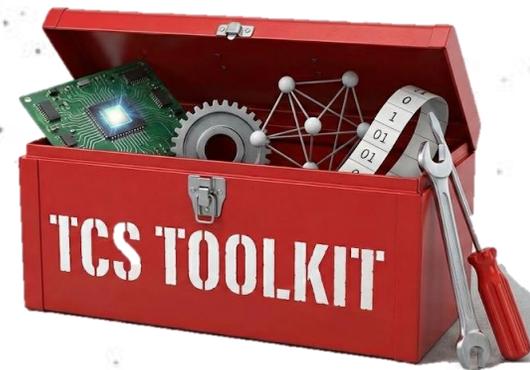


# CS 58500 – Theoretical Computer Science Toolkit

Lecture 11 (02/24)

Gradient Descent

[https://ruizhezhang.com/course\\_spring\\_2026.html](https://ruizhezhang.com/course_spring_2026.html)



# Today's Lecture

- Lipschitz Optimization
- Smooth Optimization
- Well-Conditioned Optimization

# Lipschitz Optimization

- $f: \mathcal{X} \rightarrow \mathbb{R}$  is  **$L$ -Lipschitz** if

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq L \cdot \|\mathbf{x} - \mathbf{y}\|_2 \quad \forall \mathbf{x}, \mathbf{y} \in \mathcal{X}$$

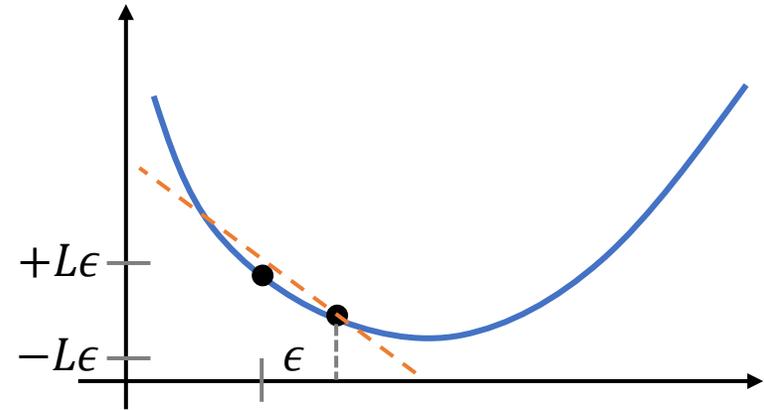
- $\mathbf{g}$  is a **subgradient** of  $f$  at  $\mathbf{x} \in \mathcal{X}$  (denoted as  $\mathbf{g} \in \partial f(\mathbf{x})$ ) if

$$f(\mathbf{x}') \geq f(\mathbf{x}) + \langle \mathbf{g}, \mathbf{x}' - \mathbf{x} \rangle \quad \forall \mathbf{x}' \in \mathcal{X}$$

**Lemma (Subgradient bound).** Let  $f: \mathcal{X} \rightarrow \mathbb{R}$  be convex and  $L$ -Lipschitz for  $\mathcal{X} \subseteq \mathbb{R}^n$ . Then, for  $\mathbf{x} \in \text{relint}(\mathcal{X})$  and  $\mathbf{g} \in \partial f(\mathbf{x})$  contained in the lowest-dimensional subspace containing  $\mathcal{X}$ , we have  $\|\mathbf{g}\|_2 \leq L$ . Moreover, the converse direction holds as well

*Proof.*

- Suppose  $\exists \mathbf{x} \in \mathcal{X}, \mathbf{g} \in \partial f(\mathbf{x})$  such that  $\|\mathbf{g}\|_2 > L$
- Consider  $\mathbf{y} := \mathbf{x} + \epsilon \mathbf{g} \in \mathcal{X}$ .
- $L$ -Lipschitz implies that  $f(\mathbf{y}) \leq f(\mathbf{x}) + L\|\mathbf{y} - \mathbf{x}\|_2 = f(\mathbf{x}) + L\epsilon\|\mathbf{g}\|_2$
- Subgradient implies that  $f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \mathbf{g}, \mathbf{y} - \mathbf{x} \rangle = f(\mathbf{x}) + \langle \mathbf{g}, \epsilon \mathbf{g} \rangle > f(\mathbf{x}) + L\epsilon\|\mathbf{g}\|_2$  (contradiction!)



# Lipschitz Optimization: Hardness

**Theorem 1 (Lipschitz convex lower bound).** Let  $\epsilon, L, R > 0$ . No algorithm  $\mathcal{A}$  which accesses a target  $L$ -Lipschitz, convex function  $f: RB_2^n \rightarrow \mathbb{R}$  using a subgradient oracle  $O$  and produces iterates  $\{\mathbf{x}_t\}_{0 \leq t < T}$  can optimize  $f$  to additive error  $\epsilon$  using  $T < \min \left\{ n, \left( \frac{LR}{4\epsilon} \right)^2 \right\}$  queries, subject to the restriction that

$$\mathbf{x}_0 = \mathbf{0}, \quad \mathbf{x}_t \in \text{span}(\{O(\mathbf{x}_s)\}_{0 \leq s < t}) \quad \forall t \in [T]$$

*Proof.*

- Define  $f(\mathbf{x}) := \gamma \max_{i \in [T]} \mathbf{x}_i + \frac{\alpha}{2} \|\mathbf{x}\|_2^2$  with  $\alpha, \gamma > 0$  to be chosen later
- $f$  is convex and  $(\gamma + \alpha R)$ -Lipschitz ( $\gamma + \alpha R \leq L$ )
- $\partial f(\mathbf{x}) = \gamma \cdot \text{conv}(\{\mathbf{e}_i : i \in \arg \max_{j \in [T]} \mathbf{x}_j\}) + \alpha \mathbf{x}$
- Let  $O(\mathbf{x}) = \gamma \mathbf{e}_i + \alpha \mathbf{x}$  with the smallest  $i \in \arg \max_{j \in [T]} \mathbf{x}_j$

The subgradient of a max function at  $\mathbf{x}$  is the convex hull of the gradients of the “active” functions at  $\mathbf{x}$

# Lipschitz Optimization: Hardness

**Theorem 1 (Lipschitz convex lower bound).** Let  $\epsilon, L, R > 0$ . No algorithm  $\mathcal{A}$  which accesses a target  $L$ -Lipschitz, convex function  $f: RB_2^n \rightarrow \mathbb{R}$  using a subgradient oracle  $O$  and produces iterates  $\{\mathbf{x}_t\}_{0 \leq t < T}$  can optimize  $f$  to additive error  $\epsilon$  using  $T < \min \left\{ n, \left( \frac{LR}{4\epsilon} \right)^2 \right\}$  queries, subject to the restriction that

$$\mathbf{x}_0 = \mathbf{0}, \quad \mathbf{x}_t \in \text{span}(\{O(\mathbf{x}_s)\}_{0 \leq s < t}) \quad \forall t \in [T]$$

