

# Recitation 1

18-460/18-660 Optimization

Chaoyi Pan

Carnegie Mellon University

# Recap from Calculus

$f : \mathbb{R}^n \rightarrow \mathbb{R}$  is continuously differentiable.

*Gradient:*

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \partial_{x_1} f(\mathbf{x}) \\ \partial_{x_2} f(\mathbf{x}) \\ \vdots \\ \partial_{x_n} f(\mathbf{x}) \end{bmatrix}$$

# Recap from Calculus

$f : \mathbb{R}^n \rightarrow \mathbb{R}$  is continuously differentiable.

*Gradient:*

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \partial_{x_1} f(\mathbf{x}) \\ \partial_{x_2} f(\mathbf{x}) \\ \vdots \\ \partial_{x_n} f(\mathbf{x}) \end{bmatrix}$$

**Example.**  $f(\mathbf{x}) = 3x_1^2 + 5 \log x_2$ .  $\nabla f(\mathbf{x}) = \left[ \quad \right]$

## Recap from Calculus

$f : \mathbb{R}^n \rightarrow \mathbb{R}$  is continuously differentiable.

*Gradient:*

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \partial_{x_1} f(\mathbf{x}) \\ \partial_{x_2} f(\mathbf{x}) \\ \vdots \\ \partial_{x_n} f(\mathbf{x}) \end{bmatrix}$$

**Example.**  $f(\mathbf{x}) = 3x_1^2 + 5 \log x_2$ .  $\nabla f(\mathbf{x}) = \begin{bmatrix} 6x_1 \\ \end{bmatrix}$

## Recap from Calculus

$f : \mathbb{R}^n \rightarrow \mathbb{R}$  is continuously differentiable.

*Gradient:*

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \partial_{x_1} f(\mathbf{x}) \\ \partial_{x_2} f(\mathbf{x}) \\ \vdots \\ \partial_{x_n} f(\mathbf{x}) \end{bmatrix}$$

**Example.**  $f(\mathbf{x}) = 3x_1^2 + 5 \log x_2$ .  $\nabla f(\mathbf{x}) = \begin{bmatrix} 6x_1 \\ 5x_2^{-1} \end{bmatrix}$

# Recap from Calculus

$f : \mathbb{R}^n \rightarrow \mathbb{R}$  is *twice* continuously differentiable.

*Hessian:*

$$\nabla_f^2(\mathbf{x}) = \begin{bmatrix} \partial_{x_1 x_1}^2 f(\mathbf{x}) & \partial_{x_1 x_2}^2 f(\mathbf{x}) & \cdots & \partial_{x_1 x_n}^2 f(\mathbf{x}) \\ \partial_{x_2 x_1}^2 f(\mathbf{x}) & \partial_{x_2 x_2}^2 f(\mathbf{x}) & \cdots & \partial_{x_2 x_n}^2 f(\mathbf{x}) \\ \vdots & \vdots & \ddots & \vdots \\ \partial_{x_n x_1}^2 f(\mathbf{x}) & \partial_{x_n x_2}^2 f(\mathbf{x}) & \cdots & \partial_{x_n x_n}^2 f(\mathbf{x}) \end{bmatrix}$$

# Recap from Calculus

$f : \mathbb{R}^n \rightarrow \mathbb{R}$  is *twice* continuously differentiable.

*Hessian:*

$$\nabla_f^2(\mathbf{x}) = \begin{bmatrix} \partial_{x_1 x_1}^2 f(\mathbf{x}) & \partial_{x_1 x_2}^2 f(\mathbf{x}) & \cdots & \partial_{x_1 x_n}^2 f(\mathbf{x}) \\ \partial_{x_2 x_1}^2 f(\mathbf{x}) & \partial_{x_2 x_2}^2 f(\mathbf{x}) & \cdots & \partial_{x_2 x_n}^2 f(\mathbf{x}) \\ \vdots & \vdots & \ddots & \vdots \\ \partial_{x_n x_1}^2 f(\mathbf{x}) & \partial_{x_n x_2}^2 f(\mathbf{x}) & \cdots & \partial_{x_n x_n}^2 f(\mathbf{x}) \end{bmatrix}$$

Because  $f \in \mathbf{C}^2$ ,  $\partial_{x_i x_j}^2 f = \partial_{x_j x_i}^2 f$ , and  $\nabla_f^2(\mathbf{x})$  is symmetric.

# Recap from Calculus

$f : \mathbb{R}^n \rightarrow \mathbb{R}$  is *twice* continuously differentiable.

*Hessian:*

$$\nabla_f^2(\mathbf{x}) = \begin{bmatrix} \partial_{x_1 x_1}^2 f(\mathbf{x}) & \partial_{x_1 x_2}^2 f(\mathbf{x}) & \cdots & \partial_{x_1 x_n}^2 f(\mathbf{x}) \\ \partial_{x_2 x_1}^2 f(\mathbf{x}) & \partial_{x_2 x_2}^2 f(\mathbf{x}) & \cdots & \partial_{x_2 x_n}^2 f(\mathbf{x}) \\ \vdots & \vdots & \ddots & \vdots \\ \partial_{x_n x_1}^2 f(\mathbf{x}) & \partial_{x_n x_2}^2 f(\mathbf{x}) & \cdots & \partial_{x_n x_n}^2 f(\mathbf{x}) \end{bmatrix}$$

Because  $f \in \mathbf{C}^2$ ,  $\partial_{x_i x_j}^2 f = \partial_{x_j x_i}^2 f$ , and  $\nabla_f^2(\mathbf{x})$  is symmetric.

**Example.**  $f(\mathbf{x}) = x_1^3 + 5x_2 + 7x_1 x_2$ .  $\nabla_f^2(\mathbf{x}) = \left[ \quad \right]$

# Recap from Calculus

$f : \mathbb{R}^n \rightarrow \mathbb{R}$  is *twice* continuously differentiable.

