \in \mathbb{R}^d : \gamma_C(\mathbf{r}) = 0\}$.
- 1227 4. If $\mathbf{0} \in \text{relint}(C)$, then $\text{relint}(C) = \{\mathbf{x} \in \mathbb{R}^d : \gamma_C(\mathbf{x}) < 1\}$.

1228 *Proof.* 1. Although 1. can be proved directly from the definition of the gauge, we postpone its proof until
 1229 we speak of *support functions* below.

1230 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 $1 \in \{\lambda > 0 : \mathbf{x} \in \lambda C\}$
 1231 and therefore, $\inf\{\lambda > 0 : \mathbf{x} \in \lambda C\} \leq 1$.

1232 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
 1233 since $\{\lambda > 0 : \mathbf{x} \in \lambda C\}$ is an unbounded interval by Lemma 3.38, this means that either $1 \in \{\lambda > 0 : \mathbf{x} \in \lambda C\}$,
 1234 and thus $\mathbf{x} \in C$ or $1 = \inf\{\lambda > 0 : \mathbf{x} \in \lambda C\} = \gamma_C(\mathbf{x})$. By Proposition 3.39, we have that
 1235 $1 \cdot \mathbf{x} \in C$.

1236 3. Since $\{\lambda > 0 : \mathbf{r} \in \lambda C\}$ is convex by Lemma 3.38, we observe that $\gamma_C(\mathbf{r}) = 0$ if and only if $\frac{1}{\lambda}\mathbf{r} \in C$ for
 1237 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, $\mathbf{0} + t\mathbf{r} \in C$
 1238 for all $t \geq 0$. This is equivalent to saying that \mathbf{r} satisfies Definition 2.43 of $\text{rec}(C)$.

1239 4. Consider any $\mathbf{x} \in \text{relint}(C)$. By definition of relative interior, there exists $\lambda > 1$ such that $\lambda\mathbf{x} \in C$.
 1240 By part 2. above, $\gamma_C(\lambda\mathbf{x}) \leq 1$ and by part 1. above, γ_C is positively homogeneous, and thus,
 1241 $\gamma_C(\mathbf{x}) \leq \frac{1}{\lambda} < 1$.

1242 Now suppose $\mathbf{x} \in \mathbb{R}^d$ such that $\gamma_C(\mathbf{x}) < 1$. If $\gamma_C(\mathbf{x}) = 0$, then $\mathbf{x} \in \text{rec}(C)$ by part 3. above. Since
 1243 $\mathbf{0} \in \text{relint}(C)$, we also have $\mathbf{x} = \mathbf{0} + \mathbf{x} \in \text{relint}(C)$. Now suppose $0 < \gamma_C(\mathbf{x}) < 1$. By part 2. above,
 1244 $\mathbf{x} \in C$. Suppose to the contrary that $\mathbf{x} \notin \text{relint}(C)$. By Theorem 2.40, \mathbf{x} is contained in a proper face
 1245 F of C . Since $\mathbf{0} \in \text{relint}(C)$, $\mathbf{0}$ is not contained in F . Also, $\gamma_C(\frac{\mathbf{x}}{\gamma_C(\mathbf{x})}) = 1$ by positive homogeneity
 1246 of γ_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
 1247 $\gamma_C(\mathbf{x}) < 1$ and $\mathbf{0} \notin F$, this would contradict the fact that F is a face.
 1248 □

1249 We derive some immediate consequences.

1250 **Corollary 3.41.** Let $C \subseteq \mathbb{R}^d$ be a closed, convex set containing the origin. Then C is compact if and only
 1251 if $\gamma(\mathbf{r}) > 0$ for all $\mathbf{r} \in \mathbb{R}^d \setminus \{\mathbf{0}\}$.

1252 **Corollary 3.42.** [Uniqueness of the gauge] Let C be a compact convex set containing the origin in its
 1253 interior, i.e., $\mathbf{0} \in \text{int}(C)$. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be any sublinear function. Then $C = \{\mathbf{x} \in \mathbb{R}^d : f(\mathbf{x}) \leq 1\}$ if and
 1254 only if $f = \gamma_C$.

1255 *Proof.* The sufficiency follows from Theorem 3.40, part 2. For the necessity, suppose to the contrary that
 1256 $f(\mathbf{x}) \neq \gamma_C(\mathbf{x})$ 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 Proposition 3.33.

1257 First suppose $f(\mathbf{x}) > \gamma_C(\mathbf{x})$. Since C is compact, we know that $\gamma_C(\mathbf{x}) > 0$ by Corollary 3.41. Con-
 1258 sider that point $\frac{1}{\gamma_C(\mathbf{x})}\mathbf{x}$. By Proposition 3.39, $\mathbf{x} \in \text{relbd}(C)$. However, since f is positively homogeneous,
 1259 $f(\frac{1}{\gamma_C(\mathbf{x})}\mathbf{x}) = \frac{1}{\gamma_C(\mathbf{x})}f(\mathbf{x}) > 1$ because $f(\mathbf{x}) > \gamma_C(\mathbf{x})$. This contradicts that $C = \{\mathbf{x} \in \mathbb{R}^d : f(\mathbf{x}) \leq 1\}$.

1260 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,
 1261 $\lambda\mathbf{x} \in C$ for all $\lambda \geq 0$ by the assumption that $C = \{\mathbf{x} \in \mathbb{R}^d : f(\mathbf{x}) \leq 1\}$. This means that $\mathbf{x} \in \text{rec}(C)$ which
 1262 contradicts the fact that C is compact (see Theorem 2.47). Thus, we may assume that $f(\mathbf{x}) > 0$.

1263 Now let $\mathbf{y} = \frac{1}{f(\mathbf{x})}\mathbf{x}$. By positive homogeneity of γ_C , we obtain that $\gamma_C(\mathbf{y}) = \gamma_C(\frac{1}{f(\mathbf{x})}\mathbf{x}) = \frac{\gamma_C(\mathbf{x})}{f(\mathbf{x})} > 1$.
 1264 Therefore, $\mathbf{y} \notin C$ by Theorem 3.40, part 2. However, $f(\mathbf{y}) = 1$, which contradicts the assumption that
 1265 $C = \{\mathbf{x} \in \mathbb{R}^d : f(\mathbf{x}) \leq 1\}$. □

1266 The proof of Corollary 3.42 can be massaged to obtain the following.

1267 **Corollary 3.43.** [Uniqueness of the gauge-II] Let C be a closed, convex set (not necessarily compact)
 1268 containing the origin in its interior, i.e., $\mathbf{0} \in \text{int}(C)$. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be any *nonnegative*, sublinear function.
 1269 Then $C = \{\mathbf{x} \in \mathbb{R}^d : f(\mathbf{x}) \leq 1\}$ if and only if $f = \gamma_C$.

1270 Consequently, for every nonnegative, sublinear function f , there exists a closed, convex set C such that
 1271 $f = \gamma_C$.

1272 We also make the following observation on when the gauge function can take $+\infty$ as a value.

1273 **Lemma 3.44.** Let C be a closed, convex set with $\mathbf{0} \in C$. Then the gauge γ_C is finite valued everywhere
 1274 (i.e., $\gamma_C(\mathbf{x}) < \infty$ for all $\mathbf{x} \in \mathbb{R}^d$) if and only if $\mathbf{0} \in \text{int}(C)$.

1275 *Proof.* (\implies) Suppose $\mathbf{0}$ is not in the interior, i.e., $\mathbf{0}$ is on the boundary of C . By the Supporting Hyperplane
 1276 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, $\delta = 0$.
 1277 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 $\lambda C \subseteq H^-(\mathbf{a}, 0)$
 1278 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
 1279 shows that γ_C takes value ∞ on the entire “open” halfspace $\{\mathbf{r} \in \mathbb{R}^d : \langle \mathbf{a}, \mathbf{r} \rangle > 0\}$.

1280 (\impliedby) Assume $\mathbf{0} \in \text{int}(C)$ and consider any $\mathbf{x} \in \mathbb{R}^d$. Since $\mathbf{0} \in \text{int}(C)$, there exists $\epsilon > 0$ such that $\epsilon\mathbf{x} \in C$.
 1281 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, $\gamma_C(\mathbf{x}) < \infty$
 1282 for all $\mathbf{x} \in \mathbb{R}^d$. □

1283 We can now finally settle the correspondence between norms and centrally symmetric, compact convex
 1284 sets.

1285 **Theorem 3.45.** Let $N : \mathbb{R}^d \rightarrow \mathbb{R}$ be a norm. Then $B_N(\mathbf{0}, 1) = \{\mathbf{x} \in \mathbb{R}^d : N(\mathbf{x}) \leq 1\}$ is a centrally
 1286 symmetric, compact convex set with $\mathbf{0}$ in its interior. Moreover, $\gamma_{B_N(\mathbf{0}, 1)} = N$.

1287 Conversely, let B be a centrally symmetric, compact convex set containing $\mathbf{0}$ in its interior. Then γ_B is
 1288 a norm on \mathbb{R}^d and $B = B_{\gamma_B}(\mathbf{0}, 1)$.

1289 *Proof.* For the first part, since N is sublinear, it is convex (by Proposition 3.32). By definition, $B_N(\mathbf{0}, 1) =$
 1290 $\{\mathbf{x} \in \mathbb{R}^d : N(\mathbf{x}) \leq 1\}$ is a sublevel set for N , and is thus a convex set. It is closed, since N is continuous by
 1291 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
 1292 that $\text{rec}(B_N(\mathbf{0}, 1)) = \{\mathbf{0}\}$; this will imply that it is compact by Theorem 2.47. Consider any nonzero vector
 1293 \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
 1294 a recession direction for $B_N(\mathbf{0}, 1)$.

1295 We verify that $\mathbf{0} \in \text{int}(B_N(\mathbf{0}, 1))$. If not, then by the Supporting Hyperplane Theorem 2.23, there exists
 1296 $\mathbf{a} \in \mathbb{R}^d \setminus \{\mathbf{0}\}$ 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 $\mathbf{a} \neq \mathbf{0}$,
 1297 $N(\mathbf{a}) > 0$. 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
 1298 contradicts the fact that $B_N(\mathbf{0}, 1) \subseteq H^-(\mathbf{a}, 0)$. Therefore, from Corollary 3.42, we obtain that $N = \gamma_{B_N(\mathbf{0}, 1)}$.

1299 For the second part, we know that γ_B is sublinear, and since B is compact, $\gamma_B(\mathbf{r}) > 0$ for all $\mathbf{r} \neq \mathbf{0}$ by
 1300 Corollary 3.41. Since $\mathbf{0} \in \text{int}(B)$, Lemma 3.44 implies that γ_C is finite valued everywhere. To confirm that
 1301 γ_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
 1302 that $\gamma_B(\mathbf{x}) > \gamma_B(-\mathbf{x})$ (note that this is without loss of generality). This implies that $\gamma_B(\frac{1}{\gamma_B(-\mathbf{x})}\mathbf{x}) > 1$.
 1303 Therefore, $\frac{1}{\gamma_B(-\mathbf{x})}\mathbf{x} \notin B$ by Theorem 3.40, 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
 1304 that $-\frac{1}{\gamma_B(-\mathbf{x})}\mathbf{x} \in B$ by Theorem 3.40, part 2. This contradicts the fact that B is centrally symmetric. Thus,
 1305 γ_B is a norm on \mathbb{R}^d . Moreover, by Theorem 3.40, part 2., $B = \{\mathbf{x} \in \mathbb{R}^d : \gamma_B(\mathbf{x}) \leq 1\} = B_{\gamma_B}(\mathbf{0}, 1)$. \square

1306 Let us build towards a more computational approach to the gauge. First, let's give an explicit formula
 1307 for the gauge of a halfspace containing the origin.

Example 3.46. 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 δ 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

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

1308 The above calculation, along with the next theorem, gives powerful computational tools for gauge func-
 1309 tions.

Theorem 3.47. Let $C_i, i \in I$ be a (not necessarily finite) family of closed, convex sets containing the origin,
 and let $C = \bigcap_{i \in I} C_i$. Then

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

1310 *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 $A = \{\lambda >$
 1311 $0 : \mathbf{r} \in \lambda C\}$. Observe that $A = \bigcap A_i$. If any A_i is empty, then $\gamma_{C_i}(\mathbf{r}) = \infty$, and A is empty and therefore
 1312 $\gamma_C(\mathbf{r}) = \infty$, and the equality holds. Now suppose all A_i 's are nonempty, and so by Lemma 3.38, each A_i is of
 1313 the form (a_i, ∞) or $[a_i, \infty)$. If $A = \emptyset$, then it must mean that $a_i \rightarrow \infty$. Since $\gamma_{C_i}(\mathbf{r}) = \inf A_i = a_i$, this shows
 1314 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
 1315 A is nonempty. Then since $A = \bigcap A_i$, A must be of the form $(a, +\infty)$ or $[a, +\infty)$ where $a := \sup_{i \in I} a_i$. Then
 1316 $\gamma_C(\mathbf{r}) = a = \sup_{i \in I} a_i = \sup_{i \in I} \gamma_{C_i}(\mathbf{r})$. \square

1317 This shows that gauge functions for polyhedra can be computed very easily.

Corollary 3.48. 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

$$P = \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{a}^i, \mathbf{x} \rangle \leq 1 \quad i = 1, \dots, m\}.$$

Then

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

1318 *Proof.* Use the formula from 3.46 and Theorem 3.47. □

1319 **Support functions.** While gauges are good in the sense that they are a nice generalization of norms from
 1320 centrally symmetric convex bodies to asymmetric convex bodies, there is a drawback. Gauges are a strict
 1321 subset of sublinear functions because they are always nonnegative, while there are many sublinear functions
 1322 that take negative values. We would like to establish a one-to-one correspondence between sublinear functions
 1323 and all closed, convex sets. Note that the correspondence via the epigraph only establishes a correspondence
 1324 with closed, convex cones, and that too not all closed, convex cones are covered. The right definition, it
 1325 turns out, is inspired by optimization of linear functions over closed, convex sets.

Definition 3.49. Let $S \subseteq \mathbb{R}^d$ be any set. The *support function* for S is a function on \mathbb{R}^d defined as

$$\sigma_S(\mathbf{r}) = \sup_{\mathbf{x} \in S} \langle \mathbf{r}, \mathbf{x} \rangle.$$

1326 The following is easy to verify, and aspects of it were already explored in the midterm and HWs.

Proposition 3.50. Let $S \subseteq \mathbb{R}^d$. Then

$$\sigma_S = \sigma_{\text{cl}(S)} = \sigma_{\text{conv}(S)} = \sigma_{\text{cl}(\text{conv}(S))}.$$

1327 **Proposition 3.51.** Let $S \subseteq \mathbb{R}^d$. Then σ_S is a closed, sublinear function, i.e., its epigraph is a closed, convex
 1328 cone.

Proof. We first check that σ_S is sublinear. We check positive homogeneity. For any $\mathbf{r} \in \mathbb{R}^d$ and $\lambda > 0$,

$$\sigma_S(\lambda \mathbf{r}) = \sup_{\mathbf{x} \in S} \langle \lambda \mathbf{r}, \mathbf{x} \rangle = \sup_{\mathbf{x} \in S} \lambda \langle \mathbf{r}, \mathbf{x} \rangle = \lambda \sup_{\mathbf{x} \in S} \langle \mathbf{r}, \mathbf{x} \rangle = \lambda \sigma_S(\mathbf{r}).$$

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

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

1329 Since σ_S is the supremum of linear functions $\langle \mathbf{x}, \mathbf{r} \rangle$, $\mathbf{x} \in S$, $\text{epi}(f)$ is the intersection of closed halfspaces,
 1330 which shows that it is closed. The fact that it is a convex cone follows from Proposition 3.34. □

1331 We now establish a fundamental correspondence between gauges and support functions via polarity.

Theorem 3.52. 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

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

By Theorem 3.47 and Example 3.46, we obtain that

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

1332 Since $\mathbf{0} \in C^\circ$, the last term above can be written as $\sup_{\mathbf{a} \in C^\circ} \langle \mathbf{a}, \mathbf{r} \rangle = \sigma_{C^\circ}(\mathbf{r})$. □

Example 3.53. Consider the polyhedron

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

From Corollary 3.48, we obtain that

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

and by Theorem 3.40 part 2., we obtain that $P = \{\mathbf{x} \in \mathbb{R}^2 : \gamma_P(\mathbf{x}) \leq 1\}$. Now consider the function

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

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

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

1333 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
 1334 set $S = \{(-1, -1), (\frac{1}{2}, -1), (-1, \frac{1}{2})\}$. This shows that Corollary 3.42 really breaks down if the assumption
 1335 of compactness is removed. Even so, given a closed, convex set C , any sublinear function that has a set C
 1336 as its 1-sublevel set must match the gauge on $\mathbb{R}^d \setminus \text{int}(\text{rec}(C))$ (see Problem 7 from “HW for Week IX”). If
 1337 you are interested in learning more about representing closed, convex sets as the sublevel sets of sublinear
 1338 functions, please see [1] on exciting new results.

1339

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

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

Proof. Consider any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$. Since C is compact, $\gamma_C(\mathbf{x}) > 0$ by Corollary 3.41, and $\sigma_C(\mathbf{y}) < +\infty$. By Proposition 3.39, $\frac{\mathbf{x}}{\gamma_C(\mathbf{x})} \in C$, and therefore,

1340

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

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

Corollary 3.55. 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 3.54 and 3.52. □

The above corollary generalizes Holder’s inequality by recalling that when $\frac{1}{p} + \frac{1}{q} = 1$, then the ℓ^p and ℓ^q unit balls are polars of each other. Note that Theorem 3.54 and Corollary 3.55 have no assumption of centrally symmetric sets, so they strictly generalize the norm inequalities of Holder and Cauchy-Schwarz.

1341 **One-to-one correspondence between closed, convex sets and closed, sublinear functions.** Propo-
 1342 sition 3.51 shows that support functions are closed, sublinear functions. Proposition 3.50 shows that two
 1343 different sets, e.g., S and $\text{conv}(S)$, may give rise to the same sublinear function $\sigma_S = \sigma_{\text{conv}(S)}$ via the support
 1344 function construction. In other words, if we consider the mapping $S \rightarrow \sigma_S$ as a mapping from the family of
 1345 subsets of \mathbb{R}^d to the family of closed, sublinear functions, this mapping is not injective. But if we restrict to
 1346 closed, convex sets, it can be shown that this mapping is injective.

1347 **Exercise 8.** Let C_1, C_2 be closed, convex sets. Then $\sigma_{C_1} = \sigma_{C_2}$ if and only if $C_1 = C_2$.

1348 A natural question now is whether the mapping $C \rightarrow \sigma_C$ from the family of closed, convex sets to the
 1349 family of closed, sublinear functions is *onto*. The answer is yes! **Thus, all closed, sublinear functions**
 1350 **are support functions and vice versa.**

1351 **Theorem 3.56.** Let $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ be a sublinear function that is also closed. Then the set

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

1352 is a closed, convex set. Moreover, $\sigma_{C_f} = f$.

1353 Conversely, if C is a closed, convex set, then $C_{\sigma_C} = C$.

1354 *Proof.* We will prove the assertion when f is finite valued everywhere; the proof for general f is more tedious
 1355 and does not provide any additional insight, in our opinion, and will be skipped here.

