# From derivatives to Jacobians

## Derivative: scalar input, scalar output

$f: R \rightarrow R$

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

← how much the "output" changes

← how much the "input" changes

$$f'(x) = \frac{\partial f(x)}{\partial x} = \frac{\partial y}{\partial x} = \lim_{h \to 0} \frac{f(x+h) - f(x)}{h}$$
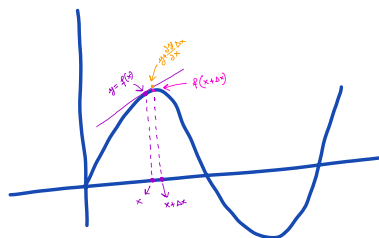
$$f(x+h) \cong f'(x) h + f(x)$$
$$f(x + \Delta x) \cong f'(x) \Delta x + f(x)$$

$$f(x + \Delta x) \cong \frac{\partial y}{\partial x} \Delta x + f(x)$$

$$y^{NEW} \cong \frac{\partial y}{\partial x} \Delta x + y^{OLD}$$

$$x^{NEW} \longrightarrow x^{old} + \Delta x \implies y^{NEW} \overset{\sim}{\longrightarrow} y^{old} + \frac{\partial y}{\partial x} \Delta x$$
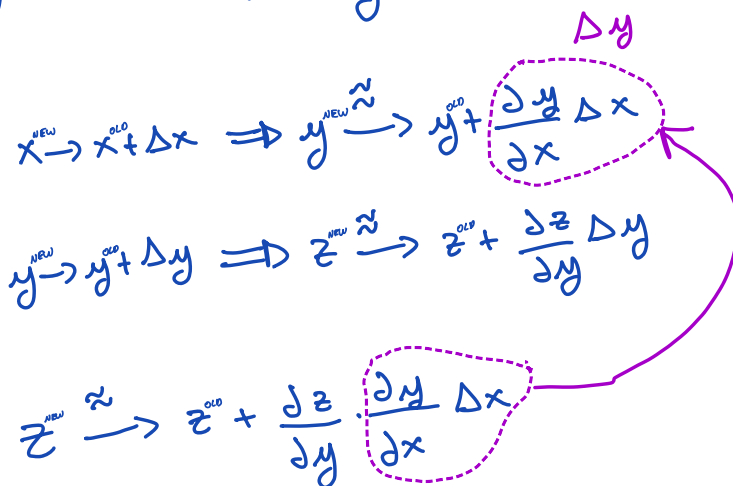
"If x is changed by $\Delta x$, then y will change approximately

by $\frac{\partial y}{\partial x} \cdot \Delta x$"

# Chain Rule

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

$$x \xrightarrow{g} y \xrightarrow{f} z$$

$$x^{NEW} \to x^{OLD} + \Delta x \implies y^{NEW} \xrightarrow{\approx} y^{OLD} + \underbrace{\frac{\partial y}{\partial x} \Delta x}_{\Delta y}$$

$$y^{NEW} \to y^{OLD} + \Delta y \implies z^{NEW} \xrightarrow{\approx} z^{OLD} + \frac{\partial z}{\partial y} \Delta y$$

$$\implies z^{NEW} \xrightarrow{\approx} z^{OLD} + \underbrace{\frac{\partial z}{\partial y} \cdot \frac{\partial y}{\partial x} \Delta x}$$

$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \cdot \frac{\partial y}{\partial x}$$

Gradient : vector input, scalar output

$$f : R^N \to R$$

$$f\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = y$$

gradient $= \left(\dfrac{\partial y}{\partial x_1}, \dfrac{\partial y}{\partial x_2}, \cdots\right)$

$$x \xrightarrow{\text{NEW}} \underbrace{x^{\text{OLD}} + \Delta x} \implies y^{\text{NEW}} \xrightarrow{\approx} y^{\text{OLD}} + \underbrace{\dfrac{\partial y}{\partial x} \cdot \Delta x}$$

vector sum

dot product

$$\dfrac{\partial y}{\partial x} \cdot \Delta x = \underbrace{\dfrac{\partial y}{\partial x_1}} \cdot \Delta x_1 + \dfrac{\partial y}{\partial x_2} \cdot \Delta x_2 + \cdots + \dfrac{\partial y}{\partial x_N} \cdot \Delta x_N$$

