

DATA8004: Optimization for Statistical Learning

Subject Lecturer: Man-Chung YUE

Lecture 1

Overview

Preliminary materials

What is optimization?

- Optimization a mathematical subject that studies techniques for finding “best” solutions/decisions.
- An optimization problem takes the form of minimizing (or maximizing) an objective function subject to constraints:

$$\begin{array}{ll}\text{Minimize} & f(x) \\ \text{subject to} & x \in \Omega.\end{array}$$

Here:

- ★ $x \in \mathbb{R}^n$ is called decision variables.
- ★ f is called the objective function.
- ★ $\Omega \subseteq \mathbb{R}^n$ is called the constraint set / feasible set / feasible region.

Example: Objectives

Objective functions:

- Linear: $f(x) = c^T x$ for some $c \in \mathbb{R}^n$.
- Quadratic: $f(x) = \frac{1}{2}x^T Gx + c^T x + \beta$ for some $c \in \mathbb{R}^n$, $\beta \in \mathbb{R}$, and **symmetric matrix** $G \in \mathbb{R}^{n \times n}$.

Example: Objectives

Objective functions:

- Linear: $f(x) = c^T x$ for some $c \in \mathbb{R}^n$.
- Quadratic: $f(x) = \frac{1}{2}x^T Gx + c^T x + \beta$ for some $c \in \mathbb{R}^n$, $\beta \in \mathbb{R}$, and **symmetric matrix** $G \in \mathbb{R}^{n \times n}$.

Note:

- The **symmetry** of G can be assume **without loss of generality**. Indeed, if $H \in \mathbb{R}^{n \times n}$ and $x \in \mathbb{R}^n$, then

$$x^T Hx = x^T H^T x = \frac{1}{2}x^T (H^T + H)x.$$

- For $f(x) = \frac{1}{2}x^T Gx + c^T x + \beta$ with $G \in \mathbb{R}^{n \times n}$ being **symmetric**, it holds that

$$\nabla f(x) = Gx + c.$$

Example: Constraints

The feasible set Ω can be specified by one or more of the following constraints.

- Equality constraints:

- ★ $x_1^2 + x_2^2 + \cdots + x_n^2 = 1$ (sphere).
- ★ $x_1 + x_2 + \cdots + x_n = 1$ (hyperplane).
- ★ In \mathbb{R}^3 : $x_3 = x_1^2 + x_2^2$ (paraboloid).

Example: Constraints

The feasible set Ω can be specified by one or more of the following constraints.

- Equality constraints:

- ★ $x_1^2 + x_2^2 + \cdots + x_n^2 = 1$ (sphere).
- ★ $x_1 + x_2 + \cdots + x_n = 1$ (hyperplane).
- ★ In \mathbb{R}^3 : $x_3 = x_1^2 + x_2^2$ (paraboloid).

- Inequality constraints:

- ★ $x_1^2 + x_2^2 + \cdots + x_n^2 \leq 1$ (ball).
- ★ $x_1 + x_2 + \cdots + x_n \leq 0$ (half-space).
- ★ In \mathbb{R}^3 : $x_3 \geq x_1^2 + x_2^2$.

Example: Constraints

The feasible set Ω can be specified by one or more of the following constraints.

- **Equality constraints:**

- ★ $x_1^2 + x_2^2 + \cdots + x_n^2 = 1$ (sphere).
- ★ $x_1 + x_2 + \cdots + x_n = 1$ (hyperplane).
- ★ In \mathbb{R}^3 : $x_3 = x_1^2 + x_2^2$ (paraboloid).

- **Inequality constraints:**

- ★ $x_1^2 + x_2^2 + \cdots + x_n^2 \leq 1$ (ball).
- ★ $x_1 + x_2 + \cdots + x_n \leq 0$ (half-space).
- ★ In \mathbb{R}^3 : $x_3 \geq x_1^2 + x_2^2$.

- **Box constraint:** $\ell \leq x \leq u$, where $\ell \in \mathbb{R}^n$ and $u \in \mathbb{R}^n$. This means

$$\ell_i \leq x_i \leq u_i \quad \forall i.$$

- The problem is said to be **unconstrained** if $\Omega = \mathbb{R}^n$.

Infimum

Definition: Let $S \subseteq \mathbb{R}$ be a nonempty set. We say that $\ell \in [-\infty, \infty)$ is the **infimum** of S if

- $s \geq \ell$ for every $s \in S$; and
- for every $\zeta > \ell$ ($\zeta \in \mathbb{R}$), one can find $s \in S$ so that $s < \zeta$.

Notation: $\ell = \inf S$. By convention, $\inf \emptyset = \infty$.

Infimum

Definition: Let $S \subseteq \mathbb{R}$ be a nonempty set. We say that $\ell \in [-\infty, \infty)$ is the **infimum** of S if

- $s \geq \ell$ for every $s \in S$; and
- for every $\zeta > \ell$ ($\zeta \in \mathbb{R}$), one can find $s \in S$ so that $s < \zeta$.

Notation: $\ell = \inf S$. By convention, $\inf \emptyset = \infty$.

Note:

- The existence and uniqueness of $\inf S$ follow from the **completeness of \mathbb{R}** .

Infimum

Definition: Let $S \subseteq \mathbb{R}$ be a nonempty set. We say that $\ell \in [-\infty, \infty)$ is the **infimum** of S if

- $s \geq \ell$ for every $s \in S$; and
- for every $\zeta > \ell$ ($\zeta \in \mathbb{R}$), one can find $s \in S$ so that $s < \zeta$.

Notation: $\ell = \inf S$. By convention, $\inf \emptyset = \infty$.

