

## Derivative

- Function  $\mathbb{R}^n \rightarrow \mathbb{R} : f(\mathbf{x}) = a_1 x_1 + a_2 x_2$

$$\nabla_{\mathbf{x}} f(\mathbf{x}) \triangleq \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \frac{\partial f(\mathbf{x})}{\partial x_2} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_n} \end{bmatrix} \in \mathbb{R}^n$$

- $\frac{\partial f(\mathbf{x})}{\partial x_i}$  : partial derivative of element i of vector  $\mathbf{x}$
- If there is no variable but  $\mathbf{x}$ ,  $\nabla_{\mathbf{x}} f(\mathbf{x})$  can be written  $\nabla f(\mathbf{x})$
- If  $\mathbf{x} \in \mathbb{R}^n$ ,  $\nabla_{\mathbf{x}} f(\mathbf{x}) \in \mathbb{R}^n$

$$\nabla^2 f(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial x_1^2} & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 x_n} \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_2 x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial x_2^2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_2 x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_n x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial x_n x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_n^2} \end{bmatrix} \in \mathbb{R}^{n \times n}$$

- Second-order gradient is also called Hessian and symmetric matrix size n
- Function  $f(X) : \mathbb{R}^{n \times m} \rightarrow \mathbb{R}$

$$\nabla f(X) = \begin{bmatrix} \frac{\partial f(X)}{\partial x_{11}} & \frac{\partial f(X)}{\partial x_{12}} & \cdots & \frac{\partial f(X)}{\partial x_{1m}} \\ \frac{\partial f(X)}{\partial x_{21}} & \frac{\partial f(X)}{\partial x_{22}} & \cdots & \frac{\partial f(X)}{\partial x_{2m}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f(X)}{\partial x_{n1}} & \frac{\partial f(X)}{\partial x_{n2}} & \cdots & \frac{\partial f(X)}{\partial x_{nm}} \end{bmatrix} \in \mathbb{R}^{n \times m}$$

- Function  $v(x) : \mathbb{R} \rightarrow \mathbb{R}^n$

$$v(x) = \begin{bmatrix} v_1(x) \\ v_2(x) \\ \vdots \\ v_n(x) \end{bmatrix}$$

- Derivative of the function by  $x$  is the **row vector**

$$\nabla_x v(x) \triangleq \begin{bmatrix} \frac{\partial v_1(x)}{\partial x} & \frac{\partial v_2(x)}{\partial x} & \dots & \frac{\partial v_n(x)}{\partial x} \end{bmatrix} \in \mathbb{R}^{1 \times n}$$

$$\nabla_x^2 v(x) \triangleq \begin{bmatrix} \frac{\partial^2 v_1(x)}{\partial x^2} & \frac{\partial^2 v_2(x)}{\partial x^2} & \dots & \frac{\partial^2 v_n(x)}{\partial x^2} \end{bmatrix}$$

- Example: Given  $\mathbf{a} \in \mathbb{R}^n$  and vector-valued  $v(x) = x\mathbf{a}$ , 
$$\begin{cases} \nabla v(x) = \mathbf{a}^T \\ \nabla^2 v(x) = 0 \end{cases}$$

- Function  $h(\mathbf{x}): \mathbb{R}^k \rightarrow \mathbb{R}^n$

$$h(\mathbf{x}) = \begin{bmatrix} h_1(\mathbf{x}): \mathbb{R}^k \rightarrow \mathbb{R} \\ h_2(\mathbf{x}): \mathbb{R}^k \rightarrow \mathbb{R} \\ \vdots \\ h_n(\mathbf{x}): \mathbb{R}^k \rightarrow \mathbb{R} \end{bmatrix}$$

$$\nabla h(\mathbf{x}) \triangleq \begin{bmatrix} \frac{\partial h_1(\mathbf{x})}{\partial x_1} & \frac{\partial h_2(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial h_n(\mathbf{x})}{\partial x_1} \\ \frac{\partial h_1(\mathbf{x})}{\partial x_2} & \frac{\partial h_2(\mathbf{x})}{\partial x_2} & \dots & \frac{\partial h_n(\mathbf{x})}{\partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial h_1(\mathbf{x})}{\partial x_k} & \frac{\partial h_2(\mathbf{x})}{\partial x_k} & \dots & \frac{\partial h_n(\mathbf{x})}{\partial x_k} \end{bmatrix} \in \mathbb{R}^{k \times n}$$

## Derivative Properties

- Product Rule
  - Assume matrix  $X$  is the input and all functions are in proper size to multiply

$$\nabla \left( f(X)^T g(X) \right) = \left( \nabla f(X) \right) g(X) + \left( \nabla g(X) \right) f(X)$$

- Above rule is just another way to represent:  $(f(x)g(x))' = f'(x)g(x) + g'(x)f(x)$

- Chain Rule

$$\nabla_X g(f(X)) = \left( \nabla_X f \right)^T \left( \nabla_f g \right)$$

## Derivative of common function

- $f(\mathbf{x}) = \mathbf{a}^T \mathbf{x}$

- Assume  $\mathbf{a}, \mathbf{x} \in \mathbb{R}^n$ , we rewrite  $f(\mathbf{x}) = \mathbf{a}^T \mathbf{x} = a_1 x_1 + a_2 x_2 + \dots + a_n x_n$

