

Foundations of Machine Learning  
African Masters in Machine Intelligence



**AIMS** | African Institute for  
Mathematical Sciences  
RWANDA

**Imperial College  
London**

# Vector Calculus

**Marc Deisenroth**

Quantum Leap Africa  
African Institute for Mathematical  
Sciences, Rwanda

Department of Computing  
Imperial College London



@mpd37

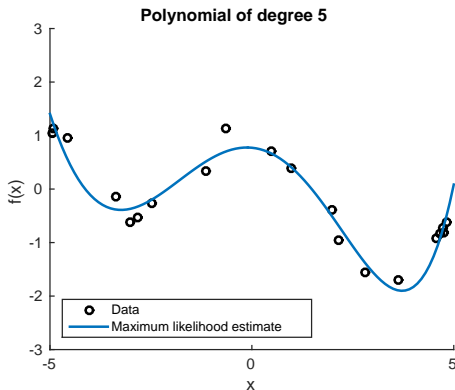
mdeisenroth@aimsammi.org

September 26, 2018

# Reference

Deisenroth et al.: Mathematics for Machine Learning, Chapter 5  
<https://mml-book.com>

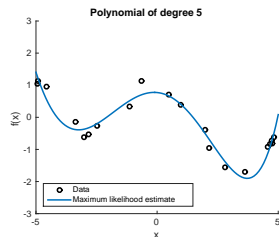
# Curve Fitting (Regression) in Machine Learning (1)



- ▶ Setting: Given inputs  $x$ , predict outputs/targets  $y$
- ▶ **Model**  $f$  that depends on parameters  $\theta$ . Examples:
  - ▶ Linear model:  $f(x, \theta) = \theta^\top x$ ,  $x, \theta \in \mathbb{R}^D$
  - ▶ Neural network:  $f(x, \theta) = NN(x, \theta)$

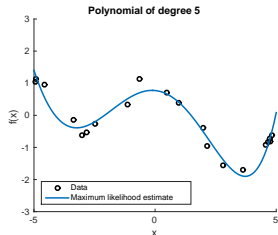
## Curve Fitting (Regression) in Machine Learning (2)

- ▶ Training data, e.g.,  $N$  pairs  $(x_i, y_i)$  of inputs  $x_i$  and observations  $y_i$
- ▶ **Training the model** means finding parameters  $\theta^*$ , such that  $f(x_i, \theta^*) \approx y_i$



## Curve Fitting (Regression) in Machine Learning (2)

- ▶ Training data, e.g.,  $N$  pairs  $(x_i, y_i)$  of inputs  $x_i$  and observations  $y_i$
- ▶ **Training the model** means finding parameters  $\theta^*$ , such that  $f(x_i, \theta^*) \approx y_i$



- ▶ Define a **loss function**, e.g.,  $\sum_{i=1}^N (y_i - f(x_i, \theta))^2$ , which we want to optimize
- ▶ Typically: Optimization based on some form of **gradient descent**
  - ▶ Differentiation required

# Types of Differentiation

1. Scalar differentiation:  $f : \mathbb{R} \rightarrow \mathbb{R}$

$y \in \mathbb{R}$  w.r.t.  $x \in \mathbb{R}$

2. Multivariate case:  $f : \mathbb{R}^N \rightarrow \mathbb{R}$

$y \in \mathbb{R}$  w.r.t. vector  $x \in \mathbb{R}^N$

3. Vector fields:  $f : \mathbb{R}^N \rightarrow \mathbb{R}^M$

vector  $y \in \mathbb{R}^M$  w.r.t. vector  $x \in \mathbb{R}^N$

4. General derivatives:  $f : \mathbb{R}^{M \times N} \rightarrow \mathbb{R}^{P \times Q}$

matrix  $y \in \mathbb{R}^{P \times Q}$  w.r.t. matrix  $X \in \mathbb{R}^{M \times N}$

## Scalar Differentiation $f : \mathbb{R} \rightarrow \mathbb{R}$

- ▶ Derivative defined as the limit of the difference quotient

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

- ▶▶ Slope of the secant line through  $f(x)$  and  $f(x+h)$

## Some Examples

$$f(x) = x^n$$

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

$$f(x) = \tanh(x)$$

$$f(x) = \exp(x)$$

$$f(x) = \log(x)$$

$$f'(x) = nx^{n-1}$$

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

$$f'(x) = 1 - \tanh^2(x)$$

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

$$f'(x) = \frac{1}{x}$$



# Rules

- ▶ Sum Rule

$$(f(x) + g(x))' = f'(x) + g'(x) = \frac{df(x)}{dx} + \frac{dg(x)}{dx}$$

# Rules

- ▶ Sum Rule

$$(f(x) + g(x))' = f'(x) + g'(x) = \frac{df(x)}{dx} + \frac{dg(x)}{dx}$$

- ▶ Product Rule

$$(f(x)g(x))' = f'(x)g(x) + f(x)g'(x) = \frac{df(x)}{dx}g(x) + f(x)\frac{dg(x)}{dx}$$

# Rules

- ▶ Sum Rule

$$(f(x) + g(x))' = f'(x) + g'(x) = \frac{df(x)}{dx} + \frac{dg(x)}{dx}$$

- ▶ Product Rule

$$(f(x)g(x))' = f'(x)g(x) + f(x)g'(x) = \frac{df(x)}{dx}g(x) + f(x)\frac{dg(x)}{dx}$$

- ▶ Chain Rule

$$(g \circ f)'(x) = (g(f(x)))' = g'(f(x))f'(x) = \frac{dg(f(x))}{df} \frac{df(x)}{dx}$$

# Rules

- ▶ Sum Rule

$$(f(x) + g(x))' = f'(x) + g'(x) = \frac{df(x)}{dx} + \frac{dg(x)}{dx}$$

- ▶ Product Rule

$$(f(x)g(x))' = f'(x)g(x) + f(x)g'(x) = \frac{df(x)}{dx}g(x) + f(x)\frac{dg(x)}{dx}$$

- ▶ Chain Rule

$$(g \circ f)'(x) = (g(f(x)))' = g'(f(x))f'(x) = \frac{dg(f(x))}{df} \frac{df(x)}{dx}$$

- ▶ Quotient Rule

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

## Example: Scalar Chain Rule

$$(g \circ f)'(x) = (g(f(x)))' = g'(f(x))f'(x) = \frac{dg}{df} \frac{df}{dx}$$

### Beginner

$$g(z) = 6z + 3$$

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

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

### Advanced

$$g(z) = \tanh(z)$$

$$z = f(x) = x^n$$

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

**Work it out with your neighbors**

## Example: Scalar Chain Rule

$$(g \circ f)'(x) = (g(f(x)))' = g'(f(x))f'(x) = \frac{dg}{df} \frac{df}{dx}$$

### Beginner

$$g(z) = 6z + 3$$

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

$$\begin{aligned}(g \circ f)'(x) &= \underbrace{(6)}_{dg/df} \underbrace{(-2)}_{df/dx} \\ &= -12\end{aligned}$$

