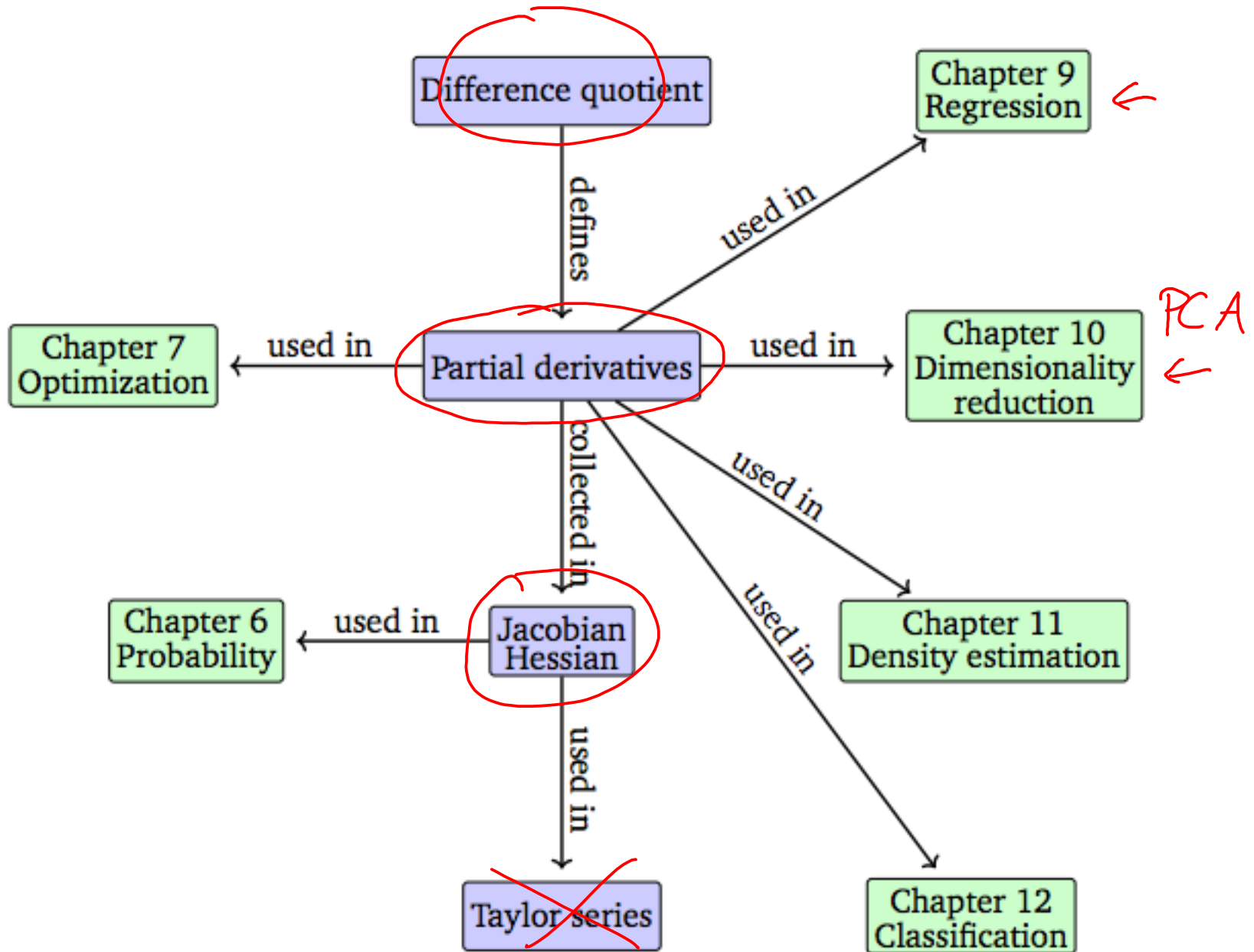


# Vector Calculus

Liang Zheng

Australian National University

[liang.zheng@anu.edu.au](mailto:liang.zheng@anu.edu.au)



# 5 Vector Calculus 向量微积分

- We discuss functions

$$f : \mathbb{R}^D \rightarrow \mathbb{R}$$

*domain*      *codomain*  
 $x \mapsto f(x)$       *norm*       $\|x\|_2$

where  $\mathbb{R}^D$  is the domain of  $f$ , and the function values  $f(x)$  are the image/codomain of  $f$ .

$$f(x) = c \cdot x \quad \mathbb{R}^n \rightarrow \mathbb{R}$$

- Example (dot product)
- Previously, we write dot product as

$$\begin{aligned} & \rightarrow f(\underline{x}) = x^T x, \quad x \in \underline{\mathbb{R}^2} \\ & \hookrightarrow f(x_1, x_2) = x_1^2 + x_2^2 \end{aligned}$$

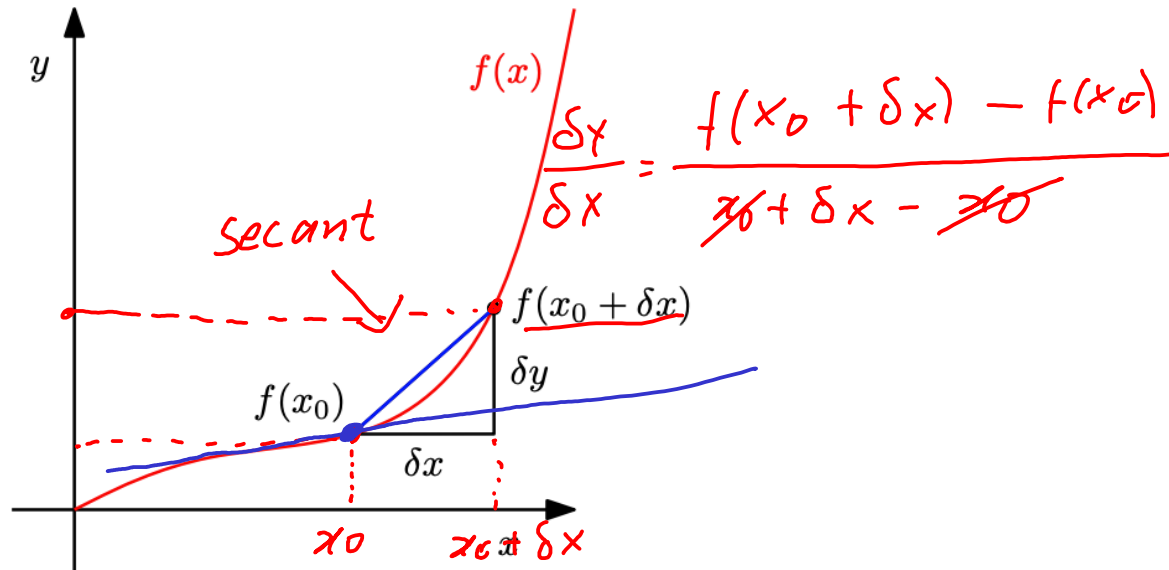
- In this chapter, we write it as

$$\begin{aligned} f : \mathbb{R}^2 &\rightarrow \mathbb{R} \\ x &\mapsto x_1^2 + x_2^2 \end{aligned}$$

# 5.1 Differentiation of Univariate Functions

微分

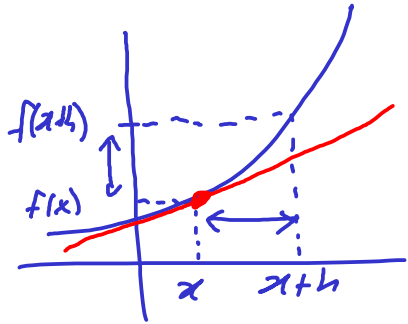
- Given  $y = f(x)$ , the **difference quotient** is defined as
$$\frac{\delta y}{\delta x} := \frac{f(x + \delta x) - f(x)}{\delta x}$$
- It computes the slope of the secant line through two points on the graph of  $f$ . In this figure, these are the points with  $x$ -coordinates  $x_0$  and  $x_0 + \delta x$ .
- In the limit for  $\delta x \rightarrow 0$ , we obtain the tangent of  $f$  at  $x$  (if  $f$  is differentiable). The tangent is then the derivative of  $f$  at  $x$ .



# 5.1 Differentiation of Univariate Functions

- For  $h > 0$ , the **derivative** of  $f$  at  $x$  is defined as the limit

$$\frac{df}{dx} := \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$



- The derivative of  $f$  points in the direction of steepest ascent of  $f$ .