Note:

- The existence and uniqueness of $\inf S$ follow from the **completeness of \mathbb{R}** .
- Roughly speaking, $\ell = \inf S$ is the **largest number** that is smaller than everything in S . However, it is **not necessary that $\ell \in S$** ! e.g., $\inf\{e^{-x} : x \in \mathbb{R}\} = 0$, but there is no $a \in \mathbb{R}$ so that $e^{-a} = 0$.

Infimum

Definition: Let $S \subseteq \mathbb{R}$ be a nonempty set. We say that $\ell \in [-\infty, \infty)$ is the **infimum** of S if

- $s \geq \ell$ for every $s \in S$; and
- for every $\zeta > \ell$ ($\zeta \in \mathbb{R}$), one can find $s \in S$ so that $s < \zeta$.

Notation: $\ell = \inf S$. By convention, $\inf \emptyset = \infty$.

Note:

- The existence and uniqueness of $\inf S$ follow from the **completeness of \mathbb{R}** .
- Roughly speaking, $\ell = \inf S$ is the **largest number** that is smaller than everything in S . However, it is **not necessary that $\ell \in S$** !
e.g., $\inf\{e^{-x} : x \in \mathbb{R}\} = 0$, but there is no $a \in \mathbb{R}$ so that $e^{-a} = 0$.
- For optimization problems, we refer to the infimum of the set $\{f(x) : x \in \Omega\}$ as the **optimal value**.
e.g., “Minimize e^{-x} subject to $x \in \mathbb{R}$ ” has optimal value 0.

Norm

Definition: A function $\|\cdot\| : \mathbb{R}^n \rightarrow \mathbb{R}$ is called a (vector) **norm** if

- $\|x\| \geq 0$ for all $x \in \mathbb{R}^n$.
- $\|x\| = 0$ if and only if $x = 0$.
- $\|\alpha x\| = |\alpha| \|x\|$ for any $\alpha \in \mathbb{R}$ and $x \in \mathbb{R}^n$.
- $\|x + y\| \leq \|x\| + \|y\|$ for any $x, y \in \mathbb{R}^n$.

Norm

Definition: A function $\| \cdot \| : \mathbb{R}^n \rightarrow \mathbb{R}$ is called a (vector) **norm** if

- $\|x\| \geq 0$ for all $x \in \mathbb{R}^n$.
- $\|x\| = 0$ if and only if $x = 0$.
- $\|\alpha x\| = |\alpha| \|x\|$ for any $\alpha \in \mathbb{R}$ and $x \in \mathbb{R}^n$.
- $\|x + y\| \leq \|x\| + \|y\|$ for any $x, y \in \mathbb{R}^n$.

Note:

- The following are some commonly used norms:
 - ★ ℓ_1 norm: $\|x\|_1 := \sum_{i=1}^n |x_i|$.
 - ★ ℓ_2 norm: $\|x\|_2 := \sqrt{\sum_{i=1}^n |x_i|^2}$.
 - ★ ℓ_∞ norm: $\|x\|_\infty := \max_{1 \leq i \leq n} |x_i|$.

Norm

Definition: A function $\| \cdot \| : \mathbb{R}^n \rightarrow \mathbb{R}$ is called a (vector) **norm** if

- $\|x\| \geq 0$ for all $x \in \mathbb{R}^n$.
- $\|x\| = 0$ if and only if $x = 0$.
- $\|\alpha x\| = |\alpha| \|x\|$ for any $\alpha \in \mathbb{R}$ and $x \in \mathbb{R}^n$.
- $\|x + y\| \leq \|x\| + \|y\|$ for any $x, y \in \mathbb{R}^n$.

Note:

- The following are some commonly used norms:
 - ★ ℓ_1 norm: $\|x\|_1 := \sum_{i=1}^n |x_i|$.
 - ★ ℓ_2 norm: $\|x\|_2 := \sqrt{\sum_{i=1}^n |x_i|^2}$.
 - ★ ℓ_∞ norm: $\|x\|_\infty := \max_{1 \leq i \leq n} |x_i|$.
- For instance, if $x = [3 \quad -4 \quad 5]^T$, then $\|x\|_1 = 12$, $\|x\|_2 = \sqrt{50}$ and $\|x\|_\infty = 5$.

Norm cont.

Theorem 1.1: Let $\|\cdot\|$ be a norm. Then there exist positive numbers C_1 and C_2 so that for all $x \in \mathbb{R}^n$,

$$C_1 \sum_{i=1}^n |x_i| \leq \|x\| \leq C_2 \sum_{i=1}^n |x_i|$$

Norm cont.

Theorem 1.1: Let $\|\cdot\|$ be a norm. Then there exist positive numbers C_1 and C_2 so that for all $x \in \mathbb{R}^n$,

$$C_1 \sum_{i=1}^n |x_i| \leq \|x\| \leq C_2 \sum_{i=1}^n |x_i|$$

Proof: We will only prove the **second inequality**. The first inequality is a consequence of **compactness** and is left as an exercise later.

To prove the **second inequality**, notice that for any $x \in \mathbb{R}^n$, we have

$$\|x\| = \left\| \sum_{i=1}^n x_i e_i \right\| \leq \sum_{i=1}^n \|x_i e_i\| = \sum_{i=1}^n |x_i| \|e_i\| \leq C_2 \sum_{i=1}^n |x_i|,$$

where $C_2 := \max_{1 \leq i \leq n} \|e_i\|$.

Convergence and norm

Definition: Let $\{x^k\} \subset \mathbb{R}^n$ be a sequence and $x^* \in \mathbb{R}^n$. We say that $\lim_{k \rightarrow \infty} x^k = x^*$ if

$$\lim_{k \rightarrow \infty} x_i^k = x_i^* \quad \forall i.$$

Convergence and norm

Definition: Let $\{x^k\} \subset \mathbb{R}^n$ be a sequence and $x^* \in \mathbb{R}^n$. We say that $\lim_{k \rightarrow \infty} x^k = x^*$ if

$$\lim_{k \rightarrow \infty} x_i^k = x_i^* \quad \forall i.$$

Corollary 1.1: Let $\|\cdot\|$ be a norm, $\{x^k\} \subset \mathbb{R}^n$ be a sequence and $x^* \in \mathbb{R}^n$. Then $\lim_{k \rightarrow \infty} x^k = x^*$ if and only if $\lim_{k \rightarrow \infty} \|x^k - x^*\| = 0$.

