



성균관대학교
SUNGKYUNKWAN UNIVERSITY

Deep Learning

- Backpropagation and Automatic Differentiation -

Eunbyung Park

Assistant Professor

School of Electronic and Electrical Engineering

[Eunbyung Park \(silverbottlep.github.io\)](https://github.com/silverbottlep)

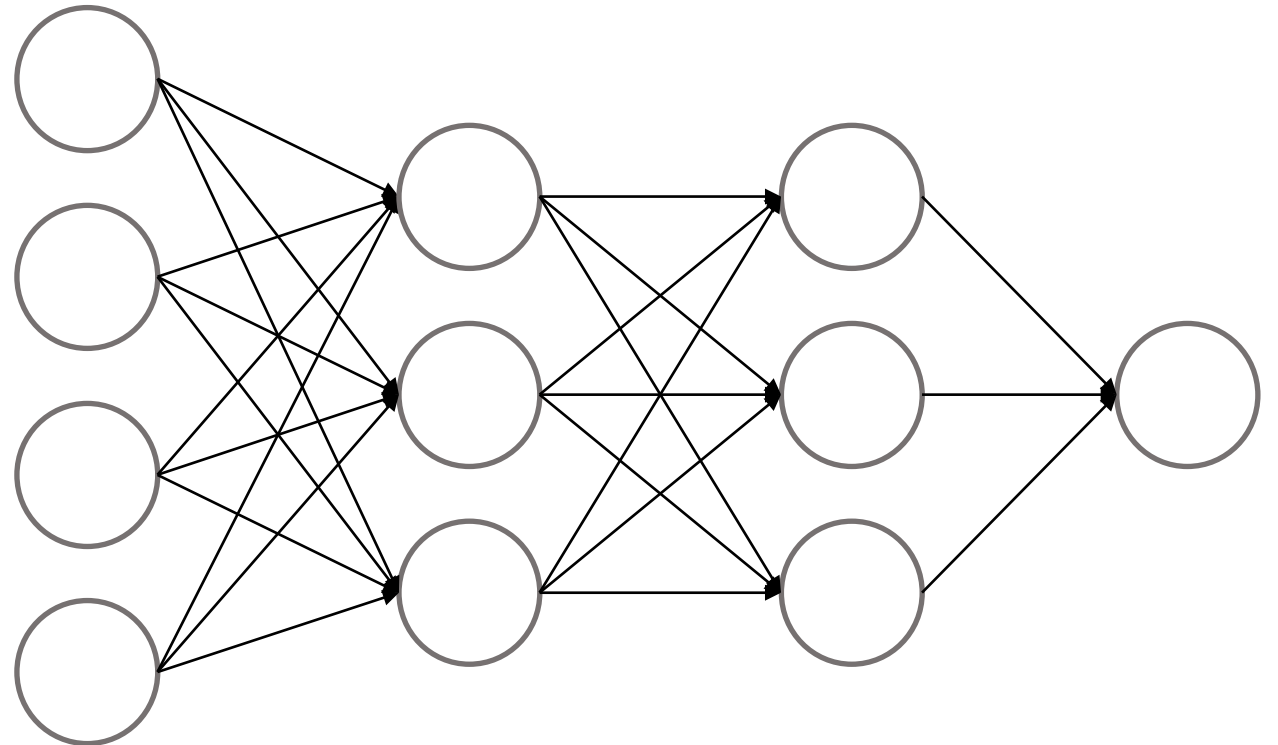
The Chain Rule

Gradient Descent

- We are using *gradient descent* for training deep neural networks
- An algorithm *to compute the derivatives of a loss function* for deep neural networks

$$W := W - \alpha \left(\frac{\partial L}{\partial W} \right)$$

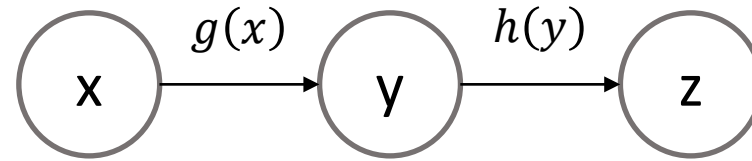
(descent) (step-size) (gradient)



The Chain Rule

- A single variable chain rule

$$f, g, h: \mathbb{R} \rightarrow \mathbb{R}$$



$$f: h \circ g$$

$$f'(x) = h'(g(x))g'(x)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

$$y = g(x), z = h(y)$$

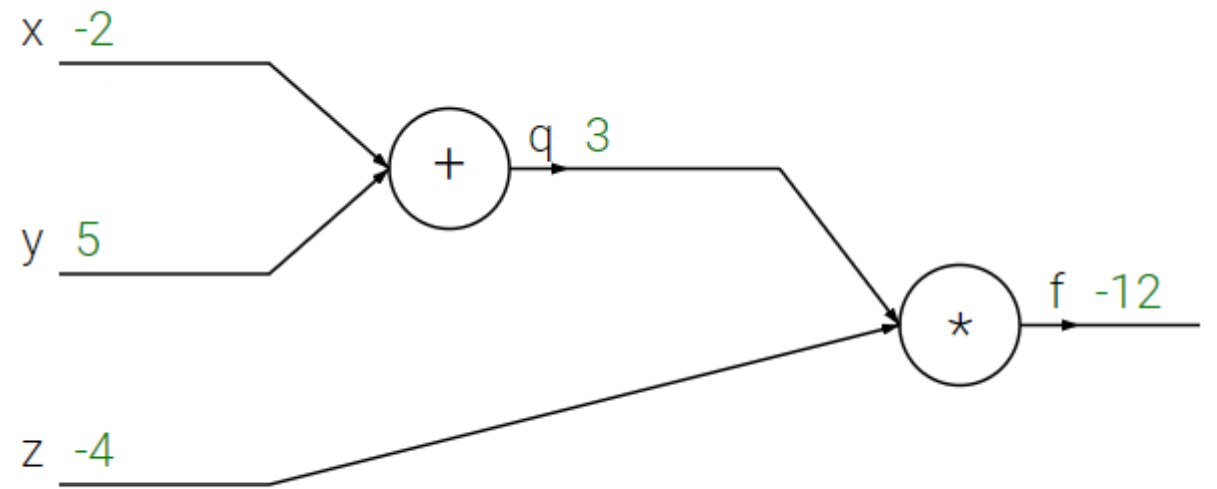
Simple Example

$$f(x, y, z) = (x + y)z$$

$$q = x + y, \quad f = qz$$

$$\frac{\partial q}{\partial x} = 1, \quad \frac{\partial q}{\partial y} = 1$$

$$\frac{\partial f}{\partial q} = z, \quad \frac{\partial f}{\partial z} = q$$



Simple Example

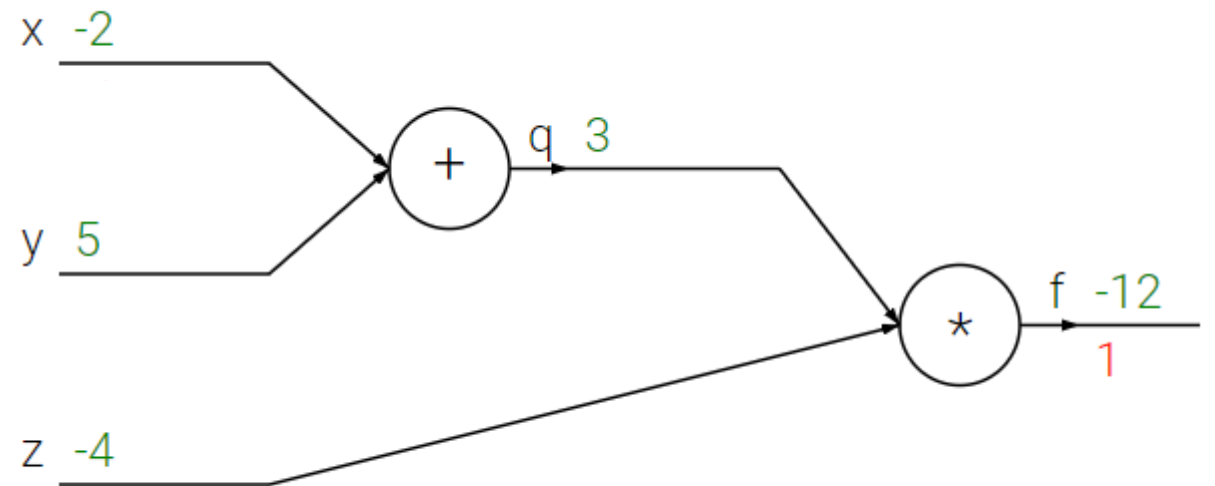
$$f(x, y, z) = (x + y)z$$

$$q = x + y, \quad f = qz$$

$$\frac{\partial q}{\partial x} = 1, \quad \frac{\partial q}{\partial y} = 1$$

$$\frac{\partial f}{\partial q} = z, \quad \frac{\partial f}{\partial z} = q$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$



Simple Example

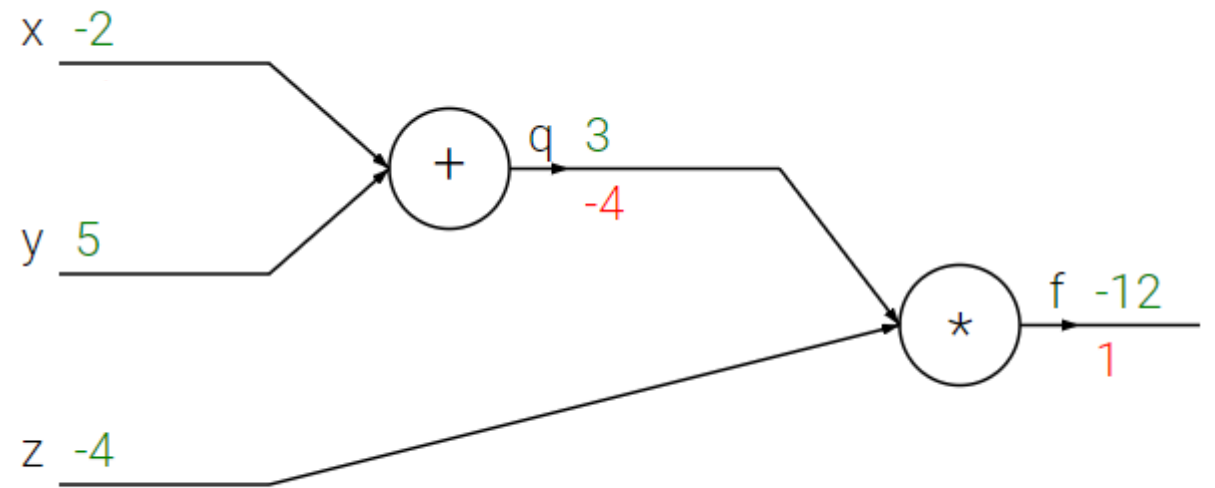
$$f(x, y, z) = (x + y)z$$

$$q = x + y, \quad f = qz$$

$$\frac{\partial q}{\partial x} = 1, \quad \frac{\partial q}{\partial y} = 1$$

$$\frac{\partial f}{\partial q} = z, \quad \frac{\partial f}{\partial z} = q$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$



