

# CS257 Linear and Convex Optimization

## Lecture 8

Bo Jiang

John Hopcroft Center for Computer Science  
Shanghai Jiao Tong University

October 26, 2020

## Recap: QP, QCQP

### Quadratic program (QP)

$$\begin{array}{ll}\min_{\mathbf{x}} & \frac{1}{2}\mathbf{x}^T \mathbf{Q}\mathbf{x} + \mathbf{c}^T \mathbf{x} \\ \text{s.t.} & \mathbf{B}\mathbf{x} \leq \mathbf{d} \\ & \mathbf{A}\mathbf{x} = \mathbf{b}\end{array}$$

### Quadratically constrained quadratic program (QCQP)

$$\begin{array}{ll}\min_{\mathbf{x}} & \frac{1}{2}\mathbf{x}^T \mathbf{Q}\mathbf{x} + \mathbf{c}^T \mathbf{x} \\ \text{s.t.} & \frac{1}{2}\mathbf{x}^T \mathbf{Q}_i \mathbf{x} + \mathbf{c}_i^T \mathbf{x} + \mathbf{d}_i \leq 0, \quad i = 1, 2, \dots, m \\ & \mathbf{A}\mathbf{x} = \mathbf{b}\end{array}$$

# Recap: LS, Lasso, and Ridge Regressions

Given  $\mathbf{y} \in \mathbb{R}^n$ ,  $\mathbf{X} \in \mathbb{R}^{n \times p}$ .

## Linear least squares regression

$$\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$$

## Lasso

$$\begin{aligned} \min_{\mathbf{w}} \quad & \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 \\ \text{s. t.} \quad & \|\mathbf{w}\|_1 \leq t \end{aligned}$$

## Ridge regression

$$\begin{aligned} \min_{\mathbf{w}} \quad & \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 \\ \text{s. t.} \quad & \|\mathbf{w}\|_2^2 \leq t \end{aligned}$$

# Recap: GP

## Posynomial form

$$\begin{aligned} \min_{\mathbf{x}} \quad & \sum_{k=1}^{p_0} \gamma_{0k} x_1^{a_{0k1}} x_2^{a_{0k2}} \cdots x_n^{a_{0kn}} \\ \text{s. t.} \quad & \sum_{k=1}^{p_i} \gamma_{ik} x_1^{a_{ik1}} x_2^{a_{ik2}} \cdots x_n^{a_{ikn}} \leq 1, \quad i = 1, \dots, m \\ & \eta_j x_1^{c_{j1}} x_2^{c_{j2}} \cdots x_n^{c_{jn}} = 1, \quad j = 1, \dots, r \end{aligned}$$

## Convex form

$$\begin{aligned} \min_{\mathbf{y}} \quad & \log \left( \sum_{k=1}^{p_0} e^{\mathbf{a}_{0k}^T \mathbf{y} + b_{0k}} \right) \\ \text{s. t.} \quad & \log \left( \sum_{k=1}^{p_i} e^{\mathbf{a}_{ik}^T \mathbf{y} + b_{ik}} \right) \leq 0, \quad i = 1, \dots, m \\ & \mathbf{c}_j^T \mathbf{y} + d_j = 0, \quad j = 1, \dots, r \end{aligned}$$

# Contents

## 1. Gradient Descent

# Unconstrained Optimization Problems

Consider an unconstrained, smooth convex optimization problem

$$\min_{\mathbf{x}} f(\mathbf{x})$$

where  $f$  is convex and differentiable on  $\mathbb{R}^n$ .

The optimal solution satisfies the first-order optimality condition

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

In some rare cases, this yields closed-form solutions, e.g.

$$\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$$

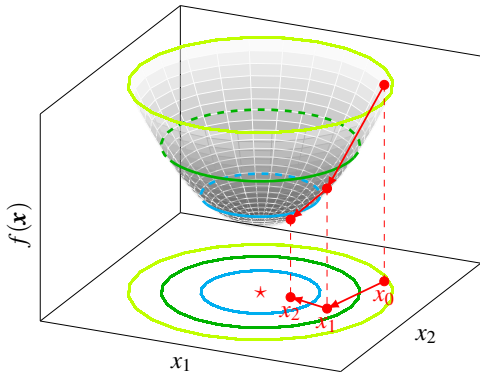
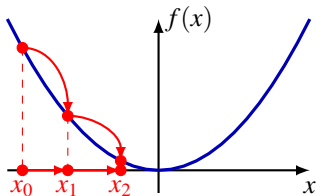
has closed-form solution

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

But in most cases we need numerical algorithms.

