



Lecture Agenda

- **Basics of Matrix Analysis**
- Taylor Expansion & Lipschitz Function
- Optimality Conditions
- Rate of Convergence

1



Example: Matrix Completion



- **Netflix Challenge:** Netflix provides highly incomplete ratings from 0.5 million users for $\approx 17,770$ movies
- How to predict unseen user ratings for movies?

1/24/24

2

UNIVERSITY OF MICHIGAN

Example: Matrix Completion

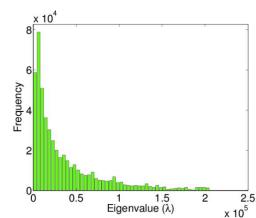


More unknowns than observations (under-determined system)

3



Example: Matrix Completion



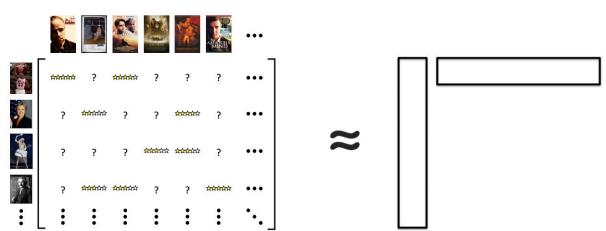
A few factors explain most of the data: **low-rank** approximation

1/24/24

4

UNIVERSITY OF MICHIGAN

Example: Matrix Completion



How to exploit **low-rank structures** in predictions?

5



Example: Matrix Completion

Given observed entries $M_{i,j}$ with $(i, j) \in \Omega$, complete the matrix via rank minimization

$$\min_{\mathbf{X}} \text{rank}(\mathbf{X}), \quad \text{s.t. } X_{i,j} = M_{i,j}, \quad (i, j) \in \Omega,$$

or, equivalently, we can rewrite it as

$$\min_{\mathbf{X}} \text{rank}(\mathbf{X}), \quad \text{s.t. } \mathcal{P}_{\Omega}(\mathbf{X}) = \mathcal{P}_{\Omega}(\mathbf{M}).$$

1/24/24

6

UNIVERSITY OF MICHIGAN

Example: Matrix Completion

$$\min_{\mathbf{X}} \text{rank}(\mathbf{X}), \quad \text{s.t. } \mathcal{P}_{\Omega}(\mathbf{X}) = \mathcal{P}_{\Omega}(\mathbf{M}),$$

where $\mathcal{P}_{\Omega}(\cdot)$ is the orthogonal projection onto the subspace of matrices support on Ω .

$$\mathcal{P}_{\Omega}(\mathbf{X}) := \begin{cases} X_{ij}, & \text{if } (i, j) \in \Omega, \\ 0, & \text{if } (i, j) \notin \Omega. \end{cases}$$

7



Example: Matrix Completion

$$\min_{\mathbf{X}} \text{rank}(\mathbf{X}), \quad \text{s.t. } \mathcal{P}_{\Omega}(\mathbf{X}) = \mathcal{P}_{\Omega}(\mathbf{M}),$$

- Similar to ℓ_0 -norm minimization, the rank minimization is also **NP-hard**
- Remedy: convex relaxation via *nuclear norm*:

$$\text{rank}(\mathbf{X}) = \|\boldsymbol{\sigma}(\mathbf{X})\|_0, \quad \|\mathbf{X}\|_* = \|\boldsymbol{\sigma}(\mathbf{X})\|_1,$$

where $\|\mathbf{X}\|_* = \sum_{i=1}^n \sigma_i(\mathbf{X})$ is a *convex surrogate* of $\text{rank}(\mathbf{X})$

Matrix Analysis Basics

- Matrix (Frobenius) inner product:**

$$\langle \mathbf{x}, \mathbf{z} \rangle = \sum_{i=1}^n x_i z_i \implies \langle \mathbf{X}, \mathbf{Z} \rangle := \sum_{i=1}^m \sum_{j=1}^n X_{ij} Z_{ij}$$

- Matrix trace:** for $\mathbf{M} \in \mathbb{R}^{n \times n}$, $\text{tr}(\mathbf{M}) := \sum_{i=1}^n M_{ii}$

$$\langle \mathbf{X}, \mathbf{Z} \rangle = \text{tr}(\mathbf{X}^\top \mathbf{Z}) = \text{tr}(\mathbf{X} \mathbf{Z}^\top)$$

Matrix Analysis Basics

- Matrix (Frobenius) inner product:**

$$\langle \mathbf{x}, \mathbf{z} \rangle = \sum_{i=1}^n x_i z_i \implies \langle \mathbf{X}, \mathbf{Z} \rangle := \sum_{i=1}^m \sum_{j=1}^n X_{ij} Z_{ij}$$

- Matrix trace:** for $\mathbf{M} \in \mathbb{R}^{n \times n}$, $\text{tr}(\mathbf{M}) := \sum_{i=1}^n M_{ii}$

$$\text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA})$$

$$\text{tr}(\mathbf{A}_1 \mathbf{A}_2 \cdots \mathbf{A}_n) = \text{tr}(\mathbf{A}_{\pi(1)} \mathbf{A}_{\pi(2)} \cdots \mathbf{A}_{\pi(n)})$$