Convergence and norm

Definition: Let $\{x^k\} \subset \mathbb{R}^n$ be a sequence and $x^* \in \mathbb{R}^n$. We say that $\lim_{k \rightarrow \infty} x^k = x^*$ if

$$\lim_{k \rightarrow \infty} x_i^k = x_i^* \quad \forall i.$$

Corollary 1.1: Let $\|\cdot\|$ be a norm, $\{x^k\} \subset \mathbb{R}^n$ be a sequence and $x^* \in \mathbb{R}^n$. Then $\lim_{k \rightarrow \infty} x^k = x^*$ if and only if $\lim_{k \rightarrow \infty} \|x^k - x^*\| = 0$.

Proof: Note that $\lim_{k \rightarrow \infty} x_i^k = x_i^*$ for all i is the same as $\lim_{k \rightarrow \infty} |x_i^k - x_i^*| = 0$ for all i , which in turn is equivalent to

$$\lim_{k \rightarrow \infty} \sum_{i=1}^n |x_i^k - x_i^*| = 0.$$

The conclusion now follows from this and [Theorem 1.1](#).

Matrix norm

Definition: A function $\| \cdot \| : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}$ is called a **matrix norm** if

- $\|A\| \geq 0$ for all $A \in \mathbb{R}^{n \times n}$.
- $\|A\| = 0$ if and only if $A = 0$.
- $\|\alpha A\| = |\alpha| \|A\|$ for any $\alpha \in \mathbb{R}$ and $A \in \mathbb{R}^{n \times n}$.
- $\|A + B\| \leq \|A\| + \|B\|$ for any $A, B \in \mathbb{R}^{n \times n}$.
- $\|AB\| \leq \|A\| \|B\|$ for any $A, B \in \mathbb{R}^{n \times n}$.

Matrix norm

The following theorem provides a large source of matrix norms.

Theorem 1.2: Let $\| \cdot \|$ be a norm. Then the following function defines a matrix norm

$$\|A\| := \max_{\|x\|=1} \|Ax\|.$$

Matrix norm

The following theorem provides a large source of matrix norms.

Theorem 1.2: Let $\| \cdot \|$ be a norm. Then the following function defines a matrix norm

$$\|A\| := \max_{\|x\|=1} \|Ax\|.$$

Remarks:

- A matrix norm taking the above form is said to be **induced by the norm $\| \cdot \|$** , or simply an **induced matrix norm**.

Matrix norm

The following theorem provides a large source of matrix norms.

Theorem 1.2: Let $\| \cdot \|$ be a norm. Then the following function defines a matrix norm

$$\|A\| := \max_{\|x\|=1} \|Ax\|.$$

Remarks:

- A matrix norm taking the above form is said to be **induced by the norm $\| \cdot \|$** , or simply an **induced matrix norm**.
- The maximum is actually **attained** at some x satisfying $\|x\| = 1$. We will need this fact below.

Matrix norm

The following theorem provides a large source of matrix norms.

Theorem 1.2: Let $\|\cdot\|$ be a norm. Then the following function defines a matrix norm

$$\|A\| := \max_{\|x\|=1} \|Ax\|.$$

Remarks:

- A matrix norm taking the above form is said to be **induced by the norm $\|\cdot\|$** , or simply an **induced matrix norm**.
- The maximum is actually **attained** at some x satisfying $\|x\| = 1$. We will need this fact below.
- The attainment is due to **compactness** of the set $\{x : \|x\| = 1\}$: the continuous function $x \mapsto \|Ax\|$ attains its maximum over the compact set $\{x : \|x\| = 1\}$.

Matrix norm

Proof of Theorem 1.2: Properties 1, 3 and 4 of the matrix norm are straightforward to verify.

Matrix norm

Proof of Theorem 1.2: Properties 1, 3 and 4 of the matrix norm are straightforward to verify.

Property 2: If $A = 0$, then clearly $\|A\| = 0$. Conversely, if $\|A\| = 0$, then $Ax = 0$ whenever $\|x\| = 1$, and hence $Ax = 0$ for all x . Thus, $A = 0$.

Matrix norm

Proof of Theorem 1.2: Properties 1, 3 and 4 of the matrix norm are straightforward to verify.

Property 2: If $A = 0$, then clearly $\|A\| = 0$. Conversely, if $\|A\| = 0$, then $Ax = 0$ whenever $\|x\| = 1$, and hence $Ax = 0$ for all x . Thus, $A = 0$.

Property 5: By the definition of $\|\cdot\|$, we have for all x with $\|x\| = 1$ that

$$\|Ax\| \leq \|A\|.$$

Consider any $x \neq 0$. Then $\|\frac{x}{\|x\|}\| = 1$ and hence $\|A\frac{x}{\|x\|}\| \leq \|A\|$. Thus, $\|Ax\| \leq \|A\|\|x\|$ for any $x \neq 0$, and hence for all x since the inequality holds trivially for $x = 0$.

Matrix norm

Proof of Theorem 1.2: Properties 1, 3 and 4 of the matrix norm are straightforward to verify.

Property 2: If $A = 0$, then clearly $\|A\| = 0$. Conversely, if $\|A\| = 0$, then $Ax = 0$ whenever $\|x\| = 1$, and hence $Ax = 0$ for all x . Thus, $A = 0$.