1356 Since C_f is defined as the intersection of a family of halfspaces (indexed by \mathbb{R}^d), C_f is a closed, convex
 1357 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
 1358 $\sigma_{C_f}(\mathbf{r}) = \sup_{\mathbf{x} \in C_f} \langle \mathbf{r}, \mathbf{x} \rangle \leq f(\mathbf{r})$. To show that $\sigma_{C_f}(\mathbf{r}) \geq f(\mathbf{r})$, it suffices to exhibit $\mathbf{y} \in C_f$ such that
 1359 $\langle \mathbf{r}, \mathbf{y} \rangle = f(\mathbf{r})$. Consider $\text{epi}(f)$, which by Proposition 3.34, is a closed convex cone (since f is assumed to be
 1360 closed). By Theorem 2.23, there exists a supporting hyperplane for $\text{epi}(f)$ at $(\mathbf{r}, f(\mathbf{r}))$. Let this hyperplane
 1361 be defined by $(\mathbf{y}, \eta) \in \mathbb{R}^d \times \mathbb{R}$ and $\alpha \in \mathbb{R}$ such that $\text{epi}(f) \subseteq H^-((\mathbf{y}, \eta), \alpha)$. Using Problems 8 and 9
 1362 from “HW for Week IX”, one can assume that $\alpha = 0$ and $\eta < 0$. After normalizing, this means that
 1363 $\text{epi}(f) \subseteq H^-((\mathbf{y}/-\eta, -1), 0)$. This implies that for every $\mathbf{r}' \in \mathbb{R}^d$, $(\mathbf{r}', f(\mathbf{r}')) \in H^-((\mathbf{y}/-\eta, -1), 0)$, which
 1364 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
 1365 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.

1366 We now show that $C_{\sigma_C} = C$ for any closed, convex set C . Consider any $\mathbf{x} \in C$. Then $\langle \mathbf{r}, \mathbf{x} \rangle \leq$
 1367 $\sup_{\mathbf{y} \in C} \langle \mathbf{r}, \mathbf{y} \rangle = \sigma_C(\mathbf{r})$. Therefore, $\mathbf{x} \in H^-(\mathbf{r}, \sigma_C(\mathbf{r}))$ for all $\mathbf{r} \in \mathbb{R}^d$. This shows that $\mathbf{x} \in C_{\sigma_C}$, and therefore,
 1368 $C \subseteq C_{\sigma_C}$. To show the reverse inclusion, consider any $\mathbf{y} \notin C$. Since C is a closed, convex set, there
 1369 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
 1370 $\sigma_C(\mathbf{a}) = \sup_{\mathbf{x} \in C} \langle \mathbf{a}, \mathbf{x} \rangle \leq \delta$. Since C_{σ_C} has $\langle \mathbf{a}, \mathbf{x} \rangle \leq \sigma_C(\mathbf{a})$ as a defining halfspace, and $\langle \mathbf{a}, \mathbf{y} \rangle > \delta \geq \sigma_C(\mathbf{a})$,
 1371 we observe that $\mathbf{y} \notin C_{\sigma_C}$. \square

1372 One can associate a nice picture with the above construction of C_f associated with the sublinear function
 1373 f , which corresponds to the following proposition.

1374 **Proposition 3.57.** Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a sublinear function, and let C_f be defined as in Theorem 3.56.
 1375 Then $\mathbf{y} \in C_f$ if and only if $(\mathbf{y}, -1) \in \text{epi}(f)^\circ$. In other words, $C_f = \{\mathbf{y} \in \mathbb{R}^d : (\mathbf{y}, -1) \in \text{epi}(f)^\circ\}$.

Proof. We simply observe the following equivalences.

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

1376 \square

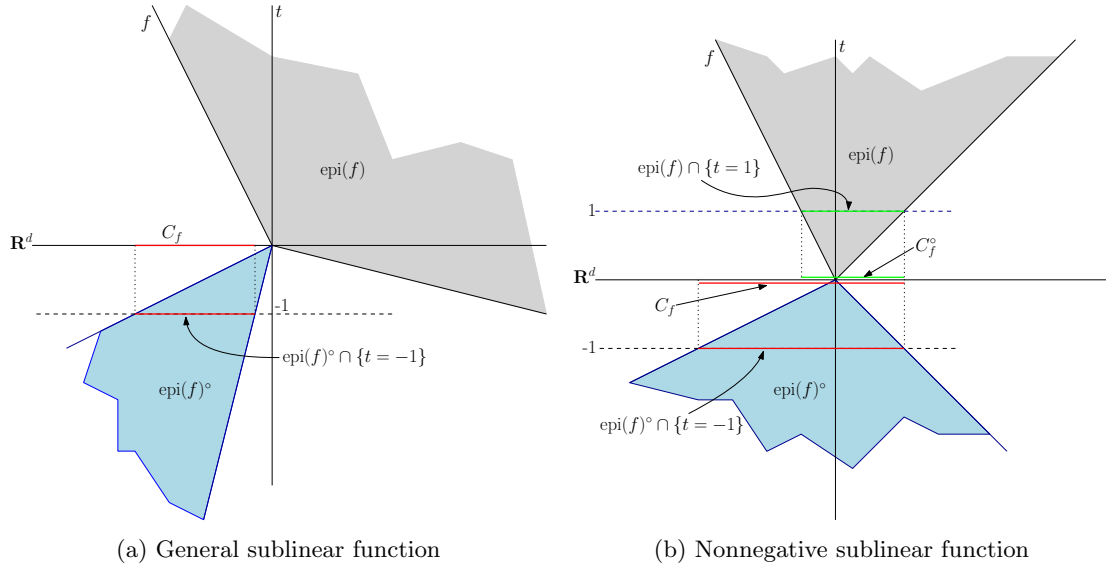


Figure 1: Illustration of Propositions 3.57 and 3.58

1377 When f is a nonnegative sublinear function, even more can be said.

1378 **Proposition 3.58.** Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a sublinear function that is nonnegative everywhere, and let C_f
 1379 be defined as in Theorem 3.56. Then $f = \gamma_{(C_f)^\circ}$, i.e., f is the gauge function for $(C_f)^\circ$. Consequently,
 1380 $(C_f)^\circ = \{\mathbf{y} \in \mathbb{R}^d : (\mathbf{y}, 1) \in \text{epi}(f)\} = \{\mathbf{y} \in \mathbb{R}^d : f(\mathbf{y}) \leq 1\}$.

1381 *Proof.* Since $f \geq 0$, $\text{epi}(f) \subseteq \{(\mathbf{r}, t) : t \geq 0\}$. Therefore, $(\mathbf{0}, -1) \in \text{epi}(f)^\circ$. By Proposition 3.57, $\mathbf{0} \in C_f$.
 1382 Moreover, by Theorems 3.56 and 3.52, $f = \sigma_{C_f} = \gamma_{(C_f)^\circ}$. By Theorem 3.40 part 2., this shows that
 1383 $(C_f)^\circ = \{\mathbf{y} \in \mathbb{R}^d : f(\mathbf{y}) \leq 1\}$. By Problem 10 from the “HW for Week XI”, we have that $(C_f)^\circ = \{\mathbf{y} \in \mathbb{R}^d :$
 1384 $(\mathbf{y}, 1) \in \text{epi}(f)\} = \{\mathbf{y} \in \mathbb{R}^d : f(\mathbf{y}) \leq 1\}$. \square

1385 3.6 Directional derivatives, subgradients and subdifferential calculus

1386 Let us look at directional derivatives of convex functions more closely. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be any function and
 1387 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:

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

1388 if that limit exists. We will be speaking of $f'(\mathbf{x}; \cdot)$ as a function from $\mathbb{R}^d \rightarrow \mathbb{R}$. When the function f is
 1389 convex, this function has very nice properties.

1390 **Lemma 3.59.** If $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex, the expression $\frac{f(\mathbf{x}+t\mathbf{r})-f(\mathbf{x})}{t}$ is a non-decreasing function of t , for
 1391 $t \neq 0$.

1392 *Proof.* By Proposition 3.26, the function $\phi(t) = f(\mathbf{x} + t\mathbf{r})$ is a convex function. By Proposition 3.24, we
 1393 observe that $\frac{\phi(t)-\phi(0)}{t}$ is a non-decreasing function of t . \square

1394 **Proposition 3.60.** Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a convex function, and let $\mathbf{x} \in \mathbb{R}^d$. Then the limit in (3.5) exists and
 1395 is finite for all $\mathbf{r} \in \mathbb{R}^d$ and the function $f'(\mathbf{x}; \cdot) : \mathbb{R}^d \rightarrow \mathbb{R}$ is sublinear.

1396 *Proof.* By Proposition 3.26, the function $\phi(t) = f(\mathbf{x} + t\mathbf{r})$ is a convex function, and $f'(\mathbf{x}; \mathbf{r}) = \lim_{t \downarrow 0} \frac{\phi(t)-\phi(0)}{t}$.
 1397 By Lemma 3.59, we observe that $\frac{\phi(t)-\phi(0)}{t}$ is a non-decreasing function of t for $t \neq 0$, and restricting to $t > 0$,

1398 $\frac{\phi(t)-\phi(0)}{t}$ 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
 1399 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 $\lambda \in (0, 1)$.

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

1400 where the inequality follows from convexity of f . By Proposition 3.32, the function f is sublinear. \square

1401 There is a nice connection with subgradients and subdifferentials – recall Definition 3.17. Also recall the
 1402 construction of the closed, convex set C_f from a sublinear function f from Theorem 3.56.

Theorem 3.61. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a convex function, and let $\mathbf{x} \in \mathbb{R}^d$. Then

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

1403 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

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

Thus, we have the following equivalences.

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

1404 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
 1405 Lemma 3.59, and the last equivalence follows from the definition of $C_{f'(\mathbf{x}; \cdot)}$ in (3.4). \square

1406 A characterization of differentiability for convex functions can be obtained using these concepts.

1407 **Theorem 3.62.** Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a convex function, and let $\mathbf{x} \in \mathbb{R}^d$. Then the following are equivalent.

1408 (i) f is differentiable at \mathbf{x} .

1409 (ii) $f'(\mathbf{x}; \cdot)$ is a linear function given by $f'(\mathbf{x}; \mathbf{r}) = \langle \mathbf{a}_{\mathbf{x}}, \mathbf{r} \rangle$ for some $\mathbf{a}_{\mathbf{x}} \in \mathbb{R}^d$.

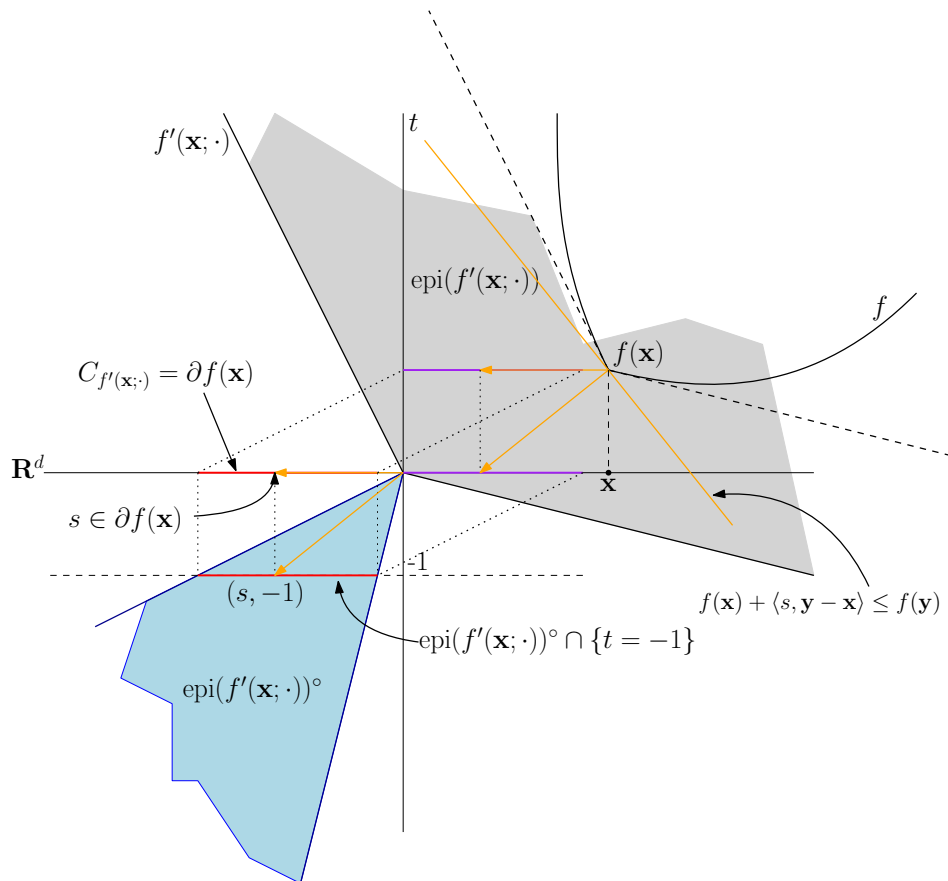


Figure 2: A picture illustrating the relationship between the sublinear function $f'(\mathbf{x}; \cdot)$, the set $C_{f'(\mathbf{x}; \cdot)}$, the subdifferential $\partial f(\mathbf{x})$, and an affine support hyperplane given by an element $s \in \partial f(\mathbf{x})$. Recall the relationships from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

1410 (iii) $\partial f(\mathbf{x})$ is a singleton, i.e., there is a unique subgradient for f at \mathbf{x} .

1411 Moreover, if any of the above conditions hold then $\nabla f(\mathbf{x}) = \mathbf{a}_\mathbf{x} = \mathbf{s}$, where \mathbf{s} is the unique subgradient in
1412 $\partial f(\mathbf{x})$.

1413 *Proof.* (i) \implies (ii). If f is differentiable, then it is well-known from calculus that $f'(\mathbf{x}; \mathbf{r}) = \langle \nabla f(\mathbf{x}), \mathbf{r} \rangle$;
1414 thus, setting $\mathbf{a}_\mathbf{x} = \nabla f(\mathbf{x})$ suffices.

(ii) \implies (iii). By Theorem 3.61 and (3.4), we obtain that

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

1415 We now observe that if $\langle \mathbf{s}, \mathbf{r} \rangle \leq \langle \mathbf{a}_\mathbf{x}, \mathbf{r} \rangle$ for all $\mathbf{r} \in \mathbb{R}^d$, then we must have $\mathbf{s} = \mathbf{a}_\mathbf{x}$. Therefore, $\partial f(\mathbf{x}) = \{\mathbf{a}_\mathbf{x}\}$.

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

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

1416 thus showing that f is differentiable at \mathbf{x} with gradient \mathbf{s} . In other words, given any $\delta > 0$, we must find
1417 $\epsilon > 0$ 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 $k \geq 1$ there exists \mathbf{h}_k such that $\|\mathbf{h}_k\| =: t_k \leq \frac{1}{k}$
and $\frac{|f(\mathbf{x} + \mathbf{h}_k) - f(\mathbf{x}) - \langle \mathbf{s}, \mathbf{h}_k \rangle|}{t_k} \geq \delta$. 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 $\{\frac{\mathbf{h}_k}{t_k}\}_{k=1}^\infty$ is the convergent sequence. Using Theorem 3.21, there exists a constant $L := L(B(\mathbf{0}, 1))$
such that $|f(\mathbf{y}) - f(\mathbf{z})| \leq L\|\mathbf{y} - \mathbf{z}\|$ for all $\mathbf{y}, \mathbf{z} \in B(\mathbf{0}, 1)$. Noting that \mathbf{h}_k and $t_k \mathbf{r}$ for all $k \geq 1$ are in the
unit ball $B(\mathbf{0}, 1)$ (since $t_k \leq \frac{1}{k}$),

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

1418 where the Cauchy-Schwartz inequality is used in the last inequality. We now let $k \rightarrow \infty$. The first term
1419 in the last expression above goes to 0, since $\frac{\mathbf{h}_k}{t_k}$ converges to \mathbf{r} . In the second term, $\frac{|f(\mathbf{x} + t_k \mathbf{r}) - f(\mathbf{x}) - \langle \mathbf{s}, t_k \mathbf{r} \rangle|}{t_k}$ goes to
1420 its limit which is the directional derivative $f'(\mathbf{x}; \mathbf{r})$. By Theorem 3.61, $f'(\mathbf{x}; \mathbf{r}) = \sup_{\mathbf{y} \in \partial f(\mathbf{x})} \langle \mathbf{y}, \mathbf{r} \rangle = \langle \mathbf{s}, \mathbf{r} \rangle$,
1421 because by assumption $\partial f(\mathbf{x}) = \{\mathbf{s}\}$. Thus, the second term in the last expression above also goes to 0. This
1422 contradicts $\delta > 0$. \square

1423 The following rules for manipulating subgradients and subdifferentials will be useful from an algorithmic
1424 perspective when we discuss optimization in the next section.

1425 **Theorem 3.63. Subdifferential calculus.** The following are all true.

1. Let $f_1, f_2 : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex functions and let $t_1, t_2 \geq 0$. 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 $\mathbf{b} \in \mathbb{R}^m$ and let $T(\mathbf{x}) = A\mathbf{x} + \mathbf{b}$ be the corresponding affine map from $\mathbb{R}^d \rightarrow \mathbb{R}^m$
and let $g : \mathbb{R}^m \rightarrow \mathbb{R}$ be a convex function. Then

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

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

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

1426 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
1427 relation, if one can impose a metric space structure on J such that $J(\mathbf{x})$ is a compact set.

1428 4 Optimization

1429 We now begin our study of the general convex optimization problem

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

1430 where $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a convex function, and C is a closed, convex set. We first observe that local minimizers
1431 are global minimizers for convex optimization problems.

1432 **Definition 4.1.** Let $g : \mathbb{R}^d \rightarrow \mathbb{R}$ be any function (not necessarily convex) and let $X \subseteq \mathbb{R}^d$ be any set (not
1433 necessarily convex). Then $\mathbf{x}^* \in X$ is said to be a *local minimizer* for the problem $\inf_{\mathbf{x} \in X} g(\mathbf{x})$ if there exists
1434 $\epsilon > 0$ such that $g(\mathbf{y}) \geq g(\mathbf{x}^*)$ for all $\mathbf{y} \in B(\mathbf{x}^*, \epsilon) \cap X$.

1435 $\mathbf{x}^* \in X$ is said to be a *global minimizer* if $g(\mathbf{y}) \geq g(\mathbf{x}^*)$ for all $\mathbf{y} \in X$.

1436 Note that if C is a compact, convex set, then (4.1) has a global minimizer by Weierstrass' Theorem
1437 (Theorem 1.11), because convex functions are continuous over the relative interior of their domain (Theo-
1438 rem 3.21).

1439 **Theorem 4.2.** Any local minimizer for (4.1) is a global minimizer.