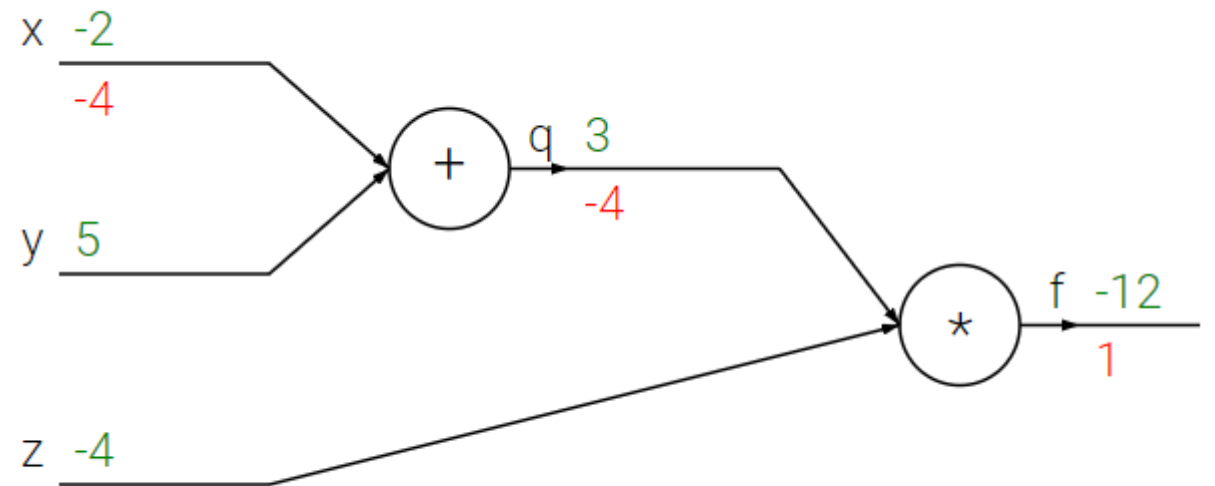
Simple Example

$$f(x, y, z) = (x + y)z$$

$$q = x + y, \quad f = qz$$

$$\frac{\partial q}{\partial x} = 1, \quad \frac{\partial q}{\partial y} = 1$$
$$\frac{\partial f}{\partial q} = z, \quad \frac{\partial f}{\partial z} = q$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$



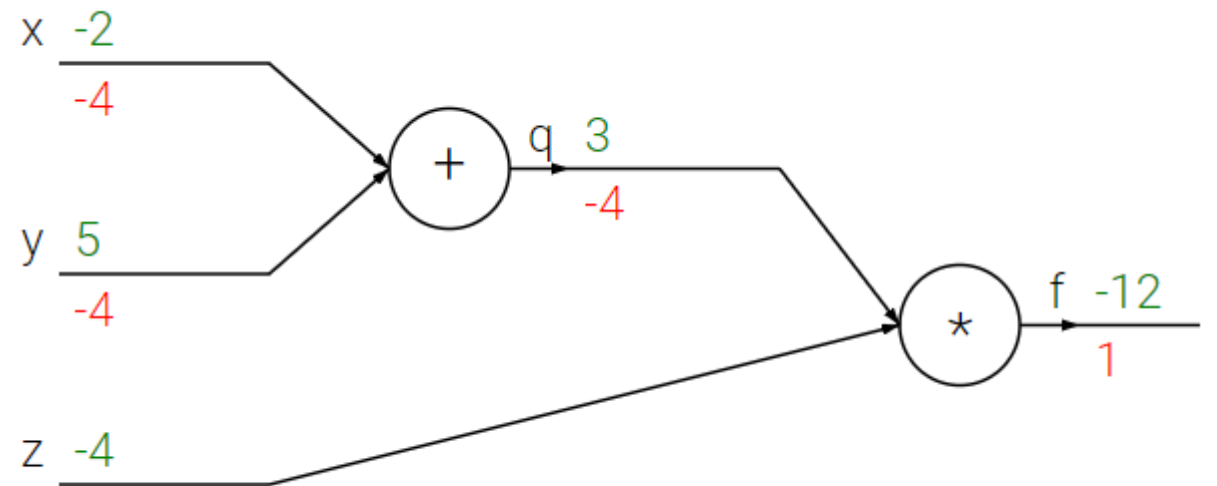
Simple Example

$$f(x, y, z) = (x + y)z$$

$$q = x + y, \quad f = qz$$

$$\frac{\partial q}{\partial x} = 1, \quad \frac{\partial q}{\partial y} = 1$$
$$\frac{\partial f}{\partial q} = z, \quad \frac{\partial f}{\partial z} = q$$

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y}$$



Simple Example

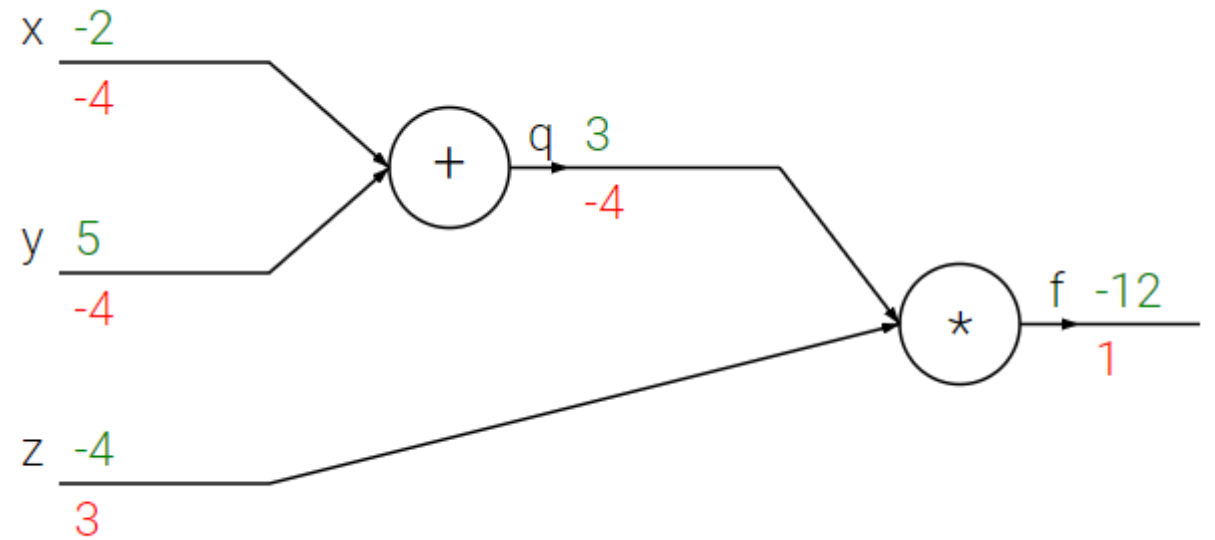
$$f(x, y, z) = (x + y)z$$

$$q = x + y, \quad f = qz$$

$$\frac{\partial q}{\partial x} = 1, \quad \frac{\partial q}{\partial y} = 1$$

$$\frac{\partial f}{\partial q} = z, \quad \frac{\partial f}{\partial z} = q$$

$$\frac{\partial f}{\partial z} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial z}$$



Sigmoid Example

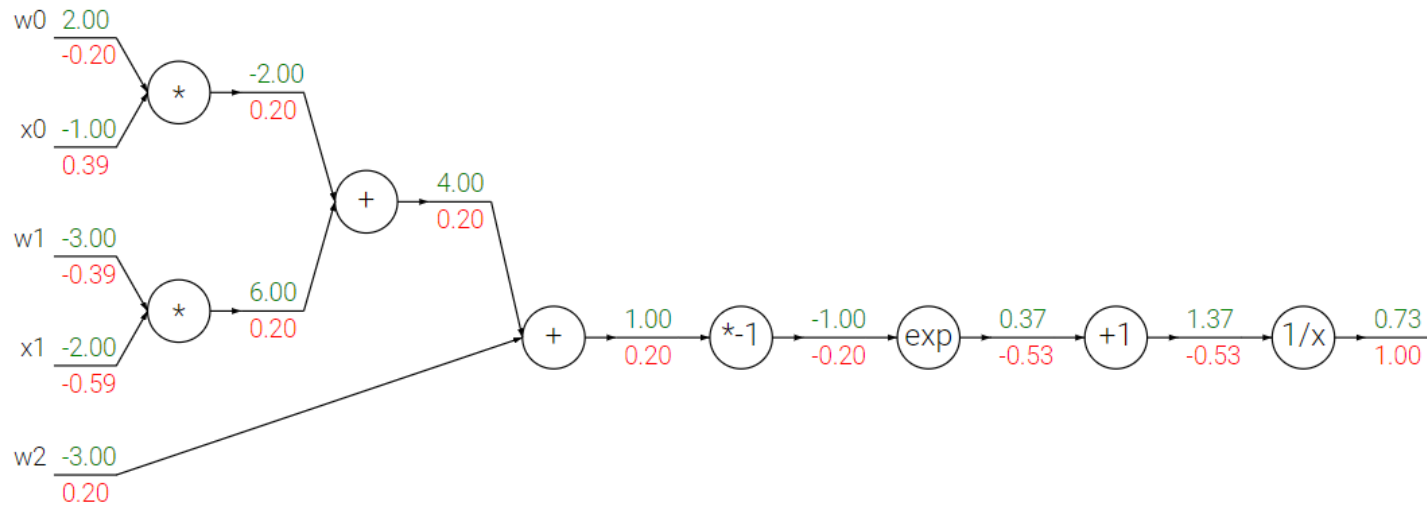
$$\sigma(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

$$f(x) = \frac{1}{x}, \quad g(x) = 1 + x, \quad h(x) = e^{-x}, \quad i(x) = w_0x_0 + w_1x_1 + w_2$$

