

Foundations of Machine Learning (ECE 5984)

- Neural Networks -

Eunbyung Park

Assistant Professor

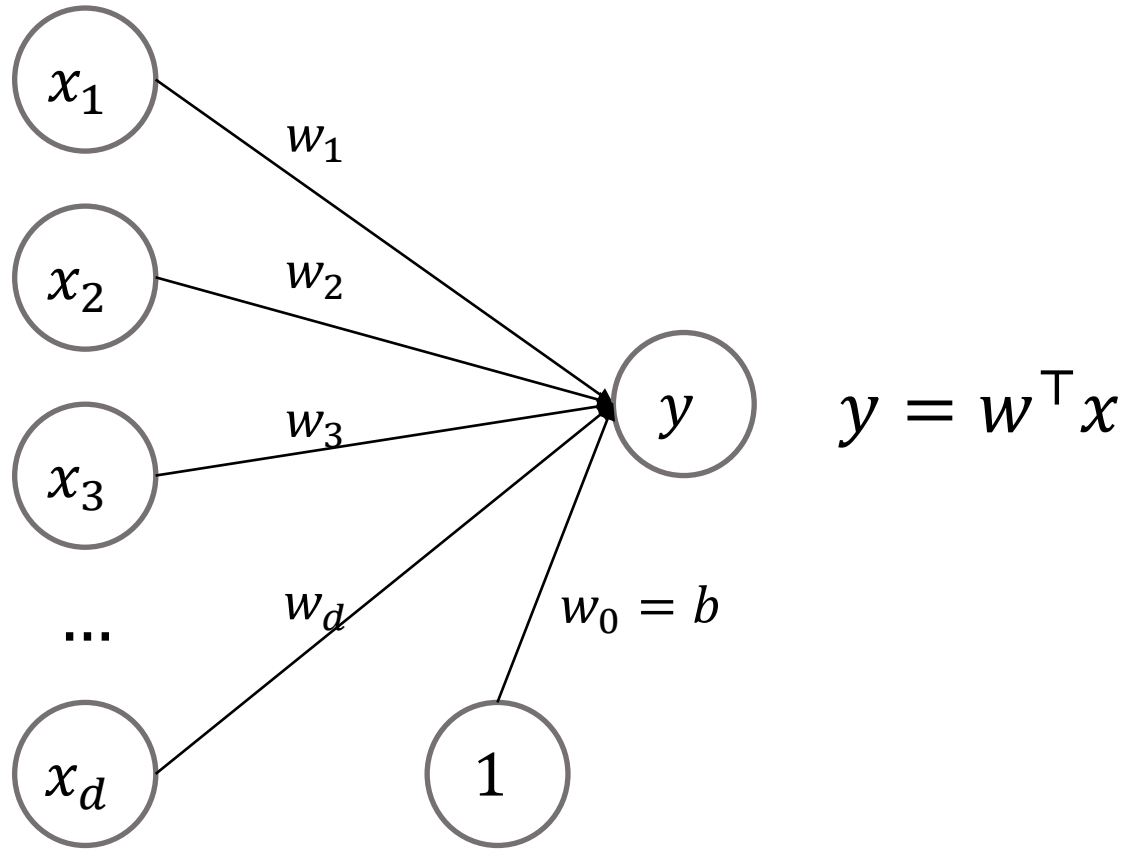
School of Electronic and Electrical Engineering

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

Multi-Layer Perceptron

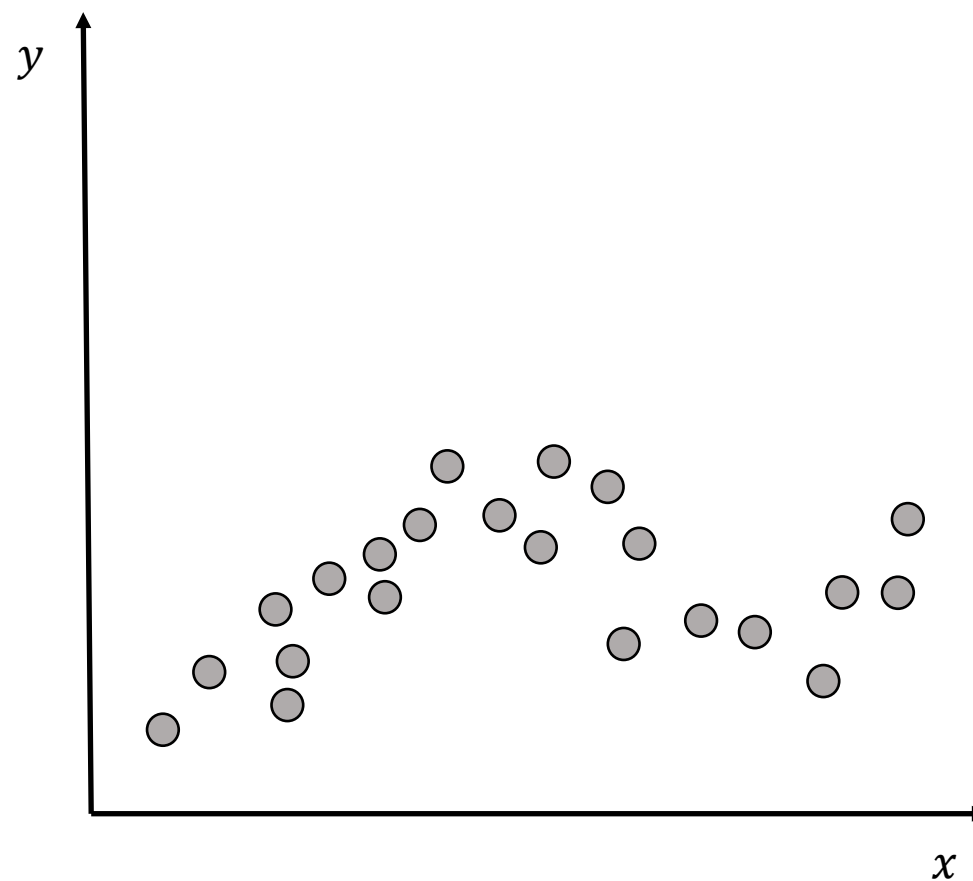
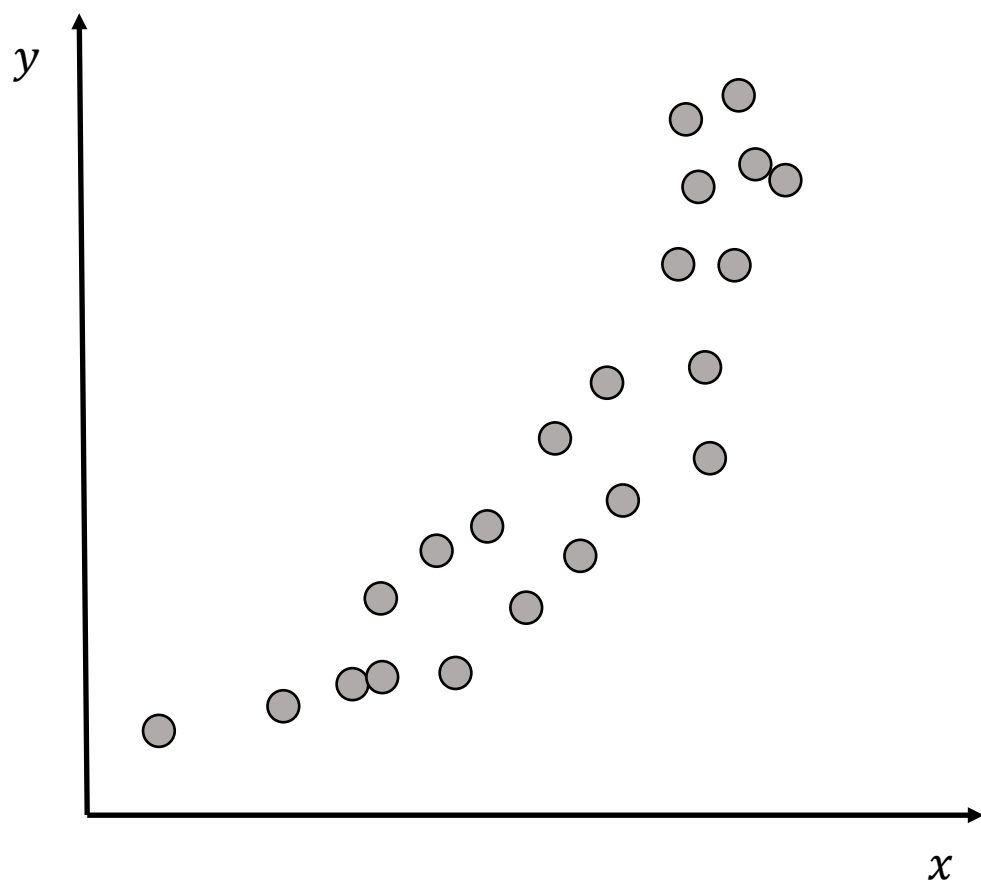
Linear Models as Shallow Neural Networks