*Hessian:*

$$\nabla_f^2(\mathbf{x}) = \begin{bmatrix} \partial_{x_1 x_1}^2 f(\mathbf{x}) & \partial_{x_1 x_2}^2 f(\mathbf{x}) & \cdots & \partial_{x_1 x_n}^2 f(\mathbf{x}) \\ \partial_{x_2 x_1}^2 f(\mathbf{x}) & \partial_{x_2 x_2}^2 f(\mathbf{x}) & \cdots & \partial_{x_2 x_n}^2 f(\mathbf{x}) \\ \vdots & \vdots & \ddots & \vdots \\ \partial_{x_n x_1}^2 f(\mathbf{x}) & \partial_{x_n x_2}^2 f(\mathbf{x}) & \cdots & \partial_{x_n x_n}^2 f(\mathbf{x}) \end{bmatrix}$$

Because  $f \in \mathbf{C}^2$ ,  $\partial_{x_i x_j}^2 f = \partial_{x_j x_i}^2 f$ , and  $\nabla_f^2(\mathbf{x})$  is symmetric.

**Example.**  $f(\mathbf{x}) = x_1^3 + 5x_2 + 7x_1 x_2$ .  $\nabla_f^2(\mathbf{x}) = \begin{bmatrix} 6x_1 & \\ & \end{bmatrix}$

# Recap from Calculus

$f : \mathbb{R}^n \rightarrow \mathbb{R}$  is *twice* continuously differentiable.

*Hessian:*

$$\nabla_f^2(\mathbf{x}) = \begin{bmatrix} \partial_{x_1 x_1}^2 f(\mathbf{x}) & \partial_{x_1 x_2}^2 f(\mathbf{x}) & \cdots & \partial_{x_1 x_n}^2 f(\mathbf{x}) \\ \partial_{x_2 x_1}^2 f(\mathbf{x}) & \partial_{x_2 x_2}^2 f(\mathbf{x}) & \cdots & \partial_{x_2 x_n}^2 f(\mathbf{x}) \\ \vdots & \vdots & \ddots & \vdots \\ \partial_{x_n x_1}^2 f(\mathbf{x}) & \partial_{x_n x_2}^2 f(\mathbf{x}) & \cdots & \partial_{x_n x_n}^2 f(\mathbf{x}) \end{bmatrix}$$

Because  $f \in \mathbf{C}^2$ ,  $\partial_{x_i x_j}^2 f = \partial_{x_j x_i}^2 f$ , and  $\nabla_f^2(\mathbf{x})$  is symmetric.

**Example.**  $f(\mathbf{x}) = x_1^3 + 5x_2 + 7x_1 x_2$ .  $\nabla_f^2(\mathbf{x}) = \begin{bmatrix} 6x_1 & 7 \\ 7 & \end{bmatrix}$

# Recap from Calculus

$f : \mathbb{R}^n \rightarrow \mathbb{R}$  is twice continuously differentiable.

Hessian:

$$\nabla_f^2(\mathbf{x}) = \begin{bmatrix} \partial_{x_1 x_1}^2 f(\mathbf{x}) & \partial_{x_1 x_2}^2 f(\mathbf{x}) & \cdots & \partial_{x_1 x_n}^2 f(\mathbf{x}) \\ \partial_{x_2 x_1}^2 f(\mathbf{x}) & \partial_{x_2 x_2}^2 f(\mathbf{x}) & \cdots & \partial_{x_2 x_n}^2 f(\mathbf{x}) \\ \vdots & \vdots & \ddots & \vdots \\ \partial_{x_n x_1}^2 f(\mathbf{x}) & \partial_{x_n x_2}^2 f(\mathbf{x}) & \cdots & \partial_{x_n x_n}^2 f(\mathbf{x}) \end{bmatrix}$$

Because  $f \in \mathbf{C}^2$ ,  $\partial_{x_i x_j}^2 f = \partial_{x_j x_i}^2 f$ , and  $\nabla_f^2(\mathbf{x})$  is symmetric.

**Example.**  $f(\mathbf{x}) = x_1^3 + 5x_2 + 7x_1 x_2$ .  $\nabla_f^2(\mathbf{x}) = \begin{bmatrix} 6x_1 & 7 \\ 7 & 0 \end{bmatrix}$

## Recap from Calculus

$f : \mathbb{R}^n \rightarrow \mathbb{R}$  is continuously differentiable.

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

## Recap from Calculus

$f : \mathbb{R}^n \rightarrow \mathbb{R}$  is continuously differentiable.

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

Integration of *projection* of gradients along the line through  $\mathbf{x}$  and  $\mathbf{y}$ .

## Recap from Calculus

$f : \mathbb{R}^n \rightarrow \mathbb{R}$  is continuously differentiable.

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

Integration of *projection* of gradients along the line through  $\mathbf{x}$  and  $\mathbf{y}$ .

**Proof.** Let  $\mathbf{z}(t) = t\mathbf{y} + (1-t)\mathbf{x}$  and  $h(t) = f(\mathbf{z}(t))$ .

## Recap from Calculus

$f : \mathbb{R}^n \rightarrow \mathbb{R}$  is continuously differentiable.

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

Integration of *projection* of gradients along the line through  $\mathbf{x}$  and  $\mathbf{y}$ .

**Proof.** Let  $\mathbf{z}(t) = t\mathbf{y} + (1-t)\mathbf{x}$  and  $h(t) = f(\mathbf{z}(t))$ .

$$\mathbf{z}(0) = \mathbf{x}, \quad \mathbf{z}(1) = \mathbf{y}.$$

## Recap from Calculus

$f : \mathbb{R}^n \rightarrow \mathbb{R}$  is continuously differentiable.

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

Integration of *projection* of gradients along the line through  $\mathbf{x}$  and  $\mathbf{y}$ .

**Proof.** Let  $\mathbf{z}(t) = t\mathbf{y} + (1-t)\mathbf{x}$  and  $h(t) = f(\mathbf{z}(t))$ .

$\mathbf{z}(0) = \mathbf{x}$ ,  $\mathbf{z}(1) = \mathbf{y}$ .

$$h'(t) = \nabla f(\mathbf{z}(t))^\top \begin{bmatrix} z'_1(t) \\ \vdots \\ z'_n(t) \end{bmatrix}$$

## Recap from Calculus

$f : \mathbb{R}^n \rightarrow \mathbb{R}$  is continuously differentiable.