### Advanced

$$g(z) = \tanh(z)$$

$$z = f(x) = x^n$$

$$(g \circ f)'(x) = \underbrace{(1 - \tanh^2(x^n))}_{dg/df} \underbrace{nx^{n-1}}_{df/dx}$$

# Multivariate Differentiation $f : \mathbb{R}^N \rightarrow \mathbb{R}$

$$y = f(\mathbf{x}), \quad \mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} \in \mathbb{R}^N$$

- ▶ **Partial derivative** (change one coordinate at a time):

$$\frac{\partial f}{\partial x_i} = \lim_{h \rightarrow 0} \frac{f(x_1, \dots, x_{i-1}, x_i + h, x_{i+1}, \dots, x_N) - f(\mathbf{x})}{h}$$

# Multivariate Differentiation $f : \mathbb{R}^N \rightarrow \mathbb{R}$

$$y = f(\mathbf{x}), \quad \mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} \in \mathbb{R}^N$$

- ▶ **Partial derivative** (change one coordinate at a time):

$$\frac{\partial f}{\partial x_i} = \lim_{h \rightarrow 0} \frac{f(x_1, \dots, x_{i-1}, x_i + h, x_{i+1}, \dots, x_N) - f(\mathbf{x})}{h}$$

- ▶ **Jacobian** vector (**gradient**) collects all partial derivatives:

$$\frac{df}{d\mathbf{x}} = \left[ \frac{\partial f}{\partial x_1} \quad \dots \quad \frac{\partial f}{\partial x_N} \right] \in \mathbb{R}^{1 \times N}$$

Note: This is a row vector.



# Example: Multivariate Differentiation

## Beginner

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

$$f(x_1, x_2) = x_1^2 x_2 + x_1 x_2^3 \in \mathbb{R}$$

## Advanced

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

$$f(x_1, x_2) = (x_1 + 2x_2^3)^2 \in \mathbb{R}$$

Partial derivatives?

**Work it out with your neighbors**

# Example: Multivariate Differentiation

## Beginner

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

$$f(x_1, x_2) = x_1^2 x_2 + x_1 x_2^3 \in \mathbb{R}$$

## Advanced

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

$$f(x_1, x_2) = (x_1 + 2x_2^3)^2 \in \mathbb{R}$$

### Partial derivatives

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

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

$$\frac{\partial f(x_1, x_2)}{\partial x_1} = 2(x_1 + 2x_2^3) \underbrace{\frac{\partial}{\partial x_1}(x_1 + 2x_2^3)}_{(1)}$$

$$\frac{\partial f(x_1, x_2)}{\partial x_2} = 2(x_1 + 2x_2^3) \underbrace{(6x_2^2)}_{\frac{\partial}{\partial x_2}(x_1 + 2x_2^3)}$$

# Example: Multivariate Differentiation

## Beginner

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

$$f(x_1, x_2) = x_1^2 x_2 + x_1 x_2^3 \in \mathbb{R}$$

## Advanced

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

$$f(x_1, x_2) = (x_1 + 2x_2^3)^2 \in \mathbb{R}$$

### Partial derivatives

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

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

$$\frac{\partial f(x_1, x_2)}{\partial x_1} = 2(x_1 + 2x_2^3) \underbrace{\frac{\partial}{\partial x_1}(x_1 + 2x_2^3)}_{(1)}$$

$$\frac{\partial f(x_1, x_2)}{\partial x_2} = 2(x_1 + 2x_2^3) \underbrace{(6x_2^2)}_{\frac{\partial}{\partial x_2}(x_1 + 2x_2^3)}$$

Gradient  $\frac{df}{dx} = \left[ \frac{\partial f(x_1, x_2)}{\partial x_1} \quad \frac{\partial f(x_1, x_2)}{\partial x_2} \right] \in \mathbb{R}^{1 \times 2}$

$$\frac{df}{dx} = [2x_1 x_2 + x_2^3 \quad x_1^2 + 3x_1 x_2^2]$$

$$\frac{df}{dx} = [2(x_1 + 2x_2^3) \quad 12(x_1 + 2x_2^3)x_2^2]$$

## Example: Multivariate Chain Rule

- ▶ Consider the function

$$L(\mathbf{e}) = \frac{1}{2} \|\mathbf{e}\|^2 = \frac{1}{2} \mathbf{e}^\top \mathbf{e}$$

$$\mathbf{e} = \mathbf{y} - \mathbf{A}\mathbf{x}, \quad \mathbf{x} \in \mathbb{R}^N, \mathbf{A} \in \mathbb{R}^{M \times N}, \mathbf{e}, \mathbf{y} \in \mathbb{R}^M$$

- ▶ Compute the gradient  $\frac{dL}{dx}$ . What is the dimension/size of  $\frac{dL}{dx}$ ?

**Work it out with your neighbors**

## Example: Multivariate Chain Rule

- ▶ Consider the function

$$L(\mathbf{e}) = \frac{1}{2} \|\mathbf{e}\|^2 = \frac{1}{2} \mathbf{e}^\top \mathbf{e}$$

$$\mathbf{e} = \mathbf{y} - \mathbf{A}\mathbf{x}, \quad \mathbf{x} \in \mathbb{R}^N, \mathbf{A} \in \mathbb{R}^{M \times N}, \mathbf{e}, \mathbf{y} \in \mathbb{R}^M$$

- ▶ Compute the gradient  $\frac{dL}{d\mathbf{x}}$ . What is the dimension/size of  $\frac{dL}{d\mathbf{x}}$ ?

$$\frac{dL}{d\mathbf{x}} = \frac{\partial L}{\partial \mathbf{e}} \frac{\partial \mathbf{e}}{\partial \mathbf{x}}$$
$$\frac{\partial L}{\partial \mathbf{e}} = \mathbf{e}^\top \in \mathbb{R}^{1 \times M} \tag{1}$$

$$\frac{\partial \mathbf{e}}{\partial \mathbf{x}} = -\mathbf{A} \in \mathbb{R}^{M \times N} \tag{2}$$

$$\implies \frac{dL}{d\mathbf{x}} = \mathbf{e}^\top (-\mathbf{A}) = -(\mathbf{y} - \mathbf{A}\mathbf{x})^\top \mathbf{A} \in \mathbb{R}^{1 \times N}$$

# Vector Field Differentiation $f : \mathbb{R}^N \rightarrow \mathbb{R}^M$

$$\mathbf{y} = f(\mathbf{x}) \in \mathbb{R}^M, \quad \mathbf{x} \in \mathbb{R}^N$$
$$\begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} f_1(x_1, \dots, x_N) \\ \vdots \\ f_M(x_1, \dots, x_N) \end{bmatrix}$$

# Vector Field Differentiation $f : \mathbb{R}^N \rightarrow \mathbb{R}^M$

$$\mathbf{y} = f(\mathbf{x}) \in \mathbb{R}^M, \quad \mathbf{x} \in \mathbb{R}^N$$
$$\begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} f_1(x_1, \dots, x_N) \\ \vdots \\ f_M(x_1, \dots, x_N) \end{bmatrix}$$