1440 *Proof.* Let \mathbf{x}^* be a local minimizer, i.e., there exists $\epsilon > 0$ such that $f(\mathbf{y}) \geq f(\mathbf{x}^*)$ for all $\mathbf{y} \in B(\mathbf{x}^*, \epsilon) \cap C$.
1441 Suppose to the contrary that there exists $\bar{\mathbf{y}} \in C$ such that $f(\bar{\mathbf{y}}) < f(\mathbf{x}^*)$. Then $\bar{\mathbf{y}} \notin B(\mathbf{x}^*, \epsilon)$; otherwise, it
1442 would contradict $f(\mathbf{y}) \geq f(\mathbf{x}^*)$ for all $\mathbf{y} \in B(\mathbf{x}^*, \epsilon) \cap C$. Consider the line segment $[\mathbf{x}^*, \bar{\mathbf{y}}]$. It must intersect
1443 $B(\mathbf{x}^*, \epsilon)$ in a point other than \mathbf{x}^* . Therefore, there exists $1 > \lambda > 0$ such that $\bar{\mathbf{x}} = \lambda \mathbf{x}^* + (1 - \lambda)\bar{\mathbf{y}}$ is in
1444 $B(\mathbf{x}^*, \epsilon)$. By convexity of f , $f(\bar{\mathbf{x}}) \leq \lambda f(\mathbf{x}^*) + (1 - \lambda)f(\bar{\mathbf{y}})$. Since $\lambda \in (0, 1)$ and $f(\bar{\mathbf{y}}) < f(\mathbf{x}^*)$, this implies
1445 that $f(\bar{\mathbf{x}}) < f(\mathbf{x}^*)$. Moreover, since C is convex, $\bar{\mathbf{x}} \in C$, and so $\bar{\mathbf{x}} \in B(\mathbf{x}^*, \epsilon) \cap C$. This contradicts that
1446 $f(\mathbf{y}) \geq f(\mathbf{x}^*)$ for all $\mathbf{y} \in B(\mathbf{x}^*, \epsilon) \cap C$. \square

1447 We now give a characterization of global minimizers of (4.1) in terms of the local geometry of C and
1448 the first order properties of f , i.e., its subdifferential ∂f . We first need some concepts related to the local
1449 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

$$F_C(\mathbf{x}) = \{\mathbf{r} \in \mathbb{R}^d : \exists \epsilon > 0 \text{ such that } \mathbf{x} + \epsilon \mathbf{r} \in C\}.$$

1450 $F_C(\mathbf{x})$ may not be a closed cone – consider C as the unit circle in \mathbb{R}^2 and $\mathbf{x} = (-1, 0)$; then $F_C(\mathbf{x}) =$
1451 $\{\mathbf{r} \in \mathbb{R}^2 : \mathbf{r}_1 > 0\} \cup \{\mathbf{0}\}$. It is much nicer to work with its closure.

1452 **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 $T_C(\mathbf{x}) :=$
1453 $\text{cl}(F_C(\mathbf{x}))$.

1454 The final concept related to the local geometry of closed, convex sets will be the *normal cone*.

1455 **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 $N_C(\mathbf{x}) := \{\mathbf{r} \in$
1456 $\mathbb{R}^d : \langle \mathbf{r}, \mathbf{x} \rangle \geq \langle \mathbf{r}, \mathbf{y} \rangle \ \forall \mathbf{y} \in C\}$.

1457 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
 1458 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
 1459 immediate that N_C is a closed, convex cone. Note that any nonzero vector $\mathbf{r} \in N_C(\mathbf{x})$ defines a supporting
 1460 hyperplane $H(\mathbf{r}, \langle \mathbf{r}, \mathbf{x} \rangle)$ at \mathbf{x} .
 1461

1462 **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
 1463 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
 1464 the normal cone are polars of each other.

1465 *Proof.* See Problem 4 in “HW for Week X”. □

1466 We are now ready to state the characterization of a global minimizer of (4.1), in terms of the local
 1467 geometry of C and the first-order information of f .

1468 **Theorem 4.7.** Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a convex function, and C be a closed, convex set. Then the following are
 1469 all equivalent.

- 1470 1. \mathbf{x}^* is a global minimizer of (4.1).
- 1471 2. $f'(\mathbf{x}^*; \mathbf{y} - \mathbf{x}^*) \geq 0$ for all $\mathbf{y} \in C$.
- 1472 3. $f'(\mathbf{x}^*; \mathbf{r}) \geq 0$ for all $\mathbf{r} \in T_C(\mathbf{x}^*)$.
- 1473 4. $\mathbf{0} \in \partial f(\mathbf{x}^*) + N_C(\mathbf{x}^*)$.

1474 *Proof.* 1. \implies 2. Since $f(\mathbf{z}) \geq f(\mathbf{x}^*)$ for all $\mathbf{z} \in C$, in particular this holds for $\mathbf{z} = \mathbf{x}^* + t(\mathbf{y} - \mathbf{x}^*)$ for all
 1475 $0 \leq t \leq 1$. Therefore, $\frac{f(\mathbf{x}^* + t(\mathbf{y} - \mathbf{x}^*)) - f(\mathbf{x}^*)}{t} \geq 0$ for all $t \in (0, 1)$. Taking the limit as $t \rightarrow 0$, we obtain that
 1476 $f'(\mathbf{x}^*; \mathbf{y} - \mathbf{x}^*) \geq 0$.

1477 2. \implies 3. We first show that $f'(\mathbf{x}^*; \mathbf{r}) \geq 0$ for all $\mathbf{r} \in F_C(\mathbf{x}^*)$. Let $\epsilon > 0$ such that $\mathbf{y} = \mathbf{x}^* + \epsilon \mathbf{r} \in C$.
 1478 By assumption, $0 \leq f'(\mathbf{x}^*; \mathbf{y} - \mathbf{x}^*) = f'(\mathbf{x}^*; \epsilon \mathbf{r}) = \epsilon f'(\mathbf{x}^*; \mathbf{r})$, using the positive homogeneity of $f'(\mathbf{x}^*; \cdot)$, since
 1479 $f'(\mathbf{x}^*; \cdot)$ is sublinear by Proposition 3.60. Dividing by ϵ , we obtain that $f'(\mathbf{x}^*; \mathbf{r}) \geq 0$ for all $\mathbf{r} \in F_C(\mathbf{x}^*)$.
 1480 Since $f'(\mathbf{x}^*; \cdot)$ is sublinear, it is convex by Proposition 3.32, and thus, it is continuous by Theorem 3.21.
 1481 Consequently, it must be nonnegative on $T_C(\mathbf{x}^*) = \text{cl}(F_C(\mathbf{x}^*))$, because it is nonnegative on $F_C(\mathbf{x}^*)$.

1482 3. \implies 4. Suppose to the contrary that $\mathbf{0} \notin \partial f(\mathbf{x}^*) + N_C(\mathbf{x}^*)$. Since f is assumed to be finite-valued
 1483 everywhere, $\text{dom}(f) = \mathbb{R}^d$. Thus, by Problem 15 in “HW for Week IX”, $\partial f(\mathbf{x}^*)$ is a compact, convex set.
 1484 Moreover, $N_C(\mathbf{x}^*)$ is a closed, convex cone by Proposition 4.6. Therefore, by Problem 6 in “HW for Week
 1485 II”, $\partial f(\mathbf{x}^*) + N_C(\mathbf{x}^*)$ is a closed, convex set. By the separating hyperplane theorem (Theorem 2.20), there
 1486 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 $\mathbf{v} \in \partial f(\mathbf{x}^*) + N_C(\mathbf{x}^*)$.

1487 First, we claim that $\langle \mathbf{a}, \mathbf{n} \rangle \leq 0$ for all $\mathbf{n} \in N_C(\mathbf{x}^*)$. Otherwise, consider $\bar{\mathbf{n}} \in N_C(\mathbf{x}^*)$ such that $\langle \mathbf{a}, \bar{\mathbf{n}} \rangle > 0$.
 1488 Since $N_C(\mathbf{x}^*)$ is a convex cone, $\lambda \bar{\mathbf{n}} \in N_C(\mathbf{x}^*)$ for all $\lambda \geq 0$. But then consider any $\mathbf{s} \in \partial f(\mathbf{x}^*)$ (which is
 1489 nonempty by Problem 15 in “HW for Week IX”) and the set of points $\mathbf{s} + \lambda \bar{\mathbf{n}}$. Since $\langle \mathbf{a}, \bar{\mathbf{n}} \rangle > 0$, we can find
 1490 $\lambda \geq 0$ 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 $\mathbf{v} \in \partial f(\mathbf{x}^*) + N_C(\mathbf{x}^*)$.

1491 Since $\langle \mathbf{a}, \mathbf{n} \rangle \leq 0$ for all $\mathbf{n} \in N_C(\mathbf{x}^*)$, we obtain that $\mathbf{a} \in N_C(\mathbf{x}^*)^\circ = T_C(\mathbf{x}^*)$, by Proposition 4.6. Now
 1492 we use the fact that $\partial f(\mathbf{x}^*) \subseteq \partial f(\mathbf{x}^*) + N_C(\mathbf{x}^*)$, since $\mathbf{0} \in N_C(\mathbf{x}^*)$. This implies that $\langle \mathbf{a}, \mathbf{s} \rangle \leq \delta < 0$
 1493 for all $\mathbf{s} \in \partial f(\mathbf{x}^*)$. Since $\partial f(\mathbf{x}^*)$ is a compact, convex set, this implies that $\sup_{\mathbf{s} \in \partial f(\mathbf{x}^*)} \langle \mathbf{a}, \mathbf{s} \rangle < 0$. From
 1494 Theorem 3.61, $f'(\mathbf{x}^*; \mathbf{a}) = \sigma_{\partial f(\mathbf{x}^*)}(\mathbf{a}) = \sup_{\mathbf{s} \in \partial f(\mathbf{x}^*)} \langle \mathbf{a}, \mathbf{s} \rangle < 0$. This contradicts the assumption of 3., because
 1495 we showed above that $\mathbf{a} \in T_C(\mathbf{x}^*)$.

4. \implies 1. Consider any $\mathbf{y} \in C$. Since $\mathbf{0} \in \partial f(\mathbf{x}^*) + N_C(\mathbf{x}^*)$, there exist $\mathbf{s} \in \partial f(\mathbf{x}^*)$ and $\mathbf{n} \in N_C(\mathbf{x}^*)$ such
 that $\mathbf{0} = \mathbf{s} + \mathbf{n}$. Now, $\mathbf{y} - \mathbf{x}^* \in T_C(\mathbf{x}^*)$ and so $\langle \mathbf{y} - \mathbf{x}^*, \mathbf{n} \rangle \leq 0$ by Proposition 4.6. Since we have

$$0 = \langle \mathbf{y} - \mathbf{x}^*, \mathbf{0} \rangle = \langle \mathbf{y} - \mathbf{x}^*, \mathbf{s} \rangle + \langle \mathbf{y} - \mathbf{x}^*, \mathbf{n} \rangle,$$

1496 this implies that $\langle \mathbf{y} - \mathbf{x}^*, \mathbf{s} \rangle \geq 0$. By definition of subgradient, $f(\mathbf{y}) \geq f(\mathbf{x}^*) + \langle \mathbf{s}, \mathbf{y} - \mathbf{x}^* \rangle \geq f(\mathbf{x}^*)$. Since
 1497 the choice of $\mathbf{y} \in C$ was arbitrary, this shows that \mathbf{x}^* is a global minimizer. □

1498 **Corollary 4.8.** Let \mathbf{x}^* be a minimizer for (4.1). If $\mathbf{x}^* \in \text{int}(C)$, then $\mathbf{0} \in \partial f(\mathbf{x}^*)$. In particular, if f is real-
 1499 valued everywhere and $C = \mathbb{R}^d$, then a minimizer of f must contain $\mathbf{0}$ in its subdifferential. Consequently,
 1500 if f is differentiable everywhere, the gradient at the minimizer must be $\mathbf{0}$.

1501 *Proof.* This follows from the fact that for any convex set C and any $\mathbf{y} \in \text{int}(C)$, $N_C(\mathbf{y}) = \{\mathbf{0}\}$ (Why?). \square

1502 **Algorithmic setup: First-order oracles.** To tackle the problem (4.1) computationally, we have to set
 1503 up a precise way to access the values/subgradients of the function f and test if given points belong to the
 1504 set C or not. To make this algorithmically clean, we define *first-order oracles*.

1505 **Definition 4.9.** A *first order oracle* for a convex function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is an oracle/algorithm/black-box that
 1506 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
 1507 $C \subseteq \mathbb{R}^d$ is an oracle/algorithm/black-box that takes as input any $\mathbf{x} \in \mathbb{R}^d$ and either correctly reports that
 1508 $\mathbf{x} \in C$ 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
 1509 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*.

1510 4.1 Subgradient algorithm

1511 To build up towards an algorithm that assumes only first-order oracles for f and C , we will first look at the
 1512 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
 1513 report the closest point in C to \mathbf{x} (assuming C is nonempty). Recall that in the proof of Theorem 2.20, we
 1514 had shown that such a closest point always exists as long as C is a nonempty, closed, convex set. In fact,
 1515 the proof holds even for a closed set; convexity was not used to show the existence of a closest point. We
 1516 now strengthen the observation by showing that under the additional assumption of convexity, the closest
 1517 point is unique.

1518 **Proposition 4.10.** Let $C \subseteq \mathbb{R}^d$ be a nonempty, closed, convex set and let $\mathbf{x} \in \mathbb{R}^d$. Then there is a unique
 1519 point $\mathbf{x}^* \in C$ such that $\|\mathbf{x} - \mathbf{x}^*\| \leq \|\mathbf{x} - \mathbf{y}\|$ for all $\mathbf{y} \in C$.

Proof. If $\mathbf{x} \in C$, then the conclusion is true by setting $\mathbf{x}^* = \mathbf{x}$. So we assume $\mathbf{x} \notin C$. Following the proof
 of Theorem 2.20, there exists a closest point $\mathbf{x}^* \in C$ and $\mathbf{a} = \mathbf{x} - \mathbf{x}^*$ satisfies $\langle \mathbf{a}, \mathbf{y} - \mathbf{x}^* \rangle \leq 0$ for all $\mathbf{y} \in C$.
 Thus,

$$\|\mathbf{x} - \mathbf{y}\|^2 = \|\mathbf{a} + (\mathbf{x}^* - \mathbf{y})\|^2 = \|\mathbf{a}\|^2 + \langle \mathbf{a}, \mathbf{x}^* - \mathbf{y} \rangle + \|\mathbf{x}^* - \mathbf{y}\|^2 > \|\mathbf{x}^* - \mathbf{y}\|^2,$$

1520 where the last inequality follows from the fact that $\mathbf{a} \neq \mathbf{0}$ and $\langle \mathbf{a}, \mathbf{x}^* - \mathbf{y} \rangle \geq 0$. \square

1521 **Definition 4.11.** $\text{Proj}_C(\mathbf{x})$ will denote the unique closest point (under the standard Euclidean norm) in C
 1522 to \mathbf{x} .

1523 Note that an oracle that reports $\text{Proj}_C(\mathbf{x})$ for any $\mathbf{x} \in \mathbb{R}^d$ is stronger than a separation oracle for C ,
 1524 because $\text{Proj}_C(\mathbf{x}) = \mathbf{x}$ if and only if $\mathbf{x} \in C$, and when $\text{Proj}_C(\mathbf{x}) \neq \mathbf{x}$, then one can use $\mathbf{a} = \mathbf{x} - \text{Proj}_C(\mathbf{x})$
 1525 and $\delta = \langle \mathbf{a}, \text{Proj}_C(\mathbf{x}) \rangle$ as a separating hyperplane; see the proof of Theorem 2.20. Even so, for “simple”
 1526 sets C , computing $\text{Proj}_C(\mathbf{x})$ is not a difficult task. For example, when $C = \mathbb{R}_+^d$, then $\text{Proj}_C(\mathbf{x}) = \mathbf{y}$, where
 1527 $y_i = \max\{0, x_i\}$ for all $i = 1, \dots, d$.

1528 We now give a simple and elegant algorithm to solve the problem (4.1) when one has access to an oracle
 1529 that can output $\text{Proj}_C(\mathbf{x})$ for any $\mathbf{x} \in \mathbb{R}^d$, and a first-order oracle for $f : \mathbb{R}^d \rightarrow \mathbb{R}$. The algorithm does not
 1530 assume any properties beyond convexity for the function f (e.g., differentiability). Note that, in particular,
 1531 when we have no constraints, i.e., $C = \mathbb{R}^n$, then $\text{Proj}_C(\mathbf{x}) = \mathbf{x}$ for all $\mathbf{x} \in \mathbb{R}^n$. Therefore, this algorithm can
 1532 be used for *unconstrained optimization of general convex functions* with only a first-order oracle for f .

1533 **Subgradient Algorithm.**

- 1534 1. Choose any sequence h_0, h_1, \dots , of strictly positive numbers. Let $\mathbf{x}^0 \in C$ (which can be found by
 1535 taking an arbitrary point in \mathbb{R}^d and projecting to C).
- 1536 2. For $i = 0, 1, 2, \dots$, do
- 1537 (a) Use the first-order oracle for f to get some $\mathbf{s}^i \in \partial f(\mathbf{x}^i)$. If $\mathbf{s}^i = \mathbf{0}$, then stop and report \mathbf{x}^i as the
 1538 optimal point.
- 1539 (b) Set $\mathbf{x}^{i+1} = \text{Proj}_C(\mathbf{x}^i - h_i \frac{\mathbf{s}^i}{\|\mathbf{s}^i\|})$.

1540 The points $\mathbf{x}^0, \mathbf{x}^1, \dots$ will be called the *iterates* of the Subgradient Algorithm. We now do a simple
 1541 convergence analysis for the algorithm. First, a simple observation about the point $\text{Proj}_C(\mathbf{x})$.