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

Integration of *projection* of gradients along the line through  $\mathbf{x}$  and  $\mathbf{y}$ .

**Proof.** Let  $\mathbf{z}(t) = t\mathbf{y} + (1-t)\mathbf{x}$  and  $h(t) = f(\mathbf{z}(t))$ .

$$\mathbf{z}(0) = \mathbf{x}, \quad \mathbf{z}(1) = \mathbf{y}.$$

$$h'(t) = \nabla f(\mathbf{z}(t))^\top \begin{bmatrix} z'_1(t) \\ \vdots \\ z'_n(t) \end{bmatrix} = \nabla f(\mathbf{z}(t))^\top \begin{bmatrix} y_1 - x_1 \\ \vdots \\ y_n - x_n \end{bmatrix}$$

## Recap from Calculus

$f : \mathbb{R}^n \rightarrow \mathbb{R}$  is continuously differentiable.

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

Integration of *projection* of gradients along the line through  $\mathbf{x}$  and  $\mathbf{y}$ .

**Proof.** Let  $\mathbf{z}(t) = t\mathbf{y} + (1-t)\mathbf{x}$  and  $h(t) = f(\mathbf{z}(t))$ .

$$\mathbf{z}(0) = \mathbf{x}, \quad \mathbf{z}(1) = \mathbf{y}.$$

$$\begin{aligned} h'(t) &= \nabla f(\mathbf{z}(t))^\top \begin{bmatrix} z'_1(t) \\ \vdots \\ z'_n(t) \end{bmatrix} = \nabla f(\mathbf{z}(t))^\top \begin{bmatrix} y_1 - x_1 \\ \vdots \\ y_n - x_n \end{bmatrix} = \\ &\nabla f(\mathbf{z}(t))^\top (\mathbf{y} - \mathbf{x}). \end{aligned}$$

## Recap from Calculus

$f : \mathbb{R}^n \rightarrow \mathbb{R}$  is continuously differentiable.

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

Integration of *projection* of gradients along the line through  $\mathbf{x}$  and  $\mathbf{y}$ .

**Proof.** Let  $\mathbf{z}(t) = t\mathbf{y} + (1-t)\mathbf{x}$  and  $h(t) = f(\mathbf{z}(t))$ .

$$\mathbf{z}(0) = \mathbf{x}, \quad \mathbf{z}(1) = \mathbf{y}.$$

$$\begin{aligned} h'(t) &= \nabla f(\mathbf{z}(t))^\top \begin{bmatrix} z'_1(t) \\ \vdots \\ z'_n(t) \end{bmatrix} = \nabla f(\mathbf{z}(t))^\top \begin{bmatrix} y_1 - x_1 \\ \vdots \\ y_n - x_n \end{bmatrix} = \\ &\nabla f(\mathbf{z}(t))^\top (\mathbf{y} - \mathbf{x}). \end{aligned}$$

$$f(\mathbf{y}) - f(\mathbf{x}) = h(1) - h(0) = \int_0^1 h'(t) dt = \int_0^1 \nabla f(t\mathbf{y} + (1-t)\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) dt.$$

## Recap from Calculus

$f : \mathbb{R} \rightarrow \mathbb{R}$  is *twice* continuously differentiable.

For any  $x, y \in \mathbb{R}$ , there exists  $t \in [0, 1]$  such that

$$f(y) = f(x) + f'(x)(y - x) + \frac{1}{2}f''(ty + (1 - t)x)(y - x)^2.$$

## Recap from Calculus

$f : \mathbb{R} \rightarrow \mathbb{R}$  is twice continuously differentiable.

For any  $x, y \in \mathbb{R}$ , there exists  $t \in [0, 1]$  such that

$$f(y) = f(x) + f'(x)(y - x) + \frac{1}{2}f''(ty + (1 - t)x)(y - x)^2.$$

What if  $f$  is multi-variate?

## Recap from Calculus

$f : \mathbb{R} \rightarrow \mathbb{R}$  is twice continuously differentiable.

For any  $x, y \in \mathbb{R}$ , there exists  $t \in [0, 1]$  such that

$$f(y) = f(x) + f'(x)(y - x) + \frac{1}{2}f''(ty + (1 - t)x)(y - x)^2.$$

What if  $f$  is multi-variate?

**Hint:** Define  $h(s) = f(s\mathbf{y} + (1 - s)\mathbf{x})$ , and do Taylor's expansion of  $h(1)$  on  $s = 0$ .

## Recap from Calculus

$f : \mathbb{R} \rightarrow \mathbb{R}$  is twice continuously differentiable.

For any  $x, y \in \mathbb{R}$ , there exists  $t \in [0, 1]$  such that

$$f(y) = f(x) + f'(x)(y - x) + \frac{1}{2}f''(ty + (1 - t)x)(y - x)^2.$$

What if  $f$  is multi-variate?

**Hint:** Define  $h(s) = f(s\mathbf{y} + (1 - s)\mathbf{x})$ , and do Taylor's expansion of  $h(1)$  on  $s = 0$ .

$$f(\mathbf{y}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{1}{2}(\mathbf{y} - \mathbf{x})^\top \nabla_f^2(t\mathbf{y} + (1 - t)\mathbf{x})(\mathbf{y} - \mathbf{x})$$

# Cauchy-Schwarz Inequality

$\mathbf{x}, \mathbf{y} \in \mathbb{R}^n.$

$$\mathbf{x}^T \mathbf{y} \leq \|\mathbf{x}\|_2 \|\mathbf{y}\|_2.$$

## Cauchy-Schwarz Inequality

$\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ .

$$\mathbf{x}^\top \mathbf{y} \leq \|\mathbf{x}\|_2 \|\mathbf{y}\|_2.$$

**Proof.** Consider (univariate) quadratic function

$$h(t) = \frac{1}{2}(\mathbf{t}\mathbf{x} + \mathbf{y})^\top (\mathbf{t}\mathbf{x} + \mathbf{y}) = \frac{1}{2}\|\mathbf{x}\|_2^2 t^2 + \mathbf{x}^\top \mathbf{y} t + \frac{1}{2}\|\mathbf{y}\|_2^2.$$

