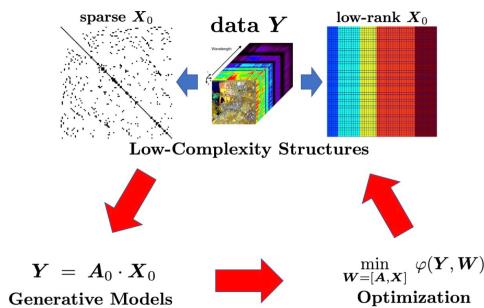


## Lecture Agenda

- Motivating Examples
- Subgradient and Subdifferential Revisit
- Subgradient Methods

### Recover Low-Dimensional Structures



2/13/24

3

UNIVERSITY OF MICHIGAN

### Convex Nonsmooth Problems

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

- Here, the function  $f(\mathbf{x}) : \mathbb{R}^n \mapsto \mathbb{R}$  is *convex*, and the set  $\mathcal{C}$  is a *convex set* or just  $\mathcal{C} = \mathbb{R}^n$ ;
- The function  $f(\mathbf{x}) : \mathbb{R}^n \mapsto \mathbb{R}$  is *nonsmooth* (i.e., the gradient is not well-defined everywhere).

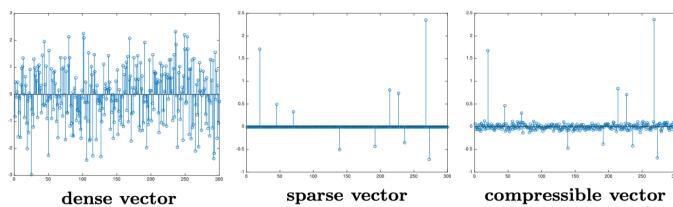
2/13/24

2/13/24

4

UNIVERSITY OF MICHIGAN

### Example I: Sparse Recovery

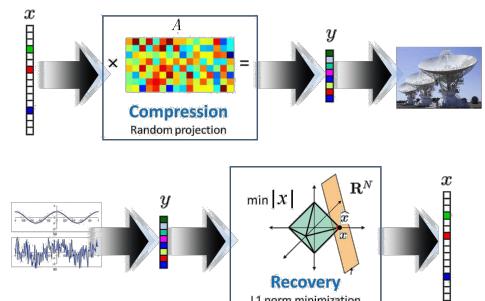


2/13/24

5

UNIVERSITY OF MICHIGAN

### Example I: Sparse Recovery



2/13/24

2/13/24

6

UNIVERSITY OF MICHIGAN

### Example I: Sparse Recovery

$$\mathbf{y}_i = \mathbf{y}_{\text{clean}} + \mathbf{z}_i = \underset{\text{patch dictionary}}{\mathbf{A}} \cdot \underset{\text{sparse coefficient}}{\mathbf{x}_i} + \mathbf{z}_i$$



noisy image

denoised image

dictionary for image patches

### Example I: Sparse Recovery

**Problem:** given measurement  $\mathbf{y} \in \mathbb{R}^m$  and sensing matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  ( $m \ll n$ ), recover the **sparsest**  $\mathbf{x} \in \mathbb{R}^n$ :

$$\min_{\mathbf{x}} \|\mathbf{x}\|_0, \quad \text{s.t. } \mathbf{y} = \mathbf{A}\mathbf{x}.$$

where we have

$$\|\mathbf{x}\|_0 := \#\{i \mid x(i) \neq 0\} = \sum_i \mathbb{1}_{x(i) \neq 0};$$

## Example I: Sparse Recovery

Convex Relaxation via *basis pursuit* (BP)

$$\min_{\mathbf{x}} \|\mathbf{x}\|_1, \quad \text{s.t.} \quad \mathbf{y} = \mathbf{A}\mathbf{x},$$

Here,  $\|\mathbf{x}\|_1 = \sum_{k=1}^m |x_k|$  denotes the  $\ell^1$ -norm of  $\mathbf{x}$ .

## Example I: Sparse Recovery

*Basis pursuit denoising* (BPDN)

$$\min_{\mathbf{x}} \|\mathbf{x}\|_1, \quad \text{s.t.} \quad \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2 \leq \delta,$$

which is equivalent to

$$\min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \cdot \|\mathbf{x}\|_1,$$

for some properly chosen  $\lambda > 0$ . The problem is termed as *Lasso* (least absolute shrinkage and selection operator).

## Example I: Sparse Recovery

*Basis pursuit denoising* (BPDN)

$$\min_{\mathbf{x}} \|\mathbf{x}\|_1, \quad \text{s.t.} \quad \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2 \leq \delta,$$

which is equivalent to

$$\min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \cdot \|\mathbf{x}\|_1,$$

for some properly chosen  $\lambda > 0$ . The  $\ell_1$ -norm introduces **nonsmoothness** into the optimization problems.

## Example II: 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?

## Example II: Matrix Completion

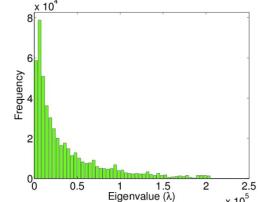
$$\begin{bmatrix} \checkmark & ? & ? & ? & \checkmark & ? \\ ? & ? & \checkmark & \checkmark & ? & ? \\ \checkmark & ? & ? & \checkmark & ? & ? \\ ? & ? & \checkmark & ? & ? & \checkmark \\ \checkmark & ? & ? & ? & ? & ? \\ ? & \checkmark & ? & ? & \checkmark & ? \\ ? & ? & \checkmark & \checkmark & ? & ? \end{bmatrix}$$

More unknowns than observations (under-determined system)

## Example II: Matrix Completion

$$\begin{bmatrix} \checkmark & ? & ? & ? & ? & ? \\ ? & \dots & ? & \dots & ? & ? \\ ? & ? & ? & \dots & ? & ? \\ ? & \dots & \dots & ? & ? & ? \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix} \quad \begin{array}{c} \text{Frequency} \\ \text{Eigenvalue } (\lambda) \end{array}$$