Property 5: By the definition of $\|\cdot\|$, we have for all x with $\|x\| = 1$ that

$$\|Ax\| \leq \|A\|.$$

Consider any $x \neq 0$. Then $\|\frac{x}{\|x\|}\| = 1$ and hence $\|A\frac{x}{\|x\|}\| \leq \|A\|$. Thus, $\|Ax\| \leq \|A\|\|x\|$ for any $x \neq 0$, and hence for all x since the inequality holds trivially for $x = 0$.

Then

$$\|AB\| = \max_{\|x\|=1} \|A(Bx)\| \leq \max_{\|x\|=1} \|A\|\|Bx\| = \|A\|\|B\|.$$

Example 1

Example: The following functions are matrix norms:

- $\|A\|_1 = \max_j \sum_{i=1}^n |a_{ij}|$ (maximum of the ℓ_1 norms of columns).
Moreover, this is an **induced matrix norm**:

$$\|A\|_1 = \max_{\|x\|_1=1} \|Ax\|_1.$$

Example 1

Example: The following functions are matrix norms:

- $\|A\|_1 = \max_j \sum_{i=1}^n |a_{ij}|$ (maximum of the ℓ_1 norms of columns).
Moreover, this is an **induced matrix norm**:

$$\|A\|_1 = \max_{\|x\|_1=1} \|Ax\|_1.$$

- $\|A\|_2 = \sqrt{\lambda_{\max}(A^T A)}$. Moreover, this is an **induced matrix norm**:

$$\|A\|_2 = \max_{\|x\|_2=1} \|Ax\|_2.$$

Example 1

Example: The following functions are matrix norms:

- $\|A\|_1 = \max_j \sum_{i=1}^n |a_{ij}|$ (maximum of the ℓ_1 norms of columns).
Moreover, this is an **induced matrix norm**:

$$\|A\|_1 = \max_{\|x\|_1=1} \|Ax\|_1.$$

- $\|A\|_2 = \sqrt{\lambda_{\max}(A^T A)}$. Moreover, this is an **induced matrix norm**:

$$\|A\|_2 = \max_{\|x\|_2=1} \|Ax\|_2.$$

- $\|A\|_\infty = \max_i \sum_{j=1}^n |a_{ij}|$ (maximum of the ℓ_1 norms of rows).
Moreover, this is an **induced matrix norm**:

$$\|A\|_\infty = \max_{\|x\|_\infty=1} \|Ax\|_\infty.$$

Example 1

Example: The following functions are matrix norms:

- $\|A\|_1 = \max_j \sum_{i=1}^n |a_{ij}|$ (maximum of the ℓ_1 norms of columns).
Moreover, this is an **induced matrix norm**:

$$\|A\|_1 = \max_{\|x\|_1=1} \|Ax\|_1.$$

- $\|A\|_2 = \sqrt{\lambda_{\max}(A^T A)}$. Moreover, this is an **induced matrix norm**:

$$\|A\|_2 = \max_{\|x\|_2=1} \|Ax\|_2.$$

- $\|A\|_\infty = \max_i \sum_{j=1}^n |a_{ij}|$ (maximum of the ℓ_1 norms of rows).
Moreover, this is an **induced matrix norm**:

$$\|A\|_\infty = \max_{\|x\|_\infty=1} \|Ax\|_\infty.$$

- $\|A\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^n |a_{ij}|^2}$. This is known as the Fröbenius norm.

Example 2

Example: Consider

$$A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

Then

- $\|A\|_1 = \max\{4, 6\} = 6$.

Example 2

Example: Consider

$$A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

Then

- $\|A\|_1 = \max\{4, 6\} = 6$.
- $A^T A = \begin{bmatrix} 10 & 14 \\ 14 & 20 \end{bmatrix}$, and the eigenvalues of $A^T A$ are $15 \pm \sqrt{221}$.

Hence

$$\|A\|_2 = \sqrt{15 + \sqrt{221}}.$$

Example 2

Example: Consider

$$A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

Then

- $\|A\|_1 = \max\{4, 6\} = 6$.
- $A^T A = \begin{bmatrix} 10 & 14 \\ 14 & 20 \end{bmatrix}$, and the eigenvalues of $A^T A$ are $15 \pm \sqrt{221}$.

Hence

$$\|A\|_2 = \sqrt{15 + \sqrt{221}}.$$

- $\|A\|_\infty = \max\{3, 7\} = 7$.

Example 2

Example: Consider

$$A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

Then

- $\|A\|_1 = \max\{4, 6\} = 6$.
- $A^T A = \begin{bmatrix} 10 & 14 \\ 14 & 20 \end{bmatrix}$, and the eigenvalues of $A^T A$ are $15 \pm \sqrt{221}$.

Hence

$$\|A\|_2 = \sqrt{15 + \sqrt{221}}.$$

- $\|A\|_\infty = \max\{3, 7\} = 7$.
- $\|A\|_F = \sqrt{1^2 + 2^2 + 3^2 + 4^2} = \sqrt{30}$.

Lower Semicontinuity

Definition: A function $f : \Omega \rightarrow \mathbb{R}$ is said to be **lower semicontinuous** at a point $x \in \Omega$ if for any sequence $\{x^k\}$ converging to x , it holds that

$$f(x) \leq \liminf_{i \rightarrow \infty} f(x^k).$$

If f is lower semicontinuous at every $x \in \Omega$, then we say that f is lower semicontinuous on Ω .

Lower Semicontinuity

Definition: A function $f : \Omega \rightarrow \mathbb{R}$ is said to be **lower semicontinuous** at a point $x \in \Omega$ if for any sequence $\{x^k\}$ converging to x , it holds that

$$f(x) \leq \liminf_{i \rightarrow \infty} f(x^k).$$

If f is lower semicontinuous at every $x \in \Omega$, then we say that f is lower semicontinuous on Ω .