Lemma 4.12. Let $C \subseteq \mathbb{R}^d$ be a closed, convex set, let $\mathbf{x}^* \in C$ and $\mathbf{x} \in \mathbb{R}^d$ (not necessarily in C). Then

$$\|\text{Proj}_C(\mathbf{x}) - \mathbf{x}^*\| \leq \|\mathbf{x} - \mathbf{x}^*\|.$$

Proof. The interesting case is when $\mathbf{x} \notin C$. The proof of Theorem 2.20 shows that if we set $\mathbf{a} = \mathbf{x} - \text{Proj}_C(\mathbf{x})$, then $\langle \mathbf{a}, \text{Proj}_C(\mathbf{x}) - \mathbf{y} \rangle \geq 0$ for all $\mathbf{y} \in C$; in particular, $\langle \mathbf{a}, \text{Proj}_C(\mathbf{x}) - \mathbf{x}^* \rangle \geq 0$. We now observe that

$$\begin{aligned} \|\mathbf{x} - \mathbf{x}^*\|^2 &= \|\mathbf{x} - \text{Proj}_C(\mathbf{x}) + \text{Proj}_C(\mathbf{x}) - \mathbf{x}^*\|^2 \\ &= \|\mathbf{a} + \text{Proj}_C(\mathbf{x}) - \mathbf{x}^*\|^2 \\ &= \|\mathbf{a}\|^2 + \|\text{Proj}_C(\mathbf{x}) - \mathbf{x}^*\|^2 + 2\langle \mathbf{a}, \text{Proj}_C(\mathbf{x}) - \mathbf{x}^* \rangle \\ &\geq \|\text{Proj}_C(\mathbf{x}) - \mathbf{x}^*\|^2, \end{aligned}$$

1542 since $\langle \mathbf{a}, \text{Proj}_C(\mathbf{x}) - \mathbf{x}^* \rangle \geq 0$. □

Theorem 4.13. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a convex function, and let $\mathbf{x}^* \in \arg \min_{\mathbf{x} \in C} f(\mathbf{x})$ (i.e, we assume a minimizer exists for the problem). Suppose $\mathbf{x}^0 \in B(\mathbf{x}^*, R)$ for some real number $R \geq 0$. Let $M := M(B(\mathbf{x}^*, R))$ be a Lipschitz constant for f , guaranteed to exist by Theorem 3.21, i.e., $|f(\mathbf{x}) - f(\mathbf{y})| \leq M\|\mathbf{x} - \mathbf{y}\|$ for all $\mathbf{x}, \mathbf{y} \in B(\mathbf{x}^*, R)$. Let $\mathbf{x}^0, \mathbf{x}^1, \dots$ be the sequence of iterates obtained by the Subgradient Algorithm above and assume that a zero subgradient was not reported in any iteration. Then, for every $k \geq 0$,

$$\min_{i=0, \dots, k} f(\mathbf{x}^i) \leq f(\mathbf{x}^*) + M \left(\frac{R^2 + \sum_{i=0}^k h_i^2}{2 \sum_{i=0}^k h_i} \right).$$

Proof. Define $r_i = \|\mathbf{x}^i - \mathbf{x}^*\|$ and $v_i = \frac{\langle \mathbf{s}^i, \mathbf{x}^i - \mathbf{x}^* \rangle}{\|\mathbf{s}^i\|}$ for $i = 0, 1, 2, \dots$. Note that $v_i \geq 0$ for all $i \geq 0$ since \mathbf{s}^i is a subgradient at \mathbf{x}^i and \mathbf{x}^* is the minimizer (Verify!!). We next observe that

$$\begin{aligned} r_{i+1}^2 &= \|\text{Proj}_C(\mathbf{x}^i - h_i \frac{\mathbf{s}^i}{\|\mathbf{s}^i\|}) - \mathbf{x}^*\|^2 \\ &\leq \|\mathbf{x}^i - h_i \frac{\mathbf{s}^i}{\|\mathbf{s}^i\|} - \mathbf{x}^*\|^2 && \text{by Lemma 4.12} \\ &= \|\mathbf{x}^i - \mathbf{x}^*\|^2 + h_i^2 - 2h_i v_i \\ &= r_i^2 + h_i^2 - 2h_i v_i \end{aligned}$$

1543 Adding these inequalities for $i = 0, 1, \dots, k$, we obtain that

$$r_{k+1}^2 \leq r_0^2 + \sum_{i=0}^k h_i^2 - 2 \sum_{i=0}^k h_i v_i. \tag{4.2}$$

Let $v_{\min} = \min_{i=0, \dots, k} v_i$ and let i^{\min} be such that $v_{\min} = v_{i^{\min}}$. Using the fact that $r_0^2 = \|\mathbf{x}^0 - \mathbf{x}^*\|^2 \leq R^2$, and that $r_{k+1}^2 \geq 0$, we obtain from (4.2) that

$$v_{\min} (2 \sum_{i=0}^k h_i) \leq 2 \sum_{i=0}^k h_i v_i \leq R^2 + \sum_{i=0}^k h_i^2.$$

1544 Consequently,

$$v_{\min} \leq \frac{R^2 + \sum_{i=0}^k h_i^2}{2 \sum_{i=0}^k h_i}. \quad (4.3)$$

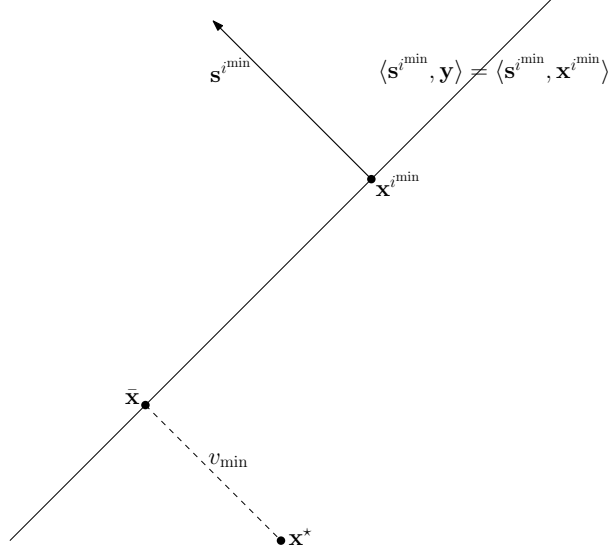


Figure 3: Using v_{\min} to bound the function value. The line through $\bar{\mathbf{x}}$ and $\mathbf{x}^{i^{\min}}$ represents the hyperplane $H := H(\mathbf{s}^{i^{\min}}, \langle \mathbf{s}^{i^{\min}}, \mathbf{x}^{i^{\min}} \rangle)$.

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}^* ; see Figure 3. By Problem 12 in “HW for Week IX”, $v_{\min} = \|\bar{\mathbf{x}} - \mathbf{x}^*\|$. Moreover, $v_{\min} \leq v_0 \leq \|\mathbf{x}^0 - \mathbf{x}^*\| \leq R$. Therefore, $\bar{\mathbf{x}} \in B(\mathbf{x}^*, R)$. Using the Lipschitz constant M , we obtain that $f(\bar{\mathbf{x}}) \leq f(\mathbf{x}^*) + Mv_{\min}$. 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 $\bar{\mathbf{x}}, \mathbf{x}^{i^{\min}} \in H$. Therefore, we obtain

$$\min_{i=0, \dots, k} f(\mathbf{x}^i) \leq f(\mathbf{x}^{i^{\min}}) \leq f(\bar{\mathbf{x}}) \leq f(\mathbf{x}^*) + Mv_{\min} \leq f(\mathbf{x}^*) + M \left(\frac{R^2 + \sum_{i=0}^k h_i^2}{2 \sum_{i=0}^k h_i} \right),$$

1545 where the last inequality follows from (4.3). □

1546 If we fix the number of steps of the algorithm to be $N \in \mathbb{N}$, then the choice of h_0, \dots, h_N that minimizes
 1547 $\frac{R^2 + \sum_{i=0}^N h_i^2}{2 \sum_{i=0}^N h_i}$ is where $h_i = \frac{R}{\sqrt{N+1}}$ for all $i = 0, \dots, N$, which yields the following corollary.

Corollary 4.14. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a convex function, and let $\mathbf{x}^* \in \arg \min_{\mathbf{x} \in C} f(\mathbf{x})$. Suppose $\mathbf{x}^0 \in B(\mathbf{x}^*, R)$ for some real number $R \geq 0$. Let $M := M(B(\mathbf{x}^*, R))$ 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 $i = 0, \dots, N$. Then the iterates of the Subgradient Algorithm, with this choice of h_i , satisfy

$$\min_{i=0, \dots, N} f(\mathbf{x}^i) \leq f(\mathbf{x}^*) + \frac{MR}{\sqrt{N+1}}.$$

1548 Turning this around, if we want to be within ϵ of the optimal value $f(\mathbf{x}^*)$ for some $\epsilon > 0$, we should run
 1549 the Subgradient Algorithm for $\frac{M^2 R^2}{\epsilon^2}$ iterates, with $h_i = \frac{\epsilon}{M}$.

1550 If we theoretically let the algorithm run for infinitely many steps, we would hope to make the difference
 1551 between $\min_i f(\mathbf{x}^i)$ and $f(\mathbf{x}^*)$ go to 0 in the limit. This, of course, depends on the choice of the sequence
 1552 h_0, h_1, \dots so that the expression $\frac{R^2 + \sum_{i=0}^k h_i^2}{2 \sum_{i=0}^k h_i} \rightarrow 0$ as $k \rightarrow \infty$. There is a general sufficient condition that
 1553 guarantees this.

Proposition 4.15. Let $\{h_i\}_{i=0}^\infty$ be a sequence of strictly positive real numbers such that $\lim_{i \rightarrow \infty} h_i = 0$ and $\sum_{i=1}^\infty h_i = \infty$ (e.g., $h_i = \frac{1}{i}$). Then, for any real number R ,

$$\lim_{k \rightarrow \infty} \frac{R^2 + \sum_{i=0}^k h_i^2}{2 \sum_{i=0}^k h_i} = 0.$$

1554 **Remark 4.16.** Corollary 4.14 shows that the subgradient algorithm has a convergence that is *independent*
 1555 of the dimension! Now matter how large d is, as long as one can access subgradients for f and project to
 1556 C , the number of iterations needed to converge to within ϵ is $O(\frac{1}{\epsilon^2})$. This is important to keep in mind for
 1557 applications where the dimension is extremely large.

1558 4.2 Generalized inequalities and convex mappings

1559 We first review the notion of a partially ordered set.

1560 **Definition 4.17.** Let X be any set. A *partial order* on X is a binary relation on X , i.e., a subset $\mathcal{R} \subseteq X \times X$
 1561 that satisfies certain conditions. We will denote $x \preceq y$ for $x, y \in X$ if $(x, y) \in \mathcal{R}$. The conditions are as
 1562 follows:

- 1563 1. $x \preceq x$ for all $x \in X$.
- 1564 2. $x \preceq y$ and $y \preceq z$ implies $x \preceq z$.
- 1565 3. $x \preceq y$ and $y \preceq x$ if and only if $x = y$.

1566 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
 1567 mindful of the vector space structure of \mathbb{R}^m .

1568 **Definition 4.18.** We will say that a binary relation on \mathbb{R}^m is a *generalized inequality*, if it satisfies the
 1569 following conditions.

- 1570 1. $\mathbf{x} \preceq \mathbf{x}$ for all $\mathbf{x} \in \mathbb{R}^m$.
- 1571 2. $\mathbf{x} \preceq \mathbf{y}$ and $\mathbf{y} \preceq \mathbf{z}$ implies $\mathbf{x} \preceq \mathbf{z}$.
- 1572 3. $\mathbf{x} \preceq \mathbf{y}$ and $\mathbf{y} \preceq \mathbf{x}$ if and only if $\mathbf{x} = \mathbf{y}$.
- 1573 4. $\mathbf{x} \preceq \mathbf{y}$ implies $\mathbf{x} + \mathbf{z} \preceq \mathbf{y} + \mathbf{z}$ for all $\mathbf{z} \in \mathbb{R}^m$.
- 1574 5. $\mathbf{x} \preceq \mathbf{y}$ implies $\lambda \mathbf{x} \preceq \lambda \mathbf{y}$ for all $\lambda \geq 0$.

1575 Generalized inequalities have an elegant geometric characterization.

1576 **Proposition 4.19.** Let $K \subseteq \mathbb{R}^m$ be a closed, convex, pointed cone. Then, the relation on \mathbb{R}^m defined by
 1577 $\mathbf{x} \preceq_K \mathbf{y}$ if and only if $\mathbf{y} - \mathbf{x} \in K$, is a generalized inequality. In this case, we say that \preceq_K is the generalized
 1578 inequality *induced by* K .

1579 Conversely, any generalized inequality \preceq is induced by a unique convex cone given by $K_\preceq = \{\mathbf{x} \in \mathbb{R}^d :$
 1580 $\mathbf{0} \preceq \mathbf{x}\}$. In other words, \preceq is the same relation as \preceq_{K_\preceq} .

1581 *Proof.* Left as an exercise. □

1582 **Example 4.20.** Here are some examples of generalized inequalities.

- 1583 1. $K = \mathbb{R}_+^m$ induces the generalized inequality $\mathbf{x} \preceq_K \mathbf{y}$ if and only if $\mathbf{x}_i \leq \mathbf{y}_i$ for all $i = 1 \dots, m$. This is
1584 often abbreviated to $\mathbf{x} \leq \mathbf{y}$, and is sometimes called the “canonical” generalized inequality on \mathbb{R}^m .
- 1585 2. $K = \{\mathbf{x} \in \mathbb{R}^d : \sqrt{\mathbf{x}_1^2 + \dots + \mathbf{x}_{d-1}^2} \leq \mathbf{x}_d\}$. This cone is called the *Lorentz cone*, and the corresponding
1586 generalized inequality is called a *second order cone constraints (SOCC)*.
- 1587 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
1588 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
1589 matrices that are positive semidefinite; see Definition 1.19. The corresponding generalized inequality
1590 on \mathbb{R}^{n^2} is called the *positive semidefinite cone constraint*.

1591 We would like to extend the notion of convex functions to vector valued maps, for which we will use the
1592 notion of generalized inequalities.

Definition 4.21. Let \preceq_K be a generalized inequality on \mathbb{R}^m induced by the cone K . We say that $G : \mathbb{R}^d \rightarrow \mathbb{R}^m$ is a *K-convex mapping* if

$$G(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) \preceq_K \lambda G(\mathbf{x}) + (1 - \lambda)G(\mathbf{y})$$

1593 for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ and $\lambda \in (0, 1)$.

1594 **Example 4.22.** Here are some examples of *K-convex mappings*.

- 1595 1. Let $K \subseteq \mathbb{R}^m$ be any closed, convex, pointed cone. If $G : \mathbb{R}^d \rightarrow \mathbb{R}^m$ is an affine map, i.e., there exist a
1596 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 \preceq be the *positive semidefinite cone constraint* from part 3. of Example 4.20, i.e., induced by the cone K of positive semidefinite matrices. Let A_0, A_1, \dots, A_d be fixed $p \times n$ matrices, for some $p \in \mathbb{N}$ (not necessarily equal to n). Define $G : \mathbb{R}^d \times \mathbb{R} \rightarrow \mathbb{R}^{n^2}$ to be the mapping

$$G(\mathbf{x}, s) = (A_0 + \mathbf{x}_1 A_1 + \dots + \mathbf{x}_d A_d)^T (A_0 + \mathbf{x}_1 A_1 + \dots + \mathbf{x}_d A_d) - B,$$

1597 where B is an arbitrary $n \times n$ matrix. Then G is a *K-convex mapping*.

- 1598 3. Let $K = \mathbb{R}_+^m$, and let $g_1, \dots, g_m : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex functions. Let $G : \mathbb{R}^d \rightarrow \mathbb{R}^m$ be defined as
1599 $G(\mathbf{x}) = (g_1(\mathbf{x}), \dots, g_m(\mathbf{x}))$, then G is a *K-convex mapping*.

1600 4.3 Convex optimization with generalized inequalities

1601 We can now define a very general framework for convex optimization problems, which is more concrete than
1602 the abstraction level of black-box first-order oracles, but is still flexible enough to incorporate the majority
1603 of convex optimization problems that show up in practice.

1604 **Definition 4.23.** Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a convex function, let $K \subseteq \mathbb{R}^m$ be a closed, convex, pointed cone,
1605 and let $G : \mathbb{R}^d \rightarrow \mathbb{R}^m$ be a *K-convex mapping*. Then f, K, G define a *convex optimization problem with*
1606 *generalized constraints* given as follows

$$\inf\{f(\mathbf{x}) : G(\mathbf{x}) \preceq_K \mathbf{0}\}. \tag{4.4}$$

1607 Problem 3 in “HW for Week XI” shows that the set $C = \{\mathbf{x} \in \mathbb{R}^d : G(\mathbf{x}) \preceq_K \mathbf{0}\}$ is a convex set, when G
1608 is a *K-convex mapping*. Thus, (4.4) is a special case of (4.1).

1609 **Example 4.24.** 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 $G : \mathbb{R}^d \rightarrow \mathbb{R}^m$ 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

$$\inf\{\langle \mathbf{c}, \mathbf{x} \rangle : A\mathbf{x} \leq \mathbf{b}\}$$

1610 which is the problem of minimizing a linear function over a polyhedron. This is more commonly known
1611 as a *linear program*, in accordance with the fact that the objective and the constraints are all linear.

1612 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
1613 convex function (see Problem 14 from “HW for Week IX”). With K and G as above, (4.4) is called a
1614 *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 $\mathbf{c} \in \mathbb{R}^{n^2}$, let $K \subseteq \mathbb{R}^{n^2}$ be the positive semidefinite cone, including the positive semidefinite cone constraint, and let $G : \mathbb{R}^d \rightarrow \mathbb{R}^{n^2}$ be an affine map, i.e., there exist $n \times n$ matrices F_0, F_1, \dots, F_d such that $G(\mathbf{x}) = F_0 + \mathbf{x}_1 F_1 + \dots + \mathbf{x}_d F_d$. Then (4.4) becomes

$$\inf\{\langle \mathbf{c}, \mathbf{x} \rangle : -F_0 - \mathbf{x}_1 F_1 - \dots - \mathbf{x}_d F_d \text{ is a PSD matrix}\}.$$

1615 This is known as a *semidefinite program*.

3. **Convex optimization with explicit constraints.** Let $f, g_1, \dots, g_m : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex functions. Define $K = \mathbb{R}_+^m$ and define $G : \mathbb{R}^d \rightarrow \mathbb{R}^m$ as $G(\mathbf{x}) = (g_1(\mathbf{x}), \dots, g_m(\mathbf{x}))$, which is the K -convex mapping from Example 4.22. Then (4.4) becomes

$$\inf\{f(\mathbf{x}) : g_1(\mathbf{x}) \leq 0, \dots, g_m(\mathbf{x}) \leq 0\}.$$

1616 4.3.1 Lagrangian duality for convex optimization with generalized constraints

1617 Given that the Subgradient Algorithm is a simple and elegant method for solving unconstrained problems,
1618 or problems with “simple” constraint sets C (i.e., when one can compute $\text{Proj}_C(\mathbf{x})$ efficiently), we will try to
1619 transform convex optimization problems with more complicated constraints into ones with simple constraints.
1620 This is the motivation for what is known as *Lagrangian duality*.

1621 Note that problem (4.4) is equivalent to the problem

$$\inf_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + I_{-K}(G(\mathbf{x})), \tag{4.5}$$

1622 where I_{-K} is the indicator function for the cone $-K$. It can be shown that the function $I_{-K} \circ G$ is a
1623 convex function – see Problem 4 from “HW for Week XI”. Thus, problem (4.5) is an unconstrained convex
1624 optimization problem. However, indicator functions are nasty to deal with because they are not finite valued,
1625 and thus, obtaining subgradients at all points becomes impossible. Thus, we try to replace I_{-K} with a “nicer”
1626 penalty function $p : \mathbb{R}^m \rightarrow \mathbb{R}$, which is not that wildly discontinuous, and is finite-valued everywhere. So we
1627 would be looking at the problem

$$\inf_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + p(G(\mathbf{x})). \tag{4.6}$$

1628 What properties should we require from our penalty function? First we would like problem (4.6) to be a
1629 convex problem, thus, we impose that

$$p \circ G : \mathbb{R}^d \rightarrow \mathbb{R} \text{ is a convex function.} \tag{4.7}$$

1630 Next, from an optimization perspective, we would like to have a guaranteed relationship between the function
1631 $f(\mathbf{x}) + I_{-K}(G(\mathbf{x}))$ and the function $f(\mathbf{x}) + p(G(\mathbf{x}))$. It turns out that a nice property to have is the guarantee
1632 that $f(\mathbf{x}) + p(G(\mathbf{x})) \leq f(\mathbf{x}) + I_{-K}(G(\mathbf{x}))$ for all $\mathbf{x} \in \mathbb{R}^d$. This can be achieved by imposing that

$$p \text{ is an underestimator of } I_{-K}, \text{ i.e., } p \leq I_{-K}. \tag{4.8}$$

1633 Lagrangian duality theory is the study of penalty functions p that are *linear* on \mathbb{R}^m , and satisfy the two
 1634 conditions highlighted above. Now a function $p : \mathbb{R}^m \rightarrow \mathbb{R}$ is linear if and only if there exists $\mathbf{c} \in \mathbb{R}^m$ such
 1635 that $p(\mathbf{z}) = \langle \mathbf{c}, \mathbf{z} \rangle$. The following proposition characterizes linear functions that satisfy the two conditions
 1636 above.