- ▶ **Jacobian** matrix (collection of all partial derivatives)

$$\begin{bmatrix} \frac{dy_1}{dx} \\ \vdots \\ \frac{dy_M}{dx} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_N} \\ \vdots & & \vdots \\ \frac{\partial f_M}{\partial x_1} & \cdots & \frac{\partial f_M}{\partial x_N} \end{bmatrix} \in \mathbb{R}^{M \times N}$$

## Example: Vector Field Differentiation

$$f(\mathbf{x}) = \mathbf{Ax}, \quad f(\mathbf{x}) \in \mathbb{R}^M, \quad \mathbf{A} \in \mathbb{R}^{M \times N}, \quad \mathbf{x} \in \mathbb{R}^N$$

$$\begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1N}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{MN}x_N \end{bmatrix}$$

- Compute the gradient  $\frac{df}{dx}$



## Example: Vector Field Differentiation

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

$$\begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1N}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{MN}x_N \end{bmatrix}$$

► Compute the gradient  $\frac{df}{dx}$

► Gradient:

$$f_i(\mathbf{x}) = \sum_{j=1}^N A_{ij}x_j \quad \implies \quad \frac{\partial f_i}{\partial x_j} = A_{ij}$$

## Example: Vector Field Differentiation

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

$$\begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1N}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{MN}x_N \end{bmatrix}$$

► Compute the gradient  $\frac{df}{dx}$

► Gradient:

$$f_i(\mathbf{x}) = \sum_{j=1}^N A_{ij}x_j \quad \implies \quad \frac{\partial f_i}{\partial x_j} = A_{ij}$$

## Example: Vector Field Differentiation

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

$$\begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(x) \\ \vdots \\ f_M(x) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1N}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{MN}x_N \end{bmatrix}$$

► Compute the gradient  $\frac{df}{dx}$

► Gradient:

$$f_i(x) = \sum_{j=1}^N A_{ij}x_j \quad \implies \quad \frac{\partial f_i}{\partial x_j} = A_{ij}$$
$$\implies \frac{df}{dx} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_N} \\ \vdots & & \vdots \\ \frac{\partial f_M}{\partial x_1} & \cdots & \frac{\partial f_M}{\partial x_N} \end{bmatrix} = \begin{bmatrix} A_{11} & \cdots & A_{1N} \\ \vdots & & \vdots \\ A_{M1} & \cdots & A_{MN} \end{bmatrix} = A \in \mathbb{R}^{M \times N}$$

# Dimensionality of the Gradient

In general: A function  $f : \mathbb{R}^N \rightarrow \mathbb{R}^M$  has a gradient that is an  $M \times N$ -matrix with

$$\frac{df}{dx} \in \mathbb{R}^{M \times N}, \quad df[m, n] = \frac{\partial f_m}{\partial x_n}$$

Gradient dimension: # target dimensions  $\times$  # input dimensions

# Chain Rule

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

## Example: Chain Rule

- ▶ Consider  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $\mathbf{x} : \mathbb{R} \rightarrow \mathbb{R}^2$

$$f(\mathbf{x}) = f(x_1, x_2) = x_1^2 + 2x_2,$$

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

## Example: Chain Rule

- ▶ Consider  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $\mathbf{x} : \mathbb{R} \rightarrow \mathbb{R}^2$

$$f(\mathbf{x}) = f(x_1, x_2) = x_1^2 + 2x_2,$$

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

- ▶ What are the dimensions of  $\frac{df}{dx}$  and  $\frac{dx}{dt}$ ?

**Work it out with your neighbors**

## Example: Chain Rule

- ▶ Consider  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $\mathbf{x} : \mathbb{R} \rightarrow \mathbb{R}^2$

$$f(\mathbf{x}) = f(x_1, x_2) = x_1^2 + 2x_2,$$

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

- ▶ What are the dimensions of  $\frac{df}{dx}$  and  $\frac{dx}{dt}$ ?  
 $1 \times 2$  and  $2 \times 1$
- ▶ Compute the gradient  $\frac{df}{dt}$  using the chain rule:



## Example: Chain Rule

- ▶ Consider  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $\mathbf{x} : \mathbb{R} \rightarrow \mathbb{R}^2$

$$f(\mathbf{x}) = f(x_1, x_2) = x_1^2 + 2x_2,$$

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

- ▶ What are the dimensions of  $\frac{df}{d\mathbf{x}}$  and  $\frac{d\mathbf{x}}{dt}$ ?

$$1 \times 2 \text{ and } 2 \times 1$$

- ▶ Compute the gradient  $\frac{df}{dt}$  using the chain rule:

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

# Derivatives with Respect to Matrices

- ▶ Recall: A function  $f : \mathbb{R}^N \rightarrow \mathbb{R}^M$  has a gradient that is an  $M \times N$ -matrix with

$$\frac{df}{dx} \in \mathbb{R}^{M \times N}, \quad df[m, n] = \frac{\partial f_m}{\partial x_n}$$

Gradient dimension: # target dimensions  $\times$  # input dimensions

# Derivatives with Respect to Matrices

- ▶ Recall: A function  $f : \mathbb{R}^N \rightarrow \mathbb{R}^M$  has a gradient that is an  $M \times N$ -matrix with

$$\frac{df}{dx} \in \mathbb{R}^{M \times N}, \quad df[m, n] = \frac{\partial f_m}{\partial x_n}$$

Gradient dimension: # target dimensions  $\times$  # input dimensions

- ▶ This generalizes to when the inputs ( $N$ ) or targets ( $M$ ) are **matrices**

# Derivatives with Respect to Matrices

- ▶ Recall: A function  $f : \mathbb{R}^N \rightarrow \mathbb{R}^M$  has a gradient that is an  $M \times N$ -matrix with

$$\frac{df}{dx} \in \mathbb{R}^{M \times N}, \quad df[m, n] = \frac{\partial f_m}{\partial x_n}$$

Gradient dimension: # target dimensions  $\times$  # input dimensions

- ▶ This generalizes to when the inputs ( $N$ ) or targets ( $M$ ) are **matrices**
- ▶ Function  $f : \mathbb{R}^{M \times N} \rightarrow \mathbb{R}^{P \times Q}$ , has a gradient that is a  $(P \times Q) \times (M \times N)$  object (tensor)

$$\frac{df}{dX} \in \mathbb{R}^{(P \times Q) \times (M \times N)}, \quad df[p, q, m, n] = \frac{\partial f_{pq}}{\partial X_{mn}}$$

## Example 1: Derivatives with Respect to Matrices

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

$$\begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1N}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{MN}x_N \end{bmatrix}$$

$$\frac{d\mathbf{f}}{d\mathbf{A}} \in \mathbb{R}^?$$

## Example 1: Derivatives with Respect to Matrices

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

$$\begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(x) \\ \vdots \\ f_M(x) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1N}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{MN}x_N \end{bmatrix}$$

$$\frac{df}{dA} \in \mathbb{R}^{\# \text{ target dim} \times \# \text{ input dim}} = M \times (M \times N)$$

$$\frac{df}{dA} = \begin{bmatrix} \frac{\partial f_1}{\partial A} \\ \vdots \\ \frac{\partial f_M}{\partial A} \end{bmatrix}, \quad \frac{\partial f_i}{\partial A} \in \mathbb{R}^{1 \times (M \times N)}$$