Sigmoid Example

$$\sigma(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

$$f(x) = \frac{1}{x}, \quad g(x) = 1 + x, \quad h(x) = e^{-x}, \quad i(x) = w_0x_0 + w_1x_1 + w_2$$



Backpropagation Algorithm

Gradient

- In vector calculus, the *gradient* of a *scalar-valued* differentiable function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ at the point x

$$\nabla f: \mathbb{R}^n \rightarrow \mathbb{R}^n$$

$$\nabla f = \frac{\partial f}{\partial x} = \left[\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right]$$

Jacobian

- In vector calculus, the *Jacobian* of a *vector-valued* differentiable function is the matrix of all its first-order partial derivatives.

$$f: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$J_{ij} = \frac{\partial f_i}{\partial x_j}$$

$$J = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

Matrix Calculus

$$X \in \mathbb{R}^{n \times m}, y \in \mathbb{R}$$

$$f: \mathbb{R}^{n \times m} \rightarrow \mathbb{R}$$

$$y = f(x)$$

$$\frac{\partial y}{\partial X} = \begin{bmatrix} \frac{\partial y}{\partial X_{11}} & \cdots & \frac{\partial y}{\partial X_{1m}} \\ \vdots & \ddots & \vdots \\ \frac{\partial y}{\partial X_{n1}} & \cdots & \frac{\partial y}{\partial X_{nm}} \end{bmatrix} \in \mathbb{R}^{n \times m}$$

Matrix Calculus

$$X \in \mathbb{R}^{n \times m}, y \in \mathbb{R}^l$$

$$f: \mathbb{R}^{n \times m} \rightarrow \mathbb{R}^l$$

$$y = f(x)$$

$$\frac{\partial y_1}{\partial X} = \begin{bmatrix} \frac{\partial y_1}{\partial X_{11}} & \cdots & \frac{\partial y_1}{\partial X_{1m}} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial X_{n1}} & \cdots & \frac{\partial y_1}{\partial X_{nm}} \end{bmatrix} \in \mathbb{R}^{n \times m}$$

$$\frac{\partial y}{\partial X} \in \mathbb{R}^{l \times n \times m} \quad (3 \text{ dim tensor})$$

Finite Difference

- Numerical method to compute the gradients based on the definition of gradients

$$\frac{df}{dx} \approx \frac{f(x + \Delta x) - f(x)}{\Delta x}$$

Forward
difference

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

$$\frac{df}{dx} \approx \frac{f(x) - f(x - \Delta x)}{\Delta x}$$

Backward
difference

$$\frac{df}{dx} \approx \frac{f(x + \Delta x) - f(x - \Delta x)}{2\Delta x}$$

Central
difference

Finite Difference

- Numerical method to compute the gradients based on the definition of gradients

$$\frac{df}{dx} \approx \frac{f(x + \Delta x) - f(x)}{\Delta x}$$

Forward
difference

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

$$\frac{df}{dx} \approx \frac{f(x) - f(x - \Delta x)}{\Delta x}$$

Backward
difference

What's wrong with this approach?

$$\frac{df}{dx} \approx \frac{f(x + \Delta x) - f(x - \Delta x)}{2\Delta x}$$

Central
difference

The Chain Rule

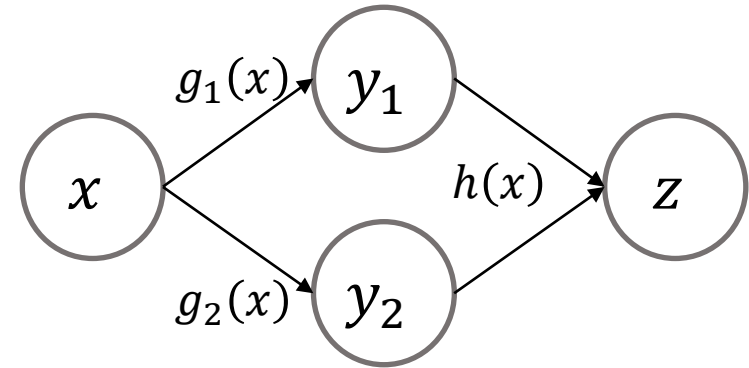
- Multi-variable chain rule

$$f, g_1, g_2: \mathbb{R} \rightarrow \mathbb{R}, \quad h: \mathbb{R}^2 \rightarrow \mathbb{R}$$

$$y_1 = g_1(x), \quad y_2 = g_2(x)$$

$$z = h(y_1, y_2)$$

$$\frac{dz}{dx} = \frac{dz}{dy_1} \frac{dy_1}{dx} + \frac{dz}{dy_2} \frac{dy_2}{dx}$$



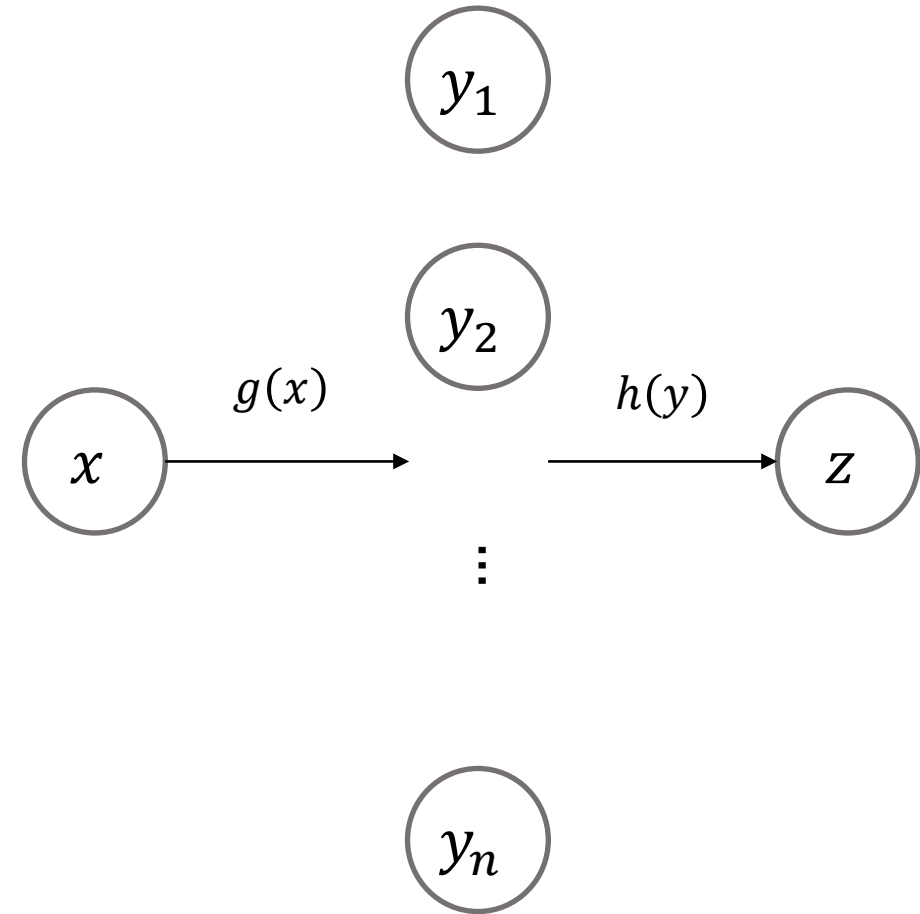
The Chain Rule

- Multi-variable chain rule

$$x \in \mathbb{R}, y \in \mathbb{R}^n, z \in \mathbb{R}$$

$$g: \mathbb{R} \rightarrow \mathbb{R}^n, \quad y = g(x)$$

$$h: \mathbb{R}^n \rightarrow \mathbb{R}, \quad z = h(y)$$



$$\frac{\partial z}{\partial x} = \sum_{i=1}^n \frac{\partial z}{\partial y_i} \frac{dy_i}{dx} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$$

$\in \mathbb{R}^{n \times 1}$

$\in \mathbb{R}^{1 \times n}$

The Chain Rule

- Multi-variable chain rule

$$x \in \mathbb{R}^n, y \in \mathbb{R}^m, z \in \mathbb{R}$$

$$g: \mathbb{R}^n \rightarrow \mathbb{R}^m, \quad y = g(x)$$

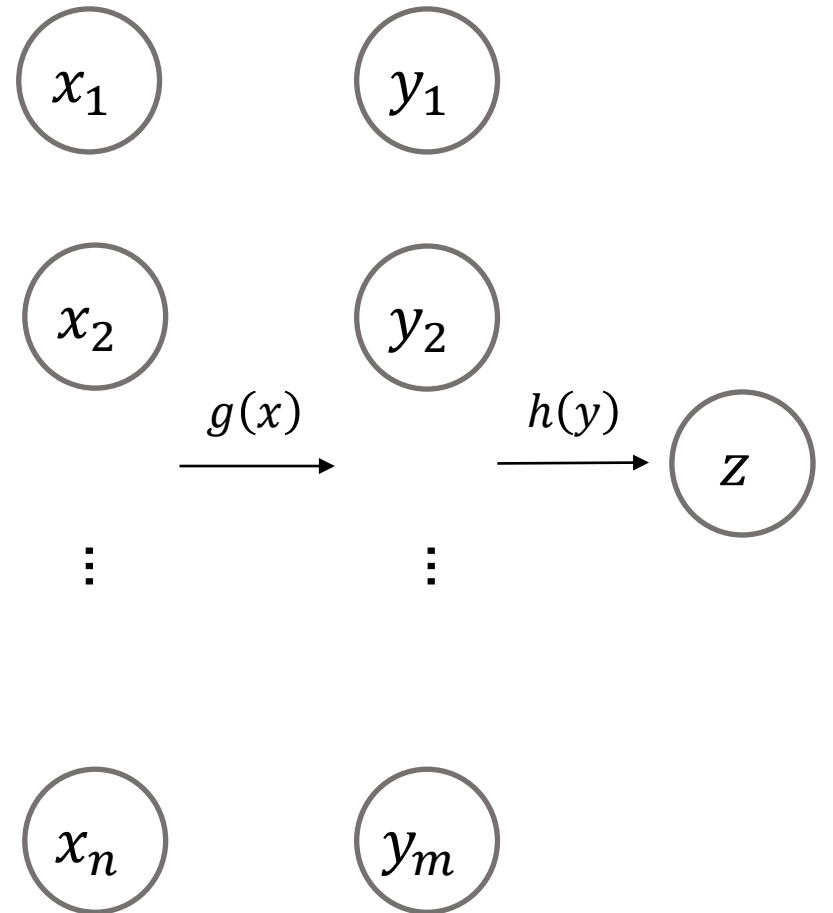
$$h: \mathbb{R}^m \rightarrow \mathbb{R}, \quad z = h(y)$$

$$\frac{\partial z}{\partial x_j} = \sum_{i=1}^m \frac{\partial z}{\partial y_i} \frac{\partial y_i}{\partial x_j} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x_j}$$