- It is a single layer neural network



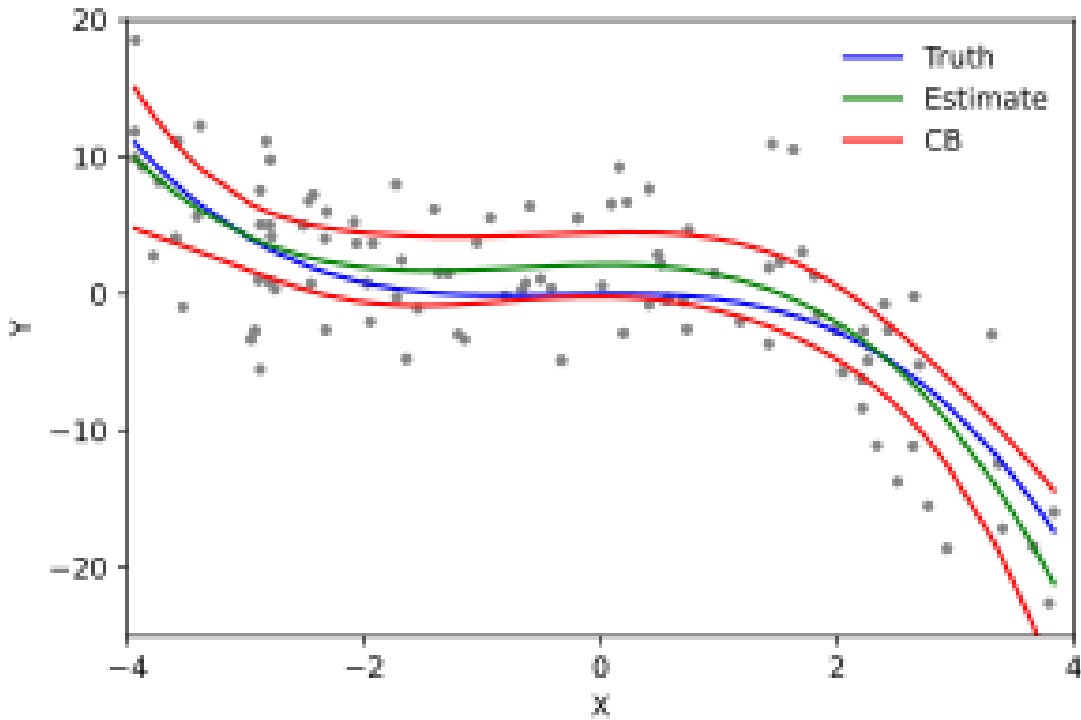
Linear Models

- Is linear model a good for all?



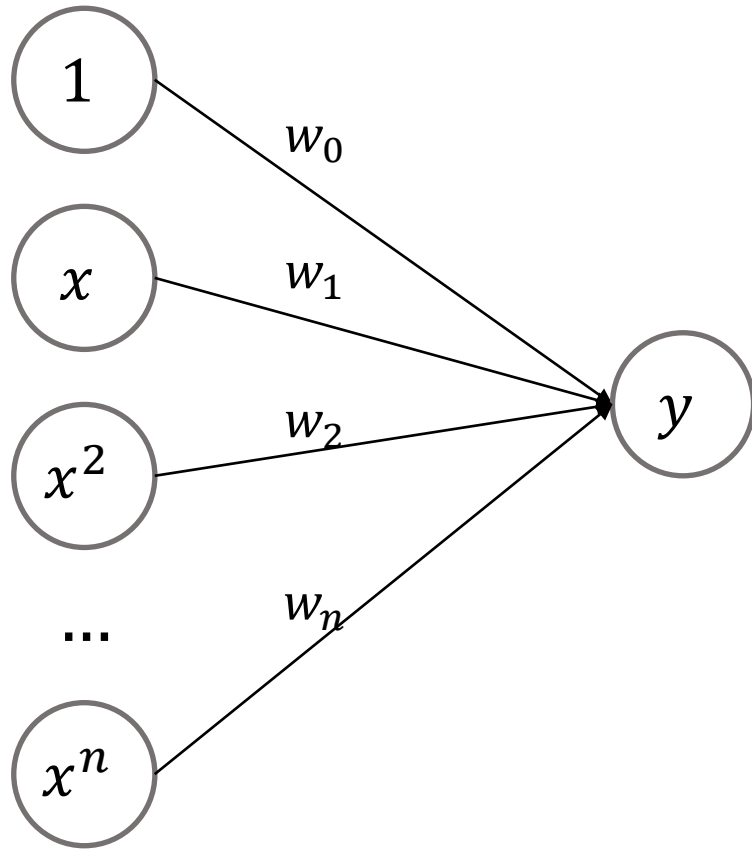
Nonlinear Models

- nth-degree Polynomial regression



$$f(x) = w_0 + w_1x + w_2x^2 + w_3x^3 + \cdots + w_nx^n$$

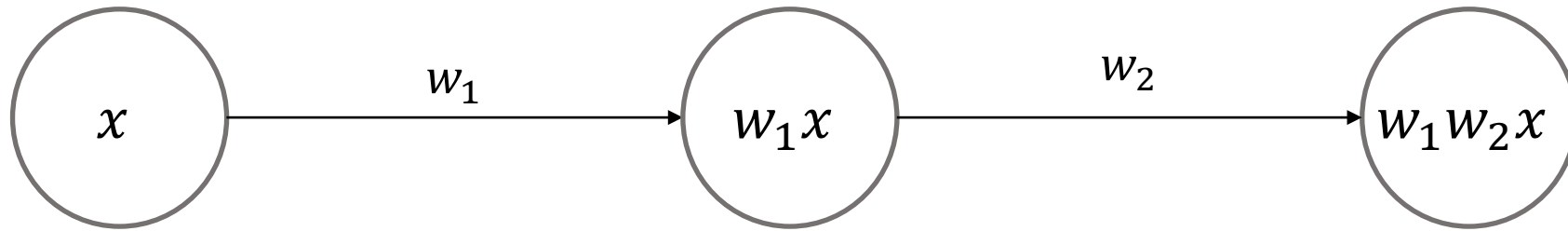
Polynomials as Neural Network



$$f(x) = w_0 + w_1x + w_2x^2 + w_3x^3 + \cdots + w_nx^n$$

- Feature engineering is hard
- Can we make it non-linear w/o feature engineering?

Feed-Forward Neural Network

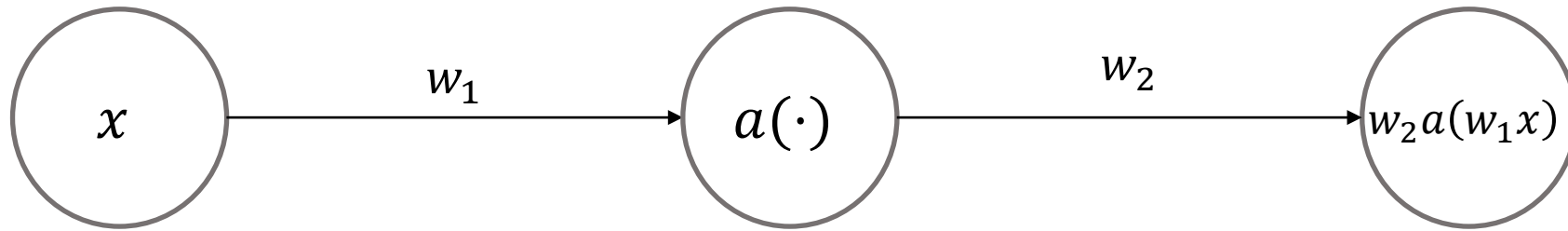


$$f(x) = w_1 w_2 x$$

Is it non-linear in x ?