*Proof.*

- Define  $f(\mathbf{x}) := \gamma \max_{i \in [T]} x_i + \frac{\alpha}{2} \|\mathbf{x}\|_2^2$  with  $\alpha, \gamma > 0$  to be chosen later
- Let  $O(\mathbf{x}) = \gamma \mathbf{e}_i + \alpha \mathbf{x}$  with the smallest  $i \in \arg \max_{j \in [T]} x_j$
- $O(\mathbf{x}_0) = (\gamma, 0, \dots, 0)$ ,  $\mathbf{x}_1 = (?, 0, \dots, 0)$ ,  $O(\mathbf{x}_1) = (?, ?, 0, \dots, 0)$
- For all  $0 \leq t < T$ ,  $\mathbf{x}_t$  is only supported in the first  $t$  coordinate and  $f(\mathbf{x}_t) \geq 0$

# Lipschitz Optimization: Hardness

**Theorem 1 (Lipschitz convex lower bound).** Let  $\epsilon, L, R > 0$ . No algorithm  $\mathcal{A}$  which accesses a target  $L$ -Lipschitz, convex function  $f: RB_2^n \rightarrow \mathbb{R}$  using a subgradient oracle  $O$  and produces iterates  $\{\mathbf{x}_t\}_{0 \leq t < T}$  can optimize  $f$  to additive error  $\epsilon$  using  $T < \min \left\{ n, \left( \frac{LR}{4\epsilon} \right)^2 \right\}$  queries, subject to the restriction that

$$\mathbf{x}_0 = \mathbf{0}, \quad \mathbf{x}_t \in \text{span}(\{O(\mathbf{x}_s)\}_{0 \leq s < t}) \quad \forall t \in [T]$$

*Proof.*

- Define  $f(\mathbf{x}) := \gamma \max_{i \in [T]} x_i + \frac{\alpha}{2} \|\mathbf{x}\|_2^2$  with  $\alpha, \gamma > 0$  to be chosen later
- $\mathbf{x}^* := \sum_{i \in [T]} -\frac{\gamma}{\alpha T} \mathbf{e}_i$
- $f(\mathbf{x}^*) = -\frac{\gamma^2}{2\alpha T}$  and  $\partial f(\mathbf{x}^*) = \sum_{i \in [T]} -\frac{\gamma}{T} \mathbf{e}_i + \gamma \cdot \text{conv}(\{\mathbf{e}_i : i \in [T]\}) \ni \mathbf{0}$
- Thus,  $\mathbf{x}^*$  is a **minimizer** if  $\|\mathbf{x}^*\|_2 \leq R \iff \gamma^2 \leq R^2 \alpha^2 T$

# Lipschitz Optimization: Hardness

**Theorem 1 (Lipschitz convex lower bound).** Let  $\epsilon, L, R > 0$ . No algorithm  $\mathcal{A}$  which accesses a target  $L$ -Lipschitz, convex function  $f: RB_2^n \rightarrow \mathbb{R}$  using a subgradient oracle  $O$  and produces iterates  $\{\mathbf{x}_t\}_{0 \leq t < T}$  can optimize  $f$  to additive error  $\epsilon$  using  $T < \min \left\{ n, \left( \frac{LR}{4\epsilon} \right)^2 \right\}$  queries, subject to the restriction that

$$\mathbf{x}_0 = \mathbf{0}, \quad \mathbf{x}_t \in \text{span}(\{O(\mathbf{x}_s)\}_{0 \leq s < t}) \quad \forall t \in [T]$$

*Proof.*

- Let  $\gamma := \frac{L}{2}$  and  $\alpha := \frac{L}{2R\sqrt{T}}$ , which satisfy the conditions
- $f(\mathbf{x}^*) = -\frac{\gamma^2}{2\alpha T} = -\frac{LR}{4\sqrt{T}}$ , while for all  $0 \leq t < T$ ,  $f(\mathbf{x}_t) \geq 0$ . That is, the error of  $\mathcal{A}$  is  $\geq \frac{LR}{4\sqrt{T}}$
- For  $\epsilon$ -error, we need  $\frac{LR}{4\sqrt{T}} < \epsilon$ , which gives  $T > \left( \frac{LR}{4\epsilon} \right)^2$



# Lipschitz Optimization: Hardness

**Theorem 1 (Lipschitz convex lower bound).** Let  $\epsilon, L, R > 0$ . No algorithm  $\mathcal{A}$  which accesses a target  $L$ -Lipschitz, convex function  $f: RB_2^n \rightarrow \mathbb{R}$  using a subgradient oracle  $O$  and produces iterates  $\{\mathbf{x}_t\}_{0 \leq t < T}$  can optimize  $f$  to additive error  $\epsilon$  using  $T < \min \left\{ n, \left( \frac{LR}{4\epsilon} \right)^2 \right\}$  queries, subject to the restriction that

$$\mathbf{x}_0 = \mathbf{0}, \quad \mathbf{x}_t \in \text{span}(\{O(\mathbf{x}_s)\}_{0 \leq s < t}) \quad \forall t \in [T]$$

- The **cutting-plane method** has query complexity  $n \log(n/\epsilon)$ , matching the first part of the lower bound
- **Projected gradient descent** can match the second part of the lower bound

# Lipschitz Optimization: Projected gradient descent

Update rule of PGD:

$$\mathbf{g}_t \in \partial f(\mathbf{x}_t), \quad \mathbf{x}_{t+1} = \Pi_{\mathcal{X}}(\mathbf{x}_t - \eta \mathbf{g}_t) \quad \forall 0 \leq t < T$$

Here,  $\eta > 0$  is the step size, and  $\Pi_{\mathcal{X}}$  is the Euclidean projection to  $\mathcal{X}$

**Theorem (PGD).** Let  $f: \mathcal{X} \rightarrow \mathbb{R}$  be convex and  $L$ -Lipschitz for  $\mathcal{X} \subseteq \mathbb{R}B_2^n$ . Suppose  $\mathbf{x}_0 \leftarrow \mathbf{0}$ ,  $\mathbf{x}_{t+1} \leftarrow \Pi_{\mathcal{X}}(\mathbf{x}_t - \eta \mathbf{g}_t)$  with  $\eta = \frac{R}{L\sqrt{T}}$ , and let  $\bar{\mathbf{x}} := \frac{1}{T} \sum_{0 \leq t < T} \mathbf{x}_t$ . Then,

$$f(\bar{\mathbf{x}}) - \min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}) \leq \frac{LR}{\sqrt{T}}$$

➤ For  $\epsilon$ -error,  $T > \left(\frac{LR}{\epsilon}\right)^2$ , which is optimal by the lower bound

# Lipschitz Optimization: Projected gradient descent

**Theorem (PGD).** Let  $f: \mathcal{X} \rightarrow \mathbb{R}$  be convex and  $L$ -Lipschitz for  $\mathcal{X} \subseteq \mathbb{R}B_2^n$ . Suppose  $\mathbf{x}_0 \leftarrow \mathbf{0}$ ,  $\mathbf{x}_{t+1} \leftarrow \Pi_{\mathcal{X}}(\mathbf{x}_t - \eta \mathbf{g}_t)$  with  $\eta = \frac{R}{L\sqrt{T}}$ , and let  $\bar{\mathbf{x}} := \frac{1}{T} \sum_{0 \leq t < T} \mathbf{x}_t$ . Then,

$$f(\bar{\mathbf{x}}) - \min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}) \leq \frac{LR}{\sqrt{T}}$$

*Proof.*

- Our goal is to upper bound  $f(\bar{\mathbf{x}}) - f(\mathbf{x}^*)$ . By convexity, it suffices to upper bound

