

# Lecture 5: Vector Calculus

Mathematics for Machine Learning

# Summary

- Machine learning is about solving an optimization problem whose variables are the parameters of a given model.
- Solving optimization problems require gradient information.
- Central to this chapter is the concept of the function, which we often write

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

$$\boldsymbol{x} \mapsto f(\boldsymbol{x})$$

1 Differentiation of Univariate Functions

2 Partial Differentiation and Gradients

3 Gradients of Vector-Valued Functions

4 Gradients of Matrices

# Difference Quotient and Derivative

- **Difference Quotient.** The average slope of  $f$  between  $x$  and  $x + \partial x$

$$\frac{\partial y}{\partial x} := \frac{f(x + \partial x) - f(x)}{\partial x}$$

- **Derivative.** Pointing in the direction of steepest ascent of  $f$ .

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

- Unless confusion arises, we often use  $f' = \frac{df}{dx}$ .

# Taylor Series

- Representation of a function as an infinite sum of terms, using derivatives of evaluated at  $x_0$ .
- **Taylor polynomial.** The Taylor polynomial of degree  $n$  of  $f : \mathbb{R} \mapsto \mathbb{R}$  at  $x_0$  is:

$$T_n(x) := \sum_{k=0}^n \frac{f^{(k)}(x_0)}{k!} (x - x_0)^k, \text{ where } f^{(k)}(x_0) \text{ is the } k\text{th derivative of } f \text{ at } x_0.$$

- **Taylor Series.** For a smooth function  $f \in \mathcal{C}^\infty$ , the Taylor series of  $f$  at  $x_0$  is:

$$T_\infty(x) := \sum_{k=0}^{\infty} \frac{f^{(k)}(x_0)}{k!} (x - x_0)^k.$$

- If  $f(x) = T_\infty(x)$ ,  $f$  is called **analytic**.

# Differentiation Rules

- **Product rule.**  $(f(x)g(x))' = f'(x)g(x) + f(x)g'(x)$
- **Quotient rule.**  $\left(\frac{f(x)}{g(x)}\right)' = \frac{f'(x)g(x) - f(x)g'(x)}{(g(x))^2}$
- **Sum rule.**  $(f(x) + g(x))' = f'(x) + g'(x)$
- **Chain rule.**  $(g(f(x)))' = g'(f(x))f'(x)$

1 Differentiation of Univariate Functions

2 Partial Differentiation and Gradients

3 Gradients of Vector-Valued Functions

4 Gradients of Matrices

# Gradient

- Now,  $f : \mathbb{R}^n \mapsto \mathbb{R}$ .
- Gradient of  $f$  w.r.t.  $\mathbf{x}$   $\nabla_{\mathbf{x}} f$ : Varying one variable at a time and keeping the others constant.

**Partial Derivative.** For  $f : \mathbb{R}^n \mapsto \mathbb{R}$ ,

$$\begin{aligned}\frac{\partial f}{\partial x_1} &= \lim_{h \rightarrow 0} \frac{f(x_1 + h, x_2, \dots, x_n) - f(\mathbf{x})}{h} \\ &\vdots \\ \frac{\partial f}{\partial x_n} &= \lim_{h \rightarrow 0} \frac{f(x_1, x_2, \dots, x_n + h) - f(\mathbf{x})}{h}\end{aligned}$$

**Gradient.** Get the partial derivatives and collect them in the row vector.

$$\nabla_{\mathbf{x}} f = \frac{df}{d\mathbf{x}} = \left( \frac{\partial f(\mathbf{x})}{\partial x_1} \quad \dots \quad \frac{\partial f(\mathbf{x})}{\partial x_n} \right) \in \mathbb{R}^{1 \times n}$$



# Example

- **Example.**  $f(x, y) = (x + 2y^3)^2$

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

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

- **Example.**  $f(x_1, x_2) = x_1^2 x_2 + x_1 x_2^3$

$$\nabla_{(x_1, x_2)} f = \frac{df}{dx} = \left( \frac{\partial f(x_1, x_2)}{\partial x_1} \quad \frac{\partial f(x_1, x_2)}{\partial x_2} \right) = \left( 2x_1 x_2 + x_2^3 \quad x_1^2 + 3x_1 x_2^2 \right)$$