Feed-Forward Neural Network

- Using non-linear activation function



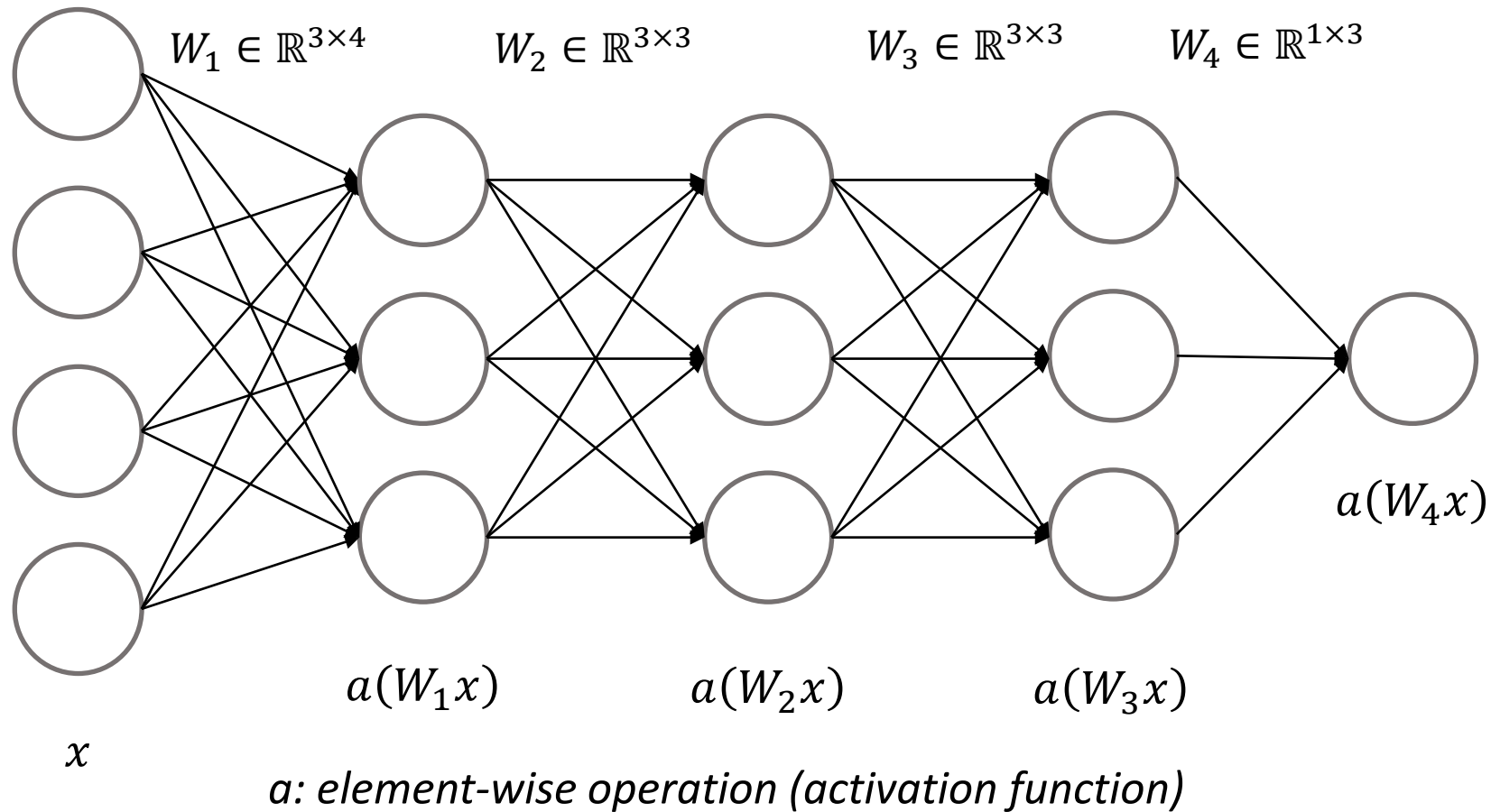
$$f(x) = w_2 a(w_1 x)$$

$$a(x) = \max(0, x) \quad (\text{Rectifier Linear Unit})$$

$$a(x) = \frac{1}{1 + e^{-x}} \quad (\text{Sigmoid})$$

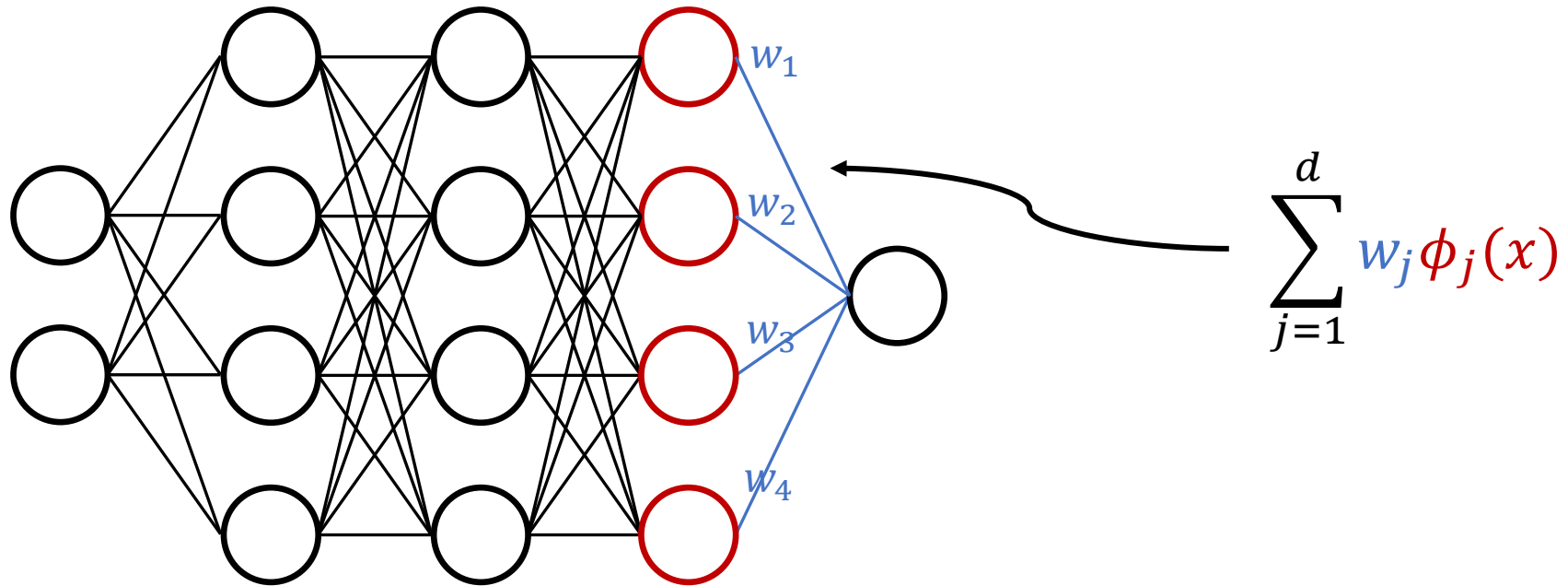
Feed-Forward Neural Network

- AKA, Multi-Layer Perceptron



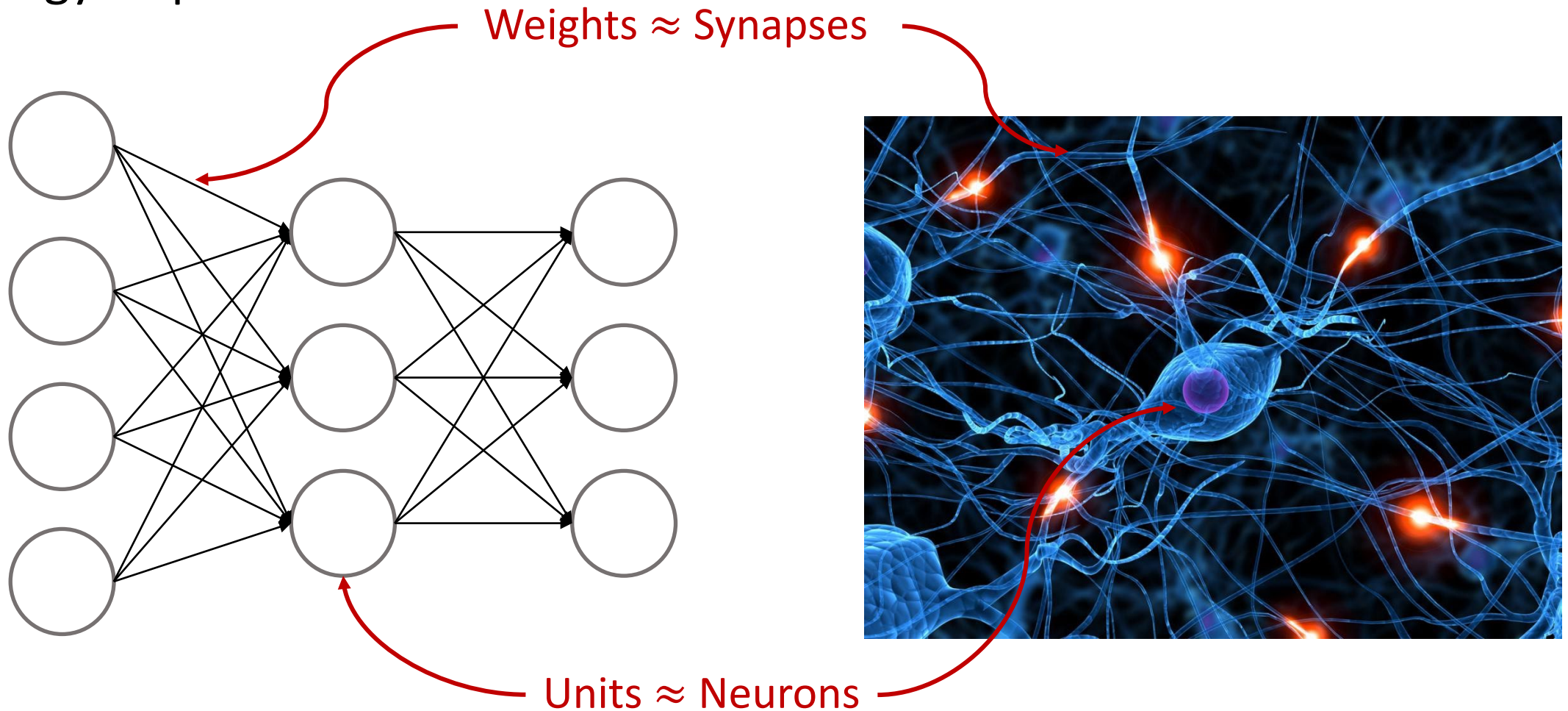
Feed-Forward Neural Network

- Connection to the Kernel Method



Feed-Forward Neural Network

- Biology Inspired



Feed-Forward Neural Network

- Regression with two layers MLP

$$D = \{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}$$

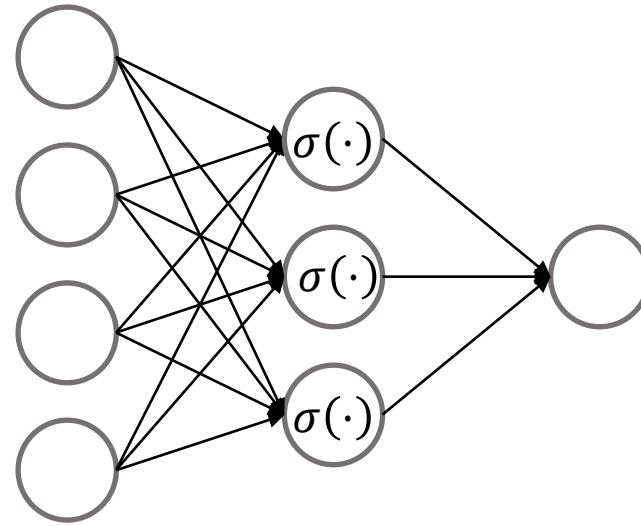
$$x^{(i)} \in \mathbb{R}^d, y^{(i)} \in \mathbb{R}, X \in \mathbb{R}^{N \times d}, Y \in \mathbb{R}^N$$

