

**A GRADUATE COURSE
IN
OPTIMIZATION**

A GRADUATE COURSE
IN
OPTIMIZATION
CIE6010 Notebook

Prof. Yin Zhang

The Chinese University of Hong Kong, Shenzhen



香港中文大學(深圳)

The Chinese University of Hong Kong, Shenzhen

Contents

| | |
|--|-----------|
| Acknowledgments | vii |
| Notations | ix |
| 1 Week1 | 1 |
| 1.1 Monday | 1 |
| 1.1.1 Introduction to Optimizaiton | 1 |
| 1.2 Wednesday | 2 |
| 1.2.1 Reviewing for Linear Algebra | 2 |
| 1.2.2 Reviewing for Calculus | 2 |
| 1.2.3 Introduction to Optimization | 3 |
| 2 Week2 | 7 |
| 2.1 Monday | 7 |
| 2.1.1 Reviewing and Announments | 7 |
| 2.1.2 Quadratic Function Case Study | 8 |
| 2.2 Wednesday | 11 |
| 2.2.1 Convex Analysis | 11 |
| 3 Week3 | 17 |
| 3.1 Wednesday | 17 |
| 3.1.1 Convex Analysis | 17 |
| 3.1.2 Iterative Method | 18 |
| 3.2 Thursday | 22 |
| 3.2.1 Announcement | 22 |
| 3.2.2 Sparse Large Scale Optimization | 22 |

Acknowledgments

This book is from the CIE6010 in fall semester, 2018.

CUHK(SZ)

Notations and Conventions

| | |
|---|---|
| X | Set |
| $\inf X \subseteq \mathbb{R}$ | Infimum over the set X |
| $\mathbb{R}^{m \times n}$ | set of all $m \times n$ real-valued matrices |
| $\mathbb{C}^{m \times n}$ | set of all $m \times n$ complex-valued matrices |
| x_i | i th entry of column vector \mathbf{x} |
| a_{ij} | (i, j) th entry of matrix \mathbf{A} |
| \mathbf{a}_i | i th column of matrix \mathbf{A} |
| \mathbf{a}_i^T | i th row of matrix \mathbf{A} |
| \mathbb{S}^n | set of all $n \times n$ real symmetric matrices, i.e., $\mathbf{A} \in \mathbb{R}^{n \times n}$ and $a_{ij} = a_{ji}$ for all i, j |
| \mathbb{H}^n | set of all $n \times n$ complex Hermitian matrices, i.e., $\mathbf{A} \in \mathbb{C}^{n \times n}$ and $\bar{a}_{ij} = a_{ji}$ for all i, j |
| \mathbf{A}^T | transpose of \mathbf{A} , i.e, $\mathbf{B} = \mathbf{A}^T$ means $b_{ji} = a_{ij}$ for all i, j |
| \mathbf{A}^H | Hermitian transpose of \mathbf{A} , i.e, $\mathbf{B} = \mathbf{A}^H$ means $b_{ji} = \bar{a}_{ij}$ for all i, j |
| $\text{trace}(\mathbf{A})$ | sum of diagonal entries of square matrix \mathbf{A} |
| $\mathbf{1}$ | A vector with all 1 entries |
| $\mathbf{0}$ | either a vector of all zeros, or a matrix of all zeros |
| \mathbf{e}_i | a unit vector with the nonzero element at the i th entry |
| $\mathcal{C}(\mathbf{A})$ | the column space of \mathbf{A} |
| $\mathcal{R}(\mathbf{A})$ | the row space of \mathbf{A} |
| $\mathcal{N}(\mathbf{A})$ | the null space of \mathbf{A} |
| $\text{Proj}_{\mathcal{M}}(\mathbf{A})$ | the projection of \mathbf{A} onto the set \mathcal{M} |

Chapter 1

Week1

1.1. Monday

1.1.1. Introduction to Optimizaiton

The usual optimization formulation is given by:

$$\begin{aligned} \min f(\mathbf{x}), \quad & \text{where } f: \mathbb{R}^n \mapsto \mathbb{R} \\ \text{such that } \mathbf{x} \in X \subseteq \mathbb{R}^n \end{aligned}$$

One example of the set X is given by:

$$X = \left\{ \mathbf{x} \in \mathbb{R}^n \left| \begin{array}{l} C_i(\mathbf{x}) = \mathbf{0}, i = 1, 2, \dots, m \leq n \\ h_i(\mathbf{x}) \geq \mathbf{0}, i = 1, 2, \dots, p \end{array} \right. \right\}$$

Linear programming can be easily solved, but Integer linear programming is much harder. The equivalent LP formulation is given by:

$$\begin{aligned} \min \quad & \mathbf{c}^T \mathbf{x} \\ \text{s.t.} \quad & \mathbf{Ax} = \mathbf{b} \\ & \mathbf{c} \leq \mathbf{Bx} \leq \mathbf{c}' \end{aligned}$$

1.2. Wednesday

1.2.1. Reviewing for Linear Algebra

Questions:

- What is the necessary and sufficient condition for the linear system $\mathbf{A}\mathbf{x} = \mathbf{b}$ to have a solution \mathbf{x} ?

Answer: $\mathbf{b} \in \mathcal{C}(\mathbf{A})$.

- For $\mathbf{A} \in \mathbb{S}^n$, what is the necessary and sufficient condition for $\mathbf{A} \succeq 0$?

Answer: $\mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$ for $\forall \mathbf{x} \in \mathbb{R}^n$; or $\lambda_i(\mathbf{A}) \geq 0$ for all i .