$\in \mathbb{R}^{m \times 1}$ (above $\frac{\partial z}{\partial y_i}$)
 $\in \mathbb{R}^{1 \times m}$ (below $\frac{\partial y}{\partial x_j}$)

$$\frac{\partial z}{\partial x} = \left[\sum_{i=1}^m \frac{\partial z}{\partial y_i} \frac{\partial y_i}{\partial x_1}, \dots, \sum_{i=1}^m \frac{\partial z}{\partial y_i} \frac{\partial y_i}{\partial x_n} \right] = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$$

$\in \mathbb{R}^{m \times n}$ (above $\frac{\partial y}{\partial x}$)
 $\in \mathbb{R}^{1 \times m}$ (below $\frac{\partial z}{\partial y}$)



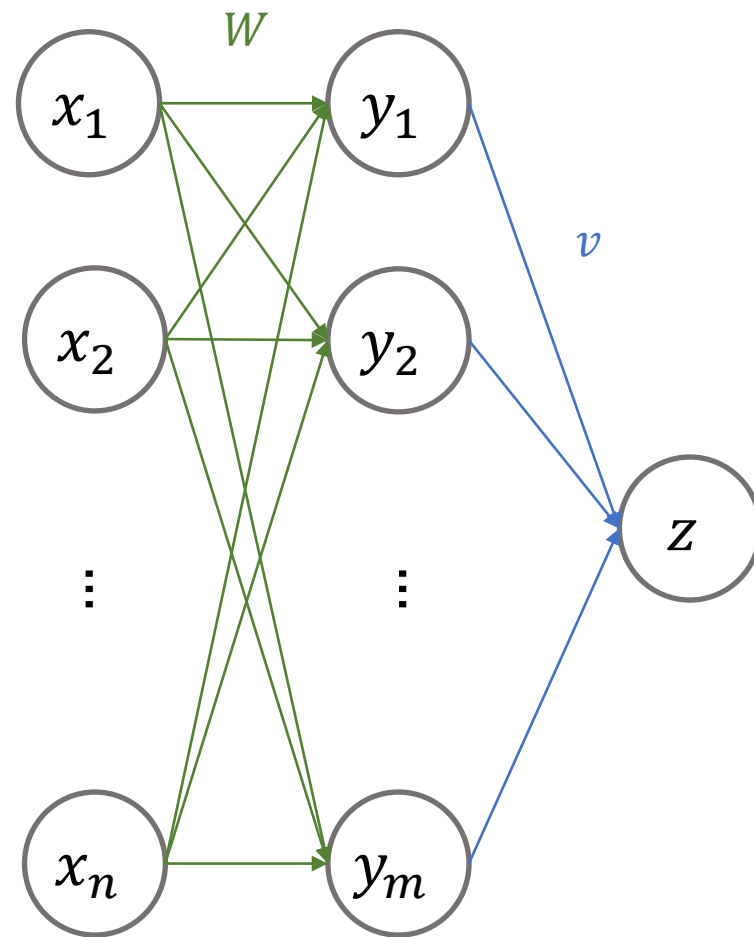
Two Layers MLP

$$x \in \mathbb{R}^n, y \in \mathbb{R}^m, z \in \mathbb{R}, W \in \mathbb{R}^{m \times n}, v \in \mathbb{R}^m$$

$$y = Wx \quad z = \sum_{i=1}^m v_i y_i = v^T y$$

$$\frac{\partial z}{\partial x_j} = \sum_{i=1}^m \frac{\partial z}{\partial y_i} \frac{\partial y_i}{\partial x_j} = \underbrace{\frac{\partial z}{\partial y}}_{\in \mathbb{R}^{1 \times m}} \underbrace{\frac{\partial y}{\partial x_j}}_{\in \mathbb{R}^{m \times 1}}$$

$$\frac{\partial z}{\partial x} = \left[\sum_{i=1}^m \frac{\partial z}{\partial y_i} \frac{\partial y_i}{\partial x_1}, \dots, \sum_{i=1}^m \frac{\partial z}{\partial y_i} \frac{\partial y_i}{\partial x_n} \right] = \underbrace{\frac{\partial z}{\partial y}}_{\in \mathbb{R}^{1 \times m}} \underbrace{\frac{\partial y}{\partial x}}_{\in \mathbb{R}^{m \times n}} = v^T W$$



Derivatives of Linear Layer

$$y = Wx \quad z = \sum_{i=1}^m v_i y_i = v^\top y$$

$$\frac{\partial z}{\partial y}$$

$$\frac{\partial y}{\partial x}$$

Two Layers MLP

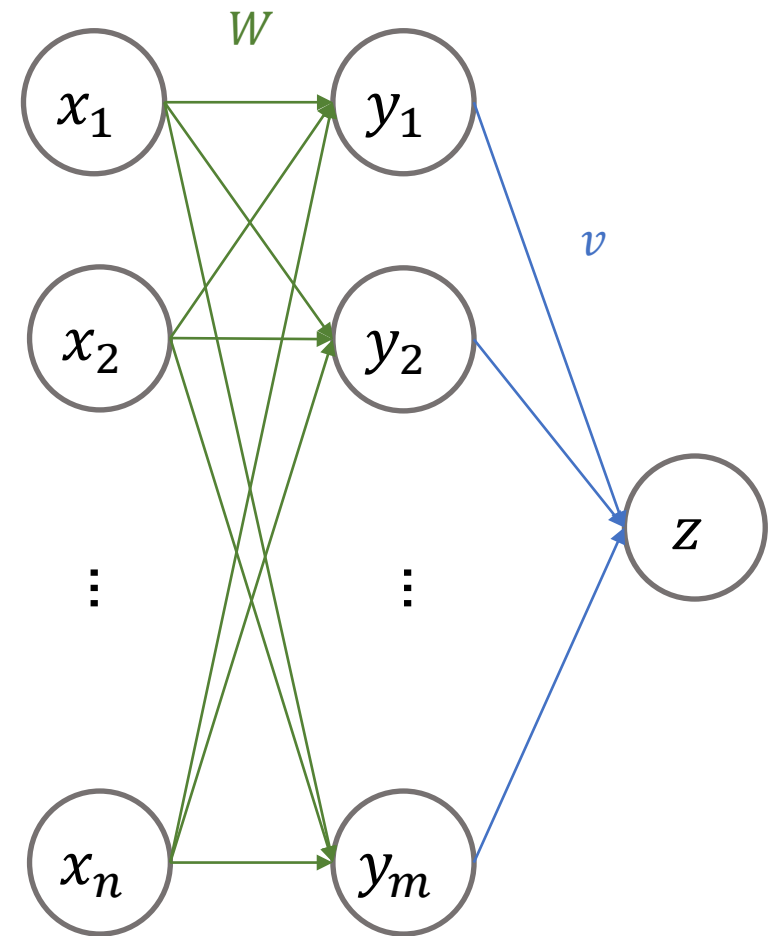
$$x \in \mathbb{R}^n, y \in \mathbb{R}^m, z \in \mathbb{R}, W \in \mathbb{R}^{m \times n}, v \in \mathbb{R}^m$$

$$y = Wx \quad z = \sum_{i=1}^m v_i y_i = v^T y$$

$$\frac{\partial z}{\partial W} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial W}$$

$\in \mathbb{R}^{m \times n}$
 $\in \mathbb{R}^{m \times m \times n}$
 $\in \mathbb{R}^{1 \times m}$

Tensor Product
(n-mode product)

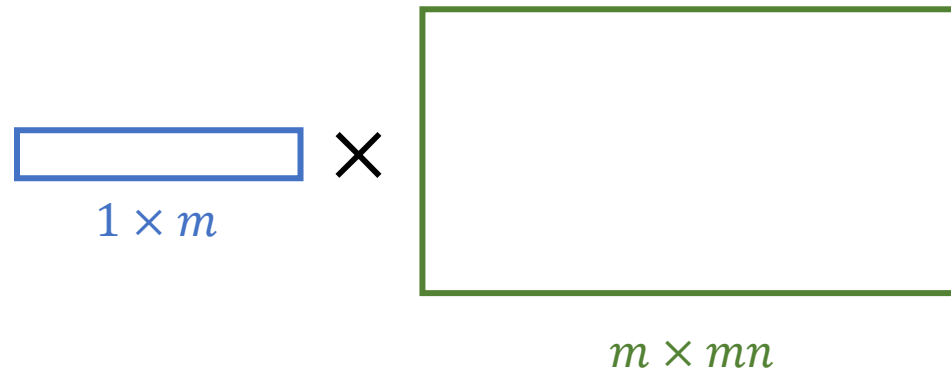
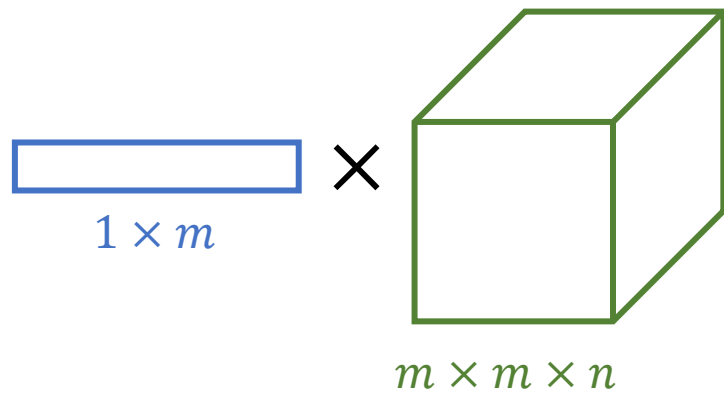