## Cauchy-Schwarz Inequality

$\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ .

$$\mathbf{x}^\top \mathbf{y} \leq \|\mathbf{x}\|_2 \|\mathbf{y}\|_2.$$

**Proof.** Consider (univariate) quadratic function

$$h(t) = \frac{1}{2}(\mathbf{t}\mathbf{x} + \mathbf{y})^\top (\mathbf{t}\mathbf{x} + \mathbf{y}) = \frac{1}{2}\|\mathbf{x}\|_2^2 t^2 + \mathbf{x}^\top \mathbf{y} t + \frac{1}{2}\|\mathbf{y}\|_2^2.$$

$$h(t) \geq 0, \quad \forall t \in \mathbb{R}.$$

## Cauchy-Schwarz Inequality

$\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ .

$$\mathbf{x}^\top \mathbf{y} \leq \|\mathbf{x}\|_2 \|\mathbf{y}\|_2.$$

**Proof.** Consider (univariate) quadratic function

$$h(t) = \frac{1}{2}(\mathbf{t}\mathbf{x} + \mathbf{y})^\top (\mathbf{t}\mathbf{x} + \mathbf{y}) = \frac{1}{2}\|\mathbf{x}\|_2^2 t^2 + \mathbf{x}^\top \mathbf{y} t + \frac{1}{2}\|\mathbf{y}\|_2^2.$$

$$h(t) \geq 0, \quad \forall t \in \mathbb{R}.$$

Reminder: quadratic function  $ax^2 + bx + c$  has less than 2 distinct real roots  $\Leftrightarrow \Delta = b^2 - 4ac \leq 0$ .

# Cauchy-Schwarz Inequality

$\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ .

$$\mathbf{x}^\top \mathbf{y} \leq \|\mathbf{x}\|_2 \|\mathbf{y}\|_2.$$

**Proof.** Consider (univariate) quadratic function

$$h(t) = \frac{1}{2}(\mathbf{t}\mathbf{x} + \mathbf{y})^\top (\mathbf{t}\mathbf{x} + \mathbf{y}) = \frac{1}{2}\|\mathbf{x}\|_2^2 t^2 + \mathbf{x}^\top \mathbf{y} t + \frac{1}{2}\|\mathbf{y}\|_2^2.$$

$$h(t) \geq 0, \quad \forall t \in \mathbb{R}.$$

Reminder: quadratic function  $ax^2 + bx + c$  has less than 2 distinct real roots  $\Leftrightarrow \Delta = b^2 - 4ac \leq 0$ .

$$0 \geq \Delta = (\mathbf{x}^\top \mathbf{y})^2 - \|\mathbf{x}\|_2^2 \|\mathbf{y}\|_2^2.$$

## Lipschitz Functions

$f : \mathcal{D} \rightarrow \mathbb{R}$  is  $L$ -Lipschitz continuous if for every  $\mathbf{x}, \mathbf{y} \in \mathcal{D}$ ,

$$|f(\mathbf{y}) - f(\mathbf{x})| \leq L\|\mathbf{y} - \mathbf{x}\|_2.$$

## Lipschitz Functions

$f : \mathcal{D} \rightarrow \mathbb{R}$  is  $L$ -Lipschitz continuous if for every  $\mathbf{x}, \mathbf{y} \in \mathcal{D}$ ,

$$|f(\mathbf{y}) - f(\mathbf{x})| \leq L\|\mathbf{y} - \mathbf{x}\|_2.$$

A differentiable  $f : \mathcal{D} \rightarrow \mathbb{R}$  has  $L$ -Lipschitz continuous gradients if for every  $\mathbf{x}, \mathbf{y} \in \mathcal{D}$ ,

$$\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\|_2 \leq L\|\mathbf{y} - \mathbf{x}\|_2.$$

## Lipschitz Functions

$f : \mathcal{D} \rightarrow \mathbb{R}$  is  $L$ -Lipschitz continuous if for every  $\mathbf{x}, \mathbf{y} \in \mathcal{D}$ ,

$$|f(\mathbf{y}) - f(\mathbf{x})| \leq L\|\mathbf{y} - \mathbf{x}\|_2.$$

A differentiable  $f : \mathcal{D} \rightarrow \mathbb{R}$  has  $L$ -Lipschitz continuous gradients if for every  $\mathbf{x}, \mathbf{y} \in \mathcal{D}$ ,

$$\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\|_2 \leq L\|\mathbf{y} - \mathbf{x}\|_2.$$

*Gradient doesn't have unbounded difference between two close points.*

## Lipschitz Functions

$f : \mathcal{D} \rightarrow \mathbb{R}$  is  $L$ -Lipschitz continuous if for every  $\mathbf{x}, \mathbf{y} \in \mathcal{D}$ ,

$$|f(\mathbf{y}) - f(\mathbf{x})| \leq L\|\mathbf{y} - \mathbf{x}\|_2.$$

A differentiable  $f : \mathcal{D} \rightarrow \mathbb{R}$  has  $L$ -Lipschitz continuous gradients if for every  $\mathbf{x}, \mathbf{y} \in \mathcal{D}$ ,

$$\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\|_2 \leq L\|\mathbf{y} - \mathbf{x}\|_2.$$

*Gradient doesn't have unbounded difference between two close points.*

**Note.** Differentiable  $f : \mathcal{D} \rightarrow \mathbb{R}$  is  $m$ -strongly convex  $\implies$  for every  $\mathbf{x}, \mathbf{y} \in \mathcal{D}$ ,

$$\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\|_2 \geq m\|\mathbf{y} - \mathbf{x}\|_2.$$