$$\theta = \{W_1, W_2\}, W_1 \in \mathbb{R}^{h \times d}, W_2 \in \mathbb{R}^{1 \times h}$$

$$f_{\theta}(x) = W_2 \sigma(W_1 x)$$

$$f_{\theta}: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$L(\theta) = \frac{1}{2} \sum_{i=1}^N (y^{(i)} - f_{\theta}(x^{(i)}))^2 = \frac{1}{2} (Y - \sigma(W_1 X^T)^T W_2^T)^T (Y - \sigma(W_1 X^T)^T W_2^T)$$



Feed-Forward Neural Network

- Regression with two layers MLP

$$D = \{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}$$

$$x^{(i)} \in \mathbb{R}^d, y^{(i)} \in \mathbb{R}, X \in \mathbb{R}^{N \times d}, Y \in \mathbb{R}^N$$

$$\theta = \{W_1, W_2\}, W_1 \in \mathbb{R}^{h \times d}, W_2 \in \mathbb{R}^{1 \times h}$$

$$f_{\theta}(x) = W_2 \sigma(W_1 x)$$

$$f_{\theta}: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$L(\theta) = \frac{1}{2} \sum_{i=1}^N (y^{(i)} - f_{\theta}(x^{(i)}))^2 = \frac{1}{2} (Y - \sigma(W_1 X^T)^T W_2^T)^T (Y - \sigma(W_1 X^T)^T W_2^T)$$

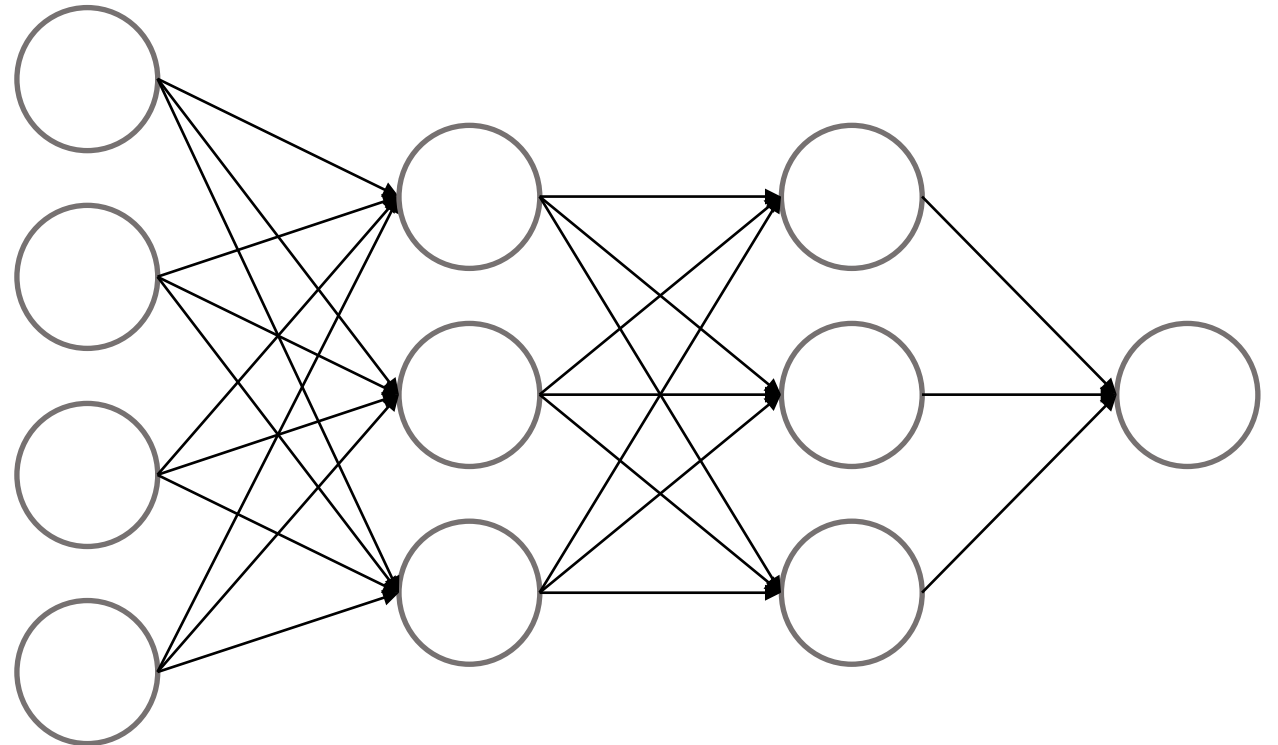
1. Can you take the gradients w.r.t θ ?
2. Does it have a closed form solution?
3. Is it a convex function?

Gradient Descent

- We are using *gradient descent* for training deep neural networks

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

(descent) (step-size) (gradient)



The Universal Approximator

The Universal Approximation Theorem

- A single hidden layer neural network can approximate any continuous function arbitrarily well, given enough hidden units.
- This holds for many different activation functions, e.g. sigmoid, tanh, ReLU, etc.

Cybenko Theorem

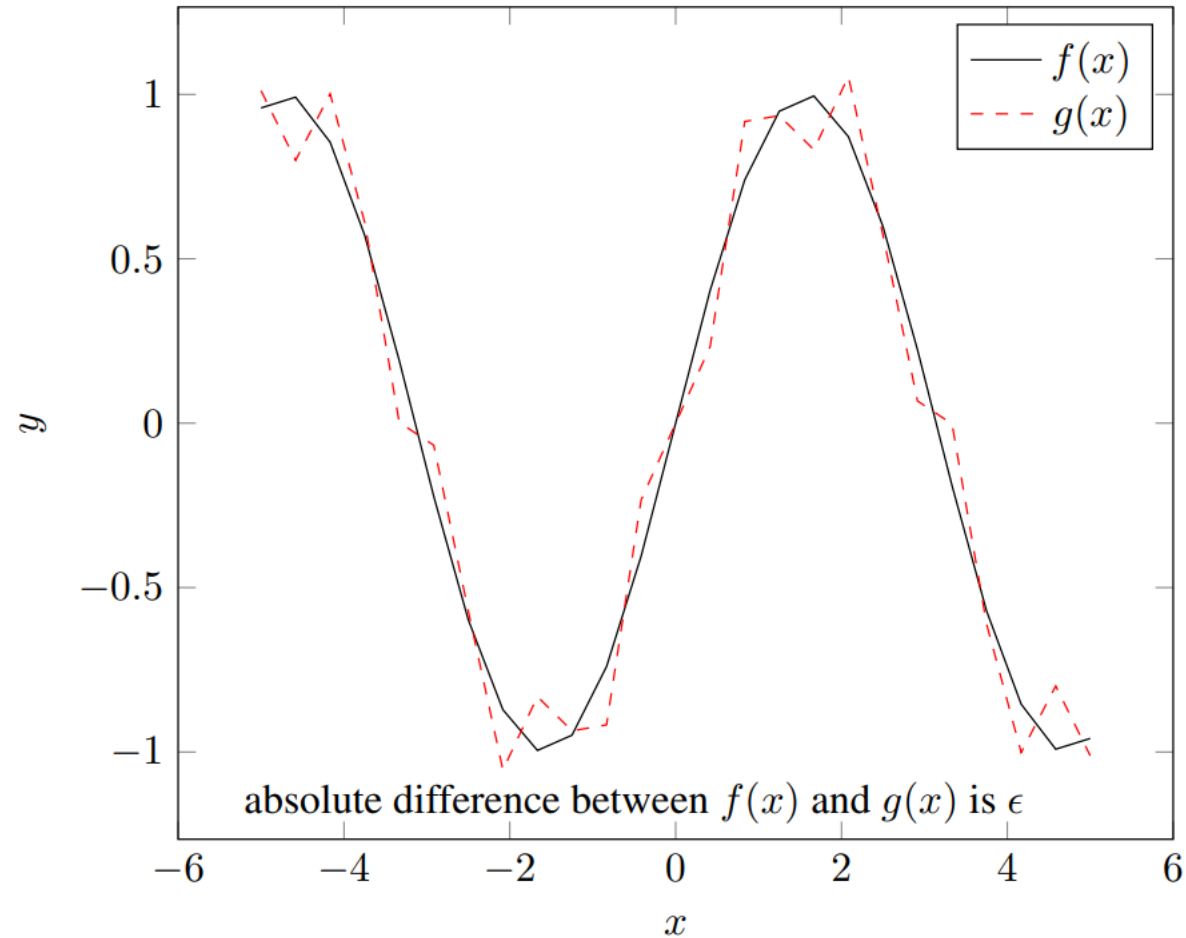
- Cybenko Approximation by Superposition of Sigmoidal Function