# Descent Method

- 1: choose initial point  $\mathbf{x}_0 \in \mathbb{R}^n$
- 2: **repeat**
- 3:     choose **descent direction**  $\mathbf{d}_k \in \mathbb{R}^n$  and **step size**  $t_k > 0$
- 4:      $\mathbf{x}_{k+1} = \mathbf{x}_k + t_k \mathbf{d}_k$    s.t.    $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$
- 5: **until** stopping criterion is satisfied



## Questions

- How to choose  $\mathbf{d}_k$  and  $t_k$ ?
- Does  $\mathbf{x}_k$  converge to  $\mathbf{x}^*$ ?

## Descent Direction

$\mathbf{d}_k$  is a **descent direction** at  $\mathbf{x}_k$  if for all small enough  $t > 0$

$$g(t) \triangleq f(\mathbf{x}_k + t\mathbf{d}_k) < f(\mathbf{x}_k) = g(0)$$

For differentiable  $f$  (not necessarily convex),

- if  $\mathbf{d}_k$  is a descent direction, then  $g'(0) = \mathbf{d}_k^T \nabla f(\mathbf{x}_k) \leq 0$ ;
- if  $g'(0) = \mathbf{d}_k^T \nabla f(\mathbf{x}_k) < 0$ , then  $\mathbf{d}_k$  is a descent direction.

For convex  $f$ , by the first-order condition for convexity,

$$f(\mathbf{x}_k) > f(\mathbf{x}_k + t\mathbf{d}_k) \geq f(\mathbf{x}_k) + t\mathbf{d}_k^T \nabla f(\mathbf{x}_k).$$

$\mathbf{d}_k^T \nabla f(\mathbf{x}_k) < 0$  is also **necessary** for  $\mathbf{d}_k$  to be a descent direction.

For convex differentiable  $f$ ,

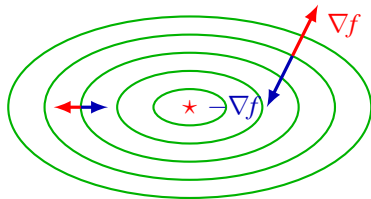
$$\mathbf{d}_k \text{ is a descent direction} \iff \mathbf{d}_k^T \nabla f(\mathbf{x}_k) < 0$$



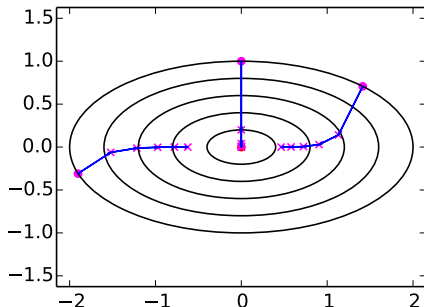
# Gradient Descent

Choose  $\mathbf{d}_k = -\nabla f(\mathbf{x}_k)$ ,  $\mathbf{d}_k^T \nabla f(\mathbf{x}_k) = -\|\nabla f(\mathbf{x}_k)\|_2^2 < 0$  unless  $\nabla f(\mathbf{x}_k) = \mathbf{0}$ .

$$\mathbf{x}_{k+1} = \mathbf{x}_k - t_k \nabla f(\mathbf{x}_k)$$



level curves of  $f(x_1, x_2) = \frac{x_1^2}{4} + x_2^2$



**Question.** What happens if  $\nabla f(\mathbf{x}_k) = \mathbf{0}$ ?

# Max-rate Descending Direction

$-\nabla f(\mathbf{x}_k)$  is the max-rate descending direction

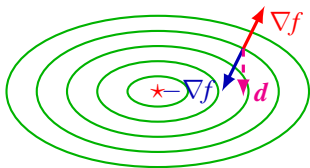
- If  $\|\mathbf{d}_k\|_2 = 1$ , the rate of change of  $f$  at  $\mathbf{x}_k$  along the direction  $\mathbf{d}_k$  is

$$\nabla_{\mathbf{d}_k} f(\mathbf{x}_k) = g'(0) = \lim_{t \downarrow 0} \frac{f(\mathbf{x}_k + t\mathbf{d}_k) - f(\mathbf{x}_k)}{t} = \mathbf{d}_k^T \nabla f(\mathbf{x}_k)$$

