

Convex Optimization

Shan-Hung Wu
shwu@cs.nthu.edu.tw

Department of Computer Science,
National Tsing Hua University, Taiwan

NetDB-ML, Fall 2014

Outline

1 Optimization Problems

- Standard Forms and Terminology
- Problem Classes

2 Convexity

- Convex Sets
- Convex Functions

3 Convex Optimization

- Optimality
- Disciplined Convex Programming and CVX
- LP and QP

4 Algorithms

- Unconstrained Problems
- Constrained Problems
- Large-Scale Problems**

5 Duality

- Weak Duality
- Strong Duality

Outline

1 Optimization Problems

- Standard Forms and Terminology
- Problem Classes

2 Convexity

- Convex Sets
- Convex Functions

3 Convex Optimization

- Optimality
- Disciplined Convex Programming and CVX
- LP and QP

4 Algorithms

- Unconstrained Problems
- Constrained Problems
- Large-Scale Problems**

5 Duality

- Weak Duality
- Strong Duality

Functional Form

- An **optimization problem** is to minimize an **objective** (or cost) function $f : \mathcal{D} \subseteq \mathbb{R}^n \rightarrow \mathbb{R}$:

$$\begin{aligned} & \min_{\mathbf{x}} f(\mathbf{x}) \\ & \text{subject to } \mathbf{x} \in C \end{aligned}$$

where $C \subseteq \mathbb{R}^n$ is called the **feasible set** containing **feasible points** (or variables)

- If $C = \mathbb{R}^n$, we say the optimization problem is unconstrained
- Maximizing f equals to minimizing $-f$
- C can be a set of function **constrains**, i.e.,
 $C = \{\mathbf{x} : g_i(\mathbf{x}) \leq 0, i = 1, \dots, m\}$
 - Sometimes, we single out equality constrains
 $C = \{\mathbf{x} : g_i(\mathbf{x}) \leq 0, h_j(\mathbf{x}) = 0, i = 1, \dots, m, j = 1, \dots, p\}$
 - Each equality constrain can be written as two inequality constrains

Epigraph form

- We can always assume that the objective is a linear function of the variables, via the *epigraph* ($\text{epi}(f) := \{(\mathbf{x}, t) \in \mathbb{R}^{n+1} : \mathbf{x} \in \mathbb{R}^n, t \geq f(\mathbf{x})\}$) representation of the problem

$$\begin{aligned} & \min_{\mathbf{x}, t} t \\ & \text{subject to } f(\mathbf{x}) - t \leq 0, \mathbf{x} \in C \end{aligned}$$

- The objective function is $\hat{A} : \mathbb{R}^{n+1} \rightarrow \mathbb{R}$, with values $\hat{A}(\mathbf{x}, t) = t$
- Consider the t -sublevel set of \hat{A} (i.e., $\{\mathbf{x} : t \geq \hat{A}(\mathbf{x})\}$), the problem amounts to finding the smallest t for which the corresponding sub-level set intersects the set of points satisfying the constraints

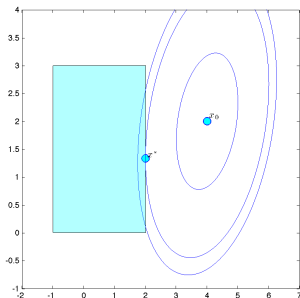
Geometric View

Functional form:

$$\min_{\mathbf{x}} 0.9x_1^2 - 0.4x_1x_2 - 0.6x_2^2 - 6.4x_1 - 0.8x_2 : -1 \leq x_1 \leq 2, 0 \leq x_2 \leq 3$$

Epigraph form:

$$\min_{\mathbf{x}, t} t : t \geq 0.9x_1^2 - 0.4x_1x_2 - 0.6x_2^2 - 6.4x_1 - 0.8x_2, -1 \leq x_1 \leq 2, 0 \leq x_2 \leq 3$$



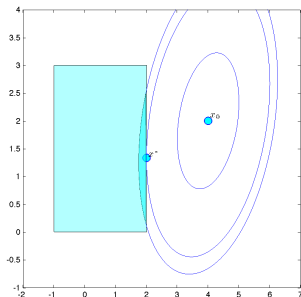
The level sets of the objective function are shown as blue lines, and the feasible set is the light-blue box. The problem amounts to find the smallest value of t such that $t = f(\mathbf{x})$ for some feasible \mathbf{x} . The two dots are the unconstrained and constrained optimal values respectively

Terminology (1)

- $p^* := \inf_{\mathbf{x}} f(\mathbf{x}) : \mathbf{x} \in C$ is called the **optimal value**, which
 - may not exist if the problem is infeasible
 - may not be attained (e.g., in $\min_x e^{-x}$, $p^* = 0$ is attained only when $x \rightarrow \infty$)
- We allow p^* to take on the values ∞ and $-\infty$ when the problem is either
 - infeasible (the feasible set is empty), or
 - unbounded below (there exists feasible points such that $f(\mathbf{x}) \rightarrow -\infty$), respectively
- A feasible point \mathbf{x}^* is called the **optimal point** if $f(\mathbf{x}^*) = p^*$
- The **optimal set** X^* is the set of all optimal points, i.e.,
$$X^* := \{\mathbf{x} \in C : f(\mathbf{x}) = p^*\} = \operatorname{argmin}_{\mathbf{x}} f(\mathbf{x}) : \mathbf{x} \in C$$
- We say the problem is **attained** iff $C \neq \emptyset$ and p^* is attained (or equivalently, $X^* \neq \emptyset$)

Terminology (2)

- The **ϵ -suboptimal set** X^ϵ is defined as $X^\epsilon := \{\mathbf{x} \in C : f(\mathbf{x}) \leq p^* + \epsilon\}$



An ϵ -suboptimal set is marked in darker color. This corresponds to the set of feasible points that achieves an objective value less or equal than $p^* + \epsilon$

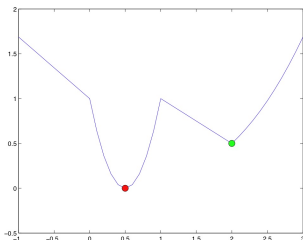
- In practice, we may be only interested in suboptimal solutions

Local vs. Global Optimality

- A point z is **locally optimal** if there is a value $\delta > 0$ such that z is optimal for problem (with new objective $\tilde{f}(\mathbf{x}, z) = f(\mathbf{x})$)

$$\min_{\mathbf{x}} f(\mathbf{x}) : z, \mathbf{x} \in C, \|\mathbf{x} - z\| \leq \delta$$

- That is, a local minimizer minimizes f , but only for its nearby points in the feasible set



Minima of a nonlinear function. The value at a local minimizer is not necessarily the (global) optimal value of the problem, unless f is a “convex” function (in the sense that $\text{epi}(f)$ is a “convex” set)

Outline

1 Optimization Problems

- Standard Forms and Terminology
- Problem Classes

2 Convexity

- Convex Sets
- Convex Functions

3 Convex Optimization

- Optimality
- Disciplined Convex Programming and CVX
- LP and QP

4 Algorithms

- Unconstrained Problems
- Constrained Problems
- Large-Scale Problems**

5 Duality

- Weak Duality
- Strong Duality

Linear Programming

- **Linear Programming (LP)** has the form:¹

$$\begin{aligned} & \min_{\mathbf{x}} \mathbf{c}^T \mathbf{x} \\ & \text{subject to } \mathbf{G}\mathbf{x} \leq \mathbf{h}, \mathbf{A}\mathbf{x} = \mathbf{b} \end{aligned}$$

where $\mathbf{c} \in \mathbb{R}^n$, $\mathbf{G} \in \mathbb{R}^{m \times n}$, $\mathbf{h} \in \mathbb{R}^m$, $\mathbf{A} \in \mathbb{R}^{p \times n}$, and $\mathbf{b} \in \mathbb{R}^p$

- The objective and the $m+p$ constrain functions are **all affine** (i.e., translated linear)
 - Note $\min_{\mathbf{x}} \mathbf{c}^T \mathbf{x} + d$ for some fixed $d \in \mathbb{R}$ amounts to $\min_{\mathbf{x}} \mathbf{c}^T \mathbf{x}$

¹The term “programming” has nothing to do with computer programs. It is named so due to historical reasons.

Quadratic Programming

- *Quadratic Programming (QP)* has the form:

$$\begin{aligned} & \min_{\mathbf{x}} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{c}^T \mathbf{x} \\ & \text{subject to } \mathbf{G} \mathbf{x} \leq \mathbf{h}, \mathbf{A} \mathbf{x} = \mathbf{b} \end{aligned}$$

where $\mathbf{Q} \in \mathbb{R}^{n \times n}$, $\mathbf{c} \in \mathbb{R}^n$, $\mathbf{G} \in \mathbb{R}^{m \times n}$, $\mathbf{h} \in \mathbb{R}^m$, $\mathbf{A} \in \mathbb{R}^{p \times n}$, and $\mathbf{b} \in \mathbb{R}^p$

- The objective is a quadratic function, and the $m + p$ constraint functions are affine

Convex Optimization

- A *convex optimization* problem is of the form:

$$\begin{aligned} & \min_{\mathbf{x}} f(\mathbf{x}) \\ & \text{subject to } \mathbf{x} \in \mathcal{C} \end{aligned}$$

where f is a *convex function*, and \mathcal{C} is a *convex set*

- In particular, with constraints

$$\mathcal{C} = \{\mathbf{x} : g_i(\mathbf{x}) \leq 0, h_j(\mathbf{x}) = 0, i = 1, \dots, m, j = 1, \dots, p\}$$

- g_i must be *convex* functions
- h_j must be *affine* functions (since h_j can be expressed as two g 's, the only way to make both g 's convex is by letting h_j affine)
- Includes LP, QP with positive semidefinite \mathbf{Q} , and more