Example:

- The following function is lower semicontinuous:

$$f(x) = \begin{cases} x^2 & \text{if } x \leq 0, \\ x^2 + 1 & \text{if } x > 0. \end{cases}$$

- But the following function is not lower semicontinuous:

$$f(x) = \begin{cases} x^2 & \text{if } x < 0, \\ x^2 + 1 & \text{if } x \geq 0. \end{cases}$$

Compactness

Definition: A set $\Omega \subseteq \mathbb{R}^n$ is said to be **closed** if it contains all the limits of convergent sequences of points in Ω .

Compactness

Definition: A set $\Omega \subseteq \mathbb{R}^n$ is said to be **closed** if it contains all the limits of convergent sequences of points in Ω .

Example:

- The set $(0, 1)$ is not closed in \mathbb{R} , the set $[0, 1]$ is closed in \mathbb{R} .
- The set $\{x : \|x\|_2 \leq 1\}$ is closed in \mathbb{R}^n but $\{x : \|x\|_2 < 1\}$ is not.

Compactness

Definition: A set $\Omega \subseteq \mathbb{R}^n$ is said to be **closed** if it contains all the limits of convergent sequences of points in Ω .

Example:

- The set $(0, 1)$ is not closed in \mathbb{R} , the set $[0, 1]$ is closed in \mathbb{R} .
- The set $\{x : \|x\|_2 \leq 1\}$ is closed in \mathbb{R}^n but $\{x : \|x\|_2 < 1\}$ is not.

Definition. A set $\Omega \subseteq \mathbb{R}^n$ is said to be **bounded** if there exists $K > 0$ so that $\Omega \subseteq \{x : \|x\|_2 \leq K\}$.

Compactness

Definition: A set $\Omega \subseteq \mathbb{R}^n$ is said to be **closed** if it contains all the limits of convergent sequences of points in Ω .

Example:

- The set $(0, 1)$ is not closed in \mathbb{R} , the set $[0, 1]$ is closed in \mathbb{R} .
- The set $\{x : \|x\|_2 \leq 1\}$ is closed in \mathbb{R}^n but $\{x : \|x\|_2 < 1\}$ is not.

Definition. A set $\Omega \subseteq \mathbb{R}^n$ is said to be **bounded** if there exists $K > 0$ so that $\Omega \subseteq \{x : \|x\|_2 \leq K\}$.

Theorem 1.3: (Bolzano-Weierstrass)

Let $\Omega \subset \mathbb{R}^n$ be bounded. If $\{x^k\} \subseteq \Omega$, then there exists a convergent subsequence $\{x^{k_i}\}$, i.e., for some $x^* \in \mathbb{R}^n$, we have

$$\lim_{i \rightarrow \infty} x^{k_i} = x^*.$$

Note: A closed and bounded set in \mathbb{R}^n is called a **compact set**.

Existence of minimizers

Theorem 1.4: (Existence of minimizers)

Let $\Omega \subset \mathbb{R}^n$ be a **nonempty compact set** and $f : \Omega \rightarrow \mathbb{R}$ be **lower semicontinuous on Ω** . Then f achieves its infimum value over Ω , i.e., there exists $x^* \in \Omega$ so that $f(x^*) = \inf\{f(x) : x \in \Omega\}$.

Proof: Let $\ell := \inf\{f(x) : x \in \Omega\}$ and let $\{\lambda_k\} \subset \mathbb{R}$ be a **strictly decreasing sequence** converging to ℓ .

By the **definition of infimum**, for each λ_k , $k = 1, 2, \dots$, there exists $x^k \in \Omega$ so that

$$\ell \leq f(x^k) < \lambda_k.$$

Since $\{x^k\} \subseteq \Omega$ and Ω is bounded, by **Bolzano-Weierstrass theorem** there exists a subsequence $\{x^{k_i}\}$ converging to some $x^* \in \mathbb{R}^n$. Since $\{x^{k_i}\}$ is itself a sequence in Ω and Ω is closed, $x^* \in \Omega$. Thus,

$$\ell \leq f(x^*) \leq \liminf_{i \rightarrow \infty} f(x^{k_i}) \leq \lim_{i \rightarrow \infty} \lambda_{k_i} = \ell,$$

showing that f achieves ℓ at $x^* \in \Omega$.

Positive semidefinite matrices

Definition: (Positive semidefinite matrices)

Let $A \in \mathbb{R}^{n \times n}$ be symmetric. We say that A is positive semidefinite if $x^T A x \geq 0$ for all $x \in \mathbb{R}^n$.

Notation: $A \succeq 0$. The set of $n \times n$ positive semidefinite matrices is denoted by S_+^n .

Positive semidefinite matrices

Definition: (Positive semidefinite matrices)

Let $A \in \mathbb{R}^{n \times n}$ be symmetric. We say that A is positive semidefinite if $x^T A x \geq 0$ for all $x \in \mathbb{R}^n$.

Notation: $A \succeq 0$. The set of $n \times n$ positive semidefinite matrices is denoted by S_+^n .

Example: The matrix

$$A = \begin{bmatrix} 3 & 1 \\ 1 & 2 \end{bmatrix}$$

is positive semidefinite. To see this, note that

$$\begin{aligned} x^T A x &= 3x_1^2 + 2x_1x_2 + 2x_2^2 \\ &= 2x_1^2 + (x_1 + x_2)^2 + x_2^2 \geq 0. \end{aligned}$$

Positive semidefinite matrices

Definition: (Positive semidefinite matrices)

Let $A \in \mathbb{R}^{n \times n}$ be symmetric. We say that A is positive semidefinite if $x^T A x \geq 0$ for all $x \in \mathbb{R}^n$.

Notation: $A \succeq 0$. The set of $n \times n$ positive semidefinite matrices is denoted by S_+^n .