partial derivative

Note: the chain rule applies in the same way as in the scalar case

Jacobian: vector input, vector output

$$f : \mathbb{R}^N \longrightarrow \mathbb{R}^M$$

$$f\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

$$\text{Jacobian} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_1}{\partial x_N} \\ \vdots & & \vdots \\ \frac{\partial y_M}{\partial x_1} & \cdots & \frac{\partial y_M}{\partial x_N} \end{bmatrix}$$

$$x^{NEW} \longrightarrow x^{OLD} + \Delta x \implies y^{NEW} \xrightarrow{\approx} y^{OLD} + \underbrace{\frac{\partial y}{\partial x} \Delta x}_{\substack{\text{Matrix-vector} \\ \text{product}}}$$

$$(M \times N) \times (N \times 1) = (M \times 1)$$
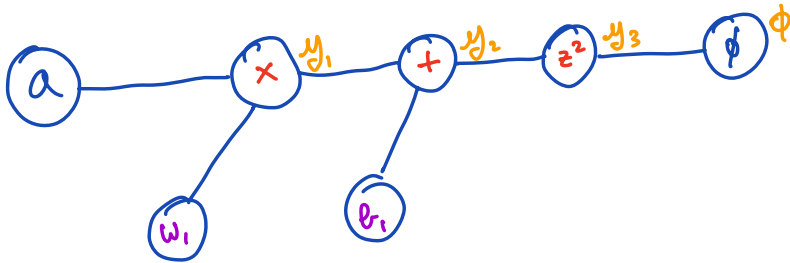
# Generalized Jacobian: tensor input, tensor output

$$f: \mathbb{R}^{N_1 \times \dots \times N_{D_x}} \longrightarrow \mathbb{R}^{M_1 \times \dots M_{D_y}}$$

$$f\left(\begin{array}{c} D_x\text{-dimensional} \\ \text{tensor} \end{array}\right) = \begin{array}{c} D_y\text{-dimensional} \\ \text{tensor} \end{array}$$

$$x^{NEW} \longrightarrow \underbrace{x^{old} + \Delta x}_{\substack{\text{tensor} \\ \text{sum}}} \Longrightarrow y^{NEW} \xrightarrow{\approx} y^{old} + \underbrace{\frac{\partial y}{\partial x} \Delta x}_{\substack{\text{Tensor product} \\ (M_1 \times \dots M_{D_y}) \times (N_1 \times \dots N_{D_x})}}$$