## ~~Example - Derivative of a Polynomial~~

- ~~Compute the derivative of  $f(x) = x^n$ ,  $n \in \mathbb{N}$ . (From our high school knowledge, the derivative is  $nx^{n-1}$ .)~~

~~$$\frac{df}{dx} := \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$~~

~~$$= \lim_{h \rightarrow 0} \frac{(x+h)^n - x^n}{h}$$~~

~~$$= \lim_{h \rightarrow 0} \frac{\sum_{i=0}^n \binom{n}{i} x^{n-i} h^i - x^n}{h}$$~~

~~we see that  $x^n = \binom{n}{0} x^{n-0} h^0$ . By starting the sum at 1, the  $x^n$  cancels.~~

$$\frac{df}{dx} := \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

Example. Prove  $\frac{d}{dx}(x^2) = 2x$

$$\begin{aligned} \frac{d}{dx}(x^2) &= \lim_{h \rightarrow 0} \frac{(x+h)^2 - x^2}{h} \\ &= \lim_{h \rightarrow 0} \frac{\cancel{x^2} + 2xh + \cancel{h^2}}{h} \end{aligned}$$

$$= \lim_{h \rightarrow 0} \frac{2xh}{h} + \frac{\cancel{h^2}}{h} h$$

$$= \lim_{h \rightarrow 0} 2x + h$$

$$= 2x + \lim_{h \rightarrow 0} h$$

$$= 2x$$

$$\frac{df}{dx} := \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

Power Law

$$\frac{d}{dx} (x^6) = 6x^5$$

Example. Prove  $\frac{d}{dx} (x^n) = nx^{n-1}$

Binomial Theorem

$$(x+h)^n = \sum_{i=0}^n \binom{n}{i} x^{n-i} h^i$$

$$\frac{d}{dx} (x^n) = \lim_{h \rightarrow 0} \frac{(x+h)^n - x^n}{h}$$

$$= \lim_{h \rightarrow 0} \frac{\sum_{i=0}^n \binom{n}{i} x^{n-i} h^i - x^n}{h}$$

$$\cancel{x^n} x^n$$

$$= \lim_{h \rightarrow 0} \frac{\cancel{x^n} + \left( \sum_{i=1}^n \binom{n}{i} x^{n-i} h^i \right) - \cancel{x^n}}{h}$$

$$= \lim_{h \rightarrow 0} \binom{n}{1} x^{n-1} + \sum_{i=2}^n \binom{n}{i} x^{n-i} \cancel{h^i} \cdot \frac{h^i}{h} = h^{i-1}$$

$$= n \cdot x^{n-1} + \lim_{h \rightarrow 0} \sum_{i=2}^n \binom{n}{i} x^{n-i} h^{i-1}$$

$$= n \cdot x^{n-1}$$

Optional Exercise. Prove  $\frac{d}{dx} e^x = e^x$

You will need the Taylor Expansion of  $\exp(x)$  to do this:  $e^x = \sum_{n=0}^{\infty} \frac{x^n}{n!}$

and that  $e^{x+y} = e^x \cdot e^y$

Why care about derivatives?

$$(y_1, x_1), \dots, (y_n, x_n)$$

Model  $\hat{y} = f(x | \theta)$

Annotations:  
 $\hat{y}$ : guess  
 $x$ : data  
 $f$ : model  
 $\theta$ : param.

Loss

$$\mathcal{L}(X; \theta) = \sum_i \|y_i - \hat{y}_i\|_2^2$$

$$f(x) = x^5 - x^2 + x$$

$$f'(x) = 5x^4 - 2x + 1 = 0 \quad ?$$

$\frac{\partial \mathcal{L}}{\partial \theta} = 0$  (easy)

Solve for  $\theta$  (hard)

Numerically. 1. Guess  $\theta_0$

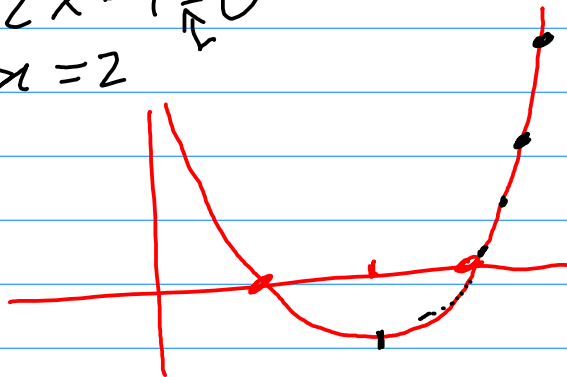
Learning Rate (0.1 to 0.01)  
 $\eta$  (Hyperparameter)

2.  $\theta_{t+1} = \theta_t - \eta \cdot \frac{\partial \mathcal{L}}{\partial \theta}$

Minimise  $f(x) = x^2 - 4x + 3$

$$f'(x) = 2x - 4 = 0$$

$$\Rightarrow x = 2$$



guess  $x_0 = 4$   $\eta = 0.1$

$$x_1 = x_0 - \eta \cdot f'(x_0)$$

$$= 4 - 0.1 \cdot (2 \cdot 4 - 4)$$

$$= 3.6$$

$$x_2 = x_1 - \eta \cdot f'(x_1)$$

$$= 3.6 - 0.1 (2 \times 3.6 - 4)$$

$$= 3.28$$

$$x_3 = 3.024, \dots$$



# Estimating Intractable Derivatives

$$\frac{\partial \mathcal{L}}{\partial \theta} = ?$$

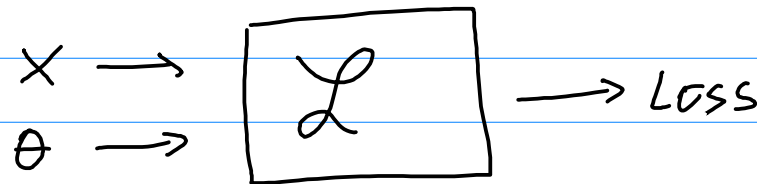
$\mathcal{L} = \text{loss}$

$$f'(x) \approx \frac{f(x+\epsilon) - f(x)}{\epsilon} \quad \text{for small } \epsilon > 0$$

NOT OFTEN

eg estimate derivative of  $f(x) = x^2$  at  $x = 2$ .

True Answer:  $f'(x) = 2x$        $f'(2) = 4$



Estimate: 
$$\frac{f(2+\epsilon) - f(2)}{\epsilon} = \frac{(2+\epsilon)^2 - 2^2}{\epsilon}$$

$\epsilon = 0.1$       
$$\frac{2.1^2 - 2^2}{0.1} = 4.1$$

$\epsilon = 0.01$       
$$\frac{2.01^2 - 2^2}{0.01} = 4.01$$

# 5.1.2 Differentiation Rules

$$\frac{d}{dx} (f(x)g(x)) = \frac{df}{dx} \cdot g + f \cdot \frac{dy}{dx} \quad f'(x) = \frac{d(f(x))}{dx}$$

- Product rule

$$\underline{(f(x)g(x))'} = \underline{f'(x)} \underline{g(x)} + \underline{f(x)} \underline{g'(x)}$$

- Quotient rule:

$$\left(\frac{f(x)}{g(x)}\right)' = \frac{f'(x)g(x) - f(x)g'(x)}{(g(x))^2}$$

- Sum rule:

Differentiation is linear

$$(f(x) + g(x))' = f'(x) + g'(x)$$

$$(x^2 + 5x)' = 2x + 5$$

- Chain rule:

$$\underline{(g(f(x)))' = (g \circ f)'(x) = g'(f(x))f'(x)}$$

Here,  $g \circ f$  denotes function composition  $g(f(x))$

$$\left(\frac{f(x)}{g(x)}\right)' = \left(f(x) \cdot \frac{1}{g(x)}\right)'$$

Prove Sum Rule

$$(f+g)' = f' + g'$$

$$\frac{df}{dx} := \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

$$\frac{d}{dx} (f(x) + g(x)) = \lim_{h \rightarrow 0} \frac{(f(x+h) + g(x+h)) - (f(x) + g(x))}{h}$$

$$= \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} + \frac{g(x+h) - g(x)}{h}$$

$$= \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} + \lim_{h \rightarrow 0} \frac{g(x+h) - g(x)}{h}$$

$$= f'(x) + g'(x)$$

Prove Product Rule  $(fg)' = fg' + f'g$   $\frac{df}{dx} := \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$

$$\frac{d}{dx} (f(x)g(x)) = \lim_{h \rightarrow 0} \frac{f(x+h)g(x+h) - f(x)g(x)}{h}$$

$$= \lim_{h \rightarrow 0} \frac{\boxed{f(x+h)g(x+h)} + \boxed{f(x+h)g(x)} - \boxed{f(x+h)g(x)} - \boxed{f(x)g(x)}}{h}$$

$$= \lim_{h \rightarrow 0} f(x+h) \frac{g(x+h) - g(x)}{h} + g(x) \frac{f(x+h) - f(x)}{h}$$

$$= \lim_{h \rightarrow 0} f(x+h) \frac{g(x+h) - g(x)}{h} + \lim_{h \rightarrow 0} g(x) \frac{f(x+h) - f(x)}{h}$$

$$\left( \lim_{h \rightarrow 0} f(x+h) \right) \left( \lim_{h \rightarrow 0} \frac{g(x+h) - g(x)}{h} \right) + g(x) \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

$$f(x) \quad g'(x) \quad + g(x) \cdot f'(x)$$

Exercise:

$$= f(x) \cdot g'(x) + g(x) \cdot f'(x) \quad \text{Prove Chain Rule.}$$

$$\text{Exercise: Prove } \frac{d}{dx} (cf(x)) = c \frac{d}{dx} f(x)$$

## Examples of Derivative Rules

$$\begin{aligned}(fg)' &= f'g + fg' \\ (f+g)' &= f' + g' \\ (f(g(x)))' &= f'(g(x)) \cdot g'(x)\end{aligned}$$