Example: The matrix

$$A = \begin{bmatrix} 3 & 1 \\ 1 & 2 \end{bmatrix}$$

is positive semidefinite. To see this, note that

$$\begin{aligned} x^T A x &= 3x_1^2 + 2x_1x_2 + 2x_2^2 \\ &= 2x_1^2 + (x_1 + x_2)^2 + x_2^2 \geq 0. \end{aligned}$$

Question: Easier way to test for positive semidefiniteness?

Positive semidefinite matrices cont.

Theorem 1.5: Let $A \in \mathbb{R}^{n \times n}$ be **symmetric**. The following statements are equivalent.

1. All eigenvalues of A are **nonnegative**.
2. There exists $M \in \mathbb{R}^{n \times n}$ so that $A = M^T M$.
3. A is **positive semidefinite**.

Theorem 1.5 proof sketch:

(1) \Rightarrow (2): Since A is symmetric, there exist an **orthogonal matrix** U and a **diagonal matrix** D so that $A = UDU^T$.

Positive semidefinite matrices cont.

Theorem 1.5: Let $A \in \mathbb{R}^{n \times n}$ be **symmetric**. The following statements are equivalent.

1. All eigenvalues of A are **nonnegative**.
2. There exists $M \in \mathbb{R}^{n \times n}$ so that $A = M^T M$.
3. A is **positive semidefinite**.

Theorem 1.5 proof sketch:

(1) \Rightarrow (2): Since A is symmetric, there exist an **orthogonal matrix** U and a **diagonal matrix** D so that $A = UDU^T$.

Since all eigenvalues of A are **nonnegative**, we have $D_{ii} \geq 0$ for all i .

Let $W \in \mathbb{R}^{n \times n}$ be the matrix so that

$$W_{ij} = \begin{cases} \sqrt{D_{ii}} & \text{if } i = j, \\ 0 & \text{otherwise.} \end{cases}$$

Positive semidefinite matrices cont.

Theorem 1.5 proof sketch cont.:

Then $W = W^T$ and

$$A = U(WW)U^T = (WU^T)^T(WU^T).$$

Thus, (2) holds with $M = WU^T$.

Positive semidefinite matrices cont.

Theorem 1.5 proof sketch cont.:

Then $W = W^T$ and

$$A = U(WW)U^T = (WU^T)^T(WU^T).$$

Thus, (2) holds with $M = WU^T$.

(2) \Rightarrow (3): Let $x \in \mathbb{R}^n$ and $y := Mx$. Then

$$x^T Ax$$

Positive semidefinite matrices cont.

Theorem 1.5 proof sketch cont.:

Then $W = W^T$ and

$$A = U(WW)U^T = (WU^T)^T(WU^T).$$

Thus, (2) holds with $M = WU^T$.

(2) \Rightarrow (3): Let $x \in \mathbb{R}^n$ and $y := Mx$. Then

$$x^T Ax = x^T M^T Mx = (Mx)^T (Mx) = y^T y \geq 0.$$

Positive semidefinite matrices cont.

Theorem 1.5 proof sketch cont.:

Then $W = W^T$ and

$$A = U(WW)U^T = (WU^T)^T(WU^T).$$

Thus, (2) holds with $M = WU^T$.

(2) \Rightarrow (3): Let $x \in \mathbb{R}^n$ and $y := Mx$. Then

$$x^T Ax = x^T M^T Mx = (Mx)^T (Mx) = y^T y \geq 0.$$

(3) \Rightarrow (1): Let λ be an eigenvalue of A with a corresponding eigenvector v , i.e.,

$$v \neq 0 \text{ and } Av = \lambda v.$$

Then $v^T v > 0$ and

$$\lambda v^T v = v^T Av \geq 0.$$

Thus, it follows that $\lambda \geq 0$.

Positive definite matrices

Definition: A symmetric matrix $A \in \mathbb{R}^{n \times n}$ is called **positive definite** if $x^T A x > 0$ for all $x \in \mathbb{R}^n \setminus \{0\}$.

Notation: $A \succ 0$.

Theorem 1.6: For a symmetric matrix $A \in \mathbb{R}^{n \times n}$, the following statements are equivalent:

- All eigenvalues of A are **positive**.
- There exists an **invertible** matrix $M \in \mathbb{R}^{n \times n}$ so that $A = M^T M$.
- A is **positive definite**.

Note: Let $A \succ 0$, then

- $A^{-1} \succ 0$ and $\lambda_{\min}(A) = \inf\{x^T A x : \|x\|_2 = 1\}$.
- $\|A\|_2 = \lambda_{\max}(A) = [\lambda_{\min}(A^{-1})]^{-1}$.

Block matrix multiplication

Let $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{n \times p}$ be **partitioned** so that

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \quad \text{and} \quad B = \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix},$$

where

- $A_{11} \in \mathbb{R}^{m_1 \times n_1}$, $A_{12} \in \mathbb{R}^{m_1 \times n_2}$, $A_{21} \in \mathbb{R}^{m_2 \times n_1}$ and $A_{22} \in \mathbb{R}^{m_2 \times n_2}$;
- $B_{11} \in \mathbb{R}^{n_1 \times p_1}$, $B_{12} \in \mathbb{R}^{n_1 \times p_2}$, $B_{21} \in \mathbb{R}^{n_2 \times p_1}$ and $B_{22} \in \mathbb{R}^{n_2 \times p_2}$.