## Example 2: Derivatives with Respect to Matrices

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

$$\begin{bmatrix} y_1 \\ \vdots \\ y_i \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_i(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1N}x_N \\ \vdots \\ A_{i1}x_1 + A_{i2}x_2 + \cdots + A_{iN}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{MN}x_N \end{bmatrix}$$

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

$$\frac{\partial f_i}{\partial A_{i,:}} = ?$$

$$\frac{\partial f_i}{\partial A_{k \neq i,:}} = ?$$

$$\frac{\partial f_i}{\partial \mathbf{A}} = ?$$

## Example 2: Derivatives with Respect to Matrices

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

$$\begin{bmatrix} y_1 \\ \vdots \\ y_i \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_i(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1N}x_N \\ \vdots \\ A_{i1}x_1 + A_{i2}x_2 + \cdots + A_{iN}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{MN}x_N \end{bmatrix}$$

$$\frac{\partial f_i}{\partial A_{iq}} = \underbrace{x_q}_{\in \mathbb{R}} \quad \frac{\partial f_i}{\partial A_{i,:}} = ? \quad \frac{\partial f_i}{\partial A_{k \neq i,:}} = ? \quad \frac{\partial f_i}{\partial \mathbf{A}} = ?$$



## Example 2: Derivatives with Respect to Matrices

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

$$\begin{bmatrix} y_1 \\ \vdots \\ y_i \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_i(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1N}x_N \\ \vdots \\ A_{i1}x_1 + A_{i2}x_2 + \cdots + A_{iN}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{MN}x_N \end{bmatrix}$$

$$\frac{\partial f_i}{\partial A_{iq}} = \underbrace{x_q}_{\in \mathbb{R}}$$

$$\frac{\partial f_i}{\partial A_{i,:}} = \underbrace{\mathbf{x}^\top}_{\in \mathbb{R}^{1 \times 1 \times N}}$$

$$\frac{\partial f_i}{\partial A_{k \neq i,:}} = ?$$

$$\frac{\partial f_i}{\partial \mathbf{A}} = ?$$

## Example 2: Derivatives with Respect to Matrices

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

$$\begin{bmatrix} y_1 \\ \vdots \\ y_i \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_i(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1N}x_N \\ \vdots \\ A_{i1}x_1 + A_{i2}x_2 + \cdots + A_{iN}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{MN}x_N \end{bmatrix}$$

$$\frac{\partial f_i}{\partial A_{iq}} = \underbrace{x_q}_{\in \mathbb{R}} \quad \frac{\partial f_i}{\partial A_{i:}} = \underbrace{\mathbf{x}^\top}_{\in \mathbb{R}^{1 \times 1 \times N}} \quad \frac{\partial f_i}{\partial A_{k \neq i,:}} = \underbrace{\mathbf{0}^\top}_{\in \mathbb{R}^{1 \times 1 \times N}} \quad \frac{\partial f_i}{\partial \mathbf{A}} = ?$$

## Example 2: Derivatives with Respect to Matrices

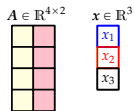
$$f_i = \sum_{j=1}^N A_{ij}x_j, \quad i = 1, \dots, M$$

$$\begin{bmatrix} y_1 \\ \vdots \\ y_i \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_i(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1N}x_N \\ \vdots \\ A_{i1}x_1 + A_{i2}x_2 + \cdots + A_{iN}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{MN}x_N \end{bmatrix}$$

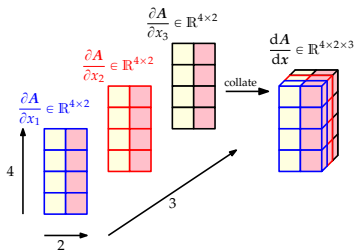
$$\frac{\partial f_i}{\partial A_{iq}} = \underbrace{x_q}_{\in \mathbb{R}} \quad \frac{\partial f_i}{\partial A_{i:}} = \underbrace{\mathbf{x}^\top}_{\in \mathbb{R}^{1 \times 1 \times N}} \quad \frac{\partial f_i}{\partial A_{k \neq i,:}} = \underbrace{\mathbf{0}^\top}_{\in \mathbb{R}^{1 \times 1 \times N}} \quad \frac{\partial f_i}{\partial \mathbf{A}} = \underbrace{\begin{bmatrix} \mathbf{0}^\top \\ \vdots \\ \mathbf{x}^\top \\ \vdots \\ \mathbf{0}^\top \end{bmatrix}}_{\in \mathbb{R}^{1 \times (M \times N)}}$$

# Gradient Computation: Two Alternatives

- ▶ Consider  $f : \mathbb{R}^3 \rightarrow \mathbb{R}^{4 \times 2}$ ,  $f(\mathbf{x}) = \mathbf{A} \in \mathbb{R}^{4 \times 2}$  where the entries  $A_{ij}$  depend on a vector  $\mathbf{x} \in \mathbb{R}^3$
- ▶ We can compute  $\frac{d\mathbf{A}(\mathbf{x})}{d\mathbf{x}} \in \mathbb{R}^{4 \times 2 \times 3}$  in two equivalent ways:

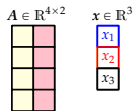


Partial derivatives:

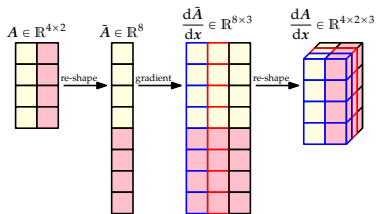
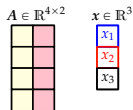
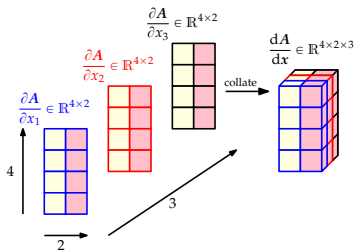


# Gradient Computation: Two Alternatives

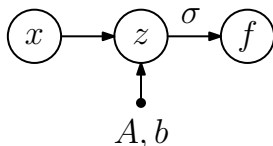
- ▶ Consider  $f : \mathbb{R}^3 \rightarrow \mathbb{R}^{4 \times 2}$ ,  $f(\mathbf{x}) = \mathbf{A} \in \mathbb{R}^{4 \times 2}$  where the entries  $A_{ij}$  depend on a vector  $\mathbf{x} \in \mathbb{R}^3$
- ▶ We can compute  $\frac{d\mathbf{A}(\mathbf{x})}{d\mathbf{x}} \in \mathbb{R}^{4 \times 2 \times 3}$  in two equivalent ways:



Partial derivatives:



# Gradients of a Single-Layer Neural Network



$$f = \tanh(\underbrace{Ax + b}_{=:z \in \mathbb{R}^M}) \in \mathbb{R}^M, \quad x \in \mathbb{R}^N, A \in \mathbb{R}^{M \times N}, b \in \mathbb{R}^M$$

