

Vector Calculus

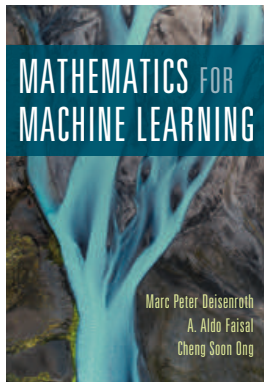
Marc Deisenroth
Centre for Artificial Intelligence
Department of Computer Science
University College London

AIMS Rwanda and AIMS Ghana
March/April, 2020

 @mpd37

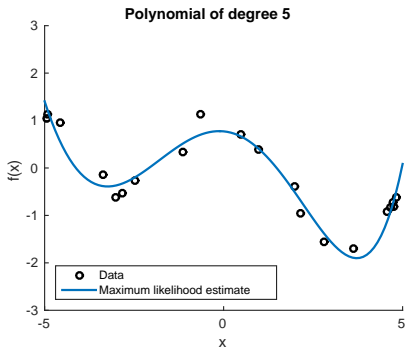
`m.deisenroth@ucl.ac.uk`

`https://deisenroth.cc`



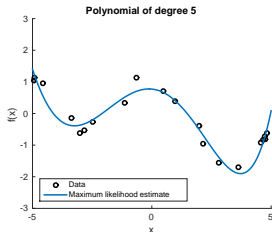
Chapter 5

<https://mml-book.com>

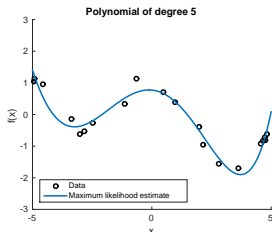


- Setting: Inputs $x \in \mathbb{R}^D$, outputs/targets $y \in \mathbb{R}$
- Goal: Find a function that models the relationship between x and y (regression/curve fitting)
- Model f that depends on parameters θ . Examples:
 - Linear model: $f(x, \theta) = \theta^\top x$, $x, \theta \in \mathbb{R}^D$
 - Neural network: $f(x, \theta) = NN(x, \theta)$

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



- Training data, e.g., N pairs (x_i, y_i) of inputs x_i and observations y_i
- **Training the model** means finding parameters θ^* , 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

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

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

2 Scalar differentiation of a vector: $f : \mathbb{R} \rightarrow \mathbb{R}^N$

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

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

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

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

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

5 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 $X \in \mathbb{R}^{M \times N}$

- 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 $f(x+h)$

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

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

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

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

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

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

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

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

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

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

■ Sum Rule

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

■ Sum Rule

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

■ Product Rule

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

■ 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

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

■ 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

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

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

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

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

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

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

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

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

$$\mathbf{f}(x) = \begin{bmatrix} f_1(x) \\ \vdots \\ f_N(x) \end{bmatrix} \in \mathbb{R}^N, \quad x \in \mathbb{R}$$

- Here, f_n are different functions, e.g.,

$$\mathbf{f}(x) = \begin{bmatrix} f_1(x) \\ f_2(x) \end{bmatrix} = \begin{bmatrix} x^2 \\ \sin(x) \end{bmatrix}$$

- Differentiation: Compute derivatives of each f_n :

$$\frac{d\mathbf{f}}{dx} = \begin{bmatrix} \frac{df_1}{dx} \\ \vdots \\ \frac{df_N}{dx} \end{bmatrix} \in \mathbb{R}^N$$

- Derivative of a (column) vector w.r.t. a scalar input is a (column) vector

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

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

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

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

- **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**.

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? Gradient?

Work it out with your neighbors

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

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

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

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

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

Beginner

Advanced

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

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

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

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

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

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

- 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

- 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 $\frac{dL}{d\mathbf{x}}$?

$$\frac{dL}{d\mathbf{x}} = \frac{\partial L}{\partial \mathbf{e}} \frac{\partial \mathbf{e}}{\partial \mathbf{x}}$$

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

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

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

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

$$\begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} f_1(x_1, \dots, x_N) \\ \vdots \\ f_M(x_1, \dots, x_N) \end{bmatrix}$$

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

$$\begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} f_1(x_1, \dots, x_N) \\ \vdots \\ f_M(x_1, \dots, x_N) \end{bmatrix}$$

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

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

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

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

- Compute the gradient $\frac{df}{dx}$
 - Dimension of $\frac{df}{dx}$:

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

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

■ Compute the gradient $\frac{d\mathbf{f}}{d\mathbf{x}}$

■ Dimension of $\frac{d\mathbf{f}}{d\mathbf{x}}$: Since $\mathbf{f} : \mathbb{R}^N \rightarrow \mathbb{R}^M$ it follows that $\frac{d\mathbf{f}}{d\mathbf{x}} \in \mathbb{R}^{M \times N}$

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

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

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

■ Dimension of $\frac{df}{dx}$: Since $\mathbf{f} : \mathbb{R}^N \rightarrow \mathbb{R}^M$ it follows that $\frac{df}{dx} \in \mathbb{R}^{M \times N}$

■ Gradient:

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

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

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

■ Compute the gradient $\frac{d\mathbf{f}}{d\mathbf{x}}$

■ Dimension of $\frac{d\mathbf{f}}{d\mathbf{x}}$: Since $\mathbf{f} : \mathbb{R}^N \rightarrow \mathbb{R}^M$ it follows that $\frac{d\mathbf{f}}{d\mathbf{x}} \in \mathbb{R}^{M \times N}$

■ Gradient:

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

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

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

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

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

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

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

- 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 $\frac{d\mathbf{x}}{dt}$?

Work it out with your neighbors

- 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:

- 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:

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

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

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

Gradient dimension: # target dimensions \times # input dimensions

- 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**

- 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

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

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

$$\frac{d\mathbf{f}}{d\mathbf{A}} \in \mathbb{R}^{\# \text{ target dim} \times \# \text{ input dim}} = M \times (M \times N)$$

$$\frac{d\mathbf{f}}{d\mathbf{A}} = \begin{bmatrix} \frac{\partial f_1}{\partial \mathbf{A}} \\ \vdots \\ \frac{\partial f_M}{\partial \mathbf{A}} \end{bmatrix}, \quad \frac{\partial f_m}{\partial \mathbf{A}} \in \mathbb{R}^{1 \times (M \times N)}$$

$$f_m = \sum_{n=1}^N A_{mn}x_n, \quad m = 1, \dots, M$$

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

$$\frac{\partial f_m}{\partial A_{mq}} = ?$$

$$\frac{\partial f_m}{\partial A_{m,:}} = ?$$

$$\frac{\partial f_m}{\partial A_{k \neq m,:}} = ?$$

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

$$f_m = \sum_{n=1}^N A_{mn}x_n, \quad m = 1, \dots, M$$

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

$$\frac{\partial f_m}{\partial A_{mq}} = \underbrace{x_q}_{\in \mathbb{R}} \quad \frac{\partial f_m}{\partial A_{m,:}} = ? \quad \frac{\partial f_m}{\partial A_{k \neq m,:}} = ? \quad \frac{\partial f_m}{\partial \mathbf{A}} = ?$$

$$f_m = \sum_{n=1}^N A_{mn}x_n, \quad m = 1, \dots, M$$

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

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

$$f_m = \sum_{n=1}^N A_{mn}x_n, \quad m = 1, \dots, M$$

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

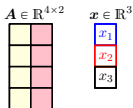
$$\frac{\partial f_m}{\partial A_{mq}} = \underbrace{x_q}_{\in \mathbb{R}} \quad \frac{\partial f_m}{\partial A_{m,:}} = \underbrace{\mathbf{x}^\top}_{\in \mathbb{R}^{1 \times 1 \times N}} \quad \frac{\partial f_m}{\partial A_{k \neq m, :}} = \underbrace{\mathbf{0}^\top}_{\in \mathbb{R}^{1 \times 1 \times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

$$f_m = \sum_{n=1}^N A_{mn}x_n, \quad m = 1, \dots, M$$

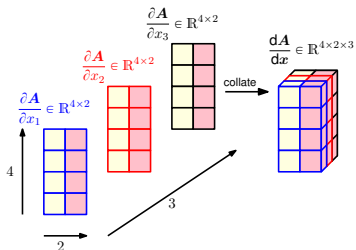
$$\begin{bmatrix} y_1 \\ \vdots \\ y_i \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_i(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1N}x_N \\ \vdots \\ A_{i1}x_1 + A_{i2}x_2 + \cdots + A_{iN}x_N \\ \vdots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

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

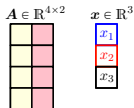
- Consider $f : \mathbb{R}^3 \rightarrow \mathbb{R}^{4 \times 2}$, $f(\mathbf{x}) = \mathbf{A} \in \mathbb{R}^{4 \times 2}$ where the entries A_{ij} depend on a vector $\mathbf{x} \in \mathbb{R}^3$
- We can compute $\frac{d\mathbf{A}(\mathbf{x})}{d\mathbf{x}} \in \mathbb{R}^{4 \times 2 \times 3}$ in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