Where π is a cyclic permutation on $\{1, 2, \dots, n\}$

Matrix Analysis Basics

- Matrix (Frobenius) inner product:**

$$\langle \mathbf{x}, \mathbf{z} \rangle = \sum_{i=1}^n x_i z_i \implies \langle \mathbf{X}, \mathbf{Z} \rangle := \sum_{i=1}^m \sum_{j=1}^n X_{ij} Z_{ij}$$

- Matrix trace:** for $\mathbf{M} \in \mathbb{R}^{n \times n}$, $\text{tr}(\mathbf{M}) := \sum_{i=1}^n M_{ii}$

$$\langle \mathbf{X}, \mathbf{Z} \rangle = \text{tr}(\mathbf{X}^\top \mathbf{Z}) = \text{tr}(\mathbf{X} \mathbf{Z}^\top)$$

- Frobenius norm:**

$$\|\mathbf{X}\|_F = \sqrt{\text{tr}(\mathbf{X}^* \mathbf{X})} = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |X_{ij}|^2}$$

Eigenvectors & Eigenvalues

Let $\mathbf{A} \in \mathbb{C}^{n \times n}$. We say that $\lambda \in \mathbb{C}$ is an *eigenvalue* of \mathbf{A} if there exists some $\mathbf{v} \in \mathbb{C}^n \setminus \{\mathbf{0}\}$, such that

$$\mathbf{A} \cdot \mathbf{v} = \begin{array}{c} \lambda \\ \text{eigenvalue} \end{array} \cdot \begin{array}{c} \mathbf{v} \\ \text{eigenvector} \end{array}.$$

Fact. $\lambda \in \mathbb{C}$ is an eigenvalue of $\mathbf{A} \in \mathbb{C}^{n \times n}$ iff it is a root of the *characteristic polynomial* (i.e., $\xi(\lambda) = 0$)

$$\xi(\lambda) := \det(\mathbf{A} - \lambda \mathbf{I}).$$

Note: even for real $\mathbf{A} \in \mathbb{R}^{n \times n}$, its eigenvalue λ is *not* necessarily real.

Eigen Decomposition of Symmetric Matrices

Fact. Let $\mathbf{A} \in \mathbb{R}^{n \times n}$ be *symmetric*. Then there exist *orthonormal* vectors $\{\mathbf{v}_i\}_{i=1}^n \in \mathbb{R}^n$ and *real* scalars $\lambda_1 \geq \dots \geq \lambda_n$,

$$\mathbf{V} = [\mathbf{v}_1 \ \dots \ \mathbf{v}_n] \in O(n), \quad \Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n).$$

such that

$$\mathbf{A} = \mathbf{V} \Lambda \mathbf{V}^\top = \sum_{i=1}^n \lambda_i \mathbf{v}_i \mathbf{v}_i^\top.$$

Matrix Analysis Basics

- Positive (semi)definiteness**

Definition. A symmetric matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is *positive definite* if $\mathbf{x}^\top \mathbf{A} \mathbf{x} > 0$ for all nonzero $\mathbf{x} \in \mathbb{R}^n$.

It is *positive semidefinite (p.s.d.)* if $\mathbf{x}^\top \mathbf{A} \mathbf{x} \geq 0$ for all $\mathbf{x} \in \mathbb{R}^n$.

$$\mathbf{A} \in \mathbb{R}^{n \times n} \text{ positive definite:} \quad \mathbf{A} \succ \mathbf{0}$$

$$\mathbf{A} \in \mathbb{R}^{n \times n} \text{ positive semidefinite:} \quad \mathbf{A} \succeq \mathbf{0}$$

Positive (Semi)definiteness

Definition. A symmetric matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is *positive definite* if $\mathbf{x}^\top \mathbf{A} \mathbf{x} > 0$ for all nonzero $\mathbf{x} \in \mathbb{R}^n$.

It is *positive semidefinite (p.s.d.)* if $\mathbf{x}^\top \mathbf{A} \mathbf{x} \geq 0$ for all $\mathbf{x} \in \mathbb{R}^n$.

$$\mathbf{A} \in \mathbb{R}^{n \times n} \text{ positive definite:} \quad \mathbf{A} \succ \mathbf{0}$$

$$\mathbf{A} \in \mathbb{R}^{n \times n} \text{ positive semidefinite:} \quad \mathbf{A} \succeq \mathbf{0}$$

We also write $\mathbf{A} \succeq \mathbf{B}$ if $\mathbf{A} - \mathbf{B}$ is p.s.d., i.e., $\mathbf{A} - \mathbf{B} \succeq \mathbf{0}$.

Positive (Semi)definiteness

$\mathbf{A} \in \mathbb{R}^{n \times n}$ positive definite: $\mathbf{A} \succ \mathbf{0}$

$\mathbf{A} \in \mathbb{R}^{n \times n}$ positive semidefinite: $\mathbf{A} \succeq \mathbf{0}$

Fact. A symmetric matrix \mathbf{A} is positive definite (resp. semidefinite) iff all of its eigenvalues $\{\lambda_i\}_{i=1}^n$ are positive (resp. nonnegative).