1.2.2. Reviewing for Calculus

For function $f : \mathbb{R}^n \mapsto \mathbb{R}$:

- We use notation $f \in \mathcal{C}^n$ to denote f is **continuously differentiable to n th order**. This course will basically deal with such functions.
- We use notation $\nabla f(x)$ to denote the **Gradient** of f at x ; and $\nabla^2 f(x)$ denotes the second order derivative of f at x . Note that $\nabla^2 f(x) \in \mathbb{S}^n$ for $f \in \mathcal{C}^1$.
- We use notation \mathbb{S}^n to denote the set of all symmetric $n \times n$ matrices, i.e.,

$$\mathbb{S}^n = \{\mathbf{X} \in \mathbb{R}^{n \times n} \mid \mathbf{X}^T = \mathbf{X}\}$$

Moreover, \mathbb{S}_+^n denotes the set of all symmetric $n \times n$ matrices with all eigenvalues non-negative:

$$\mathbb{S}_+^n = \{\mathbf{X} \in \mathbb{R}^{n \times n} \mid \mathbf{X}^T = \mathbf{X} \succeq 0\}$$

1.2.3. Introduction to Optimization

The usual optimization formulation is given by:

$$\begin{aligned} \min f(\mathbf{x}), \quad & \text{where } f: \mathbb{R}^n \mapsto \mathbb{R} \\ \text{such that } \mathbf{x} \in X \subseteq \mathbb{R}^n \end{aligned}$$

- The simplest case for the constraint is $X = \mathbb{R}^n$, which leads to **unconstrained** optimization problem.
- Or $X = P$ is a **polyhedron**, i.e., the boundaries for the region are all lines.

Definition 1.1 [Constraint Regions] In space \mathbb{R}^n ,

- the hyper-plane is defined as:

$$\{\mathbf{x} \mid \mathbf{a}^T \mathbf{x} = \beta\}$$

with constants $\mathbf{a} \in \mathbb{R}^n$ and $\beta \in \mathbb{R}$

- the half-space is defined as

$$\{\mathbf{x} \mid \mathbf{a}^T \mathbf{x} \leq \beta\}$$

- the polyhedron is defined as the **intersection** of a **finite** number of hyperplanes or half-spaces

Next, we give the definition for the basic optimization problem:

Definition 1.2 [Linear Programming] The Linear Programming is given by:

$$\begin{aligned} \min \quad & \mathbf{c}^T \mathbf{x}, \\ \text{such that } \mathbf{x} \in P(\text{polyhedron}) \end{aligned}$$

Or it can be reformulated as:

$$\begin{aligned} \min \quad & \mathbf{c}^T \mathbf{x}, \\ \text{such that} \quad & \mathbf{A}_I \mathbf{x} \leq \mathbf{b}_I \\ & \mathbf{A}_E \mathbf{x} = \mathbf{b}_E \in \mathbb{R}^m, \quad m < n. \end{aligned}$$

Definition 1.3 [Optimality] \mathbf{x}^* is said to be :

- the **local minimum** of $f(\mathbf{x})$ if there exists small ϵ such that

$$f(\mathbf{x}^*) \leq f(\mathbf{x}), \quad \forall \mathbf{x} \in \mathcal{B}(\mathbf{x}^*, \epsilon) \cap X := \{\mathbf{x} \mid \|\mathbf{x} - \mathbf{x}^*\| \leq \epsilon\} \cap X$$

- the **global minimum** if

$$f(\mathbf{x}^*) \leq f(\mathbf{x}), \quad \forall \mathbf{x} \in X$$

R Unless specified, when we want to minimize a non-convex function, it usually means we only find its **local minimum**. This is because usually the local minimum is good enough.

The optimization task is essentially find \mathbf{x}^* such that

$$\mathbf{x}^* = \arg \min_{\mathbf{x} \in X} f(\mathbf{x}) \in \mathbb{R}^n.$$

philosophy (optimization sufficient and necessity). philosophy of relaxation (convex nulls)

The Optimality conditions are the **most important** theoretical tools for optimization.

Theorem 1.1 — Optimality condition. The optimality condition contains

1. Necessary Condition (exclude non-optimal points):

$$n = 1 \text{ special case: } \begin{cases} \text{1st order: } f'(x) = 0 \\ \text{2nd order: } f''(x) \geq 0 \end{cases} \implies \begin{cases} \text{1st order: } \nabla f(x) = 0 \\ \text{2nd order: } \nabla^2 f(x) \succeq 0 \end{cases}$$

2. Sufficient Condition (may identify optimal solutions)

$$n = 1 \text{ special case: } \begin{cases} \text{1st order: } f'(x) = 0 \\ \text{2nd order: } f''(x) > 0 \end{cases} \implies \begin{cases} \text{1st order: } \nabla f(x) = 0 \\ \text{2nd order: } \nabla^2 f(x) \succ 0 \end{cases}$$

Proof. The $n = 1$ special case can imply the general case for optimality condition. For multivariate f , we set $\mathbf{x} = \mathbf{x}^* + td$ with t to be the stepsize and d to be the direction. For fixed t and d , we define $h(t) = f(\mathbf{x}) = f(\mathbf{x}^* + td)$. It follows that

$$h'(t) = \nabla^T f(\mathbf{x}^* + td)d$$

We find $h'(0) = \nabla^T f(\mathbf{x}^*)d$ for $\forall d$, which implies $\nabla f(\mathbf{x}^*) = 0$. ■

Note that there is a gap between necessary and sufficient conditions, which puts us in an embarrassing position. However, the convex condition can save us:

Theorem 1.2 If f is convex in \mathcal{C}^1 , then $\nabla f(\mathbf{x}) = 0$ is the **necessary** and **sufficient** condition.