Derive  $\sin(2x+1)$

$$\sin(2x+1) = f(g(x))$$

$$f(x) = \sin x \quad g(x) = \underline{2x+1}$$

$$\begin{aligned}(f(g(x)))' &= \underline{f'(g(x)) \cdot g'(x)}\end{aligned}$$

$$f'(x) = \cos x$$

$$f'(g(x)) = \cos(2x+1)$$

$$g'(x) = 2$$

$$= 2 \cos(2x+1)$$

Derive  $(\sin(x^4))^5$

$$h(f(g(x)))$$

$$h(x) = x^5$$

$$f(x) = \sin x$$

$$g(x) = x^4$$

$$(h(f(g(x))))' = h'(f(g(x))) \cdot (f(g(x)))'$$

$$= h'(f(g(x))) \cdot f'(g(x)) \cdot g'(x)$$

$$h'(x) = 5x^4 \quad f'(x) = \cos x \quad g'(x) = 4x^3$$

$$= 5 (\sin(x^4))^4 \cdot \cos(x^4) \cdot 4x^3$$

$$= 20 (\sin(x^4))^4 \cdot \cos(x^4)$$

## Product Rule

$$(fg)' = f'g + fg'$$
$$(f+g)' = f' + g'$$
$$(f(g(x)))' = f'(g(x)) \cdot g'(x)$$

$$\frac{d}{dx} (x^3 \cdot \sin x)$$

$$= \left( \frac{d}{dx} x^3 \right) \sin x + x^3 \left( \frac{d}{dx} \sin x \right)$$

$$= 3x^2 \cdot \sin x + x^3 \cos x.$$

An alternative view of chain rule

$$\frac{dy}{dx} = \frac{dy}{du} \cdot \frac{du}{dx}$$

Derive  $y = (x^2 + 5x + 6)^{10} = u^{10}$

let  $u = x^2 + 5x + 6$

$$\frac{dy}{dx} = \frac{dy}{du} \cdot \frac{du}{dx}$$

$$= 10u^9 \cdot (2x + 5)$$

$$= 10(x^2 + 5x + 6)^9 (2x + 5)$$

# 5.2 Partial Differentiation and Gradients

- Instead of considering  $x \in \mathbb{R}$ , we consider  $x \in \mathbb{R}^n$ , e.g.,  $f(x) = f(x_1, x_2) \in \mathbb{R}^2$
- The generalization of the derivative to functions of several variables is the gradient.   
 $\nabla_x f$  "nabla" "del"  $\text{grad } f$    
 $\left( \frac{df}{dx} \right)$
- We find the gradient of the function  $f$  with respect to  $x$  by
  - varying one variable at a time and keeping the others constant.
  - The gradient is the collection of these partial derivatives.
- For a function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $x \mapsto f(x)$ ,  $x \in \mathbb{R}^n$  of  $n$  variables  $x_1, \dots, x_n$ , we define the partial derivatives as

$$\frac{\partial f}{\partial x_1} := \lim_{h \rightarrow 0} \frac{f(x_1 + h, x_2, \dots, x_n) - f(x)}{h}$$

$$\vdots$$

$$\frac{\partial f}{\partial x_n} := \lim_{h \rightarrow 0} \frac{f(x_1, \dots, x_{n-1}, x_n + h) - f(x)}{h}$$

and collect them in the row vector

$$\nabla_x f = \text{grad } f = \frac{df}{dx} = \left[ \frac{\partial f(x)}{\partial x_1} \quad \frac{\partial f(x)}{\partial x_n} \quad \dots \quad \frac{\partial f(x)}{\partial x_n} \right] \in \mathbb{R}^{1 \times n}$$



## 5.2 Partial Differentiation and Gradients

$$f: \mathbb{R}^n \rightarrow \mathbb{R} \quad \mathbb{R}^{1 \times n} \quad \downarrow \quad \frac{\partial f(x)}{\partial x_1}: \mathbb{R}^n \rightarrow \mathbb{R} \quad \boxed{\mathbb{R}} \quad \text{codomain} \times \text{domain}$$

- $\nabla_x f = \text{grad} f = \frac{df}{dx} = \left[ \frac{\partial f(x)}{\partial x_1} \quad \frac{\partial f(x)}{\partial x_n} \quad \dots \quad \frac{\partial f(x)}{\partial x_n} \right] \in \boxed{\mathbb{R}^{1 \times n}} \quad (\mathbb{R}^n \rightarrow \mathbb{R})^{1 \times n}$
- $n$  is the number of variables and  $1$  is the dimension of the image/range/codomain of  $f$
- The row vector  $\nabla_x f \in \mathbb{R}^{1 \times n}$  is called the **gradient** of  $f$  or the **Jacobian**.

- Example - Partial Derivatives Using the Chain Rule

- For  $f(x, y) = (x + 2y^3)^2$ , we obtain the partial derivatives

$$f: \mathbb{R}^2 \rightarrow \mathbb{R} \quad \frac{\partial f(x, y)}{\partial x} = 2(x + 2y^3) \frac{\partial}{\partial x} (x + 2y^3) = 2(x + 2y^3)$$

$$\frac{\partial f(x, y)}{\partial y} = 2(x + 2y^3) \frac{\partial}{\partial y} (x + 2y^3) = 12(x + 2y^3)y^2$$

$x = (x, y)$

$$\nabla_x f = \left[ \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right] = \left[ 2(x + 2y^3), 12(x + 2y^3)y^2 \right]^T \in \mathbb{R}^{1 \times 2}$$

## 5.2 Partial Differentiation and Gradients

$$f: \mathbb{R}^2 \rightarrow \mathbb{R}$$

↓

- For  $f(x_1, x_2) = \underline{x_1^2 x_2} + \underline{x_1 x_2^3} \in \mathbb{R}$ , the partial derivatives (i.e., the derivatives of  $f$  with respect to  $x_1$  and  $x_2$ ) are

$$\frac{\partial f(x_1, x_2)}{\partial x_1} = 2x_1 x_2 + x_2^3$$

$$\frac{\partial f(x_1, x_2)}{\partial x_2} = x_1^2 + 3x_1 x_2^2$$

- and the gradient is then

$$\frac{df}{d\mathbf{x}} = \left[ \frac{\partial f(x_1, x_2)}{\partial x_1} \quad \frac{\partial f(x_1, x_2)}{\partial x_2} \right] = [ 2x_1 x_2 + x_2^3 \quad x_1^2 + 3x_1 x_2^2 ] \in \mathbb{R}^{1 \times 2}$$

## 5.2.1 Basic Rules of Partial Differentiation

$$\boxed{f: \mathbb{R}^n \rightarrow \mathbb{R}}$$

$$g(f(x))$$

- Product rule:

$$\frac{\partial}{\partial x}(f(x)g(x)) = \frac{\partial f}{\partial x}g(x) + f(x)\frac{\partial g}{\partial x}$$

- Sum rule:

$$\frac{\partial}{\partial x}(f(x) + g(x)) = \frac{\partial f}{\partial x} + \frac{\partial g}{\partial x}$$

$$g \circ f: \mathbb{R}^n \rightarrow \mathbb{R}$$

- Chain rule:

$$\frac{\partial}{\partial x}(g \circ f)(x) = \frac{\partial}{\partial x}(g(f(x))) = \frac{\partial g}{\partial f} \frac{\partial f}{\partial x}$$

$$\boxed{g: \mathbb{R} \rightarrow \mathbb{R}}$$

$$\frac{\partial g}{\partial f}: \underline{\mathbb{R}} \rightarrow \mathbb{R} \quad (\mathbb{R}?)$$

$$\frac{\partial f}{\partial x}: \underline{\mathbb{R}^n} \rightarrow \mathbb{R}^{1 \times n} \quad (\mathbb{R}^{1 \times n})$$

## 5.2.2 Chain Rule

$$f: \mathbb{R}^2 \rightarrow \mathbb{R} \quad \text{codomain } \mathbb{R} \times \text{domain } \mathbb{R}^2 \quad \frac{\partial f}{\partial \mathbf{x}} \in \mathbb{R}^{1 \times 2} \quad \frac{\partial \mathbf{x}}{\partial t} \in \mathbb{R}^{2 \times 1}$$

$$\mathbf{x}(t) = \{x_1(t), x_2(t)\}$$

$$\mathbf{x}: \mathbb{R} \rightarrow \mathbb{R}^2$$

- Consider a function  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$  of two variables  $x_1$  and  $x_2$ .
- $x_1(t)$  and  $x_2(t)$  are themselves functions of  $t$ .
- To compute the gradient of  $f$  with respect to  $t$ , we apply the chain rule:

$$\boxed{\frac{df}{dt}} = \frac{\partial f}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial t} = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} \end{bmatrix} \begin{bmatrix} \frac{\partial x_1(t)}{\partial t} \\ \frac{\partial x_2(t)}{\partial t} \end{bmatrix} = \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial t} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial t}$$

Where  $d$  denotes the gradient and  $\partial$  partial derivatives.