Tensor Product

- N-mode product
 - Matricization -> matrix multiplication

$$\frac{\partial z}{\partial W} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial W}$$

$\in \mathbb{R}^{m \times n}$
 $\in \mathbb{R}^{m \times m \times n}$
 $\in \mathbb{R}^{1 \times m}$



Vector Jacobian Product (VJP)

- Jacobian is very sparse and explicit formation of it is too expensive

$$x \in \mathbb{R}^n, y \in \mathbb{R}^m, z \in \mathbb{R}, W \in \mathbb{R}^{m \times n}, v \in \mathbb{R}^m$$

$$y = Wx$$

$$\frac{\partial y_1}{\partial W} \in \mathbb{R}^{m \times n} = \begin{bmatrix} \frac{\partial y_1}{\partial W_{11}} & \dots & \frac{\partial y_1}{\partial W_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial W_{m1}} & \dots & \frac{\partial y_1}{\partial W_{mn}} \end{bmatrix} = \begin{bmatrix} \frac{\partial y_1}{\partial W_{11}} & \dots & \frac{\partial y_1}{\partial W_{1n}} \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\frac{\partial y_2}{\partial W} \in \mathbb{R}^{m \times n} = \begin{bmatrix} \frac{\partial y_2}{\partial W_{11}} & \dots & \frac{\partial y_2}{\partial W_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_2}{\partial W_{m1}} & \dots & \frac{\partial y_2}{\partial W_{mn}} \end{bmatrix} = \begin{bmatrix} 0 & 0 & 0 \\ \frac{\partial y_2}{\partial W_{21}} & \dots & \frac{\partial y_2}{\partial W_{2n}} \\ 0 & 0 & 0 \end{bmatrix}$$

Vector Jacobian Product (VJP)

- Jacobian is very sparse and explicit formation of it is too expensive

$$\frac{\partial z}{\partial y} \text{reshape} \left(\frac{\partial y}{\partial W} \right) =$$

$$\begin{bmatrix} \frac{\partial z}{\partial y_1} & \frac{\partial z}{\partial y_2} & \dots & \frac{\partial z}{\partial y_m} \end{bmatrix} \begin{bmatrix} \frac{\partial y_1}{\partial W_{11}} & \dots & \frac{\partial y_1}{\partial W_{1n}} & 0 & 0 & 0 & 0 & 0 & 0 & \dots \\ 0 & 0 & 0 & \frac{\partial y_2}{\partial W_{21}} & \dots & \frac{\partial y_2}{\partial W_{2n}} & 0 & 0 & 0 & \dots \\ 0 & 0 & 0 & 0 & 0 & 0 & \frac{\partial y_3}{\partial W_{31}} & \dots & \frac{\partial y_3}{\partial W_{3n}} & \dots \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \dots \end{bmatrix}$$

$$= \begin{bmatrix} \frac{\partial z}{\partial y_1} \frac{\partial y_1}{\partial W_{11}} & \dots & \frac{\partial z}{\partial y_1} \frac{\partial y_1}{\partial W_{1n}} & \frac{\partial z}{\partial y_2} \frac{\partial y_2}{\partial W_{21}} & \dots & \frac{\partial z}{\partial y_2} \frac{\partial y_2}{\partial W_{2n}} & \dots \end{bmatrix}$$

Vector Jacobian Product (VJP)

- Jacobian is very sparse and explicit formation of it is too expensive

$$\begin{aligned} \frac{\partial z}{\partial W} &= \frac{\partial z}{\partial y} \frac{\partial y}{\partial W} = \text{reshape} \left(\frac{\partial z}{\partial y} \text{reshape} \left(\frac{\partial y}{\partial W} \right) \right) \\ &= \text{reshape} \left(\left[\begin{array}{ccc} \frac{\partial z}{\partial y_1} \frac{\partial y_1}{\partial W_{11}} & \cdots & \frac{\partial z}{\partial y_1} \frac{\partial y_1}{\partial W_{1n}} \\ \frac{\partial z}{\partial y_2} \frac{\partial y_2}{\partial W_{21}} & \cdots & \frac{\partial z}{\partial y_2} \frac{\partial y_2}{\partial W_{2n}} \\ \vdots & \ddots & \vdots \end{array} \right] \right) \\ &= \left[\begin{array}{ccc} \frac{\partial z}{\partial y_1} \frac{\partial y_1}{\partial W_{11}} & \cdots & \frac{\partial z}{\partial y_1} \frac{\partial y_1}{\partial W_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial z}{\partial y_m} \frac{\partial y_m}{\partial W_{m1}} & \cdots & \frac{\partial z}{\partial y_m} \frac{\partial y_m}{\partial W_{mn}} \end{array} \right] = \left[\begin{array}{ccc} \frac{\partial z}{\partial y_1} x_1 & \cdots & \frac{\partial z}{\partial y_1} x_n \\ \vdots & \ddots & \vdots \\ \frac{\partial z}{\partial y_m} x_1 & \cdots & \frac{\partial z}{\partial y_m} x_n \end{array} \right] = \left(\frac{\partial z}{\partial y} \right)^\top x^\top \end{aligned}$$

Vector Jacobian Product (VJP)

- Explicit formation of Jacobian is too expensive

$$x \in \mathbb{R}^n, y \in \mathbb{R}^m, z \in \mathbb{R}, W \in \mathbb{R}^{m \times n}, v \in \mathbb{R}^m$$

$$y = Wx \quad z = \sum_{i=1}^m v_i y_i = v^T y$$

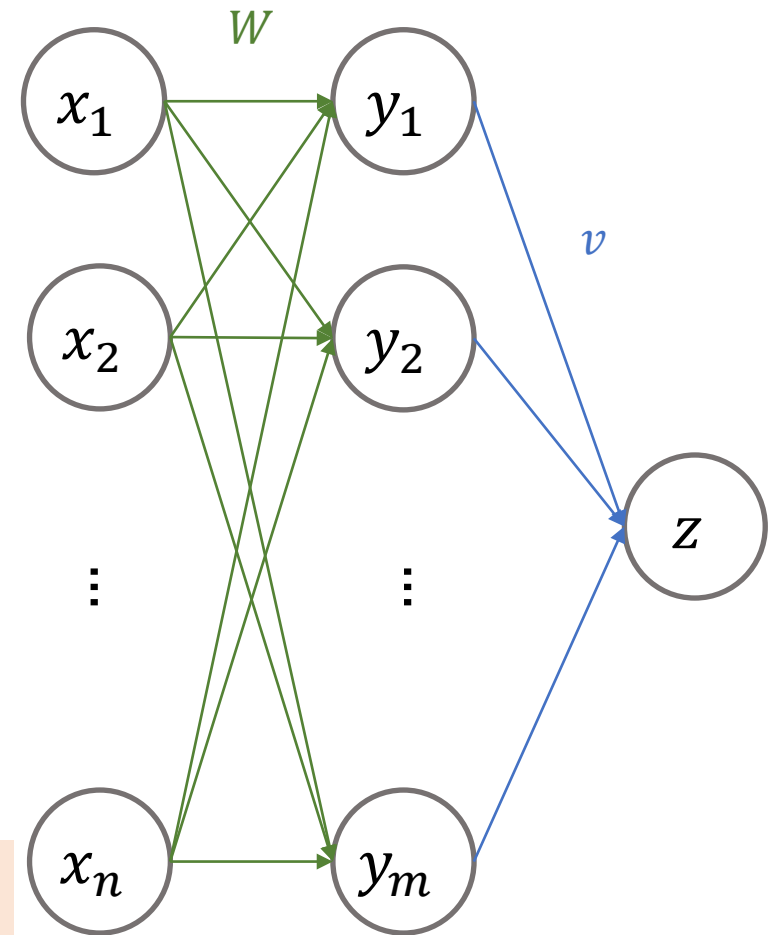
$$\frac{\partial z}{\partial W} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial W}$$

$\in \mathbb{R}^{m \times n}$ (red)
 $\in \mathbb{R}^{1 \times m}$ (blue)
 $\in \mathbb{R}^{m \times m \times n}$ (green)

$$\frac{\partial z}{\partial W} = \left(\frac{\partial z}{\partial y} \right)^T x^T$$

$\in \mathbb{R}^{m \times n}$ (red)
 $\in \mathbb{R}^{m \times 1}$ (blue)
 $\in \mathbb{R}^{1 \times n}$ (green)

We almost never explicitly construct Jacobians $\left(\frac{\partial y}{\partial W} \right)$. We instead directly compute vector-Jacobian product (VJP, $\frac{\partial z}{\partial y} \frac{\partial y}{\partial W}$) in more efficient way $\left(\left(\frac{\partial z}{\partial y} \right)^T x^T \right)$



Vector Jacobian Product (VJP)