Then it holds that

$$AB = \begin{bmatrix} A_{11}B_{11} + A_{12}B_{21} & A_{11}B_{12} + A_{12}B_{22} \\ A_{21}B_{11} + A_{22}B_{21} & A_{21}B_{12} + A_{22}B_{22} \end{bmatrix}.$$

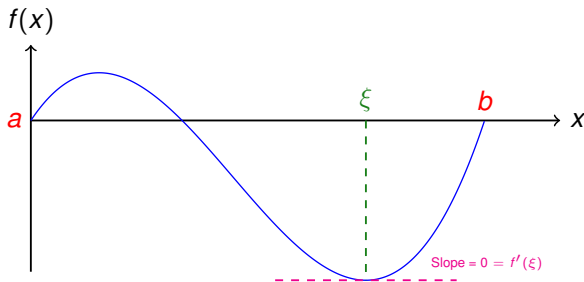
Roughly speaking, whenever the sizes match, matrix blocks can be multiplied as if they were numbers.

Mean value theorem

Theorem 1.7. (Rolle's mean value theorem)

Let f be continuous on $[a, b]$ and differentiable in (a, b) . If $f(b) = f(a)$, then there exists $\xi \in (a, b)$ so that

$$f'(\xi) = 0.$$



Taylor's theorem

Theorem 1.8. (Taylor's theorem with remainder term)

Suppose that f is $(n + 1)$ times differentiable on an open interval containing $[a, b]$. Then

$$f(b) = f(a) + f'(a)(b - a) + \frac{f''(a)}{2!}(b - a)^2 + \cdots + \frac{f^{(n)}(a)}{n!}(b - a)^n \\ + \frac{f^{(n+1)}(\xi)}{(n + 1)!}(b - a)^{n+1}$$

for some $\xi \in (a, b)$.

Taylor's theorem cont.

Proof of Theorem 1.8.: Define

$$T_n(x) = f(a) + f'(a)(x - a) + \frac{f''(a)}{2!}(x - a)^2 + \cdots + \frac{f^{(n)}(a)}{n!}(x - a)^n$$

and define K so that

$$f(b) = T_n(b) + K(b - a)^{n+1}.$$

Taylor's theorem cont.

Proof of Theorem 1.8.: Define

$$T_n(x) = f(a) + f'(a)(x - a) + \frac{f''(a)}{2!}(x - a)^2 + \cdots + \frac{f^{(n)}(a)}{n!}(x - a)^n$$

and define K so that

$$f(b) = T_n(b) + K(b - a)^{n+1}.$$

We need to show that K is given by $\frac{f^{(n+1)}(\xi)}{(n+1)!}$ for some $\xi \in (a, b)$.

Taylor's theorem cont.

Proof of Theorem 1.8.: Define

$$T_n(x) = f(a) + f'(a)(x - a) + \frac{f''(a)}{2!}(x - a)^2 + \cdots + \frac{f^{(n)}(a)}{n!}(x - a)^n$$

and define K so that

$$f(b) = T_n(b) + K(b - a)^{n+1}.$$

We need to show that K is given by $\frac{f^{(n+1)}(\xi)}{(n+1)!}$ for some $\xi \in (a, b)$.

To this end, consider

$$g(x) = f(x) - T_n(x) - K(x - a)^{n+1}.$$

Taylor's theorem cont.

Proof of Theorem 1.8.: Define

$$T_n(x) = f(a) + f'(a)(x - a) + \frac{f''(a)}{2!}(x - a)^2 + \cdots + \frac{f^{(n)}(a)}{n!}(x - a)^n$$

and define K so that

$$f(b) = T_n(b) + K(b - a)^{n+1}.$$

We need to show that K is given by $\frac{f^{(n+1)}(\xi)}{(n+1)!}$ for some $\xi \in (a, b)$.

To this end, consider

$$g(x) = f(x) - T_n(x) - K(x - a)^{n+1}.$$

Note: $g(a) = 0$ and $g(b) = 0$. Thus, **Rolle's mean value theorem** gives the existence of $a < \xi_1 < b$ with $g'(\xi_1) = 0$.

Taylor's theorem cont.

Proof of Theorem 1.8. cont.:

Note that

$$\begin{aligned} g'(x) &= f'(x) - T'_n(x) - K(n+1)(x-a)^n \\ &= f'(x) - f'(a) - f''(a)(x-a) - \cdots - \frac{f^{(n)}(a)}{(n-1)!}(x-a)^{n-1} \\ &\quad - K(n+1)(x-a)^n. \end{aligned}$$

Hence $g'(a) = g'(\xi_1) = 0$. Thus, again by **Rolle's mean value theorem**, there exists $a < \xi_2 < \xi_1$ so that $g''(\xi_2) = 0$.

Taylor's theorem cont.

Proof of Theorem 1.8. cont.:

Note that

$$\begin{aligned}g'(x) &= f'(x) - T'_n(x) - K(n+1)(x-a)^n \\&= f'(x) - f'(a) - f''(a)(x-a) - \cdots - \frac{f^{(n)}(a)}{(n-1)!}(x-a)^{n-1} \\&\quad - K(n+1)(x-a)^n.\end{aligned}$$

Hence $g'(a) = g'(\xi_1) = 0$. Thus, again by **Rolle's mean value theorem**, there exists $a < \xi_2 < \xi_1$ so that $g''(\xi_2) = 0$.

Proceeding **inductively**, there exist $a < \xi_n < \xi_{n-1} < \cdots < \xi_1 < b$ so that

$$g'(\xi_1) = g''(\xi_2) = \cdots = g^{(n)}(\xi_n) = 0.$$

Taylor's theorem cont.

Proof of Theorem 1.8. cont.: Finally, notice that

$$\begin{aligned} g^{(n)}(x) &= f^{(n)}(x) - T_n^{(n)}(x) - K(n+1)!(x-a) \\ &= f^{(n)}(x) - f^{(n)}(a) - K(n+1)!(x-a). \end{aligned}$$

Taylor's theorem cont.