Combinatorial Optimization

- In combinatorial optimization, some (or all) the variables are Boolean or integers, reflecting discrete choices to be made
 - E.g., Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ be an incidence matrix of a directed graph where $A_{i,j}$ equals to 1 if the arc j starts at node i ; -1 if j ends at i ; 0 otherwise. The problem of finding the shortest path between nodes 1 and m can be expressed as

$$\min_{\mathbf{x}} \mathbf{1}^\top \mathbf{x} : \mathbf{A}\mathbf{x} = [1, 0, \dots, 0, -1]^\top, \mathbf{x} \in \{0, 1\}^n$$

- E.g., the traveling salesman problem
- Generally, extremely hard to solve
- However, they can often be approximately solved with linear or convex programming
 - E.g., the LP-*relaxed* single-pair shortest path problem:

$$\min_{\mathbf{x}} \mathbf{1}^\top \mathbf{x} : \mathbf{A}\mathbf{x} = [1, 0, \dots, 0, -1]^\top, \mathbf{0} \leq \mathbf{x} \in \mathbb{R}^n \leq \mathbf{1}$$

Hard vs. Easy Problems

- We say a problem is hard if cannot be solved in a reasonable amount of time and/or memory space
- Roughly speaking, *convex problems are easy*; non-convex ones are hard
- Of course, not all convex problems are easy, but a (reasonably large) subset
 - E.g., LP and QP with positive semidefinite Q
- Conversely, some non-convex problems are actually easy
 - E.g., the LP-relaxed single-pair shortest path problem has optimal points turn out to be Boolean, so these points are also optimal to the original problem

Outline

1 Optimization Problems

- Standard Forms and Terminology
- Problem Classes

2 Convexity

- Convex Sets
- Convex Functions

3 Convex Optimization

- Optimality
- Disciplined Convex Programming and CVX
- LP and QP

4 Algorithms

- Unconstrained Problems
- Constrained Problems
- Large-Scale Problems**

5 Duality

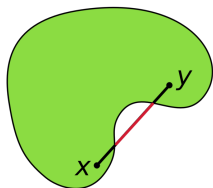
- Weak Duality
- Strong Duality

Convex Sets

Definition (Convex Set)

A set C of points is **convex** iff for any $x, y \in C$ and $\theta \in [0, 1]$, we have $(1 - \theta)x + \theta y \in C$.

- The point $(1 - \theta)x + \theta y$ is called the **convex combination** of points x and y
- Non-convex set:
- Any convex set you know?

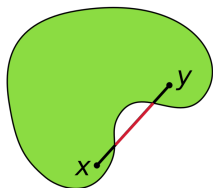


Convex Sets

Definition (Convex Set)

A set C of points is **convex** iff for any $x, y \in C$ and $\theta \in [0, 1]$, we have $(1 - \theta)x + \theta y \in C$.

- The point $(1 - \theta)x + \theta y$ is called the **convex combination** of points x and y
- Non-convex set:
- Any convex set you know? \mathbb{R}^n , non-negative orthant \mathbb{R}_+^n , \emptyset , $\{x\}$, line segments, etc.



- A set is said to be a **convex cone** if it is convex, and has the property that if $x \in C$, then $\theta x \in C$ for every $\theta \geq 0$
 - E.g., \mathbb{R}^n , \mathbb{R}_+^n , union of scalings of a convex set (must contain $\mathbf{0}$)

More Examples

- Subspaces and affine subspaces such as lines, hyperplanes, and higher-dimensional “flat” sets
- Half-spaces, linear varieties (polyhedra, intersections of half-spaces)
- The **convex hulls** of a set of points $\{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ is a convex set:

$$\text{Co}(\mathbf{x}_1, \dots, \mathbf{x}_m) := \left\{ \sum_{i=1}^m \theta_i \mathbf{x}_i : \theta_i \geq 0, \forall i, \sum_{i=1}^m \theta_i = 1 \right\}$$

- Norm balls: $N = \{\mathbf{x} : \|\mathbf{x}\| \leq 1\}$, where $\|\cdot\|$ is some norm on \mathbb{R}^n
 - As for any $\mathbf{x}, \mathbf{y} \in N$,
 $\|(1-\theta)\mathbf{x} + \theta\mathbf{y}\| \leq \|(1-\theta)\mathbf{x}\| + \|\theta\mathbf{y}\| = (1-\theta)\|\mathbf{x}\| + \theta\|\mathbf{y}\| \leq 1$
- The set of all (symmetric) positive semidefinite matrices, denoted by $\mathbb{S}_+^n \subset \mathbb{R}^{n \times n}$, is a convex cone
 - For any $\mathbf{A}, \mathbf{B} \in \mathbb{S}_+^n$ and $\mathbf{x} \in \mathbb{R}^n$,
 $\mathbf{x}^\top ((1-\theta)\mathbf{A} + \theta\mathbf{B})\mathbf{x} = \mathbf{x}^\top (1-\theta)\mathbf{A}\mathbf{x} + \mathbf{x}^\top \theta\mathbf{B}\mathbf{x} \geq 0$

Operations That Preserve Convexity

- Given a convex set $C_1, C_2 \subseteq \mathbb{R}^n$,
 - Scaling: $\beta C = \{\beta \mathbf{x} : \mathbf{x} \in C\}$ is convex for any $\beta \in \mathbb{R}$
 - Sum: $C_1 + C_2 = \{\mathbf{x}_1 + \mathbf{x}_2 : \mathbf{x}_1 \in C_1, \mathbf{x}_2 \in C_2\}$ is convex
 - Augmentation: $\{(\mathbf{x}_1, \mathbf{x}_2) : \mathbf{x}_1 \in C_1, \mathbf{x}_2 \in C_2\} \subseteq \mathbb{R}^{2n}$ is convex
 - **Intersection:** $C_1 \cap C_2$ is convex [Homework]
- Affine transformation: if a map $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is affine, and C is convex, then the set

$$f(C) := \{f(\mathbf{x}) : \mathbf{x} \in C\}$$

is convex [Proof]

- In particular, the projection of a convex set on a subspace is convex

Outline

1 Optimization Problems

- Standard Forms and Terminology
- Problem Classes

2 Convexity

- Convex Sets
- Convex Functions

3 Convex Optimization

- Optimality
- Disciplined Convex Programming and CVX
- LP and QP

4 Algorithms

- Unconstrained Problems
- Constrained Problems
- Large-Scale Problems**

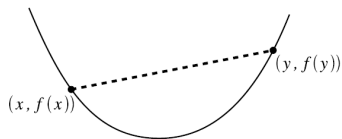
5 Duality

- Weak Duality
- Strong Duality

Convex Functions

Definition (Convex Function)

A function $f : \mathcal{D} \subseteq \mathbb{R}^n \rightarrow \mathbb{R}$ is **convex** iff a) \mathcal{D} is convex; and b) for any $\mathbf{x}, \mathbf{y} \in \mathcal{D}$ and $\theta \in [0, 1]$, we have $f((1-\theta)\mathbf{x} + \theta\mathbf{y}) \leq (1-\theta)f(\mathbf{x}) + \theta f(\mathbf{y})$

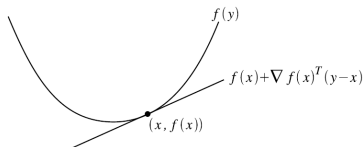


- We say that a function f is
 - **strictly convex** if $f((1-\theta)\mathbf{x} + \theta\mathbf{y}) < (1-\theta)f(\mathbf{x}) + \theta f(\mathbf{y})$ for $\mathbf{x} \neq \mathbf{y}$
 - **concave** if $-f$ is convex

- Condition a) is necessary (what if \mathcal{D} is union of two line segments?)
- Alternatively, f is **convex** iff its epigraph $\text{epi}(f) := \{(\mathbf{x}, t) \in \mathbb{R}^{n+1} : \mathbf{x} \in \mathbb{R}^n, t \geq f(\mathbf{x})\}$ is convex

More Alternate Definitions

- First-order condition: if $f \in \mathcal{C}^1$ is differentiable (that is, \mathcal{D} is open and the gradient exists everywhere on \mathcal{D}), then f is convex iff for any \mathbf{x} and \mathbf{y} ,
$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})$$
 - I.e., the graph of f is bounded below everywhere by anyone of its tangent planes
- Restriction to a line: f is convex iff its restriction to **any** line is convex, i.e., for every $\mathbf{x}_0, \mathbf{v} \in \mathbb{R}^n$, the function $g(t) := f(\mathbf{x}_0 + t\mathbf{v})$ is convex when $\mathbf{x}_0 + t\mathbf{v} \in \mathcal{D}$
- Second-order condition: If f is twice differentiable, then it is convex iff its Hessian $\nabla^2 f$ is positive semidefinite everywhere on \mathcal{D} ; i.e., for any $\mathbf{x} \in \mathcal{D}$, $\nabla^2 f(\mathbf{x}) \succeq \mathbf{O}$



Examples

- $f(x) = e^{ax}$ for $a \in \mathbb{R}$, $f(x) = |x|$, $f(x) = -\log x$ on \mathbb{R}_{++} (strict positive real numbers), negative entropy $f(x) = x \log x$ on \mathbb{R}_{++}
- Affine functions $f(\mathbf{x}) = \mathbf{A}\mathbf{x} + \mathbf{b}$
- Quadratic functions $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{A}\mathbf{x} + \mathbf{b}\mathbf{x} + c$ with positive semidefinite \mathbf{A}
- Function $\lambda_{\max}(\mathbf{X})$ that maps an $n \times n$ symmetric matrix \mathbf{X} to its maximum eigenvalue λ_{\max}
 - Since the condition $\lambda_{\max}(\mathbf{X}) \leq t$ is equivalent to the condition that $t\mathbf{I} - \mathbf{X} \in \mathbb{S}_+^n$, the epigraph is convex
- Norms
 - As $\|(1-\theta)\mathbf{x} + \theta\mathbf{y}\| \leq \|(1-\theta)\mathbf{x}\| + \|\theta\mathbf{y}\| = (1-\theta)\|\mathbf{x}\| + \theta\|\mathbf{y}\|$
- Log-sum-exp $f(\mathbf{x}) = \log \sum_i e^{x_i}$ (a smooth approximation to $f(\mathbf{x}) = \max\{x_i\}$)

Convexity of Sublevel Sets

- Convex functions give rise to a particularly important type of convex set, the t -sublevel set:

Theorem

Given a convex function $f : \mathcal{D} \rightarrow \mathbb{R}$ and $t \in \mathbb{R}$. The t -sublevel set (i.e., $\{\mathbf{x} \in \mathcal{D} : f(\mathbf{x}) \leq t\}$) is Convex.