- $\frac{\partial f(\mathbf{x})}{\partial x_i} = a_i, \forall i = 1, 2, 3, \dots, n$

$$\rightarrow \nabla f(\mathbf{x}) = [a_1 \quad a_2 \quad \dots \quad a_n]^T = \mathbf{a}$$

- $f(\mathbf{x}) = A\mathbf{x}$

- This is the function  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$  with  $\mathbf{x} \in \mathbb{R}^n, A \in \mathbb{R}^{m \times n}$ , so  $A\mathbf{x} = \begin{bmatrix} \mathbf{a}_1 \mathbf{x} \\ \mathbf{a}_2 \mathbf{x} \\ \vdots \\ \mathbf{a}_m \mathbf{x} \end{bmatrix}$
  - $\nabla_{\mathbf{x}}(A\mathbf{x}) = \begin{bmatrix} \mathbf{a}_1^T & \mathbf{a}_2^T & \cdots & \mathbf{a}_m^T \end{bmatrix} = A^T$

- $f(\mathbf{x}) = \mathbf{x}^T A \mathbf{x}$

- Assume  $\mathbf{x} \in \mathbb{R}^n, A \in \mathbb{R}^{n \times n}$ , we have

$$\begin{aligned} \nabla_{\mathbf{x}}^T A \mathbf{x} &= \nabla f(\mathbf{x}) = \nabla \left( (\mathbf{x}^T)(A\mathbf{x}) \right) \\ &= (\nabla(\mathbf{x}) A \mathbf{x}) + (\nabla(A\mathbf{x})) \mathbf{x} \\ &= I A \mathbf{x} + A^T \mathbf{x} \\ &= (A + A^T) \mathbf{x} \end{aligned}$$

$$\nabla_{\mathbf{x}}^2 \mathbf{x}^T A \mathbf{x} = A + A^T$$

- $f(\mathbf{x}) = \|A\mathbf{x} - b\|_2^2$

$$\left. \begin{aligned} \nabla(A\mathbf{x} - b) &= A^T \\ \nabla \|A\mathbf{x} - b\|_2^2 &= 2A\mathbf{x} - 2b \end{aligned} \right\} \rightarrow \nabla \|A\mathbf{x} - b\|_2^2 = 2A^T(A\mathbf{x} - b)$$

- $f(\mathbf{x}) = \mathbf{a}^T \mathbf{x} \mathbf{x}^T \mathbf{b}$

- Rewrite as  $f(\mathbf{x}) = (\mathbf{x}^T \mathbf{b})(\mathbf{a}^T \mathbf{x})$  and apply product rule:

$$\nabla(\mathbf{a}^T \mathbf{x} \mathbf{x}^T \mathbf{b}) = \mathbf{b} \mathbf{a}^T \mathbf{x} + \mathbf{x}^T \mathbf{b} \mathbf{a} = \mathbf{b} \mathbf{a}^T \mathbf{x} + \mathbf{b}^T \mathbf{x} \mathbf{a} = (\mathbf{a} \mathbf{b}^T + \mathbf{b} \mathbf{a}^T) \mathbf{x}$$

- $f(X) = \text{trace}(AX)$

- Assume  $A \in \mathbb{R}^{n \times m}, X \in \mathbb{R}^{m \times n}, B = AX \in \mathbb{R}^{n \times n}$ , based on the definition of trace:

$$f(X) = \text{trace}(AX) = \text{trace}(B) = \sum_{j=1}^n b_{jj} = \sum_{j=1}^n \sum_{i=1}^m a_{ji} x_{ij}$$

$$\rightarrow \frac{\partial f(X)}{\partial x_{ij}} = a_{ji} \rightarrow \nabla_{\mathbf{x}} \text{trace}(AX) = A^T$$

- $f(X) = \mathbf{a}^T X \mathbf{b}$

- Assume  $\mathbf{a} \in \mathbb{R}^m, X \in \mathbb{R}^{m \times n}, \mathbf{b} \in \mathbb{R}^n$ , after multiplying, the result:  $f(X) = \sum_{i=1}^m \sum_{j=1}^n x_{ij} a_i b_j$

$$\rightarrow \nabla_x (\mathbf{a}^T \mathbf{X} \mathbf{b}) = \begin{bmatrix} a_1 b_1 & a_1 b_2 & \cdots & a_1 b_n \\ a_2 b_1 & a_2 b_2 & \cdots & a_2 b_n \\ \vdots & \vdots & \ddots & \vdots \\ a_m b_1 & a_m b_2 & \cdots & a_m b_n \end{bmatrix} = \mathbf{a} \mathbf{b}^T$$