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



## Example II: Matrix Completion

$$\begin{bmatrix} \checkmark & ? & ? & ? & ? & ? \\ ? & \dots & ? & \dots & ? & ? \\ ? & ? & ? & \dots & ? & ? \\ ? & \dots & \dots & ? & ? & ? \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix} \quad \approx \quad \begin{array}{c} \text{Matrix} \\ \text{rank} \end{array}$$

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

## Example II: Matrix Completion

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

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

or, equivalently, we can rewrite it as

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

## Example II: Matrix Completion

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

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

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

2/13/24

17

UNIVERSITY OF MICHIGAN

## Example II: Matrix Completion

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

- Similar to  $\ell_0$ -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  $\text{rank}(\mathbf{X})$

2/13/24

18

UNIVERSITY OF MICHIGAN

## More General: Low-rank Matrix Recovery

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

The matrix completion problem can be extended to a more general form

$$\min_{\mathbf{X}} \|\mathbf{X}\|_*, \quad \text{s.t.} \quad \mathbf{y} = \mathcal{A}(\mathbf{X}).$$

with given linear sensing matrices  $\mathbf{A}_i (1 \leq i \leq m)$

$$\mathbf{y} = \mathcal{A}(\mathbf{M}) := \begin{bmatrix} \text{tr}(\mathbf{A}_1^\top \mathbf{M}) \\ \vdots \\ \text{tr}(\mathbf{A}_m^\top \mathbf{M}) \end{bmatrix}$$

2/13/24

19

UNIVERSITY OF MICHIGAN

## More General: Low-rank Matrix Recovery

Low-rank matrix recovery under noise

$$\min_{\mathbf{X}} \|\mathbf{X}\|_*, \quad \text{s.t.} \quad \|\mathbf{y} - \mathcal{A}(\mathbf{X})\|_2 \leq \delta,$$

is equivalent to

$$\min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{X})\|_2^2 + \lambda \cdot \|\mathbf{X}\|_*,$$

for some properly chosen  $\lambda > 0$ .

- Because  $\|\cdot\|_*$  is **nonsmooth**, both problems are nonsmooth.

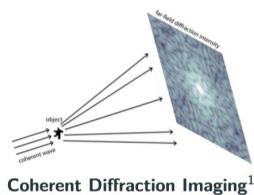
2/13/24

20

UNIVERSITY OF MICHIGAN

## Example III: Phase Retrieval

Electric field  $x(t_1, t_2) \rightarrow$  Fourier transform  $\hat{x}(f_1, f_2)$



**Applications:** X-ray crystallography, diffraction imaging (left), optics, astronomical imaging, and microscopy

**Phase retrieval:** given intensity  $\mathbf{y} = |\mathcal{F}(\mathbf{x})|$ , recover  $\mathbf{x}$ .

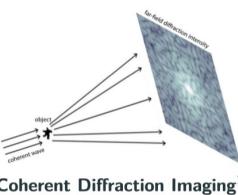
2/13/24

21

UNIVERSITY OF MICHIGAN

## Example III: Phase Retrieval

Electric field  $x(t_1, t_2) \rightarrow$  Fourier transform  $\hat{x}(f_1, f_2)$



**Applications:** X-ray crystallography, diffraction imaging (left), optics, astronomical imaging, and microscopy

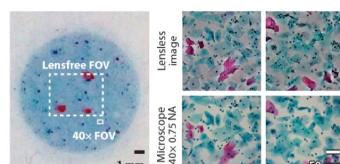
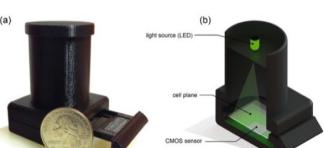
**Generalized phase retrieval:** given intensity  $\mathbf{y} = |\mathbf{A}\mathbf{x}|$  recover  $\mathbf{x}$ , where  $\mathbf{A}$  can be a general sensing matrix.

2/13/24

22

UNIVERSITY OF MICHIGAN

## Example III: Phase Retrieval



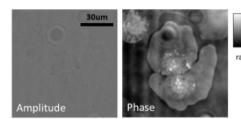
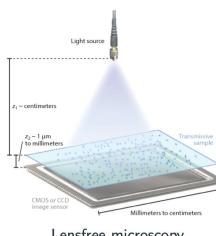
Lensfree microscopy imaging

2/13/24

23

UNIVERSITY OF MICHIGAN

## Example III: Phase Retrieval



Phase contrast microscopy: phase is important<sup>1</sup>.

2/13/24

24

UNIVERSITY OF MICHIGAN

### Example III: Phase Retrieval

$$A \quad x \quad Ax \quad \Rightarrow \quad y = |Ax|^2$$

Solve for  $x \in \mathbb{R}^n$  in  $m$  quadratic equations

$$y_k = |\mathbf{a}_k^\top \mathbf{x}|^2, \quad k = 1, \dots, m,$$

$$\text{or } \mathbf{y} = |\mathbf{Ax}|^2, \quad \text{where } |\mathbf{x}|^2 := [|\mathbf{x}_1|^2, \dots, |\mathbf{x}_m|^2]^\top.$$

2/13/24

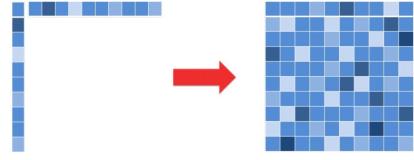
25

UNIVERSITY OF MICHIGAN

### Example III: Phase Retrieval

**Lifting:** introduce  $\mathbf{X} = \mathbf{x}\mathbf{x}^*$  to linearize the problem

$$y_k = |\mathbf{a}_k^* \mathbf{x}|^2 = \underbrace{\mathbf{a}_k^* (\mathbf{x}\mathbf{x}^*)}_{\mathbf{X}} \mathbf{a}_k \implies y_k = \langle \mathbf{a}_k \mathbf{a}_k^*, \mathbf{X} \rangle$$