Let $C([0,1]^n)$ denote the set of all continuous function $[0,1]^n \rightarrow \mathbb{R}$, let σ be any sigmoidal activation function then the finite sum of the form $f(x) = \sum_{i=1}^N \alpha_i \sigma(w_i^\top x + b_i)$ is dense in $C([0,1]^n)$

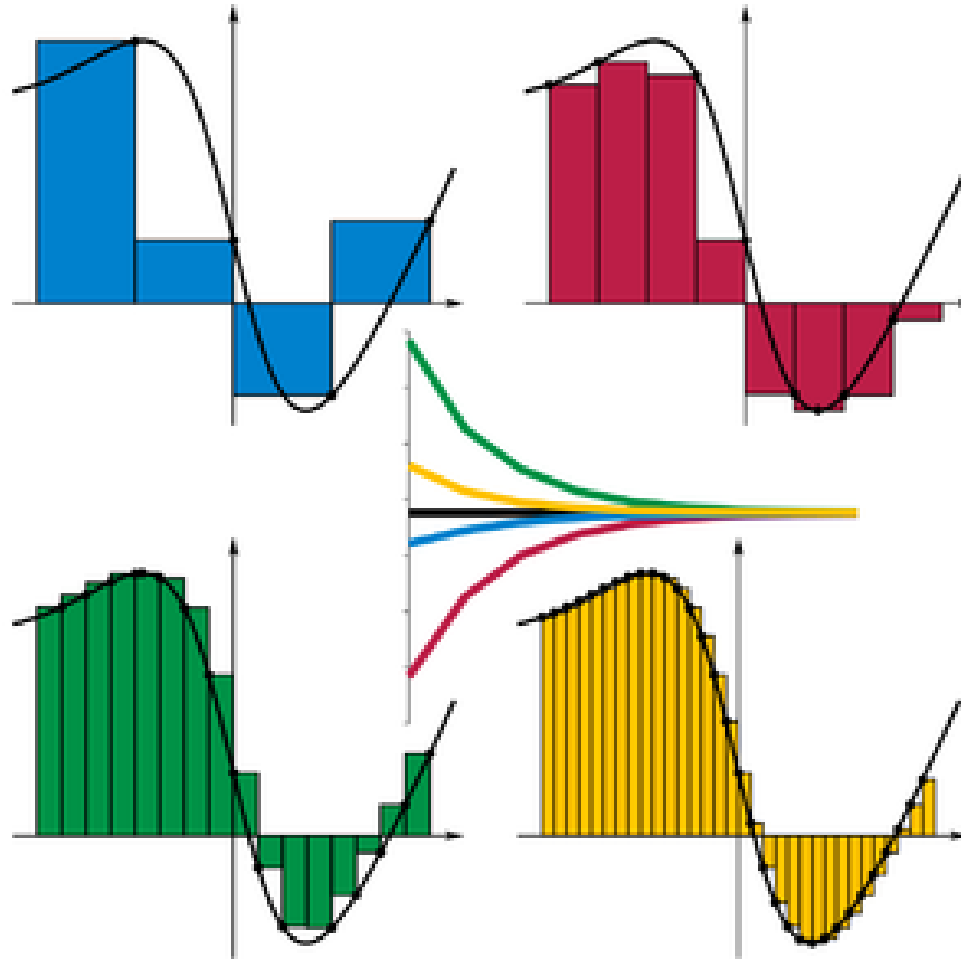
For any $g \in C([0,1]^n)$ and any $\epsilon > 0$, there exists $f: x \rightarrow \sum_{i=1}^N \alpha_i \sigma(w_i^\top x + b_i)$, such that $|f(x) - g(x)| < \epsilon$ for all $x \in [0,1]^n$.

Cybenko Theorem

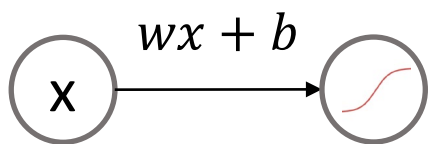
- Cybenko Approximation by Superposition of Sigmoidal Function