# Rules for Partial Differentiation

- Product rule

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

- Sum rule

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

- Chain rule

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

## More about Chain Rule

- $f : \mathbb{R}^2 \mapsto \mathbb{R}$  of two variables  $x_1$  and  $x_2$ .  $x_1(t)$  and  $x_2(t)$  are functions of  $t$ .

$$\frac{df}{dt} = \begin{pmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} \end{pmatrix} \begin{pmatrix} \frac{\partial x_1(t)}{\partial t} \\ \frac{\partial x_2(t)}{\partial t} \end{pmatrix} = \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial t} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial t}$$

- **Example.**  $f(x_1, x_2) = x_1^2 + 2x_2$ , where  $x_1(t) = \sin(t)$ ,  $x_2(t) = \cos(t)$

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

- $f : \mathbb{R}^2 \mapsto \mathbb{R}$  of two variables  $x_1$  and  $x_2$ .  $x_1(s, t)$  and  $x_2(s, t)$  are functions of  $s, t$ .

$$\left. \begin{aligned} \frac{\partial f}{\partial s} &= \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial s} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial s} \\ \frac{\partial f}{\partial t} &= \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial t} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial t} \end{aligned} \right| \frac{df}{d(s, t)} = \frac{\partial f}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial (s, t)} = \begin{pmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} \end{pmatrix} \begin{pmatrix} \frac{\partial x_1}{\partial s} & \frac{\partial x_1}{\partial t} \\ \frac{\partial x_2}{\partial s} & \frac{\partial x_2}{\partial t} \end{pmatrix}$$

1 Differentiation of Univariate Functions

2 Partial Differentiation and Gradients

3 Gradients of Vector-Valued Functions

4 Gradients of Matrices

$$\mathbf{f} : \mathbb{R}^n \mapsto \mathbb{R}^m$$

- For a function  $\mathbf{f} : \mathbb{R}^n \mapsto \mathbb{R}^m$  and vector  $\mathbf{x} = (x_1 \ \dots \ x_n)^\top \in \mathbb{R}^n$ , the vector-valued function is:

$$\mathbf{f}(\mathbf{x}) = \begin{pmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_m(\mathbf{x}) \end{pmatrix}$$

- Partial derivative w.r.t.  $x_i$  is a column vector:  $\frac{\partial \mathbf{f}}{\partial x_i} = \begin{pmatrix} \frac{\partial f_1}{\partial x_i} \\ \vdots \\ \frac{\partial f_m}{\partial x_i} \end{pmatrix}$
- Gradient (or Jacobian):  $\frac{\mathrm{d}\mathbf{f}(\mathbf{x})}{\mathrm{d}\mathbf{x}} = \left( \frac{\partial \mathbf{f}(\mathbf{x})}{\partial x_1} \ \dots \ \frac{\partial \mathbf{f}(\mathbf{x})}{\partial x_n} \right)$

$$\begin{aligned}\mathbf{J} &= \nabla_{\mathbf{x}} \mathbf{f} = \frac{d\mathbf{f}(\mathbf{x})}{d\mathbf{x}} = \left( \frac{\partial \mathbf{f}(\mathbf{x})}{\partial x_1} \quad \dots \quad \frac{\partial \mathbf{f}(\mathbf{x})}{\partial x_n} \right) \\ &= \begin{pmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial f_1(\mathbf{x})}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_m(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial f_m(\mathbf{x})}{\partial x_n} \end{pmatrix}\end{aligned}$$

- For a  $\mathbb{R}^n \mapsto \mathbb{R}^m$  function, its Jacobian is a  $m \times n$  matrix.