## Lipschitz Functions

$f : \mathcal{D} \rightarrow \mathbb{R}$  is  $L$ -Lipschitz continuous if for every  $\mathbf{x}, \mathbf{y} \in \mathcal{D}$ ,

$$|f(\mathbf{y}) - f(\mathbf{x})| \leq L\|\mathbf{y} - \mathbf{x}\|_2.$$

A differentiable  $f : \mathcal{D} \rightarrow \mathbb{R}$  has  $L$ -Lipschitz continuous gradients if for every  $\mathbf{x}, \mathbf{y} \in \mathcal{D}$ ,

$$\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\|_2 \leq L\|\mathbf{y} - \mathbf{x}\|_2.$$

*Gradient doesn't have unbounded difference between two close points.*

**Note.** Differentiable  $f : \mathcal{D} \rightarrow \mathbb{R}$  is  $m$ -strongly convex  $\implies$  for every  $\mathbf{x}, \mathbf{y} \in \mathcal{D}$ ,

$$\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\|_2 \geq m\|\mathbf{y} - \mathbf{x}\|_2.$$

**Proof.**

$$\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\|_2 = \left\| \overset{(1)}{\left( \nabla f(\mathbf{y}) - \nabla f(\mathbf{x}) \right)^T} \right\|_2$$

## Lipschitz Functions

$f : \mathcal{D} \rightarrow \mathbb{R}$  is  $L$ -Lipschitz continuous if for every  $\mathbf{x}, \mathbf{y} \in \mathcal{D}$ ,

$$|f(\mathbf{y}) - f(\mathbf{x})| \leq L\|\mathbf{y} - \mathbf{x}\|_2.$$

A differentiable  $f : \mathcal{D} \rightarrow \mathbb{R}$  has  $L$ -Lipschitz continuous gradients if for every  $\mathbf{x}, \mathbf{y} \in \mathcal{D}$ ,

$$\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\|_2 \leq L\|\mathbf{y} - \mathbf{x}\|_2.$$

*Gradient doesn't have unbounded difference between two close points.*

**Note.** Differentiable  $f : \mathcal{D} \rightarrow \mathbb{R}$  is  $m$ -strongly convex  $\implies$  for every  $\mathbf{x}, \mathbf{y} \in \mathcal{D}$ ,

$$\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\|_2 \geq m\|\mathbf{y} - \mathbf{x}\|_2.$$

**Proof.**

$$\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\|_2 = \left\| \overset{(1)}{\left( \nabla f(\mathbf{y}) - \nabla f(\mathbf{x}) \right)^T} \overset{(2)}{\left( \mathbf{y} - \mathbf{x} \right)} \right\|_2$$

## Eigenvalues and Eigenvectors

Let  $A \in \mathbb{R}^{n \times n}$ . A nonzero vector  $v$  is an eigenvector with eigenvalue  $\lambda \in \mathbb{C}$  if

$$Av = \lambda v.$$

# Eigenvalues and Eigenvectors

Let  $A \in \mathbb{R}^{n \times n}$ . A nonzero vector  $v$  is an eigenvector with eigenvalue  $\lambda \in \mathbb{C}$  if

$$Av = \lambda v.$$

If  $A$  is real symmetric ( $A = A^T$ ), then

- ▶ All eigenvalues are real.
- ▶ There exists an orthonormal basis of eigenvectors.
- ▶  $A$  admits the eigendecomposition  $A = Q\Lambda Q^T$  with  $Q$  orthogonal,  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ .

## Rayleigh Quotient and Inequalities (Symmetric $A$ )

For  $A = A^T$  and any nonzero  $x \in \mathbb{R}^n$ , the Rayleigh quotient satisfies

$$\lambda_{\min}(A) \leq \frac{x^T A x}{x^T x} \leq \lambda_{\max}(A).$$

## Rayleigh Quotient and Inequalities (Symmetric $A$ )

For  $A = A^T$  and any nonzero  $x \in \mathbb{R}^n$ , the Rayleigh quotient satisfies

$$\lambda_{\min}(A) \leq \frac{x^T A x}{x^T x} \leq \lambda_{\max}(A).$$

Equivalent statements:

- ▶  $A \succeq 0$  (positive semidefinite, PSD)  $\Leftrightarrow x^T A x \geq 0$  for all  $x$ .
- ▶  $A \preceq B \Leftrightarrow B - A \succeq 0$  (matrix Loewner order).

## Positive Semidefinite (PSD) and Notation $A \succeq B$

- ▶  $A \succeq 0$  (PSD)  $\Leftrightarrow x^T A x \geq 0$  for all  $x \in \mathbb{R}^n$ .
- ▶  $A \succ 0$  (PD)  $\Leftrightarrow x^T A x > 0$  for all nonzero  $x$ .
- ▶  $A \succeq B$  means  $A - B \succeq 0$ .

## Positive Semidefinite (PSD) and Notation $A \succeq B$

- ▶  $A \succeq 0$  (PSD)  $\Leftrightarrow x^T A x \geq 0$  for all  $x \in \mathbb{R}^n$ .
- ▶  $A \succ 0$  (PD)  $\Leftrightarrow x^T A x > 0$  for all nonzero  $x$ .
- ▶  $A \succeq B$  means  $A - B \succeq 0$ .

For symmetric  $A$ :

$$A \succeq 0 \Leftrightarrow \lambda_i(A) \geq 0 \ \forall i \Leftrightarrow A = Q \Lambda Q^T, \Lambda \succeq 0.$$

# Singular Value Decomposition (SVD)

For any  $A \in \mathbb{R}^{m \times n}$ ,

$$A = U\Sigma V^T,$$

where  $U \in \mathbb{R}^{m \times m}$  and  $V \in \mathbb{R}^{n \times n}$  are orthogonal, and  $\Sigma \in \mathbb{R}^{m \times n}$  is diagonal with nonnegative entries  $\sigma_1 \geq \dots \geq \sigma_r > 0$  (singular values).