Chapter 2

Week2

2.1. Monday

2.1.1. Reviewing and Announments

Tutorial: Thursday 7:00pm -9:00pm, ChengDao 208

Homework is due every Monday.

The first homework has been uploaded.

To proof the optimality condition in \mathbb{R}^n , we set $h(t) = f(x^* + td)$ for fixed x^* and d .

It follows that

$$h'(t) = \nabla^T f(x^* + td)d$$

and

$$h''(t) = d^T \nabla^2 f(x^* + td)d$$

By the optimality condition for \mathbb{R} , we derive the necessary condition:

$$\begin{cases} h'(0) = \nabla^T f(x^*)d = 0 \text{ for } \forall d \implies \nabla f(x^*) = 0; \\ h''(0) = d^T \nabla^2 f(x^*)d = 0 \text{ for } \forall d \implies \nabla^2 f(x^*) \succeq 0 \end{cases}$$

together with the sufficient condition:

$$\begin{cases} \nabla f(x^*) = 0; \\ \nabla^2 f(x^*) \succ 0 \end{cases}$$

2.1.2. Quadratic Function Case Study

Given a quadratic function

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{b}^T \mathbf{x}$$

w.l.o.g., assume the matrix \mathbf{Q} is symmetric (recall the quadratic section studied in linear algebra).

Definition 2.1 [Stationarity] A point \mathbf{x}^* is said to be the stationary point of $f(\mathbf{x})$ if $\nabla f(\mathbf{x}^*) = \mathbf{0}$. ■

To minimize such a function without constraint, we apply the optimality condition:

1. The first order optimality condition is given by:

$$\nabla f(\mathbf{x}) = \mathbf{Q} \mathbf{x} + \mathbf{b} = \mathbf{0}$$

The stationary point of the quadratic function $f(\mathbf{x})$ exists iff $\mathbf{b} \in \mathcal{C}(\mathbf{Q})$.

2. The second order necessary condition should be:

$$\nabla^2 f(\mathbf{x}) = \mathbf{Q} \succeq 0$$

For this special case, if $\mathbf{Q} \succeq 0$, then $f(\mathbf{x})$ is convex, the solutions to $\nabla f(\mathbf{x}) = \mathbf{0}$ are local minimum points. Furthermore, they are global minimum points (prove by Taylor Expansion). However, for general functions, we cannot obtain such good results.

Least Squares Problem. Such a problem has been well-studied in statistics given by:

$$\min_{\mathbf{x}} f(\mathbf{x}) := \frac{1}{2} \|\mathbf{A} \mathbf{x} - \mathbf{b}\|_2^2$$

The first order derivative of the minimizer should satisfy:

$$\nabla f(\mathbf{x}) = \mathbf{A}^T (\mathbf{A} \mathbf{x} - \mathbf{b})$$

Note that $\mathbf{A}^T \mathbf{b} \in \mathcal{C}(\mathbf{A}^T \mathbf{A})$, thus the least squares problem always has a solution. However, such a solution is not unique unless \mathbf{A} is full rank.

A Non-trivial Quadratic Function. To minimize the function

$$f(x, y) = \frac{1}{2}(\alpha x^2 + \beta y^2) - x$$

We take the first order derivative to be zero:

$$\nabla f(x, y) = \begin{bmatrix} \alpha x - 1 \\ \beta y \end{bmatrix} = \mathbf{0}$$

The second order derivative is given by:

$$\nabla^2 f(x, y) = \begin{bmatrix} \alpha & 0 \\ 0 & \beta \end{bmatrix}$$

The optimal solutions depend on the value of α and β : (although we haven't introduce the definition for convex formally)

- If $\alpha, \beta > 0$, then this problem is **strongly convex**. By the necessary and sufficient optimality condition for convex problem, we find that $(\frac{1}{\alpha}, 0)$ is the unique local minimum (It is also the global minimum by plotting the figure).
- If $\alpha = 0$, this problem has no solution. The objective value $f(x, y) \rightarrow -\infty$ as $x \rightarrow \infty$.
- If $\beta = 0, \alpha > 0$, this problem is convex. By the necessary and sufficient optimality condition for convex problem, $\{(\frac{1}{\alpha}, \xi) \mid \xi \in \mathbb{R}\}$ is the set of local minimum. (By plotting the graph, we find that such set is the set of global minimum points)
- For $\alpha > 0, \beta < 0$ case, this problem is non-convex. Actually, $f(x, y) \rightarrow -\infty$ as $y \rightarrow \infty$. Hence, this problem has no global minimum point.