- $f(X) = \|X\|_F^2$

- Assume  $X \in \mathbb{R}^{n \times n}$ , rewrite  $\|X\|_F^2 = \sum_{i=1}^n \sum_{j=1}^n x_{ij}^2$ , so  $\frac{\partial f}{\partial x_{ij}} = 2x_{ij}$

$$\rightarrow \nabla \|X\|_F^2 = 2X$$

$f(\mathbf{x})$	$\nabla f(\mathbf{x})$	$f(X)$	$\nabla_x f(X)$
$\mathbf{x}$	$\mathbf{I}$	$\text{trace}(X)$	$\mathbf{I}$
$\mathbf{a}^T \mathbf{x}$	$\mathbf{a}$	$\text{trace}(A^T X)$	$A$
$\mathbf{x}^T \mathbf{A} \mathbf{x}$	$(A + A^T) \mathbf{x}$	$\text{trace}(X^T A X)$	$(A + A^T) X$
$\mathbf{x}^T \mathbf{x} = \ \mathbf{x}\ _2^2$	$2\mathbf{x}$	$\text{trace}(X^T X) = \ X\ _F^2$	$2X$
$\ \mathbf{A} \mathbf{x} - \mathbf{b}\ _2^2$	$2A^T (\mathbf{A} \mathbf{x} - \mathbf{b})$	$\ AX - B\ _F^2$	$2A^T (AX - B)$
$\mathbf{a}^T \mathbf{x}^T \mathbf{x} \mathbf{b}$	$2\mathbf{a}^T \mathbf{b} \mathbf{x}$	$\mathbf{a}^T X \mathbf{b}$	$\mathbf{a} \mathbf{b}^T$
$\mathbf{a}^T \mathbf{x} \mathbf{x}^T \mathbf{b}$	$(\mathbf{a} \mathbf{b}^T + \mathbf{b} \mathbf{a}^T) \mathbf{x}$	$\text{trace}(A^T X B)$	$AB^T$

## Numerical Method of Derivative

- Based on definition,  $f'(x) = \lim_{\varepsilon \rightarrow 0} \frac{f(x+\varepsilon) - f(x)}{\varepsilon}$ ,  $\varepsilon$  usually gets very small or  $= 10^{-6}$ ,  $f'(x) \approx \frac{f(x+\varepsilon) - f(x-\varepsilon)}{2\varepsilon}$
- Above  $f'(x)$  can be solve by either Analysis (Taylor Series) or Geometry

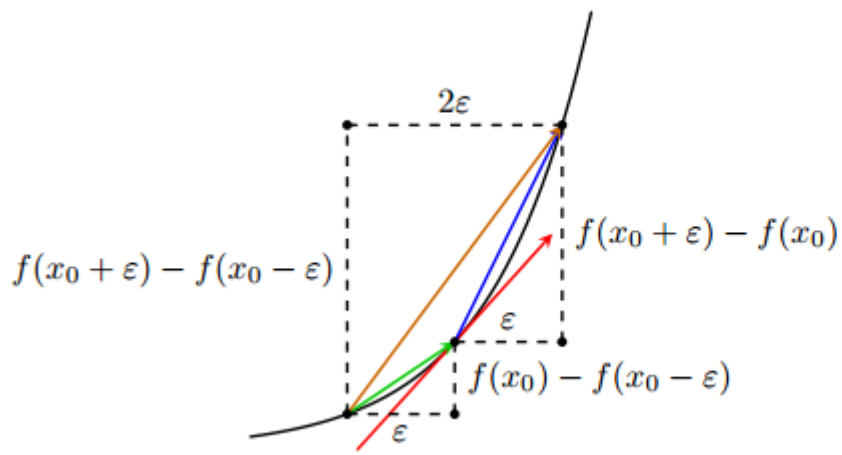
- Taylor Series 
$$\begin{cases} f(x+\varepsilon) \approx f(x) + f'(x)\varepsilon + \frac{f''(x)}{2}\varepsilon^2 + \frac{f^{(3)}(x)}{6}\varepsilon^3 + \dots \\ f(x-\varepsilon) \approx f(x) - f'(x)\varepsilon + \frac{f''(x)}{2}\varepsilon^2 - \frac{f^{(3)}(x)}{6}\varepsilon^3 + \dots \end{cases}$$

$$\rightarrow \begin{cases} \frac{f(x+\varepsilon) - f(x)}{\varepsilon} \approx f'(x) + \frac{f''(x)}{2}\varepsilon + \dots = f'(x) + O(\varepsilon) & (2.21) \\ \frac{f(x+\varepsilon) - f(x-\varepsilon)}{2\varepsilon} \approx f'(x) + \frac{f^{(3)}(x)}{6}\varepsilon^2 + \dots = f'(x) + O(\varepsilon^2) & (2.22) \end{cases}$$

(2.21) has the error of  $O(\varepsilon)$ , but 2.22 just has the error of  $O(\varepsilon^2)$ , and when  $\varepsilon \rightarrow 0$ ,  $O(\varepsilon^2) \ll O(\varepsilon)$ , so (2.22) is much better

- Geometry

**Hình 2.1:** Giải thích cách xấp xỉ đạo hàm bằng hình học.



Easily see that  $f(x_0 + \varepsilon) - f(x_0 - \varepsilon)$  is more similar with the real derivative than  $f(x_0) - f(x_0 - \varepsilon)$  and  $f(x_0 + \varepsilon) - f(x_0)$