Proof.

[Homework] □

- Consider an inequality constraint $g \leq 0$ in a convex optimization problem, if g is a convex function, then it defines a convex feasible set, the 0-sublevel set
 - When there are multiple inequality constraints, the final feasible set is the intersection of multiple convex sets, which is still convex

Operations That Preserve Convexity (1)

- Composition with an affine function: if \mathbf{A} in $\mathbb{R}^{m \times n}$, \mathbf{b} in \mathbb{R}^m and $f: \mathbb{R}^m \rightarrow \mathbb{R}$ is convex, then the function $g: \mathbb{R}^n \rightarrow \mathbb{R}$ with values $g(\mathbf{x}) = f(\mathbf{A}\mathbf{x} + \mathbf{b})$ is convex
- **Point-wise maximum**: the pointwise maximum of a family of convex functions is convex—if $\{f_i\}_{i \in \mathcal{A}}$ is a family of convex functions, then the function $f(\mathbf{x}) := \max_{i \in \mathcal{A}} f_i(\mathbf{x})$ is convex
 - E.g., $f(\mathbf{x}) = \max\{x_i\}$, induced matrix norm $\|\mathbf{A}\| = \max_{\mathbf{x}: \|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\|$ is convex
 - Extension: $\sup_{\mathbf{y} \in \mathcal{A}} f(\mathbf{x}, \mathbf{y})$ is convex if for each $\mathbf{y} \in \mathcal{A}$, $f(\mathbf{x}, \mathbf{y})$ is convex in \mathbf{x}
- Nonnegative weighted sum of convex functions is convex
 - E.g., entropy $f(\mathbf{x}) = -\sum_{i=1}^n x_i \log x_i$ for a distribution $\mathbf{x} \in [0, 1]^n$ and $\mathbf{1}^\top \mathbf{x} = 1$ is concave
- Partial minimum: If f is a convex function in (\mathbf{y}, \mathbf{z}) , then the function $g(\mathbf{y}) := \min_{\mathbf{z}} f(\mathbf{y}, \mathbf{z})$ is convex
 - Note that joint convexity in (\mathbf{y}, \mathbf{z}) is essential

Operations That Preserve Convexity (2)

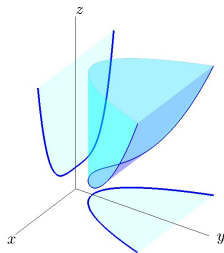
- **Composition with monotone convex functions:** if $f(\mathbf{x}) = h(g_1(\mathbf{x}), \dots, g_k(\mathbf{x}))$, with $g_i: \mathbb{R}^n \rightarrow \mathbb{R}$ convex, $h: \mathbb{R}^k \rightarrow \mathbb{R}$ convex and non-decreasing in each variable, then f is convex
 - For simplicity, assume $k = 1$ and $h, g \in \mathcal{C}^2$. The above conditions ensure that $\nabla^2 g_1(\mathbf{x}) \in \mathbb{R}^{n \times n} \succeq \mathbf{O}$, $h''(\mathbf{y}) \in \mathbb{R}^n \geq 0$, and $h'(\mathbf{y}) \in \mathbb{R}^n \geq 0$
 - Then for any $\mathbf{x} \in \mathcal{D}$, (remember the chain and product rules?)

$$\begin{aligned}\nabla^2 f(\mathbf{x}) &= (\nabla f)'(\mathbf{x})^\top = \left\{ [\nabla g_1(\mathbf{x}) h'(g_1(\mathbf{x}))]' \right\}^\top \\ &= \{ \nabla g_1(\mathbf{x}) h''(g_1(\mathbf{x})) g_1'(\mathbf{x}) + (\nabla g_1)'(\mathbf{x}) h'(g_1(\mathbf{x})) \}^\top \\ &= h''(g_1(\mathbf{x})) \{ \nabla g_1(\mathbf{x}) \nabla g_1(\mathbf{x})^\top \} + h'(g_1(\mathbf{x})) \{ \nabla^2 g_1(\mathbf{x}) \} \\ &\succeq \mathbf{O}\end{aligned}$$

- E.g., $\log \sum_i \exp(g_i)$ is convex if g_i is

Operations That Preserve Convexity (3)

- Let $g(x) = x^2$, $h(y) = y^2$ for $y \geq 0$, and $f(x) = h \circ g(x) = x^4$
- To show that $\text{epi}(f)$ is convex, observe first that $f(x) \leq z$ is equivalent to the existence of y such that $h(y) \leq z$ and $g(x) \leq y$
- The above conditions ensure that the set $\{(x, y, z) : h(y) \leq z, g(x) \leq y\}$ in the space of (x, y, z) -variables is convex
- Hence, $\text{epi}(f)$, the projection of that convex set onto the space of (x, z) -variables, is convex



Outline

1 Optimization Problems

- Standard Forms and Terminology
- Problem Classes

2 Convexity

- Convex Sets
- Convex Functions

3 Convex Optimization

- Optimality
- Disciplined Convex Programming and CVX
- LP and QP

4 Algorithms

- Unconstrained Problems
- Constrained Problems
- Large-Scale Problems**

5 Duality

- Weak Duality
- Strong Duality

Problem Revisited

- Form:

$$\begin{aligned} & \min_{\mathbf{x}} f(\mathbf{x}) \\ & \text{subject to } g_i(\mathbf{x}) \leq 0, h_j(\mathbf{x}) = 0, i = 1, \dots, m, j = 1, \dots, p \end{aligned}$$

where f is a **convex function**, g_i are **convex** functions, and h_j are **affine** functions

- $\text{epi}(f)$ is a convex set
- $C = \{\mathbf{x} : g_i(\mathbf{x}) \leq 0, h_j(\mathbf{x}) = 0, i = 1, \dots, m, j = 1, \dots, p\}$ is a convex set
 - g_i 's are convex implies that the 0-sublevel sets $\{\mathbf{x} : g_i(\mathbf{x}) \leq 0\}$ are convex sets
 - C is the intersection of convex sublevel sets and hyperplanes
- The problem amounts to finding the “lowest” point in the set $\text{epi}(f) \cap \{(\mathbf{x}, t) : \mathbf{x} \in C, t \in \mathbb{R}\}$, which is convex
 - Local optimal points are also global optima

Global vs. Local Optima in Convex Optimization

Theorem

For convex problems with objective $f : \mathcal{D} \rightarrow \mathbb{R}$, any locally optimal point is globally optimal. In addition, the optimal set is convex.

Proof.

Let \mathbf{y} and \mathbf{x}^* be a point and a local minimizer of f on the intersection of feasible set C and \mathcal{D} . We need to prove that $f(\mathbf{y}) \geq f(\mathbf{x}^*) = p^*$. By convexity of f and C , we have $\mathbf{x}_\theta := \theta\mathbf{y} + (1-\theta)\mathbf{x}^*$, and:

$$f(\mathbf{x}_\theta) - f(\mathbf{x}^*) \leq \theta f(\mathbf{y}) + (1-\theta)f(\mathbf{x}^*) - f(\mathbf{x}^*) = \theta(f(\mathbf{y}) - f(\mathbf{x}^*)).$$

Since \mathbf{x}^* is a local minimizer, the left-hand side in this inequality is nonnegative for all small enough values of $\theta > 0$. We conclude that the right hand side is nonnegative, i.e., $f(\mathbf{y}) \geq f(\mathbf{x}^*) = p^*$ as claimed.

Also, the optimal set is convex, since it can be written as

$X^* = \{\mathbf{x} \in C \cap \mathcal{D} : f(\mathbf{x}^*) \leq p^*\}$. This ends our proof. □

Outline

1 Optimization Problems

- Standard Forms and Terminology
- Problem Classes

2 Convexity

- Convex Sets
- Convex Functions

3 Convex Optimization

- Optimality
- Disciplined Convex Programming and CVX
- LP and QP

4 Algorithms

- Unconstrained Problems
- Constrained Problems
- Large-Scale Problems**

5 Duality

- Weak Duality
- Strong Duality

Disciplined Convex Programming and CVX

- A convex optimization software can solve a convex optimization problem efficiently
 - E.g., CVX, optimization toolbox in Matlab (for LP and QP)
- But it cannot identify whether a problem, in an arbitrary form, is convex or not
 - Don't expect it to accept any problem you give, and tell you the problem is not convex
- Discipline convex optimization defines
 - A library of convex functions
 - The rule sets corresponding to operations that preserve convexity. E.g., sum, affine composition, pointwise maximum, partial minimization, composition with monotone convex functions, etc.

