# Mathematics for AI

## Lecture 18
Jacobian Matrix, Chain Rule, Automatic Differentiation, Backpropagation

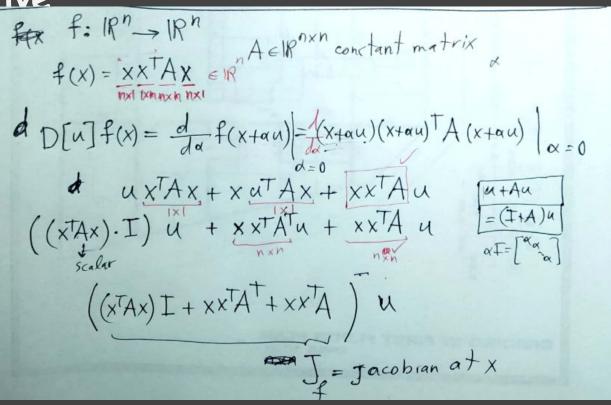# Remember: The Jacobian Matrix

# Example: Derive Jacobian by directional derivative



$f: \mathbb{R}^n \to \mathbb{R}^n$    $A \in \mathbb{R}^{n \times n}$ constant matrix

$$f(x) = \underset{n \times 1}{x} \underset{1 \times n}{x^T} \underset{n \times n}{A} \underset{n \times 1}{x} \in \mathbb{R}^n$$

$$D[u] f(x) = \frac{d}{d\alpha} f(x + \alpha u) \Big| = \frac{d}{d\alpha} (x + \alpha u)(x + \alpha u)^T A (x + \alpha u) \Big|_{\alpha = 0}$$
$$\qquad \alpha = 0$$

$$u \underset{1 \times 1}{x^T A x} + x \underset{1 \times 1}{u^T A x} + \boxed{x x^T A u}$$

$$\qquad \boxed{\begin{array}{l} u + A u \\ = (I + A) u \end{array}}$$

$$((x^T A x) \cdot I) u + x \underset{n \times n}{x^T A^T} u + \underset{n \times n}{x x^T A} u$$
$$\downarrow$$
$$\text{scalar}$$

$$\alpha I = \begin{bmatrix} \alpha & & \\ & \alpha & \\ & & \alpha \end{bmatrix}$$

$$\left( (x^T A x) I + x x^T A^T + x x^T A \right) u$$

$$\underbrace{\qquad\qquad\qquad\qquad\qquad}$$

$$J_f = \text{Jacobian at } x$$

# Multi-layer neural nets

# Multi-layer neural nets



$$f(x) = \sigma(Ax + b)$$

$$y_p = f_p(y_{p-1}, \theta_p) = f_p(f_{p-1}(y_{p-2}, \theta_{p-1}), \theta_p)$$

$$y_p = f_p\left(\cdots f_3\left(f_2\left(f_1(x, \theta_1), \theta_2\right), \theta_3\right) \cdots, \theta_p\right)$$

$$C(\theta_1, \theta_2 \cdots \theta_p) = \sum_{i=1}^{N} d\left(y_p^i, f_p(f_{p-1}(\cdots f_2(f_1(\underline{x}^i, \theta_1), \theta_2) \cdots)\right)$$

$$\frac{\partial C}{\partial \theta_i} = \mathbb{0} \quad \nabla_C^{\theta_i}$$

# Derivative of Composition of functions



$$f: \mathbb{R}^n \longrightarrow \mathbb{R}^m$$
$$g: \mathbb{R}^m \longrightarrow \mathbb{R}^P$$

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

$$g(\underbrace{f}_{\mathbb{R}^m}(\theta))$$

$$\underset{\mathbb{R}^P}{\underbrace{\underset{\mathbb{R}^n}{\downarrow}}}$$

$$n = m = p = 1 \implies \frac{\partial}{\partial \theta} = \frac{d}{d\theta} g(f(\theta)) = g'(f(\theta)) \, f(\theta)$$
$$= g'\big|_{f(\theta)} \, f\big|_{\theta}$$

Chain Rule

# Derivative of Composition of functions



$f: \mathbb{R}^n \to \mathbb{R}^m \implies J_f \in \mathbb{R}^{m \times n}$

$g: \mathbb{R}^m \to \mathbb{R}^p \implies J_g \in \mathbb{R}^{p \times m}$

$J_{g \circ f} = ?$

Intuitive ~~Proof~~ Argument

$D[u](g \circ f)(x) = \dfrac{d}{d\alpha}$

$D[u](g \circ f)(x) = \dfrac{d}{d\alpha} g(f(x + \alpha u))\Big|_{\alpha = 0}$

$\cong \dfrac{d}{d\alpha} g(f(x) + J_f(\alpha u))$

$= \dfrac{d}{d\alpha} g(f(x) + \alpha J_f u)\Big|_{\alpha = 0}$

$= D[J_f u] g\Big|_{f(x)}$

$= J_g\Big|_{f(x)} J_f\Big|_x u$

$= J_g(f(x)) \underbrace{J_f(x)}_{J_{f \circ g}(x)} u$

# Derivative of Composition of functions

$$C(\theta) = f_n\left(f_{n-1}\left(\cdots f_2\left(f_1(\theta)\right)\right)\right)$$

$$J_C^\theta = J_{f_n}\Big|_{f_n} J_{f_{n-1}}\Big|_{f_{n-1}} \cdots J_{f_2}\Big|_{f_1(\theta)} J_{f_1}\Big|_\theta \qquad \text{chain Rule}$$

$$J = \frac{\partial f_n}{\partial f_{n-1}} \cdots \frac{\partial f_2}{\partial f_1} \frac{\partial f_1}{\partial \theta}$$

# Numeric Differentiation



(f(x)) Numeric ~~Computation~~ Differentiation

$$h(x) = h\left(\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}\right)$$

$$\frac{\partial h}{\partial x_1} \simeq \frac{h\left(\begin{bmatrix} x_1 + \varepsilon \\ x_2 \\ \vdots \\ x_n \end{bmatrix}\right) - h\left(\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}\right)}{\varepsilon}$$

$$h(x) = f_n(f_{n-1}(\cdots f_2(f_1(x))\cdots))$$

$$\nabla = \begin{bmatrix} \dfrac{\partial h}{\partial x_1} \\ \dfrac{\partial h}{\partial x_2} \\ \vdots \end{bmatrix}$$

# Symbolic Differentiation



Symbolic Differentiation

$f(x)$

$\dfrac{x * x}{2x}$ $\longrightarrow$ too slow

# Algorithmic Differentiation (Automatic Differentiation)



https://en.wikipedia.org/wiki/Automatic_differentiation