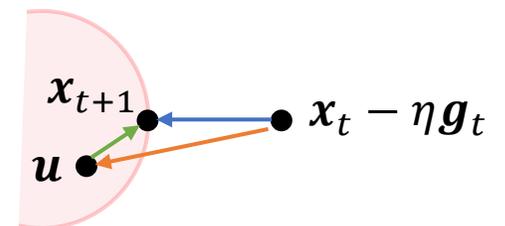
$$f(\bar{\mathbf{x}}) - f(\mathbf{x}^*) \leq \frac{1}{T} \sum_{0 \leq t < T} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq \frac{1}{T} \sum_{0 \leq t < T} \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \quad \text{“potential”}$$

- According to the update rule,

$$\mathbf{x}_{t+1} = \arg \min_{\mathbf{x} \in \mathcal{X}} \frac{1}{2} \|\mathbf{x} - (\mathbf{x}_t - \eta \mathbf{g}_t)\|_2^2 = \arg \min_{\mathbf{x} \in \mathcal{X}} \langle \eta \mathbf{g}_t, \mathbf{x} \rangle + \frac{1}{2} \|\mathbf{x} - \mathbf{x}_t\|_2^2$$

- The **first-order optimality condition** implies that

$$\langle \eta \mathbf{g}_t + (\mathbf{x}_{t+1} - \mathbf{x}_t), \mathbf{x}_{t+1} - \mathbf{u} \rangle \leq 0 \quad \forall \mathbf{u} \in \mathcal{X}$$



# Lipschitz Optimization: Projected gradient descent

$$\langle \eta \mathbf{g}_t + (\mathbf{x}_{t+1} - \mathbf{x}_t), \mathbf{x}_{t+1} - \mathbf{u} \rangle \leq 0 \quad \forall \mathbf{u} \in \mathcal{X}$$

- We have 
$$\begin{aligned} \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle &= \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}_{t+1} \rangle + \langle \mathbf{g}_t, \mathbf{x}_{t+1} - \mathbf{x}^* \rangle \\ &\leq \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}_{t+1} \rangle + \eta^{-1} \langle \mathbf{x}_t - \mathbf{x}_{t+1}, \mathbf{x}_{t+1} - \mathbf{x}^* \rangle \\ &\leq \|\mathbf{g}_t\|_2 \|\mathbf{x}_t - \mathbf{x}_{t+1}\|_2 + \eta^{-1} \langle \mathbf{x}_t - \mathbf{x}_{t+1}, \mathbf{x}_{t+1} - \mathbf{x}^* \rangle \end{aligned}$$

- You can check that the following identity holds for any  $\mathbf{x}_t, \mathbf{x}_{t+1}, \mathbf{u}$ :

$$\langle \mathbf{x}_t - \mathbf{x}_{t+1}, \mathbf{x}_{t+1} - \mathbf{u} \rangle = \frac{1}{2} \|\mathbf{x}_t - \mathbf{u}\|_2^2 - \frac{1}{2} \|\mathbf{x}_{t+1} - \mathbf{u}\|_2^2 - \frac{1}{2} \|\mathbf{x}_t - \mathbf{x}_{t+1}\|_2^2$$

- Thus, we have

$$\begin{aligned} \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle &\leq \|\mathbf{g}_t\|_2 \|\mathbf{x}_t - \mathbf{x}_{t+1}\|_2 - \frac{1}{2\eta} \|\mathbf{x}_t - \mathbf{x}_{t+1}\|_2^2 + \frac{1}{2\eta} \|\mathbf{x}_t - \mathbf{x}^*\|_2^2 - \frac{1}{2\eta} \|\mathbf{x}_{t+1} - \mathbf{x}^*\|_2^2 \\ &\leq \frac{\eta}{2} \|\mathbf{g}_t\|_2^2 + \frac{1}{2\eta} \|\mathbf{x}_t - \mathbf{x}^*\|_2^2 - \frac{1}{2\eta} \|\mathbf{x}_{t+1} - \mathbf{x}^*\|_2^2 \\ &\leq \frac{\eta L^2}{2} + \frac{1}{2\eta} (\|\mathbf{x}_t - \mathbf{x}^*\|_2^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|_2^2) \end{aligned}$$

# Lipschitz Optimization: Projected gradient descent

$$\langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \leq \frac{\eta L^2}{2} + \frac{1}{2\eta} (\|\mathbf{x}_t - \mathbf{x}^*\|_2^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|_2^2)$$

- Now, we can sum for  $0 \leq t < T$ :

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle &\leq \frac{\eta L^2}{2} + \frac{1}{2\eta T} \sum_{t=0}^{T-1} (\|\mathbf{x}_t - \mathbf{x}^*\|_2^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|_2^2) \\ &= \frac{\eta L^2}{2} + \frac{1}{2\eta T} (\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 - \|\mathbf{x}_T - \mathbf{x}^*\|_2^2) \\ &\leq \frac{\eta L^2}{2} + \frac{1}{2\eta T} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 \leq \frac{\eta L^2}{2} + \frac{R^2}{2\eta T} \end{aligned}$$

- If we take  $\eta = \frac{R}{L\sqrt{T}}$ , we obtain that LHS  $\leq LR/\sqrt{T}$
- Thus,  $f(\bar{\mathbf{x}}) - \min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}) \leq LR/\sqrt{T}$



# Lipschitz Optimization

Let the target function  $f: RB_2^n \rightarrow \mathbb{R}$  be convex and  $L$ -Lipschitz

- Lower bound:  $\min \left\{ n, \Omega \left( \frac{L^2 R^2}{\epsilon^2} \right) \right\}$
- Upper bound:
  - CPM:  $\sim n \log(n/\epsilon)$  (high-accuracy)
  - PGD:  $\sim L^2 R^2 / \epsilon^2$  (low-accuracy)
- The lower bound relies on the assumption that the algorithm can only query points within the span of subgradients returned by the oracle
- This assumption can be removed and obtain an information-theoretic lower bound against any randomized algorithm (Agarwal-Bartlett-Ravikumar-Wainwright '12)

# Today's Lecture

- Lipschitz Optimization
- **Smooth Optimization**
- Well-Conditioned Optimization

# Smooth Optimization

- Consider an infinitesimal version of the update rule:

$$\frac{d}{dt} \mathbf{x}_t = \mathbf{v}_t$$

- By chain rule, we have

$$\frac{d}{dt} f(\mathbf{x}_t) = \langle \nabla f(\mathbf{x}_t), \mathbf{v}_t \rangle = -\|\nabla f(\mathbf{x}_t)\|_2^2 \quad \text{if } \mathbf{v}_t = -\nabla f(\mathbf{x}_t)$$

- It indicates that  $-\nabla f(\mathbf{x}_t)$  is the direction of **steepest descent**
- We obtain the ODE (so-called the **gradient flow**):

$$\frac{d}{dt} \mathbf{x}_t = -\nabla f(\mathbf{x}_t)$$

- However, we cannot implement this continuous-time dynamics

# Smooth Optimization

$$\frac{d}{dt} \mathbf{x}_t = -\nabla f(\mathbf{x}_t)$$

- For time-discretization, one simple method is the **forward Euler method**:

$$\frac{d}{dt} \mathbf{x}_t = -\nabla f(\mathbf{x}_{t_0}) \quad \forall t \in [t_0, t_0 + \eta] \quad \Rightarrow \quad \mathbf{x}_{t_0+\eta} = \mathbf{x}_{t_0} - \eta \nabla f(\mathbf{x}_{t_0})$$