Outline

1 Optimization Problems

- Standard Forms and Terminology
- Problem Classes

2 Convexity

- Convex Sets
- Convex Functions

3 Convex Optimization

- Optimality
- Disciplined Convex Programming and CVX
- LP and QP

4 Algorithms

- Unconstrained Problems
- Constrained Problems
- Large-Scale Problems**

5 Duality

- Weak Duality
- Strong Duality

Outline

1 Optimization Problems

- Standard Forms and Terminology
- Problem Classes

2 Convexity

- Convex Sets
- Convex Functions

3 Convex Optimization

- Optimality
- Disciplined Convex Programming and CVX
- LP and QP

4 Algorithms

- Unconstrained Problems
- Constrained Problems
- Large-Scale Problems**

5 Duality

- Weak Duality

• Strong Duality

Unconstrained Problems

- Form:

$$\min_{\mathbf{x}} f(\mathbf{x})$$

where f is convex

- For simplicity, here we assume $f \in \mathcal{C}^1$
- Optimality condition: \mathbf{x}^* is optimal iff $\nabla f(\mathbf{x}^*) = 0$
- For general f (other than affine or quadratic), we may not be able to solve \mathbf{x}^* in a close form
- In practice, suboptimal solutions may be acceptable
- There exist iterative algorithms that yield suboptimal points much faster

Iterative Algorithms

- Assumption: the problem is attained (i.e., $C \neq \emptyset$ and p^* is attained)

Algorithm 4.1: General Descent Method

Input: $\mathbf{x}^{(0)}$, an initial guess from \mathcal{D}

1 **repeat**

2 Determine a **search direction** $\mathbf{d}^{(t)} \in \mathbb{R}^n$;

3 **Line search:** Choose a **step size** $\eta^{(t)}$ such that
 $f(\mathbf{x}^{(t)} + \eta^{(t)} \mathbf{d}^{(t)}) < f(\mathbf{x}^{(t)})$;

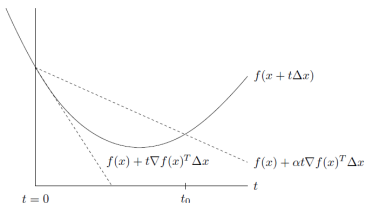
4 **Update rule:** $\mathbf{x}^{(t+1)} \leftarrow \mathbf{x}^{(t)} + \eta^{(t)} \mathbf{d}^{(t)}$;

5 **until** *convergence criterion is satisfied*;

- Convergence criterion: $\|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\| \leq \epsilon$, $\|\nabla f(\mathbf{x}^{(t+1)})\| \leq \epsilon$, etc.
- Line search could be exact: $\eta^{(t)} \leftarrow \arg \min_{\eta > 0} \phi(\eta) := f(\mathbf{x}^{(t)} + \eta \mathbf{d}^{(t)})$, which minimizes f along the ray $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} + \eta \mathbf{d}^{(t)}$, $\forall \eta \in \mathbb{R} > 0$

Backtracking Line Search

- In practice, $\eta^{(t)}$ is usually obtained by another iterations called *backtracking linear search*



($\eta = t$ here)

Algorithm 4.2: Backtracking Line Search

Input: $\alpha \in (0, 0.5)$, $\beta \in (0, 1)$

- 1 $\eta \leftarrow 1$;
- 2 **while** $\mathbf{x}^{(t)} + \eta \mathbf{d}^{(t)} \notin \mathcal{D}$ **do**
- 3 | $\eta \leftarrow \beta \eta$;
- 4 **end**
- 5 **while** $f(\mathbf{x}^{(t)} + \eta \mathbf{d}^{(t)}) = \phi(\eta) > \phi(0) + \alpha \phi'(0) \eta =$
 $f(\mathbf{x}^{(t)}) + \alpha \nabla f(\mathbf{x}^{(t)})^\top \mathbf{d}^{(t)} \eta$ **do**
- 6 | $\eta \leftarrow \beta \eta$;
- 7 **end**

- α , typically in $[0.01, 0.3]$, indicates how much relaxation we accept to the descent direction predicted by the linear extrapolation
- β , typically in $[0.1, 0.8]$, determines how fine-grained the search is

Newton's Method (1)

- Recall that when $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{Q}\mathbf{x} + \mathbf{c}^\top \mathbf{x}$ is quadratic and $\mathbf{Q} \succeq \mathbf{0}$, we can obtain \mathbf{x}^* by solving $\mathbf{Q}\mathbf{x}^* = -\mathbf{c}$
 - No solution if $\mathbf{c} \notin \mathcal{R}(\mathbf{Q})$; otherwise $X^* = \{-\mathbf{Q}^\dagger \mathbf{c} + \mathbf{z} : \mathbf{z} \in \mathcal{N}(\mathbf{Q})\}$
(remember how to solve linear equations using SVD?)
 - When $\mathbf{Q} \succ \mathbf{0}$, $\mathbf{x}^* = -\mathbf{Q}^{-1}\mathbf{c}$ is unique
 - Complexity?

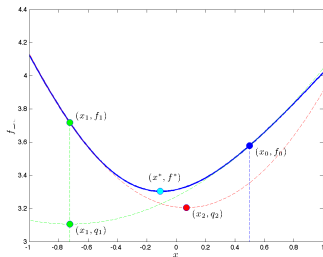
Newton's Method (1)

- Recall that when $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{Q}\mathbf{x} + \mathbf{c}^\top \mathbf{x}$ is quadratic and $\mathbf{Q} \succeq \mathbf{0}$, we can obtain \mathbf{x}^* by solving $\mathbf{Q}\mathbf{x}^* = -\mathbf{c}$
 - No solution if $\mathbf{c} \notin \mathcal{R}(\mathbf{Q})$; otherwise $X^* = \{-\mathbf{Q}^\dagger \mathbf{c} + \mathbf{z} : \mathbf{z} \in \mathcal{N}(\mathbf{Q})\}$
(remember how to solve linear equations using SVD?)
 - When $\mathbf{Q} \succ \mathbf{0}$, $\mathbf{x}^* = -\mathbf{Q}^{-1}\mathbf{c}$ is unique
 - Complexity? $O(n^3)$
- We can leverage the quadratic approximation of a general f to give an iterative algorithm

Newton's Method (2)

- Assumption: $f \in \mathcal{C}^2$ and is strictly convex (i.e., $\nabla^2 f(\mathbf{x}) \succ \mathbf{0}$ everywhere)

Update rule: $\mathbf{x}^{(t+1)} \leftarrow \mathbf{x}^{(t)} - (\nabla^2 f(\mathbf{x}^{(t)}))^{-1} \nabla f(\mathbf{x}^{(t)});$



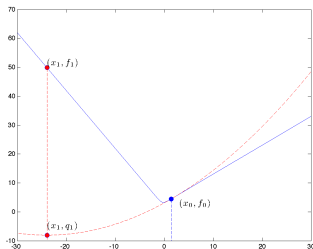
- Based on a **local quadratic approximation** of the the function at the current point \mathbf{x}_t :
$$\tilde{f}(\mathbf{x}) := f(\mathbf{x}^{(t)}) + \nabla f(\mathbf{x}^{(t)})^\top (\mathbf{x} - \mathbf{x}^{(t)}) + \frac{1}{2} (\mathbf{x} - \mathbf{x}^{(t)})^\top \nabla^2 f(\mathbf{x}^{(t)}) (\mathbf{x} - \mathbf{x}^{(t)})$$
- $\mathbf{x}^{(t+1)}$ is set to be a solution to the problem of minimizing \tilde{f}

Remarks (1)

- Pros:
 - No need for line search (although in practice, we often set $\mathbf{d}^{(n)} = -(\nabla^2 f(\mathbf{x}^{(t)}))^{-1} \nabla f(\mathbf{x}^{(t)})$ and perform linear search)
 - Converges fast (1 iteration for quadratic f)
- Cons:
 - Computing $(\nabla^2 f(\mathbf{x}_t))^{-1}$ may be too costly for large-scale problems
 - $\nabla^2 f(\mathbf{x}_t)$ may be singular or ill-conditioned (try $\mathbf{d}^{(n)} = -[\nabla^2 f(\mathbf{x}^{(t)}) + \mu \mathbf{I}]^{-1} \nabla f(\mathbf{x}^{(t)})$ instead)

Remarks (2)

- Might fail to converge for some convex functions
 - Works best for self-concordant functions, whose the Hessians do not vary too fast



- Failure of the Newton method. x_0 is chosen in a region where the function is almost linear. As a result, the quadratic approximation is almost a straight line, and the Hessian is close to zero, sending x_1 to a relatively large negative value. The method quickly diverges in this case

Gradient Descent (1)

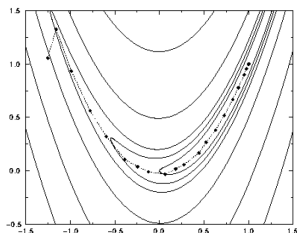
- Assumption: $f \in \mathcal{C}^1$
- Recall that at a given point \mathbf{x} , $\nabla f(\mathbf{x})$ points to the steepest ascend direction

Search direction: $\mathbf{d}^{(t)} = -\nabla f(\mathbf{x}^{(t)});$

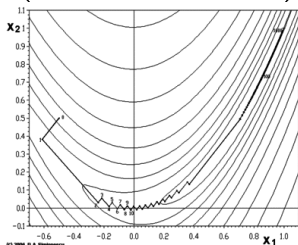
- Since $\nabla f(\mathbf{x}^{(t)} + \eta^{(t)} \mathbf{d}^{(t)})^\top \mathbf{d}^{(t)} = 0$, the next gradient $\nabla f(\mathbf{x}^{(t+1)})$ is orthogonal to the current descent direction $\mathbf{d}^{(t)} = -\nabla f(\mathbf{x}^{(t)})$

Remarks

- Pros:
 - Easy to implement
 - Requires only the first order information on f (computing each iteration is cheap)
- Cons:
 - Much more iterations (as compared to the Newton's method) to convergence
 - “Zig-zagging” around a narrow valley with flat bottom
 - E.g., Rosenbrock's banana:
$$f(\mathbf{x}) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$



(Newton vs. Gradient)