# Gradients of a Single-Layer Neural Network

$$\mathbf{f} = \tanh(\underbrace{\mathbf{A}\mathbf{x} + \mathbf{b}}_{=: \mathbf{z} \in \mathbb{R}^M}) \in \mathbb{R}^M, \quad \mathbf{x} \in \mathbb{R}^N, \mathbf{A} \in \mathbb{R}^{M \times N}, \mathbf{b} \in \mathbb{R}^M$$

$$\frac{\partial \mathbf{f}}{\partial \mathbf{b}} =$$

$$\frac{\partial \mathbf{f}}{\partial \mathbf{A}} =$$

# Gradients of a Single-Layer Neural Network

$$f = \tanh(\underbrace{\mathbf{Ax} + \mathbf{b}}_{=: \mathbf{z} \in \mathbb{R}^M}) \in \mathbb{R}^M, \quad \mathbf{x} \in \mathbb{R}^N, \mathbf{A} \in \mathbb{R}^{M \times N}, \mathbf{b} \in \mathbb{R}^M$$

$$\frac{\partial f}{\partial \mathbf{b}} = \underbrace{\frac{\partial f}{\partial \mathbf{z}}}_{M \times M} \underbrace{\frac{\partial \mathbf{z}}{\partial \mathbf{b}}}_{M \times M} \in \mathbb{R}^{M \times M}$$

$$\frac{\partial f}{\partial \mathbf{A}} =$$



# Gradients of a Single-Layer Neural Network

$$f = \tanh(\underbrace{\mathbf{Ax} + \mathbf{b}}_{=: \mathbf{z} \in \mathbb{R}^M}) \in \mathbb{R}^M, \quad \mathbf{x} \in \mathbb{R}^N, \mathbf{A} \in \mathbb{R}^{M \times N}, \mathbf{b} \in \mathbb{R}^M$$

$$\frac{\partial f}{\partial \mathbf{b}} = \underbrace{\frac{\partial f}{\partial \mathbf{z}}}_{M \times M} \underbrace{\frac{\partial \mathbf{z}}{\partial \mathbf{b}}}_{M \times M} \in \mathbb{R}^{M \times M}$$

$$\frac{\partial f}{\partial \mathbf{A}} = \underbrace{\frac{\partial f}{\partial \mathbf{z}}}_{M \times M} \underbrace{\frac{\partial \mathbf{z}}{\partial \mathbf{A}}}_{M \times (M \times N)} \in \mathbb{R}^{M \times (M \times N)}$$

# Gradients of a Single-Layer Neural Network

$$f = \tanh(\underbrace{\mathbf{Ax} + \mathbf{b}}_{=: \mathbf{z} \in \mathbb{R}^M}) \in \mathbb{R}^M, \quad \mathbf{x} \in \mathbb{R}^N, \mathbf{A} \in \mathbb{R}^{M \times N}, \mathbf{b} \in \mathbb{R}^M$$

$$\frac{\partial f}{\partial \mathbf{b}} = \underbrace{\frac{\partial f}{\partial \mathbf{z}}}_{M \times M} \underbrace{\frac{\partial \mathbf{z}}{\partial \mathbf{b}}}_{M \times M} \in \mathbb{R}^{M \times M}$$

$$\frac{\partial f}{\partial \mathbf{A}} = \underbrace{\frac{\partial f}{\partial \mathbf{z}}}_{M \times M} \underbrace{\frac{\partial \mathbf{z}}{\partial \mathbf{A}}}_{M \times (M \times N)} \in \mathbb{R}^{M \times (M \times N)}$$

$$\frac{\partial f}{\partial \mathbf{z}} = \underbrace{\text{diag}(1 - \tanh^2(\mathbf{z}))}_{\in \mathbb{R}^{M \times M}}$$

# Gradients of a Single-Layer Neural Network

$$f = \tanh(\underbrace{\mathbf{Ax} + \mathbf{b}}_{=: \mathbf{z} \in \mathbb{R}^M}) \in \mathbb{R}^M, \quad \mathbf{x} \in \mathbb{R}^N, \mathbf{A} \in \mathbb{R}^{M \times N}, \mathbf{b} \in \mathbb{R}^M$$

$$\frac{\partial f}{\partial \mathbf{b}} = \underbrace{\frac{\partial f}{\partial \mathbf{z}}}_{M \times M} \underbrace{\frac{\partial \mathbf{z}}{\partial \mathbf{b}}}_{M \times M} \in \mathbb{R}^{M \times M}$$

$$\frac{\partial f}{\partial \mathbf{A}} = \underbrace{\frac{\partial f}{\partial \mathbf{z}}}_{M \times M} \underbrace{\frac{\partial \mathbf{z}}{\partial \mathbf{A}}}_{M \times (M \times N)} \in \mathbb{R}^{M \times (M \times N)}$$

$$\frac{\partial f}{\partial \mathbf{z}} = \underbrace{\text{diag}(1 - \tanh^2(\mathbf{z}))}_{\in \mathbb{R}^{M \times M}} \quad \frac{\partial \mathbf{z}}{\partial \mathbf{b}} = \underbrace{\mathbf{I}}_{\in \mathbb{R}^{M \times M}}$$

# Gradients of a Single-Layer Neural Network

$$f = \tanh(\underbrace{\mathbf{Ax} + \mathbf{b}}_{=: \mathbf{z} \in \mathbb{R}^M}) \in \mathbb{R}^M, \quad \mathbf{x} \in \mathbb{R}^N, \mathbf{A} \in \mathbb{R}^{M \times N}, \mathbf{b} \in \mathbb{R}^M$$

$$\frac{\partial f}{\partial \mathbf{b}} = \underbrace{\frac{\partial f}{\partial \mathbf{z}}}_{M \times M} \underbrace{\frac{\partial \mathbf{z}}{\partial \mathbf{b}}}_{M \times M} \in \mathbb{R}^{M \times M}$$

$$\frac{\partial f}{\partial \mathbf{A}} = \underbrace{\frac{\partial f}{\partial \mathbf{z}}}_{M \times M} \underbrace{\frac{\partial \mathbf{z}}{\partial \mathbf{A}}}_{M \times (M \times N)} \in \mathbb{R}^{M \times (M \times N)}$$

$$\frac{\partial f}{\partial \mathbf{z}} = \underbrace{\text{diag}(1 - \tanh^2(\mathbf{z}))}_{\in \mathbb{R}^{M \times M}} \quad \frac{\partial \mathbf{z}}{\partial \mathbf{b}} = \underbrace{\mathbf{I}}_{\in \mathbb{R}^{M \times M}} \quad \frac{\partial \mathbf{z}}{\partial \mathbf{A}} = \underbrace{\begin{bmatrix} \mathbf{x}^\top & \cdot & \mathbf{0}^\top & \cdot & \mathbf{0}^\top \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \mathbf{0}^\top & \cdot & \mathbf{x}^\top & \cdot & \mathbf{0}^\top \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \mathbf{0}^\top & \cdot & \mathbf{0}^\top & \cdot & \mathbf{x}^\top \end{bmatrix}}_{\in \mathbb{R}^{M \times (M \times N)}}$$