A Non-trivial Function Study. To minimize the function

$$\begin{aligned} \min \quad & f(\mathbf{y}) = e^{y_1} + \cdots + e^{y_n} \\ \text{such that} \quad & y_1 + \cdots + y_n = S \end{aligned}$$

We can transform such a constrained optimization problem into unconstrained. Let $y_n = S - y_1 - \cdots - y_{n-1}$ and substitute it into the objective function, it suffices to solve

$$\min e^{y_1} + \cdots + e^{y_{n-1}} + e^{S-y_1-\cdots-y_{n-1}}$$

The stationary point should satisfy:

$$e^{y_i} = e^{S-y_1-\cdots-y_{n-1}}, \quad i = 1, 2, \dots, n-1$$


Or equivalently, $y_1 = y_2 = \cdots = y_{n-1} = y_n$. Hence we derive the unique stationary point:

$$y_1^* = y_2^* = \cdots = y_n^* = \frac{S}{n}$$

The value on the stationary point is $f(\mathbf{y}^*) = ne^{S/n}$. By checking the second order sufficient optimality condition,

$$\frac{f}{\partial y_i \partial y_j} = \begin{cases} e^{y_i} + e^{S-y_1-\cdots-y_{n-1}} & i = j \\ e^{S-y_1-\cdots-y_{n-1}} & i \neq j \end{cases} \implies \nabla^2 f = e^{S-y_1-\cdots-y_{n-1}} \mathbf{E} + \text{diag}(e^{y_1}, \dots, e^{y_{n-1}})$$

where \mathbf{E} is a matrix with entries all ones. Thus $\nabla^2 f \succ 0$ for any stationary point. By the second order sufficient optimality condition, this stationary point is local minimum. Actually, for this special problem, this unique local minimum point is the global minimum.

-  In this problem, we find that this stationary point is the unique local minimum point, but the unique local minimum point is not necessarily the global minimum point, unless the function is **coercive** or the feasible region is compact. Here is the counter-example: $f(x) = x^2 - x^4$. We will discuss the definition for coercive in the future.

2.2. Wednesday

2.2.1. Convex Analysis

This lecture will study the convex analysis.

Definition 2.2 [Convex] The subset $\mathcal{C} \subseteq \mathbb{R}^n$ is convex if

$$x, y \in \mathcal{C} \implies \{\lambda x + (1 - \lambda)y \mid \lambda \in [0, 1]\} \subset \mathcal{C},$$

i.e., the line segment between arbitrarily two elements lies in \mathcal{C} ■

R Intersections of convex sets are convex. Empty set is assumed to be convex.

Definition 2.3 [Convex] The function $f : \mathbb{R}^n \mapsto \mathbb{R}$ is convex if $\text{dom } f$ is convex and

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$$

for $\forall x, y \in \text{dom } f$ and $\forall \lambda \in [0, 1]$, i.e., the function evaluated in the line segment is lower than secant line between x and y (f lies below secant line). ■

R

- f is convex iff $-f$ is concave. (The concave definition simply changes the inequality direction in Def.(2.3))
- Affines are both convex and concave.
- The convexity depends on the domain of the function.

For a second order differentiable function, we have a much easier way to determine its convexity.

Theorem 2.1 If $f \in \mathcal{C}^1$, then the followings are equivalent:

1. f is convex

2. $f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla^T f(\mathbf{x})(\mathbf{y} - \mathbf{x})$ for $\forall \mathbf{x}, \mathbf{y} \in \text{dom } f$, i.e., f lies above the tangent line.

Proof. 1. From the definition for convexity,

$$f(\mathbf{y}) - f(\mathbf{x}) \geq \frac{f(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) - f(\mathbf{x})}{1 - \lambda}$$

Letting $\lambda \rightarrow 1$, the RHS becomes a direction derivative:

$$f(\mathbf{y}) - f(\mathbf{x}) \geq \nabla^T f(\mathbf{x})(\mathbf{y} - \mathbf{x})$$

2. To show the converse, we let $\mathbf{z} = \lambda \mathbf{x} + (1 - \lambda)\mathbf{y}$. By applying the inequality in (2.1) twice, we have

$$f(\mathbf{x}) \geq f(\mathbf{z}) + \nabla^T f(\mathbf{z})(\mathbf{x} - \mathbf{z}) \quad (2.1)$$

$$f(\mathbf{y}) \geq f(\mathbf{z}) + \nabla^T f(\mathbf{z})(\mathbf{y} - \mathbf{z}) \quad (2.2)$$

Letting Eq.(2.1) times λ add Eq.(2.2) times $(1 - \lambda)$, we derive that f is convex. ■

Theorem 2.2 If $f \in \mathcal{C}^2$, then the followings are equivalent:

1. f is convex
2. $\nabla^2 f(\mathbf{x}) \succeq 0$ for $\forall \mathbf{x} \in \text{dom } f$.

Proof. We rewrite $f(\mathbf{y})$ by applying Taylor expansion:

$$f(\mathbf{y}) = f(\mathbf{x}) + \nabla^T f(\mathbf{x})(\mathbf{y} - \mathbf{x}) + \frac{1}{2}(\mathbf{y} - \mathbf{x})^T \nabla^2 f(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))(\mathbf{y} - \mathbf{x}), \quad (2.3)$$

for some $t \in [0, 1]$.

1. If f is convex, from Theorem(2.1) and Eq.(2.3), we derive

$$(\mathbf{y} - \mathbf{x})^T \nabla^2 f(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))(\mathbf{y} - \mathbf{x}) \geq 0 \implies \frac{(\mathbf{y} - \mathbf{x})^T}{\|\mathbf{y} - \mathbf{x}\|} \nabla^2 f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) \frac{(\mathbf{y} - \mathbf{x})}{\|\mathbf{y} - \mathbf{x}\|} \geq 0$$

Set $d := \frac{(\mathbf{y}-\mathbf{x})}{\|\mathbf{y}-\mathbf{x}\|}$ and let $\mathbf{y} \rightarrow \mathbf{x}$, we derive

$$\mathbf{d}^T \nabla^2 f(\mathbf{x}) \mathbf{d} \geq 0,$$

which implies $\nabla^2 f(\mathbf{x}) \succeq 0$ since \mathbf{d} could have an arbitrary direction.