## Example: Gradient of Vector-Valued Function

- $\mathbf{f}(\mathbf{x}) = \mathbf{A}\mathbf{x}$ ,  $\mathbf{f} : \mathbb{R}^n \mapsto \mathbb{R}^m$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{x} \in \mathbb{R}^n$
- Partial derivatives:  $f_i(\mathbf{x}) = \sum_{j=1}^n A_{ij}x_j \implies \frac{\partial f_i}{\partial x_j} = A_{ij}$
- Gradient

$$\frac{d\mathbf{f}}{d\mathbf{x}} = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{pmatrix} = \begin{pmatrix} A_{11} & \cdots & A_{1n} \\ \vdots & & \vdots \\ A_{m1} & \cdots & A_{mn} \end{pmatrix} = \mathbf{A}$$

## Example: Chain Rule

- $h : \mathbb{R} \mapsto \mathbb{R}$ ,  $h(t) = (f \circ g)(t)$  with

$$f : \mathbb{R}^2 \mapsto \mathbb{R}, f(\mathbf{x}) = \exp(x_1 x_2^2), \quad g : \mathbb{R} \mapsto \mathbb{R}^2, \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = g(t) = \begin{pmatrix} t \cos(t) \\ t \sin(t) \end{pmatrix}$$

- (Note)  $\frac{\partial f}{\partial \mathbf{x}} \in \mathbb{R}^{1 \times 2}$  and  $\frac{\partial g}{\partial t} \in \mathbb{R}^{2 \times 1}$
- Using the chain rule,

$$\begin{aligned} \frac{dh}{dt} &= \frac{\partial f}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial t} = \begin{pmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} \end{pmatrix} \begin{pmatrix} \frac{\partial x_1}{\partial t} \\ \frac{\partial x_2}{\partial t} \end{pmatrix} \\ &= \left( \exp(x_1 x_2^2) x_2^2 \quad 2 \exp(x_1 x_2^2) x_1 x_2 \right) \begin{pmatrix} \cos(t) - t \sin(t) \\ \sin(t) + t \cos(t) \end{pmatrix} \end{aligned}$$



## Example: Least-Square Loss (1)

- A linear model:  $\mathbf{y} = \Phi\boldsymbol{\theta}$
- $\boldsymbol{\theta} \in \mathbb{R}^D$ : parameter vector
- $\Phi \in \mathbb{R}^{N \times D}$ : input features
- $\mathbf{y} \in \mathbb{R}^N$ : observations
- Goal: Find a good parameter vector that provides the best-fit, formulated by minimizing the following loss  $L : \mathbb{R}^D \mapsto \mathbb{R}$  over the parameter vector  $\boldsymbol{\theta}$ .

$$L(\mathbf{e}) := \|\mathbf{e}\|^2, \quad \text{where } \mathbf{e}(\boldsymbol{\theta}) = \mathbf{y} - \Phi\boldsymbol{\theta}$$

## Example: Least-Square Loss (2)

- $\frac{\partial L}{\partial \theta} = \frac{\partial L}{\partial e} \frac{\partial e}{\partial \theta}$
- **Note.**  $\frac{\partial L}{\partial \theta} \in \mathbb{R}^{1 \times D}$ ,  $\frac{\partial L}{\partial e} \in \mathbb{R}^{1 \times N}$ ,  $\frac{\partial e}{\partial \theta} \in \mathbb{R}^{N \times D}$
- Using that  $\|e\|^2 = e^T e$ ,  $\frac{\partial L}{\partial e} = 2e^T \in \mathbb{R}^{1 \times N}$  and  $\frac{\partial e}{\partial \theta} = -\Phi \in \mathbb{R}^{N \times D}$

Finally, we get: 
$$\frac{\partial L}{\partial \theta} = 2e^T(-\Phi) = -\underbrace{2(y^T - \theta^T \Phi^T)}_{1 \times N} \underbrace{\Phi}_{N \times D}$$