- Elementwise activation functions

$$y \in \mathbb{R}^m, \hat{y} \in \mathbb{R}^m$$

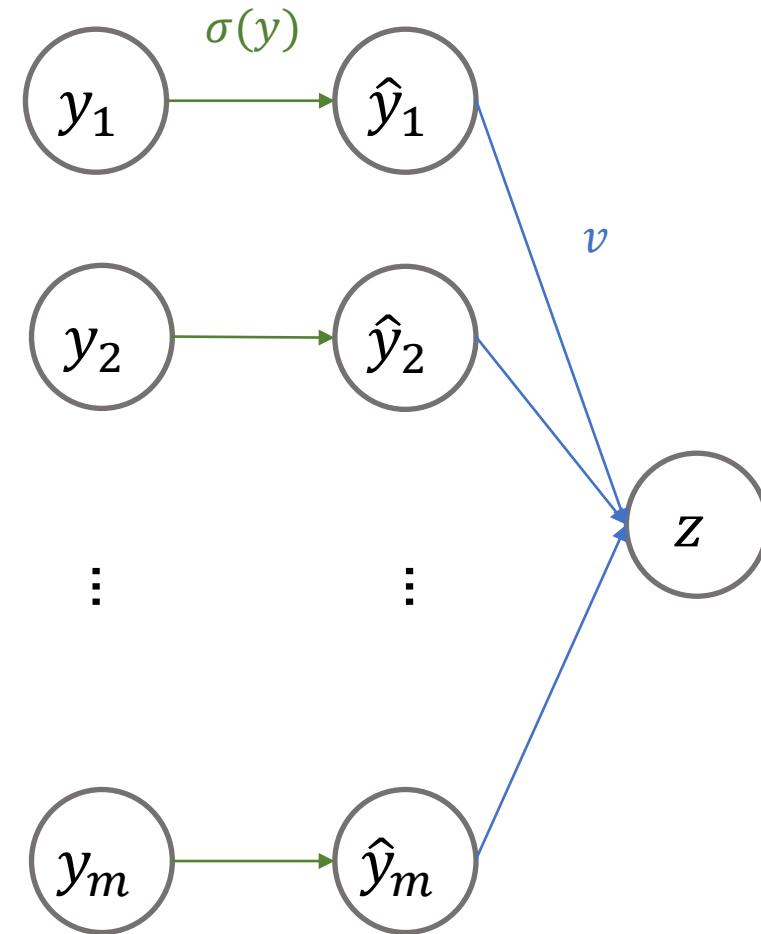
$$\hat{y} = \sigma(y) \quad z = \sum_{i=1}^m v_i \hat{y}_i = v^\top \hat{y}$$

$$\frac{\partial z}{\partial y} = \frac{\partial z}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial y}$$

$\in \mathbb{R}^{1 \times m}$ (top left)
 $\in \mathbb{R}^{m \times m}$ (top right)
 $\in \mathbb{R}^{1 \times m}$ (bottom)

$$\frac{\partial \hat{y}}{\partial y} = \begin{bmatrix} \frac{\partial \hat{y}_1}{\partial y_1} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \frac{\partial \hat{y}_m}{\partial y_m} \end{bmatrix}$$

$$= \begin{bmatrix} \sigma(y_1)(1 - \sigma(y_1)) & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \sigma(y_m)(1 - \sigma(y_m)) \end{bmatrix}$$



Vector Jacobian Product (VJP)

- Elementwise activation functions

$$\frac{\partial z}{\partial y} = \frac{\partial z}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial y} \quad \frac{\partial \hat{y}}{\partial y} = \begin{bmatrix} \frac{\partial \hat{y}_1}{\partial y_1} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \frac{\partial \hat{y}_m}{\partial y_m} \end{bmatrix} = \begin{bmatrix} \sigma(y_1)(1 - \sigma(y_1)) & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \sigma(y_m)(1 - \sigma(y_m)) \end{bmatrix}$$

$\in \mathbb{R}^{1 \times m}$ $\in \mathbb{R}^{m \times m}$ $\in \mathbb{R}^{1 \times m}$

Element-wise product

$$\frac{\partial z}{\partial y} = \frac{\partial z}{\partial \hat{y}} \odot \left(\sigma(y)(1 - \sigma(y)) \right)^T$$

$\in \mathbb{R}^{1 \times m}$ $\in \mathbb{R}^{1 \times m}$

Automatic Differentiation

Automatic Differentiation (AD)

- A procedure for automatic evaluation of derivatives of arbitrary algebraic functions
- Backpropagation == reverse-mode AD

Reverse-Mode AD (a.k.a Backpropagation)

$$f: \mathbb{R}^{n_1} \rightarrow \mathbb{R}^{n_2}$$

$$b = f(a)$$

$$g: \mathbb{R}^{n_2} \rightarrow \mathbb{R}^{n_3}$$

$$c = g(b)$$

$$d = h(c)$$

$$h: \mathbb{R}^{n_3} \rightarrow \mathbb{R}^{n_4}$$

$$e = i(d)$$

$$i: \mathbb{R}^{n_4} \rightarrow \mathbb{R}$$

$$\frac{\partial e}{\partial a} ?$$

Loss function:
scalar function

Reverse-Mode AD (a.k.a Backpropagation)

$$f: \mathbb{R}^{n_1} \rightarrow \mathbb{R}^{n_2}$$

$$b = f(a)$$

$$g: \mathbb{R}^{n_2} \rightarrow \mathbb{R}^{n_3}$$

$$c = g(b)$$

$$d = h(c)$$

$$h: \mathbb{R}^{n_3} \rightarrow \mathbb{R}^{n_4}$$

$$e = i(d)$$

$$i: \mathbb{R}^{n_4} \rightarrow \mathbb{R}$$

Loss function:
scalar function

$$\frac{\partial e}{\partial d} = \overset{\in \mathbb{R}^{1 \times 1}}{\frac{\partial e}{\partial e}} \overset{\in \mathbb{R}^{1 \times n_4}}{\frac{\partial e}{\partial d}} = \mathbf{1} \frac{\partial e}{\partial d}$$

Reverse-Mode AD (a.k.a Backpropagation)

$$f: \mathbb{R}^{n_1} \rightarrow \mathbb{R}^{n_2}$$

$$b = f(a)$$

$$g: \mathbb{R}^{n_2} \rightarrow \mathbb{R}^{n_3}$$

$$c = g(b)$$

$$h: \mathbb{R}^{n_3} \rightarrow \mathbb{R}^{n_4}$$

$$d = h(c)$$

$$e = i(d)$$

$$i: \mathbb{R}^{n_4} \rightarrow \mathbb{R}$$

Loss function:
scalar function

$$\frac{\partial e}{\partial d} = \frac{\overset{\in \mathbb{R}^{1 \times 1}}{\partial e}}{\overset{\in \mathbb{R}^{1 \times n_4}}{\partial d}} = 1 \frac{\partial e}{\partial d}$$

$$\frac{\partial e}{\partial c} = \frac{\overset{\in \mathbb{R}^{1 \times n_4}}{\partial e}}{\overset{\in \mathbb{R}^{n_4 \times n_3}}{\partial d}} \frac{\partial d}{\partial c} = \frac{\partial e}{\partial d} \frac{\partial d}{\partial c}$$

Reverse-Mode AD (a.k.a Backpropagation)

$$f: \mathbb{R}^{n_1} \rightarrow \mathbb{R}^{n_2}$$

$$b = f(a)$$

$$g: \mathbb{R}^{n_2} \rightarrow \mathbb{R}^{n_3}$$

$$c = g(b)$$

$$h: \mathbb{R}^{n_3} \rightarrow \mathbb{R}^{n_4}$$

$$d = h(c)$$

$$e = i(d)$$

$$i: \mathbb{R}^{n_4} \rightarrow \mathbb{R}$$

Loss function:
scalar function

$$\frac{\partial e}{\partial d} = \frac{\overset{\in \mathbb{R}^{1 \times 1}}{\partial e}}{\overset{\in \mathbb{R}^{1 \times n_4}}{\partial d}} = 1 \frac{\partial e}{\partial d}$$

$$\frac{\partial e}{\partial c} = \frac{\overset{\in \mathbb{R}^{1 \times n_4}}{\partial e}}{\overset{\in \mathbb{R}^{n_4 \times n_3}}{\partial d}} \frac{\partial d}{\partial c} = \frac{\partial e}{\partial d} \frac{\partial d}{\partial c}$$

$$\frac{\partial e}{\partial b} = \frac{\overset{\in \mathbb{R}^{1 \times n_3}}{\partial e}}{\overset{\in \mathbb{R}^{n_3 \times n_2}}{\partial d}} \frac{\partial d}{\partial c} \frac{\partial c}{\partial b} = \frac{\partial e}{\partial c} \frac{\partial c}{\partial b}$$

Reverse-Mode AD (a.k.a Backpropagation)

$$f: \mathbb{R}^{n_1} \rightarrow \mathbb{R}^{n_2}$$

$$b = f(a)$$

$$g: \mathbb{R}^{n_2} \rightarrow \mathbb{R}^{n_3}$$

$$c = g(b)$$

$$h: \mathbb{R}^{n_3} \rightarrow \mathbb{R}^{n_4}$$

$$d = h(c)$$

$$e = i(d)$$

$$i: \mathbb{R}^{n_4} \rightarrow \mathbb{R}$$

Loss function:
scalar function

$$\frac{\partial e}{\partial d} = \frac{\partial e}{\partial e} \frac{\partial e}{\partial d} = 1 \frac{\partial e}{\partial d}$$

$\in \mathbb{R}^{1 \times 1} \quad \in \mathbb{R}^{1 \times n_4}$

Vector-Jacobian
Product (VJP)

$$\frac{\partial e}{\partial c} = \frac{\partial e}{\partial e} \frac{\partial e}{\partial d} \frac{\partial d}{\partial c} = \frac{\partial e}{\partial d} \frac{\partial d}{\partial c}$$

$\in \mathbb{R}^{1 \times n_4} \quad \in \mathbb{R}^{n_4 \times n_3}$

$$\frac{\partial e}{\partial b} = \frac{\partial e}{\partial e} \frac{\partial e}{\partial d} \frac{\partial d}{\partial c} \frac{\partial c}{\partial b} = \frac{\partial e}{\partial c} \frac{\partial c}{\partial b}$$

$\in \mathbb{R}^{1 \times n_3} \quad \in \mathbb{R}^{n_3 \times n_2}$

$$\frac{\partial e}{\partial a} = \frac{\partial e}{\partial e} \frac{\partial e}{\partial d} \frac{\partial d}{\partial c} \frac{\partial c}{\partial b} \frac{\partial b}{\partial a} = \frac{\partial e}{\partial b} \frac{\partial b}{\partial a}$$