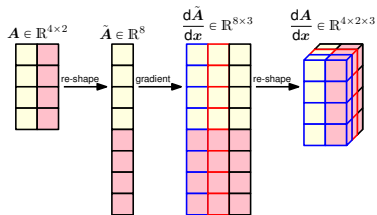
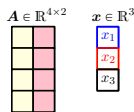
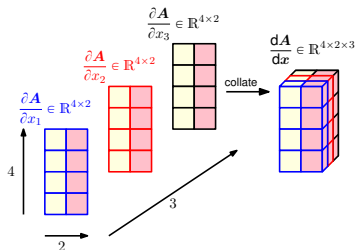
Partial derivatives:

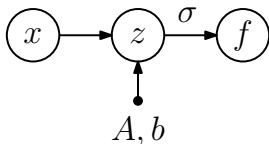


- Consider $f : \mathbb{R}^3 \rightarrow \mathbb{R}^{4 \times 2}$, $f(\mathbf{x}) = \mathbf{A} \in \mathbb{R}^{4 \times 2}$ where the entries A_{ij} depend on a vector $\mathbf{x} \in \mathbb{R}^3$
- We can compute $\frac{d\mathbf{A}(\mathbf{x})}{d\mathbf{x}} \in \mathbb{R}^{4 \times 2 \times 3}$ in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



Partial derivatives:





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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

- 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

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

- 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

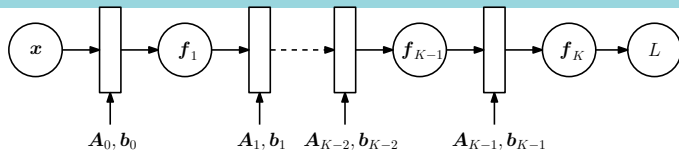
$$L(\boldsymbol{\theta}) = \frac{1}{2} \|\mathbf{e}\|^2 \in \mathbb{R}, \quad \mathbf{e} = \mathbf{y} - \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{z}) \in \mathbb{R}^M$$

is minimized

Partial derivatives:

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

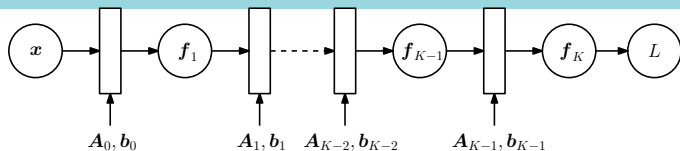
$$\begin{aligned} \frac{\partial L}{\partial \mathbf{e}} &= \underbrace{\mathbf{e}^\top}_{\in \mathbb{R}^{1 \times M}} & \frac{\partial \mathbf{e}}{\partial \mathbf{f}} &= \underbrace{-\mathbf{I}}_{\in \mathbb{R}^{M \times M}} & \frac{\partial \mathbf{f}}{\partial \mathbf{z}} &= \underbrace{\text{diag}(1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



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

$$f_0 = x$$

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



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

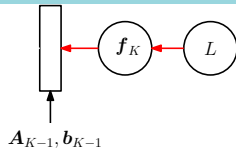
$$\mathbf{f}_0 = \mathbf{x}$$

$$\mathbf{f}_i = \sigma_i(\mathbf{A}_{i-1} \mathbf{f}_{i-1} + \mathbf{b}_{i-1}), \quad i = 1, \dots, K$$

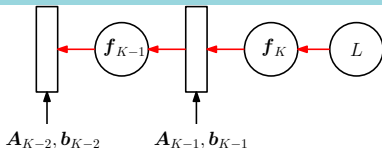
- Find $\mathbf{A}_j, \mathbf{b}_j$ for $j = 0, \dots, K - 1$, such that the squared loss

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

is minimized, where $\boldsymbol{\theta} = \{\mathbf{A}_j, \mathbf{b}_j\}$, $j = 0, \dots, K - 1$