1 Differentiation of Univariate Functions

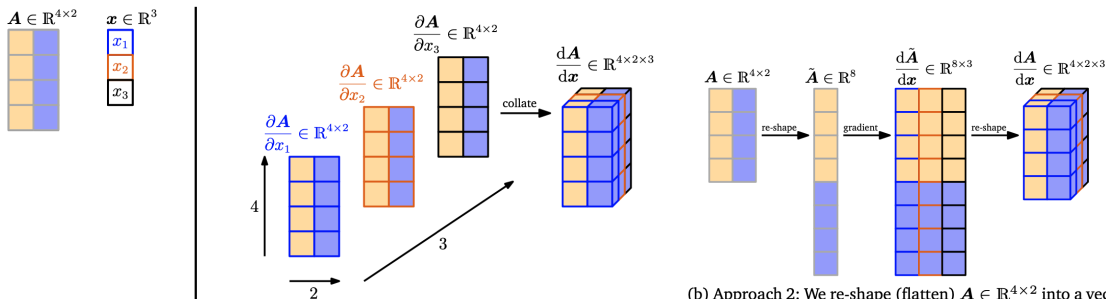
2 Partial Differentiation and Gradients

3 Gradients of Vector-Valued Functions

4 Gradients of Matrices

# Gradients of matrices

- Gradient of matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  w.r.t. matrix  $\mathbf{B} \in \mathbb{R}^{p \times q}$
- Jacobian: A four-dimensional tensor<sup>1</sup>  $\mathbf{J} = \frac{d\mathbf{A}}{d\mathbf{B}} \in \mathbb{R}^{(m \times n) \times (p \times q)}$



(a) Approach 1: We compute the partial derivative  $\frac{\partial \mathbf{A}}{\partial x_1}$ ,  $\frac{\partial \mathbf{A}}{\partial x_2}$ ,  $\frac{\partial \mathbf{A}}{\partial x_3}$ , each of which is a  $4 \times 2$  matrix, and collate them in a  $4 \times 2 \times 3$  tensor.

(b) Approach 2: We re-shape (flatten)  $\mathbf{A} \in \mathbb{R}^{4 \times 2}$  into a vector  $\tilde{\mathbf{A}} \in \mathbb{R}^8$ . Then, we compute the gradient  $\frac{d\tilde{\mathbf{A}}}{d\mathbf{x}} \in \mathbb{R}^{8 \times 3}$ . We obtain the gradient tensor by re-shaping this gradient as illustrated above.

<sup>1</sup>A multidimensional array

## Example: Gradient of Vectors for Matrices (1)

- $\mathbf{f}(\mathbf{x}) = \mathbf{A}\mathbf{x}$ ,  $\mathbf{f} \in \mathbb{R}^m$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{x} \in \mathbb{R}^n$ . What is  $\frac{d\mathbf{f}}{d\mathbf{A}}$ ?
- Dimension: If we consider  $\mathbf{f} : \mathbb{R}^{m \times n} \mapsto \mathbb{R}^m$ ,  $\frac{d\mathbf{f}}{d\mathbf{A}} \in \mathbb{R}^{m \times (m \times n)}$

- Partial derivatives:  $\frac{\partial f_i}{\partial \mathbf{A}} \in \mathbb{R}^{1 \times (m \times n)}$ ,  $\frac{d\mathbf{f}}{d\mathbf{A}} = \begin{pmatrix} \frac{\partial f_1}{\partial \mathbf{A}} \\ \vdots \\ \frac{\partial f_m}{\partial \mathbf{A}} \end{pmatrix}$

$$f_i = \sum_{j=1}^n A_{ij}x_j, \quad i = 1, \dots, m \implies \frac{\partial f_i}{\partial A_{iq}} = x_q,$$

$$\frac{\partial f_i}{\partial A_{i\cdot}} = \mathbf{x}^\top \in \mathbb{R}^{1 \times 1 \times n} \quad (\text{for } i\text{th row vector})$$

$$\frac{\partial f_i}{\partial A_{k \neq i \cdot}} = \mathbf{0}^\top \in \mathbb{R}^{1 \times 1 \times n} \quad (\text{for } k\text{th row vector, } k \neq i)$$

$$\frac{\partial f_i}{\partial \mathbf{A}} = \begin{pmatrix} \mathbf{0}^\top \\ \vdots \\ \mathbf{0}^\top \\ \mathbf{x}^\top \\ \mathbf{0}^\top \\ \vdots \\ \mathbf{0}^\top \end{pmatrix} \in \mathbb{R}^{1 \times (m \times n)}$$