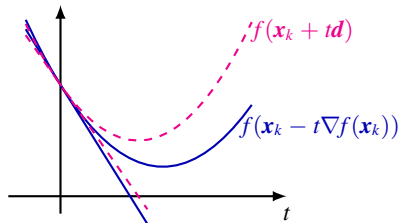
- by the Cauchy-Schwarz inequality,

$$\mathbf{d}_k^T \nabla f(\mathbf{x}_k) \geq -\|\mathbf{d}_k\|_2 \cdot \|\nabla f(\mathbf{x}_k)\|_2 = -\|\nabla f(\mathbf{x}_k)\|_2$$

with equality iff  $\mathbf{d}_k = -\nabla f(\mathbf{x}_k) / \|\nabla f(\mathbf{x}_k)\|$



level curves of  $f(x_1, x_2) = \frac{x_1^2}{4} + x_2^2$



# Gradient Descent Algorithm

- 1: initialization  $\mathbf{x} \leftarrow \mathbf{x}_0 \in \mathbb{R}^n$
- 2: **while**  $\|\nabla f(\mathbf{x})\| > \delta$  **do**
- 3:      $\mathbf{x} \leftarrow \mathbf{x} - t \nabla f(\mathbf{x})$
- 4: **end while**
- 5: **return**  $\mathbf{x}$

Step size (aka **learning rate** in machine learning)

- the above algorithm uses constant step size  $t$  for all iterations
- there are other methods for choosing  $t$  for each iteration, e.g. **exact line search**, **backtracking line search**

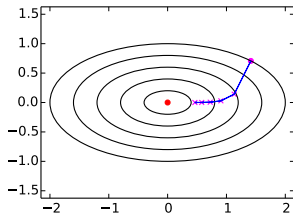
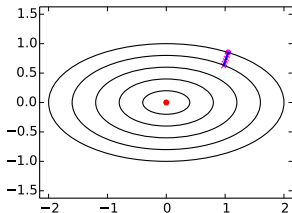
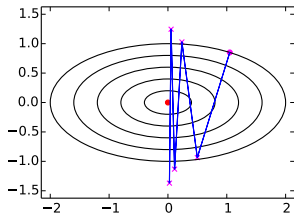
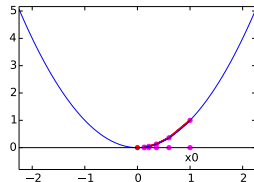
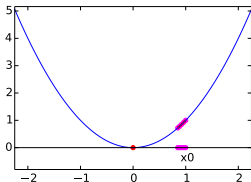
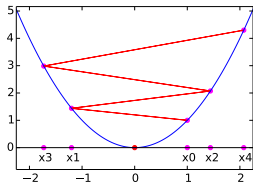
Stopping criterion

- ideally, stop if  $\nabla f(\mathbf{x}) = 0$  (optimality condition), but impractical
- more practical: stop when  $\|\nabla f(\mathbf{x})\| \leq \delta$  for some small  $\delta$
- other criteria:  $|f(\mathbf{x}_{\text{new}}) - f(\mathbf{x}_{\text{old}})| \leq \delta$ ,  $\frac{|f(\mathbf{x}_{\text{new}}) - f(\mathbf{x}_{\text{old}})|}{|f(\mathbf{x}_{\text{old}})|} \leq \delta$ , ...
- in practice, also stop if maximum # of iterations is reached

# Large vs. Small Step Size

Consider constant step size. How large should the step size be?

- Too large: may oscillate and diverge
- Too small: may be too slow
- “Just right”: fast convergence



# 1D Example

Consider  $f(x) = \frac{1}{2}ax^2$ , where  $a > 0$ .