- Example
- Consider  $f(x_1, x_2) = x_1^2 + 2x_2$ , where  $x_1 = \sin t$  and  $x_2 = \cos t$ , then

$$\begin{aligned} \frac{df}{dt} &= \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial t} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial t} \\ &= 2 \sin t \frac{\partial \sin t}{\partial t} + 2 \frac{\partial \cos t}{\partial t} \\ &= 2 \sin t \cos t - 2 \sin t = 2 \sin t (\cos t - 1) \end{aligned}$$

The above is the corresponding derivative of  $f$  with respect to  $t$

$$f: \mathbb{R}^2 \rightarrow \mathbb{R}$$

$$x_1: \mathbb{R} \rightarrow \mathbb{R} \quad x_2: \mathbb{R} \rightarrow \mathbb{R}$$

$$f(x_1, x_2) = x_1^2 x_2 + 3x_1 + 6x_2$$

$$x_1(t) = \sin t$$

$$x_2(t) = e^{2t}$$

$$X(t) =$$

$$\frac{d}{dt} = \frac{\partial}{\partial x} \frac{\partial}{\partial t}$$

$$X: \mathbb{R} \rightarrow \mathbb{R}^2$$

dom = 1      codom = 2

$$\frac{\partial f}{\partial x} = \left[ \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2} \right] = [2x_1 x_2 + 3, x_1^2 + 6] \in \mathbb{R}^{1 \times 2}$$

$$\frac{\partial X}{\partial t} = \begin{bmatrix} \frac{\partial x_1}{\partial t} \\ \frac{\partial x_2}{\partial t} \end{bmatrix} = \begin{bmatrix} \cos t \\ 2e^{2t} \end{bmatrix} \in \mathbb{R}^{2 \times 1}$$

$$\frac{d}{dx}(e^{f(x)}) = f'(x) e^{f(x)}$$

$$\mathbb{R}^{\text{codom} \times \text{dom}} = \mathbb{R}^{2 \times 1} \approx \mathbb{R}^2$$

$$\frac{\partial f}{\partial t} = \frac{\partial f}{\partial x} \frac{\partial X}{\partial t} = (2x_1 x_2 + 3) \cos t + (x_1^2 + 6) 2e^{2t}$$

Alternatively:

$$f(t) = \underbrace{(\sin t)^2}_{\text{product rule}} e^{2t} + 3 \sin t + 6 e^{2t} : \mathbb{R} \rightarrow \mathbb{R}$$

## 5.2.2 Chain Rule

- If  $f(x_1, x_2)$  is a function of  $(x_1)$  and  $(x_2)$ , where  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $x_1(s, t)$  and  $x_2(s, t)$  are themselves functions of two variables  $s$  and  $t$ , the chain rule yields the partial derivatives

$$\begin{aligned} \frac{\partial f}{\partial s} &= \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial s} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial s} = \frac{\partial f}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial s} \\ \frac{\partial f}{\partial t} &= \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial t} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial t} = \frac{\partial f}{\partial \mathbf{x}} \cdot \frac{\partial \mathbf{x}}{\partial t} \end{aligned}$$

- The gradient can be obtained by the matrix multiplication

$$\mathbf{x}: \mathbb{R}^2 \rightarrow \mathbb{R}^2$$

$$\begin{aligned} \nabla_{(s,t)} f &= \frac{df}{d(s,t)} = \frac{\partial f}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial (s,t)} = \underbrace{\begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} \end{bmatrix}}_{= \frac{\partial f}{\partial \mathbf{x}}} \underbrace{\begin{bmatrix} \frac{\partial x_1}{\partial s} & \frac{\partial x_1}{\partial t} \\ \frac{\partial x_2}{\partial s} & \frac{\partial x_2}{\partial t} \end{bmatrix}}_{= \frac{\partial \mathbf{x}}{\partial (s,t)}} \\ \frac{\partial \mathbf{x}}{\partial (s,t)} &\in \mathbb{R}^{2 \times 2} \end{aligned}$$

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$

In general, if  $f(x_1, \dots, x_n)$  and each  $x_i$  is a function  ~~$x_i(u_1, \dots, u_k)$~~  then  $x_i: \mathbb{R}^k \rightarrow \mathbb{R}$   
 $x_i(u_1, u_2, \dots, u_k)$

$$\frac{\partial f}{\partial u_i} = \sum_{j=1}^n \frac{\partial f}{\partial x_j} \boxed{\frac{\partial x_j}{\partial u_i}} \quad \leftarrow$$

$$= \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial u_i} + \dots + \frac{\partial f}{\partial x_n} \frac{\partial x_n}{\partial u_i}$$

$$x_1(u, v) = u + v$$

~~$$x_2(u) = 2u$$~~

$$x_2(u, v) = 2u$$

$$\frac{\partial f}{\partial t} = \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial t} + \cancel{\frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial t}}$$

$$f(x_1, x_2)$$

$$\boxed{\begin{aligned} x_1(s, t) &= s \cdot t \\ x_2(s, p, q) &= sp + 2 \end{aligned}}$$

$$x_1(s, t, \overset{d}{p}, q) = st$$

$$\boxed{x_2(s, t, p, q) = sp + q}$$

$$f: \mathbb{R}^3 \rightarrow \mathbb{R} \quad \downarrow$$

$$x, y, z: \mathbb{R}^3 \rightarrow \mathbb{R}$$

Example  $f(x, y, z) = 3xy + yz^2$

$$x(p, q, r) = pqr$$

$$y(p, q, r) = 2p^2 + q + 5r$$

$$z(p, q, r) = p \cos(q + 2r)$$



$$f = \underbrace{3pqr}_x (\underbrace{2p^2 + q + 5r}_y) + \dots$$

compute  $\frac{\partial f}{\partial p} = \frac{\partial f}{\partial(x, y, z)} \cdot \frac{\partial(x, y, z)}{\partial p}$

$$= \frac{\partial f}{\partial x} \cdot \frac{\partial x}{\partial p} + \frac{\partial f}{\partial y} \frac{\partial y}{\partial p} + \frac{\partial f}{\partial z} \cdot \frac{\partial z}{\partial p}$$

$$= (3y)(qr) + (3x + z^2)(4p) + (2yz)(\cos(q + 2r))$$

Exercise:  $\frac{\partial f}{\partial q} \quad \& \quad \frac{\partial f}{\partial r}$



$$f(\mathbf{X}) : \mathbb{R}^{n \times m} \rightarrow \mathbb{R} \quad F(\mathbf{x}) : \mathbb{R} \rightarrow \mathbb{R}^{n \times m}$$

## 5.3 Gradients of Vector-Valued Functions

$$f(x) : \mathbb{R} \rightarrow \mathbb{R} \quad f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R} \quad \mathbf{f}(x) : \mathbb{R} \rightarrow \mathbb{R}^n \quad \mathbf{f}(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

- We discussed partial derivatives and gradients of function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$
- We will generalize the concept of the gradient to vector-valued functions (vector fields)  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ , where  $n \geq 1$  and  $m > 1$ .

- For a function  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  and a vector  $\mathbf{x} = [x_1, \dots, x_n]^T \in \mathbb{R}^n$ , the corresponding vector of function values is given as

$$\mathbf{f}(\mathbf{x}) = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_m(\mathbf{x}) \end{bmatrix} \in \mathbb{R}^m \quad \left( \mathbb{R}^n \rightarrow \mathbb{R} \right)^m$$

$\left. \begin{array}{l} \mathbb{R}^n \rightarrow \mathbb{R} \\ \mathbb{R}^n \rightarrow \mathbb{R} \\ \vdots \\ \mathbb{R}^n \rightarrow \mathbb{R} \end{array} \right\} m$

- Writing the vector-valued function in this way allows us to view a vector valued function  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  as a vector of functions  $[f_1, \dots, f_m]^T$ ,  $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$  that map onto  $\mathbb{R}$ .
- The differentiation rules for every  $f_i$  are exactly the ones we discussed before.