Generalized Jacobian

# Autograd with derivatives



$$\phi = y_3 = (y_2)^2 = (y_1 + b_1)^2 = (aw_1 + b_1)^2$$

$$\frac{\partial \phi}{\partial w_1} = 2(aw_1 + b_1)(a) = 2a(aw_1 + b_1)$$

$$\frac{\partial \phi}{\partial w_1} = \frac{\partial \phi}{\partial y_3} \cdot \frac{\partial y_3}{\partial y_2} \cdot \frac{\partial y_2}{\partial y_1} \cdot \frac{\partial y_1}{\partial w_1} =$$
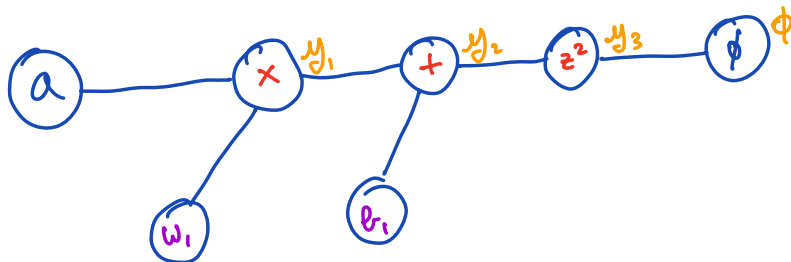
$$= 1 \cdot 2y_2 \cdot 1 \cdot a =$$
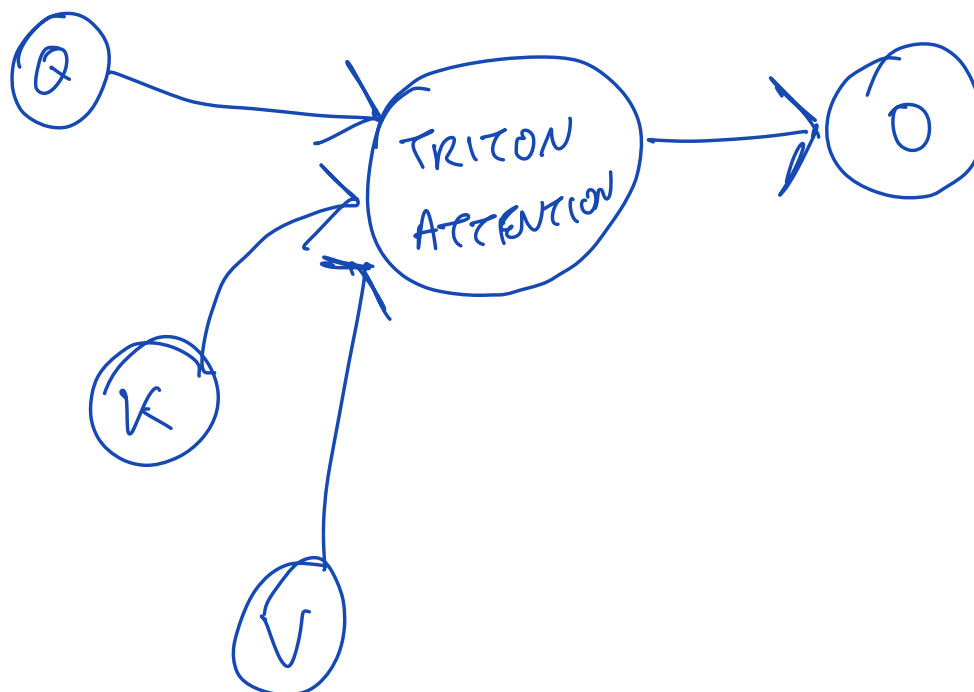
$$= 2ay_2 = 2a(aw_1 + b_1)$$

$$\frac{\partial \phi}{\partial w_1} = \underbrace{\overbrace{\frac{\partial \phi}{\partial y_3} \cdot \frac{\partial y_3}{\partial y_2}}^{\textstyle \frac{\partial \phi}{\partial y_1}} \cdot \frac{\partial y_2}{\partial y_1}}_{\textstyle \frac{\partial \phi}{\partial y_2}} \cdot \frac{\partial y_1}{\partial w_1}$$

$$\frac{\partial \phi}{y_2} = \frac{\partial \phi}{\partial y_3} \cdot \frac{\partial y_3}{\partial y_2}$$

$$\frac{\partial \phi}{y_1} = \frac{\partial \phi}{y_2} \cdot \frac{\partial y_2}{\partial y_1}$$

$$\frac{\partial \phi}{\partial w_1} = \frac{\partial \phi}{y_1} \cdot \frac{\partial y_1}{\partial w_1}$$

$$Q \rightarrow \text{TRITON ATTENTION} \rightarrow O$$

$$K \nearrow$$

$$V \nearrow$$

$N = 1024 \qquad 1024$

$\lceil N, M \rceil \qquad \downarrow \leftarrow \qquad \lceil N, D, M \rceil \qquad \sqcup 8$

$\lceil N, D \rceil$

$[...,...] \qquad [N, M] \qquad [...,...] \to 20^{...}$

$Y = XW$

UPSTREAM GRADIENT

$$\frac{\partial \phi}{\partial X} = \frac{\partial \phi}{\partial X} \cdot \frac{\partial Y}{\partial X}$$

$(N, M) \times (N, D)$

$1024 \times 2048 \times 1024 \times 1024$

LOCAL JACOBIAN

DOWNSTREAM GRADIENT



$$\begin{bmatrix} [.\,.\,.\,.\,.\, \quad ] \\ [.\,.\,-\,-\, \quad ] \\ [-\,-\,-\,-\, \quad ] \\ [.\,.\,.\,-\, \quad ] \end{bmatrix} \times \begin{bmatrix} \quad \\ \quad \end{bmatrix} =$$

$(N, D) \qquad\qquad (D, M)$

$$\begin{bmatrix} [.\,.\,.\,- \quad ] \\ [.\,-\,-\, \quad ] \\ [.\,-\,-\, \quad ] \\ [.\,.\,.\,.\, \quad ] \end{bmatrix}$$