Singular Value Decomposition (SVD)

Fact (Compact SVD). Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ with $\text{rank}(\mathbf{A}) = r$. There exist $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$, and matrices

$$\mathbf{U} \in \mathbb{R}^{m \times r}, \mathbf{V} \in \mathbb{R}^{n \times r},$$

with orthonormal columns ($\mathbf{U}^\top \mathbf{U} = \mathbf{I}$, $\mathbf{V}^\top \mathbf{V} = \mathbf{I}$) such that

$$\mathbf{A} = \mathbf{U} \Sigma \mathbf{V}^\top = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^\top, \quad \Sigma = \text{diag}(\sigma_1, \dots, \sigma_r).$$

The $\{\sigma_i\}_{i=1}^r$ are called the *singular values* of \mathbf{A} ;

The columns of \mathbf{U} and \mathbf{V} are called the (left & right) *singular vectors*.

Singular Value Decomposition (SVD)

Fact (properties of compact SVD).

Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ with $\text{rank}(\mathbf{A}) = r$ and compact SVD $\mathbf{A} = \mathbf{U} \Sigma \mathbf{V}^\top$

Let us denote:

$$\text{range}(\mathbf{A}) := \{\mathbf{Ax} \mid \mathbf{x} \in \mathbb{R}^n\} = \text{col}(\mathbf{A}).$$

Then we have

- $\text{range}(\mathbf{A}) = \text{range}(\mathbf{U})$;
- $\text{range}(\mathbf{A}^\top) = \text{range}(\mathbf{V})$.

Singular Value Decomposition (SVD)

Relationship between eigen decomposition and SVD:

- The columns of \mathbf{V} are eigenvectors of $\mathbf{A}^\top \mathbf{A}$;
- The columns of \mathbf{U} are eigenvectors of $\mathbf{A} \mathbf{A}^\top$.

$$\sigma_i(\mathbf{A}) = \sqrt{\lambda_i(\mathbf{A} \mathbf{A}^\top)} = \sqrt{\lambda_i(\mathbf{A}^\top \mathbf{A})}$$

Singular Value Decomposition (SVD)

Theorem (best rank- r approximation).

Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ with SVD, $\mathbf{A} = \sum_{i=1}^{\min\{m,n\}} \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$. Then an optimal solution to the best rank- r approximation problem

$$\min_{\mathbf{X}} \|\mathbf{X} - \mathbf{A}\|_2, \quad \text{s.t. } \text{rank}(\mathbf{X}) \leq r$$

is the truncated SVD of the following form

$$\widehat{\mathbf{A}}_r = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^\top, \quad \text{if } \sigma_r > \sigma_{r+1}, \text{ it is unique.}$$

Singular Value Decomposition (SVD)

Proof. Suppose $m \geq n$, then we observe that

$$\left\| \widehat{\mathbf{A}}_r - \mathbf{A} \right\|_2 = \left\| \sum_{i=r+1}^n \sigma_i \mathbf{u}_i \mathbf{v}_i^\top \right\|_2 = \sigma_{r+1}$$

Second, to prove our result, we want to show that

$$\|\mathbf{B}_r - \mathbf{A}\|_2 \geq \sigma_{r+1}, \quad \forall \mathbf{B}_r \text{ with } \text{rank}(\mathbf{B}_r) = r.$$

We suppose $\mathbf{B}_r = \widetilde{\mathbf{U}} \widetilde{\mathbf{V}}^\top$, $\widetilde{\mathbf{U}} \in \mathbb{R}^{m \times r}$, $\widetilde{\mathbf{V}} \in \mathbb{R}^{n \times r}$.

$$\begin{aligned} \|\mathbf{B}_r - \mathbf{A}\|_2^2 &= \sup_{\|\mathbf{z}\|_2 \leq 1} \|(\mathbf{B}_r - \mathbf{A}) \mathbf{z}\|_2^2 \\ &\geq \|(\mathbf{B}_r - \mathbf{A}) \mathbf{w}\|_2^2 \end{aligned}$$

Singular Value Decomposition (SVD)

Since $\widetilde{\mathbf{V}}$ has r columns, then there must be a nontrivial linear combination of the first $r+1$ columns of \mathbf{V} , i.e.,

$$\mathbf{w} = \gamma_1 \mathbf{v}_1 + \dots + \gamma_{r+1} \mathbf{v}_{r+1}$$

such that $\widetilde{\mathbf{V}}^\top \mathbf{w} = \mathbf{0}$. Wlog, we can scale $\|\mathbf{w}\|_2 = 1$ and $\sum_{i=1}^{r+1} \gamma_i^2 = 1$

- Therefore, we have

$$\|\mathbf{A} - \mathbf{B}_k\|_2^2 \geq \|(\mathbf{A} - \mathbf{B}_k) \mathbf{w}\|_2^2 = \|\mathbf{A} \mathbf{w}\|_2^2 = \gamma_1^2 \sigma_1^2 + \dots + \gamma_{r+1}^2 \sigma_{r+1}^2 \geq \sigma_{r+1}^2.$$

Matrix & Operator Norms

Definition. A vector norm is any real-valued function $\|\cdot\|$ that satisfies the following properties

- if $\mathbf{x} \neq \mathbf{0}$, then $\|\mathbf{x}\| > 0$
- for any $\alpha \neq 0$, then $\|\alpha \mathbf{x}\| = |\alpha| \|\mathbf{x}\|$
- triangle inequality: $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$.