- gradient at  $x_k$  is  $f'(x_k) = ax_k$ , so  $d_k = -ax_k$  in gradient descent
- update

$$x_{k+1} = x_k - tf'(x_k) = (1 - at)x_k$$

- in order for  $f(x_{k+1}) < f(x_k)$ , the step size should satisfy

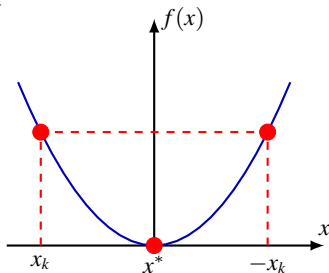
$$|x_{k+1}| < |x_k| \implies |1 - at| < 1 \implies 0 < t < \frac{2}{a}$$

- $x_k \rightarrow x^* = 0$  geometrically for such  $t$

Note  $f$  satisfies

- $|f'(x) - f'(y)| = a|x - y|$
- $f''(x) = a$

$f'$  is Lipschitz continuous



# Lipschitz Continuity

A function  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is **Lipschitz continuous** with **Lipschitz constant**  $L > 0$ , or simply  **$L$ -Lipschitz**, if

$$\|f(\mathbf{x}) - f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|, \quad \forall \mathbf{x}, \mathbf{y}$$

**Note.** Lipschitz continuity can be defined with respect to any norms. But we will assume the norms in the above definition are the 2-norms in  $\mathbb{R}^n$  and  $\mathbb{R}^m$ , respectively, unless stated otherwise.

**Note.** Lipschitz continuity implies uniform continuity.

**Example.**  $f(x) = ax$  is  $|a|$ -Lipschitz,  $|f(x) - f(y)| = |a| \cdot |x - y|$

**Example.**  $f(x) = |x|$  is 1-Lipschitz,  $|f(x) - f(y)| = ||x| - |y|| \leq |x - y|$

**Example.**  $f(\mathbf{x}) = \mathbf{a}^T \mathbf{x}$  is  $\|\mathbf{a}\|$ -Lipschitz,  $|\mathbf{a}^T \mathbf{x} - \mathbf{a}^T \mathbf{y}| \leq \|\mathbf{a}\| \cdot \|\mathbf{x} - \mathbf{y}\|$  by the Cauchy-Schwarz inequality.

## Lipschitz Continuity (cont'd)

**Example.** Let  $\mathbf{Q} = \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}$ .  $f(\mathbf{x}) = \mathbf{Q}\mathbf{x} = (x_1, 2x_2)^T$  is 2-Lipschitz.

$$f(\mathbf{x}) - f(\mathbf{y}) = (x_1 - y_1, 2x_2 - 2y_2)^T = (d_1, 2d_2)^T$$

$$\|f(\mathbf{x}) - f(\mathbf{y})\| = \sqrt{d_1^2 + 4d_2^2} \leq 2\sqrt{d_1^2 + d_2^2} = 2\|\mathbf{x} - \mathbf{y}\|$$

More generally,  $f(\mathbf{x}) = \mathbf{Q}\mathbf{x}$  with  $\mathbf{Q} \succeq \mathbf{O}$  is  $\lambda_{\max}(\mathbf{Q})$ -Lipschitz, where  $\lambda_{\max}(\mathbf{Q})$  is the largest eigenvalue of  $\mathbf{Q}$ <sup>1</sup>.

**Proof.** Let  $\mathbf{d} = \mathbf{x} - \mathbf{y}$ .

$$\|f(\mathbf{x}) - f(\mathbf{y})\| = \|\mathbf{Q}\mathbf{d}\| = \sqrt{\mathbf{d}^T \mathbf{Q}^2 \mathbf{d}} \leq \sqrt{\lambda_{\max}(\mathbf{Q}^2) \|\mathbf{d}\|^2} = \lambda_{\max}(\mathbf{Q}) \|\mathbf{x} - \mathbf{y}\|$$

The last equality uses the fact  $\lambda_{\max}(\mathbf{Q}^2) = \lambda_{\max}^2(\mathbf{Q})$ .

---

<sup>1</sup>if we do not assume  $\mathbf{Q} \succeq \mathbf{O}$ , then we should replace  $\lambda_{\max}(\mathbf{Q})$  by  $\sigma_{\max}(\mathbf{Q})$ , the largest singular value of  $\mathbf{Q}$

## Appendix: Bounds on Quadratic Forms

**Proposition.** For a symmetric matrix  $\mathbf{Q} \in \mathbb{R}^{n \times n}$ ,

$$\lambda_{\min} \|\mathbf{x}\|_2^2 \leq \mathbf{x}^T \mathbf{Q} \mathbf{x} \leq \lambda_{\max} \|\mathbf{x}\|_2^2, \quad \forall \mathbf{x} \in \mathbb{R}^n$$

where  $\lambda_{\max}$  and  $\lambda_{\min}$  are the largest and the smallest eigenvalues of  $\mathbf{Q}$ , respectively.