2/13/24

2/13/24

26

UNIVERSITY OF MICHIGAN

### Example III: Phase Retrieval

$$\text{find } \mathbf{X} \succeq \mathbf{0}, \quad \text{s.t. } y_k = \langle \mathbf{a}_k \mathbf{a}_k^*, \mathbf{X} \rangle, \quad k = 1, \dots, m \\ \text{rank}(\mathbf{X}) = 1$$

- Convex relaxation:

$$\min_{\mathbf{X}} \|\mathbf{X}\|_*, \quad \text{s.t. } \mathbf{y} = \mathcal{A}(\mathbf{X}), \quad \mathbf{X} \succeq \mathbf{0}.$$

- Stable version:

$$\min_{\mathbf{X} \succeq \mathbf{0}} \frac{1}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{X})\|_2^2 + \lambda \cdot \|\mathbf{X}\|_*.$$

2/13/24

27

UNIVERSITY OF MICHIGAN

### Example III: Sparse Phase Retrieval

- Convex relaxation:

$$\min_{\mathbf{X}} \|\mathbf{X}\|_* + \lambda \|\mathbf{X}\|_1, \quad \text{s.t. } \mathbf{y} = \mathcal{A}(\mathbf{X}), \quad \mathbf{X} \succeq \mathbf{0}.$$

- Stable version:

$$\min_{\mathbf{X} \succeq \mathbf{0}} \|\mathbf{X}\|_* + \lambda \|\mathbf{X}\|_1 + \frac{\mu}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{X})\|_2^2.$$

2/13/24

2/13/24

28

UNIVERSITY OF MICHIGAN

### Example IV: Robust PCA

$$\text{observation } \mathbf{Y} = \text{low-rank } \mathbf{L} + \text{sparse } \mathbf{S}$$

observation  $\mathbf{Y}$

low-rank  $\mathbf{L}$

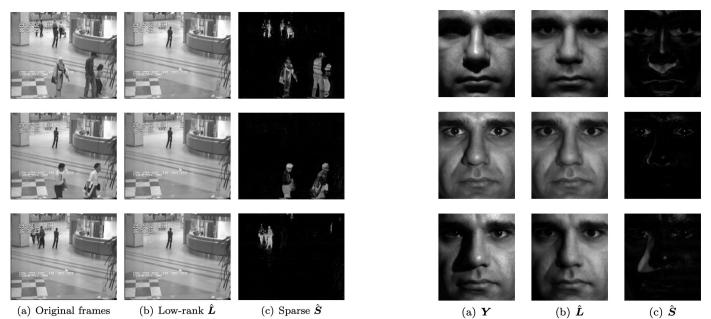
sparse  $\mathbf{S}$

2/13/24

29

UNIVERSITY OF MICHIGAN

### Example IV: Robust PCA



2/13/24

2/13/24

30

UNIVERSITY OF MICHIGAN

### Example IV: Robust PCA

- Convex relaxation:

$$\min_{\mathbf{L}, \mathbf{S}} \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1, \quad \text{s.t. } \mathbf{Y} = \mathbf{L} + \mathbf{S}.$$

- Stable version:

$$\min_{\mathbf{L}, \mathbf{S}} \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1 + \frac{\mu}{2} \|\mathbf{Y} - \mathbf{L} - \mathbf{S}\|_F^2.$$

2/13/24

31

UNIVERSITY OF MICHIGAN

### Further Readings

- *High-Dimensional Data Analysis with Low-Dimensional Models: Principles, Computation, and Applications*. John Wright, Yi Ma. (Chapter 3-5)

32

UNIVERSITY OF MICHIGAN

# Lecture Agenda

- Motivating Examples
- Subgradient and Subdifferential Revisit
- Subgradient Methods

2/13/24

33

UNIVERSITY OF MICHIGAN

## (Projected) Subgradient Methods

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

Practically, a natural choice is “subgradient-based methods”

$$\mathbf{x}_{k+1} = \mathcal{P}_{\mathcal{C}}(\mathbf{x}_k - \tau_k \cdot \mathbf{g}_k),$$

where  $\mathbf{g}_k \in \partial f(\mathbf{x}_k)$  is any subgradient of  $f(\cdot)$  at  $\mathbf{x}_k$ .

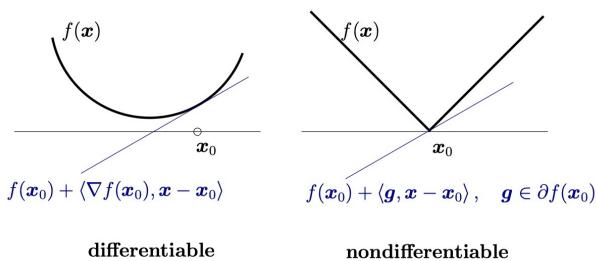
- **Caution:** unlike gradient descent method, a subgradient direction might *not* be a descent direction on  $f$ .

2/13/24

34

UNIVERSITY OF MICHIGAN

## Basics: Subdifferential & Subgradient



2/13/24

35

UNIVERSITY OF MICHIGAN

## Basics: Subgradient & Subdifferential

- **Subgradient.** Let  $f : \mathbb{R}^n \mapsto \mathbb{R}$  be *convex*. A *subgradient* of  $f$  at  $\mathbf{x}_0$  is any  $\mathbf{u}$  satisfying

$$f(\mathbf{x}) \geq f(\mathbf{x}_0) + \langle \mathbf{u}, \mathbf{x} - \mathbf{x}_0 \rangle, \quad \forall \mathbf{x}.$$

- **Subdifferential.** It is the *set of all subgradients* of  $f$  at  $\mathbf{x}_0$ 

$$\partial f(\mathbf{x}_0) := \{\mathbf{u} \mid f(\mathbf{x}) \geq f(\mathbf{x}_0) + \langle \mathbf{u}, \mathbf{x} - \mathbf{x}_0 \rangle, \forall \mathbf{x} \in \mathbb{R}^n\}.$$

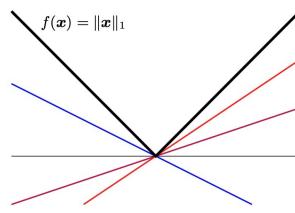
If  $f$  is differentiable at  $\mathbf{x}_0$ , then  $\partial f(\mathbf{x}) = \{\nabla f(\mathbf{x}_0)\}$ .