There are many norms, for example

$$\|\mathbf{x}\|_p = (\sum_i |x_i|^p)^{1/p}, \quad \text{for any } p \geq 1.$$

Matrix & Operator Norms

Definition (matrix operator norm). Let $\mathbf{A} \in \mathbb{R}^{m \times n}$. If $\|\cdot\|_a$ and $\|\cdot\|_b$ are norms on \mathbb{R}^n and \mathbb{R}^m , we have

$$\|\mathbf{A}\|_{a \rightarrow b} := \sup_{\|\mathbf{x}\|_a \leq 1} \|\mathbf{A}\mathbf{x}\|_b.$$

It satisfies the *three criteria* for norm, and is *submultiplicative*

$$\|\mathbf{AB}\|_{a \rightarrow b} \leq \|\mathbf{A}\|_{a \rightarrow b} \cdot \|\mathbf{B}\|_{a \rightarrow b}.$$

Matrix & Operator Norms

Definition (matrix operator norm). Let $\mathbf{A} \in \mathbb{R}^{m \times n}$. If $\|\cdot\|_a$ and $\|\cdot\|_b$ are norms on \mathbb{R}^n and \mathbb{R}^m , we have

$$\|\mathbf{A}\|_{a \rightarrow b} := \sup_{\|\mathbf{x}\|_a \leq 1} \|\mathbf{A}\mathbf{x}\|_b.$$

- $\|\mathbf{A}\|_{2 \rightarrow 2} = \sigma_1(\mathbf{A})$ (**spectral norm**, write as $\|\mathbf{A}\|$);
- $\|\mathbf{A}\|_{1 \rightarrow b} = \max_{j=1, \dots, n} \|\mathbf{A}\mathbf{e}_j\|_b$
- $\|\mathbf{A}\|_{a \rightarrow \infty} = \max_{i=1, \dots, m} \|\mathbf{e}_i^* \mathbf{A}\|_b^*$, $\|\mathbf{v}\|_b^* := \sup_{\|\mathbf{u}\|_b \leq 1} \langle \mathbf{u}, \mathbf{v} \rangle$.

Matrix & Operator Norms

Definition (unitary invariant matrix norms).

Let $\mathbf{A} \in \mathbb{R}^{m \times n}$. We say the matrix norm is unitary invariant if

$$\|\mathbf{A}\|_{\sharp} = \|\mathbf{PAQ}\|_{\sharp}, \quad \forall \mathbf{P} \in O(m), \mathbf{Q} \in O(n).$$

For example:

- **Spectral norm.** $\|\mathbf{A}\|_{2 \rightarrow 2} = \sigma_1(\mathbf{A}) = \|\sigma(\mathbf{A})\|_{\infty}$
- **Frobenius norm.**

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^{\min\{m,n\}} \sigma_i^2(\mathbf{A})} = \|\sigma(\mathbf{A})\|_2$$

Matrix & Operator Norms

Definition (Schatten p -norm).

Let $\mathbf{A} \in \mathbb{R}^{m \times n}$. For any $p \in [1, +\infty)$, the function

$$\|\mathbf{A}\|_{S_p} := \|\sigma(\mathbf{A})\|_p$$

is a norm on $\mathbb{R}^{m \times n}$.

Special case of great interest: **nuclear/trace norm**

$$\|\mathbf{A}\|_* := \|\mathbf{A}\|_{S_1} = \sum_i \sigma_i(\mathbf{A}) = \|\sigma(\mathbf{A})\|_1$$

Lecture Agenda

- Basics of Matrix Analysis
- **Taylor Expansion & Lipschitz Function**
- Optimality Conditions
- Rate of Convergence

Lipschitz Continuity

Suppose that

- $f : \mathcal{X} \mapsto \mathcal{Y}$ with \mathcal{X} and \mathcal{Y} being open sets;
- $\|\cdot\|_{\mathcal{X}}$ and $\|\cdot\|_{\mathcal{Y}}$ are norms on \mathcal{X} and \mathcal{Y} , respectively.