# Gradients of a Single-Layer Neural Network

$$f = \tanh(\underbrace{\mathbf{Ax} + \mathbf{b}}_{=: \mathbf{z} \in \mathbb{R}^M}) \in \mathbb{R}^M, \quad \mathbf{x} \in \mathbb{R}^N, \mathbf{A} \in \mathbb{R}^{M \times N}, \mathbf{b} \in \mathbb{R}^M$$

$$\frac{\partial f}{\partial \mathbf{b}} = \underbrace{\frac{\partial f}{\partial \mathbf{z}}}_{M \times M} \underbrace{\frac{\partial \mathbf{z}}{\partial \mathbf{b}}}_{M \times M} \in \mathbb{R}^{M \times M}$$

$$\frac{\partial f}{\partial \mathbf{b}}[i, j] = \sum_{l=1}^M \frac{\partial f}{\partial z}[i, l] \frac{\partial z}{\partial \mathbf{b}}[l, j]$$

$$\frac{\partial f}{\partial \mathbf{A}} = \underbrace{\frac{\partial f}{\partial \mathbf{z}}}_{M \times M} \underbrace{\frac{\partial \mathbf{z}}{\partial \mathbf{A}}}_{M \times (M \times N)} \in \mathbb{R}^{M \times (M \times N)}$$

$$\frac{\partial f}{\partial \mathbf{z}} = \underbrace{\text{diag}(1 - \tanh^2(\mathbf{z}))}_{\in \mathbb{R}^{M \times M}} \quad \frac{\partial \mathbf{z}}{\partial \mathbf{b}} = \underbrace{\mathbf{I}}_{\in \mathbb{R}^{M \times M}} \quad \frac{\partial \mathbf{z}}{\partial \mathbf{A}} = \underbrace{\begin{bmatrix} \mathbf{x}^\top & \cdot & \mathbf{0}^\top & \cdot & \mathbf{0}^\top \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \mathbf{0}^\top & \cdot & \mathbf{x}^\top & \cdot & \mathbf{0}^\top \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \mathbf{0}^\top & \cdot & \mathbf{0}^\top & \cdot & \mathbf{x}^\top \end{bmatrix}}_{\in \mathbb{R}^{M \times (M \times N)}}$$

# Gradients of a Single-Layer Neural Network

$$f = \tanh(\underbrace{\mathbf{Ax} + \mathbf{b}}_{=: \mathbf{z} \in \mathbb{R}^M}) \in \mathbb{R}^M, \quad \mathbf{x} \in \mathbb{R}^N, \mathbf{A} \in \mathbb{R}^{M \times N}, \mathbf{b} \in \mathbb{R}^M$$

$$\frac{\partial f}{\partial \mathbf{b}} = \underbrace{\frac{\partial f}{\partial \mathbf{z}}}_{M \times M} \underbrace{\frac{\partial \mathbf{z}}{\partial \mathbf{b}}}_{M \times M} \in \mathbb{R}^{M \times M}$$

$$\frac{\partial f}{\partial \mathbf{b}}[i, j] = \sum_{l=1}^M \frac{\partial f}{\partial z}[i, l] \frac{\partial z}{\partial \mathbf{b}}[l, j]$$

$$\frac{\partial f}{\partial \mathbf{A}} = \underbrace{\frac{\partial f}{\partial \mathbf{z}}}_{M \times M} \underbrace{\frac{\partial \mathbf{z}}{\partial \mathbf{A}}}_{M \times (M \times N)} \in \mathbb{R}^{M \times (M \times N)}$$

$$\frac{\partial f}{\partial \mathbf{A}}[i, j, k] = \sum_{l=1}^M \frac{\partial f}{\partial z}[i, l] \frac{\partial z}{\partial \mathbf{A}}[l, j, k]$$

$$\frac{\partial f}{\partial \mathbf{z}} = \underbrace{\text{diag}(1 - \tanh^2(\mathbf{z}))}_{\in \mathbb{R}^{M \times M}}$$

$$\frac{\partial \mathbf{z}}{\partial \mathbf{b}} = \underbrace{\mathbf{I}}_{\in \mathbb{R}^{M \times M}}$$

$$\frac{\partial \mathbf{z}}{\partial \mathbf{A}}$$

$$= \underbrace{\begin{bmatrix} \mathbf{x}^\top & \cdot & \mathbf{0}^\top & \cdot & \mathbf{0}^\top \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \mathbf{0}^\top & \cdot & \mathbf{x}^\top & \cdot & \mathbf{0}^\top \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \mathbf{0}^\top & \cdot & \mathbf{0}^\top & \cdot & \mathbf{x}^\top \end{bmatrix}}_{\in \mathbb{R}^{M \times (M \times N)}}$$

# Putting Things Together

- ▶ Inputs  $\mathbf{x} \in \mathbb{R}^N$

# Putting Things Together

- ▶ Inputs  $\mathbf{x} \in \mathbb{R}^N$
- ▶ Observed outputs  $\mathbf{y} = f_{\theta}(\mathbf{z}) + \boldsymbol{\epsilon} \in \mathbb{R}^M, \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$



# Putting Things Together

- ▶ Inputs  $\mathbf{x} \in \mathbb{R}^N$
- ▶ Observed outputs  $\mathbf{y} = f_{\boldsymbol{\theta}}(\mathbf{z}) + \boldsymbol{\epsilon} \in \mathbb{R}^M, \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$
- ▶ Train single-layer neural network with

$$f_{\boldsymbol{\theta}}(\mathbf{z}) = \tanh(\mathbf{z}) \in \mathbb{R}^M, \quad \mathbf{z} = \mathbf{A}\mathbf{x} + \mathbf{b} \in \mathbb{R}^M, \quad \boldsymbol{\theta} = \{\mathbf{A}, \mathbf{b}\}$$

# Putting Things Together

- ▶ Inputs  $\mathbf{x} \in \mathbb{R}^N$
- ▶ Observed outputs  $\mathbf{y} = f_{\boldsymbol{\theta}}(\mathbf{z}) + \boldsymbol{\epsilon} \in \mathbb{R}^M, \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$
- ▶ Train single-layer neural network with

$$f_{\boldsymbol{\theta}}(\mathbf{z}) = \tanh(\mathbf{z}) \in \mathbb{R}^M, \quad \mathbf{z} = \mathbf{A}\mathbf{x} + \mathbf{b} \in \mathbb{R}^M, \quad \boldsymbol{\theta} = \{\mathbf{A}, \mathbf{b}\}$$