# Singular Value Decomposition (SVD)

For any  $A \in \mathbb{R}^{m \times n}$ ,

$$A = U\Sigma V^T,$$

where  $U \in \mathbb{R}^{m \times m}$  and  $V \in \mathbb{R}^{n \times n}$  are orthogonal, and  $\Sigma \in \mathbb{R}^{m \times n}$  is diagonal with nonnegative entries  $\sigma_1 \geq \dots \geq \sigma_r > 0$  (singular values). Useful facts:

- ▶  $\sigma_i(A) = \sqrt{\lambda_i(A^T A)}$ .
- ▶  $\|A\|_2 = \sigma_{\max}(A)$ ,  $\|A\|_F^2 = \sum_i \sigma_i^2$ .
- ▶ Best rank- $k$  approximation (Eckart–Young):  
 $A_k = U_{:,1:k} \Sigma_{1:k,1:k} V_{:,1:k}^T$  minimizes  $\|A - X\|_F$  over  $\text{rank}(X) \leq k$ .

# Common Inequalities

- ▶ Cauchy–Schwarz:  $|x^T y| \leq \|x\|_2 \|y\|_2$ .
- ▶ For PSD  $A$ :  $|x^T A y| \leq \sqrt{x^T A x} \sqrt{y^T A y}$ .
- ▶ Submultiplicativity:  $\|A x\|_2 \leq \|A\|_2 \|x\|_2$ .
- ▶ For symmetric  $A$ :  $\lambda_{\min}(A) \|x\|_2^2 \leq x^T A x \leq \lambda_{\max}(A) \|x\|_2^2$ .

## Use Hessian Matrix: Example 1

Prove that  $f(x, y) = xy + a(x^2 + y^2)$  is convex if and only if  $a \geq 1/2$ .

## Use Hessian Matrix: Example 1

Prove that  $f(x, y) = xy + a(x^2 + y^2)$  is convex if and only if  $a \geq 1/2$ .

$$\nabla_f^2(x, y) = \begin{bmatrix} 2a & 1 \\ 1 & 2a \end{bmatrix}.$$

## Use Hessian Matrix: Example 1

Prove that  $f(x, y) = xy + a(x^2 + y^2)$  is convex if and only if  $a \geq 1/2$ .

$$\nabla_f^2(x, y) = \begin{bmatrix} 2a & 1 \\ 1 & 2a \end{bmatrix}.$$

Let  $\lambda$  be its eigenvalue.

## Use Hessian Matrix: Example 1

Prove that  $f(x, y) = xy + a(x^2 + y^2)$  is convex if and only if  $a \geq 1/2$ .

$$\nabla_f^2(x, y) = \begin{bmatrix} 2a & 1 \\ 1 & 2a \end{bmatrix}.$$

Let  $\lambda$  be its eigenvalue.

$$0 = \begin{vmatrix} 2a - \lambda & 1 \\ 1 & 2a - \lambda \end{vmatrix} = (2a - \lambda)^2 - 1.$$

## Use Hessian Matrix: Example 1

Prove that  $f(x, y) = xy + a(x^2 + y^2)$  is convex if and only if  $a \geq 1/2$ .

$$\nabla_f^2(x, y) = \begin{bmatrix} 2a & 1 \\ 1 & 2a \end{bmatrix}.$$

Let  $\lambda$  be its eigenvalue.

$$0 = \begin{vmatrix} 2a - \lambda & 1 \\ 1 & 2a - \lambda \end{vmatrix} = (2a - \lambda)^2 - 1.$$

$$\lambda = 2a - 1, 2a + 1.$$

## Use Hessian Matrix: Example 1

Prove that  $f(x, y) = xy + a(x^2 + y^2)$  is convex if and only if  $a \geq 1/2$ .

$$\nabla_f^2(x, y) = \begin{bmatrix} 2a & 1 \\ 1 & 2a \end{bmatrix}.$$

Let  $\lambda$  be its eigenvalue.

$$0 = \begin{vmatrix} 2a - \lambda & 1 \\ 1 & 2a - \lambda \end{vmatrix} = (2a - \lambda)^2 - 1.$$

$$\lambda = 2a - 1, 2a + 1.$$

Both  $\lambda$ 's are non-negative iff  $a \geq 1/2$ .

## Use Hessian Matrix: Example 2

Prove that  $f(\beta_1, \beta_2) = \log(1 + e^{\beta_1 - 2\beta_2})$  is convex over  $\mathbb{R}^2$ .

## Use Hessian Matrix: Example 2

Prove that  $f(\beta_1, \beta_2) = \log(1 + e^{\beta_1 - 2\beta_2})$  is convex over  $\mathbb{R}^2$ .

$$\nabla f(\boldsymbol{\beta}) = \left[ \frac{e^{\beta_1 - 2\beta_2}}{1 + e^{\beta_1 - 2\beta_2}} \right]^\top.$$

## Use Hessian Matrix: Example 2

Prove that  $f(\beta_1, \beta_2) = \log(1 + e^{\beta_1 - 2\beta_2})$  is convex over  $\mathbb{R}^2$ .

$$\nabla f(\beta) = \begin{bmatrix} \frac{e^{\beta_1 - 2\beta_2}}{1 + e^{\beta_1 - 2\beta_2}} & -\frac{2e^{\beta_1 - 2\beta_2}}{1 + e^{\beta_1 - 2\beta_2}} \end{bmatrix}^\top.$$

## Use Hessian Matrix: Example 2

Prove that  $f(\beta_1, \beta_2) = \log(1 + e^{\beta_1 - 2\beta_2})$  is convex over  $\mathbb{R}^2$ .

$$\nabla f(\beta) = \begin{bmatrix} \frac{e^{\beta_1 - 2\beta_2}}{1 + e^{\beta_1 - 2\beta_2}} & -\frac{2e^{\beta_1 - 2\beta_2}}{1 + e^{\beta_1 - 2\beta_2}} \end{bmatrix}^\top.$$

$$\nabla_f^2(\beta) = \begin{bmatrix} \frac{e^{\beta_1 - 2\beta_2}}{(1 + e^{\beta_1 - 2\beta_2})^2} \\ \\ \end{bmatrix}$$

## Use Hessian Matrix: Example 2

Prove that  $f(\beta_1, \beta_2) = \log(1 + e^{\beta_1 - 2\beta_2})$  is convex over  $\mathbb{R}^2$ .