The Universal Approximation Theorem



The Universal Approximation Theorem



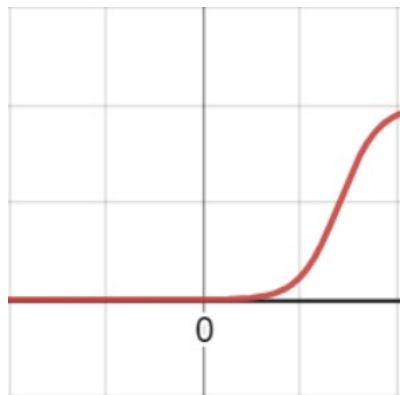
$$w = 5, b = 0$$



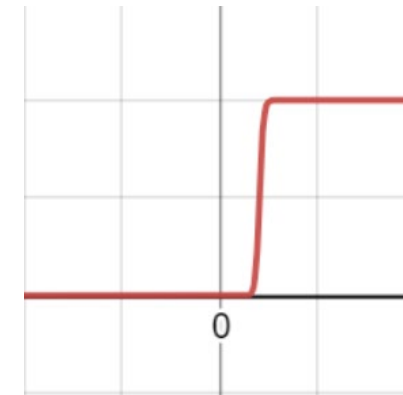
$$w = 5, b = 3$$



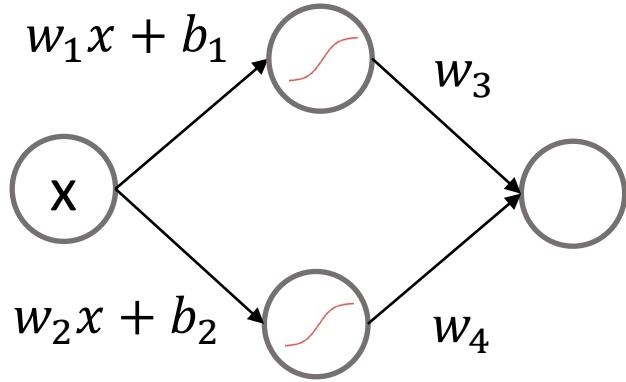
$$w = 10, b = -7$$



$$w = 100, b = -20$$



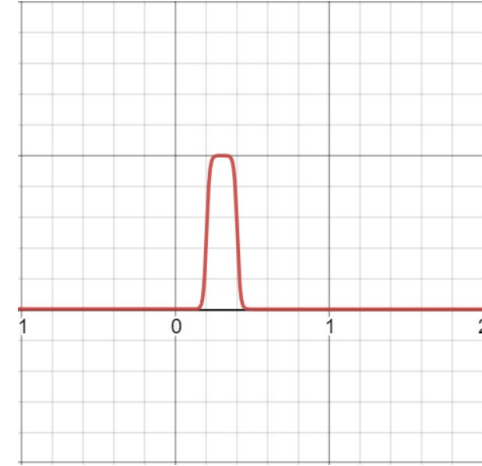
The Universal Approximation Theorem



$$w_1 = 100, b_1 = -20$$

$$w_2 = 100, b_2 = -40$$

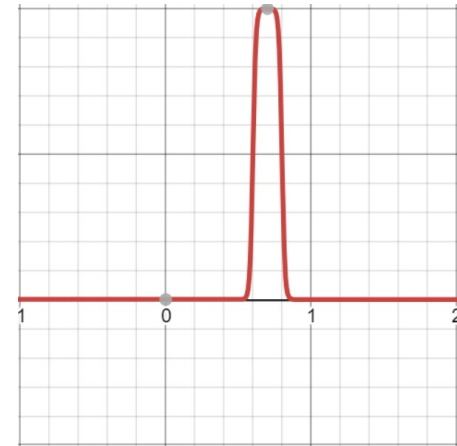
$$w_3 = 1, w_4 = -1$$



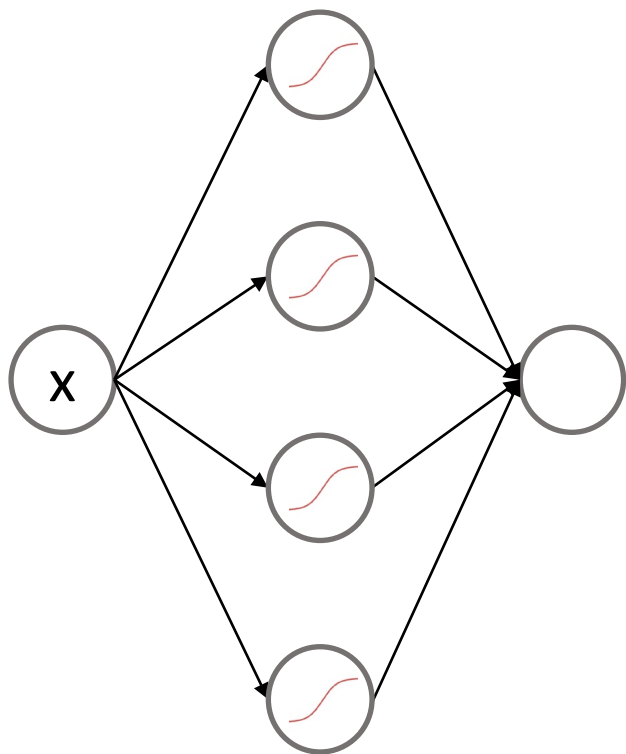
$$w_1 = 100, b_1 = -60$$

$$w_2 = 100, b_2 = -80$$

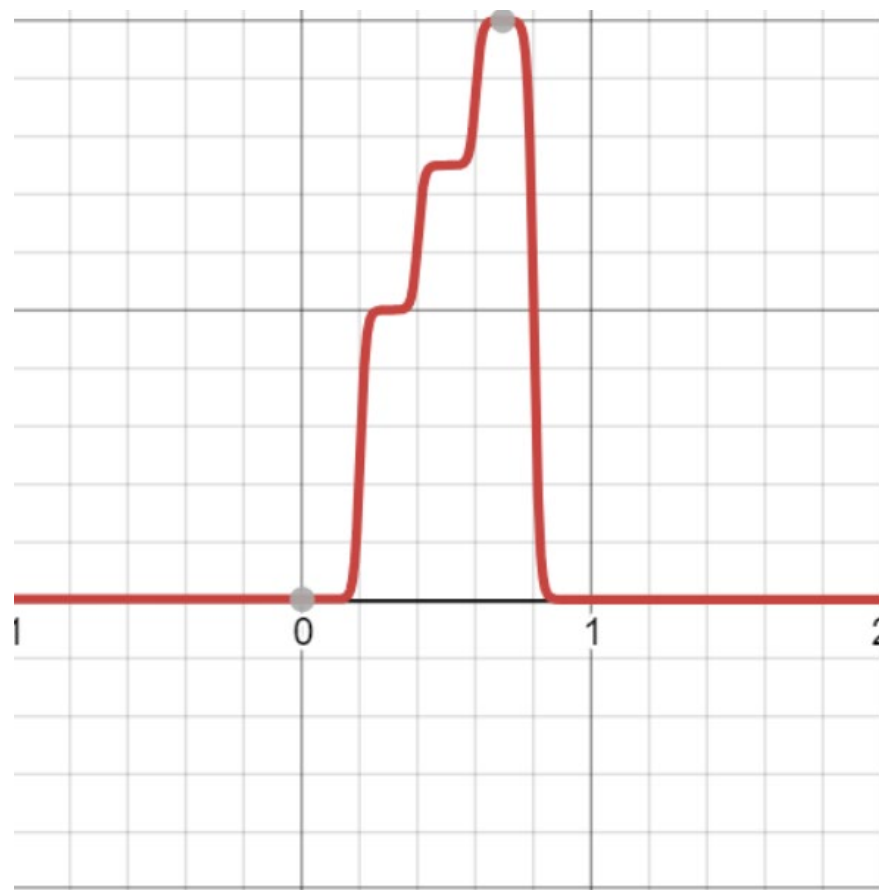
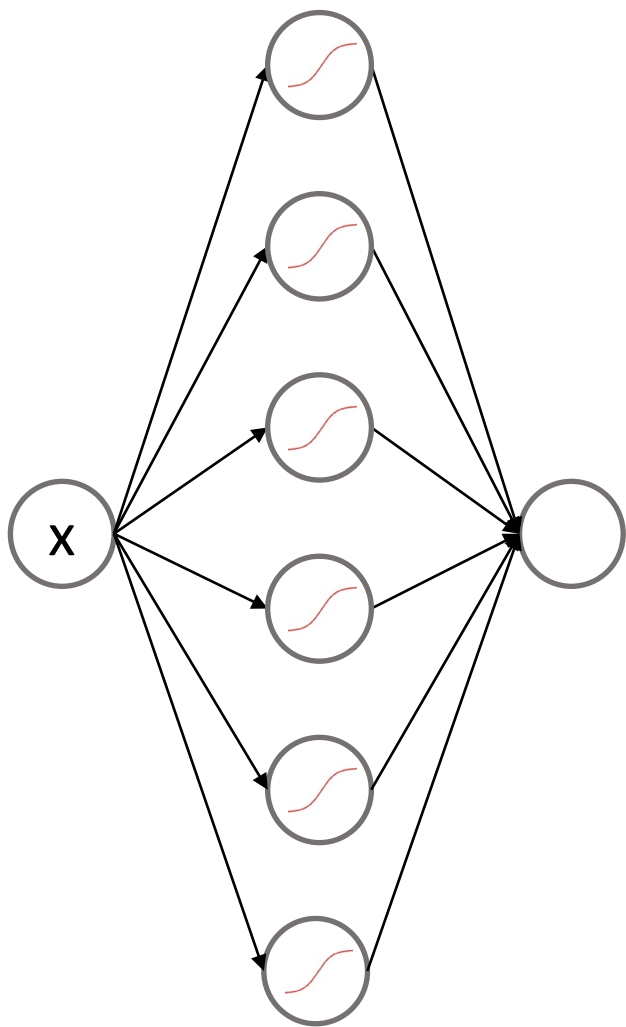
$$w_3 = 2, w_4 = -2$$



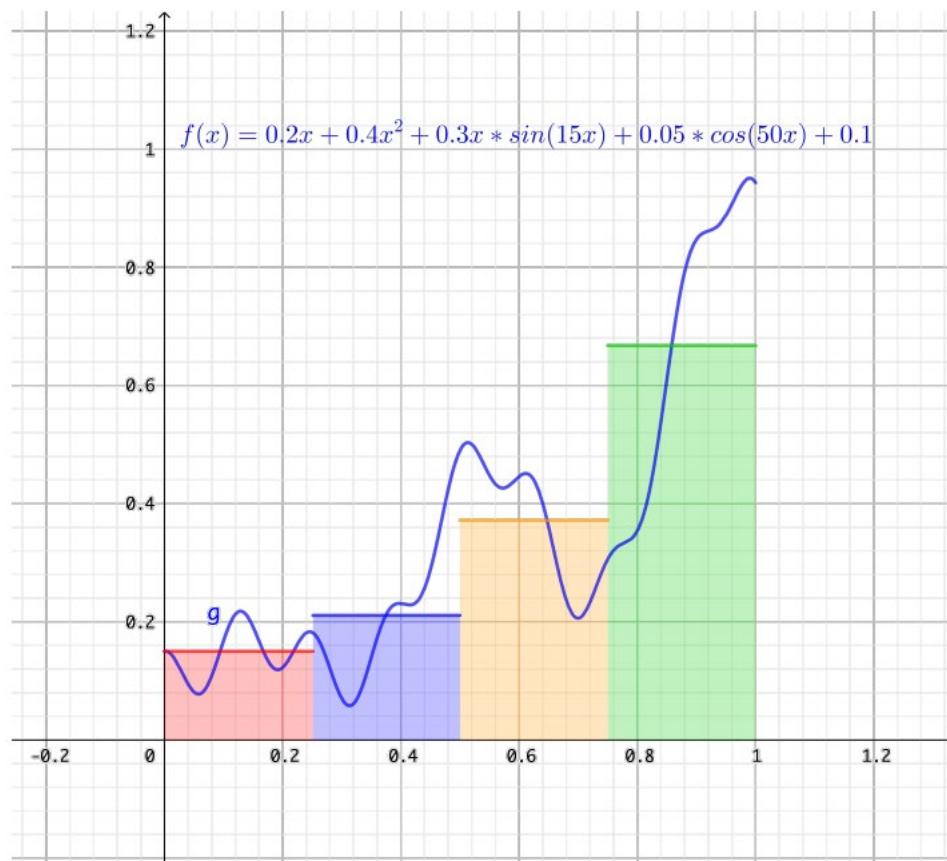
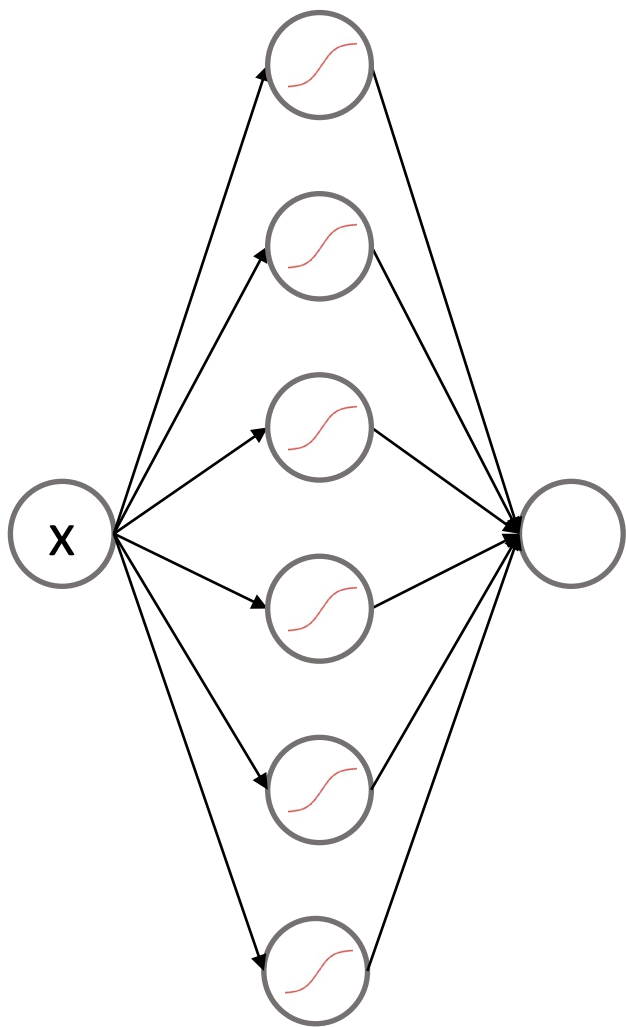
The Universal Approximation Theorem