- ▶ Find  $\mathbf{A}, \mathbf{b}$ , such that the squared loss

$$L(\boldsymbol{\theta}) = \frac{1}{2} \|\mathbf{e}\|^2 \in \mathbb{R}, \quad \mathbf{e} = \mathbf{y} - f_{\boldsymbol{\theta}}(\mathbf{z}) \in \mathbb{R}^M$$

is minimized

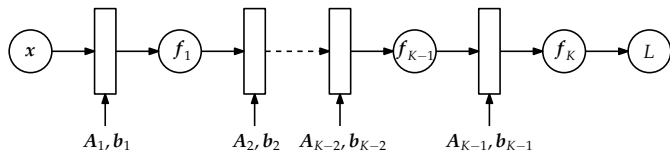
# Putting Things Together

Partial derivatives:

$$\begin{aligned}\frac{\partial L}{\partial \mathbf{A}} &= \frac{\partial L}{\partial \mathbf{e}} \frac{\partial \mathbf{e}}{\partial \mathbf{f}} \frac{\partial \mathbf{f}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{A}} \\ \frac{\partial L}{\partial \mathbf{b}} &= \frac{\partial L}{\partial \mathbf{e}} \frac{\partial \mathbf{e}}{\partial \mathbf{f}} \frac{\partial \mathbf{f}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{b}}\end{aligned}$$

$$\begin{aligned}\frac{\partial L}{\partial \mathbf{e}} &= \underbrace{\mathbf{e}^\top}_{\in \mathbb{R}^{1 \times M}} & \frac{\partial \mathbf{e}}{\partial \mathbf{f}} &= \underbrace{-\mathbf{I}}_{\in \mathbb{R}^{M \times M}} & \frac{\partial \mathbf{f}}{\partial \mathbf{z}} &= \underbrace{\text{diag}(1 - \tanh^2(\mathbf{z}))}_{\in \mathbb{R}^{M \times M}} \\ \frac{\partial \mathbf{z}}{\partial \mathbf{A}} &= \underbrace{\begin{bmatrix} \mathbf{x}^\top & \cdot & \mathbf{0}^\top & \cdot & \mathbf{0}^\top \\ \cdot & & \cdot & & \cdot \\ \mathbf{0}^\top & \cdot & \mathbf{x}^\top & \cdot & \mathbf{0}^\top \\ \cdot & & \cdot & & \cdot \\ \mathbf{0}^\top & \cdot & \mathbf{0}^\top & \cdot & \mathbf{x}^\top \end{bmatrix}}_{\in \mathbb{R}^{M \times (M \times N)}} & \frac{\partial \mathbf{z}}{\partial \mathbf{b}} &= \underbrace{\mathbf{I}}_{\in \mathbb{R}^{M \times M}}\end{aligned}$$

# Gradients of a Multi-Layer Neural Network

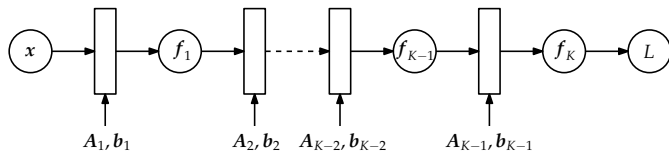


- ▶ Inputs  $x$ , observed outputs  $y$
- ▶ Train multi-layer neural network with

$$f_0 = x$$

$$f_i = \sigma_i(A_{i-1}f_{i-1} + b_{i-1}), \quad i = 1, \dots, K$$

# Gradients of a Multi-Layer Neural Network



- ▶ Inputs  $x$ , observed outputs  $y$
- ▶ Train multi-layer neural network with

$$f_0 = x$$

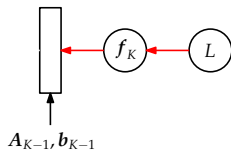
$$f_i = \sigma_i(A_{i-1}f_{i-1} + b_{i-1}), \quad i = 1, \dots, K$$

- ▶ Find  $A_j, b_j$  for  $j = 0, \dots, K-1$ , such that the squared loss

$$L(\theta) = \|\mathbf{y} - \mathbf{f}_{K,\theta}(x)\|^2$$

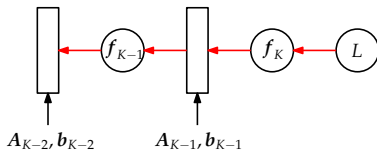
is minimized, where  $\theta = \{A_j, b_j\}$ ,  $j = 0, \dots, K-1$

# Gradients of a Multi-Layer Neural Network



$$\frac{\partial L}{\partial \boldsymbol{\theta}_{K-1}} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial \boldsymbol{\theta}_{K-1}}$$

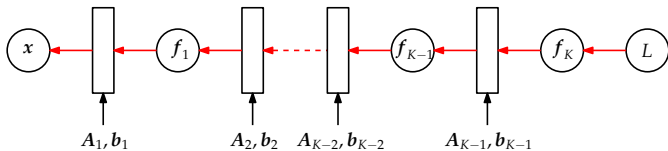
# Gradients of a Multi-Layer Neural Network



$$\frac{\partial L}{\partial \theta_{K-1}} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial \theta_{K-1}}$$

$$\frac{\partial L}{\partial \theta_{K-2}} = \frac{\partial L}{\partial f_K} \boxed{\frac{\partial f_K}{\partial f_{K-1}} \frac{\partial f_{K-1}}{\partial \theta_{K-2}}}$$

# Gradients of a Multi-Layer Neural Network



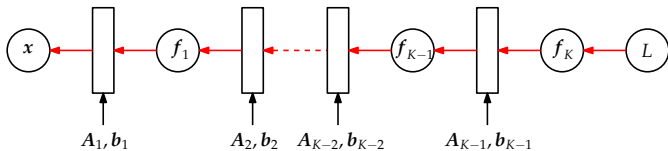
$$\frac{\partial L}{\partial \theta_{K-1}} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial \theta_{K-1}}$$

$$\frac{\partial L}{\partial \theta_{K-2}} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial f_{K-1}} \frac{\partial f_{K-1}}{\partial \theta_{K-2}}$$

$$\frac{\partial L}{\partial \theta_{K-3}} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial f_{K-1}} \frac{\partial f_{K-1}}{\partial f_{K-2}} \frac{\partial f_{K-2}}{\partial \theta_{K-3}}$$



# Gradients of a Multi-Layer Neural Network



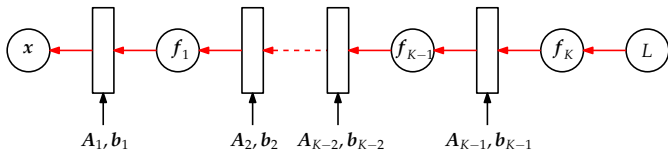
$$\frac{\partial L}{\partial \theta_{K-1}} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial \theta_{K-1}}$$

$$\frac{\partial L}{\partial \theta_{K-2}} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial f_{K-1}} \frac{\partial f_{K-1}}{\partial \theta_{K-2}}$$

$$\frac{\partial L}{\partial \theta_{K-3}} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial f_{K-1}} \frac{\partial f_{K-1}}{\partial f_{K-2}} \frac{\partial f_{K-2}}{\partial \theta_{K-3}}$$

$$\frac{\partial L}{\partial \theta_i} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial f_{K-1}} \dots \frac{\partial f_{i+2}}{\partial f_{i+1}} \frac{\partial f_{i+1}}{\partial \theta_i}$$

# Gradients of a Multi-Layer Neural Network



$$\frac{\partial L}{\partial \theta_{K-1}} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial \theta_{K-1}}$$

$$\frac{\partial L}{\partial \theta_{K-2}} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial f_{K-1}} \frac{\partial f_{K-1}}{\partial \theta_{K-2}}$$

$$\frac{\partial L}{\partial \theta_{K-3}} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial f_{K-1}} \frac{\partial f_{K-1}}{\partial f_{K-2}} \frac{\partial f_{K-2}}{\partial \theta_{K-3}}$$

$$\frac{\partial L}{\partial \theta_i} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial f_{K-1}} \dots \frac{\partial f_{i+2}}{\partial f_{i+1}} \frac{\partial f_{i+1}}{\partial \theta_i}$$