Definition. $f(\cdot)$ is Lipschitz continuous over \mathcal{X} if $\exists L < \infty$, such that

$$\|f(\mathbf{y}) - f(\mathbf{x})\|_{\mathcal{Y}} \leq L(\mathbf{x}) \|\mathbf{y} - \mathbf{x}\|_{\mathcal{X}}, \quad \forall \mathbf{x}, \mathbf{y} \in \mathcal{X}.$$

Example I: ℓ_1 -norm

ℓ^1 -norm: $f(\mathbf{x}) = \|\mathbf{x}\|_1$

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq \|\mathbf{x} - \mathbf{y}\|_1 \leq \sqrt{n} \|\mathbf{x} - \mathbf{y}\|_2$$

- A subgradient of $f(\mathbf{x}) = \|\mathbf{x}\|_1$ is not Lipschitz;
- The function $f(\mathbf{x}) = \|\mathbf{x}\|_1$ is *not* continuously differentiable, but Lipschitz continuous

Lipschitz Continuity

Let f and g be Lipschitz continuous functions with best Lipschitz constants L_f and L_g , respectively.

$h(\mathbf{x})$	L_h	
$\alpha f(\mathbf{x}) + \beta$	$ \alpha L_f$	scale/shift
$f(\mathbf{x} - \mathbf{x}_0)$	L_f	translate
$f(\mathbf{x}) + g(\mathbf{x})$	$\leq L_f + L_g$	add
$f(g(\mathbf{x}))$	$\leq L_f L_g$	compose (HW)
$\mathbf{Ax} + \mathbf{b}$	$\ \mathbf{A}\ $	affine (for same norm on \mathbb{F}^M and \mathbb{F}^N)
$f(\mathbf{x})g(\mathbf{x})$?	multiply

Smooth Functions

Definition. A differentiable function $f(\mathbf{x})$ is called **smooth** iff it has a *Lipschitz continuous gradient*, i.e., iff $L < +\infty$ such that

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{z})\|_2 \leq L \|\mathbf{x} - \mathbf{z}\|_2, \quad \forall \mathbf{x}, \mathbf{z} \in \mathbb{R}^n$$

Note: Lipschitz continuity of ∇f is a stronger condition than mere continuity, so any differentiable function whose gradient is Lipschitz continuous is a **continuously differentiable** function.

Mean Value Theorem and Lipschitzness

- If $f : \mathbb{R}^n \mapsto \mathbb{R}$ is continuously differentiable,

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq \sup_{\mathbf{z} \in \mathbb{R}^n} \|\nabla f(\mathbf{z})\|_2 \|\mathbf{x} - \mathbf{y}\|_2$$

so that the Lipschitz constant of f is $\sup_{\mathbf{z} \in \mathbb{R}^n} \|\nabla f(\mathbf{z})\|_2$

- If $f : \mathbb{R}^n \mapsto \mathbb{R}$ is twice continuously differentiable,

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq \sup_{\mathbf{z} \in \mathbb{R}^n} \|\nabla^2 f(\mathbf{z})\|_2 \|\mathbf{x} - \mathbf{y}\|_2$$

so that the Lipschitz constant of ∇f is $\sup_{\mathbf{z} \in \mathbb{R}^n} \|\nabla^2 f(\mathbf{z})\|_2$

Mean Value Theorem I

Theorem. Let $f : \mathbb{R}^n \mapsto \mathbb{R}$ be continuously differentiable. For any fixed \mathbf{x} and \mathbf{y} , we have

$$f(\mathbf{y}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{z}(t_L)), \mathbf{y} - \mathbf{x} \rangle$$

for some $\mathbf{z}(t_L) = (1 - t_L) \cdot \mathbf{x} + t_L \cdot \mathbf{y}$ with $t_L \in (0, 1)$.

Mean Value Theorem and Lipschitzness

- If $f : \mathbb{R}^n \mapsto \mathbb{R}$ is continuously differentiable,

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq \sup_{\mathbf{z} \in \mathbb{R}^n} \|\nabla f(\mathbf{z})\|_2 \|\mathbf{x} - \mathbf{y}\|_2$$

so that the Lipschitz constant of f is $\sup_{\mathbf{z} \in \mathbb{R}^n} \|\nabla f(\mathbf{z})\|_2$

- If $f : \mathbb{R}^n \mapsto \mathbb{R}$ is twice continuously differentiable,

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq \sup_{\mathbf{z} \in \mathbb{R}^n} \|\nabla^2 f(\mathbf{z})\|_2 \|\mathbf{x} - \mathbf{y}\|_2$$

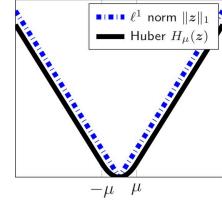
so that the Lipschitz constant of ∇f is $\sup_{\mathbf{z} \in \mathbb{R}^n} \|\nabla^2 f(\mathbf{z})\|_2$

Example II: Huber Function

Huber function: $H_\mu(\mathbf{x}) = \sum_{i=1}^n h_\mu(x_i)$, $h_\mu(x) = \begin{cases} |x| & |x| \geq \mu \\ \frac{x^2}{2\mu} + \frac{\mu}{2} & |x| < \mu \end{cases}$

By mean value theory,

$$\begin{aligned} |h_\mu(x) - h_\mu(y)| &\leq \sup_z |h'_\mu(z)| |x - y| \\ &\leq |x - y| \\ |h'_\mu(x) - h'_\mu(y)| &\leq \sup_z \left| \frac{\partial^2}{\partial z^2} h_\mu(z) \right| |x - y| \\ &\leq \frac{1}{\mu} |x - y| \end{aligned}$$