# 5.3 Gradients of Vector-Valued Functions

- The partial derivative of a vector-valued function  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$  with respect to  $x_i \in \mathbb{R}, i = 1, \dots, n$ , is given as the vector

$$\left( \frac{\partial f}{\partial x_i} \right) = \begin{bmatrix} \frac{\partial f_1}{\partial x_i} \\ \vdots \\ \frac{\partial f_m}{\partial x_i} \end{bmatrix} = \begin{bmatrix} \lim_{h \rightarrow 0} \frac{f_1(x_1, \dots, x_{i-1}, x_i + h, x_{i+1}, \dots, x_n) - f_1(x)}{h} \\ \vdots \\ \lim_{h \rightarrow 0} \frac{f_m(x_1, \dots, x_{i-1}, x_i + h, x_{i+1}, \dots, x_n) - f_m(x)}{h} \end{bmatrix} \in \mathbb{R}^m$$

- In above, every partial derivative  $\frac{\partial f}{\partial x_i}$  is a column vector
- Recall that the gradient of  $f$  with respect to a vector is the row vector of the partial derivatives
- Therefore, we obtain the gradient of  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$  with respect to  $x \in \mathbb{R}^n$ , by collecting these partial derivatives:

$$\frac{df(x)}{dx} = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} & \dots & \frac{\partial f(x)}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1(x)}{\partial x_1} & \dots & \frac{\partial f_1(x)}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_m(x)}{\partial x_1} & \dots & \frac{\partial f_m(x)}{\partial x_n} \end{bmatrix}$$

Jacobian

$J_{ij} = \frac{\partial f_i(x)}{\partial x_j}$

$\in \mathbb{R}^{m \times n}$

codomain  $\rightarrow$  domain

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$

$\mathbb{R}$

$\mathbb{R}$

$$\frac{\partial f}{\partial x_i} = \lim_{h \rightarrow 0} \frac{f(x_1, x_2, \dots, x_{i-1}, \underline{x_i + h}, x_{i+1}, \dots, x_n) - f(x_1, \dots, x_n)}{h}$$

$(h) \mathbb{R}$

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$

Prove  $\nabla_x g(x) = c^T$

$$f(x) = x^T x = \|x\|_2^2 \quad g: \mathbb{R}^n \rightarrow \mathbb{R} \quad g(x) = c^T x$$

$$\nabla_x f(x) = 2x^T$$

Kronecker

~~Kronecker~~  
Delta

$$\delta_{ij} = \begin{cases} 1 & i=j \\ 0 & i \neq j \end{cases}$$

$$\frac{\partial x_i}{\partial x_j} = \begin{cases} 1 & i=j \\ 0 & i \neq j \end{cases} = \delta_{ij}$$

$$\nabla_x g(x) = \left[ \frac{\partial (c^T x)}{\partial x_1}, \dots, \frac{\partial (c^T x)}{\partial x_n} \right]^T \in \mathbb{R}^{1 \times n}$$

$$c^T x = \sum_i c_i x_i$$

$$\frac{\partial}{\partial x_j} \left( \sum_{i=1}^n c_i x_i \right) = \sum_{i=1}^n c_i \frac{\partial x_i}{\partial x_j}$$

$$= \sum_{i=1}^n c_i \delta_{ij} = c_j$$

$$\nabla_x g(x) = [c_1, c_2, \dots, c_n]^T = c^T$$

$$\nabla_x (c^T x) = c^T$$

$$\frac{d}{dx} (cx) = c$$

$$\frac{\partial f}{\partial x_1} \quad \frac{\partial^2 f}{\partial x_1^2}$$

$$= \frac{\partial}{\partial x_1} \left( \frac{\partial f}{\partial x_1} \right)$$

Second order.

$$\frac{\partial^2 f}{\partial x_1 \partial x_2} = \frac{\partial}{\partial x_1} \left( \frac{\partial f}{\partial x_2} \right)$$

$$= \|x\|_2^2$$

Prove  $\nabla_x (x^T x) = 2x^T$

$$x \in \mathbb{R}^n$$

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\nabla_x (x^T x) = \left[ \frac{\partial x^T x}{\partial x_1}, \dots, \frac{\partial x^T x}{\partial x_n} \right] \in \mathbb{R}^{1 \times n}$$

$$\frac{\partial}{\partial x_j} (x^T x) = \frac{\partial}{\partial x_j} \left( \sum_i x_i^2 \right) = \sum_i \frac{\partial x_i^2}{\partial x_j} \begin{cases} 2x_j & j=i \\ 0 & j \neq i \end{cases}$$

$$= \sum_i 2x_i \delta_{ij} = 2x_j$$

$$\nabla_x \|x\|_2^2$$

$$\nabla_x (x^T x) = [2x_1, \dots, 2x_n] = (2x^T)$$

$$\nabla_x (c^T x) = c^T$$

## 5.3 Gradients of Vector-Valued Functions

- The collection of all first-order partial derivatives of a vector-valued function  $\mathbf{f}: \mathbb{R}^n \rightarrow \mathbb{R}^m$  is called the **Jacobian**. The Jacobian  $\mathbf{J}$  is an  $m \times n$  matrix, which we define and arrange as follows:

$$\begin{aligned}\mathbf{J} = \nabla_{\mathbf{x}} \mathbf{f} &= \frac{d\mathbf{f}(\mathbf{x})}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial \mathbf{f}(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial \mathbf{f}(\mathbf{x})}{\partial x_n} \end{bmatrix} \\ &= \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial f_1(\mathbf{x})}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_m(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial f_m(\mathbf{x})}{\partial x_n} \end{bmatrix} \\ \mathbf{x} &= \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \quad \mathbf{J}(i,j) = \frac{\partial f_i}{\partial x_j}\end{aligned}$$

- The elements of  $\mathbf{f}$  define the rows and the elements of  $\mathbf{x}$  define the columns of the corresponding Jacobian
- Special case: for a function  $\mathbf{f}: \mathbb{R}^n \rightarrow \mathbb{R}^1$  which maps a vector  $\mathbf{x} \in \mathbb{R}^n$  onto a scalar, i.e.,  $m = 1$ , the Jacobian is a row vector of dimension  $1 \times n$ .

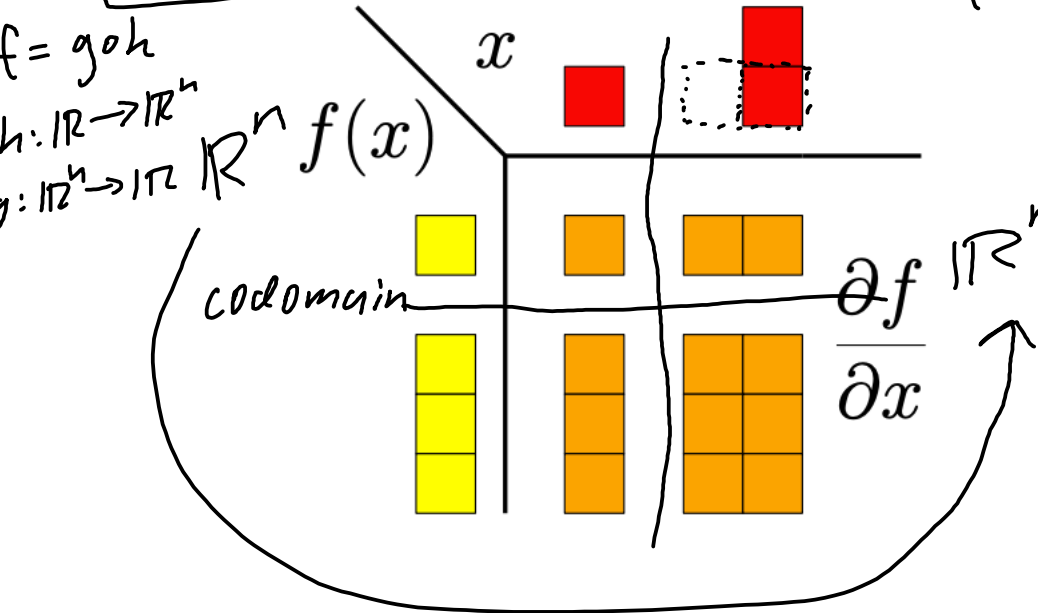
# 5.3 Gradients of Vector-Valued Functions