## Example: Gradient of Matrices for Matrices (2)

- $\mathbf{R} \in \mathbb{R}^{m \times n}$  and  $\mathbf{f} : \mathbb{R}^{m \times n} \mapsto \mathbb{R}^{n \times n}$  with  $\mathbf{f}(\mathbf{R}) = \mathbf{K} := \mathbf{R}^\top \mathbf{R} \in \mathbb{R}^{n \times n}$ . What is  $\frac{d\mathbf{K}}{d\mathbf{R}} \in \mathbb{R}^{(n \times n) \times (m \times n)}$ ?
- $\frac{dK_{pq}}{d\mathbf{R}} \in \mathbb{R}^{1 \times m \times n}$ . Let  $\mathbf{r}_i$  be the  $i$ th column of  $\mathbf{R}$ . Then  $K_{pq} = \mathbf{r}_p^\top \mathbf{r}_q = \sum_{k=1}^m R_{kp} R_{kq}$ .
- Partial derivative  $\frac{\partial K_{pq}}{\partial R_{ij}}$

$$\frac{\partial K_{pq}}{\partial R_{ij}} = \sum_{k=1}^m \frac{\partial}{\partial R_{ij}} R_{kp} R_{kq} = \partial_{pqij}, \quad \partial_{pqij} = \begin{cases} R_{iq} & \text{if } j = p, p \neq q \\ R_{ip} & \text{if } j = q, p \neq q \\ 2R_{iq} & \text{if } j = p, p = q \\ 0 & \text{otherwise} \end{cases}$$

## Useful Identities

$$\frac{\partial}{\partial \mathbf{X}} \mathbf{f}(\mathbf{X})^\top = \left( \frac{\partial \mathbf{f}(\mathbf{X})}{\partial \mathbf{X}} \right)^\top \quad (5.99)$$

$$\frac{\partial}{\partial \mathbf{X}} \text{tr}(\mathbf{f}(\mathbf{X})) = \text{tr} \left( \frac{\partial \mathbf{f}(\mathbf{X})}{\partial \mathbf{X}} \right) \quad (5.100)$$

$$\frac{\partial}{\partial \mathbf{X}} \det(\mathbf{f}(\mathbf{X})) = \det(\mathbf{f}(\mathbf{X})) \text{tr} \left( \mathbf{f}(\mathbf{X})^{-1} \frac{\partial \mathbf{f}(\mathbf{X})}{\partial \mathbf{X}} \right) \quad (5.101)$$

$$\frac{\partial}{\partial \mathbf{X}} \mathbf{f}(\mathbf{X})^{-1} = -\mathbf{f}(\mathbf{X})^{-1} \frac{\partial \mathbf{f}(\mathbf{X})}{\partial \mathbf{X}} \mathbf{f}(\mathbf{X})^{-1} \quad (5.102)$$

$$\frac{\partial \mathbf{a}^\top \mathbf{X}^{-1} \mathbf{b}}{\partial \mathbf{X}} = -(\mathbf{X}^{-1})^\top \mathbf{a} \mathbf{b}^\top (\mathbf{X}^{-1})^\top \quad (5.103)$$

$$\frac{\partial \mathbf{x}^\top \mathbf{a}}{\partial \mathbf{x}} = \mathbf{a}^\top \quad (5.104)$$

$$\frac{\partial \mathbf{a}^\top \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a}^\top \quad (5.105)$$

$$\frac{\partial \mathbf{a}^\top \mathbf{X} \mathbf{b}}{\partial \mathbf{X}} = \mathbf{a} \mathbf{b}^\top \quad (5.106)$$

$$\frac{\partial \mathbf{x}^\top \mathbf{B} \mathbf{x}}{\partial \mathbf{x}} = \mathbf{x}^\top (\mathbf{B} + \mathbf{B}^\top) \quad (5.107)$$