Mean Value Theorem II

Theorem. Let $f : \mathbb{R}^n \mapsto \mathbb{R}$ be twice continuously differentiable. For any fixed \mathbf{x} and \mathbf{y} , we have

$$\nabla f(\mathbf{y}) = \nabla f(\mathbf{x}) + \int_0^1 \nabla^2 f(\mathbf{z}(t)) \cdot (\mathbf{y} - \mathbf{x}) dt$$

for some $\mathbf{z}(t) = (1 - t)\mathbf{x} + t\mathbf{y}$ with $t \in (0, 1)$, and that

$$f(\mathbf{y}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{1}{2} (\mathbf{y} - \mathbf{x})^\top \nabla^2 f(\mathbf{z}(t_Q)) (\mathbf{y} - \mathbf{x})$$

for some $t_Q \in (0, 1)$.

Taylor Approximation Theory I

Theorem (first-order approximation).

Let $f : \mathbb{R}^n \mapsto \mathbb{R}$ be continuously differentiable, and $\nabla f(\mathbf{x}_0)$ is $\gamma_L(\mathbf{x}_0)$ -Lipschitz continuous at \mathbf{x}_0 , then

$$|f(\mathbf{x}) - \hat{f}_L(\mathbf{x}; \mathbf{x}_0)| \leq \frac{\gamma_L(\mathbf{x}_0)}{2} \|\mathbf{x} - \mathbf{x}_0\|_2^2,$$

where we define

$$\hat{f}_L(\mathbf{x}; \mathbf{x}_0) := f(\mathbf{x}_0) + \langle \nabla f(\mathbf{x}_0), \mathbf{x} - \mathbf{x}_0 \rangle.$$

Taylor Approximation Theory II

Theorem (second-order approximation).

Let $f : \mathbb{R}^n \mapsto \mathbb{R}$ be twice continuously differentiable, and suppose $\nabla^2 f(\mathbf{x})$ is $\gamma_Q(\mathbf{x}_0)$ -Lipschitz continuous at \mathbf{x}_0 , then

$$|f(\mathbf{x}) - \hat{f}_Q(\mathbf{x}; \mathbf{x}_0)| \leq \frac{\gamma_Q(\mathbf{x}_0)}{6} \|\mathbf{x} - \mathbf{x}_0\|_2^3,$$

where we define

$$\begin{aligned} \hat{f}_Q(\mathbf{x}; \mathbf{x}_0) &:= f(\mathbf{x}_0) + \langle \nabla f(\mathbf{x}_0), \mathbf{x} - \mathbf{x}_0 \rangle \\ &+ \frac{1}{2} (\mathbf{x} - \mathbf{x}_0)^\top \nabla^2 f(\mathbf{x}_0) (\mathbf{x} - \mathbf{x}_0). \end{aligned}$$

Lecture Agenda

- Basics of Matrix Analysis
- Taylor Expansion & Lipschitz Function
- Optimality Conditions**
- Rate of Convergence

Optimality: Unconstrained Problem

Consider a continuously differentiable function $f : \mathbb{R}^n \mapsto \mathbb{R}$

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$$

Definition. A point \mathbf{x}_* is a *stationary point* of $f(\cdot)$, if

$$\nabla f(\mathbf{x}_*) = \mathbf{0}.$$

Here, \mathbf{x}_* is also known as a *critical point*.

40

UNIVERSITY OF MICHIGAN

41

UNIVERSITY OF MICHIGAN

Optimality: Unconstrained Problem

- Second-order *necessary* condition for optimality

If \mathbf{x}_* is a local minimizer of $f : \mathbb{R}^n \mapsto \mathbb{R}$, and f is *twice* continuously differentiable in an open neighborhood around \mathbf{x}_* , then we must have

- (i) $\nabla f(\mathbf{x}_*) = \mathbf{0};$
- (ii) $\nabla^2 f(\mathbf{x}_*) \succeq \mathbf{0}.$

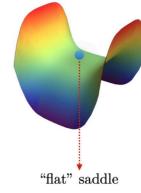
42

UNIVERSITY OF MICHIGAN

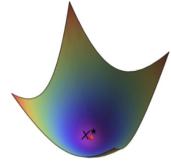
43

UNIVERSITY OF MICHIGAN

Optimality: Unconstrained Problem



$$f(\mathbf{x}) = x_1^3 - x_2^3$$



$$f(\mathbf{x}) = \mathbf{x}^\top \mathbf{A} \mathbf{x}, \quad \mathbf{A} \succ \mathbf{0}.$$

Optimality: Unconstraint Problem

- Second-order *sufficient* condition for optimality

If $f : \mathbb{R}^n \mapsto \mathbb{R}$ is *twice* continuously differentiable in an open neighborhood around \mathbf{x}_* and that

- (i) $\nabla f(\mathbf{x}_*) = \mathbf{0};$
- (ii) $\nabla^2 f(\mathbf{x}_*) \succ \mathbf{0},$

44

UNIVERSITY OF MICHIGAN

45

UNIVERSITY OF MICHIGAN

Global Optimality of Convex Functions

Let $f : \mathbb{R}^n \mapsto \mathbb{R}$ be a *convex* function, then

- A local minimizer of f is also its global minimizer. If f is *strictly convex*, the global minimizer is *unique*.
- A point is a global minimizer of f iff
 $\mathbf{0} \in \partial f(\mathbf{x}_*).$

If $f \in \mathcal{C}^1$, then $\nabla f(\mathbf{x}_*) = \mathbf{0}$ implies that \mathbf{x}_* is a *global* minimizer.