2. To show the converse, due to the semidefiniteness of $\nabla^2 f(\mathbf{x})$, we obtain a new inequality from Eq.(2.3):

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla^T f(\mathbf{x})(\mathbf{y} - \mathbf{x})$$

From Theorem(2.1) we imply f is convex. ■

Definition 2.4 [Epigraph] The Epigraph of f is given by:

$$\text{Epi}(f) := \{(x, t) \in \mathbb{R}^n \times \mathbb{R} \mid x \in \text{dom } f, t \geq f(x)\} \subseteq \mathbb{R}^{n+1}$$

Theorem 2.3 f is convex iff $\text{Epi}(f)$ is convex.

Proof. 1. Suppose f is convex. For any $(x, t), (y, s) \in \text{Epi}(f)$, it suffices to show

$$(\lambda x + (1 - \lambda)y, \lambda t + (1 - \lambda)s) \in \text{Epi}(f) \iff \lambda t + (1 - \lambda)s \geq f(\lambda x + (1 - \lambda)y).$$

The convexity of f implies that

$$\begin{aligned} f(\lambda x + (1 - \lambda)y) &\leq \lambda f(x) + (1 - \lambda)f(y) \\ &\leq \lambda t + (1 - \lambda)s. \end{aligned}$$

2. The reverse direction is obvious by applying definitions. ■

Definition 2.5 [Strict Convex] The function $f : \mathbb{R}^n \mapsto \mathbb{R}$ is strict convex if $\text{dom } f$ is convex and

$$f(\lambda x + (1 - \lambda)y) < \lambda f(x) + (1 - \lambda)f(y)$$

for $\forall x \neq y, x, y \in \text{dom } f$ and $\forall \lambda \in (0, 1)$ ■

R Strict convex implies the uniqueness of minimum

However, for function $f(x) = \frac{1}{x}$, the curvature becomes more and more flat. We want to exclude such kind of functions.

Definition 2.6 The function $f : \mathbb{R}^n \mapsto \mathbb{R}$ is said to be strongly convex if $\text{dom } f$ is convex and $\exists \alpha > 0$ such that $f(\mathbf{x}) - \alpha \mathbf{x}^T \mathbf{x}$ is convex; or equivalently,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla^T f(\mathbf{x})(\mathbf{y} - \mathbf{x}) + \frac{\alpha}{2} \|\mathbf{y} - \mathbf{x}\|^2$$

R The strong convexity places a quadratic lower bound in the curvature of the function, i.e., the function must rise up at least as fast as a quadratic function. How fast it rises depends on the parameter α .

The convexity properties are extremely useful in forcing optimization algorithms to rapidly converge to optima. However, most functions are not convex. The most important result that requires convexity is given below:

Theorem 2.4 If f is convex in \mathcal{C}^1 , then $\nabla f(\mathbf{x}) = 0$ is the **necessary** and **sufficient** condition for **global** minimum.

Note that convex function does not have a local minimum that is not global minimum.

Proof. If $f \in \mathcal{C}^1$ is convex, recall the Theorem(2.1) that

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla^T f(\mathbf{x})(\mathbf{y} - \mathbf{x}) \tag{2.4}$$

1. If $\nabla f(\mathbf{x}) = \mathbf{0}$, then Eq.(2.3) implies $f(\mathbf{y}) \geq f(\mathbf{x})$ for $\forall \mathbf{y}$.
2. If \mathbf{x} is the global minimum, recall the optimality condition, $\nabla f(\mathbf{x}) = \mathbf{0}$.

■

In practice, we cannot solve all convex optimization problems. So we need to carefully study the structure of every problem we have faced.

Chapter 3

Week3

3.1. Wednesday

Assignment 2 posted.

CIE6010: Exercise 1.2.9 and 1.3.9; together with MATLAB project.

3.1.1. Convex Analysis

Last time we have shown that for a unconstrained problem, $\nabla f(\mathbf{x}) = 0$ is the necessary and sufficient condition for global minimum ensurance. However, the case for constrained problem will be different.

Proposition 3.1 For the **constrained** problem

$$\begin{aligned} \min \quad & f(x) \\ & \mathbf{x} \in X \subseteq \mathbb{R} \\ & f \text{ is convex in } \mathcal{C}^1 \end{aligned}$$

\mathbf{x} is a global minimum iff

$$\langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \geq 0$$

for $\forall \mathbf{y} \in X$.

Proof. Since f is convex, the inequality below holds:

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle, \forall \mathbf{x}, \mathbf{y} \in X$$

Note that \mathbf{x} is a global minimum iff $f(\mathbf{y}) \geq f(\mathbf{x})$, $\forall \mathbf{y} \in X$. Combining the inequality above, the proof is complete. ■

R

- Such a condition is not so useful unless \mathbf{y} lies in the whole space, at that time we have no choice but $\nabla f(\mathbf{x}) = \mathbf{0}$. (otherwise we can construct a \mathbf{y} to let the inner product negative.)
- An equivalent version of the condition is that every **feasible** direction is **ascending**.

Definition 3.1 [Descending Direction] The vector $\mathbf{d} \in \mathbb{R}^n$ is said to be a **descending direction** of f at \mathbf{x} if

$$\langle \nabla f(\mathbf{x}), \mathbf{d} \rangle < 0.$$

This definition is the motivation of descent method.