$\in \mathbb{R}^{1 \times n_2} \quad \in \mathbb{R}^{n_2 \times n_1}$

Reverse-Mode AD (a.k.a Backpropagation)

$$f: \mathbb{R}^{n_1} \rightarrow \mathbb{R}^{n_2}$$

$$b = f(a)$$

$$g: \mathbb{R}^{n_2} \rightarrow \mathbb{R}^{n_3}$$

$$c = g(b)$$

$$h: \mathbb{R}^{n_3} \rightarrow \mathbb{R}^{n_4}$$

$$d = h(c)$$

$$e = i(d)$$

$$i: \mathbb{R}^{n_4} \rightarrow \mathbb{R}^3$$

What if i is
vector-valued function?

$$\frac{\partial e}{\partial d} = \frac{\partial e}{\partial c} \frac{\partial c}{\partial d} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \frac{\partial e}{\partial d}$$

$\in \mathbb{R}^{3 \times 3} \quad \in \mathbb{R}^{3 \times n_4}$

3 X
Computation

$$\frac{\partial e}{\partial c} = \frac{\partial e}{\partial d} \frac{\partial d}{\partial c} = \frac{\partial e}{\partial d} \frac{\partial d}{\partial c}$$

$\in \mathbb{R}^{3 \times n_4} \quad \in \mathbb{R}^{n_4 \times n_3}$

$$\frac{\partial e}{\partial b} = \frac{\partial e}{\partial c} \frac{\partial c}{\partial d} \frac{\partial d}{\partial c} \frac{\partial c}{\partial b} = \frac{\partial e}{\partial c} \frac{\partial c}{\partial b}$$

$\in \mathbb{R}^{3 \times n_3} \quad \in \mathbb{R}^{n_3 \times n_2}$

$$\frac{\partial e}{\partial a} = \frac{\partial e}{\partial c} \frac{\partial c}{\partial d} \frac{\partial d}{\partial c} \frac{\partial c}{\partial b} \frac{\partial b}{\partial a} = \frac{\partial e}{\partial c} \frac{\partial b}{\partial a}$$

$\in \mathbb{R}^{3 \times n_2} \quad \in \mathbb{R}^{n_2 \times n_1}$

Reverse-Mode AD (a.k.a Backpropagation)

$$f: \mathbb{R}^{n_1} \rightarrow \mathbb{R}^{n_2}$$

$$b = f(a)$$

$$g: \mathbb{R}^{n_2} \rightarrow \mathbb{R}^{n_3}$$

$$c = g(b)$$

$$h: \mathbb{R}^{n_3} \rightarrow \mathbb{R}^{n_4}$$

$$d = h(c)$$

$$e = i(d)$$

$$i: \mathbb{R}^{n_4} \rightarrow \mathbb{R}^3$$

$$\frac{\partial e}{\partial d} = \frac{\partial e}{\partial e} \frac{\partial e}{\partial d} = \begin{matrix} \in \mathbb{R}^{3 \times 3} & \in \mathbb{R}^{3 \times n_4} \\ \boxed{\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}} \end{matrix} \frac{\partial e}{\partial d}$$

$\frac{\partial e_1}{\partial d}$

Forward-Mode AD

single variable

$$f: \mathbb{R} \rightarrow \mathbb{R}^{n_1}$$

$$b = f(a)$$

$$g: \mathbb{R}^{n_1} \rightarrow \mathbb{R}^{n_2}$$

$$c = g(b)$$

$$h: \mathbb{R}^{n_2} \rightarrow \mathbb{R}^{n_3}$$

$$d = h(c)$$

$$e = L(d)$$

$$L: \mathbb{R}^{n_3} \rightarrow \mathbb{R}^{n_4}$$

$$\frac{\partial b}{\partial a} = \frac{\partial b}{\partial a} \frac{\partial a}{\partial a} = \frac{\partial b}{\partial a} \mathbf{1}$$

$\in \mathbb{R}^{n_1 \times 1} \quad \in \mathbb{R}^{1 \times 1}$

Forward-Mode AD

single variable

$$f: \mathbb{R} \rightarrow \mathbb{R}^{n_1}$$

$$b = f(a)$$

$$g: \mathbb{R}^{n_1} \rightarrow \mathbb{R}^{n_2}$$

$$c = g(b)$$

$$h: \mathbb{R}^{n_2} \rightarrow \mathbb{R}^{n_3}$$

$$d = h(c)$$

$$e = L(d)$$

$$L: \mathbb{R}^{n_3} \rightarrow \mathbb{R}^{n_4}$$

$$\frac{\partial b}{\partial a} = \frac{\partial b}{\partial a} \frac{\partial a}{\partial a} = \frac{\partial b}{\partial a} \mathbf{1}$$

$\in \mathbb{R}^{n_1 \times 1} \quad \in \mathbb{R}^{1 \times 1}$

$$\frac{\partial c}{\partial a} = \frac{\partial c}{\partial b} \frac{\partial b}{\partial a} \frac{\partial a}{\partial a} = \frac{\partial c}{\partial b} \frac{\partial b}{\partial a}$$

$\in \mathbb{R}^{n_2 \times n_1} \quad \in \mathbb{R}^{n_1 \times 1}$

Forward-Mode AD

single variable

$$f: \mathbb{R} \rightarrow \mathbb{R}^{n_1}$$

$$b = f(a)$$

$$g: \mathbb{R}^{n_1} \rightarrow \mathbb{R}^{n_2}$$

$$c = g(b)$$

$$h: \mathbb{R}^{n_2} \rightarrow \mathbb{R}^{n_3}$$

$$d = h(c)$$

$$e = L(d)$$

$$L: \mathbb{R}^{n_3} \rightarrow \mathbb{R}^{n_4}$$

$$\frac{\partial b}{\partial a} = \frac{\partial b}{\partial a} \frac{\partial a}{\partial a} = \frac{\partial b}{\partial a} \mathbf{1}$$

$\in \mathbb{R}^{n_1 \times 1} \in \mathbb{R}^{1 \times 1}$

$$\frac{\partial c}{\partial a} = \frac{\partial c}{\partial b} \frac{\partial b}{\partial a} \frac{\partial a}{\partial a} = \frac{\partial c}{\partial b} \frac{\partial b}{\partial a}$$

$\in \mathbb{R}^{n_2 \times n_1} \in \mathbb{R}^{n_1 \times 1}$

$$\frac{\partial d}{\partial a} = \frac{\partial d}{\partial c} \frac{\partial c}{\partial b} \frac{\partial b}{\partial a} \frac{\partial a}{\partial a} = \frac{\partial d}{\partial c} \frac{\partial c}{\partial a}$$

$\in \mathbb{R}^{n_3 \times n_2} \in \mathbb{R}^{n_2 \times 1}$

Forward-Mode AD

single variable

$$f: \mathbb{R} \rightarrow \mathbb{R}^{n_1}$$

$$b = f(a)$$

$$g: \mathbb{R}^{n_1} \rightarrow \mathbb{R}^{n_2}$$

$$c = g(b)$$

$$h: \mathbb{R}^{n_2} \rightarrow \mathbb{R}^{n_3}$$

$$d = h(c)$$

$$L: \mathbb{R}^{n_3} \rightarrow \mathbb{R}^{n_4}$$

$$e = L(d)$$

$$\frac{\partial b}{\partial a} = \frac{\partial b}{\partial a} \frac{\partial a}{\partial a} = \frac{\partial b}{\partial a} \mathbf{1}$$

$\in \mathbb{R}^{n_1 \times 1} \in \mathbb{R}^{1 \times 1}$

Jacobian-Vector
Product (JVP)

$$\frac{\partial c}{\partial a} = \frac{\partial c}{\partial b} \frac{\partial b}{\partial a} \frac{\partial a}{\partial a} = \frac{\partial c}{\partial b} \frac{\partial b}{\partial a}$$

$\in \mathbb{R}^{n_2 \times n_1} \in \mathbb{R}^{n_1 \times 1}$

$$\frac{\partial d}{\partial a} = \frac{\partial d}{\partial c} \frac{\partial c}{\partial b} \frac{\partial b}{\partial a} \frac{\partial a}{\partial a} = \frac{\partial d}{\partial c} \frac{\partial c}{\partial a}$$

$\in \mathbb{R}^{n_3 \times n_2} \in \mathbb{R}^{n_2 \times 1}$

$$\frac{\partial e}{\partial a} = \frac{\partial e}{\partial d} \frac{\partial d}{\partial c} \frac{\partial c}{\partial b} \frac{\partial b}{\partial a} \frac{\partial a}{\partial a} = \frac{\partial e}{\partial d} \frac{\partial d}{\partial a}$$

$\in \mathbb{R}^{n_4 \times n_3} \in \mathbb{R}^{n_3 \times 1}$

Forward-Mode AD

$$f: \mathbb{R}^3 \rightarrow \mathbb{R}^{n_1}$$

$$b = f(a)$$

$$g: \mathbb{R}^{n_1} \rightarrow \mathbb{R}^{n_2}$$

$$c = g(b)$$

$$h: \mathbb{R}^{n_2} \rightarrow \mathbb{R}^{n_3}$$

$$d = h(c)$$

$$e = L(d)$$

$$L: \mathbb{R}^{n_3} \rightarrow \mathbb{R}^{n_4}$$

What if input is
Multi-variables?

$$\frac{\partial b}{\partial a} = \frac{\partial b}{\partial a} \frac{\partial a}{\partial a} = \frac{\partial b}{\partial a} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$\in \mathbb{R}^{n_1 \times 3} \quad \in \mathbb{R}^{3 \times 3}$

3 X
Computation

$$\frac{\partial c}{\partial a} = \frac{\partial c}{\partial b} \frac{\partial b}{\partial a} \frac{\partial a}{\partial a} = \frac{\partial c}{\partial b} \frac{\partial b}{\partial a}$$

$\in \mathbb{R}^{n_2 \times n_1} \quad \in \mathbb{R}^{n_1 \times 3}$

$$\frac{\partial d}{\partial a} = \frac{\partial d}{\partial c} \frac{\partial c}{\partial b} \frac{\partial b}{\partial a} \frac{\partial a}{\partial a} = \frac{\partial d}{\partial c} \frac{\partial c}{\partial a}$$