Optimality: Constrained Problem

Consider a smooth contrained problem with

$$\begin{aligned} \min_{\mathbf{x}} f(\mathbf{x}), \quad \text{s.t. } r_i(\mathbf{x}) &= 0, \quad 1 \leq i \leq p, \\ h_j(\mathbf{x}) &\leq 0, \quad 1 \leq j \leq q. \end{aligned}$$

Consider its *Lagrangian* function

$$\mathcal{L}(\mathbf{x}, \mathbf{u}, \mathbf{v}) = f(\mathbf{x}) + \sum_{i=1}^p u_i \cdot r_i(\mathbf{x}) + \sum_{j=1}^q v_j \cdot h_j(\mathbf{x}).$$

46

UNIVERSITY OF MICHIGAN

Optimality: Constraint Problem

First-order necessary condition (aka KKT condition)

- Stationary**

$$\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}_*, \mathbf{u}_*, \mathbf{v}_*) = \mathbf{0}$$

- Feasibility**

$$\begin{array}{ll} (\text{Primal}) & r_i(\mathbf{x}_*) = 0, \quad 1 \leq i \leq p, \\ & h_j(\mathbf{x}_*) \leq 0, \quad 1 \leq j \leq q. \end{array} \quad \begin{array}{ll} (\text{Dual}) & \mathbf{v}_* \geq \mathbf{0}. \end{array}$$

- Complimentary Slackness**

$$h_j(\mathbf{x}_*) \cdot v_{*j} = 0, \quad \forall 1 \leq j \leq q.$$

47

UNIVERSITY OF MICHIGAN

Optimality: Constraint Problem

- KKT condition is only a **necessary condition** for general constraint nonlinear problems;
- For *convex problems* under certain constraint qualifications, it is also a *sufficient* condition.

48

UNIVERSITY OF MICHIGAN

Lecture Agenda

- Basics of Matrix Analysis
- Taylor Expansion & Lipschitz Function
- Optimality Conditions
- **Rate of Convergence**

49

UNIVERSITY OF MICHIGAN

Algorithmic Convergence

$$\min_{\mathbf{x}} f(\mathbf{x}), \quad \text{s.t. } \mathbf{x} \in \mathcal{C}.$$

Solve the problem via iterative methods of optimization, which produce a sequence of points

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k, \dots$$

starting from an initialization \mathbf{x}_0 .

50

UNIVERSITY OF MICHIGAN

Algorithmic Convergence

$$\min_{\mathbf{x}} f(\mathbf{x}), \quad \text{s.t. } \mathbf{x} \in \mathcal{C}.$$

- \mathcal{X} is the *set of all optimal solutions* \mathbf{x}_* ;
- $f_* = f(\mathbf{x}_*)$ is the *optimal function value*;
- ε is the user-defined *precision value*.

51

UNIVERSITY OF MICHIGAN

Algorithmic Convergence

$$\min_{\mathbf{x}} f(\mathbf{x}), \quad \text{s.t. } \mathbf{x} \in \mathcal{C}.$$

- **Total computation time:**

$$\text{complexity} = \text{per iter cost} \times \#\text{of iterations.}$$

- **Per iteration cost:** how much computation it takes to generate the next point \mathbf{x}_{k+1} from \mathbf{x}_k .

52

UNIVERSITY OF MICHIGAN

Algorithmic Convergence

$$\text{complexity} = \text{per iter cost} \times \#\text{of iterations.}$$

Convergence rate. How quickly the sequence $\{\mathbf{x}_k\}_{k \geq 1}$ converges, measured by

Distance to a minimizer	$\ \mathbf{x}_k - \mathbf{x}_*\ _2 \leq \varepsilon$
Sub-optimality in objective	$ f(\mathbf{x}_k) - f(\mathbf{x}_*) \leq \varepsilon$
Gradient	$\ \nabla f(\mathbf{x}_k)\ _2 \leq \varepsilon$

53

UNIVERSITY OF MICHIGAN

Definitions of Convergence

- **Iterate convergence to the (set) minimizer**

$$d(\mathbf{x}_k, \mathcal{X}) = \min_{\mathbf{x}_* \in \mathcal{X}} \|\mathbf{x}_k - \mathbf{x}_*\|_2 \leq \varepsilon.$$

- **Zero-th order function value convergence**

$$|f(\mathbf{x}_k) - f_*| \leq \varepsilon$$

Function value convergence does *not* imply iterate convergence.