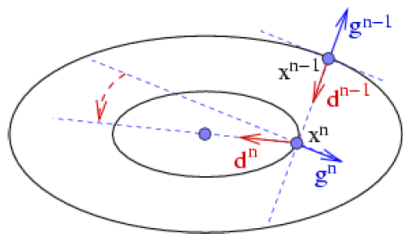


Conjugate Gradient Descent (1)

- A simple variation of the gradient descent
 - Line search and update rule are the same
 - But tilt the next search direction to better aim at the minimum of the Hessian of f

Search direction: $\mathbf{d}^{(t)} = -\nabla f(\mathbf{x}^{(t)}) + c^{(t)} \mathbf{d}^{(t-1)}$ for some constant $c^{(t)}$;

- $c^{(t)}$ can be $\frac{\|\nabla f(\mathbf{x}^{(t)})\|^2}{\|\nabla f(\mathbf{x}^{(t-1)})\|^2}$,
 $\frac{(\nabla f(\mathbf{x}^{(t)}) - \nabla f(\mathbf{x}^{(t-1)}))^\top \nabla f(\mathbf{x}^{(t)})}{\|\nabla f(\mathbf{x}^{(t-1)})\|^2}$,
etc.
- Designed to perform well on quadratic functions



$$(\mathbf{g}^{(t)} := \nabla f(\mathbf{x}^{(t)}))$$

Conjugate Gradient Descent (2)

- Suppose $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{A}\mathbf{x} + \mathbf{b}^\top \mathbf{x}$ is quadratic (so that $\nabla f(\mathbf{x}) = \mathbf{A}\mathbf{x} + \mathbf{b}$)
- Idea: instead of searching for $\mathbf{x}^{(t+1)}$ minimizing f along $\mathbf{x}^{(t)} - \eta \nabla f(\mathbf{x}^{(t)})$, seek for $\mathbf{x}^{(t+1)}$ minimizing f **in the affine space**
 $\mathcal{W}^{(t+1)} := \mathbf{x}^{(0)} + \text{span}(\mathbf{d}^{(0)}, \mathbf{d}^{(1)}, \dots, \mathbf{d}^{(t-1)}, \nabla f(\mathbf{x}^{(t)}))$

Lemma

If $\mathbf{x}^{(t+1)}$ is the minimizer of f in $\mathcal{W}^{(t+1)}$, then $\nabla f(\mathbf{x}^{(t+1)}) \perp \mathcal{W}^{(t+1)}$.

Proof.

Otherwise, we can decrease f along the projection of $\nabla f(\mathbf{x}^{(t+1)})$ onto $\mathcal{W}^{(t+1)}$, contradicting to that $\mathbf{x}^{(t+1)}$ is the minimizer. □

Conjugate Gradient Descent (3)

Lemma

Let $\mathbf{x}^{(t)}$ be the minimizer of f in $\mathcal{W}^{(t)}$. From $\mathbf{x}^{(t)}$, the direction $\mathbf{d}^{(t)}$ points to the minimizer $\mathbf{x}^{(t+1)}$ in $\mathcal{W}^{(t+1)}$ iff $\mathbf{d}^{(t)\top} \mathbf{A} \mathbf{d}^{(i)} = 0$ for $0 \leq i \leq t-1$. The direction $\mathbf{d}^{(t)}$ is said to be **conjugate** to all previous $\mathbf{d}^{(i)}$.

Proof.

By definition, we have $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} + \eta \mathbf{d}^{(t)}$ and

$$\nabla f(\mathbf{x}^{(t+1)}) = \mathbf{A} \mathbf{x}^{(t+1)} + \mathbf{b} = \nabla f(\mathbf{x}^{(t)}) + \eta \mathbf{A} \mathbf{d}^{(t)}.$$

From the above lemma $\nabla f(\mathbf{x}^{(t+1)}) \perp \mathcal{W}^{(t+1)}$ and $\nabla f(\mathbf{x}^{(t)}) \perp \mathcal{W}^{(t)}$, we have

$$0 = \nabla f(\mathbf{x}^{(t+1)})^\top \nabla f(\mathbf{x}^{(t)}) = \|\nabla f(\mathbf{x}^{(t)})\|^2 + \eta \mathbf{d}^{(t)\top} \mathbf{A} \nabla f(\mathbf{x}^{(t)}),$$

implying $\eta \neq 0$. Furthermore,

$$0 = \nabla f(\mathbf{x}^{(t+1)})^\top \mathbf{d}^{(i)} = \nabla f(\mathbf{x}^{(t)})^\top \mathbf{d}^{(i)} + \eta \mathbf{d}^{(t)\top} \mathbf{A} \mathbf{d}^{(i)} = \eta \mathbf{d}^{(t)\top} \mathbf{A} \mathbf{d}^{(i)},$$

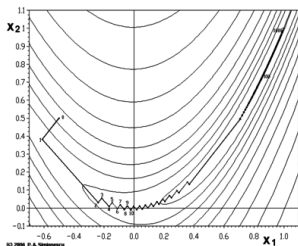
implying $\mathbf{d}^{(t)\top} \mathbf{A} \mathbf{d}^{(i)} = 0$ for all i . □

Conjugate Gradient Descent (4)

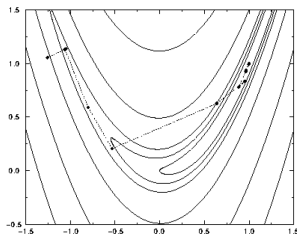
- How to find $\mathbf{d}^{(t)}$ such that it is conjugate to all $\mathbf{d}^{(i)}$?
- Notice that $\nabla f(\mathbf{x}^{(t+1)}) - \nabla f(\mathbf{x}^{(t)}) = \mathbf{A}(\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}) = \eta \mathbf{A} \mathbf{d}^{(t)}$ (see the proof of the above lemma).
- So, $\mathbf{d}^{(t)\top} \mathbf{A} \mathbf{d}^{(i)} = 0 \Rightarrow \mathbf{d}^{(t)\top} (\nabla f(\mathbf{x}^{(t+1)}) - \nabla f(\mathbf{x}^{(t)})) = 0 \Rightarrow \mathbf{d}^{(t)\top} \nabla f(\mathbf{x}^{(t+1)}) = \mathbf{d}^{(t)\top} \nabla f(\mathbf{x}^{(t)}) = \text{some constant}$
- Since $\nabla f(\mathbf{x}^{(i)})$ forms an orthogonal family, we have $\mathbf{d}^{(t)}$ a scaling of
$$\sum_{i=0}^t \frac{\nabla f(\mathbf{x}^{(i)})}{\|\nabla f(\mathbf{x}^{(i)})\|^2}$$
- Apply the above to $\mathbf{d}^{(0)}, \mathbf{d}^{(1)}, \dots, \mathbf{d}^{(t)}$, we have
$$\mathbf{d}^{(t)} = -\nabla f(\mathbf{x}^{(t)}) + c^{(t)} \mathbf{d}^{(t-1)}$$
 - You can easily verify that $c^{(t)} = \frac{\|\nabla f(\mathbf{x}^{(t)})\|^2}{\|\nabla f(\mathbf{x}^{(t-1)})\|^2}$ makes the equation holds

Remarks

- Pros:
 - Easy to implement
 - Still a first order method (same cheap iterations as in gradient descent)
 - Converges fast (at most n iterations for quadratic function $f: \mathbb{R}^n \rightarrow \mathbb{R}$)
 - Can be applied to non-quadratic f , by replacing \mathbf{A} with the Hessian of f
 - Works well if $\nabla^2 f(\mathbf{x}^{(t+1)})$ and $\nabla^2 f(\mathbf{x}^{(t)})$ do not vary too much
- Caution:
 - For general f , \mathbf{d}^n may not be a descent direction. Set it to $-\nabla f(\mathbf{x}^{(t)})$ in this case



(Gradient vs. Conjugate Gradient)



Outline

1 Optimization Problems

- Standard Forms and Terminology
- Problem Classes

2 Convexity

- Convex Sets
- Convex Functions

3 Convex Optimization

- Optimality
- Disciplined Convex Programming and CVX
- LP and QP

4 Algorithms

- Unconstrained Problems
- **Constrained Problems**
- Large-Scale Problems**

5 Duality

- Weak Duality
- Strong Duality

Constrained Problems

- Form:

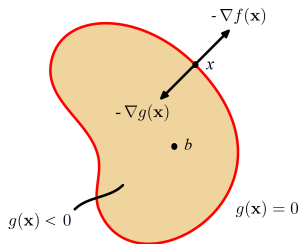
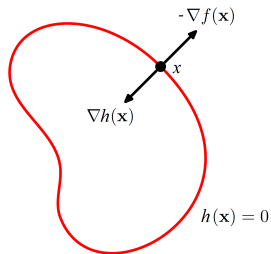
$$\begin{aligned} & \min_{\mathbf{x}} f(\mathbf{x}) \\ & \text{subject to } \mathbf{x} \in C = \{\mathbf{x} : g_i(\mathbf{x}) \leq 0, h_j(\mathbf{x}) = 0, i = 1, \dots, m, j = 1, \dots, p\} \end{aligned}$$

where f and g_i are convex, h_j are affine

- For simplicity, here we assume $f \in \mathcal{C}^1$
- Optimality condition: \mathbf{x}^* is optimal iff $\nabla f(\mathbf{x}^*)^\top (\mathbf{x} - \mathbf{x}^*) \geq 0, \forall \mathbf{x} \in C$, as $f(\mathbf{x}) \geq f(\mathbf{x}^*) + \nabla f(\mathbf{x}^*)^\top (\mathbf{x} - \mathbf{x}^*)$

Active Sets

- Define the **active set** $\mathcal{A}(\mathbf{x})$ at a point \mathbf{x} as the set of constrains θ 's such that $\theta(\mathbf{x}) = 0$, i.e., $\mathcal{A}(\mathbf{x}) := \{\theta : \theta(\mathbf{x}) = 0\}$
 - Equality constrains h_j 's are always active
- Recall for any constrain θ , the gradient $\nabla\theta(\mathbf{x})$ is orthogonal to a tangent line/space passing through the level set at \mathbf{x}
- \mathbf{x}^* occurs when
 - $\forall j, \nabla h_j(\mathbf{x}^*)$ and $-\nabla f(\mathbf{x}^*)$ are **parallel** (i.e., $-\nabla f(\mathbf{x}^*) = \nu_j \nabla h_j(\mathbf{x}^*)$ for some $\nu_j \neq 0$)
 - $\forall i$ such that g_i is active, $-\nabla g_i(\mathbf{x}^*)$ and $-\nabla f(\mathbf{x}^*)$ are **opposite** (i.e., $-\nabla f(\mathbf{x}^*) = \lambda_i \nabla g_i(\mathbf{x}^*)$ for some $\lambda_i > 0$)



Iterative Algorithms

- Assumption: the problem is attained (i.e., $C \neq \emptyset$ and p^* is attained)
- Iterative algorithms in the presence of constraints?