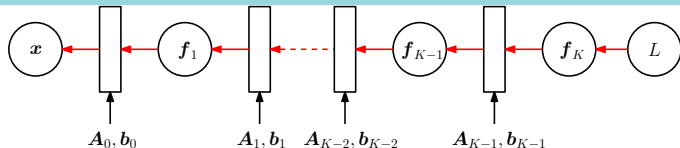


$$\frac{\partial L}{\partial \boldsymbol{\theta}_{K-1}} = \frac{\partial L}{\partial \mathbf{f}_K} \frac{\partial \mathbf{f}_K}{\partial \boldsymbol{\theta}_{K-1}}$$



$$\frac{\partial L}{\partial \theta_{K-1}} = \frac{\partial L}{\partial \mathbf{f}_K} \frac{\partial \mathbf{f}_K}{\partial \theta_{K-1}}$$

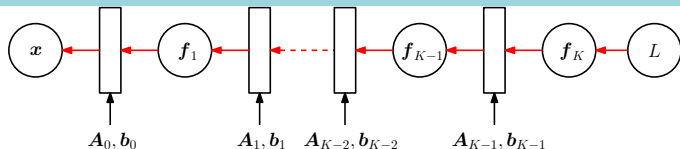
$$\frac{\partial L}{\partial \theta_{K-2}} = \frac{\partial L}{\partial \mathbf{f}_K} \boxed{\frac{\partial \mathbf{f}_K}{\partial \mathbf{f}_{K-1}} \frac{\partial \mathbf{f}_{K-1}}{\partial \theta_{K-2}}}$$



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

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

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

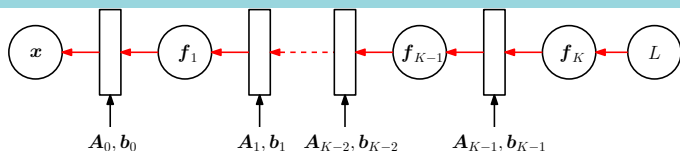


$$\frac{\partial L}{\partial \theta_{K-1}} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial \theta_{K-1}}$$

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

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

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



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

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

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

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

► Intermediate derivatives are stored during the forward pass

- Linear regression with a neural network parametrized by θ :

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

- Linear regression with a neural network parametrized by θ :

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

- Given inputs \mathbf{x}_n and corresponding (noisy) observations y_n , $n = 1, \dots, N$, find parameters θ^* 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 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 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 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)

- Linear regression with a polynomial of order M :

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

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

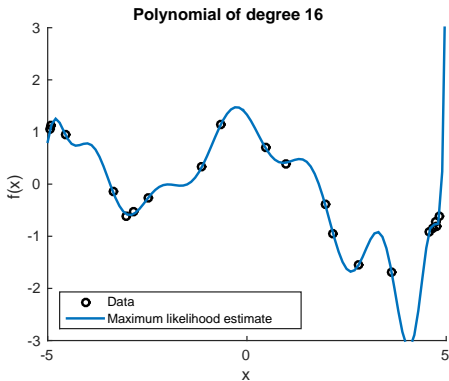
- Linear regression with a polynomial of order M :

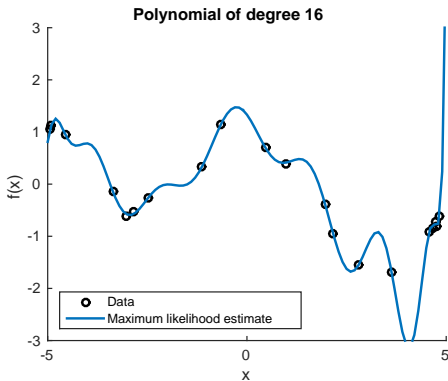
$$y = f(x, \boldsymbol{\theta}) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_\epsilon^2)$$

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

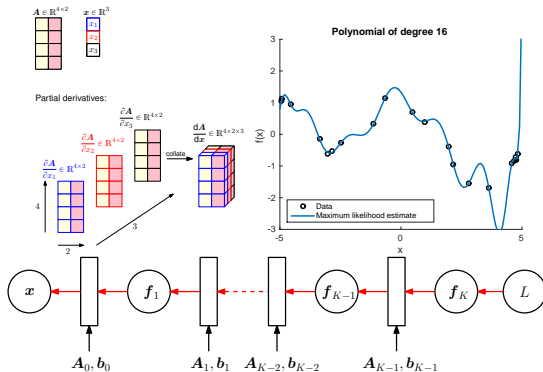
- Given inputs x_i and corresponding (noisy) observations y_i , $i = 1, \dots, N$, find parameters $\boldsymbol{\theta} = [\theta_0, \dots, \theta_M]^\top$, that minimize the squared loss (equivalently: maximize the likelihood)

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





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



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