36

UNIVERSITY OF MICHIGAN

## Example: $f(x) = |x|$

$$\partial f(x) = \begin{cases} \{1\} & x > 0, \\ [-1, 1] & x = 0, \\ \{-1\} & x < 0. \end{cases}$$

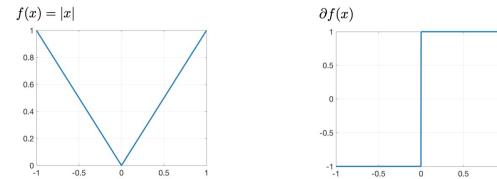


2/13/24

37

UNIVERSITY OF MICHIGAN

## Example: $f(x) = |x|$

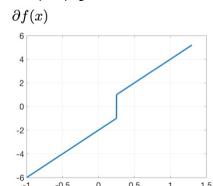
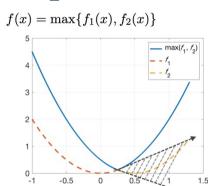


$$f(x) = |x| \quad \partial f(x) = \begin{cases} \{1\} & x > 0, \\ [-1, 1] & x = 0, \\ \{-1\} & x < 0. \end{cases}$$

38

UNIVERSITY OF MICHIGAN

## Example: $\max\{f_1(x), f_2(x)\}$



$\max\{f_1(x), f_2(x)\}$ , where  $f_1$  and  $f_2$  are differentiable

$$\partial f(x) = \begin{cases} \{f'_1(x)\}, & \text{if } f_1(x) > f_2(x) \\ [f'_1(x), f'_2(x)], & \text{if } f_1(x) = f_2(x) \\ \{f'_2(x)\}, & \text{if } f_1(x) < f_2(x) \end{cases}$$

2/13/24

39

UNIVERSITY OF MICHIGAN

## A Subgradient of Norms at 0

**Theorem.** Let  $f(\mathbf{x}) = \|\mathbf{x}\|$  for any norm  $\|\cdot\|$ , then for any  $\mathbf{g}$  obeying  $\|\mathbf{g}\|_* \leq 1$ , we have

$$\mathbf{g} \in \partial f(\mathbf{0}),$$

where  $\|\cdot\|_*$  is the dual norm of  $\|\cdot\|$  with

$$\|\mathbf{x}\|_* := \sup_{\mathbf{z}: \|\mathbf{z}\| \leq 1} \langle \mathbf{z}, \mathbf{x} \rangle$$

Here,  $\|\cdot\|_*$  is not unclear norm.

40

UNIVERSITY OF MICHIGAN

## A Subgradient of Norms at 0

**Proof.** To show this, it suffices to show that for  $f(\mathbf{z}) = \|\mathbf{z}\|$ ,

$$\forall \mathbf{z} \in \mathbb{R}^n \quad f(\mathbf{z}) \geq f(\mathbf{0}) + \langle \mathbf{g}, \mathbf{z} - \mathbf{0} \rangle, \quad \|\mathbf{g}\|_* \leq 1$$

$$\iff \|\mathbf{z}\| \geq \langle \mathbf{g}, \mathbf{z} \rangle \quad \|\mathbf{g}\|_* \leq 1$$

To prove the above, notice that

$$\|\mathbf{g}\|_* = \sup_{\|\mathbf{z}\| \leq 1} \langle \mathbf{g}, \mathbf{z} \rangle = \sup_{\mathbf{z}} \left\langle \mathbf{g}, \frac{\mathbf{z}}{\|\mathbf{z}\|} \right\rangle \geq \left\langle \mathbf{g}, \frac{\mathbf{z}}{\|\mathbf{z}\|} \right\rangle, \quad \forall \mathbf{z}.$$

$$\implies \langle \mathbf{g}, \mathbf{z} \rangle \leq \|\mathbf{g}\|_* \|\mathbf{z}\| \leq \|\mathbf{z}\|$$

2/13/24

41

UNIVERSITY OF MICHIGAN

## Basic Rules for Subdifferential

- **Scaling:**  $\partial(\alpha f) = \alpha \partial f$  (for  $\alpha > 0$ )

- **Summation:**  $\partial(f_1 + f_2) = \partial f_1 + \partial f_2$

2/13/24

2/13/24

42

UNIVERSITY OF MICHIGAN

### Example: the $\ell_1$ -Norm

$$f(\mathbf{x}) = \|\mathbf{x}\|_1 = \sum_{i=1}^n \underbrace{|x_i|}_{=: f_i(\mathbf{x})}$$

Since

$$\partial f_i(\mathbf{x}) = \begin{cases} \text{sign}(x_i) \mathbf{e}_i, & \text{if } x_i \neq 0, \\ [-1, 1] \cdot \mathbf{e}_i & \text{if } x_i = 0. \end{cases}$$

We have

$$\partial f(\mathbf{x}) = \sum_{i: x_i \neq 0} \text{sign}(x_i) \mathbf{e}_i + \sum_{j: x_j = 0} [-1, 1] \cdot \mathbf{e}_j$$

In particular,  $\sum_{i: x_i \neq 0} \text{sign}(x_i) \mathbf{e}_i \in \partial f(\mathbf{x})$ .