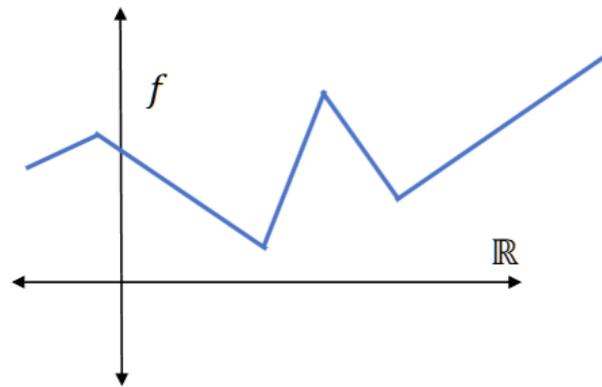
- To guarantee that the discretized dynamics (gradient descent) is close to the continuous time dynamics (gradient flow), we need  **$\nabla f(\mathbf{x})$  to be stable**

# Smooth Optimization

$f: \mathbb{R}^n \rightarrow \mathbb{R}$  is  **$L$ -smooth**, if  $f$  is differentiable, and  $\nabla f$  is  $L$ -Lipschitz

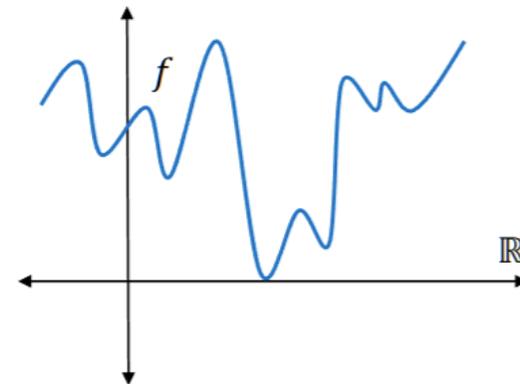
$$|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})| \leq L \cdot \|\mathbf{x} - \mathbf{y}\| \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$$

$f$  is  $L$ -Lipschitz



(bounded slope)  
(bounded 1<sup>st</sup> derivatives)

$f$  is  $L$ -smooth



(bounded curvature)  
(bounded 2<sup>nd</sup> derivative)

# Smooth Optimization

$f: \mathbb{R}^n \rightarrow \mathbb{R}$  is  **$L$ -smooth**, if  $f$  is differentiable, and  $\nabla f$  is  $L$ -Lipschitz

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

**Lemma.** If  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is differentiable and convex, then  $f$  is  $L$ -smooth if and only if

$$f(\mathbf{x}') \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{x}' - \mathbf{x}\|_2^2 \quad \forall \mathbf{x}, \mathbf{x}' \in \mathbb{R}^n$$

If  $f$  is twice-differentiable (and possibly nonconvex),  $f$  is  $L$ -smooth if and only if

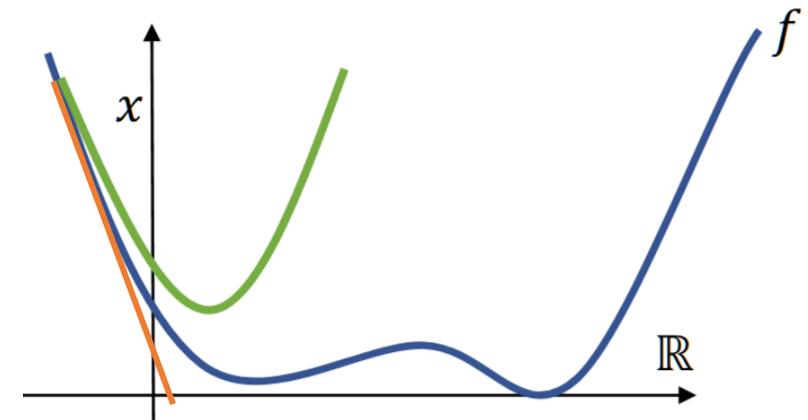
$$|\nabla^2 f(\mathbf{x})[\mathbf{v}, \mathbf{v}]| := |\mathbf{v}^\top \nabla^2 f(\mathbf{x}) \mathbf{v}| \leq L \|\mathbf{v}\|_2^2 \quad \forall \mathbf{v} \in \mathbb{R}^n$$

**Convexity  $\Rightarrow$  linear lower bound**

$$f(\mathbf{x}') \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle$$

**Smoothness  $\Rightarrow$  quadratic upper bound**

$$f(\mathbf{x}') \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{x}' - \mathbf{x}\|_2^2$$



# Smooth Optimization

**Lemma.** Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be  $L$ -smooth. Then for any  $\mathbf{x} \in \mathbb{R}^n$ , letting  $\mathbf{x}' \leftarrow \mathbf{x} - \frac{1}{L} \nabla f(\mathbf{x})$ ,

$$f(\mathbf{x}') - f(\mathbf{x}) \leq -\frac{1}{2L} \|\nabla f(\mathbf{x})\|_2^2$$

*Proof.*

- $L$ -smoothness implies an **upper bound**:

$$f(\mathbf{x}') \leq U_x(\mathbf{x}') := f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{x}' - \mathbf{x}\|_2^2$$

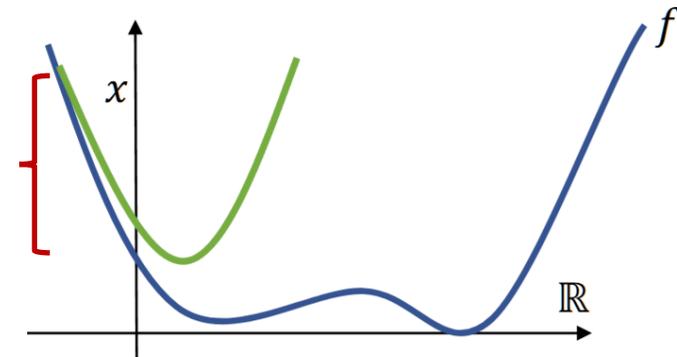
- Thus, we have  $f(\mathbf{x}') - f(\mathbf{x}) \leq U_x(\mathbf{x}') - U_x(\mathbf{x}) = \min_{\mathbf{x}'} U_x(\mathbf{x}') - U_x(\mathbf{x})$

- The **first-order optimality condition** for RHS is

$$\nabla f(\mathbf{x}) + L(\mathbf{x}' - \mathbf{x}) = 0 \iff \mathbf{x}' = \mathbf{x} - L^{-1} \nabla f(\mathbf{x})$$

- Then, we have

$$f(\mathbf{x}') \leq f(\mathbf{x}) - \frac{1}{L} \|\nabla f(\mathbf{x})\|_2^2 + \frac{1}{2L} \|\nabla f(\mathbf{x})\|_2^2 = f(\mathbf{x}) - \frac{1}{2L} \|\nabla f(\mathbf{x})\|_2^2$$