1637 **Proposition 4.25.** Let $p : \mathbb{R}^m \rightarrow \mathbb{R}$ 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
 1638 following are equivalent:

- 1639 1. p satisfies condition (4.8).
- 1640 2. $\mathbf{c} \in -K^\circ$, i.e., $-\mathbf{c}$ is in the polar of K .
- 1641 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.,

$$\begin{aligned} & \langle \mathbf{c}, \mathbf{z} \rangle \leq 0 \quad \text{for all } \mathbf{z} \in -K \\ \Leftrightarrow & \langle \mathbf{c}, -\mathbf{z} \rangle \leq 0 \quad \text{for all } \mathbf{z} \in K \\ \Leftrightarrow & \langle -\mathbf{c}, \mathbf{z} \rangle \leq 0 \quad \text{for all } \mathbf{z} \in K \\ \Leftrightarrow & -\mathbf{c} \in K^\circ \\ \Leftrightarrow & \mathbf{c} \in -K^\circ \end{aligned}$$

1642

(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 $\lambda \in (0, 1)$,

$$\begin{aligned} & \langle \mathbf{c}, \lambda G(\mathbf{x}) + (1 - \lambda)G(\mathbf{y}) - G(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) \rangle \geq 0 \\ \implies & \langle \mathbf{c}, \lambda G(\mathbf{x}) \rangle + \langle \mathbf{c}, (1 - \lambda)G(\mathbf{y}) \rangle \geq \langle \mathbf{c}, G(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) \rangle \\ \implies & \lambda \langle \mathbf{c}, G(\mathbf{x}) \rangle + (1 - \lambda)\langle \mathbf{c}, G(\mathbf{y}) \rangle \geq \langle \mathbf{c}, G(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) \rangle \\ \implies & \lambda p(G(\mathbf{x})) + (1 - \lambda)p(G(\mathbf{y})) \geq p(G(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y})) \end{aligned}$$

1643 Hence, condition (4.7) is satisfied.

1644 (3. \implies 1.) Trivial. □

1645 **Definition 4.26.** The set $-K^\circ$ is important in Lagrangian duality, and a separate notation and name has
 1646 been invented: $-K^\circ$ is called the *dual cone* of K and is denoted by K^* .

1647 The above discussions show that for any $\mathbf{y} \in K^*$, the optimal value of the (4.6), with p given by
 1648 $p(\mathbf{z}) = \langle \mathbf{y}, \mathbf{z} \rangle$, is a lower bound on the optimal value of (4.4). This motivates definition of the so-called *dual*
 1649 *function* $\mathcal{L} : \mathbb{R}^m \rightarrow \mathbb{R}$ associated with (4.4) as follows:

$$\mathcal{L}(\mathbf{y}) := \inf_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + \langle \mathbf{y}, G(\mathbf{x}) \rangle \tag{4.9}$$

1650 We state the lower bound property formally.

1651 **Proposition 4.27.** [Weak Duality] Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex, let $K \subseteq \mathbb{R}^m$ be a closed, convex, pointed
 1652 cone, and let $G : \mathbb{R}^d \rightarrow \mathbb{R}^m$ be a K -convex mapping. Let $\mathcal{L} : \mathbb{R}^m \rightarrow \mathbb{R}$ be as defined in (4.9). Then, for all
 1653 $\bar{\mathbf{x}} \in \mathbb{R}^d$ such that $G(\bar{\mathbf{x}}) \preceq_K \mathbf{0}$ and all $\bar{\mathbf{y}} \in K^*$, we must have $\mathcal{L}(\bar{\mathbf{y}}) \leq f(\bar{\mathbf{x}})$. Consequently, $\mathcal{L}(\bar{\mathbf{y}}) \leq \inf\{f(\mathbf{x}) : G(\mathbf{x}) \preceq_K \mathbf{0}\}$.
 1654

Proof. We simply follow the inequalities

$$\begin{aligned} \mathcal{L}(\bar{\mathbf{y}}) &= \inf_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + \langle \bar{\mathbf{y}}, G(\mathbf{x}) \rangle \\ &\leq f(\bar{\mathbf{x}}) + \langle \bar{\mathbf{y}}, G(\bar{\mathbf{x}}) \rangle \\ &\leq f(\bar{\mathbf{x}}), \end{aligned}$$

1655 where the last inequality holds because $G(\bar{\mathbf{x}}) \preceq_K \mathbf{0}$ and $\bar{\mathbf{y}} \in K^*$, and so $\langle \bar{\mathbf{y}}, G(\bar{\mathbf{x}}) \rangle \leq 0$. □

1656 Proposition 4.27 shows that any $\mathbf{y} \in K^*$ provides the lower bound $\mathcal{L}(\mathbf{y})$ on the optimal value of the
 1657 optimization problem (4.4). The *Lagrangian dual optimization problem* is the problem of finding the $\mathbf{y} \in K^*$
 1658 that provides the *best/largest* lower bound. In other words, the Lagrangian dual problem is defined as

$$\sup_{\mathbf{y} \in K^*} \mathcal{L}(\mathbf{y}), \tag{4.10}$$

1659 and Proposition 4.27 can be restated as

$$\sup\{\mathcal{L}(\mathbf{y}) : \mathbf{y} \in K^*\} \leq \inf\{f(\mathbf{x}) : G(\mathbf{x}) \preceq_K \mathbf{0}\}. \tag{4.11}$$

1660 If we have equality in (4.11), then to solve (4.4), one can instead solve (4.10). This merits a definition.

1661 **Definition 4.28** (Strong Duality). We say that we have a *zero duality gap* if equality holds in (4.11). In
 1662 addition, if the supremum in (4.10) is attained for some $\mathbf{y} \in K^*$, then we say that *strong duality* holds.

1663 4.3.2 Solving the Lagrangian dual problem

1664 Before we investigate conditions under which we have zero duality gap or strong duality, let us try to see
 1665 how one could use the subgradient algorithm to solve (4.10).

1666 **Proposition 4.29.** $\mathcal{L}(\mathbf{y})$ 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

$$\begin{aligned} -\mathcal{L}(\mathbf{y}) &= -\inf_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + \langle \mathbf{y}, G(\mathbf{x}) \rangle \\ &= \sup_{\mathbf{x} \in \mathbb{R}^d} -f(\mathbf{x}) + \langle \mathbf{y}, -G(\mathbf{x}) \rangle, \end{aligned}$$

1667 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
 1668 Theorem 3.12, $-\mathcal{L}(\mathbf{y})$ is convex in \mathbf{y} . □

1669 We could now use the subgradient algorithm to solve (4.10), if we had a first order oracle for $-\mathcal{L}(\mathbf{y})$
 1670 and an algorithm to project to K^* . We show that a subgradient for $-\mathcal{L}(\mathbf{y})$ can be found by solving an
 1671 unconstrained convex optimization problem.

1672 **Proposition 4.30.** 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 $-G(\bar{\mathbf{x}}) \in \partial(-\mathcal{L})(\bar{\mathbf{y}})$.

1673 *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
 1674 3. of Theorem 3.63, 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
 1675 simply $\{-G(\bar{\mathbf{x}})\}$. □

1676 Now if we have an algorithm that can compute $\text{Proj}_{K^*}(\mathbf{y})$ for all $\mathbf{y} \in \mathbb{R}^m$, then using Propositions 4.29
 1677 and 4.30, one can solve the Lagrangian dual problem (4.10), where in each iteration of the algorithm, one
 1678 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 $\bar{\mathbf{y}} \in K^*$. This can, in turn, be solved
 1679 by the subgradient algorithm if one has the appropriate first order oracles for $f(\mathbf{x})$ and $\langle \bar{\mathbf{y}}, G(\mathbf{x}) \rangle$.

1680 4.3.3 Explicit examples of the Lagrangian dual

1681 We will now explore some special settings of convex optimization problems with generalized inequalities, and
 1682 see that the Lagrangian dual has a particularly nice form.

1683 **Conic optimization.** Let $K \subseteq \mathbb{R}^m$ be a closed, convex, pointed cone. Let $G : \mathbb{R}^d \rightarrow \mathbb{R}^m$ be an affine map
 1684 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 $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a linear function given by
 1685 $f(\mathbf{x}) = \langle \mathbf{c}, \mathbf{x} \rangle$ for some $\mathbf{c} \in \mathbb{R}^d$. Then Problem (4.4) becomes

$$\inf\{\langle \mathbf{c}, \mathbf{x} \rangle : A\mathbf{x} \preceq_K \mathbf{b}\}. \quad (4.12)$$

1686 For a fixed cone K , problems of the form (4.12) with are called *conic optimization problems over the cone*
 1687 K . As we pick different data $A, \mathbf{b}, \mathbf{c}$, we get different instances of a conic optimization problem over the
 1688 cone K . A special case is when $K = \mathbb{R}_+^m$, which is known as *linear programming or linear optimization* – see
 1689 Example 4.24 – 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

$$\begin{aligned} \inf_{\mathbf{x} \in \mathbb{R}^d} \langle \mathbf{c}, \mathbf{x} \rangle + \langle \mathbf{y}, A\mathbf{x} - \mathbf{b} \rangle &= \inf_{\mathbf{x} \in \mathbb{R}^d} \langle \mathbf{c}, \mathbf{x} \rangle + \langle \mathbf{y}, A\mathbf{x} \rangle - \langle \mathbf{y}, \mathbf{b} \rangle \\ &= \inf_{\mathbf{x} \in \mathbb{R}^d} \langle \mathbf{c}, \mathbf{x} \rangle + \langle A^T \mathbf{y}, \mathbf{x} \rangle - \langle \mathbf{y}, \mathbf{b} \rangle \\ &= \inf_{\mathbf{x} \in \mathbb{R}^d} \langle \mathbf{c} + A^T \mathbf{y}, \mathbf{x} \rangle - \langle \mathbf{y}, \mathbf{b} \rangle. \end{aligned}$$

1690 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
 1691 $-\langle \mathbf{b}, \mathbf{y} \rangle$. Therefore, for (4.12), the dual function is given by

$$\mathcal{L}(\mathbf{y}) = \begin{cases} -\infty & \text{if } \mathbf{c} + A^T \mathbf{y} \neq \mathbf{0} \\ -\langle \mathbf{b}, \mathbf{y} \rangle & \text{if } \mathbf{c} + A^T \mathbf{y} = \mathbf{0} \end{cases} \quad (4.13)$$

Therefore,

$$\sup_{\mathbf{y} \in K^*} \mathcal{L}(\mathbf{y}) = \sup\{-\langle \mathbf{b}, \mathbf{y} \rangle : A^T \mathbf{y} = -\mathbf{c}, \mathbf{y} \in K^*\} = -\inf\{\langle \mathbf{b}, \mathbf{y} \rangle : A^T \mathbf{y} = -\mathbf{c}, \mathbf{y} \in K^*\}.$$

1692 To remove the slightly annoying minus sign in front of \mathbf{c} above, it is more standard to write (4.12) as
 1693 $-\sup\{\langle -\mathbf{c}, \mathbf{x} \rangle : A\mathbf{x} \preceq_K \mathbf{b}\}$, and then replace $-\mathbf{c}$ with \mathbf{c} throughout the above derivation. Thus, the
 1694 standard primal dual pairs for conic optimization problems are

$$\sup\{\langle \mathbf{c}, \mathbf{x} \rangle : A\mathbf{x} \preceq_K \mathbf{b}\} \leq \inf\{\langle \mathbf{b}, \mathbf{y} \rangle : A^T \mathbf{y} = \mathbf{c}, \mathbf{y} \in K^*\}. \quad (4.14)$$

1695 *Linear Programming/Optimization.* Specializing to the linear programming case with $K = \mathbb{R}_+^m$ and observing
 1696 that $K^* = K = \mathbb{R}_+^m$ (see Problem 2 from “HW for Week III”), we obtain the primal dual pair

$$\sup\{\langle \mathbf{c}, \mathbf{x} \rangle : A\mathbf{x} \leq \mathbf{b}\} \leq \inf\{\langle \mathbf{b}, \mathbf{y} \rangle : A^T \mathbf{y} = \mathbf{c}, \mathbf{y} \geq \mathbf{0}\}. \quad (4.15)$$

1697 *Semidefinite Programming/Optimization.* Another special case is that of semidefinite optimization. This is
 1698 the situation when $m = n^2$ and K is the cone of positive semidefinite matrices. $G : \mathbb{R}^d \rightarrow \mathbb{R}^{n^2}$ is an affine
 1699 map from \mathbb{R}^d to the space of $n \times n$ matrices. To avoid dealing with asymmetric matrices, G is always assumed
 1700 to be of the form $G(\mathbf{x}) = \mathbf{x}_1 A_1 + \dots + \mathbf{x}_d A_d - A_0$, where A_0, A_1, \dots, A_d are $n \times n$ symmetric matrices⁴. If
 1701 one works through the algebra in this case and uses the fact that the positive semidefinite cone is *self-dual*,
 1702 i.e., $K = K^*$, (4.14) becomes

$$\sup\{\langle \mathbf{c}, \mathbf{x} \rangle : \mathbf{x}_1 A_1 + \dots + \mathbf{x}_d A_d - A_0 \text{ is a PSD matrix}\} \leq \inf\{\langle A_0, Y \rangle : \langle A_i, Y \rangle = \mathbf{c}_i, Y \text{ is a PSD matrix}\},$$

1703 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 4.24, where $K = \mathbb{R}_+^m$, $f, g_1, \dots, g_m : \mathbb{R}^d \rightarrow \mathbb{R}$ are convex functions, and $G : \mathbb{R}^d \rightarrow \mathbb{R}^m$ was defined as $G(\mathbf{x}) = (g_1(\mathbf{x}), \dots, g_m(\mathbf{x}))$, giving the explicit problem

$$\inf\{f(\mathbf{x}) : g_1(\mathbf{x}) \leq 0, \dots, g_m(\mathbf{x}) \leq 0\}.$$

In this case, since $K^* = K = \mathbb{R}_+^m$ (see Problem 2 from “HW for Week III”), the dual problem is

$$\sup_{\mathbf{y} \in K^*} \mathcal{L}(\mathbf{y}) = \sup_{\mathbf{y} \geq \mathbf{0}} \inf_{\mathbf{x} \in \mathbb{R}^d} \{f(\mathbf{x}) + \mathbf{y}_1 g_1(\mathbf{x}) + \dots + \mathbf{y}_m g_m(\mathbf{x})\}.$$

⁴Dealing with asymmetric matrices is not hard, but involves little details that can be overlooked for this exposition, and don't provide any great insight.

1704 **A closer look at linear programming duality.** Consider the following linear program:

$$\begin{aligned}
 \sup \quad & 2x_1 - 1.5x_2 \\
 & x_1 + x_2 \leq 1 \\
 & x_1 - x_2 \leq 1 \\
 & -x_1 + x_2 \leq 1 \\
 & -x_1 - x_2 \leq 1
 \end{aligned} \tag{4.16}$$

1705 To solve this problem, let us make some simple observations. If we multiply the first inequality by 0.5,
 1706 the second inequality by 3.5, the third by 1.75 and the fourth by 0.25 and add all these scaled inequalities,
 1707 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
 1708 program must also satisfy this new inequality. This shows that our supremum is *at most* 6. Now if we choose
 1709 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
 1710 gives a better bound of $2 \leq 6$ on the optimal solution value. Now, consider the point $x_1 = 1, x_2 = 0$: this
 1711 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
 1712 $x_1 = 1, x_2 = 0$ is actually the optimal solution to the above linear program! Thus, we have provided the
 1713 optimal solution, and a quick certificate of its optimality. If you think about how we were deriving the upper
 1714 bounds of 6 and 2, we were looking for nonnegative multipliers y_1, y_2, y_3, y_4 such that the corresponding
 1715 combination of the inequalities gives us $2x_1 - 1.5x_2$ on the left hand side, and the upper bound was simply
 1716 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
 1717 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
 1718 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
 1719 $y_1 - y_2 + y_3 - y_4 = -1.5$, and all y_i 's are nonnegative. But this is exactly the dual problem in (4.15). We
 1720 hope this gives the reader a more “hands-on” perspective on the Lagrangian dual of a linear program.

1721 4.3.4 Strong duality: sufficient conditions and complementary slackness

1722 In the above example of the linear program in (4.16), it turned out that we could find a primal feasible
 1723 solution and a dual feasible solution that have the same value, which shows that we have strong duality, and
 1724 certifies the optimality of the two solutions. We will see below that this always happens for linear programs.
 1725 For general conic optimization problems, or a convex optimization problem with generalized inequalities,
 1726 this does not always hold and one may not even have zero duality gap. We now supply two conditions under
 1727 which strong duality is obtained. Linear programming strong duality will be a special case of the second
 1728 condition.

1729 **Slater’s condition for strong duality.** The following is perhaps the most well-known sufficient condition
 1730 in convex optimization that guarantees strong duality.

1731 **Theorem 4.31.** [Slater’s condition] Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex, let $K \subseteq \mathbb{R}^m$ be a closed, convex, pointed
 1732 cone, and let $G : \mathbb{R}^d \rightarrow \mathbb{R}^m$ be a K -convex mapping. Let $\mathcal{L} : \mathbb{R}^m \rightarrow \mathbb{R}$ be as defined in (4.9). If there exists
 1733 $\bar{\mathbf{x}}$ such that $-G(\bar{\mathbf{x}}) \in \text{int}(K)$ and $\inf\{f(\mathbf{x}) : G(\mathbf{x}) \preceq_K \mathbf{0}\}$ is a finite value, then there exists $\mathbf{y}^* \in K^*$ such
 1734 that $\sup_{\mathbf{y} \in K^*} \mathcal{L}(\mathbf{y}) = \mathcal{L}(\mathbf{y}^*) = \inf\{f(\mathbf{x}) : G(\mathbf{x}) \preceq_K \mathbf{0}\}$, i.e., strong duality holds.

1735 Before we begin the proof, we need to establish a slight variant of the separating hyperplane theorem,
 1736 that does not make any closedness or compactness assumptions.

1737 **Proposition 4.32.** Let $A, B \subseteq \mathbb{R}^d$ be convex sets (not necessarily closed) such that $A \cap B = \emptyset$. Then there
 1738 exist $\mathbf{a} \in \mathbb{R}^d \setminus \{\mathbf{0}\}$ and $\delta \in \mathbb{R}$ such that $\langle \mathbf{a}, \mathbf{x} \rangle \geq \delta$ for all $\mathbf{x} \in A, \mathbf{y} \in B$.