$$\nabla f(\beta) = \begin{bmatrix} \frac{e^{\beta_1 - 2\beta_2}}{1 + e^{\beta_1 - 2\beta_2}} & -\frac{2e^{\beta_1 - 2\beta_2}}{1 + e^{\beta_1 - 2\beta_2}} \end{bmatrix}^\top.$$

$$\nabla^2 f(\beta) = \begin{bmatrix} \frac{e^{\beta_1 - 2\beta_2}}{(1 + e^{\beta_1 - 2\beta_2})^2} & -\frac{2e^{\beta_1 - 2\beta_2}}{(1 + e^{\beta_1 - 2\beta_2})^2} \\ -\frac{2e^{\beta_1 - 2\beta_2}}{(1 + e^{\beta_1 - 2\beta_2})^2} & \frac{2e^{\beta_1 - 2\beta_2}}{(1 + e^{\beta_1 - 2\beta_2})^2} \end{bmatrix}$$

## Use Hessian Matrix: Example 2

Prove that  $f(\beta_1, \beta_2) = \log(1 + e^{\beta_1 - 2\beta_2})$  is convex over  $\mathbb{R}^2$ .

$$\nabla f(\beta) = \begin{bmatrix} \frac{e^{\beta_1 - 2\beta_2}}{1 + e^{\beta_1 - 2\beta_2}} & -\frac{2e^{\beta_1 - 2\beta_2}}{1 + e^{\beta_1 - 2\beta_2}} \end{bmatrix}^\top.$$

$$\nabla^2 f(\beta) = \begin{bmatrix} \frac{e^{\beta_1 - 2\beta_2}}{(1 + e^{\beta_1 - 2\beta_2})^2} & -\frac{2e^{\beta_1 - 2\beta_2}}{(1 + e^{\beta_1 - 2\beta_2})^2} \\ -\frac{2e^{\beta_1 - 2\beta_2}}{(1 + e^{\beta_1 - 2\beta_2})^2} & \frac{4e^{\beta_1 - 2\beta_2}}{(1 + e^{\beta_1 - 2\beta_2})^2} \end{bmatrix}$$

## Use Hessian Matrix: Example 2

Prove that  $f(\beta_1, \beta_2) = \log(1 + e^{\beta_1 - 2\beta_2})$  is convex over  $\mathbb{R}^2$ .

$$\nabla f(\beta) = \begin{bmatrix} \frac{e^{\beta_1 - 2\beta_2}}{1 + e^{\beta_1 - 2\beta_2}} & -\frac{2e^{\beta_1 - 2\beta_2}}{1 + e^{\beta_1 - 2\beta_2}} \end{bmatrix}^\top.$$

$$\nabla_f^2(\beta) = \begin{bmatrix} \frac{e^{\beta_1 - 2\beta_2}}{(1 + e^{\beta_1 - 2\beta_2})^2} & -\frac{2e^{\beta_1 - 2\beta_2}}{(1 + e^{\beta_1 - 2\beta_2})^2} \\ -\frac{2e^{\beta_1 - 2\beta_2}}{(1 + e^{\beta_1 - 2\beta_2})^2} & \frac{4e^{\beta_1 - 2\beta_2}}{(1 + e^{\beta_1 - 2\beta_2})^2} \end{bmatrix} \propto + \begin{bmatrix} 1 & -2 \\ -2 & 4 \end{bmatrix}.$$

## Use Hessian Matrix: Example 2

Prove that  $f(\beta_1, \beta_2) = \log(1 + e^{\beta_1 - 2\beta_2})$  is convex over  $\mathbb{R}^2$ .

$$\nabla f(\beta) = \begin{bmatrix} e^{\beta_1 - 2\beta_2} \\ 1 + e^{\beta_1 - 2\beta_2} \end{bmatrix}^T - \frac{2e^{\beta_1 - 2\beta_2}}{1 + e^{\beta_1 - 2\beta_2}}.$$

$$\nabla_f^2(\beta) = \begin{bmatrix} e^{\beta_1 - 2\beta_2} & -\frac{2e^{\beta_1 - 2\beta_2}}{(1 + e^{\beta_1 - 2\beta_2})^2} \\ -\frac{2e^{\beta_1 - 2\beta_2}}{(1 + e^{\beta_1 - 2\beta_2})^2} & \frac{4e^{\beta_1 - 2\beta_2}}{(1 + e^{\beta_1 - 2\beta_2})^2} \end{bmatrix} \propto + \begin{bmatrix} 1 & -2 \\ -2 & 4 \end{bmatrix}.$$

The eigenvalues are 0 and  $\frac{5e^{\beta_1 - 2\beta_2}}{(1 + e^{\beta_1 - 2\beta_2})^2} \geq 0$ .

$\nabla_f^2(\beta) \succeq 0$ .  $f$  is convex.

## Use Hessian Matrix: Example 2

$f(\beta_1, \beta_2) = \log(1 + e^{\beta_1 - 2\beta_2})$  is convex over  $\mathbb{R}^2$ .

*Extension.* Is  $f(\beta) = \log(1 + e^{\mathbf{c}^\top \beta})$  always convex for any constant vector  $\mathbf{c}$ ?

## Use Hessian Matrix: Example 2

$f(\beta_1, \beta_2) = \log(1 + e^{\beta_1 - 2\beta_2})$  is convex over  $\mathbb{R}^2$ .

*Extension.* Is  $f(\beta) = \log(1 + e^{\mathbf{c}^\top \beta})$  always convex for any constant vector  $\mathbf{c}$ ?

$$\nabla_f^2(\beta) = \frac{e^{\mathbf{c}^\top \beta}}{(1 + e^{\mathbf{c}^\top \beta})^2} \begin{bmatrix} c_1^2 & c_1 c_2 & \cdots & c_1 c_n \\ c_1 c_2 & c_2^2 & \cdots & c_2 c_n \\ \vdots & \vdots & \ddots & \vdots \\ c_1 c_n & c_2 c_n & \cdots & c_n^2 \end{bmatrix} \propto_+ \mathbf{c} \mathbf{c}^\top.$$

## Use Hessian Matrix: Example 2

$f(\beta_1, \beta_2) = \log(1 + e^{\beta_1 - 2\beta_2})$  is convex over  $\mathbb{R}^2$ .