►► Intermediate derivatives are stored during the forward pass

## Example: Linear Regression with Neural Networks

- ▶ Linear regression with a neural network parametrized by  $\theta$ ,  $f_\theta$ :

$$y = f_\theta(\mathbf{x}) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_\epsilon^2)$$

## Example: Linear Regression with Neural Networks

- ▶ Linear regression with a neural network parametrized by  $\theta$ ,  $f_\theta$ :

$$y = f_\theta(\mathbf{x}) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_\epsilon^2)$$

- ▶ Given inputs  $\mathbf{x}_n$  and corresponding (noisy) observations  $y_n$ ,  $n = 1, \dots, N$ , find parameters  $\theta^*$  that minimize the squared loss

$$L(\theta) = \sum_{n=1}^N (y_n - f_\theta(\mathbf{x}_n))^2 = \|\mathbf{y} - \mathbf{f}(\mathbf{X})\|^2$$

# Training Neural Networks as Maximum Likelihood Estimation

- ▶ Training a neural network in the above way corresponds to **maximum likelihood estimation**:
  - ▶ If  $\mathbf{y} = NN(\mathbf{x}, \boldsymbol{\theta}) + \boldsymbol{\epsilon}$ ,  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  then the **log-likelihood** is

$$\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = -\frac{1}{2}\|\mathbf{y} - NN(\mathbf{x}, \boldsymbol{\theta})\|^2$$

# Training Neural Networks as Maximum Likelihood Estimation

- ▶ Training a neural network in the above way corresponds to **maximum likelihood estimation**:

- ▶ If  $\mathbf{y} = NN(\mathbf{x}, \boldsymbol{\theta}) + \boldsymbol{\epsilon}$ ,  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  then the **log-likelihood** is

$$\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = -\frac{1}{2}\|\mathbf{y} - NN(\mathbf{x}, \boldsymbol{\theta})\|^2$$

- ▶ Find  $\boldsymbol{\theta}^*$  by **minimizing the negative log-likelihood**:

$$\begin{aligned}\boldsymbol{\theta}^* &= \arg \min_{\boldsymbol{\theta}} -\log p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) \\ &= \arg \min_{\boldsymbol{\theta}} \frac{1}{2}\|\mathbf{y} - NN(\mathbf{x}, \boldsymbol{\theta})\|^2 \\ &= \arg \min_{\boldsymbol{\theta}} L(\boldsymbol{\theta})\end{aligned}$$

# Training Neural Networks as Maximum Likelihood Estimation

- ▶ Training a neural network in the above way corresponds to **maximum likelihood estimation**:

- ▶ If  $\mathbf{y} = NN(\mathbf{x}, \boldsymbol{\theta}) + \boldsymbol{\epsilon}$ ,  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  then the **log-likelihood** is

$$\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = -\frac{1}{2}\|\mathbf{y} - NN(\mathbf{x}, \boldsymbol{\theta})\|^2$$

- ▶ Find  $\boldsymbol{\theta}^*$  by **minimizing the negative log-likelihood**:

$$\begin{aligned}\boldsymbol{\theta}^* &= \arg \min_{\boldsymbol{\theta}} -\log p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) \\ &= \arg \min_{\boldsymbol{\theta}} \frac{1}{2}\|\mathbf{y} - NN(\mathbf{x}, \boldsymbol{\theta})\|^2 \\ &= \arg \min_{\boldsymbol{\theta}} L(\boldsymbol{\theta})\end{aligned}$$

- ▶ Maximum likelihood estimation can lead to **overfitting** (interpret noise as signal)

## Example: Linear Regression (1)

- ▶ Linear regression with a polynomial of order  $M$ :

$$y = f(x, \boldsymbol{\theta}) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_\epsilon^2)$$

$$f(x, \boldsymbol{\theta}) = \theta_0 + \theta_1 x + \theta_2 x^2 + \cdots + \theta_M x^M = \sum_{i=0}^M \theta_i x^i$$



## Example: Linear Regression (1)

- ▶ Linear regression with a polynomial of order  $M$ :

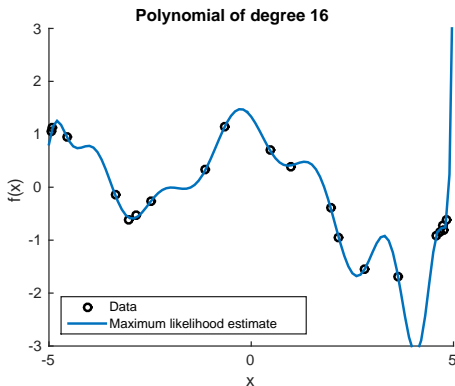
$$y = f(x, \boldsymbol{\theta}) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_\epsilon^2)$$

$$f(x, \boldsymbol{\theta}) = \theta_0 + \theta_1 x + \theta_2 x^2 + \cdots + \theta_M x^M = \sum_{i=0}^M \theta_i x^i$$

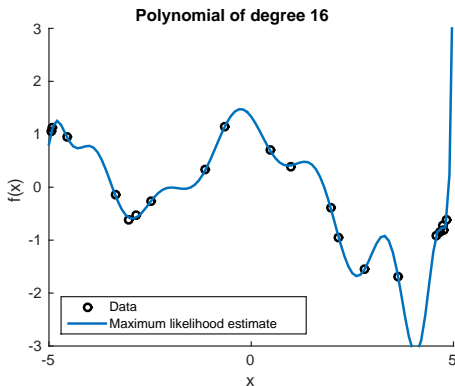
- ▶ Given inputs  $x_i$  and corresponding (noisy) observations  $y_i$ ,  $i = 1, \dots, N$ , find parameters  $\boldsymbol{\theta} = [\theta_0, \dots, \theta_M]^\top$ , that minimize the squared loss (equivalently: maximize the likelihood)

$$L(\boldsymbol{\theta}) = \sum_{i=1}^N (y_i - f(x_i, \boldsymbol{\theta}))^2$$

## Example: Linear Regression (2)



## Example: Linear Regression (2)

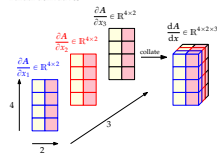


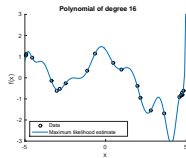
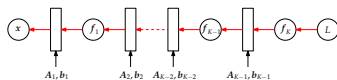
- ▶ Regularization, model selection etc. can address overfitting
- ▶ Alternative approach based on integration

# Summary

$$A \in \mathbb{R}^{4 \times 2} \quad x \in \mathbb{R}^3$$


Partial derivatives:

$$\frac{\partial A}{\partial x_1} \in \mathbb{R}^{4 \times 2} \quad \frac{\partial A}{\partial x_2} \in \mathbb{R}^{4 \times 2} \quad \frac{\partial A}{\partial x_3} \in \mathbb{R}^{4 \times 2}$$




- ▶ Vector-valued differentiation
- ▶ Chain rule
- ▶ Check the dimension of the gradients