**Proof.** Recall that  $\mathbf{Q}$  can be orthogonally diagonalized, i.e.  $\mathbf{Q} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T$ , where  $\mathbf{\Lambda} = \text{diag}\{\lambda_1, \dots, \lambda_n\}$  and  $\mathbf{U}^T \mathbf{U} = \mathbf{I}$  (Lecture 2, slide 24). Let  $\mathbf{x} = \mathbf{U} \mathbf{y}$ .

$$\mathbf{x}^T \mathbf{Q} \mathbf{x} = \mathbf{y}^T (\mathbf{U}^T \mathbf{Q} \mathbf{U}) \mathbf{y} = \mathbf{y}^T \mathbf{\Lambda} \mathbf{y} = \sum_{i=1}^n \lambda_i y_i^2 \leq \sum_{i=1}^n \lambda_{\max} y_i^2 = \lambda_{\max} \|\mathbf{y}\|_2^2$$

Then use the fact that orthogonal transformations preserves 2-norms, i.e.

$$\|\mathbf{x}\|_2^2 = \mathbf{x}^T \mathbf{x} = (\mathbf{U} \mathbf{y})^T (\mathbf{U} \mathbf{y}) = \mathbf{y}^T (\mathbf{U}^T \mathbf{U}) \mathbf{y} = \mathbf{y}^T \mathbf{y} = \|\mathbf{y}\|_2^2.$$

Similarly for  $\mathbf{x}^T \mathbf{Q} \mathbf{x} \geq \lambda_{\min} \|\mathbf{x}\|_2^2$ .



# Lipschitz Continuity of Gradient

A function is  **$L$ -smooth** if it is differentiable and its gradient is  $L$ -Lipschitz, i.e.

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|, \quad \forall \mathbf{x}, \mathbf{y}$$

**Note.**  $L$  upper bounds the rate of change of  $\nabla f$

**Example.**  $f(x) = \frac{1}{2}ax^2$  is  $|a|$ -smooth, since  $f'(x) = ax$  is  $|a|$ -Lipschitz

**Example.**  $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{Q} \mathbf{x}$  with  $\mathbf{Q} \succeq \mathbf{O}$  is  $\lambda_{\max}(\mathbf{Q})$ -smooth, since  $\nabla f(\mathbf{x}) = \mathbf{Q}\mathbf{x}$  is  $\lambda_{\max}(\mathbf{Q})$ -Lipschitz.

With  $\mathbf{Q} = \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}$ , we obtain  $f(\mathbf{x}) = \frac{1}{2}x_1^2 + x_2^2$  is 2-smooth.

**Lemma.** A twice continuously differentiable **convex**  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is  $L$ -smooth iff  $\nabla^2 f(\mathbf{x}) \preceq L\mathbf{I}$ , meaning  $L\mathbf{I} - \nabla^2 f(\mathbf{x}) \succeq \mathbf{O}$ , or equivalently  $\lambda_{\max}(\nabla^2 f(\mathbf{x})) \leq L$ .

## Lipschitz Continuity of Gradient (cont'd)

**Lemma.** A twice continuously differentiable  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is  $L$ -smooth iff for any  $\mathbf{x}$ ,  $-L\mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq L\mathbf{I}$ , or equivalently  $|\lambda| \leq L$  for all eigenvalues  $\lambda$  of  $\nabla^2 f(\mathbf{x})$ .

**Proof.** “ $\Leftarrow$ ”. Assume  $-L\mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq L\mathbf{I}$  for all  $\mathbf{x}$ .

By the Mean Value Theorem, there exists  $\mathbf{z} = \mathbf{y} + t(\mathbf{x} - \mathbf{y})$  for some  $t \in [0, 1]$  s.t.

$$\nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) = \nabla^2 f(\mathbf{z})(\mathbf{x} - \mathbf{y})$$

Since  $-L\mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq L\mathbf{I}$ ,

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| = \|\nabla^2 f(\mathbf{z})(\mathbf{x} - \mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|$$

by the proof on slide 15.