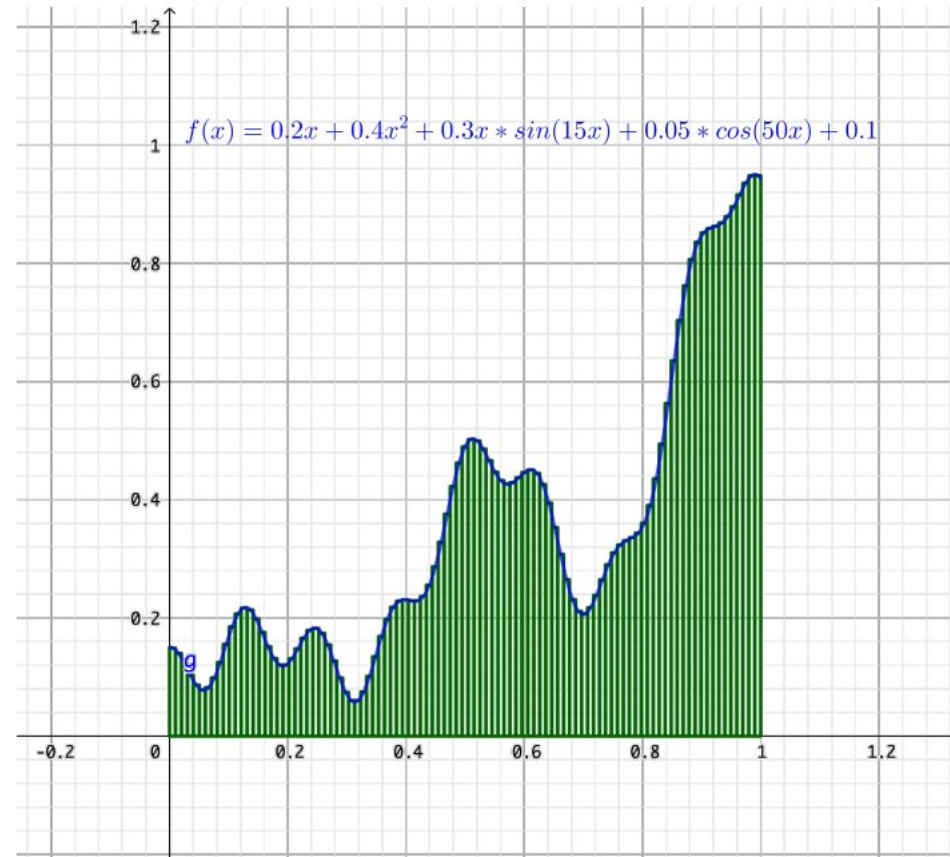
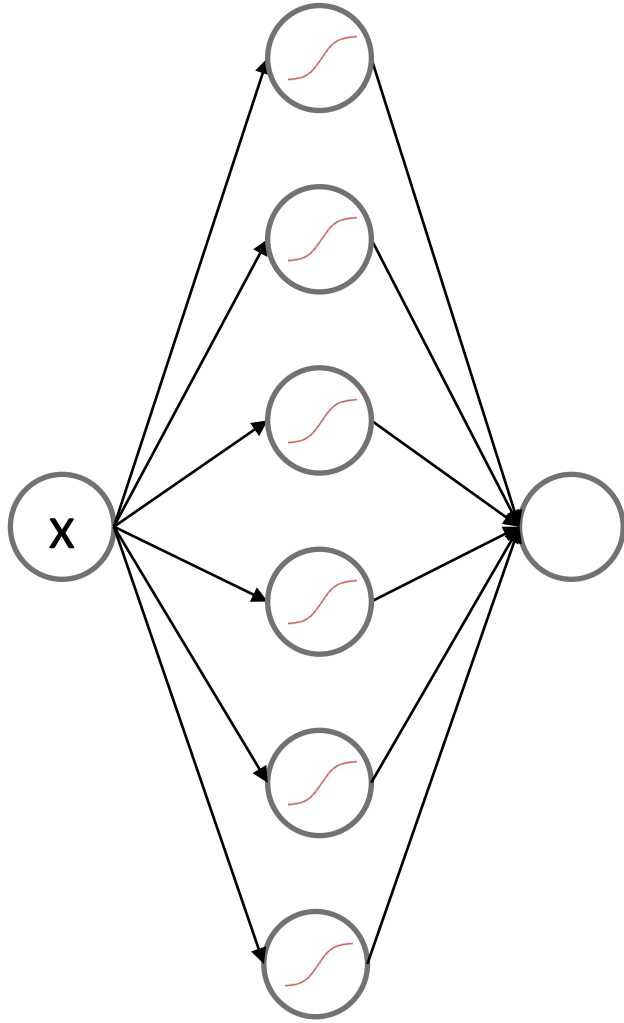
The Universal Approximation Theorem



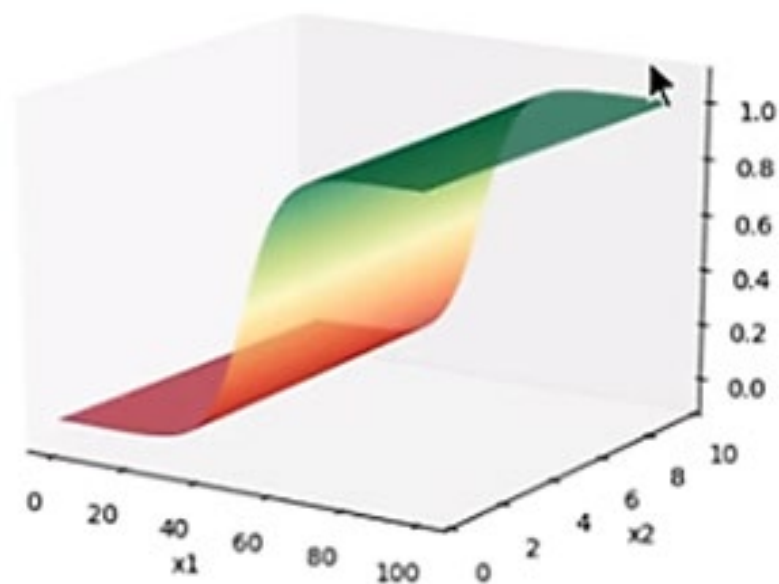
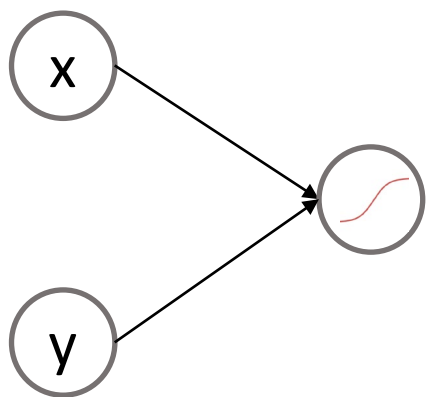
The Universal Approximation Theorem