*Extension.* Is  $f(\beta) = \log(1 + e^{\mathbf{c}^\top \beta})$  always convex for any constant vector  $\mathbf{c}$ ?

$$\nabla_f^2(\beta) = \frac{e^{\mathbf{c}^\top \beta}}{(1 + e^{\mathbf{c}^\top \beta})^2} \begin{bmatrix} c_1^2 & c_1 c_2 & \cdots & c_1 c_n \\ c_1 c_2 & c_2^2 & \cdots & c_2 c_n \\ \vdots & \vdots & \ddots & \vdots \\ c_1 c_n & c_2 c_n & \cdots & c_n^2 \end{bmatrix} \propto_+ \mathbf{c} \mathbf{c}^\top.$$

Why is  $\nabla_f^2(\beta)$  always positive semi-definite?

## Use Hessian Matrix: Example 2

$f(\beta_1, \beta_2) = \log(1 + e^{\beta_1 - 2\beta_2})$  is convex over  $\mathbb{R}^2$ .

*Extension.* Is  $f(\beta) = \log(1 + e^{\mathbf{c}^\top \beta})$  always convex for any constant vector  $\mathbf{c}$ ?

$$\nabla_f^2(\beta) = \frac{e^{\mathbf{c}^\top \beta}}{(1 + e^{\mathbf{c}^\top \beta})^2} \begin{bmatrix} c_1^2 & c_1 c_2 & \cdots & c_1 c_n \\ c_1 c_2 & c_2^2 & \cdots & c_2 c_n \\ \vdots & \vdots & \ddots & \vdots \\ c_1 c_n & c_2 c_n & \cdots & c_n^2 \end{bmatrix} \propto_+ \mathbf{c} \mathbf{c}^\top.$$

Why is  $\nabla_f^2(\beta)$  always positive semi-definite?

Because for every  $\mathbf{x}, \beta \in \mathbb{R}^n$ ,

$$\mathbf{x}^\top \nabla_f^2(\beta) \mathbf{x} = \frac{e^{\mathbf{c}^\top \beta}}{(1 + e^{\mathbf{c}^\top \beta})^2} \mathbf{x}^\top \mathbf{c} \mathbf{c}^\top \mathbf{x} = \frac{e^{\mathbf{c}^\top \beta}}{(1 + e^{\mathbf{c}^\top \beta})^2} (\mathbf{x}^\top \mathbf{c})^2 \geq 0.$$

## Use Hessian Matrix: Example 2

$f(\beta) = \log(1 + e^{\mathbf{c}^\top \beta})$  is always convex over  $\mathbb{R}^n$ .

*Extension II.* Is  $f(\beta) = \log(1 + e^{\mathbf{c}^\top \beta}) - \mathbf{m}^\top \beta + \text{Const}$  always convex for any constant vector  $\mathbf{c}, \mathbf{m}$ ?

## Use Hessian Matrix: Example 2

$f(\beta) = \log(1 + e^{\mathbf{c}^\top \beta})$  is always convex over  $\mathbb{R}^n$ .

*Extension II.* Is  $f(\beta) = \log(1 + e^{\mathbf{c}^\top \beta}) - \mathbf{m}^\top \beta + \text{Const}$  always convex for any constant vector  $\mathbf{c}, \mathbf{m}$ ?

*Application: Logistic Regression.*  $X_i \in \mathbb{R}^n, y_i \in \{0, 1\}$  are given data.

## Use Hessian Matrix: Example 2

$f(\beta) = \log(1 + e^{\mathbf{c}^\top \beta})$  is always convex over  $\mathbb{R}^n$ .

*Extension II.* Is  $f(\beta) = \log(1 + e^{\mathbf{c}^\top \beta}) - \mathbf{m}^\top \beta + \text{Const}$  always convex for any constant vector  $\mathbf{c}, \mathbf{m}$ ?

*Application: Logistic Regression.*  $X_i \in \mathbb{R}^n, y_i \in \{0, 1\}$  are given data.

We want a good  $\beta$  for predicting  $y$  with  $\hat{y} \in [0, 1]$  given another  $X$ .

## Use Hessian Matrix: Example 2

$f(\beta) = \log(1 + e^{\mathbf{c}^\top \beta})$  is always convex over  $\mathbb{R}^n$ .

*Extension II.* Is  $f(\beta) = \log(1 + e^{\mathbf{c}^\top \beta}) - \mathbf{m}^\top \beta + \text{Const}$  always convex for any constant vector  $\mathbf{c}, \mathbf{m}$ ?

*Application: Logistic Regression.*  $X_i \in \mathbb{R}^n, y_i \in \{0, 1\}$  are given data.

We want a good  $\beta$  for predicting  $y$  with  $\hat{y} \in [0, 1]$  given another  $X$ .

With a good  $\beta^*$ , we predict that

$$\hat{y} = \frac{1}{1 + e^{-X^\top \beta^*}}.$$

## Use Hessian Matrix: Example 2

$f(\beta) = \log(1 + e^{\mathbf{c}^\top \beta})$  is always convex over  $\mathbb{R}^n$ .

*Extension II.* Is  $f(\beta) = \log(1 + e^{\mathbf{c}^\top \beta}) - \mathbf{m}^\top \beta + \text{Const}$  always convex for any constant vector  $\mathbf{c}, \mathbf{m}$ ?

*Application: Logistic Regression.*  $X_i \in \mathbb{R}^n, y_i \in \{0, 1\}$  are given data.

We want a good  $\beta$  for predicting  $y$  with  $\hat{y} \in [0, 1]$  given another  $X$ .

With a good  $\beta^*$ , we predict that

$$\hat{y} = \frac{1}{1 + e^{-X^\top \beta^*}}.$$

*Maximum Likelihood Estimation (MLE)* gives that