2/13/24

43

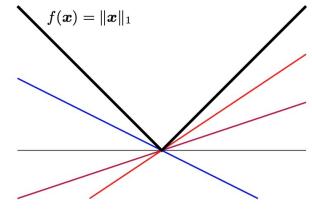
UNIVERSITY OF MICHIGAN

### Example: $f(\mathbf{x}) = \|\mathbf{x}\|_1, \mathbf{x} \in \mathbb{R}^n$

$$\partial f(\mathbf{x}) = \mathcal{J}_1 \times \cdots \times \mathcal{J}_n,$$

where for  $k = 1, 2, \dots, n$ ,

$$\mathcal{J}_k = \begin{cases} \{1\} & x_k > 0, \\ [-1, 1] & x_k = 0, \\ \{-1\} & x_k < 0. \end{cases}$$



44

UNIVERSITY OF MICHIGAN

## Basic Rules for Subdifferential

- **Scaling:**  $\partial(\alpha f) = \alpha \partial f$  (for  $\alpha > 0$ )
- **Summation:**  $\partial(f_1 + f_2) = \partial f_1 + \partial f_2$
- **Affine transformation:** if  $h(\mathbf{x}) = f(\mathbf{Ax} + \mathbf{b})$ , then  

$$\partial h(\mathbf{x}) = \mathbf{A}^\top \partial f(\mathbf{Ax} + \mathbf{b})$$

2/13/24

45

UNIVERSITY OF MICHIGAN

### Example: $h(\mathbf{x}) = \|\mathbf{Ax} + \mathbf{b}\|_1$

Let  $f(\mathbf{x}) = \|\mathbf{x}\|_1$  and  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_m]^\top$ , we have

$$\mathbf{g} = \sum_{i: \mathbf{a}_i^\top \mathbf{x} + b_i \neq 0} \text{sign}(\mathbf{a}_i^\top \mathbf{x} + b_i) \cdot \mathbf{e}_i \in \partial f(\mathbf{Ax} + \mathbf{b})$$

$$\implies \mathbf{A}^\top \mathbf{g} = \sum_{i: \mathbf{a}_i^\top \mathbf{x} + b_i \neq 0} \text{sign}(\mathbf{a}_i^\top \mathbf{x} + b_i) \cdot \mathbf{a}_i \in \partial h(\mathbf{x}).$$

46

UNIVERSITY OF MICHIGAN

## Basic Rules for Subdifferential

- **Chain rule:** suppose  $f : \mathbb{R}^n \mapsto \mathbb{R}$  is convex, and  $g : \mathbb{R} \mapsto \mathbb{R}$  is differentiable, nondecreasing, and convex. Let  $h = f \circ g$ ,  

$$\partial h(\mathbf{x}) = g'(f(\mathbf{x})) \cdot \partial f(\mathbf{x})$$
- **Composition:** suppose  $f(\mathbf{x}) = h(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$ , where  $f_i$ 's are convex, and  $h$  is differentiable, nondecreasing, and convex. Let  $\mathbf{q} = \nabla h(\mathbf{y}) \mid_{\mathbf{y}=[f_1(\mathbf{x}), \dots, f_n(\mathbf{x})]}$  and  $\mathbf{g}_i \in \partial f_i(\mathbf{x})$ , then  

$$\sum_{i=1}^n q_i \mathbf{g}_i \in \partial f(\mathbf{x}).$$

2/13/24

47

UNIVERSITY OF MICHIGAN

## Basic Rules for Subdifferential

- **Pointwise maximum:** if  $f(\mathbf{x}) = \max_{1 \leq i \leq k} f_i(\mathbf{x})$ , then

$$\partial f(\mathbf{x}) = \underbrace{\text{conv} \{ \cup \{ \partial f_i(\mathbf{x}) \mid f_i(\mathbf{x}) = f(\mathbf{x}) \} \}}_{\text{convex hull of subdifferentials of all active functions}}$$

- **Pointwise supremum:** if  $f(\mathbf{x}) = \sup_{\alpha \in \mathcal{F}} f_\alpha(\mathbf{x})$ , then

$$\partial f(\mathbf{x}) = \text{closure} (\text{conv} \{ \cup_{\alpha \in \mathcal{F}} \{ \partial f_\alpha(\mathbf{x}) \mid f_\alpha(\mathbf{x}) = f(\mathbf{x}) \} \})$$

48

UNIVERSITY OF MICHIGAN

## Example: Piece-wise Linear Functions

$$f(\mathbf{x}) = \max_{1 \leq i \leq m} \{ \mathbf{a}_i^\top \mathbf{x} + b_i \}$$

At any given  $\mathbf{x}_0$ , pick any  $\mathbf{a}_j$  such that

$$\mathbf{a}_j^\top \mathbf{x}_0 + b_j = \max_{1 \leq i \leq m} \{ \mathbf{a}_i^\top \mathbf{x}_0 + b_i \}$$

then

$$\mathbf{a}_j \in \partial f(\mathbf{x}_0).$$

2/13/24

49

UNIVERSITY OF MICHIGAN

## Example: the $\ell_\infty$ -Norm

$$f(\mathbf{x}) = \|\mathbf{x}\|_\infty = \max_{1 \leq i \leq n} |\mathbf{x}_i|.$$

At any given  $\mathbf{x}'$ , pick any  $x'_j$  such that

$$|x'_j| = \max_{1 \leq i \leq m} |x'_i|$$

then

$$\text{sign}(x'_j) \cdot \mathbf{e}_j \in \partial f(\mathbf{x}')$$

2/13/24

50

UNIVERSITY OF MICHIGAN

## Example: the Maximum Eigenvalue

$$f(\mathbf{x}) = \lambda_{\max} \left( \sum_{i=1}^n x_i \mathbf{A}_i \right)$$

where  $\mathbf{A}_1, \dots, \mathbf{A}_n$  are real symmetric matrices. Rewrite

$$f(\mathbf{x}) = \sup_{\|\mathbf{v}\|=1} \mathbf{v}^\top \left( \sum_{i=1}^n x_i \mathbf{A}_i \right) \mathbf{v}.$$