# Smooth Optimization: Finding Stationary Point

**Theorem (Finding stationary point).** Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be  $L$ -smooth, let  $\epsilon > 0$ , and suppose for  $\mathbf{x}_0 \in \mathbb{R}^n$  we have  $f(\mathbf{x}_0) - \min f(\mathbf{x}) \leq \Delta$ . Then iterating  $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t)$  for  $0 \leq t < T$  where  $T \geq \frac{2L\Delta}{\epsilon^2}$ ,

$$\min_{0 \leq t < T} \|\nabla f(\mathbf{x}_t)\|_2 \leq \epsilon$$

*Proof.*

- Suppose the conclusion is not true, i.e.,  $\|\nabla f(\mathbf{x}_t)\|_2 > \epsilon \forall t$ . Then, by the previous lemma,

$$f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) \leq -\frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|_2^2$$

- Summing over  $0 \leq t < T$ , we have

$$f(\mathbf{x}_T) - f(\mathbf{x}_0) \leq -\frac{1}{2L} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|_2^2 < -\frac{T\epsilon^2}{2L} \leq -\Delta$$



# Smooth Optimization: Finding Stationary Point

**Theorem (Finding stationary point).** Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be  $L$ -smooth, let  $\epsilon > 0$ , and suppose for  $\mathbf{x}_0 \in \mathbb{R}^n$  we have  $f(\mathbf{x}_0) - \min f(\mathbf{x}) \leq \Delta$ . Then iterating  $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t)$  for  $0 \leq t < T$  where

$$T \geq \frac{2L\Delta}{\epsilon^2},$$

$$\min_{0 \leq t < T} \|\nabla f(\mathbf{x}_t)\|_2 \leq \epsilon$$

- Finding  $\epsilon$ -stationary/ $\epsilon$ -critical point (i.e.,  $\|\nabla f(\mathbf{x})\|_2 \leq \epsilon$ ) is an important problem in **non-convex optimization**
- The complexity  $\mathcal{O}\left(\frac{L\Delta}{\epsilon^2}\right)$  is optimal (Carmon-Duchi-Hinder-Sidford'20)

# Smooth Optimization: Smooth Gradient Descent

**Theorem (Smooth GD).** Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be  $L$ -smooth and convex, and suppose for  $\mathbf{x}_0 \in \mathbb{R}^n$  we have  $\|\mathbf{x}_0 - \mathbf{x}^*\|_2 \leq R$  for  $\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$ . Then iterating  $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t)$  for  $0 \leq t < T$ ,

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{2LR^2}{T}$$

**Lemma (Contractivity).** Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be  $L$ -smooth and convex, and let  $\mathbf{y} \leftarrow \mathbf{x} - \eta \nabla f(\mathbf{x})$  for  $\eta \leq \frac{1}{L}$ . Then for  $\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$ ,  $\|\mathbf{y} - \mathbf{x}^*\|_2 \leq \|\mathbf{x} - \mathbf{x}^*\|_2$

*Proof.*

- It suffices to prove that  $\|\mathbf{x} - \mathbf{x}^*\|_2^2 \geq \|\mathbf{x} - \eta \nabla f(\mathbf{x}) - \mathbf{x}^*\|_2^2$ , which is

$$\langle \nabla f(\mathbf{x}), \mathbf{x} - \mathbf{x}^* \rangle \geq \frac{\eta}{2} \|\nabla f(\mathbf{x})\|_2^2$$

- By convexity,  $\langle \nabla f(\mathbf{x}), \mathbf{x} - \mathbf{x}^* \rangle \geq f(\mathbf{x}) - f(\mathbf{x}^*) \geq f(\mathbf{x}) - f(\mathbf{y}) \geq \frac{1}{2L} \|\nabla f(\mathbf{x})\|_2^2 \geq \frac{\eta}{2} \|\nabla f(\mathbf{x})\|_2^2$



# Smooth Optimization: Smooth Gradient Descent

**Theorem (Smooth GD).** Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be  $L$ -smooth and convex, and suppose for  $\mathbf{x}_0 \in \mathbb{R}^n$  we have  $\|\mathbf{x}_0 - \mathbf{x}^*\|_2 \leq R$  for  $\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$ . Then iterating  $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t)$  for  $0 \leq t < T$ ,

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{2LR^2}{T}$$

*Proof.*

- Let  $\Phi_t := f(\mathbf{x}_t) - f(\mathbf{x}^*)$
- By **contractivity**, we know that  $\|\mathbf{x}_t - \mathbf{x}^*\|_2 \leq R \forall 0 \leq t \leq T$
- By convexity,  $\Phi_t \leq \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}^* \rangle \leq \|\nabla f(\mathbf{x}_t)\|_2 \cdot \|\mathbf{x}_t - \mathbf{x}^*\|_2 \leq \|\nabla f(\mathbf{x}_t)\|_2 \cdot R$
- On the other hand, we have

$$\Phi_{t+1} - \Phi_t = f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) \leq -\frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|_2^2 \leq -\frac{1}{2LR^2} \Phi_t^2$$

# Smooth Optimization: Smooth Gradient Descent

**Theorem (Smooth GD).** Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be  $L$ -smooth and convex, and suppose for  $\mathbf{x}_0 \in \mathbb{R}^n$  we have  $\|\mathbf{x}_0 - \mathbf{x}^*\|_2 \leq R$  for  $\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$ . Then iterating  $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t)$  for  $0 \leq t < T$ ,

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{2LR^2}{T}$$

*Proof.*

- We have the recursion:

$$\Phi_{t+1} - \Phi_t \leq -\frac{1}{2LR^2} \Phi_t^2$$

- Note that  $\frac{1}{\Phi_{t+1}} - \frac{1}{\Phi_t} = \frac{\Phi_t - \Phi_{t+1}}{\Phi_t \Phi_{t+1}} \geq \frac{\Phi_t - \Phi_{t+1}}{\Phi_t^2} \geq \frac{1}{2LR^2}$

- Telescoping for  $T$  iterations:

$$\frac{1}{\Phi_T} \geq \frac{1}{\Phi_0} + \frac{T}{2LR^2} \geq \frac{T}{2LR^2}$$



# Smooth Optimization: Smooth Gradient Descent

**Theorem (Smooth GD).** Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be  $L$ -smooth and convex, and suppose for  $\mathbf{x}_0 \in \mathbb{R}^n$  we have  $\|\mathbf{x}_0 - \mathbf{x}^*\|_2 \leq R$  for  $\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$ . Then iterating  $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t)$  for  $0 \leq t < T$ ,

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{2LR^2}{T}$$

- The error bound can be tightened:

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{2LR^2}{T + 4}$$

- The gradient descent is provably non-optimal in this setting

# Today's Lecture

- Lipschitz Optimization
- Smooth Optimization
- **Well-Conditioned Optimization**

# Well-Conditioned Optimization

- For smooth GD, the convergence rate is  $\mathcal{O}(LR^2/\epsilon)$
- Can gradient descent achieve **linear convergence rate**, i.e.,  $\sim \log(1/\epsilon)$ , or equivalently,  
$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \exp(-\Omega(t))$$

- To get some intuition, consider the gradient flow:

$$\frac{d}{dt}(f(\mathbf{x}_t) - f(\mathbf{x}^*)) = \frac{d}{dt}f(\mathbf{x}_t) = -\|\nabla f(\mathbf{x}_t)\|_2^2$$

- Suppose the following condition holds:  $\|\nabla f(\mathbf{x}_t)\|_2^2 \geq C(f(\mathbf{x}_t) - f(\mathbf{x}^*))$
- Then, by solving the ODE (or by **Grönwall's inequality**), we get the linear convergence rate:

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq e^{-Ct}(f(\mathbf{x}_0) - f(\mathbf{x}^*))$$

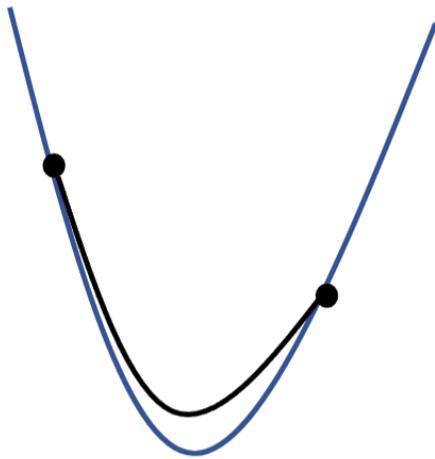
- Another way to see this is to define a **Lyapunov function**

$$V(t) := e^{Ct}(f(\mathbf{x}_t) - f(\mathbf{x}^*)), \quad dV(t) = e^{Ct}(C(f(\mathbf{x}_t) - f(\mathbf{x}^*)) - \|\nabla f(\mathbf{x}_t)\|_2^2) \leq 0$$

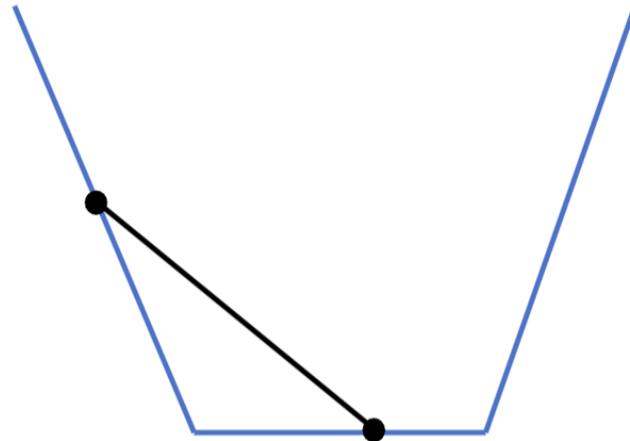
# Well-Conditioned Optimization

$f: \mathbb{R}^n \rightarrow \mathbb{R}$  is  $\mu$ -strongly convex if

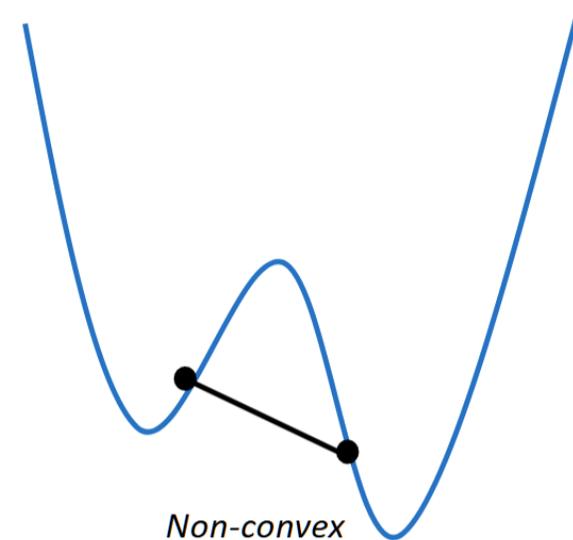
$$f(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y}) - \frac{\mu\lambda(1 - \lambda)}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 \quad \forall \mathbf{x} \neq \mathbf{y} \in \mathbb{R}^n$$



*Strongly convex*



*Convex*



*Non-convex*

# Well-Conditioned Optimization

$f: \mathbb{R}^n \rightarrow \mathbb{R}$  is  $\mu$ -strongly convex if

$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}) - \frac{\mu \lambda (1 - \lambda)}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 \quad \forall \mathbf{x} \neq \mathbf{y} \in \mathbb{R}^n$$

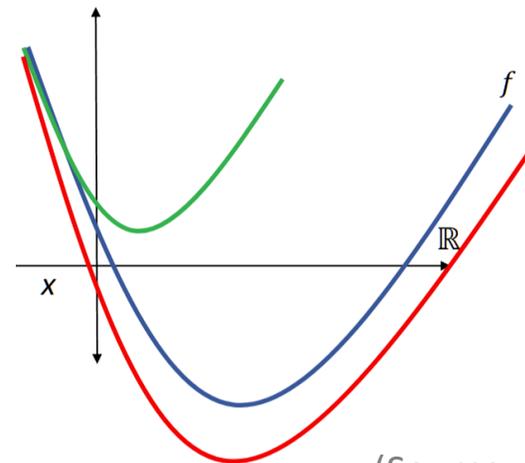
**Lemma.** If  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is differentiable, then  $f$  is  $\mu$ -strongly convex if and only if

$$f(\mathbf{x}') \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle + \frac{\mu}{2} \|\mathbf{x} - \mathbf{x}'\|_2^2 \quad \forall \mathbf{x}, \mathbf{x}' \in \mathbb{R}^n$$

If  $f$  is twice-differentiable,  $f$  is  $\mu$ -strongly convex iff  $\nabla^2 f(\mathbf{x})[\mathbf{v}, \mathbf{v}] \geq \mu \|\mathbf{v}\|_2^2 \quad \forall \mathbf{v} \in \mathbb{R}^n$

Smoothness  $\Rightarrow$  quadratic upper bound

Strong-convexity  $\Rightarrow$  quadratic lower bound



# Well-Conditioned Optimization

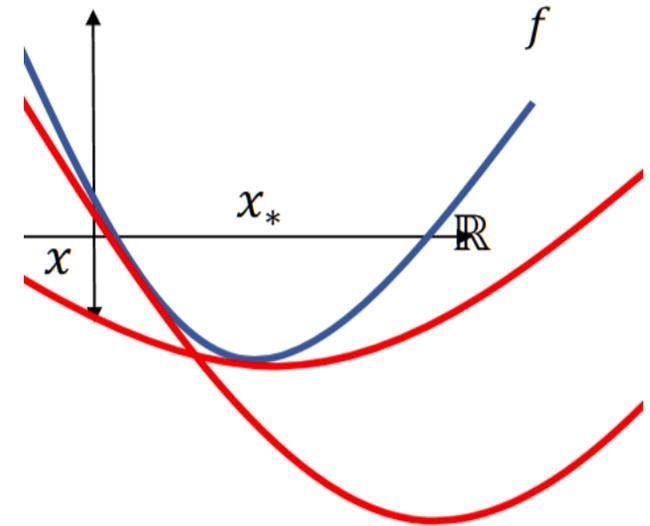
**Lemma.** Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be  $\mu$ -strongly convex, and let  $\mathbf{x}^* := \arg \min f(\mathbf{x})$ . Then

$$f(\mathbf{x}) - f(\mathbf{x}^*) \leq \frac{1}{2\mu} \|\nabla f(\mathbf{x})\|_2^2$$