1739 *Proof.* Left as an exercise. □

1740 *Proof of Theorem 4.31.* Let $\mu_0 := \inf\{f(\mathbf{x}) : G(\mathbf{x}) \preceq_K \mathbf{0}\} \in \mathbb{R}$. Define the sets

$$\begin{aligned}
 A &:= \{(\mathbf{z}, r) \in \mathbb{R}^m \times \mathbb{R} : \exists \mathbf{x} \in \mathbb{R}^d \text{ such that } f(\mathbf{x}) \leq r, G(\mathbf{x}) \preceq_K \mathbf{z}\}, \\
 B &:= \{(\mathbf{z}, r) \in \mathbb{R}^m \times \mathbb{R} : r < \mu_0, \mathbf{z} \preceq_K \mathbf{0}\}.
 \end{aligned}$$

1741 It is not hard to verify that A, B are convex. Moreover, it is also not hard to verify that $A \cap B = \emptyset$. By
 1742 Proposition 4.32, there exists $\mathbf{a} \in \mathbb{R}^m, \gamma \in \mathbb{R}$ such that

$$\langle \mathbf{a}, \mathbf{z}_1 \rangle + \gamma r_1 \geq \langle \mathbf{a}, \mathbf{z}_2 \rangle + \gamma r_2 \quad (4.17)$$

1743 for all $(\mathbf{z}_1, r_1) \in A$ and $(\mathbf{z}_2, r_2) \in B$.

1744 **Claim 3.** $\mathbf{a} \in K^*$ and $\gamma \geq 0$.

1745 *Proof of Claim.* Suppose to the contrary that $\mathbf{a} \notin K^*$. Then $\mathbf{a} \notin -K^\circ = (-K)^\circ$. Thus, there exists $\bar{\mathbf{z}} \in -K$,
 1746 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
 1747 theorem), $r_1 = f(\bar{\mathbf{x}})$, $r_2 = \mu_0 - 1$ and $\mathbf{z}_2 = \lambda \bar{\mathbf{z}}$ for all $\lambda \geq 0$. But since $\langle \mathbf{a}, \bar{\mathbf{z}} \rangle > 0$, the inequality (4.17) would
 1748 be violated for large enough λ . Thus, we must have $\mathbf{a} \in K^*$.

1749 Similarly, (4.17) holds with $\mathbf{z}_1 = G(\bar{\mathbf{x}})$, $r_1 = f(\bar{\mathbf{x}})$, $\mathbf{z}_2 = \mathbf{0}$ and all $r_2 < \mu_0$. If $\gamma < 0$, then letting $r_2 \rightarrow -\infty$
 1750 would violate (4.17). \square

We now show that, in fact, $\gamma > 0$ because of the existence of $\bar{\mathbf{x}}$ assumed in the hypothesis of the theorem. Substitute $\mathbf{z}_1 = G(\bar{\mathbf{x}})$, $r_1 = f(\bar{\mathbf{x}})$, $r_2 = \mu_0 - 1$ and $\mathbf{z}_2 = \mathbf{0}$ in (4.17). If $\gamma = 0$, then this relation becomes

$$\langle \mathbf{a}, G(\bar{\mathbf{x}}) \rangle \geq 0.$$

1751 However, $-G(\bar{\mathbf{x}}) \in \text{int}(K)$ and $\mathbf{a} \neq \mathbf{0}$ and therefore, $\langle \mathbf{a}, G(\bar{\mathbf{x}}) \rangle < 0$ (see Problem 3 from “HW for Week II”).
 1752 By Claim 3, $\gamma > 0$.

Let $\mathbf{y}^* := \frac{\mathbf{a}}{\gamma}$; by Claim 3, $\mathbf{y}^* \in K^*$. We will now show that for every $\epsilon > 0$, $\mathcal{L}(\mathbf{y}^*) \geq \mu_0 - \epsilon$. This will
 establish the result because this means $\mathcal{L}(\mathbf{y}^*) \geq \mu_0$ and since $\mathcal{L}(\mathbf{y}) \leq \mu_0$ for all $\mathbf{y} \in K^*$ by Proposition 4.27,
 we must have $\sup_{\mathbf{y} \in K^*} \mathcal{L}(\mathbf{y}) = \mathcal{L}(\mathbf{y}^*) = \mu_0$. 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 $\mathbf{z}_2 = \mathbf{0}$ 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 γ , we obtain

$$\langle \mathbf{y}^*, G(\mathbf{x}) \rangle + f(\mathbf{x}) \geq \mu_0 - \epsilon.$$

1753 This implies that $\mathcal{L}(\mathbf{y}^*) = \inf_{\mathbf{x} \in \mathbb{R}^d} \langle \mathbf{y}^*, G(\mathbf{x}) \rangle + f(\mathbf{x}) \geq \mu_0 - \epsilon$. \square

1754 **Closed cone condition for strong duality in conic optimization.** Slater’s condition applied to conic
 1755 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
 1756 very useful strong duality condition uses topological properties of the dual cone K^* .

1757 **Theorem 4.33.** [Closed cone condition] Consider the conic optimization primal dual pair (4.14). Suppose
 1758 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 $\mathbf{y} \in K^*$ such
 1759 that $A^T \mathbf{y} = \mathbf{c}$. Then we have zero duality gap. If the optimal dual value is finite, then strong duality holds
 1760 in (4.14).

1761 *Proof.* Since the dual is feasible, its optimal value is either $-\infty$ or finite. By weak duality (Proposition 4.27),
 1762 in the first case we must have zero duality gap and the primal is infeasible. So we consider the case when the
 1763 optimal value of the dual is finite, say $\mu_0 \in \mathbb{R}$. Let us label the set $S := \{(A^T \mathbf{y}, \langle \mathbf{b}, \mathbf{y} \rangle) : \mathbf{y} \in K^*\} \subseteq \mathbb{R}^d \times \mathbb{R}$.
 1764 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
 1765 $\{r \in \mathbb{R} : (\mathbf{c}, r) \in S\}$ is closed because it is topologically the same as $S \cap (\mathbf{c} \times \mathbb{R})$. Therefore the infimum in
 1766 $\inf\{r \in \mathbb{R} : (\mathbf{c}, r) \in S\}$ is over a closed subset of the real line. Hence, $(\mathbf{c}, \mu_0) \in S$ and so there exists $\mathbf{y}^* \in K^*$
 1767 such that $A^T \mathbf{y}^* = \mathbf{c}$ and $\langle \mathbf{b}, \mathbf{y}^* \rangle = \mu_0$.

1768 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
 1769 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 $\mathbf{y} \in K^*$, and $\langle \mathbf{a}, \mathbf{c} \rangle + \gamma(\mu_0 - \epsilon) >$
 1770 δ . By Problem 8 from “HW for Week IX”, we may assume $\delta = 0$. Therefore, we have

$$\langle \mathbf{a}, A^T \mathbf{y} \rangle + \gamma \cdot \langle \mathbf{b}, \mathbf{y} \rangle \leq 0 \text{ for all } \mathbf{y} \in K^*, \quad (4.18)$$

$$\langle \mathbf{a}, \mathbf{c} \rangle + \gamma(\mu_0 - \epsilon) > 0 \quad (4.19)$$

1771 Substituting \mathbf{y}^* 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
 1772 implies that $\gamma < 0$ since $\epsilon > 0$. Now (4.18) can be rewritten as $\langle \mathbf{A}\mathbf{a} + \gamma \mathbf{b}, \mathbf{y} \rangle \leq 0$ for all $\mathbf{y} \in K^*$ and (4.19) can
 1773 be rewritten as $\langle \mathbf{a}, \mathbf{c} \rangle > -\gamma(\mu_0 - \epsilon)$. Dividing through both these relations by $-\gamma > 0$, and setting $\mathbf{x} = \frac{\mathbf{a}}{-\gamma}$,
 1774 we obtain that $\langle \mathbf{A}\mathbf{x} - \mathbf{b}, \mathbf{y} \rangle \leq 0$ for all $\mathbf{y} \in K^*$ implying that $\mathbf{A}\mathbf{x} \preceq_K \mathbf{b}$, and $\langle \mathbf{x}, \mathbf{c} \rangle > \mu_0 - \epsilon$. Thus, we have
 1775 a feasible solution \mathbf{x} for the primal with value at least $\mu_0 - \epsilon$. Since $\epsilon > 0$ was chosen arbitrarily, this shows
 1776 that for every $\epsilon > 0$, the primal has optimal value better than $\mu_0 - \epsilon$. Therefore, the primal value must be
 1777 μ_0 and we have zero duality gap. The existence of \mathbf{y}^* shows that we have strong duality. \square

1778 *Linear Programming strong duality.* The closed cone condition for strong duality implies that linear programs
 1779 always enjoy strong duality when either the primal or the dual (or both) are feasible. This is because the
 1780 cone $K = \mathbb{R}_+^m$ is a polyhedral cone and also self-dual, i.e., $K^* = K = \mathbb{R}_+^m$. Since linear transformations of
 1781 polyhedral cones are polyhedral (see part 5. of Problem 1 in “HW for Week VI”), and hence closed, linear
 1782 programs always satisfy the condition in Theorem 4.33. One therefore has the following table for the possible
 1783 outcomes in the primal dual linear programming pair.

Dual Primal	Infeasible	Finite	Unbounded
Infeasible	Possible	Impossible	Possible
Finite	Impossible	Possible, Zero duality gap	Impossible
Unbounded	Possible	Impossible	Impossible

1785 An alternate proof of zero duality gap for linear programming follows from our results on polyhedral
 1786 theory. We outline it here to illustrate that linear programming duality can be approached in different
 1787 ways (although ultimately both proofs go back to the separating hyperplane theorem – Theorem 2.20). We
 1788 consider two cases:

1789 *Primal is infeasible.* In this case, we will show that if the dual is feasible, then the dual must be
 1790 unbounded. Since the primal is infeasible, the polyhedron $\mathbf{A}\mathbf{x} \leq \mathbf{b}$ is empty. By Theorem 2.88, there exists
 1791 $\hat{\mathbf{y}} \geq \mathbf{0}$ such that $\mathbf{A}^T \hat{\mathbf{y}} = \mathbf{0}$ and $\langle \mathbf{b}, \hat{\mathbf{y}} \rangle = -1$. Since the dual is feasible, consider any $\bar{\mathbf{y}} \geq \mathbf{0}$ such that $\mathbf{A}^T \bar{\mathbf{y}} = \mathbf{c}$.
 1792 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
 1793 be made to go to $-\infty$ because $\langle \mathbf{b}, \hat{\mathbf{y}} \rangle = -1$.

1794 *Primal is feasible.* If the primal is unbounded, then by weak duality, the dual must be infeasible. So let
 1795 us consider the case that the primal has a finite value μ_0 . This means that the inequality $\langle \mathbf{c}, \mathbf{x} \rangle \leq \mu_0$ is a
 1796 valid inequality for the polyhedron $\mathbf{A}\mathbf{x} \leq \mathbf{b}$. By Theorem 2.85, there exists $\hat{\mathbf{y}} \geq \mathbf{0}$ such that $\mathbf{A}^T \hat{\mathbf{y}} = \mathbf{c}$ and
 1797 $\langle \mathbf{b}, \hat{\mathbf{y}} \rangle \leq \mu_0$. Therefore the dual has a solution $\hat{\mathbf{y}}$ whose objective value is equal to the primal value μ_0 . This
 1798 guarantees strong duality.

1799 **Complementary slackness.** Complementary slackness is a useful necessary condition when we have
 1800 primal and dual optimal solutions with zero duality gap.

1801 **Theorem 4.34.** Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex, let $K \subseteq \mathbb{R}^m$ be a closed, convex, pointed cone, and let
 1802 $G : \mathbb{R}^d \rightarrow \mathbb{R}^m$ be a K -convex mapping. Let $\mathcal{L} : \mathbb{R}^m \rightarrow \mathbb{R}$ be as defined in (4.9). Let \mathbf{x}^* be such that
 1803 $G(\mathbf{x}^*) \preceq_K \mathbf{0}$ and $\mathbf{y}^* \in K^*$ such that $f(\mathbf{x}^*) = \mathcal{L}(\mathbf{y}^*)$. Then $\langle \mathbf{y}^*, G(\mathbf{x}^*) \rangle = 0$.

Proof. We simply observe that since $G(\mathbf{x}^*) \preceq_K \mathbf{0}$ and $\mathbf{y}^* \in K^*$, we must have $\langle \mathbf{y}^*, G(\mathbf{x}^*) \rangle \leq 0$. Therefore,

$$f(\mathbf{x}^*) \geq f(\mathbf{x}^*) + \langle \mathbf{y}^*, G(\mathbf{x}^*) \rangle \geq \inf_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + \langle \mathbf{y}^*, G(\mathbf{x}) \rangle = \mathcal{L}(\mathbf{y}^*).$$

1804 Since $f(\mathbf{x}^*) = \mathcal{L}(\mathbf{y}^*)$ by assumption, equality must hold throughout above giving us $\langle \mathbf{y}^*, G(\mathbf{x}^*) \rangle = 0$. \square

1805 **4.3.5 Saddle point interpretation of the Lagrangian dual**

1806 Let us go back to the original problem (4.4) and revisit the dual function $\mathcal{L}(\mathbf{y})$. Define the function

$$\hat{\mathcal{L}}(\mathbf{x}, \mathbf{y}) := f(\mathbf{x}) + \langle \mathbf{y}, G(\mathbf{x}) \rangle \quad (4.20)$$

1807 which is often called the *Lagrangian function* associated with (4.4). A characterization of a pair of optimal
1808 solutions to (4.4) and (4.10) can be obtained using saddle points of the Lagrangian function.

1809 **Theorem 4.35.** Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex, let $K \subseteq \mathbb{R}^m$ be a closed, convex, pointed cone, and let
1810 $G : \mathbb{R}^d \rightarrow \mathbb{R}^m$ be a K -convex mapping. Let $\mathcal{L} : \mathbb{R}^m \rightarrow \mathbb{R}$ be as defined in (4.9) and $\hat{\mathcal{L}} : \mathbb{R}^d \times \mathbb{R}^m \rightarrow \mathbb{R}$ be as
1811 defined in (4.20). Let \mathbf{x}^* be such that $G(\mathbf{x}^*) \preceq_K \mathbf{0}$ and $\mathbf{y}^* \in K^*$. Then the following are equivalent.

- 1812 1. $\mathcal{L}(\mathbf{y}^*) = f(\mathbf{x}^*)$.
1813 2. $\hat{\mathcal{L}}(\mathbf{x}^*, \hat{\mathbf{y}}) \leq \hat{\mathcal{L}}(\mathbf{x}^*, \mathbf{y}^*) \leq \hat{\mathcal{L}}(\hat{\mathbf{x}}, \mathbf{y}^*)$, for all $\hat{\mathbf{x}} \in \mathbb{R}^d$ and $\hat{\mathbf{y}} \in K^*$.

Proof. 1. \implies 2. Consider any $\hat{\mathbf{x}} \in \mathbb{R}^d$ and $\hat{\mathbf{y}} \in K^*$. We now derive the following chain of inequalities:

$$\begin{aligned} \hat{\mathcal{L}}(\mathbf{x}^*, \hat{\mathbf{y}}) &= f(\mathbf{x}^*) + \langle \hat{\mathbf{y}}, G(\mathbf{x}^*) \rangle \\ &\leq f(\mathbf{x}^*) && \text{since } \langle \hat{\mathbf{y}}, G(\mathbf{x}^*) \rangle \leq 0 \text{ because } \hat{\mathbf{y}} \in K^*, G(\mathbf{x}^*) \preceq_K \mathbf{0} \\ &= f(\mathbf{x}^*) + \langle \mathbf{y}^*, G(\mathbf{x}^*) \rangle = \hat{\mathcal{L}}(\mathbf{x}^*, \mathbf{y}^*) && \text{since } \langle \mathbf{y}^*, G(\mathbf{x}^*) \rangle = 0 \text{ by Theorem 4.34} \\ &= \mathcal{L}(\mathbf{y}^*) && \text{since } \mathcal{L}(\mathbf{y}^*) = f(\mathbf{x}^*) \\ &= \inf_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + \langle \mathbf{y}^*, G(\mathbf{x}) \rangle \\ &\leq f(\hat{\mathbf{x}}) + \langle \mathbf{y}^*, G(\hat{\mathbf{x}}) \rangle \\ &= \hat{\mathcal{L}}(\hat{\mathbf{x}}, \mathbf{y}^*) \end{aligned}$$

2. \implies 1. Since $\hat{\mathcal{L}}(\mathbf{x}^*, \hat{\mathbf{y}}) \leq \hat{\mathcal{L}}(\mathbf{x}^*, \mathbf{y}^*)$ for all $\hat{\mathbf{y}} \in K^*$, we have that

$$\hat{\mathcal{L}}(\mathbf{x}^*, \mathbf{y}^*) = \sup_{\mathbf{y} \in K^*} \hat{\mathcal{L}}(\mathbf{x}^*, \hat{\mathbf{y}}) = \sup_{\mathbf{y} \in K^*} f(\mathbf{x}^*) + \langle \mathbf{y}, G(\mathbf{x}^*) \rangle = f(\mathbf{x}^*),$$

where the last equality follows from the fact that $\langle \mathbf{y}, G(\mathbf{x}^*) \rangle \leq 0$ for all $\mathbf{y} \in K^*$. So the supremum is achieved
for $\mathbf{y} = \mathbf{0}$. On the other hand, since $\hat{\mathcal{L}}(\mathbf{x}^*, \mathbf{y}^*) \leq \hat{\mathcal{L}}(\hat{\mathbf{x}}, \mathbf{y}^*)$ for all $\hat{\mathbf{x}} \in \mathbb{R}^d$, we have that

$$\hat{\mathcal{L}}(\mathbf{x}^*, \mathbf{y}^*) = \inf_{\mathbf{x} \in \mathbb{R}^d} \hat{\mathcal{L}}(\mathbf{x}, \mathbf{y}^*) = \inf_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + \langle \mathbf{y}^*, G(\mathbf{x}) \rangle = \mathcal{L}(\mathbf{y}^*).$$

1814 Thus, we obtain that $f(\mathbf{x}^*) = \hat{\mathcal{L}}(\mathbf{x}^*, \mathbf{y}^*) = \mathcal{L}(\mathbf{y}^*)$. □

1815 Theorem 4.35 says that \mathbf{x}^* and \mathbf{y}^* are solutions for the primal problem (4.4) and dual problem (4.10)
1816 respectively, if and only if $(\mathbf{x}^*, \mathbf{y}^*)$ is a saddle point for the function $\hat{\mathcal{L}}(\mathbf{x}, \mathbf{y})$ of the type that \mathbf{x}^* is the
1817 minimizer when \mathbf{y} is fixed at \mathbf{y}^* and \mathbf{y}^* is the maximizer when \mathbf{x} is fixed at \mathbf{x}^* . This can be used to
1818 directly solve (4.4) and (4.10) simultaneously by searching for such saddle-points of the function $\hat{\mathcal{L}}(\mathbf{x}, \mathbf{y})$.
1819 This approach can be useful, if one has analytical forms for f and G (with sufficient differentiable properties)
1820 so that finding saddle-points is a reasonable option.

1821 **4.4 Cutting plane schemes**

1822 We now go back to the most general convex optimization problem (4.1). As before, we make no assumptions
1823 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
1824 returns an element from the subdifferential $\partial f(\mathbf{x})$, and if $\mathbf{x} \notin C$ then it returns a separating hyperplane.