- $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ 
 $\nabla_x f: \mathbb{R}^n \rightarrow \mathbb{R}^{m \times n}$ 

codomain

domain
- If  $f: \mathbb{R} \rightarrow \mathbb{R}$ , the gradient is a scalar
  $f: \mathbb{R}^{(n_1 \times \dots \times n_k)} \rightarrow \mathbb{R}^{(m_1 \times \dots \times m_s)}$
- If  $f: \mathbb{R}^D \rightarrow \mathbb{R}$ , the gradient is a  $1 \times D$  row vector
  $\nabla_x f: \mathbb{R}^{(n_1 \times \dots \times n_k)} \rightarrow \mathbb{R}^{(m_1 \times \dots \times m_s) \times (n_1 \times \dots \times n_k)}$
- If  $f: \mathbb{R} \rightarrow \mathbb{R}^E$ , the gradient is a  $E \times 1$  column vector
  $\nabla_x f: \mathbb{R}^{(n_1 \times \dots \times n_k)} \rightarrow \mathbb{R}^{(m_1 \times \dots \times m_s) \times (n_1 \times \dots \times n_k)}$
- If  $f: \mathbb{R}^D \rightarrow \mathbb{R}^E$ , the gradient is an  $E \times D$  matrix

$\frac{df}{dt} = \frac{df}{dx} \cdot \frac{dx}{dt}$   
 $f = g \circ h$   
 $h: \mathbb{R} \rightarrow \mathbb{R}^n$   
 $g: \mathbb{R}^n \rightarrow \mathbb{R}$



$f: \mathbb{R}^{a \times b} \rightarrow \mathbb{R}^{(c)}$   
 $\nabla_x f(x): \mathbb{R}^{a \times b} \rightarrow \mathbb{R}^{c \times (a \times b)}$   
 $\nabla_x f(x) \in \mathbb{R}^{c \times (a \times b)}$

$\mathbb{R}^n = \mathbb{R}^{n \times 1}$

# Example - Gradient of a Vector-Valued Function

$$\vec{f}: \mathbb{R}^N \rightarrow \mathbb{R}^M$$

- We are given  $\underline{f(x) = Ax}$ ,  $f(x) \in \mathbb{R}^M$ ,  $A \in \mathbb{R}^{M \times N}$ ,  $\underline{x \in \mathbb{R}^N}$ .
- To compute the gradient  $df/dx$  we first determine the dimension of  $df/dx$ : Since  $f: \mathbb{R}^N \rightarrow \mathbb{R}^M$ , it follows that  $df/dx \in \mathbb{R}^{M \times N}$ .
- Then, we determine the partial derivatives of  $f$  with respect to every  $x_j$ :

$$\underline{f_i(x) = \sum_{j=1}^N A_{ij} x_j} \Rightarrow \frac{\partial f_i}{\partial x_j} = A_{ij} \quad i \downarrow \left[ \begin{matrix} \vdots \\ \vdots \\ \vdots \end{matrix} \right] \left[ \begin{matrix} \vdots \\ \vdots \\ \vdots \end{matrix} \right] = \left[ \begin{matrix} \vdots \\ \vdots \\ \vdots \end{matrix} \right]$$

$A \quad x \quad f$

$$\left[ \frac{\partial f_i}{\partial x_j} \right] = \frac{\partial}{\partial x_j} \left( \sum_k A_{ik} x_k \right) = \sum_k A_{ik} \frac{\partial x_k}{\partial x_j} \delta_{jk} = \boxed{A_{ij}}$$

- We collect the partial derivatives in the Jacobian and obtain the gradient

$$\left[ \frac{\partial f}{\partial x} \right] = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_N} \\ \vdots & & \vdots \\ \frac{\partial f_M}{\partial x_1} & \cdots & \frac{\partial f_M}{\partial x_N} \end{bmatrix} = \begin{bmatrix} A_{11} & \cdots & A_{1N} \\ \vdots & & \vdots \\ A_{M1} & \cdots & A_{MN} \end{bmatrix} = \boxed{A} \in \mathbb{R}^{M \times N}$$

$\frac{d}{dx}(cx) = c \quad \frac{d}{d\vec{x}} A\vec{x} = A$



# Example - Chain Rule

$$(e^{f(x)})' = f'(x)e^{f(x)} \quad \frac{\partial x_1}{\partial t} = \cos t - t \sin t$$

$$\frac{\partial x_2}{\partial t} = \sin t + t \cos t$$

- Consider the function  $\underline{h: \mathbb{R} \rightarrow \mathbb{R}}$ ,  $h(t) = (f \circ g)(t)$  with

$$f: \mathbb{R}^2 \rightarrow \mathbb{R}$$

$$g: \mathbb{R} \rightarrow \mathbb{R}^2$$

$$f(x) = \exp(x_1 x_2^2)$$

$$\underline{x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = g(t) = \begin{bmatrix} t \cos t \\ t \sin t \end{bmatrix}}$$

$$\frac{\partial f}{\partial x_1} = x_2^2 \exp(x_1 x_2^2)$$

$$\frac{\partial f}{\partial x_2} = 2x_1 x_2 \exp(x_1 x_2^2)$$

- We compute the gradient of  $h$  with respect to  $t$ . Since  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$  and  $g: \mathbb{R} \rightarrow \mathbb{R}^2$  we note that

$$\frac{\partial f}{\partial x} \in \mathbb{R}^{1 \times 2}, \quad \frac{\partial g}{\partial t} \in \mathbb{R}^{2 \times 1}$$

- The desired gradient is computed by applying the chain rule: *and product rule.*

$$\left[ \frac{dh}{dt} \right] = \left[ \frac{\partial f}{\partial x} \frac{\partial x}{\partial t} \right] = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} \end{bmatrix} \begin{bmatrix} \frac{\partial x_1}{\partial t} \\ \frac{\partial x_2}{\partial t} \end{bmatrix}$$

$$= \begin{bmatrix} \exp(x_1 x_2^2) x_2^2 & 2 \exp(x_1 x_2^2) x_1 x_2 \end{bmatrix} \begin{bmatrix} \cos t - t \sin t \\ \sin t + t \cos t \end{bmatrix}$$

$$= \exp(x_1 x_2^2) \left( x_2^2 (\cos t - t \sin t) + 2x_1 x_2 (\sin t + t \cos t) \right)$$

where  $x_1 = t \cos t$  and  $x_2 = t \sin t$

$$f(x_1, x_2) = \exp(x_1 \cdot x_2^2) \quad g(t) = \begin{bmatrix} t \cos t \\ t \sin t \end{bmatrix}$$

Alternative:  $h: \mathbb{R} \rightarrow \mathbb{R}$   
 $h(t) = (f \circ g)(t) = f\left(\begin{bmatrix} t \cos t \\ t \sin t \end{bmatrix}\right) = \exp\left((t \cos t)(t \sin t)^2\right)$

$$h'(t) = \left( \underbrace{t \cos t (t \sin t)^2}_{(t^3 \cos t \sin^2 t)} \right)' \exp(t \cos t (t \sin t)^2) \quad (e^{f(x)})' = f'(x) e^{f(x)}$$

Exercise:

Same answer as before.

$$\begin{aligned} (f \circ g)' & \quad \text{Triple Product Rule} \\ &= f'(gh) + f(gh)' \\ &= f'(gh) + f(g'h + gh') \\ &= f'gh + fg'h + fgh' \end{aligned}$$

# Example - Gradient of a Least-Squares Loss in a Linear Model

- Let us consider the linear model  $\hat{y} = \Phi \theta$  where  $\theta \in \mathbb{R}^D$  is a parameter vector,  $\Phi \in \mathbb{R}^{N \times D}$  are input features and  $y \in \mathbb{R}^N$  are the corresponding observations. We define the functions

$$L(e) := \|e\|_2^2, \quad L: \mathbb{R}^N \rightarrow \mathbb{R}$$

$$e(\theta) := y - \Phi \theta, \quad e: \mathbb{R}^D \rightarrow \mathbb{R}^N$$

- We seek  $\frac{\partial L}{\partial \theta}$ , and we will use the chain rule for this purpose.  $L$  is called a least-squares loss function.
- First, we determine the dimensionality of the gradient as

$$\frac{\partial L}{\partial \theta} \in \mathbb{R}^{1 \times D} \quad (\text{as } L \circ e: \mathbb{R}^D \rightarrow \mathbb{R})$$

- The chain rule allows us to compute the gradient as

$$\frac{\partial L}{\partial \theta} = \underbrace{\frac{\partial L}{\partial e}}_{\text{blue}} \underbrace{\frac{\partial e}{\partial \theta}}_{\text{orange}} \quad \nabla_x(Ax) = A$$

$$\frac{\partial L}{\partial \vec{e}} = \frac{\partial}{\partial \vec{e}} (\vec{e}^T \vec{e}) = 2 \vec{e}^T$$

$$\frac{\partial \vec{e}}{\partial \vec{\theta}} = \frac{\partial}{\partial \vec{\theta}} (\cancel{y} - \Phi \vec{\theta}) = \frac{\partial}{\partial \vec{\theta}} (-\Phi \vec{\theta}) = -\Phi$$