## Lipschitz Continuity of Gradient (cont'd)

**Proof (cont'd).** “ $\Rightarrow$ ”. Assume  $f$  is  $L$ -smooth. Fix a direction  $\mathbf{d}$ . Let  $h(t) = \nabla f(\mathbf{x} + t\mathbf{d})$ . Since  $f$  is  $L$ -smooth,

$$\|g(t) - g(0)\| = \|\nabla f(\mathbf{x} + t\mathbf{d}) - \nabla f(\mathbf{x})\| \leq L\|t\mathbf{d}\|$$

so

$$\left\| \frac{g(t) - g(0)}{t} \right\| \leq L\|\mathbf{d}\|$$

Letting  $t \rightarrow 0$  and using the chain rule

$$\|\nabla^2 f(\mathbf{x})\mathbf{d}\| = \|g'(0)\| \leq L\|\mathbf{d}\|$$

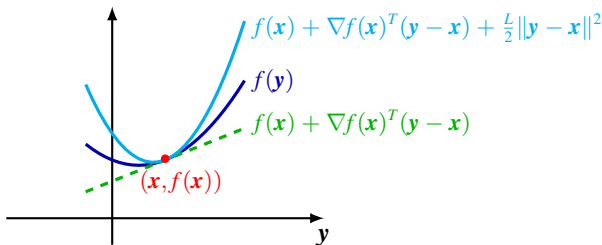
Let  $\mathbf{d}$  be an eigenvector of  $\nabla^2 f(\mathbf{x})$  with associated eigenvalue  $\lambda$ ,

$$|\lambda| \cdot \|\mathbf{d}\| = \|\lambda\mathbf{d}\| \leq L\|\mathbf{d}\| \implies |\lambda| \leq L$$

# Quadratic Upper Bound

**Lemma.** If  $f$  is  $L$ -smooth, then

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \frac{L}{2}\|\mathbf{y} - \mathbf{x}\|^2$$



**Note.** The upper bound does **not** assume the convexity of  $f$ .

If  $\nabla^2 f(\mathbf{x}) \preceq L\mathbf{I}$ , this is intuitive from the second-order Taylor expansion

$$f(\mathbf{y}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \frac{1}{2}(\mathbf{y} - \mathbf{x})^T \nabla^2 f(\mathbf{z})(\mathbf{y} - \mathbf{x})$$

for some  $\mathbf{z}$  on the line segment between  $\mathbf{x}$  and  $\mathbf{y}$ . (Check  $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{Q}\mathbf{x}$ )

## Proof

Let  $\mathbf{z}(t) = \mathbf{x} + t(\mathbf{y} - \mathbf{x})$  and  $g(t) = f(\mathbf{z}(t))$ . Then  $g(0) = f(\mathbf{x})$ ,  $g(1) = f(\mathbf{y})$ ,  $g'(t) = \nabla f(\mathbf{z}(t))^T(\mathbf{y} - \mathbf{x})$ ,  $g'(0) = \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x})$ .

1. By the Newton-Leibniz formula,

$$g(1) - g(0) - g'(0) = \int_0^1 [g'(t) - g'(0)] dt \leq \int_0^1 |g'(t) - g'(0)| dt$$

2. By the Cauchy-Schwarz inequality and  $L$ -smoothness

$$\begin{aligned} |g'(t) - g'(0)| &= |[\nabla f(\mathbf{z}(t)) - \nabla f(\mathbf{x})]^T(\mathbf{y} - \mathbf{x})| \\ &\leq \|\nabla f(\mathbf{z}(t)) - \nabla f(\mathbf{x})\| \cdot \|\mathbf{y} - \mathbf{x}\| \\ &\leq L\|\mathbf{z}(t) - \mathbf{x}\| \cdot \|\mathbf{y} - \mathbf{x}\| \\ &= tL\|\mathbf{x} - \mathbf{y}\|^2 \end{aligned}$$

3. Plugging  $|g'(t) - g'(0)| \leq tL\|\mathbf{x} - \mathbf{y}\|^2$  into step 1,

$$f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) = g(1) - g(0) - g'(0) \leq \frac{L}{2}\|\mathbf{x} - \mathbf{y}\|^2$$

## Consequence of Quadratic Upper Bound

For  $L$ -smooth  $f$ , the sequence  $\{\mathbf{x}_k\}$  produced by gradient descent satisfies

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) - t \left(1 - \frac{Lt}{2}\right) \|\nabla f(\mathbf{x}_k)\|^2$$