$\in \mathbb{R}^{n_3 \times n_2} \quad \in \mathbb{R}^{n_2 \times 3}$

$$\frac{\partial e}{\partial a} = \frac{\partial e}{\partial d} \frac{\partial d}{\partial c} \frac{\partial c}{\partial b} \frac{\partial b}{\partial a} \frac{\partial a}{\partial a} = \frac{\partial e}{\partial d} \frac{\partial d}{\partial a}$$

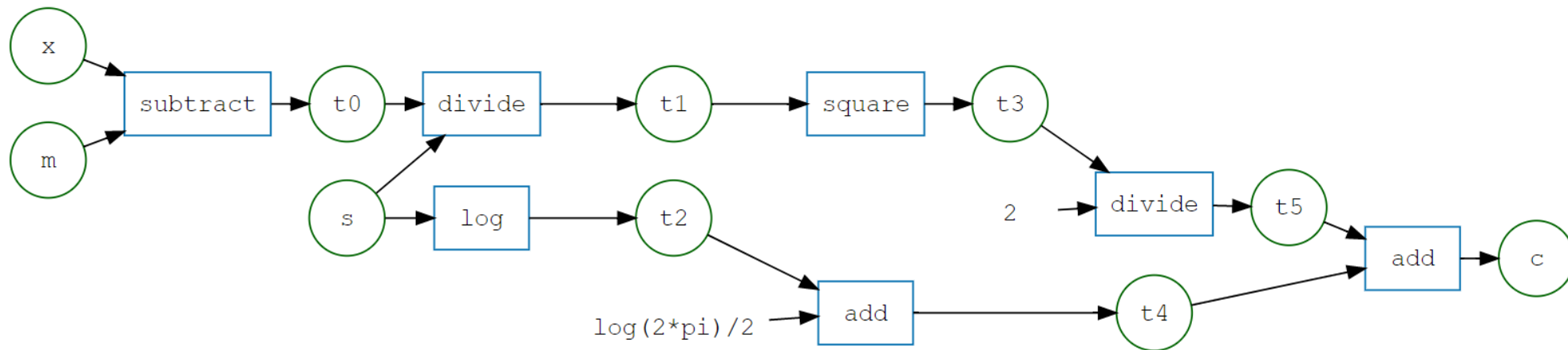
$\in \mathbb{R}^{n_4 \times n_3} \quad \in \mathbb{R}^{n_3 \times 3}$

Automatic Differentiation (AD)

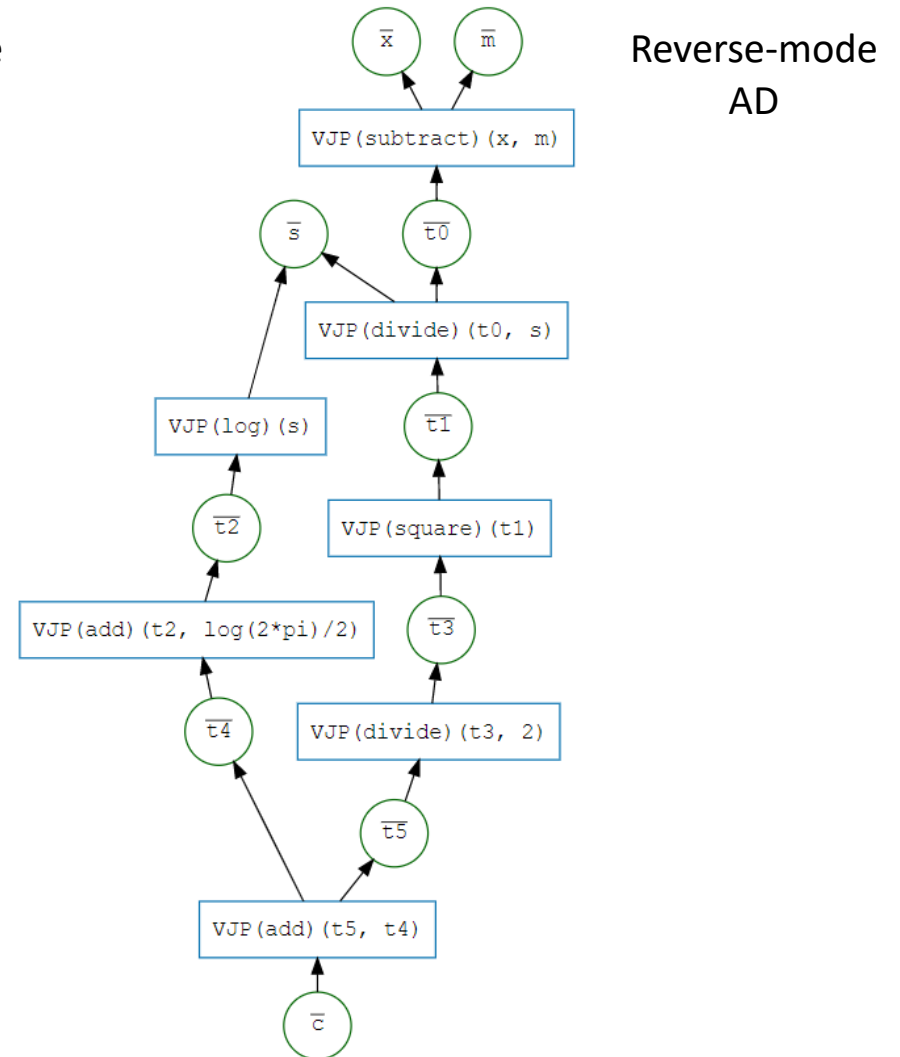
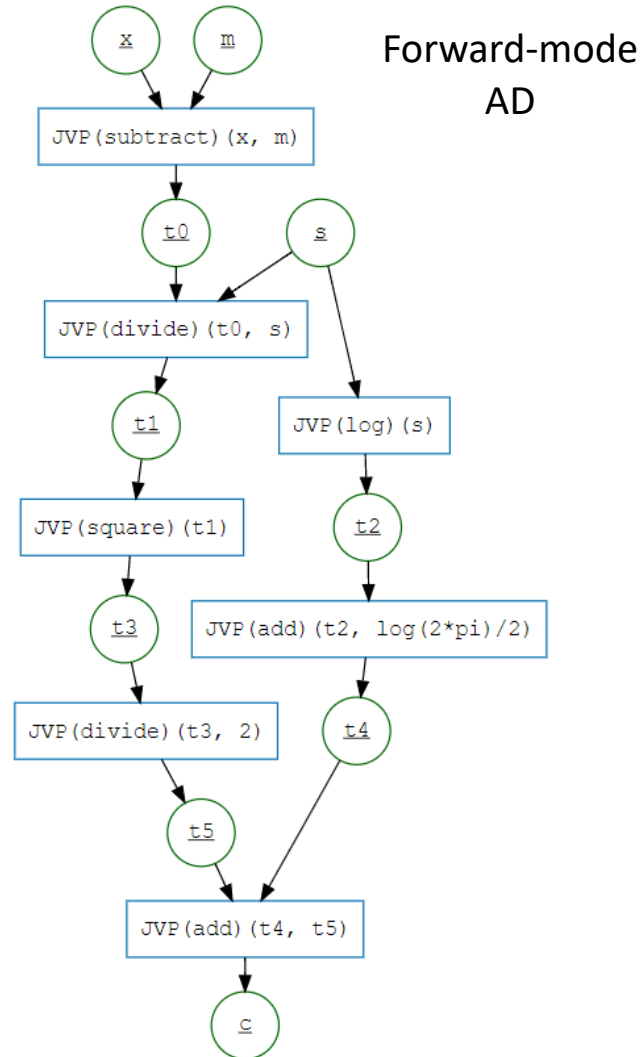
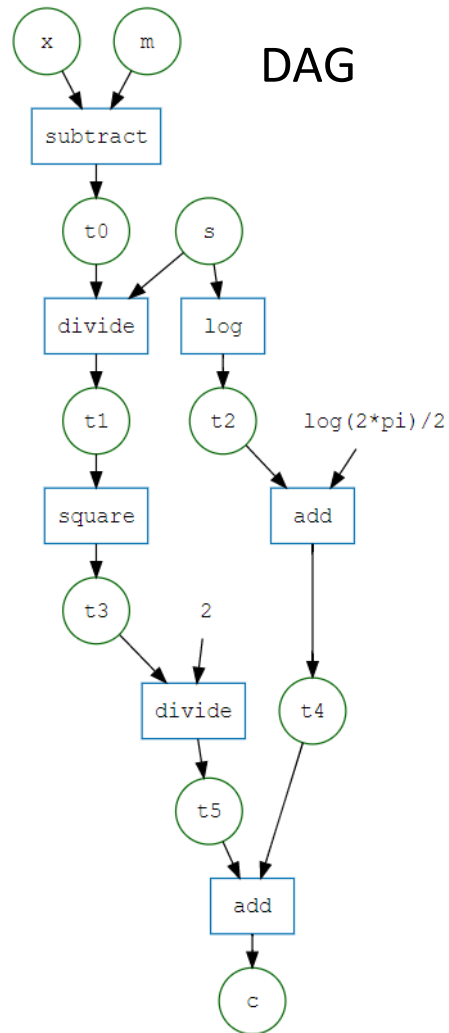
- For low dimensional outputs and high dimensional inputs
 - Objective function w/ deep neural networks
 - reverse-mode AD
- For high dimensional outputs and low dimensional inputs
 - Forward-mode AD

Computational Graph

```
t0 = x - m
t1 = t0 / s
t2 = np.log(s)
t3 = t1**2
t4 = t2 + np.log(2 * np.pi) / 2
t5 = t3 / 2
c = t4 + t5
```



Automatic Differentiation



References

- [mattjj/autodidact: A pedagogical implementation of Autograd \(github.com\)](#)
- [\[1502.05767\] Automatic differentiation in machine learning: a survey \(arxiv.org\)](#)
- [CSC321 Lecture 10: Automatic Differentiation \(toronto.edu\)](#)