3.1.2. Iterative Method

Definition 3.2 [Descent Method] At any non-stationary \mathbf{x} , i.e., $\nabla f(\mathbf{x}) \neq \mathbf{0}$, we find the descending direction \mathbf{d} , i.e., $\langle \nabla f(\mathbf{x}), \mathbf{d} \rangle < 0$. We update our old \mathbf{x} as:

$$\mathbf{x}^{r+1} \leftarrow \mathbf{x}^r + \alpha^r \mathbf{d}^r, \quad \alpha > 0.$$

The key is how to choose \mathbf{d} and α . We have a general formula for \mathbf{d} :

$$\mathbf{d} = -\mathbf{D} \cdot \nabla f(\mathbf{x}),$$

where $\mathbf{D} \in \mathbb{S}^n$ and $\mathbf{D} \succ 0$. (Verify by yourself that \mathbf{d} satisfies the descending direction definition)

1. $D = I$ implies gradient method (Steepest Descent).
2. $D = (\nabla^2 f(\mathbf{x}))^{-1}$ implies the Newton's method.

Nonlinear LS. The optimization problem is

$$\begin{aligned} \min \quad & f(\mathbf{x}) = \frac{1}{2} \sum_{i=1}^m g_i^2(\mathbf{x}) := \frac{1}{2} \|g(\mathbf{x})\|_2^2 \\ & g(\mathbf{x}) = \begin{pmatrix} g_1(\mathbf{x}) & g_2(\mathbf{x}) & \cdots & g_m(\mathbf{x}) \end{pmatrix}^T \end{aligned}$$

The gradient function is

$$\begin{aligned} \nabla f(\mathbf{x}) &= \sum_{i=1}^m g_i(\mathbf{x}) \nabla g_i(\mathbf{x}) \\ &= \underbrace{\begin{bmatrix} \nabla g_1(\mathbf{x}) & \cdots & \nabla g_m(\mathbf{x}) \end{bmatrix}}_{\nabla g(\mathbf{x}) \in \mathbb{R}^{n \times m}} \begin{bmatrix} g_1(\mathbf{x}) \\ \vdots \\ g_m(\mathbf{x}) \end{bmatrix} \\ &= \nabla g(\mathbf{x}) \cdot g(\mathbf{x}) \\ &= \langle J(\mathbf{x}), g(\mathbf{x}) \rangle, \end{aligned}$$

where $J(\mathbf{x}) \in \mathbb{R}^{m \times n}$ is said to be the Jacobian matrix of $g(\mathbf{x})$.

The second order derivative function is given as: (complete the calculation process by yourself)

$$\nabla^2 f(\mathbf{x}) = J^T(\mathbf{x})J(\mathbf{x}) + \sum_{i=1}^m g_i(\mathbf{x}) \nabla^2 g_i(\mathbf{x}),$$

the second term in RHS is complicated and hard to compute. To solve this LS problem, the Gauss-Newton method directly ignore it, which leads to the descent direction

$$\mathbf{d} = -(J^T J)^{-1} J^T g(\mathbf{x})$$

Choice of Step Length α . We apply the Limited Minimization Rule to find α , i.e., for fixed $s > 0$, choose α^r such that

$$\min_{\alpha^r \in (0, s]} f(\mathbf{x}^r + \alpha^r \mathbf{d}^r).$$

Usually this rule is too computationally expensive. The alternative ways are:

- Choose α just as a constant
- Choose $\alpha^r \rightarrow 0$ as $r \rightarrow \infty$ but also satisfies the infinite travel condition

$$\sum_{r=0}^{\infty} \alpha^r = \infty$$

Adding Lipschitz condition will make the choice of step-length easier:

Definition 3.3 [Lipschitz Continuous] ∇f is **Lipschitz continuous** with Lipschitz constant L if

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

for all \mathbf{x}, \mathbf{y} . ■



- It is useful to note that convexity places a lower bound on the growth of the function at every point; whereas Lipschitzness places an upper bound on the growth of the function that is linear in the perturbation i.e., $\|\mathbf{x} - \mathbf{y}\|_2$. Also note that Lipschitz functions need not be differentiable. However, differentiable functions with bounded gradients are always Lipschitz.
- The Lipschitz condition induces that for iterative method we have

$$f(\mathbf{x}^r) - f(\mathbf{x}^{r+1}) \geq \frac{L}{2} \|\nabla f(\mathbf{x}^r)\|^2.$$

From this inequality, we imply that the result of iterative convergence is $\nabla f(\mathbf{x}^r) \rightarrow 0$, but the minimum point is still un-guaranteed. In Deep Learning people often train the data using this way, which is not so rigorous.

Convergence Rate Analysis. We apply the Lipschitzness to analysis the rate of convergence first. Setting $h(t) = f(\mathbf{x} + t\alpha\mathbf{d})$, we find that