Iterative Algorithms

- Assumption: the problem is attained (i.e., $C \neq \emptyset$ and p^* is attained)
 - Iterative algorithms in the presence of constraints?
- 1 Transform the constrained problem into a unconstrained one, or
 - 2 Make sure that $\mathbf{x}^{(t+1)}$ falls inside the feasible set during each iteration

Exterior-Point Methods

- For equality constraints $h_j(\mathbf{x}) = 0$
- Idea: penalize non-admissible solutions
- Create “barrier functions” $\psi_j(\mathbf{x})$ such that $\psi_j(\mathbf{x}) = 0$ if $h_j(\mathbf{x}) = 0$; $\psi_j(\mathbf{x}) \gg 0$ otherwise
 - E.g., $\psi_j(\mathbf{x}) = \mu \|h_j(\mathbf{x})\|^2$ for some large μ
- Solve the unconstrained problem: $\min_{\mathbf{x}} f(\mathbf{x}) + \mu \sum_{j=1}^p \psi_j(\mathbf{x})$
 - Objective is still convex
- A solution falls outside the feasible set, an “exterior point”

Interior-Point Methods

- For inequality constraints $g_i(\mathbf{x}) \leq 0$
- Assumption: the original problem is **strictly** feasible (i.e., there exists $\mathbf{x} \in X^*$ such that $g_i(\mathbf{x}) < 0$ for all i)
- Idea: penalize non-admissible solutions
- Create barrier functions $\psi_i(\mathbf{x})$ such that $\psi_i(\mathbf{x}) = 0$ if $g_i(\mathbf{x}) \leq 0$; $\psi_i(\mathbf{x}) \gg 0$ otherwise
 - E.g., the **logarithmic barrier** $\psi_i(\mathbf{x}) = -\mu \log(-g_i(\mathbf{x}))$ for some μ
- Solve the unconstrained problem (still convex):
$$\min_{\mathbf{x}} f(\mathbf{x}) - \mu \sum_{i=1}^m \log(-g_i(\mathbf{x}))$$
- A solution falls inside the feasible set, an “interior point”

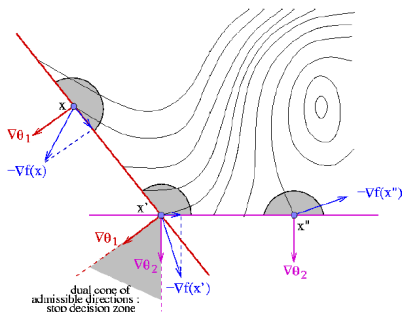
- For μ large, solving the above problem results in a point well aligned/inside the feasible set
- As $\mu \rightarrow 0$ the solution converges to a global minimizer for the original, constrained problem
 - In fact, the theory of convex optimization says that if we set $\mu = m/\epsilon$ (or $\mu = p/\epsilon$ for equality constraints), then the minimizer is ϵ -suboptimal.
- In practice, we solve the unconstrained problem several times, with μ from large to small

Projected Gradient Descent (1)

- $\mathbf{x}^{(t+1)}$ may fall outside C during an iteration
- Idea: if so, project $\mathbf{x}^{(t+1)}$ onto the boundary of C

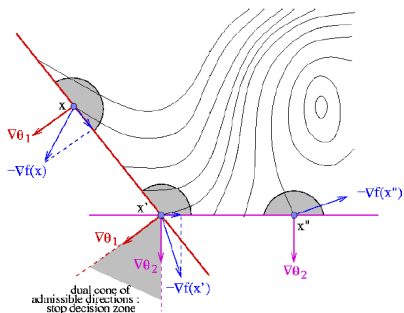
Update rule: $\mathbf{x}^{(t+1)} \leftarrow \mathbf{P}(\mathbf{x}^{(t)} - \eta^{(t)} \nabla f(\mathbf{x}^{(t)}))$ for some projector \mathbf{P} ;

- For simplicity, we consider only the affine constraints here
- Suppose $\mathbf{x}^{(t)}$ is already on the boundary of C
- We can identify the active set $\mathcal{A}(\mathbf{x}^{(t)})$ at $\mathbf{x}^{(t)}$
- Define the **tangent space of active constraints** at $\mathbf{x}^{(t)}$:
$$\bigcap_{\theta \in \mathcal{A}(\mathbf{x}^{(t)})} \{\mathbf{x} : \nabla \theta(\mathbf{x}^{(t)})^\top (\mathbf{x} - \mathbf{x}^{(t)}) = 0\}$$
- We seek for the projection of $\mathbf{x}^{(t+1)}$ onto that tangent space



Projected Gradient Descent (2)

- Since $\mathbf{x}^{(t)}$ is already in the tangent space, the update rule can be written as
$$\mathbf{x}^{(t+1)} \leftarrow (\mathbf{x}^{(t)} - \eta^{(t)} \mathbf{P} \nabla f(\mathbf{x}^{(t)}))$$
(recall $\mathbf{P}^2 = \mathbf{P}$)
- $\nabla \theta(\mathbf{x}^{(t)})^\top (\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}) = 0$ implies
$$\nabla \theta(\mathbf{x}^{(t)})^\top (-\eta^{(t)} \mathbf{P} \nabla f(\mathbf{x}^{(t)})) = 0$$
- Let
$$\Theta = [\nabla \theta_1(\mathbf{x}^{(t)}), \dots, \nabla \theta_a(\mathbf{x}^{(t)})] \in \mathbb{R}^{n \times a}, \text{ where } a = |\mathcal{A}(\mathbf{x}^{(t)})|$$
- We instead seek for the projection of $-\nabla f(\mathbf{x}^{(t)})$ onto $\{\mathbf{x} : \Theta^\top \mathbf{x} = \mathbf{0}\}$



Projected Gradient Descent (3)

- Target: $-P\nabla f(\mathbf{x}^{(t)}) \in \{\mathbf{x} : \Theta^\top \mathbf{x} = \mathbf{0}\}$. How to find P ?

Projected Gradient Descent (3)

- Target: $-P\nabla f(\mathbf{x}^{(t)}) \in \{\mathbf{x} : \Theta^\top \mathbf{x} = \mathbf{0}\}$. How to find P ?
- Recall from the fundamental theorem of linear algebra that $\{\mathbf{x} : \Theta^\top \mathbf{x} = \mathbf{0}\} = \mathcal{R}(\Theta)^\perp = \text{span}(\nabla\theta_1(\mathbf{x}^{(t)}), \dots, \nabla\theta_a(\mathbf{x}^{(t)}))^\perp$
- Also, recall that the projection of any point \mathbf{y} onto $\mathcal{R}(\Theta)$ is $\Theta\mathbf{x}^*$, where $\mathbf{x}^* = (\Theta^\top \Theta)^{-1} \Theta^\top \mathbf{y}$ is the solution to the least square problem

$$\arg \min_{\mathbf{x}} \|\Theta \mathbf{x} - \mathbf{y}\|^2$$

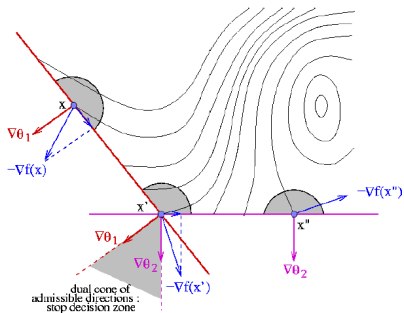
- Let $\mathbf{Q} = \Theta(\Theta^\top \Theta)^{-1} \Theta^\top$, the projection of \mathbf{y} onto $\mathcal{R}(\Theta)^\perp$ is $\mathbf{y} - \mathbf{Q}\mathbf{y} = (\mathbf{I} - \mathbf{Q})\mathbf{y}$, so $\mathbf{P} = \mathbf{I} - \mathbf{Q}$

The Changing Active Sets

- We may encounter $-\mathbf{P}\nabla f(\mathbf{x}^{(t)}) = \mathbf{0}$ during an iteration. Should we stop?