*Proof.*

- Strong-convexity implies that for any  $\mathbf{y} \in \mathbb{R}^n$ ,  
$$f(\mathbf{y}) \geq L_{\mathbf{x}}(\mathbf{y}) := f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$$
- Thus, we have

$$\begin{aligned} f(\mathbf{x}^*) &= \min_{\mathbf{y} \in \mathbb{R}^n} f(\mathbf{y}) \geq \min_{\mathbf{y} \in \mathbb{R}^n} L_{\mathbf{x}}(\mathbf{y}) \\ &= f(\mathbf{x}) + \left\langle \nabla f(\mathbf{x}), -\frac{1}{\mu} \nabla f(\mathbf{x}) \right\rangle + \frac{\mu}{2} \left\| \frac{1}{\mu} \nabla f(\mathbf{x}) \right\|_2^2 \\ &= f(\mathbf{x}) - \frac{1}{2\mu} \|\nabla f(\mathbf{x})\|_2^2 \end{aligned}$$



# Well-Conditioned Optimization: Well-conditioned Gradient Descent

**Theorem (Well-conditioned GD).** Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be  $L$ -smooth and  $\mu$ -strongly convex, and let  $x^* := \arg \min f(x)$  and  $\kappa := \frac{L}{\mu} \geq 1$ . Then iterating  $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t)$  for  $0 \leq t < T$ ,

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \left(1 - \frac{1}{\kappa}\right)^T (f(\mathbf{x}_0) - f(\mathbf{x}^*))$$

*Proof.*

- For  $0 \leq t < T$ , we have

$$f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) \leq \underbrace{-\frac{1}{2L}}_{L\text{-smooth}} \|\nabla f(\mathbf{x}_t)\|_2^2 \leq \underbrace{-\frac{\mu}{L}}_{\mu\text{-strongly convex}} (f(\mathbf{x}_t) - f(\mathbf{x}^*))$$

- Thus,

$$f(\mathbf{x}_{t+1}) - f(\mathbf{x}^*) \leq \left(1 - \frac{1}{\kappa}\right) (f(\mathbf{x}_t) - f(\mathbf{x}^*))$$



# Well-Conditioned Optimization: Well-conditioned Gradient Descent

**Theorem (Well-conditioned GD).** Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be  $L$ -smooth and  $\mu$ -strongly convex, and let  $\mathbf{x}^* := \arg \min f(x)$  and  $\kappa := \frac{L}{\mu} \geq 1$ . Then iterating  $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t)$  for  $0 \leq t < T$ ,

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \left(1 - \frac{1}{\kappa}\right)^T (f(\mathbf{x}_0) - f(\mathbf{x}^*))$$

- GD can achieve linear convergence rate  $T = \mathcal{O}(\kappa \log(1/\epsilon))$
- This result only uses the smoothness and the bound  $C(f(\mathbf{x}) - f(\mathbf{x}^*)) \leq \|\nabla f(\mathbf{x})\|_2^2$ , which can still hold even if  $f$  is non-convex
- Bounds of this form is called **Polyak-Łojasiewicz conditions**

# Well-Conditioned Optimization: Hardness

**Theorem.** Let  $\kappa \geq 1$  and  $\epsilon \in (0,1)$ . No algorithm  $\mathcal{A}$  which accesses a target  $L$ -smooth,  $\mu$ -strongly convex function  $f$  with  $\kappa = \frac{L}{\mu}$  using a gradient oracle  $O$  can optimize  $f$  to additive error  $\epsilon(f(\mathbf{x}_0) - \min f(\mathbf{x}))$  using  $T < \frac{\sqrt{\kappa}-1}{2} \log\left(\frac{1}{\kappa\epsilon}\right)$  queries, subject to the restriction that

$$\mathbf{x}_0 = \mathbf{0}, \quad \mathbf{x}_t \in \text{span}(\{O(\mathbf{x}_s)\}_{0 \leq s < t}) \quad \forall t \in [T]$$

Proof.

- Define the hard function

$$f(x) := -\frac{L-\mu}{4}x_1 + \frac{L-\mu}{8} \sum_{i=1}^{n-1} (x_i - x_{i+1})^2 + \frac{L-\mu}{8}x_1^2 + \frac{\mu}{2} \sum_{i=1}^n x_i^2 + \frac{\sqrt{L\mu}-\mu}{4}x_n^2$$

- Note that  $f = (\text{linear}) + \frac{1}{2}\mathbf{x}^\top A\mathbf{x}$ , and thus,  $\nabla^2 f(\mathbf{x}) = A$

# Well-Conditioned Optimization: Hardness

$$\mathbf{v}^2 \nabla^2 f(\mathbf{x}) \mathbf{v} = \frac{L - \mu}{4} \sum_{i=1}^{n-1} (v_i - v_{i+1})^2 + \frac{L - \mu}{4} v_1^2 + \mu \sum_{i=1}^n v_i^2 + \frac{\sqrt{L\mu} - \mu}{2} v_n^2$$

- $L$ -smooth:

$$\begin{aligned} \mathbf{v}^2 \nabla^2 f(\mathbf{x}) \mathbf{v} &\leq \frac{L - \mu}{4} \sum_{i=1}^{n-1} (2v_i^2 + 2v_{i+1}^2) + \frac{L - \mu}{4} v_1^2 + \mu \sum_{i=1}^n v_i^2 + \frac{\sqrt{L\mu} - \mu}{2} v_n^2 \\ &\leq (L - \mu) \sum_{i=1}^n v_i^2 + \mu \sum_{i=1}^n v_i^2 = L \sum_{i=1}^n v_i^2 \end{aligned}$$

- $\mu$ -strongly convex:

$$\mathbf{v}^2 \nabla^2 f(\mathbf{x}) \mathbf{v} \geq \mu \sum_{i=1}^n v_i^2$$

# Well-Conditioned Optimization: Hardness

$$f(x) := -\frac{L-\mu}{4}x_1 + \frac{L-\mu}{8}\sum_{i=1}^{n-1}(x_i - x_{i+1})^2 + \frac{L-\mu}{8}x_1^2 + \frac{\mu}{2}\sum_{i=1}^n x_i^2 + \frac{\sqrt{L\mu}-\mu}{4}x_n^2$$

- Note that  $\nabla f(\mathbf{x}_0) = \nabla f(\mathbf{0}) = (?, 0, \dots, 0)$  and  $\mathbf{x}_1 = (?, 0, \dots, 0)$  by our assumption
- Then,  $\nabla f(\mathbf{x}_1) = (?, ?, 0, \dots, 0)$
- By induction, you can show that only the first  $k$  coordinates of  $x_k$  are non-zero
- By strong convexity, we have

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \geq \frac{\mu}{2}\|\mathbf{x}_T - \mathbf{x}^*\|_2^2 \geq \frac{\mu}{2}\sum_{i=T+1}^n (x_i^*)^2$$

# Well-Conditioned Optimization: Hardness

$$f(x) := -\frac{L-\mu}{4}x_1 + \frac{L-\mu}{8}\sum_{i=1}^{n-1}(x_i - x_{i+1})^2 + \frac{L-\mu}{8}x_1^2 + \frac{\mu}{2}\sum_{i=1}^n x_i^2 + \frac{\sqrt{L\mu}-\mu}{4}x_n^2$$

- Let  $\mathbf{x}^* := \arg \min f(\mathbf{x})$ . By the first-order optimality condition,  $\nabla f(\mathbf{x}^*) = \mathbf{0}$ , i.e.,

$$-\frac{L-\mu}{4} + \frac{L-\mu}{4}(x_1^* - x_2^*) + \frac{L-\mu}{4}x_1^* + \mu x_1^* = 0$$

$$\frac{L-\mu}{4}(x_i^* - x_{i+1}^*) + \frac{L-\mu}{4}(x_i^* - x_i^*) + \mu x_i^* = 0 \quad \forall i \in \{2, \dots, n-1\}$$

$$\frac{L-\mu}{4}(x_n^* - x_{n-1}^*) + \mu x_n^* + \frac{\sqrt{L\mu}-\mu}{2}x_n^* = 0$$