$$\begin{aligned}
f(\mathbf{x} + \alpha\mathbf{d}) - f(\mathbf{x}) &= h(1) - h(0) = \int_0^1 h'(t) dt \\
&= \int_0^1 \langle \nabla f(\mathbf{x} + t \cdot \alpha\mathbf{d}), \alpha\mathbf{d} \rangle dt \\
&= \int_0^1 [\langle \nabla f(\mathbf{x} + t \cdot \alpha\mathbf{d}), \alpha\mathbf{d} \rangle - \langle \nabla f(\mathbf{x}), \alpha\mathbf{d} \rangle + \langle \nabla f(\mathbf{x}), \alpha\mathbf{d} \rangle] dt \\
&= \langle \nabla f(\mathbf{x}), \alpha\mathbf{d} \rangle + \int_0^1 \langle \nabla f(\mathbf{x} + t \cdot \alpha\mathbf{d}) - \nabla f(\mathbf{x}), \alpha\mathbf{d} \rangle dt \\
&\leq \langle \nabla f(\mathbf{x}), \alpha\mathbf{d} \rangle + \int_0^1 \|\nabla f(\mathbf{x} + t \cdot \alpha\mathbf{d}) - \nabla f(\mathbf{x})\| \cdot \|\alpha\mathbf{d}\| dt \\
&\leq \langle \nabla f(\mathbf{x}), \alpha\mathbf{d} \rangle + L \int_0^1 t\alpha^2 \|\mathbf{d}\|^2 dt \\
&= \underbrace{\langle \nabla f(\mathbf{x}), \alpha\mathbf{d} \rangle}_{\text{negative}} + \frac{L\alpha^2 \|\mathbf{d}\|^2}{2}
\end{aligned}$$

Choice of Step Length. To get the optimal step length α , differentiating the RHS w.r.t. α leads to

$$\langle \nabla f(\mathbf{x}), \alpha\mathbf{d} \rangle + L\alpha \|\mathbf{d}\|^2 = 0 \implies \alpha = -\frac{\langle \nabla f(\mathbf{x}), \alpha\mathbf{d} \rangle}{L\|\mathbf{d}\|^2} > 0,$$

which seems a reasonable choice. If \mathbf{d} is the steepest descent direction, the step-length becomes a constant:

$$\alpha = \frac{1}{L}.$$

3.2. Thursday

3.2.1. Announcement

The assignment 2 requires to do a MATLAB project. The grade usually depends on your understanding of the reading materials and the time spent on experimentation.

3.2.2. Sparse Large Scale Optimization

Given an underlying signal $\mathbf{x} \in \mathbb{R}^n$ satisfying the undermined system $\mathbf{Ax} = \mathbf{b}$, we aim to recover the desired $\hat{\mathbf{x}}$. It suffices to solve the optimization problem

$$\begin{aligned} \min \quad & \|\mathbf{D}\mathbf{x}\|_1 \\ \text{s.t.} \quad & \mathbf{Ax} = \mathbf{b} \\ & \mathbf{A} \in \mathbb{R}^{m \times n}, m < n \end{aligned}$$

with \mathbf{D} to be the difference matrix and $\min \|\mathbf{D}\mathbf{x}\|_1$ is sparsity promoting. Here we list two basic but effective ways to solve such a problem.

Linear Programming Approach. One way is to reformulate the problem into LP.

1. Define new variables $t_i = |(\mathbf{D}\mathbf{x})_i|$, we can reformulate the origin problem as:

$$\begin{aligned} \min \quad & \sum_{i=1}^n t_i \\ \text{s.t.} \quad & \mathbf{Ax} = \mathbf{b} \\ & -t_i \leq \sum_{k=1}^n d_{ik}x_k \leq t_i \\ & \mathbf{A} \in \mathbb{R}^{m \times n}, m < n \end{aligned}$$

2. Alternatively, recall what we have learnt in MAT3007. Define slack variables

$(\mathbf{D}\mathbf{x})_i = u_i - v_i$, where $u_i := (\mathbf{D}\mathbf{x})_i^+$, $v_i = (\mathbf{D}\mathbf{x})_i^-$. It suffices to solve

$$\begin{aligned} \min \quad & \sum_{i=1}^n (u_i + v_i) \\ \text{s.t.} \quad & \mathbf{A}\mathbf{x} = \mathbf{b} \\ & -t_i \leq \sum_{k=1}^n d_{ik}x_k \leq t_i \\ & \mathbf{A} \in \mathbb{R}^{m \times n}, m < n \\ & u_i, v_i \geq 0 \end{aligned}$$

However, linear programming is not the optimal way to solve large-scale problem.

Gredient-Based Approach. We can also transform it into the unconstraint minimization problem, i.e., we add the penalty for the constraint $\mathbf{A}\mathbf{x} - \mathbf{b} = \mathbf{0}$:

$$\min \|\mathbf{D}\mathbf{x}\|_1 + \frac{\mu}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2$$

You may see that this reformulation is not exactly equivalent to the origin problem. However, it is not meaningful to stress $\mathbf{A}\mathbf{x}$ should exactly equal to \mathbf{b} , as there exists some noise perturbing the equality in nature.

Another problem is that this objective function is not differentiable once there is at least zero entry from $\mathbf{D}\mathbf{x}$. Thus we do the approximation

$$|t| \approx \sqrt{t^2 + \sigma}, \text{ for small } \sigma > 0.$$

Hence, it suffices to solve

$$\min f(x) := \Theta_\sigma(\mathbf{D}\mathbf{x}) + \frac{\mu}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2 \quad (3.1)$$

where

$$\Theta_\sigma(\mathbf{y}) = \sum_{i=1}^n \sqrt{y_i^2 + \sigma}$$

Descent Direction. Since problem(3.1) is convex, taking the derivative leads to minimum point. Hence we use the gredient descent method, i.e., $\mathbf{d} = -\nabla f(\mathbf{x})$.

Although this direction is not optimal (trying another direction may be faster after

several iterations), let's assume we are short-sighted such that we just want to take the steepest direction.

