
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

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.

R 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

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 $\nabla^T f(\mathbf{x})(\mathbf{y} - \mathbf{x}) \geq 0$ for $\forall \mathbf{y} \in X$

Such a condition is not useful unless \mathbf{y} lies in the whole space, at that time we have no choice but $\nabla f(\mathbf{x}) = \mathbf{0}$.

An equivalent version of the condition is that every **feasible** direction is **ascending**.

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

$$\nabla^T f(\mathbf{x})\mathbf{d} < 0.$$

This definition is the motivation of descent method.

3.1.2. Iterative Method

Definition 3.2 [Descent Method] At ant non-stationary \mathbf{x} , i.e., $\nabla f(\mathbf{x}) \neq \mathbf{0}$, we find descent \mathbf{d} , i.e., $\nabla^T f(\mathbf{x})\mathbf{d} < 0$. 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} , which is the **descent direction**:

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

where $\mathbf{D} \in \mathbb{S}^n$ and $\mathbf{D} \succ 0$.

e.g., $\mathbf{D} = \mathbf{I}$ implies gradient method (Steepest Descent).

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

Nonlinear LS. The optimization problem is

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

The gradient function is

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

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

The second order derivative is given as:

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

the second term in RHS is complicated and hard to compute. The Gauss-Newton

method directly ignore it.

Choice of Step Length α . Exact line search

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

usually it is too expensive. Sometimes we choose α as a constant.

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} . ■

R This condition guarantees that

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

From this inequality, we imply that the result of 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.

We set $h(t) = f(\mathbf{x} + t\alpha\mathbf{d})$, then we can use Lipschitz continuous to analysis the rate of convergence and the choice of step-length:

$$\begin{aligned} f(\mathbf{x} + \alpha\mathbf{d}) - f(\mathbf{x}) &= \int_0^1 \nabla^T f(\mathbf{x} + t \cdot \alpha\mathbf{d}) \alpha\mathbf{d} \, dt \\ &= \int_0^1 \left[\nabla^T f(\mathbf{x} + t \cdot \alpha\mathbf{d}) \alpha\mathbf{d} - \nabla^T f(\mathbf{x}) \alpha\mathbf{d} + \nabla^T f(\mathbf{x}) \alpha\mathbf{d} \right] \\ &\leq \nabla^T f(\mathbf{x}) \alpha\mathbf{d} + \int_0^1 \|\nabla f(\mathbf{x} + \alpha\mathbf{d}) - \nabla f(\mathbf{x})\| \cdot \|\alpha\mathbf{d}\| \, dt \\ &\leq \nabla^T f(\mathbf{x}) \alpha\mathbf{d} + L \int_0^1 t \alpha^2 \|\mathbf{d}\|^2 \, dt \\ &= \underbrace{\nabla^T f(\mathbf{x}) \alpha\mathbf{d}}_{\text{negative}} + \frac{L\alpha^2 \|\mathbf{d}\|^2}{2} \end{aligned}$$

Differentiate the RHS w.r.t. α leads to

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

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

$$\alpha = \frac{1}{L} \text{ (constant)}$$