- Solving the linear system, we get that

$$x_i^* = \left(1 - \frac{2}{\sqrt{\kappa} + 1}\right)^i \quad \forall i \in [n]$$

# Well-Conditioned Optimization: Hardness

$$f(x) := -\frac{L-\mu}{4}x_1 + \frac{L-\mu}{8}\sum_{i=1}^{n-1}(x_i - x_{i+1})^2 + \frac{L-\mu}{8}x_1^2 + \frac{\mu}{2}\sum_{i=1}^n x_i^2 + \frac{\sqrt{L\mu}-\mu}{4}x_n^2$$

- By induction, you can show that only the first  $k$  coordinates of  $x_k$  are non-zero
- By strong convexity, we have

$$\begin{aligned} f(\mathbf{x}_T) - f(\mathbf{x}^*) &\geq \frac{\mu}{2}\|\mathbf{x}_T - \mathbf{x}^*\|_2^2 \geq \frac{\mu}{2}\sum_{i=T+1}^n (x_i^*)^2 = \frac{\mu}{2}\sum_{i=T+1}^n \left(1 - \frac{2}{\sqrt{\kappa} + 1}\right)^{2i} \\ &\underset{n \rightarrow \infty}{\approx} \frac{\mu}{2}\left(1 - \frac{2}{\sqrt{\kappa} + 1}\right)^{2T} \sum_{i=1}^n \left(1 - \frac{2}{\sqrt{\kappa} + 1}\right)^{2i} = \frac{\mu}{2}\left(1 - \frac{2}{\sqrt{\kappa} + 1}\right)^{2T} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 \\ &\geq \frac{1}{\kappa}\left(1 - \frac{2}{\sqrt{\kappa} + 1}\right)^{2T} \underbrace{(f(\mathbf{x}_0) - f(\mathbf{x}^*))}_{L\text{-smooth}} \end{aligned}$$



# Well-Conditioned Optimization: Hardness

- The hardness result indicates that the optimal query complexity in the well-conditioned setting should be  $T = \sqrt{\kappa} \log(1/\epsilon)$ , while GD can only achieve  $\mathcal{O}(\kappa \log(1/\epsilon))$ , which is sub-optimal
- To achieve the optimal rate, we need to use some “**history information**” (i.e., momentum) in the update rule, and the algorithms are called **accelerated GD** (first discovered by Nesterov)
- The hard function in the proof is

$$f(x) := -\frac{L-\mu}{4}x_1 + \frac{L-\mu}{8}\sum_{i=1}^{n-1}(x_i - x_{i+1})^2 + \frac{L-\mu}{8}x_1^2 + \frac{\mu}{2}\sum_{i=1}^n x_i^2 + \frac{\sqrt{L\mu}-\mu}{4}x_n^2$$

It is related to the **Laplacian** of a path graph



- Given any convex function, we construct the **dependence graph**  $G$  on the set of variables  $x_i$  by connecting  $x_i \sim x_j$  if  $\nabla f(x)_i$  depends on  $x_j$  or  $\nabla f(x)_j$  depends on  $x_i$
- GD can only transmit information from one vertex to another in each iteration

# Reduction between Smooth and Well-Conditioned Optimizations

**Lemma.** Suppose that algorithm  $\mathcal{A}$  can optimize an  $L$ -smooth, convex  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  with  $\mathbf{x}^* \in \arg \min f(\mathbf{x})$ , and outputs a point  $\mathbf{x}_T \in \mathbb{R}^n$  using  $T$  queries such that,

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \lesssim \frac{L \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{T^c}$$

Then there is an algorithm  $\mathcal{A}'$  that can optimize an  $L$ -smooth,  $\mu$ -strongly convex  $g: \mathbb{R}^n \rightarrow \mathbb{R}$  with  $\mathbf{y}^* = \arg \min g(\mathbf{y})$  and  $\kappa := \frac{L}{\mu}$ , and outputs a point  $\mathbf{y}_T \in \mathbb{R}^n$  such that

$$g(\mathbf{y}_T) - g(\mathbf{y}^*) \leq \epsilon (g(\mathbf{y}_0) - g(\mathbf{y}^*))$$

in  $\mathcal{O}(\kappa^{1/c} \log(1/\epsilon))$  queries

- The converse direction also holds ([Allen-Zhu-Hazan '16](#))
- GD for smooth function can achieve  $c = 1$ , and our lower bound for well-conditioned case shows that  $c \leq 2$  (and  $c = 2$  can be achieved by AGD). Thus, GD is also sub-optimal in smooth optimization

# Reduction between Smooth and Well-Conditioned Optimizations

*Proof.*

- Let  $T' := \mathcal{O}\left((4\kappa)^{1/c}\right)$ .
- We first show that by running  $\mathcal{A}$  with  $T'$  queries, the error can be decreased by a half
- By the convergence rate of  $\mathcal{A}$ , we have

$$g(\mathbf{y}_T) - g(\mathbf{y}^*) \lesssim \frac{L\|\mathbf{y}_0 - \mathbf{y}^*\|_2^2}{\left((4\kappa)^{1/c}\right)^c} = \frac{\mu}{4}\|\mathbf{y}_0 - \mathbf{y}^*\|_2^2 \leq \frac{1}{2}(g(\mathbf{y}_0) - g(\mathbf{y}^*))$$

- Then, by iterating  $\log(1/\epsilon)$  times, the error can be decreased to  $\epsilon$



# Summary

Regularity	Goal	Algorithm	Iterations
$f: B_2^n \rightarrow \mathbb{R}$ , $L$ -Lipschitz	$\epsilon$ -optimal	Any algorithm	$\Omega((L/\epsilon)^n)$
$f: RB_2^n \rightarrow \mathbb{R}$ , convex, $L$ -Lipschitz	$\epsilon$ -optimal	Any first-order method	$\Omega(\min\{n, (LR/\epsilon)^2\})$
		Cutting-plane method	$\tilde{O}(n)$
		Projected GD	$\mathcal{O}((LR/\epsilon)^2)$
$f: \mathbb{R}^n \rightarrow \mathbb{R}$ , $L$ -smooth	$\epsilon$ -critical	GD	$\mathcal{O}(L(f(\mathbf{x}_0) - f^*)/\epsilon)$
$f: \mathbb{R}^n \rightarrow \mathbb{R}$ , convex, $L$ -smooth	$\epsilon$ -optimal	GD	$\mathcal{O}(L\ \mathbf{x}_0 - \mathbf{x}^*\ _2^2/\epsilon)$
		Any first-order method	$\Omega\left(\sqrt{L\ \mathbf{x}_0 - \mathbf{x}^*\ _2^2/\epsilon}\right)$
		GD	$\tilde{O}(L/\mu)$
$f: \mathbb{R}^n \rightarrow \mathbb{R}$ , $L$ -smooth, $\mu$ -strongly convex	$\epsilon$ -optimal	GD	$\tilde{O}(L/\mu)$
		Any first-order method	$\Omega\left(\sqrt{L/\mu}\right)$