# Example - Gradient of a Least-Squares Loss in a Linear Model

$$L(\theta) : \mathbb{R}^n \rightarrow \mathbb{R}$$

- We know that  $\|e\|^2 = e^T e$  and determine

$$\frac{\partial L}{\partial e} = 2e^T \in \boxed{\mathbb{R}^{1 \times N}} \quad \vec{e} = \vec{y} - \Phi \vec{\theta}$$

- Further, we obtain

$$\frac{\partial e}{\partial \theta} = -\Phi \in \boxed{\mathbb{R}^{N \times D}} \quad e(\theta) : \mathbb{R}^D \rightarrow \mathbb{R}^N$$

- Our desired derivative is

$$\frac{\partial L}{\partial \theta} = -\underbrace{2 \underbrace{e^T}_{1 \times N} \underbrace{\Phi}_{N \times D}}_{1 \times D} = -2(\underbrace{y^T}_{1 \times N} - \underbrace{\theta^T \Phi^T}_{1 \times N}) \underbrace{\Phi}_{N \times D} \in \mathbb{R}^{1 \times D}$$

Solve  $\frac{\partial L}{\partial \theta} = 0$  for  $\theta$ , to find best  $\theta$ .

$$\cancel{-(y^T - \theta^T \Phi^T) \Phi} = 0$$
$$y^T \Phi - \theta^T \Phi^T \Phi = 0$$

$$y^T \Phi = \theta^T \Phi^T \Phi$$
$$y^T \Phi (\Phi^T \Phi)^{-1} = \theta^T$$
$$\theta = \underbrace{((\Phi^T \Phi)^T)^{-1}}_{\text{Pseudo-Inv.}} \Phi^T y = \underbrace{(\Phi^T \Phi)^{-1}}_{\text{Pseudo-Inv.}} \Phi^T y$$

# 5.4 Gradients of Matrices

domain of interest

codomain

$$f(\underline{A}, \underline{x}) : \mathbb{R}^{\underline{M \times N}} \times \mathbb{R}^N \rightarrow \mathbb{R}^{\underline{M}}$$

- Consider the following example

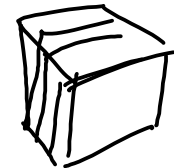
$$\underline{f} = \underline{A}\underline{x}, \quad \underline{f} \in \mathbb{R}^M, \quad \underline{A} \in \mathbb{R}^{M \times N}, \quad \underline{x} \in \mathbb{R}^N$$

- We seek the gradient  $\boxed{\frac{df}{dA}}$

$$f_x(A) : \mathbb{R}^{M \times N} \rightarrow \mathbb{R}^M$$

- First, we determine the dimension of the gradient

$$\frac{df}{dA} \in \mathbb{R}^{\underline{M} \times \underline{(M \times N)}} \quad (3\text{-tensor})$$



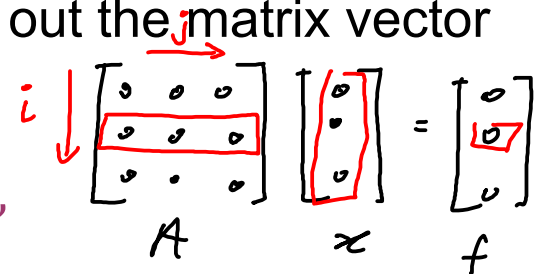
- By definition, the gradient is the collection of the partial derivatives:

$$\frac{df}{dA} = \begin{bmatrix} \frac{\partial f_1}{\partial A} \\ \vdots \\ \frac{\partial f_M}{\partial A} \end{bmatrix},$$

$$\boxed{\frac{\partial f_i}{\partial A} \in \mathbb{R}^{1 \times (M \times N)}} \approx \mathbb{R}^{M \times N}$$

- To compute the partial derivatives, we explicitly write out the matrix vector multiplication

$$\underline{f_i} = \sum_{j=1}^N \underline{A_{ij}} \underline{x_j}, \quad i = 1, \dots, M,$$



$$f_i = \sum_j A_{ij} x_j$$

$\nabla_A f_i$

compute  $\frac{\partial f_i}{\partial A_{pq}} = \frac{\partial}{\partial A_{pq}} \left( \sum_j A_{ij} x_j \right)$

$$= \sum_j x_j \frac{\partial A_{ij}}{\partial A_{pq}} = \begin{cases} 1 & \boxed{i=p, j=q} \\ 0 & \text{else} \end{cases} = \delta_{ip} \delta_{jq}$$

$$= \sum_j x_j \delta_{ip} \delta_{jq} = \underline{x_q \delta_{ip}}$$

$$\frac{\partial f_i}{\partial A_{iq}} = x_q$$

$$\frac{\partial f_i}{\partial A_{pq}} = 0$$

$p \neq i$

$$f_i = \sum_{j=1}^N A_{ij} x_j, \quad i = 1, \dots, M, \quad \frac{\partial f_i}{\partial A_{pq}}$$

- The partial derivatives are then given as

$$\frac{\partial f_i}{\partial A_{iq}} = x_q, \quad p \neq i$$

- Partial derivatives of  $f_i$  with respect to a row of  $\mathbf{A}$  are given as

$$\frac{\partial f_i}{\partial A_{i,:}} = \mathbf{x}^T \in \mathbb{R}^{1 \times 1 \times N}, \quad \frac{\partial f_i}{\partial A_{k \neq i,:}} = \mathbf{0}^T \in \mathbb{R}^{1 \times 1 \times N}$$

- Since  $f_i$  maps onto  $\mathbb{R}$  and each row of  $\mathbf{A}$  is of size  $1 \times N$ , we obtain a  $1 \times 1 \times N$  sized tensor as the partial derivative of  $f_i$  with respect to a row of  $\mathbf{A}$ .

- We stack the partial derivatives and get the desired gradient

$$\boxed{\frac{\partial f_i}{\partial \mathbf{A}}} = \begin{bmatrix} \mathbf{0}^T \\ \vdots \\ \mathbf{0}^T \\ \mathbf{x}^T \\ \mathbf{0}^T \\ \vdots \\ \mathbf{0}^T \end{bmatrix} \in \mathbb{R}^{1 \times (M \times N)} = \begin{bmatrix} 0 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 \\ x_1 & x_2 & \dots & x_n \\ 0 & 0 & \dots & 0 \end{bmatrix}$$

So  $\frac{df}{dA} = \begin{bmatrix} \begin{bmatrix} x^T \\ 0^T \\ \vdots \\ 0^T \end{bmatrix}, \begin{bmatrix} 0^T \\ x^T \\ \vdots \\ 0^T \end{bmatrix}, \dots, \begin{bmatrix} 0^T \\ 0^T \\ \vdots \\ x^T \end{bmatrix} \end{bmatrix} \in \mathbb{R}^{M \times (M \times N)}$

$\frac{df}{dA} = \begin{bmatrix} x_1, x_2, \dots, x_n \\ 0, 0, \dots, 0 \\ \vdots \\ 0, 0, \dots, 0 \end{bmatrix}$  (with dimensions  $M$  and  $N$  indicated)

$i^{th}$

$\begin{bmatrix} 0 & 0 & \dots & 0 \\ x_1 & x_2 & \dots & x_n \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix}$

$\dots$

$\begin{bmatrix} 0 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 \\ \vdots & \vdots & & \vdots \\ x_1 & x_2 & \dots & x_n \end{bmatrix}$



# Example - Gradient of Matrices with Respect to Matrices

- Consider a matrix  $\mathbf{R} \in \mathbb{R}^{M \times N}$  and  $f: \mathbb{R}^{M \times N} \rightarrow \mathbb{R}^{N \times N}$  with  $f(\mathbf{R}) = \mathbf{R}^T \mathbf{R} =: \mathbf{K} \in \mathbb{R}^{N \times N}$

- We seek the gradient  $\frac{d\mathbf{K}}{d\mathbf{R}} = \frac{d f(\mathbf{R})}{d\mathbf{R}}$

- First, the dimension of the gradient is given as

$$\frac{d\mathbf{K}}{d\mathbf{R}} \in \mathbb{R}^{(N \times N) \times (M \times N)} \quad (4\text{-tensor})$$

$$\frac{dK_{pq}}{d\mathbf{R}} \in \mathbb{R}^{1 \times M \times N} \approx \mathbb{R}^{M \times N}$$

for  $p, q = 1, \dots, N$ , where  $K_{pq}$  is the  $pq$ th entry of  $\mathbf{K} = f(\mathbf{R})$ .  $\mathbf{K} = \mathbf{R}^T \mathbf{R}$   
 $K_{pq} = (\mathbf{R}^T \mathbf{R})_{pq}$