**Proof.** Plugging in  $\mathbf{x} = \mathbf{x}_k$  and  $\mathbf{y} = \mathbf{x}_{k+1} = \mathbf{x}_k - t\nabla f(\mathbf{x}_k)$  in the quadratic upper bound,

$$\begin{aligned} f(\mathbf{x}_{k+1}) &\leq f(\mathbf{x}_k) - t\|\nabla f(\mathbf{x}_k)\|^2 + \frac{L}{2}t^2\|\nabla f(\mathbf{x}_k)\|^2 \\ &= f(\mathbf{x}_k) - t \left(1 - \frac{Lt}{2}\right) \|\nabla f(\mathbf{x}_k)\|^2 \end{aligned}$$

**Note.** If  $\nabla f(\mathbf{x}_k) \neq 0$  and  $0 < t < \frac{2}{L}$ , then  $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$ , so gradient descent with step size  $t \in (0, 2/L)$  is indeed a descent method.

**Note.** We can lower bound the decrease in function value in each step. In particular, for  $0 < t \leq \frac{1}{L}$ ,

$$f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) \geq \frac{t}{2} \|\nabla f(\mathbf{x}_k)\|^2$$

# Convergence Analysis

**Theorem.** If  $f$  is convex and  $L$ -smooth, and  $\mathbf{x}^*$  is a minimum of  $f$ , then for step size  $t \in (0, \frac{1}{L}]$ , the sequence  $\{\mathbf{x}_k\}$  produced by the gradient descent algorithm satisfies

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|^2}{2tk}$$

## Notes.

- $f(\mathbf{x}_k) \downarrow f^*$  as  $k \rightarrow \infty$ .
- Any limiting point of  $\mathbf{x}_k$  is an optimal solution.
- The rate of convergence is  $O(1/k)$ , i.e. # of iterations to guarantee  $f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \epsilon$  is  $O(1/\epsilon)$ . For  $\epsilon = 10^{-p}$ ,  $k = O(10^p)$ , exponential in the number of significant digits!
- Faster convergence with larger  $t$ ; best  $t = \frac{1}{L}$ , but  $L$  is unknown.
- Good initial guess helps.

## Proof

1. By the basic gradient step  $\mathbf{x}_{k+1} = \mathbf{x}_k - t\nabla f(\mathbf{x}_k)$ ,

$$\begin{aligned}\|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 &= \|\mathbf{x}_k - t\nabla f(\mathbf{x}_k) - \mathbf{x}^*\|^2 \\ &= \|\mathbf{x}_k - \mathbf{x}^*\|^2 + t^2\|\nabla f(\mathbf{x}_k)\|^2 + 2t\nabla f(\mathbf{x}_k)^T(\mathbf{x}^* - \mathbf{x}_k)\end{aligned}$$

2. By the first-order condition for convexity,

$$\nabla f(\mathbf{x}_k)^T(\mathbf{x}^* - \mathbf{x}_k) \leq f(\mathbf{x}^*) - f(\mathbf{x}_k)$$

3. Plugging 2 into 1,

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 \leq \|\mathbf{x}_k - \mathbf{x}^*\|^2 + t^2\|\nabla f(\mathbf{x}_k)\|^2 + 2t[f(\mathbf{x}^*) - f(\mathbf{x}_k)]$$

4. Plugging in  $\frac{t}{2}\|\nabla f(\mathbf{x}_k)\|^2 \leq f(\mathbf{x}_k) - f(\mathbf{x}_{k+1})$  from slide 21,

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 \leq \|\mathbf{x}_k - \mathbf{x}^*\|^2 + 2t[f(\mathbf{x}^*) - f(\mathbf{x}_{k+1})]$$



## Proof (cont'd)

5. Rearranging,

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*) \leq \frac{\|\mathbf{x}_k - \mathbf{x}^*\|^2 - \|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2}{2t}$$

6. Summing over  $k$  from 0 to  $N - 1$ ,

$$\sum_{k=0}^{N-1} [f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*)] \leq \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|^2 - \|\mathbf{x}_N - \mathbf{x}^*\|^2}{2t} \leq \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|^2}{2t}$$

7. Recalling the descent property  $f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k)$ ,

$$f(\mathbf{x}_N) - f(\mathbf{x}^*) \leq \frac{1}{N} \sum_{k=0}^{N-1} [f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*)] \leq \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|^2}{2tN}$$