For any given  $\mathbf{x}$ , taking  $\mathbf{v}_0$  as the leading eigenvector of  $\sum_i x_i \mathbf{A}_i$

$$[\mathbf{v}_0^\top \mathbf{A}_1 \mathbf{v}_0 \quad \dots \quad \mathbf{v}_0^\top \mathbf{A}_n \mathbf{v}_0]^\top \in \partial f(\mathbf{x})$$

2/13/24

51

UNIVERSITY OF MICHIGAN

## Example: the Nuclear Norm

Let  $\mathbf{X} \in \mathbb{R}^{m \times n}$  with SVD  $\mathbf{X} = \mathbf{U} \Sigma \mathbf{V}^\top$  and

$$f(\mathbf{X}) = \|\mathbf{X}\|_* := \sum_{i=1}^{\min\{m,n\}} \sigma_i(\mathbf{X})$$

Then we have

$$\mathbf{U} \mathbf{V}^\top \in \partial f(\mathbf{X}).$$

More generally, we can show that

$$\partial \|\mathbf{X}\|_* = \{ \mathbf{U} \mathbf{V}^\top + \mathbf{W} \mid \|\mathbf{W}\| \leq 1, \mathbf{W} \perp \text{col}(\mathbf{U}), \mathbf{W} \perp \text{row}(\mathbf{V}) \}.$$

2/13/24

52

UNIVERSITY OF MICHIGAN

## Example: the Nuclear Norm

**Proof.** First, notice the fact that

$$\|\mathbf{X}\|_* = \sup_{\|\mathbf{Z}\| \leq 1} \langle \mathbf{X}, \mathbf{Z} \rangle$$

Therefore, we have

$$\|\mathbf{X}\|_* = \sup_{\|\mathbf{X}\| \leq 1} \langle \mathbf{Z}, \mathbf{U} \Sigma_X \mathbf{V}^\top \rangle = \sup_{\|\mathbf{X}\| \leq 1} \langle \mathbf{U}^\top \mathbf{Z} \mathbf{V}, \Sigma_X \rangle$$

where the supreme is achieved when  $\mathbf{Z} = \mathbf{U} \mathbf{V}^\top$ , so that  $\|\mathbf{X}\|_* = \text{tr}(\Sigma_X)$

Thus, let  $f_Z(\mathbf{X}) = \langle \mathbf{Z}, \mathbf{X} \rangle$ , we have  $\|\mathbf{X}\|_* = \sup_{\|\mathbf{Z}\| \leq 1} f_Z(\mathbf{X})$ . By the property of subgradient, we have

$$\nabla_{\mathbf{X}} f_Z(\mathbf{X}) \mid_{\mathbf{Y}=\mathbf{U} \mathbf{V}^\top} \in \partial \|\mathbf{X}\|_* \implies \mathbf{U} \mathbf{V}^\top \in \partial \|\mathbf{X}\|_*$$

2/13/24

53

UNIVERSITY OF MICHIGAN

## Optimality Condition for Convex Nonsmooth Problem

**Theorem.** For any convex function  $f$ ,

$$f(\mathbf{x}_*) = \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \iff \mathbf{0} \in \partial f(\mathbf{x}_*)$$

that  $\mathbf{x}_*$  is a minimizer iff  $\mathbf{0}$  is a subgradient of  $f$  at  $\mathbf{x}_*$ .

**Reason:** for any  $\mathbf{y}$ , by the definition of subgradient

$$f(\mathbf{y}) \geq f(\mathbf{x}_*) + \mathbf{0}^\top (\mathbf{y} - \mathbf{x}_*)$$

2/13/24

54

UNIVERSITY OF MICHIGAN

## Soft-thresholding Operator

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) = \frac{1}{2} \|\mathbf{x} - \mathbf{z}\|_2^2 + \lambda \|\mathbf{x}\|_1$$

**Claim.** An optimal solution  $\mathbf{x}_* = \arg \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$ , is the soft-thresholding operator  $\mathbf{x}_* = S_\lambda(\mathbf{z})$  with

$$[S_\lambda(\mathbf{z})]_i = \begin{cases} z_i - \lambda & \text{if } z_i > \lambda, \\ 0 & \text{if } |z_i| \leq \lambda, \\ z_i + \lambda & \text{if } z_i \leq -\lambda. \end{cases}$$

In summary,  $S_\lambda(\mathbf{z}) = \text{sign}(\mathbf{z}) \odot (|\mathbf{z}| - \lambda)_+$ .

2/13/24

55

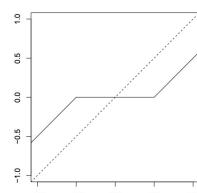
UNIVERSITY OF MICHIGAN

## Soft-thresholding Operator

**Proof.** Notice that  $\partial f(\mathbf{x}) = (\mathbf{x} - \mathbf{z}) + \lambda \partial \|\mathbf{x}\|_1$ , we find the optimal solution  $\mathbf{x}_*$  by solving  $\mathbf{0} \in \partial f(\mathbf{x}_*)$

Given that  $[\partial f(\mathbf{x}_*)]_i = [\mathbf{x}_* - \mathbf{z} + \lambda \partial \|\mathbf{x}_*\|_1]_i$

- If  $z_i > \lambda$ ,  $0 \in x_{*i} - z_i + \lambda [\partial \|\mathbf{x}_*\|_1]_i \implies x_{*i} = z_i - \lambda$
- If  $z_i < -\lambda$ ,  $0 \in x_{*i} - z_i + \lambda [\partial \|\mathbf{x}_*\|_1]_i \implies x_{*i} = z_i + \lambda$
- If  $|z_i| \leq \lambda$ ,  $0 \in x_{*i} - z_i + \lambda [\partial \|\mathbf{x}_*\|_1]_i \implies x_{*i} = 0$



2/13/24

56

UNIVERSITY OF MICHIGAN

# Singular Value Thresholding

**Theorem.** The unique solution  $\mathbf{X}_*$  to the program

$$\min_{\mathbf{X}} \frac{1}{2} \|\mathbf{X} - \mathbf{Z}\|_F^2 + \lambda \|\mathbf{X}\|_*,$$

is given by

$$\mathbf{X}_* = \mathbf{U} \mathbf{S}_\lambda[\Sigma] \mathbf{V}^\top,$$

where  $\mathbf{Z} = \mathbf{U} \Sigma \mathbf{V}^\top$  and  $\mathbf{S}_\lambda(\cdot)$  is the soft-thresholding operator.