The Changing Active Sets

- We may encounter $-\mathbf{P}\nabla f(\mathbf{x}^{(t)}) = \mathbf{0}$ during an iteration. Should we stop?
- No, some constrains θ in $\mathcal{A}(\mathbf{x}^{(t)})$ may be “unnecessary,” i.e., we cannot find $\eta > 0$ such that $\mathbf{x}^{(t)} - \eta\mathbf{P}_\theta\nabla f(\mathbf{x}^{(t)})$ is on the boundary of C ,
 - \mathbf{P}_θ projects $\mathbf{d}^{(t)}$ onto $\{\mathbf{x} : \nabla\theta(\mathbf{y}^{(t)})^\top \mathbf{x} = 0\}$
 - We can obtain η by first solving $g(\mathbf{x}^{(t)} - \eta\mathbf{P}_\theta\nabla f(\mathbf{x}^{(t)})) = 0$ for each another constrain $g \in C$, and then take the minimum of the solutions that are in $(0, \infty)$
- Remove all such constrains θ 's in $\mathcal{A}(\mathbf{x}^{(t)})$. Stop only if $\mathcal{A}(\mathbf{x}^{(t)}) = \emptyset$



Algorithm 4.3: Projected Gradient Descent Method

Input: $\mathbf{x}^{(0)}$, an initial guess from $\mathcal{D} \cap \mathcal{C}$

```
1 repeat
2    $\mathbf{d}^{(t)} \leftarrow -\nabla f(\mathbf{x}^{(t)});$ 
3   Determine  $\eta^{(t)}$ ;
4    $\mathbf{x}^{(t+1)} \leftarrow \mathbf{x}^{(t)} + \eta^{(t)} \mathbf{d}^{(t)};$ 
5   if  $\mathbf{x}^{(t+1)} \notin \mathcal{C}$  then
6      $\mathbf{y}^{(t)} \leftarrow \mathbf{x}^{(t)} + \eta' \mathbf{d}^{(t)}$  is the intersect between  $\{\mathbf{x}^{(t)} + \eta \mathbf{d}^{(t)} : \eta > 0\}$  and
       the boundary of  $\mathcal{C}$ ;
7      $\mathcal{A}(\mathbf{y}^{(t)}) \leftarrow$  set of active constrains at  $\mathbf{y}^{(t)}$ , excluding those  $\theta$ 's such that
       there is no intersect between  $\{\mathbf{x}^{(t)} + \eta \mathbf{P}_\theta \mathbf{d}^{(t)} : \eta > 0\}$  and the boundary
       of  $\mathcal{C}$ ;
8     if  $\mathcal{A}(\mathbf{y}^{(t)}) \neq \emptyset$  then  $\mathbf{x}^{(t+1)} \leftarrow \mathbf{y}^{(t)} + (\eta^{(t)} - \eta') \mathbf{P} \mathbf{d}^{(t)}$  else  $\mathbf{x}^{(t+1)} \leftarrow \mathbf{y}^{(t)};$ 
9   end
10 until convergence criterion is satisfied;
```

Outline

1 Optimization Problems

- Standard Forms and Terminology
- Problem Classes

2 Convexity

- Convex Sets
- Convex Functions

3 Convex Optimization

- Optimality
- Disciplined Convex Programming and CVX
- LP and QP

4 Algorithms

- Unconstrained Problems
- Constrained Problems
- Large-Scale Problems**

5 Duality

- Weak Duality
- Strong Duality

Decomposition Methods

TBA

Weak and Strong Duality

- Next, we show how the notion of *weak duality* allows to develop, in a systematic way, approximations of non-convex problems based on convex optimization.
- Starting with any given minimization problem, which we call the *primal problem*, we can form a *dual problem*, which
 - Is always convex (specifically, a concave maximization problem)
 - Provides a lower bound on the values of the primal
- When the primal is convex, the *strong duality* holds—the dual problem shares the same optimal value as that of the primal
 - Gives more insights to the optimality conditions

Outline

1 Optimization Problems

- Standard Forms and Terminology
- Problem Classes

2 Convexity

- Convex Sets
- Convex Functions

3 Convex Optimization

- Optimality
- Disciplined Convex Programming and CVX
- LP and QP

4 Algorithms

- Unconstrained Problems
- Constrained Problems
- Large-Scale Problems**

5 Duality

- Weak Duality
- Strong Duality

Primal Problem

- Consider a primal problem:

$$\begin{aligned} & \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \\ & \text{subject to } g_i(\mathbf{x}) \leq 0, h_j(\mathbf{x}) = 0, i = 1, \dots, m, j = 1, \dots, p \end{aligned}$$

- f , g_i , and h_j can be arbitrary (need not be convex or affine)
- For simplicity, let $f(\mathbf{x}) = \infty$ (resp., $g_i(\mathbf{x})$ and $h_j(\mathbf{x})$) if \mathbf{x} is not in the domain of f (resp., g_i and h_j)
- $p^* := \inf_{\mathbf{x}: g_i(\mathbf{x}) \leq 0, h_j(\mathbf{x}) = 0} f(\mathbf{x})$ and \mathbf{x} are call **primal value** and **variables** respectively

Lagrange Function

- Define a **Lagrange function** (or simply **Lagrangian**) $\mathcal{L} : \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^p \rightarrow \mathbb{R}$ with values

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta}) := f(\mathbf{x}) + \sum_{i=1}^m \alpha_i g_i(\mathbf{x}) + \sum_{j=1}^p \beta_j h_j(\mathbf{x})$$

- Then the primal problem can be written as

$$\min_{\mathbf{x} \in \mathbb{R}^n} \sup_{\boldsymbol{\alpha} \geq \mathbf{0}} \mathcal{L}(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta})$$

- $p^* = \inf_{\mathbf{x} \in \mathbb{R}^n} \sup_{\boldsymbol{\alpha} \geq \mathbf{0}} \mathcal{L}(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta})$
- This creates “barriers” that penalize $g_i(\mathbf{x}) > 0$ and $h_j(\mathbf{x}) \neq 0$
- The constraints $\boldsymbol{\alpha} \geq \mathbf{0}$ are essential

Dual Problem

- Given a primal problem $\min_{\mathbf{x} \in \mathbb{R}^n} \sup_{\boldsymbol{\alpha} \geq \mathbf{0}} \mathcal{L}(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta})$, define its **dual problem** as

$$\max_{\boldsymbol{\alpha} \geq \mathbf{0}} \inf_{\mathbf{x} \in \mathbb{R}^n} \mathcal{L}(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta})$$

- $d^* := \sup_{\boldsymbol{\alpha} \geq \mathbf{0}} \inf_{\mathbf{x} \in \mathbb{R}^n} \mathcal{L}(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta})$ is called the **dual value**
- It can be easily shown that $d^* = \sup_{\boldsymbol{\alpha} \geq \mathbf{0}} \inf_{\mathbf{x} \in \mathbb{R}^n} \mathcal{L}(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta}) \leq \inf_{\mathbf{x} \in \mathbb{R}^n} \sup_{\boldsymbol{\alpha} \geq \mathbf{0}} \mathcal{L}(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = p^*$ (called max-min inequality) [Homework]
 - d^* is a lower bound of p^*
 - $p^* - d^*$ is called the **duality gap**
- $dual(\boldsymbol{\alpha}, \boldsymbol{\beta}; \mathbf{x}) := \inf_{\mathbf{x} \in \mathbb{R}^n} \mathcal{L}(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta})$ is called the **dual function**
 - Defined as a point-wise minimum (in \mathbf{x}), therefore concave
- The dual problem $\max_{\boldsymbol{\alpha} \geq \mathbf{0}} dual(\boldsymbol{\alpha}, \boldsymbol{\beta})$ is always a concave-maximization problem (convex)

Example

- Consider a primal problem:

$$\begin{aligned} \min_{\mathbf{x}} \quad & \frac{1}{2} \|\mathbf{x}\|^2 \\ \text{subject to} \quad & \mathbf{Ax} \leq \mathbf{b} \end{aligned}$$

- $dual(\boldsymbol{\alpha}; \mathbf{x}) = \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{x}\|^2 + \boldsymbol{\alpha}^\top (\mathbf{Ax} - \mathbf{b}) = -\frac{1}{2} \|\mathbf{A}^\top \boldsymbol{\alpha}\|^2 - \mathbf{b}^\top \boldsymbol{\alpha}$ [Proof]
 - $\mathbf{x}^* = \mathbf{A}^\top \boldsymbol{\alpha}$

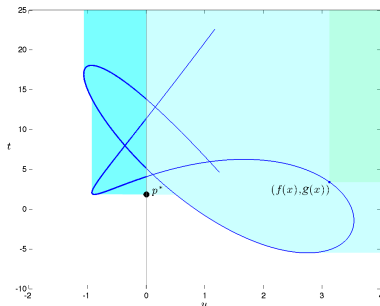
- Dual problem:

$$\begin{aligned} \max_{\boldsymbol{\alpha}} \quad & -\frac{1}{2} \|\mathbf{A}^\top \boldsymbol{\alpha}\|^2 - \mathbf{b}^\top \boldsymbol{\alpha} \\ \text{subject to} \quad & \boldsymbol{\alpha} \geq \mathbf{0} \end{aligned}$$

- Equivalent to $\min_{\boldsymbol{\alpha} \geq \mathbf{0}} \frac{1}{2} \|\mathbf{A}^\top \boldsymbol{\alpha}\|^2 + \mathbf{b}^\top \boldsymbol{\alpha}$

Geometric Interpretation (1)

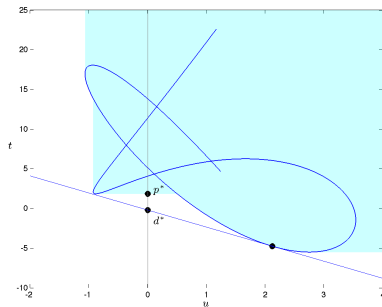
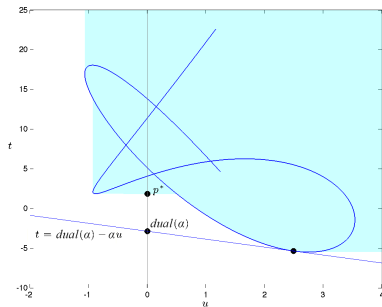
- Consider a primal problem: $\min_{\mathbf{x}} f(\mathbf{x})$ subject to $g(\mathbf{x}) \leq 0$
- Dual problem: $\max_{\alpha \geq 0} \text{dual}(\alpha) = \max_{\alpha \geq 0} \inf_{\mathbf{x}} f(\mathbf{x}) + \alpha g(\mathbf{x})$
- Let $A := \{(u, t) : u \geq g(\mathbf{x}), t \geq f(\mathbf{x})\}$, the blue area