- Denoting the  $i$ th column of  $\mathbf{R}$  by  $\mathbf{r}_i$ , every entry of  $\mathbf{K}$  is given by the dot product of two columns of  $\mathbf{R}$ , i.e.,

$$K_{pq} = \mathbf{r}_p^T \mathbf{r}_q = \sum_{m=1}^M R_{mp} R_{mq}$$

$$= \sum_m R_{pm}^T R_{mq}$$

# Example - Gradient of Matrices with Respect to Matrices

- Denoting the  $i$ th column of  $\mathbf{R}$  by  $\mathbf{r}_i$ , every entry of  $\mathbf{K}$  is given by the dot product of two columns of  $\mathbf{R}$ , i.e.,

$$K_{pq} = \mathbf{r}_p^T \mathbf{r}_q = \sum_{m=1}^M R_{mp} R_{mq}$$

$$\begin{aligned} \frac{\partial K_{pq}}{\partial R_{ij}} &= \frac{\partial}{\partial R_{ij}} \sum_m R_{mp} R_{mq} \\ &= \sum_m \left( \delta_{im} R_{mq} + R_{mp} \delta_{jm} \right) \end{aligned}$$

$i=m, j=p, j=q$

- We now compute the partial derivative  $\frac{\partial K_{pq}}{\partial R_{ij}}$ , we obtain

$$\frac{\partial K_{pq}}{\partial R_{ij}} = \sum_{m=1}^M \left( \frac{\partial}{\partial R_{ij}} R_{mp} R_{mq} \right) = \partial_{pqij}$$

Exercise.

$$\partial_{pqij} = \begin{cases} R_{iq} & \text{if } j = p, p \neq q \\ R_{ip} & \text{if } j = q, p \neq q \\ 2R_{iq} & \text{if } j = p, p = q \\ 0 & \text{otherwise} \end{cases}$$

$\leftarrow R_{mq} \quad R_{ij}$

(4-tensor)

- The desired gradient has the dimension  $(N \times N) \times (M \times N)$ , and every single entry of this tensor is given by  $\partial_{pqij}$ , where  $p, q, j = 1, \dots, N$  and  $i = 1, \dots, M$



# 5.5 Useful Identities for Computing Gradients

- Some useful gradients that are frequently required in machine learning
- $\text{tr}(\cdot)$ : trace     $\det(\cdot)$ : determinant     $f(X)^{-1}$ : the inverse of  $f(X)$

$$X \cdot Y = X^T Y$$

$$\langle x, y \rangle_A = x^T A y$$

$A$  is sym and pos. def.

$$\left\{ \begin{array}{l} \frac{\partial x^T a}{\partial x} = a^T \\ \frac{\partial a^T x}{\partial x} = a^T \end{array} \right\} \text{Already proven.}$$

$$\rightarrow \boxed{\frac{\partial a^T X b}{\partial X} = ab^T} \text{ We will prove this!}$$

$$\rightarrow \boxed{\frac{\partial x^T B x}{\partial x} = x^T (B + B^T)} \text{ Assignment 2}$$

$$\frac{\partial}{\partial s} \underline{(x - As)^T W (x - As)} = -2(x - As)^T W A \text{ for symmetric } W$$

You should be able to calculate these gradients

Prove

$$\frac{\partial (a^T X b)}{\partial X} = a b^T$$

$$a \in \mathbb{R}^{n \times 1}$$

$$b \in \mathbb{R}^{m \times 1}$$

$$X \in \mathbb{R}^{n \times m}$$

$$a^T X b \in \mathbb{R}$$

$$\textcircled{1} \xrightarrow{n \times m} \xrightarrow{m \times 1} \textcircled{1}$$

$$f(X) = \bar{a}^T X \bar{b} \\ : \mathbb{R}^{n \times m} \rightarrow \mathbb{R}'$$

$$\nabla_X (a^T X b) \in \mathbb{R}^{1 \times (n \times m)}$$

$$\simeq \mathbb{R}^{n \times m}$$

$$\nabla_X (a^T X b) = \nabla_X (a^T (X b))$$

$$a^T X b = a^T (X b) = \sum_i a_i (X b)_i = \sum_i a_i \left( \sum_j x_{ij} b_j \right)$$

$$(a b^T)$$

$$\begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \begin{bmatrix} \dots \end{bmatrix} = \begin{bmatrix} a_1 b_1 & a_2 b_1 \\ a_1 b_2 & \dots \end{bmatrix}$$

$$\sum_{i,j} a_i x_{ij} b_j$$

$$\begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \begin{bmatrix} \dots \end{bmatrix} \rightarrow \square$$

$$\frac{\partial}{\partial x_{pq}} \left( \sum_{i,j} a_i x_{ij} b_j \right)$$

$$= \sum_{i,j} a_i b_j \frac{\partial x_{ij}}{\partial x_{pq}}$$

$$\begin{cases} 1 & i=p, j=q \\ 0 & \text{else} \end{cases}$$

$$\frac{\partial x_{ij}}{\partial x_{pq}} = \delta_{ip} \delta_{jq}$$

$$\begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \begin{bmatrix} \dots \end{bmatrix} \rightarrow \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$$

$$\underline{\delta_{jq}} = \underline{\delta_{qj}}$$

$$= \sum_{i,j} a_i b_j \underline{\delta_{ip}} \underline{\delta_{jq}} \xrightarrow{i=p, j=q} \underline{a_p b_q} = \underline{(\bar{a} \bar{b}^T)_{pq}}$$

$$\boxed{\frac{\partial}{\partial x} (a^T x b)}$$

$\downarrow h(x)$

$$f(g) = a^T g$$

$$g(x) = Xb$$

$$h(x) = (f \circ g)(x)$$

$$= f(g(x))$$

$$g: \mathbb{R}^{n \times m} \rightarrow \mathbb{R}^n$$

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\frac{\partial h}{\partial x} = \frac{\partial h}{\partial g} \boxed{\frac{\partial g}{\partial x}}$$

3 tensor

{  
vector

matrix

K-Means

$$r_{nh} = \begin{cases} 1 & x_n \text{ is cluster } k \\ 0 & \text{else} \end{cases}$$

$$L(M, R) = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|_2^2$$

$$R_{nh} = r_{nh} \quad M = (\mu_1, \dots, \mu_K)$$

$$\begin{bmatrix} \mu_1 & \mu_2 & \dots & \mu_K \\ \vdots & \vdots & & \vdots \end{bmatrix}$$

Compute  $\nabla_M L = \left[ \frac{\partial L}{\partial \mu_1}, \dots, \frac{\partial L}{\partial \mu_K} \right] = 0$

$$\frac{\partial L}{\partial \mu_j} = \frac{\partial}{\partial \mu_j} \left( \sum_n \sum_k r_{nk} \|\bar{x}_n - \mu_k\|_2^2 \right)$$

$$= \sum_n \frac{\partial}{\partial \mu_j} \sum_k r_{nk} \|\bar{x}_n - \mu_k\|_2^2$$

(only non-zero term  $k=j$ )

$$= \sum_n r_{nj} \frac{\partial}{\partial \mu_j} \|\bar{x}_n - \mu_j\|_2^2$$

$$= \sum_n r_{nj} \frac{\partial}{\partial \mu_j} (-2x_n^T \mu_j + \mu_j^T \mu_j)$$

$$= \sum_n r_{nj} (-2x_n^T + 2\mu_j^T) = 0$$

$$\mu_k = \frac{\sum_n r_{nk} x_n}{\sum_n r_{nk}}$$

$$r_{n1} \|x_n - \mu_1\|_2^2 + \dots + r_{nK} \|x_n - \mu_K\|_2^2$$

$$\frac{\partial}{\partial x} (C^T \bar{x}) = C^T$$

$$\frac{\partial}{\partial \bar{x}} (\bar{x}^T \bar{x}) = 2\bar{x}^T$$

$$\|x_n - \mu_j\|_2^2 = (x_n - \mu_j)^T (x_n - \mu_j)$$

$$= \cancel{x_n^T x_n} - \mu_j^T x_n - x_n^T \mu_j + \mu_j^T \mu_j$$

$$= -2x_n^T \mu_j + \mu_j^T \mu_j$$

$$\sum_n r_{nj} (-2x_n^T + 2\mu_j^T) = 0$$

$$= \sum_n r_{nj} (-2x_n^T) + \sum_n r_{nj} (2\mu_j^T) = 0$$

$$\sum_n r_{nj} \mu_j = \sum_n r_{nj} x_n$$

$$\mu_j = \frac{\sum_n r_{nj} x_n}{\sum_n r_{nj}}$$