54

UNIVERSITY OF MICHIGAN

Definitions of Convergence

- **First-order function value convergence**

$$\|\nabla f(\mathbf{x}_k)\|_2 \leq \varepsilon.$$

- For *convex functions*, this also means convergence to the *global minimizer*
- For *nonconvex functions*, this only means convergence to a *stationary point* (e.g., a local minimizer or a saddle point).

55

UNIVERSITY OF MICHIGAN

Rate of Convergence

Definition. We say the Q -convergence is of order p (≥ 1) and with factor γ (> 0), if $\exists k_0$, such that $\forall k \geq k_0$:

$$\min_{\mathbf{x}_* \in \mathcal{X}} \|\mathbf{x}_k - \mathbf{x}_*\|_2 \leq \gamma \cdot \left(\min_{\mathbf{x}_* \in \mathcal{X}} \|\mathbf{x}_{k-1} - \mathbf{x}_*\|_2 \right)^p.$$

- The larger p is, the faster the convergence.
- The smaller γ is, the faster convergence (with p fixed).
- Typically look for the largest p and the smallest γ such that the inequality holds.

Q -Convergence

The following sequence

$$\{a_k\}_{k \geq 1} = \left\{ 2, 1, \frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{16}, \dots, \frac{1}{2^k}, \dots \right\}$$

- **Q -Linear Convergence** ($p = 1, \gamma < 1$)

$$\|\mathbf{x}_k - \mathbf{x}_*\|_2 \leq \gamma \|\mathbf{x}_{k-1} - \mathbf{x}_*\|_2$$

Q -Convergence

The following sequence is Q -linearly convergent

$$\{a_k\}_{k \geq 1} = \left\{ 2, 1, \frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{16}, \dots, \frac{1}{2^k}, \dots \right\}$$

with

$$\gamma = \lim_{k \rightarrow \infty} \frac{|1/2^{k+1} - 0|}{|1/2^k - 0|} = \lim_{k \rightarrow \infty} \frac{2^k}{2^{k+1}} = \frac{1}{2}.$$

R -Convergence

Q -convergence does not include some sequences that converges reasonably fast:

$$\{b_k\}_{k \geq 1} = \left\{ 1, 1, \frac{1}{4}, \frac{1}{4}, \frac{1}{16}, \frac{1}{16}, \dots, \frac{1}{4^{\lfloor k/2 \rfloor}}, \dots \right\}$$

- **R -convergence** is of order p (≥ 1),

$$\min_{\mathbf{x}_* \in \mathcal{X}} \|\mathbf{x}_k - \mathbf{x}_*\|_2 \leq \rho_k,$$

if the sequence $\{\rho_k\}_{k \geq 1}$ is Q -convergence of the order p .

Rate of Convergence

- **Q -Sublinear Convergence** ($p = 1, \gamma = 1$)

$$\|\mathbf{x}_k - \mathbf{x}_*\|_2 \leq k^{-1/r} \|\mathbf{x}_0 - \mathbf{x}_*\|_2$$

➤ implies k is order $\mathcal{O}(\frac{1}{\varepsilon^r})$.

Rate of Convergence

- **Q -Linear Convergence** ($p = 1, \gamma < 1$)

$$\|\mathbf{x}_k - \mathbf{x}_*\|_2 \leq \gamma \|\mathbf{x}_{k-1} - \mathbf{x}_*\|_2$$

➤ $\gamma < 1$, implies k is order $\mathcal{O}(\log \frac{1}{\varepsilon})$

Rate of Convergence

- **Q -Superlinear Convergence** ($1 < p (\leq 2)$)

$$\frac{\|\mathbf{x}_k - \mathbf{x}_*\|_2}{\|\mathbf{x}_{k-1} - \mathbf{x}_*\|_2} \rightarrow 0$$

Rate of Convergence

- **Quadratic Convergence** ($p = 2$)

$$\|\mathbf{x}_k - \mathbf{x}_*\|_2 \leq \gamma \|\mathbf{x}_{k-1} - \mathbf{x}_*\|_2^2$$

- implies k is order $\mathcal{O}(\log \log \frac{1}{\varepsilon})$;
- $\gamma > 0$ here does not need to be smaller than 1;
- implies superlinear convergence;

Rate of Convergence

Methods	Convex Nonsmooth	Convex smooth	Strongly Convex	Strict Saddle
Subgradient Method	$O(1/\sqrt{k})$			
Gradient Descent		$O(1/k^2)$	Linear	$O(1/k^2)$
Proximal Method	$O(1/k^2)$		Linear	
Quasi-Newton Method			Superlinear	
Newton Method			Quadratic	Quadratic

References & Readings

- *Numerical Optimization*, Jorge Nocedal, and Stephen Wright, Springer. ([Chapter 2](#))
- *High-Dimensional Data Analysis with Low-Dimensional Models: Principles, Computation, and Applications*. John Wright, Yi Ma. ([Appendix A-C](#))