Hence the iterative algorithm to solve this problem can be summarized into one formula: Take a initial guess \mathbf{x}^0 , then for $r = 0, 1, 2 \dots$

$$\mathbf{x}^{r+1} = \mathbf{x}^r - \alpha^r \nabla f(\mathbf{x}^r)$$

Stopping Criteria. The stopping criteria has two conditions, either one is satisfied is ok. Always keep mind of scaling for stopping criteria, i.e., how large of an objective should depend on the scale of the problem.

- First is $\|\nabla f(\mathbf{x}^k)\| \leq 10^{-2} \|\nabla f(\mathbf{x}^0)\|$, i.e., the iterative method converge to the near stationary point
- Another is $|f(\mathbf{x}^k) - f(\mathbf{x}^{k+1})| \leq 10^{-8} |f(\mathbf{x}^k)|$, i.e., the function does not change too much.

The next questions turn out that how to choose initial guess? How to choose step-length? Is steepest descent usually effective?

1. For large-scale optimization, the steepest descent is usually one of the **best** way among iterative methods.
2. To choose the initial guess, sometimes we choose the LS solution, i.e., enter the matlab command $A' A / A' b$.
3. Last lecture we tell that we can choose step-length to be the Lipschitz constant, but the disadvantage is that the constant L for large scale optimization is too small. We have a better alternative.

Armijo Condition and BB Step. The motivation is that we aim to let

$$f(\mathbf{x}^k + \alpha \mathbf{d}^k) \leq f(\mathbf{x}^k) + C_1 \alpha \langle \nabla f(\mathbf{x}), \mathbf{d}^k \rangle, (0 < C_1 < 1) \quad (3.2)$$

i.e., our updated function value should be at least less than the old function minus the descent decrease gain, i.e., it should sufficiently decrease faster than a constant times

the steepest gradient descent.

This method has a geometrically meaning:

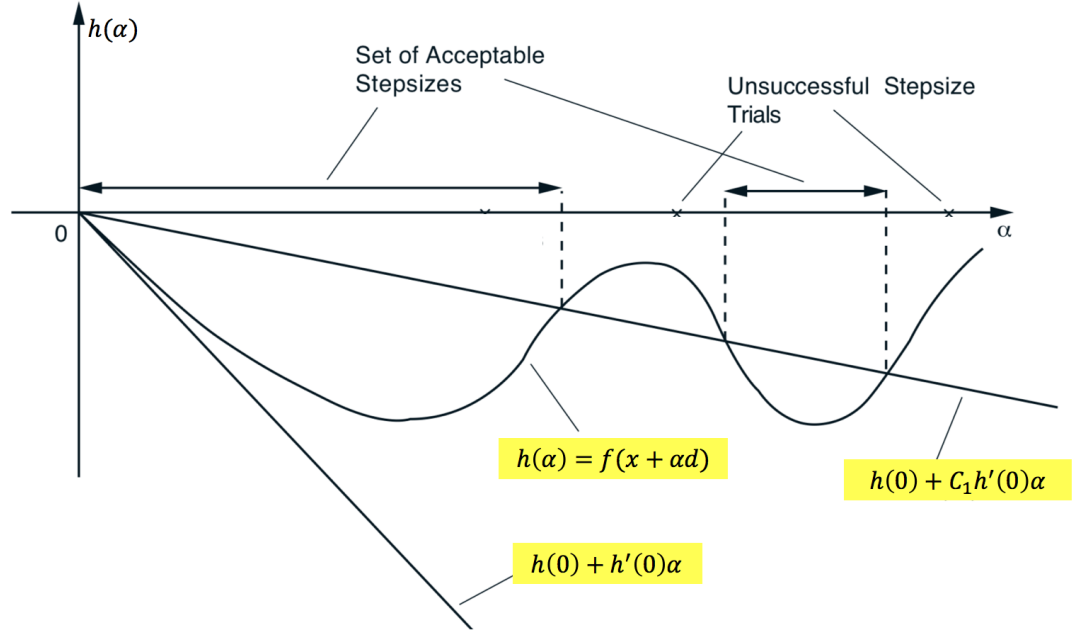


Figure 3.1: Geometric Interpretation of Armijo Condition

We set $h(\alpha) := f(\mathbf{x} + \alpha \mathbf{d})$, then $h'(0) = \langle \nabla f(\mathbf{x} + \alpha \mathbf{d}), \mathbf{d} \rangle$, thus the tangent line at $h(0)$ is given by:

$$h(0) + \alpha h'(0) := f(\mathbf{x}) + \alpha \langle \nabla f(\mathbf{x} + \alpha \mathbf{d}), \mathbf{d} \rangle$$

Geometrically we can see that no α can be chosen such that the updated function value $h(\alpha)$ is less than this tangent line. Hence we make the tangent line more flat, i.e., we want to find α such that the updated function value $h(\alpha)$ is below the line $h(0) + C_1 h'(0) \alpha$, $0 < C_1 < 1$:

$$h(\alpha) \leq h(0) + C_1 h'(0) \alpha$$

How to choose such α ? Take a initial long step-length $\bar{\alpha}$ first, if condition(3.2) is not satisfied, try step length $\beta \bar{\alpha}, \beta^2 \bar{\alpha}, \dots$ respectively. (Take a big step, if not satisfied, shorten the step.)

R However, it is not suggested to do that. Although it is mathematically true,

during the computer run, the step-length will decrease exponentially.

How to choose C_1 ? Empirically, $C_1 = 10^{-3}$ or 10^{-4} , i.e., it is very flat.

How to choose initial $\bar{\alpha}$? It depends on the scale of functionm which requires for your reading of materials.

How to choose the value of β ? Do the experiment. (0.5, 0.8 for example).