Proof of Theorem 1.8. cont.: Finally, notice that

$$\begin{aligned}g^{(n)}(x) &= f^{(n)}(x) - T_n^{(n)}(x) - K(n+1)!(x-a) \\&= f^{(n)}(x) - f^{(n)}(a) - K(n+1)!(x-a).\end{aligned}$$

Since $g^{(n)}(a) = g^{(n)}(\xi_n) = 0$, **Rolle's mean value theorem** gives the existence of $\xi_{n+1} \in (a, \xi_n) \subset (a, b)$ such that

$$0 = g^{(n+1)}(\xi_{n+1}) = f^{(n+1)}(\xi_{n+1}) - K(n+1)!,$$

which gives

$$K = \frac{f^{(n+1)}(\xi_{n+1})}{(n+1)!}.$$

Gradient and Hessian

- Let $f \in C^1(\mathbb{R}^n)$. Its gradient at an $x \in \mathbb{R}^n$ is

$$\nabla f(x) := \begin{bmatrix} \frac{\partial f}{\partial x_1}(x) \\ \frac{\partial f}{\partial x_2}(x) \\ \vdots \\ \frac{\partial f}{\partial x_n}(x) \end{bmatrix}.$$

- Let $f \in C^2(\mathbb{R}^n)$. Its Hessian at an $x \in \mathbb{R}^n$ is

$$\nabla^2 f(x) := \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2}(x) & \frac{\partial^2 f}{\partial x_2 \partial x_1}(x) & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_1}(x) \\ \frac{\partial^2 f}{\partial x_1 \partial x_2}(x) & \cdots & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_2}(x) \\ \vdots & \ddots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_1 \partial x_n}(x) & \cdots & \cdots & \frac{\partial^2 f}{\partial x_n^2}(x) \end{bmatrix}.$$

Note: Since $f \in C^2(\mathbb{R}^n)$, we have $\frac{\partial^2 f}{\partial x_i \partial x_j} = \frac{\partial^2 f}{\partial x_j \partial x_i}$ for all i and j .

High-dimensional Taylor's theorem

Theorem 1.9. (Taylor's theorem in \mathbb{R}^n with remainder term)

- Let $f \in C^1(\mathbb{R}^n)$, x and $y \in \mathbb{R}^n$. Then there exists $\xi \in \{(1-s)x + sy : s \in (0, 1)\}$ such that

$$f(y) = f(x) + [\nabla f(\xi)]^T (y - x).$$

- Let $f \in C^2(\mathbb{R}^n)$, x and $y \in \mathbb{R}^n$. Then there exists $\xi \in \{(1-s)x + sy : s \in (0, 1)\}$ such that

$$f(y) = f(x) + [\nabla f(x)]^T (y - x) + \frac{1}{2}(y - x)^T \nabla^2 f(\xi)(y - x).$$

Proof sketch: Consider the function $\psi(t) := f((1-t)x + ty)$. Observe that ψ is C^1 (resp. C^2) if f is so.

High-dimensional Taylor's theorem

Theorem 1.9. (Taylor's theorem in \mathbb{R}^n with remainder term)

- Let $f \in C^1(\mathbb{R}^n)$, x and $y \in \mathbb{R}^n$. Then there exists $\xi \in \{(1-s)x + sy : s \in (0, 1)\}$ such that

$$f(y) = f(x) + [\nabla f(\xi)]^T (y - x).$$

- Let $f \in C^2(\mathbb{R}^n)$, x and $y \in \mathbb{R}^n$. Then there exists $\xi \in \{(1-s)x + sy : s \in (0, 1)\}$ such that

$$f(y) = f(x) + [\nabla f(x)]^T (y - x) + \frac{1}{2}(y - x)^T \nabla^2 f(\xi)(y - x).$$

Proof sketch: Consider the function $\psi(t) := f((1-t)x + ty)$. Observe that ψ is C^1 (resp. C^2) if f is so. Moreover, using chain rule,

$$\psi'(t) = [\nabla f((1-t)x + ty)]^T (y-x), \psi''(t) = (y-x)^T [\nabla^2 f((1-t)x + ty)](y-x).$$

Now apply Taylor's theorem in 1 dimension to ψ .

Convention of sequence notation

Here are some conventions for **superscript** and **subscript** notation for understanding the lecture notes:

Convention of sequence notation

Here are some conventions for **superscript** and **subscript** notation for understanding the lecture notes:

For a **scalar** x :

- x^k represents x to the power k .
- x_k is term k in the sequence $\{x_k\}$.

Convention of sequence notation

Here are some conventions for **superscript** and **subscript** notation for understanding the lecture notes:

For a **scalar** x :

- x^k represents x to the power k .
- x_k is term k in the sequence $\{x_k\}$.

For a **vector** $x \in \mathbb{R}^n$, $n \geq 2$:

- x^k is term k in the sequence $\{x^k\}$.
- x_k is the k th entry of x .

Convention of sequence notation

Here are some conventions for **superscript** and **subscript** notation for understanding the lecture notes:

For a **scalar** x :

- x^k represents x to the power k .
- x_k is term k in the sequence $\{x_k\}$.

For a **vector** $x \in \mathbb{R}^n$, $n \geq 2$:

- x^k is term k in the sequence $\{x^k\}$.
- x_k is the k th entry of x .

Ambiguity: x_1^2 usually means $(x_1)^2$, and rarely means $(x^2)_1$.

Convention of sequence notation

Here are some conventions for **superscript** and **subscript** notation for understanding the lecture notes:

For a **scalar** x :

- x^k represents x to the power k .
- x_k is term k in the sequence $\{x_k\}$.

For a **vector** $x \in \mathbb{R}^n$, $n \geq 2$:

- x^k is term k in the sequence $\{x^k\}$.
- x_k is the k th entry of x .

Ambiguity: x_1^2 usually means $(x_1)^2$, and rarely means $(x^2)_1$.
You do not need to follow this convention in your writings.