2/13/24

57

UNIVERSITY OF MICHIGAN

2/13/24

58

UNIVERSITY OF MICHIGAN

# Lecture Agenda

- Motivating Examples
- Subgradient and Subdifferential Revisit
- Subgradient Methods**

## (Projected) Subgradient Methods

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

Practically, a natural choice is “subgradient-based methods”

$$\mathbf{x}_{k+1} = \mathcal{P}_{\mathcal{C}}(\mathbf{x}_k - \tau_k \cdot \mathbf{g}_k),$$

where  $\mathbf{g}_k \in \partial f(\mathbf{x}_k)$  is any subgradient of  $f(\cdot)$  at  $\mathbf{x}_k$ .

- Caution:** unlike gradient descent method, a subgradient direction might *not* be a descent direction on  $f$ .

2/13/24

59

UNIVERSITY OF MICHIGAN

## Negative Subgradient Not a Descent Direction

**Example:**  $f(\mathbf{x}) = |x_1| + 3|x_2|$

at  $\mathbf{x} = (1, 0)$ :

- $\mathbf{g}_1 = (1, 0) \in \partial f(\mathbf{x})$ , and  $-\mathbf{g}_1$  is a descent direction;
- $\mathbf{g}_2 = (1, 3) \in \partial f(\mathbf{x})$ , but  $-\mathbf{g}_2$  is *not* a descent direction!

**Reason:** lack of continuity. One can change directions significantly without violating the validity of subgradients.

60

UNIVERSITY OF MICHIGAN

## Subgradient Methods

- Since  $f(\mathbf{x}_k)$  is not necessarily monotone, we keep track of the best point

$$f_{\text{best},k} := \min_{1 \leq i \leq k} f(\mathbf{x}_i)$$

- At each step  $k$ , we set

$$f_{\text{best},k} = \min \{f_{\text{best},k-1}, f(\mathbf{x}_k)\}.$$

2/13/24

61

UNIVERSITY OF MICHIGAN

## Convergence of Subgradient Methods

$$\min_{\mathbf{x} \text{ nonsmooth}} f(\mathbf{x}), \quad \text{s.t. } \mathbf{x} \in \mathbb{R}^n.$$

**Lemma.** Suppose  $f : \mathbb{R}^n \mapsto \mathbb{R}$  is  $L$ -Lipschitz, then the iterate

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \tau_k \cdot \mathbf{g}_k,$$

satisfies

$$f_{\text{best},k} - f_* \leq \frac{R^2 + L^2 \sum_{i=1}^k \tau_i^2}{2 \sum_{i=1}^k \tau_i},$$

where  $R \geq \|\mathbf{x}_0 - \mathbf{x}_*\|_2$ .

2/13/24      Subgradient Methods, Notes for EE364b, Stanford University. Stephen Boyd. (Section 3, page 6)

62

UNIVERSITY OF MICHIGAN

## Convergence of Subgradient Methods

**Proof.** Let us denote  $f_k = f(\mathbf{x}_k)$  and  $f_* = f(\mathbf{x}_*)$

$$\begin{aligned} 0 \leq \|\mathbf{x}_{k+1} - \mathbf{x}_*\|_2^2 &= \|\mathbf{x}_k - \mathbf{x}_* - \tau_k \mathbf{g}_k\|_2^2 \\ &= \|\mathbf{x}_k - \mathbf{x}_*\|_2^2 - 2 \langle \tau_k \mathbf{g}_k, \mathbf{x}_k - \mathbf{x}_* \rangle + \tau_k^2 \|\mathbf{g}_k\|_2^2 \\ &\leq \|\mathbf{x}_k - \mathbf{x}_*\|_2^2 - 2\tau_k(f_k - f_*) + \tau_k^2 \|\mathbf{g}_k\|_2^2 \end{aligned}$$

which has been derived by the def. of subgradient

$$f_* \geq f_k + \mathbf{g}_k^\top (\mathbf{x}_* - \mathbf{x}_k) \quad \rightarrow \quad f_* - f_k \geq \mathbf{g}_k^\top (\mathbf{x}_* - \mathbf{x}_k)$$

Therefore, recursively, we have

$$0 \leq \|\mathbf{x}_{k+1} - \mathbf{x}_*\|_2^2 \leq \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2 - 2 \sum_{i=0}^k \tau_i(f_i - f_*) + \sum_{i=0}^k \tau_i^2 \|\mathbf{g}_i\|_2^2$$

2/13/24

63

UNIVERSITY OF MICHIGAN

## Convergence of Subgradient Methods

- On one hand:

$$2 \sum_{i=0}^k \tau_i(f_i - f_*) \leq R^2 + \sum_{i=0}^k \tau_i^2 \|\mathbf{g}_i\|_2^2 \leq R^2 + L^2 \sum_{i=0}^k \tau_i^2$$

- On the other hand:

$$\sum_{i=0}^k \tau_i(f_i - f_*) \geq \left( \sum_{i=0}^k \tau_i \right) \min_{0 \leq i \leq k} (f_i - f_*)$$

Combining the above together, we obtain

$$f_{\text{best},k} - f_* = \min_{0 \leq i \leq k} (f_i - f_*) \leq \frac{1}{2 \sum_{i=0}^k \tau_i} \left( R^2 + L^2 \sum_{i=0}^k \tau_i^2 \right)$$

64

UNIVERSITY OF MICHIGAN

## Step Size Rules

$$f_{\text{best},k} - f_* \leq \frac{R^2 + L^2 \sum_{i=1}^k \tau_i^2}{2 \sum_{i=1}^k \tau_i}$$