1825 The subgradient algorithm from Section 4.1 can be used to solve (4.1) if one has access to the projection
1826 operator $\text{Proj}_C(\mathbf{x})$, which is stronger than a separation oracle. *Cutting plane schemes* are a class of algo-
1827 rithms that work with just a separation oracle. Moreover, the number of oracle calls is quite different from
1828 the number of oracle calls made by the subgradient algorithm: on the one hand, they typically exhibit a

1829 logarithmic dependence of $\ln(\frac{MR}{\epsilon})$ on the initial data M, R and error guarantee ϵ as opposed to the quadratic
1830 dependence $\frac{M^2R^2}{\epsilon^2}$ of the subgradient algorithm; on the other other, cutting plane schemes have a polynomial
1831 dependence on the dimension d of the problem (typically of the order of d^2), and such a dependence does
1832 not exist for the subgradient algorithm – see Remark 4.16.

1833 We will present the algorithm and the analysis for the situation when C is compact and full-dimensional.
1834 Hence the minimizer \mathbf{x}^* exists for (4.1) since f is convex, and therefore, continuous by Theorem 3.21. There
1835 are ways to get around this assumption, but we will ignore this complication in this write-up.

1836 General cutting plane scheme

- 1837 1. Choose any $E_0 \supseteq C$.
- 1838 2. For $i = 0, 1, 2, \dots$, do
 - 1839 (a) Choose $\mathbf{x}^i \in E_i$.
 - 1840 (b) Call the separation oracle for C with \mathbf{x}^i as input.
 1841 *Case 1:* $\mathbf{x}^i \in C$. Call the first order oracle for f to get some $\mathbf{s}^i \in \partial f(\mathbf{x}^i)$.
 1842 *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 .
 - 1843 (c) Set $E_{i+1} \supseteq E_i \cap \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{s}^i, \mathbf{x} \rangle \leq \langle \mathbf{s}^i, \mathbf{x}^i \rangle\}$.

1844 The points $\mathbf{x}^0, \mathbf{x}^1, \dots$ will be called the *iterates* of the Cutting Plane scheme.

1845 **Remark 4.36.** The above general scheme actually defines a family of algorithms. We have two choices to
1846 make to get a particular algorithm out of this scheme. First, there must be a strategy/procedure to choose
1847 $\mathbf{x}^i \in E_i$ in step 2(a) in every iteration. Second, there should be a strategy to define E_{i+1} as a superset of
1848 $E_i \cap \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{s}^i, \mathbf{x} \rangle \leq \langle \mathbf{s}^i, \mathbf{x}^i \rangle\}$ in step 2(c) of the scheme. Depending on what these two strategies are, we
1849 get different variants of the general cutting plane scheme. We will look at two variants below: the *center of*
1850 *gravity method* and the *ellipsoid method*.

1851 Technically, we also have to make a choice for E_0 in Step 1, but this is usually given as part of the input
1852 to the problem: E_0 is usually a large ball or polytope containing C that is provided or known at the start.

1853 We now start our analysis of cutting plane schemes. We introduce a useful notation to denote the
1854 polyhedron defined by the halfspaces obtained during the iterations of the cutting plane scheme.

Definition 4.37. Let $\mathbf{z}^1, \dots, \mathbf{z}^k \subseteq \mathbb{R}^d$ and let $\mathbf{s}^1, \dots, \mathbf{s}^k$ 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

$$G(\mathbf{z}^1, \dots, \mathbf{z}^k) := \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{s}^i, \mathbf{x} \rangle \leq \langle \mathbf{s}^i, \mathbf{z}^i \rangle \quad i = 1, \dots, k\}.$$

1855 This polyhedron will be referred to as the *gradient polyhedron* of $\mathbf{z}^1, \dots, \mathbf{z}^k$. The name is a bit of a misnomer,
1856 because we are considering general f , so we may have no gradients, and also some of the halfspaces could
1857 correspond to separating hyperplanes which have nothing to do with gradients. Even so we stick with this
1858 terminology.

Definition 4.38. Let $\mathbf{x}^0, \mathbf{x}^1, \dots$ be the iterates of a cutting plane scheme. For any iteration $t \geq 0$, we define
 $h(t) := |C \cap \{\mathbf{x}^0, \dots, \mathbf{x}^t\}|$, i.e., $h(t)$ is the number of feasible iterates until iteration t . We also define

$$S_t = C \cap G(\mathbf{x}^0, \dots, \mathbf{x}^t).$$

1859 As we shall see below, the volume of S_t will be central in measuring our progress towards the optimal
1860 solution. We first observe in the next lemma that S_t can be described as the intersection of C and the
1861 gradient polyhedron of only the feasible iterates.

1862 **Lemma 4.39.** Let $\mathbf{x}^0, \mathbf{x}^1, \dots$ be the iterates of a cutting plane scheme. Let $t \geq 0$ be any natural number
1863 and let the feasible iterates be denoted by $\{\mathbf{x}^{i_1}, \dots, \mathbf{x}^{i_{h(t)}}\} = C \cap \{\mathbf{x}^0, \dots, \mathbf{x}^t\}$, with $0 \leq i_1 \leq i_2 \leq \dots \leq i_{h(t)}$.
1864 Then $S_t = C \cap G(\mathbf{x}^{i_1}, \dots, \mathbf{x}^{i_{h(t)}})$.

Proof. Let $X_t = \{\mathbf{x}^0, \dots, \mathbf{x}^t\}$. We derive the following relations.

$$\begin{aligned} S_t &= C \cap G(\mathbf{x}^0, \dots, \mathbf{x}^t) \\ &= C \cap G(X_t \setminus \{\mathbf{x}^{i_1}, \dots, \mathbf{x}^{i_{h(t)}}\}) \cap G(\mathbf{x}^{i_1}, \dots, \mathbf{x}^{i_{h(t)}}) \\ &= C \cap G(\mathbf{x}^{i_1}, \dots, \mathbf{x}^{i_{h(t)}}), \end{aligned}$$

1865 where the last equality follows since $C \subseteq G(X_t \setminus \{\mathbf{x}^{i_1}, \dots, \mathbf{x}^{i_{h(t)}}\})$ because each $\mathbf{z} \in X_t \setminus \{\mathbf{x}^{i_1}, \dots, \mathbf{x}^{i_{h(t)}}\}$ is
 1866 infeasible, i.e., $\mathbf{z} \notin C$, and therefore, the corresponding vector \mathbf{s} is a separating hyperplane for \mathbf{z} and C , i.e.,
 1867 $C \subseteq \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{s}, \mathbf{x} \rangle \leq \langle \mathbf{s}, \mathbf{z} \rangle\}$. \square

1868 Since our analysis will involve the volume of S_t , while our algorithm only works with the sets E_t , we need
 1869 to establish a definite relationship between these two sets.

1870 **Lemma 4.40.** Let $\mathbf{x}^0, \mathbf{x}^1, \dots$ be the iterates of a cutting plane scheme. Then $E_{t+1} \supseteq S_t$ for all $t \geq 0$.

1871 *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 $i = 0, \dots, t$. By putting all these
 1872 relationships together, we obtain that

$$E_{t+1} \supseteq E_0 \cap G(\mathbf{x}^0, \dots, \mathbf{x}^t) \supseteq C \cap G(\mathbf{x}^0, \dots, \mathbf{x}^t) = S_t, \quad (4.21)$$

1873 where the second containment follows from the assumption that $E_0 \supseteq C$. \square

1874 We now state our main structural result for the analysis of cutting plane schemes. We use $\text{dist}(\mathbf{x}, X)$ to
 1875 denote the distance of $\mathbf{x} \in \mathbb{R}^d$ from any subset $X \subseteq \mathbb{R}^d$, i.e., $\text{dist}(\mathbf{x}, X) := \inf_{\mathbf{y} \in X} \|\mathbf{x} - \mathbf{y}\|$.

Theorem 4.41. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a convex function and let C be a compact, convex set. Let \mathbf{x}^* be a
 minimizer for (4.1). Let $\mathbf{x}^0, \mathbf{x}^1, \dots$ be the iterates of any cutting plane scheme. For any $t \geq 0$, let the feasible
 iterates be denoted by $\{\mathbf{x}^{i_1}, \dots, \mathbf{x}^{i_{h(t)}}\} = C \cap \{\mathbf{x}^0, \dots, \mathbf{x}^t\}$, with $0 \leq i_1 \leq i_2 \leq \dots \leq i_{h(t)}$. Define

$$v_{\min}(t) := \min_{j=i_1, \dots, i_{h(t)}} \text{dist}(\mathbf{x}^*, H(\mathbf{s}^j, \langle \mathbf{s}^j, \mathbf{x}^j \rangle)),$$

1876 i.e., $v_{\min}(t)$ is the minimum distance of \mathbf{x}^* from the hyperplanes $\{\mathbf{x} : \langle \mathbf{s}^j, \mathbf{x} \rangle = \langle \mathbf{s}^j, \mathbf{x}^j \rangle\}$, $j = i_1, \dots, i_{h(t)}$. Let
 1877 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.

- 1878 1. For any $t \geq 0$, if $\text{vol}(E_{t+1}) < \text{vol}(C)$ then $h(t) > 0$, i.e., there is at least one feasible iterate.
- 1879 2. For any $t \geq 0$ such that $h(t) > 0$, $v_{\min}(t) \leq D \left(\frac{\text{vol}(S_t)}{\text{vol}(C)} \right)^{\frac{1}{d}} \leq D \left(\frac{\text{vol}(E_{t+1})}{\text{vol}(C)} \right)^{\frac{1}{d}}$.
- 1880 3. For any $t \geq 0$ such that $h(t) > 0$, $\min_{j=i_1, \dots, i_{h(t)}} f(\mathbf{x}^j) \leq f(\mathbf{x}^*) + M v_{\min}(t) \leq f(\mathbf{x}^*) + MD \left(\frac{\text{vol}(E_{t+1})}{\text{vol}(C)} \right)^{\frac{1}{d}}$,
 1881 where $M = L(B(\mathbf{x}^*, v_{\min}))$ is a Lipschitz constant for f over $B(\mathbf{x}^*, v_{\min})$ (see Theorem 3.21). This
 1882 provides a bound on the value of the best feasible point seen upto iteration t , in comparison to the
 1883 optimal value $f(\mathbf{x}^*)$.

1884 Theorem 4.41 shows that if we can ensure $\text{vol}(E_t) \rightarrow 0$ as $t \rightarrow \infty$, then we have a convergent algorithm.

1885 *Proof of Theorem 4.41.* 1. We prove the contrapositive. If $h(t) = 0$, then all iterates upto iteration t
 1886 are infeasible, i.e., $\mathbf{x}^i \notin C$ for all $i = 1, \dots, t$. This implies that all the vector \mathbf{s}^i are normal vectors
 1887 for separating hyperplanes. So $C \subseteq G(\mathbf{x}^0, \dots, \mathbf{x}^t)$. Since $C \subseteq E_0$, this implies that $C = E_0 \cap C \subseteq$
 1888 $E_0 \cap G(\mathbf{x}^0, \dots, \mathbf{x}^t) \subseteq E_{t+1}$, where the last containment follows from the first containment in (4.21).
 1889 Therefore, $\text{vol}(C) \leq \text{vol}(E_{t+1})$.

2. Let $\alpha = \frac{v_{\min}(t)}{D}$. Since D is the diameter of C , we must have $C \subseteq B(\mathbf{x}^*, D)$. Thus,

$$\alpha(C - \mathbf{x}^*) + \mathbf{x}^* \subseteq B(\mathbf{x}^*, \alpha D) = B(\mathbf{x}^*, v_{\min}(t)) \subseteq G(\mathbf{x}^{i_1}, \dots, \mathbf{x}^{i_{h(t)}}),$$

where the first equality follows from the definition of α and the final containment follows from the definition of $v_{\min}(t)$. Since $\mathbf{x}^* \in C$ and C is convex, we know that $\alpha(C - \mathbf{x}^*) + \mathbf{x}^* = \alpha C + (1 - \alpha)\mathbf{x}^* \subseteq C$. Therefore, $\alpha(C - \mathbf{x}^*) + \mathbf{x}^* = C \cap (\alpha(C - \mathbf{x}^*) + \mathbf{x}^*) \subseteq C \cap G(\mathbf{x}^{i_1}, \dots, \mathbf{x}^{i_{h(t)}}) = S_t$, where the last equality follows from Lemma 4.39. This implies that $\alpha^d \text{vol}(C) = \text{vol}(\alpha(C - \mathbf{x}^*)) \leq \text{vol}(S_t)$.

Rearranging and using the definition of α , we obtain that $v_{\min}(t) \leq D \left(\frac{\text{vol}(S_t)}{\text{vol}(C)} \right)^{\frac{1}{d}}$. By Lemma 4.40, $D \left(\frac{\text{vol}(S_t)}{\text{vol}(C)} \right)^{\frac{1}{d}} \leq D \left(\frac{\text{vol}(E_{t+1})}{\text{vol}(C)} \right)^{\frac{1}{d}}$.

3. It suffices to prove the first inequality; the second inequality follows from part 2. above. Let $i^{\min} \in \{i_1, i_2, \dots, i_{h(t)}\}$ be such that $v_{\min}(t) = \text{dist}(\mathbf{x}^*, H(\mathbf{s}^{i^{\min}}, \langle \mathbf{s}^{i^{\min}}, \mathbf{x}^{i^{\min}} \rangle))$. Denote by $H := H(\mathbf{s}^{i^{\min}}, \langle \mathbf{s}^{i^{\min}}, \mathbf{x}^{i^{\min}} \rangle)$ the hyperplane 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}^* . Using the Lipschitz constant M , we obtain that $f(\bar{\mathbf{x}}) \leq f(\mathbf{x}^*) + Mv_{\min}(t)$; see Figure 3 in Section 4.1. 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 $\bar{\mathbf{x}}, \mathbf{x}^{i^{\min}} \in H$ implying that $\langle \mathbf{s}^{i^{\min}}, \bar{\mathbf{x}} - \mathbf{x}^{i^{\min}} \rangle = 0$. Therefore, we obtain

$$\min_{j=i_1, \dots, i_{h(t)}} f(\mathbf{x}^j) \leq f(\mathbf{x}^{i^{\min}}) \leq f(\bar{\mathbf{x}}) \leq f(\mathbf{x}^*) + Mv_{\min}(t).$$

□

We now analyze two instantiations of the cutting plane scheme with concrete strategies to choose \mathbf{x}^i and E_{i+1} in each iteration i .

Center of Gravity Method. The first one is called the *center of gravity method*.

Definition 4.42. The *center of gravity* for any compact set $X \subseteq \mathbb{R}^d$ with non-zero volume is defined as

$$\frac{\int_X \mathbf{x} d\mathbf{x}}{\text{vol}(X)}.$$

An important property of the center gravity of compact, convex sets was established by Grünbaum [4].

Theorem 4.43. 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 $\bar{\mathbf{x}} \in H$,

$$\frac{1}{e} \leq \left(\frac{d}{d+1} \right)^d \leq \frac{\text{vol}(H^+ \cap C)}{\text{vol}(C)} \leq 1 - \left(\frac{d}{d+1} \right)^d \leq 1 - \frac{1}{e},$$

where H^+ is a halfspace with boundary H .

Theorem 4.43 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 4.43 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 4.44. In the center of gravity method, if $h(t) > 0$ for some iteration $t \geq 0$, then

$$\min_{j=i_1, \dots, i_{h(t)}} f(\mathbf{x}^j) \leq f(\mathbf{x}^*) + MD \left(1 - \frac{1}{e} \right)^{t/d} \left(\frac{\text{vol}(E_0)}{\text{vol}(C)} \right)^{1/d},$$

where D is the diameter of C and M is a Lipschitz constant for f over $B(\mathbf{x}^*, D)$.

In particular, if $E_0 = C$, then $\min_{j=i_1, \dots, i_{h(t)}} f(\mathbf{x}^j) \leq f(\mathbf{x}^*) + MD \left(1 - \frac{1}{e} \right)^{t/d}$.

1910 *Proof.* Follows from Theorem 4.41 part 3., and the fact that $B(\mathbf{x}^*, v_{\min}) \subseteq B(\mathbf{x}^*, D)$ implying that M is a
 1911 Lipschitz constant for f over $B(\mathbf{x}^*, v_{\min})$, and $\text{vol}(E_{t+1}) \leq (1 - \frac{1}{e})^t \text{vol}(E_0)$ by Theorem 4.43. \square

1912 By setting the error term $MD(1 - \frac{1}{e})^{t/d} (\frac{\text{vol}(E_0)}{\text{vol}(C)})^{1/d}$ less than equal to ϵ in Theorem 4.44, the following
 1913 is an immediate consequence.

Corollary 4.45. For any $\epsilon > 0$, after $O\left(d \ln\left(\frac{MD}{\epsilon}\right) + \ln\left(\frac{\text{vol}(E_0)}{\text{vol}(C)}\right)\right)$ iterations of the center of gravity method,

$$\min_{j=i_1, \dots, i_{h(t)}} f(\mathbf{x}^j) \leq f(\mathbf{x}^*) + \epsilon.$$

1914 In particular, if $E_0 = C$, then one needs $O(d \ln(\frac{MD}{\epsilon}))$ iterations.

1915 **Ellipsoid method.** The ellipsoid method is a cutting plane scheme where E_0 is assumed to be a large
 1916 ball with radius R around a known point \mathbf{x}^0 (typically $\mathbf{x}^0 = \mathbf{0}$) that is guaranteed to contain C . At
 1917 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
 1918 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
 1919 $\text{vol}(E_{i+1}) \leq (1 - \frac{1}{d^2+1})^{d/2} \text{vol}(E_i)$. The technical bulk of the analysis goes into showing that such an ellipsoid
 1920 E_{i+1} *always* exists.

Definition 4.46. 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 $Q \in \mathbb{R}^{d \times d}$ and a point $\mathbf{y} \in \mathbb{R}^d$, we define

$$E(Q, \mathbf{y}) := \{\mathbf{x} \in \mathbb{R}^d : (\mathbf{x} - \mathbf{y})^T Q^{-1} (\mathbf{x} - \mathbf{y}) \leq 1\}.$$

1921 The next proposition follows from unwrapping the definition. It shows that ellipsoids are simply the
 1922 image of the Euclidean unit norm ball under an invertible linear transformation.

1923 **Proposition 4.47.** Let $Q \in \mathbb{R}^{d \times d}$ be a positive definite matrix and let $Q^{-1} = X^T X$ for some invertible
 1924 matrix $X \in \mathbb{R}^{d \times d}$. Then $E(Q, \mathbf{y}) = \mathbf{y} + X^{-1}(B(\mathbf{0}, 1))$. Thus, $\text{vol}(E(Q, \mathbf{y})) = \det(X^{-1}) \text{vol}(B(\mathbf{0}, 1)) =$
 1925 $\sqrt{\det(Q)} \text{vol}(B(\mathbf{0}, 1))$.

1926 In the following, we will utilize the following relation for any $\mathbf{w}, \mathbf{z} \in \mathbb{R}^d$ and $A \in \mathbb{R}^{d \times d}$

$$(\mathbf{w} + \mathbf{z})^T A (\mathbf{w} + \mathbf{z}) = \mathbf{w}^T A \mathbf{w} + 2\mathbf{w}^T A \mathbf{z} + \mathbf{z}^T A \mathbf{z}. \quad (4.22)$$

Theorem 4.48. Let $Q \in \mathbb{R}^{d \times d}$ be a positive definite matrix and $\mathbf{y} \in \mathbb{R}^d$. Let $\mathbf{s} \in \mathbb{R}^d$ and let $E_+ = E(Q, \mathbf{y}) \cap H^-(\mathbf{s}, \langle \mathbf{s}, \mathbf{y} \rangle)$. Define

$$\mathbf{y}_+ = \mathbf{y} - \frac{1}{d+1} \cdot \frac{Q\mathbf{s}}{\sqrt{\mathbf{s}^T Q \mathbf{s}}}$$

$$Q_+ = \frac{d^2}{d^2-1} \left(Q - \frac{2}{d+1} \cdot \frac{Q\mathbf{s}\mathbf{s}^T Q}{\mathbf{s}^T Q \mathbf{s}} \right).$$

1927 Then $E_+ \subseteq E(Q_+, \mathbf{y}_+)$ and $\text{vol}(E(Q_+, \mathbf{y}_+)) \leq (1 - \frac{1}{(d+1)^2})^{d/2} \text{vol}(E(Q, \mathbf{y}))$.