- The solutions are feasible only in the dark blue area

Geometric Interpretation (2)

- $\inf_{\mathbf{x}} f(\mathbf{x}) + \alpha g(\mathbf{x})$ is attained, so we can rewrite the dual function as $dual(\alpha) = \min_{(u,t) \in A} t + \alpha u = t^* + \alpha u^*$
- Given any fixed $\alpha \geq 0$, $\{(u, t) : t = dual(\alpha) - \alpha u\}$ is a line with slope $-\alpha$ intercepting A at (t^*, u^*)
 - The line intercepts $\{(u, t) : u = 0\}$ at $(0, dual(\alpha))$
- The dual problem is to find the best line intercepting A that produce the highest intercept with $\{(u, t) : u = 0\}$



- The dual function $dual$ may not be easy to compute: it is itself an optimization problem!
 - Duality works best when $dual$ can be computed in closed form
- Even if it is possible to compute $dual$, it might not be easy to maximize: convex problems are not always easy to solve
- A lower bound might not be of great practical interest: often we need a sub-optimal solution
 - Duality does not seem at first to offer a way to compute such a primal point
- However, duality is a powerful tool in understanding the problem

Outline

1 Optimization Problems

- Standard Forms and Terminology
- Problem Classes

2 Convexity

- Convex Sets
- Convex Functions

3 Convex Optimization

- Optimality
- Disciplined Convex Programming and CVX
- LP and QP

4 Algorithms

- Unconstrained Problems
- Constrained Problems
- Large-Scale Problems**

5 Duality

- Weak Duality
- Strong Duality

Strong Duality

- Primal problem:

$$\begin{aligned} & \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \\ & \text{subject to } g_i(\mathbf{x}) \leq 0, h_j(\mathbf{x}) = 0, i = 1, \dots, m, j = 1, \dots, p \end{aligned}$$

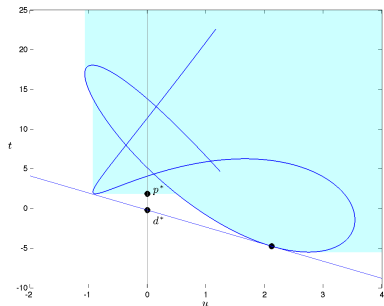
- $p^* := \inf_{\mathbf{x}: g_i(\mathbf{x}) \leq 0, h_j(\mathbf{x}) = 0} f(\mathbf{x})$
- Dual problem:

$$\max_{\boldsymbol{\alpha} \geq \mathbf{0}} \inf_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) + \sum_{i=1}^m \alpha_i g_i(\mathbf{x}) + \sum_{j=1}^p \beta_j h_j(\mathbf{x})$$

- $d^* := \sup_{\boldsymbol{\alpha} \geq \mathbf{0}} \inf_{\mathbf{x} \in \mathbb{R}^n} \mathcal{L}(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta})$
- We say that **strong duality** holds if the duality gap is zero: $d^* = p^*$

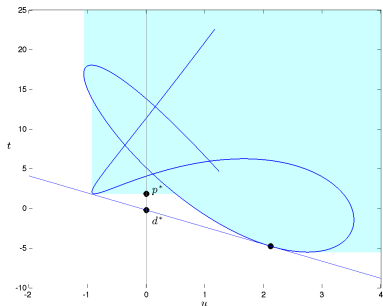
Slater's Sufficient Condition for Strong Duality (1)

- How to make $(0, d^*) = (0, p^*)$?



Slater's Sufficient Condition for Strong Duality (1)

- How to make $(0, d^*) = (0, p^*)$?
- One sufficient condition:
 - 1 $A = \{(u, t) : u \geq g(\mathbf{x}), t \geq f(\mathbf{x})\}$ (the blue area) is a convex set
 - 2 The line $\{(u, t) : t = dual(\alpha) - \alpha u\}$ is not vertical (so d^* is attained)



Slater's Sufficient Condition for Strong Duality (2)

- The above two points imply:
- ① The primal problem is **convex**
 - Since $\{u : u \geq g(\mathbf{x})\}$ and $\{t : t \geq f(\mathbf{x})\}$ are convex, so does A [Proof]
- ② **Slater condition**: the primal problem is **strictly feasible**:
 $\exists \mathbf{x} : g_i(\mathbf{x}) < 0, h_j(\mathbf{x}) = 0$
 - The interior points of $A = \{(u, t) : u \geq g(\mathbf{x}), t \geq f(\mathbf{x})\}$ (the blue area) cut into the area $\{(u, t) : u < 0\}$
 - If $g_i(\mathbf{x})$ is affine, we can relax the feasibility above by $g_i(\mathbf{x}) \leq 0$
- Sufficient condition for strong duality, but **not** necessary

Solving the Dual Problem

- Suppose the strong duality holds, then by solving the dual problem, we obtain:
 - The primal value $p^* = d^*$
 - Furthermore, \mathbf{x}^* if we can write \mathbf{x}^* in a close form with respect to $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ in $dual(\boldsymbol{\alpha}, \boldsymbol{\beta}; \mathbf{x}) := \inf_{\mathbf{x} \in \mathbb{R}^n} \mathcal{L}(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta})$
- Why solving the dual problem instead?

Solving the Dual Problem

- Suppose the strong duality holds, then by solving the dual problem, we obtain:
 - The primal value $p^* = d^*$
 - Furthermore, \mathbf{x}^* if we can write \mathbf{x}^* in a close form with respect to $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ in $dual(\boldsymbol{\alpha}, \boldsymbol{\beta}; \mathbf{x}) := \inf_{\mathbf{x} \in \mathbb{R}^n} \mathcal{L}(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta})$
- Why solving the dual problem instead?
 - We gain insights to the primal problem

Karush-Kuhn-Tucker (KKT) Conditions

Theorem

Suppose f , g_i , and h_j are continuously differentiable at \mathbf{x}^* , and the primal problem is attained, convex, and satisfies the Slater condition. Then a primal variable \mathbf{x}^* is optimal iff there exists α^* and β^* such that the following conditions, called **Karush-Kuhn-Tucker (KKT) conditions** are satisfied:

Lagrangian stationarity:

$$\nabla f(\mathbf{x}^*) + \sum_{i=1}^m \alpha_i^* \nabla g_i(\mathbf{x}^*) + \sum_{j=1}^p \beta_j^* \nabla h_j(\mathbf{x}^*) = 0$$

Primal feasibility: $g_i(\mathbf{x}^*) \leq 0$ and $h_j(\mathbf{x}^*) = 0$ for all $i = 1, \dots, m$ and $j = 1, \dots, p$

Dual feasibility: $\alpha_i^* \geq 0$ for all $i = 1, \dots, m$

Complementary slackness: $\alpha_i^* g_i(\mathbf{x}^*) = 0$ for all $i = 1, \dots, m$

Complementary Slackness

- Why $\alpha_i^* g_i(\mathbf{x}^*) = 0$ for all $i = 1, \dots, m$?
- When strong duality holds and both primal and dual problems are attained, by $(\mathbf{x}^*, \boldsymbol{\alpha}^*, \boldsymbol{\beta}^*)$, we have

$$f(\mathbf{x}^*) + \sum_{i=1}^m \alpha_i^* g_i(\mathbf{x}^*) + \sum_{j=1}^p \beta_j^* h_j(\mathbf{x}^*) = \text{dual}(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*; \mathbf{x}^*) = d^* = p^* = f(\mathbf{x}^*)$$

- Since $\boldsymbol{\alpha}^* \geq 0$, each term in $\sum_{i=1}^m \alpha_i^* g_i(\mathbf{x}^*)$ must be 0
- So what?

Complementary Slackness

- Why $\alpha_i^* g_i(\mathbf{x}^*) = 0$ for all $i = 1, \dots, m$?
- When strong duality holds and both primal and dual problems are attained, by $(\mathbf{x}^*, \boldsymbol{\alpha}^*, \boldsymbol{\beta}^*)$, we have

$$f(\mathbf{x}^*) + \sum_{i=1}^m \alpha_i^* g_i(\mathbf{x}^*) + \sum_{j=1}^p \beta_j^* h_j(\mathbf{x}^*) = \text{dual}(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*; \mathbf{x}^*) = d^* = p^* = f(\mathbf{x}^*)$$

- Since $\boldsymbol{\alpha}^* \geq 0$, each term in $\sum_{i=1}^m \alpha_i^* g_i(\mathbf{x}^*)$ must be 0
- So what? If $\alpha_i^* > 0$, then $g_i(\mathbf{x}^*) = 0$
- We can tell from the values of α_i^* 's which inequality constraint is *active*

Example

- Suppose $\mathbf{A} \in \mathbb{R}^{m \times n}$, in the primal problem:

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^n} \quad & \frac{1}{2} \|\mathbf{x}\|^2 \\ \text{subject to} \quad & \mathbf{Ax} \leq \mathbf{b} \end{aligned}$$

- Dual problem:

$$\begin{aligned} \min_{\boldsymbol{\alpha} \in \mathbb{R}^m} \quad & \frac{1}{2} \|\mathbf{A}^\top \boldsymbol{\alpha}\|^2 + \mathbf{b}^\top \boldsymbol{\alpha} \\ \text{subject to} \quad & \boldsymbol{\alpha} \geq \mathbf{0} \end{aligned}$$

- $\mathbf{x}^* = \mathbf{A}^\top \boldsymbol{\alpha}$
- We now solve m instead of n variables
 - If $n \gg m$, solving the dual problem takes less time
- Furthermore, by complementary slackness, $\boldsymbol{\alpha}^\top (\mathbf{Ax} - \mathbf{b}) = 0$
 - We can tell that the j -th constraint is active (i.e., $\mathbf{A}_{j \cdot} \mathbf{x} = b_j$) iff $\alpha_j \neq 0$