- Constant step size  $\tau_k \equiv \tau$

$$\lim_{k \rightarrow \infty} f_{\text{best},k} - f_* \leq \frac{L^2 \tau}{2}.$$

- Diminishing step size with  $\sum_k \tau_k^2 < +\infty$  and  $\sum_k \tau_k \rightarrow +\infty$

$$\lim_{k \rightarrow \infty} f_{\text{best},k} - f_* = 0.$$

2/13/24

65

UNIVERSITY OF MICHIGAN

## Optimal Step Size for Subgradient Method?

**Theorem.** Suppose  $f : \mathbb{R}^n \mapsto \mathbb{R}$  is  $L$ -Lipschitz, then choose a constant step size

$$\tau_i = R/(L\sqrt{k}), \quad i = 1, \dots, k$$

gives

$$f_{\text{best},k} - f_* \leq \frac{RL}{\sqrt{k}}$$

where  $R \geq \|\mathbf{x}_0 - \mathbf{x}_*\|_2$ .

2/13/24    Subgradient Methods, Notes for EE364b, Stanford University. Stephen Boyd. (Section 3, page 6)

UNIVERSITY OF MICHIGAN

## Optimal Step Size for Subgradient Method?

**Proof.** We already know that

$$f_{\text{best},k} - f_* \leq h(\tau) := \frac{R^2 + L^2 \sum_{i=0}^k \tau_i^2}{2 \sum_{i=0}^k \tau_i}$$

and we can show that  $h(\tau)$  is convex and symmetric in  $\tau = \{\tau_1, \dots, \tau_k\}$

The optimal solution for  $\min_\tau h(\tau)$  is  $\tau_1 = \dots = \tau_k = \frac{R}{L\sqrt{k}}$

Plugging the solution in, we obtain

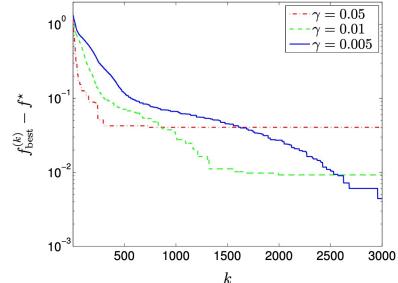
$$f_{\text{best},k} - f_* \leq h^*(\tau_*) = \frac{RL}{\sqrt{k}}$$

2/13/24

67

UNIVERSITY OF MICHIGAN

## Experiments: Subgradient Methods

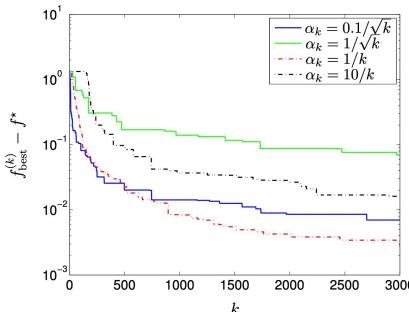


2/13/24    Subgradient Methods, Notes for EE364b, Stanford University. Stephen Boyd. (Section 4, page 12)

UNIVERSITY OF MICHIGAN

## Experiments: Subgradient Methods

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) = \|\mathbf{Ax} + \mathbf{b}\|_\infty$$



2/13/24

69

UNIVERSITY OF MICHIGAN

## Optimal Step Size with Known $f_*$

**Theorem. (Polyak's Stepsize Rule)**

Suppose  $f : \mathbb{R}^n \mapsto \mathbb{R}$  is  $L$ -Lipschitz, and the optimal is  $f_*$  known. Choose stepsizes that satisfy

$$\tau_k = \frac{f_k - f_*}{\|\mathbf{g}_k\|_2^2}.$$

Then, we have

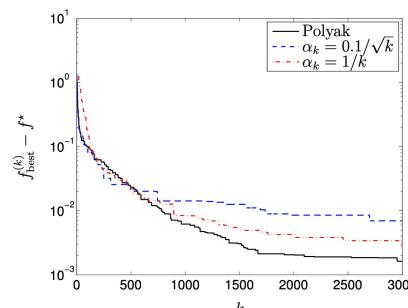
$$f_{\text{best},k} - f_* \leq \frac{RL}{\sqrt{k}}.$$

2/13/24    Subgradient Methods, Notes for EE364b, Stanford University. Stephen Boyd. (Section 4, page 12)

UNIVERSITY OF MICHIGAN

## Experiments: Subgradient Methods

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) = \|\mathbf{Ax} + \mathbf{b}\|_\infty$$



2/13/24

71

UNIVERSITY OF MICHIGAN

## Strongly Convex Functions

**Theorem.** Let  $f$  be  $\mu$ -strongly convex and  $L$ -Lipschitz continuous over  $\mathbb{R}^n$ . Choose the stepsize as  $\tau_k \equiv \tau = \frac{2}{\mu(k+1)}$ , then

$$f_{\text{best},k} - f_* \leq \frac{2L^2}{\mu} \frac{1}{k+1}.$$

2/13/24

72

UNIVERSITY OF MICHIGAN

## Summary: Subgradient Methods

	Stepsize Rule	Convergence Rate	Iteration Complexity
Convex & Lipschitz Problems	$\tau_k \asymp \frac{1}{\sqrt{k}}$	$O\left(\frac{1}{\sqrt{k}}\right)$	$O\left(\frac{1}{\varepsilon^2}\right)$
Strongly Convex & Lipschitz Problems	$\tau_k \asymp \frac{1}{k}$	$O\left(\frac{1}{k}\right)$	$O\left(\frac{1}{\varepsilon}\right)$

Similar results can be shown for the projected version.

## Further Readings

- *Subgradient Methods, Notes for EE364b, Stanford University.* Stephen Boyd.  
[https://stanford.edu/class/ee364b/lectures/subgrad\\_method\\_notes.pdf](https://stanford.edu/class/ee364b/lectures/subgrad_method_notes.pdf)