1928 *Proof.* We first prove $E_+ \subseteq E(Q_+, \mathbf{y}_+)$. Consider any $\mathbf{x} \in E_+ = E(Q, \mathbf{y}) \cap H^-(\mathbf{s}, \langle \mathbf{s}, \mathbf{y} \rangle)$. To ease notational
 1929 burden, we denote $G = Q^{-1}$ and $G_+ = Q_+^{-1}$. A direct calculation shows that $G_+ = \frac{d^2-1}{d^2} \left(G + \frac{2}{d-1} \cdot \frac{\mathbf{s}\mathbf{s}^T}{\mathbf{s}^T Q \mathbf{s}} \right)$.
 1930 Note that \mathbf{x} satisfies

$$(\mathbf{x} - \mathbf{y})^T G (\mathbf{x} - \mathbf{y}) \leq 1 \quad (4.23)$$

$$\langle \mathbf{s}, \mathbf{x} - \mathbf{y} \rangle \leq 0 \quad (4.24)$$

We now verify that

$$\begin{aligned} (\mathbf{x} - \mathbf{y}_+)^T G_+ (\mathbf{x} - \mathbf{y}_+) &= (\mathbf{x} - \mathbf{y} + \frac{1}{d+1} \cdot \frac{Q\mathbf{s}}{\sqrt{\mathbf{s}^T Q \mathbf{s}}})^T G_+ (\mathbf{x} - \mathbf{y} + \frac{1}{d+1} \cdot \frac{Q\mathbf{s}}{\sqrt{\mathbf{s}^T Q \mathbf{s}}}) \\ &= (\mathbf{x} - \mathbf{y})^T G_+ (\mathbf{x} - \mathbf{y}) + \frac{2}{d+1} (\mathbf{x} - \mathbf{y})^T G_+ \left(\frac{Q\mathbf{s}}{\sqrt{\mathbf{s}^T Q \mathbf{s}}} \right) + \left(\frac{1}{d+1} \right)^2 \cdot \frac{\mathbf{s}^T Q^T G_+ Q \mathbf{s}}{\mathbf{s}^T Q \mathbf{s}}, \end{aligned}$$

1931 where we use (4.22). Let us analyze the three terms separately. The first term can be written in terms of
1932 G, \mathbf{s} , and \mathbf{y} :

$$\begin{aligned} (\mathbf{x} - \mathbf{y})^T G_+ (\mathbf{x} - \mathbf{y}) &= (\mathbf{x} - \mathbf{y})^T \left(\frac{d^2-1}{d^2} (G + \frac{2}{d-1} \cdot \frac{\mathbf{s}\mathbf{s}^T}{\mathbf{s}^T Q \mathbf{s}}) \right) (\mathbf{x} - \mathbf{y}) \\ &= \frac{d^2-1}{d^2} \left((\mathbf{x} - \mathbf{y})^T G (\mathbf{x} - \mathbf{y}) + \frac{2}{d-1} \frac{(\mathbf{s}^T (\mathbf{x} - \mathbf{y}))^2}{\mathbf{s}^T Q \mathbf{s}} \right) \end{aligned}$$

1933 The second term simplifies to

$$\begin{aligned} \frac{2}{d+1} (\mathbf{x} - \mathbf{y})^T G_+ \left(\frac{Q\mathbf{s}}{\sqrt{\mathbf{s}^T Q \mathbf{s}}} \right) &= \frac{2}{d+1} (\mathbf{x} - \mathbf{y})^T \left(\frac{d^2-1}{d^2} (G + \frac{2}{d-1} \cdot \frac{\mathbf{s}\mathbf{s}^T}{\mathbf{s}^T Q \mathbf{s}}) \right) \left(\frac{Q\mathbf{s}}{\sqrt{\mathbf{s}^T Q \mathbf{s}}} \right) \\ &= \frac{d^2-1}{d^2} \cdot \frac{2}{d+1} \left(\frac{\mathbf{s}^T (\mathbf{x} - \mathbf{y})}{\sqrt{\mathbf{s}^T Q \mathbf{s}}} + \frac{2}{d-1} \cdot \frac{(\mathbf{x} - \mathbf{y})^T \mathbf{s} \mathbf{s}^T Q \mathbf{s}}{\mathbf{s}^T Q \mathbf{s} \cdot \sqrt{\mathbf{s}^T Q \mathbf{s}}} \right) \\ &= \frac{d^2-1}{d^2} \cdot \frac{2}{d+1} \left(\frac{\mathbf{s}^T (\mathbf{x} - \mathbf{y})}{\sqrt{\mathbf{s}^T Q \mathbf{s}}} + \frac{2}{d-1} \cdot \frac{(\mathbf{x} - \mathbf{y})^T \mathbf{s}}{\sqrt{\mathbf{s}^T Q \mathbf{s}}} \right) \\ &= \frac{d^2-1}{d^2} \cdot \frac{2}{d-1} \left(\frac{\mathbf{s}^T (\mathbf{x} - \mathbf{y})}{\sqrt{\mathbf{s}^T Q \mathbf{s}}} \right) \end{aligned}$$

1934 The third term simplifies to

$$\begin{aligned} \left(\frac{1}{d+1} \right)^2 \cdot \frac{\mathbf{s}^T Q^T G_+ Q \mathbf{s}}{\mathbf{s}^T Q \mathbf{s}} &= \left(\frac{1}{d+1} \right)^2 \cdot \frac{\mathbf{s}^T Q \left(\frac{d^2-1}{d^2} (G + \frac{2}{d-1} \cdot \frac{\mathbf{s}\mathbf{s}^T}{\mathbf{s}^T Q \mathbf{s}}) \right) Q \mathbf{s}}{\mathbf{s}^T Q \mathbf{s}} \\ &= \frac{d^2-1}{d^2} \cdot \left(\frac{1}{d+1} \right)^2 \cdot \left(\frac{\mathbf{s}^T Q \mathbf{s} + \frac{2}{d-1} (\mathbf{s}^T Q \mathbf{s})}{\mathbf{s}^T Q \mathbf{s}} \right) \\ &= \frac{d^2-1}{d^2} \left(\frac{1}{d^2-1} \right), \end{aligned}$$

1935 Putting all of it together, we obtain that

$$(\mathbf{x} - \mathbf{y}_+)^T G_+ (\mathbf{x} - \mathbf{y}_+) = \frac{d^2-1}{d^2} \left((\mathbf{x} - \mathbf{y})^T G (\mathbf{x} - \mathbf{y}) + \frac{2}{d-1} \frac{(\mathbf{s}^T (\mathbf{x} - \mathbf{y}))^2}{\mathbf{s}^T Q \mathbf{s}} + \frac{2}{d-1} \left(\frac{\mathbf{s}^T (\mathbf{x} - \mathbf{y})}{\sqrt{\mathbf{s}^T Q \mathbf{s}}} \right) + \frac{1}{d^2-1} \right) \quad (4.25)$$

1936 We now argue that $\frac{(\mathbf{s}^T (\mathbf{x} - \mathbf{y}))^2}{\mathbf{s}^T Q \mathbf{s}} + \frac{\mathbf{s}^T (\mathbf{x} - \mathbf{y})}{\sqrt{\mathbf{s}^T Q \mathbf{s}}} = \frac{\mathbf{s}^T (\mathbf{x} - \mathbf{y})}{\mathbf{s}^T Q \mathbf{s}} (\sqrt{\mathbf{s}^T Q \mathbf{s}} + \mathbf{s}^T (\mathbf{x} - \mathbf{y})) \leq 0$. Since $\mathbf{s}^T (\mathbf{x} - \mathbf{y}) \leq 0$ by
1937 (4.24), it suffices to show that $\sqrt{\mathbf{s}^T Q \mathbf{s}} + \mathbf{s}^T (\mathbf{x} - \mathbf{y}) \geq 0$; we will in fact show that $|\mathbf{s}^T (\mathbf{x} - \mathbf{y})| \leq \sqrt{\mathbf{s}^T Q \mathbf{s}}$.

1938 **Claim 4.** $|\mathbf{s}^T (\mathbf{x} - \mathbf{y})| \leq \sqrt{\mathbf{s}^T Q \mathbf{s}}$.

Proof of Claim. Let the eigendecomposition of Q be given as $Q = S\Lambda S^T$, where S is the orthonormal matrix which has the eigenvectors of Q as columns, and Λ is a diagonal matrix with the corresponding eigenvalues. Then $Q^{-1} = S\Lambda^{-1}S^T = G$. Now,

$$\begin{aligned} |\mathbf{s}^T (\mathbf{x} - \mathbf{y})| &= |\mathbf{s}^T S \Lambda^{\frac{1}{2}} \Lambda^{-\frac{1}{2}} S^T (\mathbf{x} - \mathbf{y})| \\ &= | \langle \Lambda^{\frac{1}{2}} S^T \mathbf{s}, \Lambda^{-\frac{1}{2}} S^T (\mathbf{x} - \mathbf{y}) \rangle | \\ &\leq \| \Lambda^{\frac{1}{2}} S^T \mathbf{s} \|_2 \| \Lambda^{-\frac{1}{2}} S^T (\mathbf{x} - \mathbf{y}) \|_2 \\ &= \sqrt{(\Lambda^{\frac{1}{2}} S^T \mathbf{s})^T (\Lambda^{\frac{1}{2}} S^T \mathbf{s})} \sqrt{(\Lambda^{-\frac{1}{2}} S^T (\mathbf{x} - \mathbf{y}))^T (\Lambda^{-\frac{1}{2}} S^T (\mathbf{x} - \mathbf{y}))} \\ &= \sqrt{\mathbf{s}^T S \Lambda^{\frac{1}{2}} \Lambda^{\frac{1}{2}} S^T \mathbf{s}} \sqrt{(\mathbf{x} - \mathbf{y})^T S \Lambda^{-\frac{1}{2}} \Lambda^{-\frac{1}{2}} S^T (\mathbf{x} - \mathbf{y})} \\ &= \sqrt{\mathbf{s}^T Q \mathbf{s}} \sqrt{(\mathbf{x} - \mathbf{y})^T G (\mathbf{x} - \mathbf{y})} \\ &\leq \sqrt{\mathbf{s}^T Q \mathbf{s}}, \end{aligned}$$

1939 where there first inequality is the Cauchy-Schwarz inequality, and the last inequality follows from (4.23). \square

This claim, together with (4.25), implies that

$$\begin{aligned} (\mathbf{x} - \mathbf{y}_+)^T G_+ (\mathbf{x} - \mathbf{y}_+) &\leq \frac{d^2-1}{d^2} \left((\mathbf{x} - \mathbf{y})^T G (\mathbf{x} - \mathbf{y}) + \frac{1}{d^2-1} \right) \\ &\leq \frac{d^2-1}{d^2} \left(1 + \frac{1}{d^2-1} \right) \\ &= 1, \end{aligned}$$

1940 where the second inequality follows from (4.23). This establishes that $\mathbf{x} \in E(Q_+, \mathbf{y}_+)$.

We now prove the volume claim. Let $Q = B^T B$ for some invertible matrix B . We use I_d to denote the $d \times d$ identity matrix. By Proposition 4.47,

$$\begin{aligned} \frac{\text{vol}(E(Q_+, \mathbf{y}_+))}{\text{vol}(E(Q, \mathbf{y}))} &= \sqrt{\frac{\det(Q_+)}{\det(Q)}} \\ &= \sqrt{\frac{\det\left(\frac{d^2}{d^2-1} \left(Q - \frac{2}{d+1} \cdot \frac{Q \mathbf{s} \mathbf{s}^T Q}{\mathbf{s}^T Q \mathbf{s}}\right)\right)}{\det(Q)}} \\ &= \left(\frac{d^2}{d^2-1}\right)^{\frac{d}{2}} \sqrt{\frac{\det\left(Q - \frac{2}{d+1} \cdot \frac{Q \mathbf{s} \mathbf{s}^T Q}{\mathbf{s}^T Q \mathbf{s}}\right)}{\det(Q)}} \\ &= \left(\frac{d^2}{d^2-1}\right)^{\frac{d}{2}} \sqrt{\frac{\det\left(B^T B - \frac{2}{d+1} \cdot \frac{B^T B \mathbf{s} \mathbf{s}^T B^T B}{\mathbf{s}^T B^T B \mathbf{s}}\right)}{\det(B^T B)}} \\ &= \left(\frac{d^2}{d^2-1}\right)^{\frac{d}{2}} \sqrt{\frac{\det\left(B^T \left(I_d - \frac{2}{d+1} \cdot \frac{B \mathbf{s} \mathbf{s}^T B^T}{\mathbf{s}^T B^T B \mathbf{s}}\right) B\right)}{\det(B^T) \det(B)}} \\ &= \left(\frac{d^2}{d^2-1}\right)^{\frac{d}{2}} \sqrt{\frac{\det(B^T) \det\left(I_d - \frac{2}{d+1} \cdot \frac{B \mathbf{s} \mathbf{s}^T B^T}{\mathbf{s}^T B^T B \mathbf{s}}\right) \det(B)}{\det(B^T) \det(B)}} \\ &= \left(\frac{d^2}{d^2-1}\right)^{\frac{d}{2}} \sqrt{\det\left(I_d - \frac{2}{d+1} \cdot \frac{B \mathbf{s} \mathbf{s}^T B^T}{\mathbf{s}^T B^T B \mathbf{s}}\right)} \\ &= \left(\frac{d^2}{d^2-1}\right)^{\frac{d}{2}} \cdot \left(1 - \frac{2}{d+1}\right)^{\frac{1}{2}}, \end{aligned}$$

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 $d - 1$. Now finally we observe that

$$\begin{aligned} \left(\frac{d^2}{d^2-1}\right)^{\frac{d}{2}} \cdot \left(1 - \frac{2}{d+1}\right)^{\frac{1}{2}} &= \left(\frac{d^2}{d^2-1} \cdot \left(1 - \frac{2}{d+1}\right)^{\frac{1}{d}}\right)^{\frac{d}{2}} \\ &\leq \left(\frac{d^2}{d^2-1} \cdot \left(1 - \frac{2}{d(d+1)}\right)\right)^{\frac{d}{2}} \\ &= \left(\frac{d^2(d^2+d-2)}{d(d+1)(d^2-1)}\right)^{\frac{d}{2}} \\ &= \left(1 - \frac{1}{(d+1)^2}\right)^{d/2} \end{aligned}$$

1941 This completes the proof. □

1942 This can be used to give the guarantee of the ellipsoid method as follows.

Theorem 4.49. Using the ellipsoid method with $E_0 = B(\mathbf{x}^0, R)$, if $h(t) > 0$ for some iteration $t \geq 0$, then

$$\min_{j=i_1, \dots, i_{h(t)}} f(\mathbf{x}^j) \leq f(\mathbf{x}^*) + MR \left(1 - \frac{1}{(d+1)^2}\right)^{t/2} \cdot \left(\frac{\text{vol}(E_0)}{\text{vol}(C)}\right)^{1/d} \leq MR e^{-\frac{t}{2(d+1)^2}} \cdot \left(\frac{\text{vol}(E_0)}{\text{vol}(C)}\right)^{1/d},$$

1943 where M is a Lipschitz constant for f over $B(\mathbf{x}^0, 2R)$.

1944 *Proof.* The first inequality follows from Theorem 4.41 part 3., and the fact that $B(\mathbf{x}^*, v_{\min}) \subseteq B(\mathbf{x}^0, 2R)$
 1945 implying that M is a Lipschitz constant for f over $B(\mathbf{x}^*, v_{\min})$, and $\text{vol}(E_{t+1}) \leq \left(1 - \frac{1}{(d+1)^2}\right)^{t/2} \text{vol}(E_0)$ by
 1946 Theorem 4.48. The second inequality follows from the general inequality that $(1+x) \leq e^x$ for all $x \in \mathbb{R}$. □

1947 By setting the error term $MR e^{-\frac{t}{2(d+1)^2}} \cdot \left(\frac{\text{vol}(E_0)}{\text{vol}(C)}\right)^{1/d}$ less than equal to ϵ in Theorem 4.49, the following
 1948 is an immediate consequence.

Corollary 4.50. For any $\epsilon > 0$, after $2(d+1)^2 \left(\ln\left(\frac{MR}{\epsilon}\right) + \frac{1}{d} \ln\left(\frac{\text{vol}(E_0)}{\text{vol}(C)}\right) \right)$ iterations of the ellipsoid method,

$$\min_{j=i_1, \dots, i_{h(t)}} f(\mathbf{x}^j) \leq f(\mathbf{x}^*) + \epsilon.$$

1949 In particular, if there exists $\rho > 0$ such that $B(\mathbf{z}, \rho) \subseteq C$ for some $\mathbf{z} \in C$, then after $2(d+1)^2 \ln\left(\frac{MR^2}{\epsilon\rho}\right)$
 1950 iterations of the ellipsoid method, $\min_{j=i_1, \dots, i_{h(t)}} f(\mathbf{x}^j) \leq f(\mathbf{x}^*) + \epsilon$.

1951 *Proof.* We simply use the fact that $\text{vol}(B(\mathbf{z}, \lambda)) = \lambda^d \text{vol}(B(\mathbf{0}, 1))$ for any $\mathbf{z} \in \mathbb{R}^d$ and $\lambda \geq 0$. □

1952 Because of the logarithmic dependence on the data (M, R, ρ) and the error guarantee ϵ , and the quadratic
 1953 dependence on the dimension d , the ellipsoid method is said to have *polynomial* running time for convex
 1954 optimization.

References

1955

- 1956 [1] Michele Conforti, Gérard Cornuéjols, Aris Daniilidis, Claude Lemaréchal, and Jérôme Malick. Cut-
1957 generating functions and S-free sets. *Mathematics of Operations Research*, 40(2):276–391, 2014.
- 1958 [2] Luc Devroye, László Györfi, and Gábor Lugosi. *A probabilistic theory of pattern recognition*, volume 31.
1959 Springer Science & Business Media, 2013.
- 1960 [3] P.M. Gruber. *Convex and Discrete Geometry*, volume 336 of *Grundlehren der Mathematischen Wis-*
1961 *senschaften [Fundamental Principles of Mathematical Sciences]*. Springer-Verlag, Berlin, 2007.
- 1962 [4] B. Grünbaum. Partitions of mass-distributions and of convex bodies by hyperplanes. *Pacific J. Math.*,
1963 10:1257–1261, 1960.
- 1964 [5] Vladimir N Vapnik and A Ya Chervonenkis. On the uniform convergence of relative frequencies of events
1965 to their probabilities. *Theory of Probability & Its Applications*, 16(2):264–280, 1971.