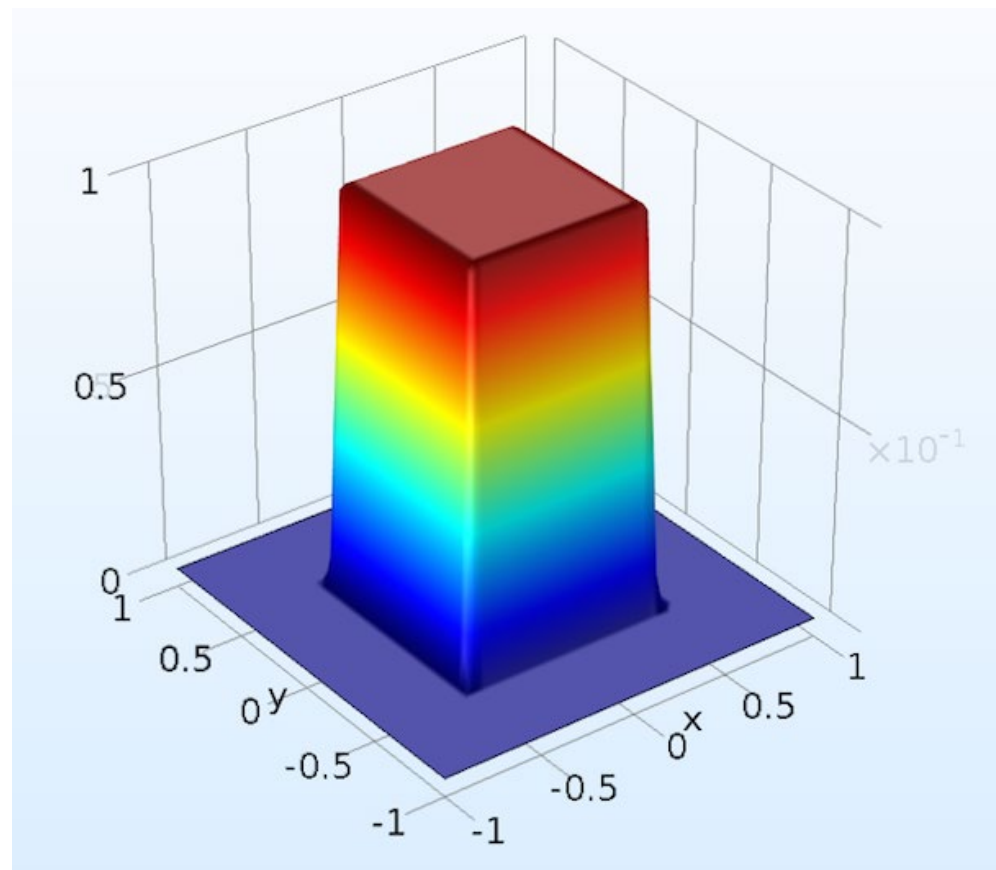
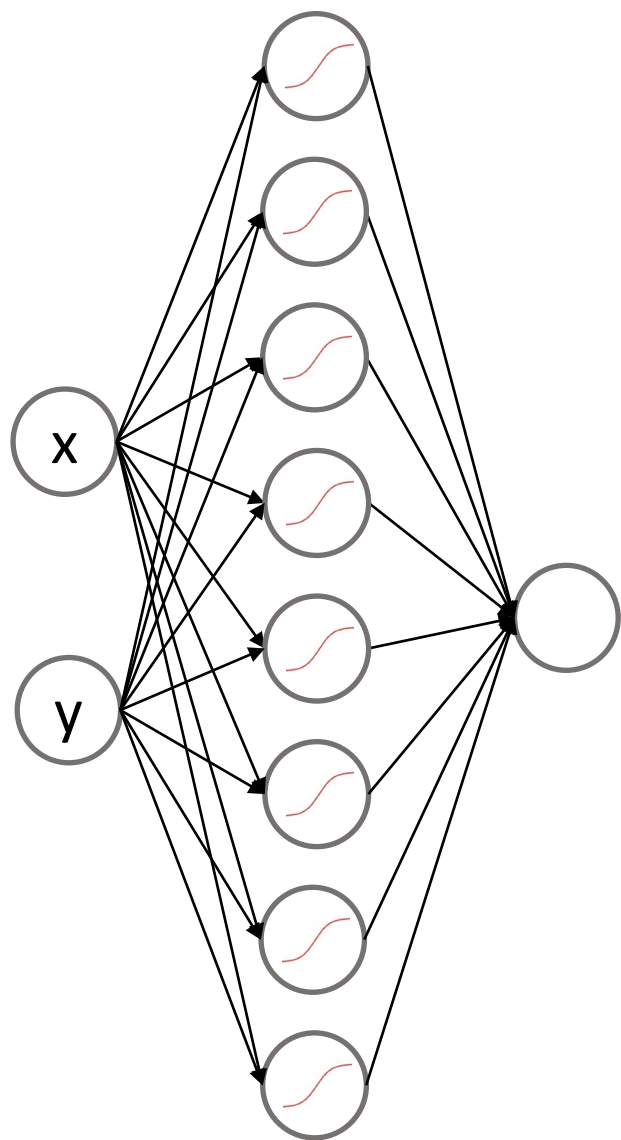
The Universal Approximation Theorem



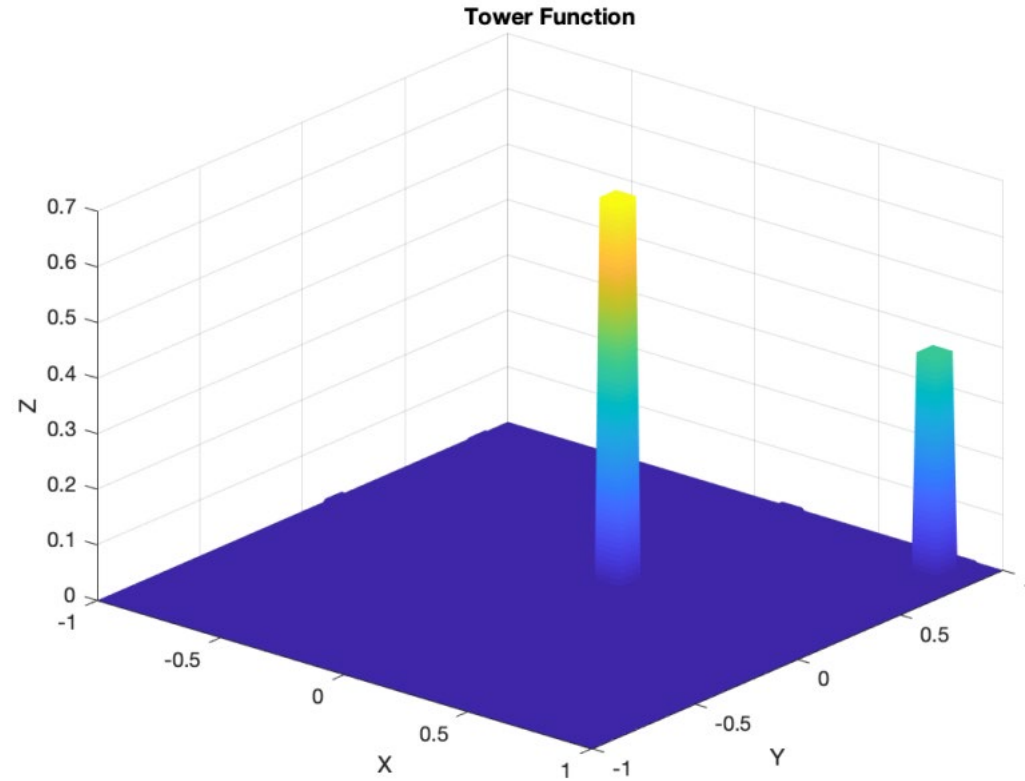
The Universal Approximator in 2D



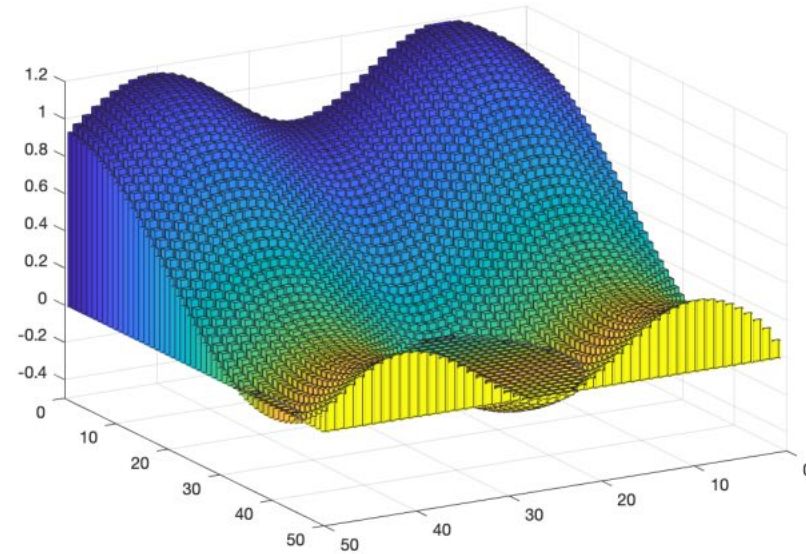
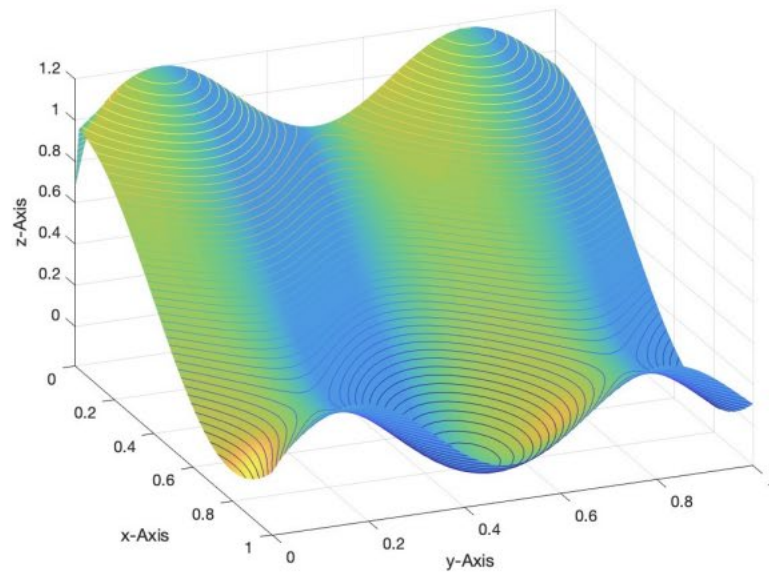
The Universal Approximator in 2D



The Universal Approximator in 2D



The Universal Approximator in 2D



The Universal Approximation Theorem

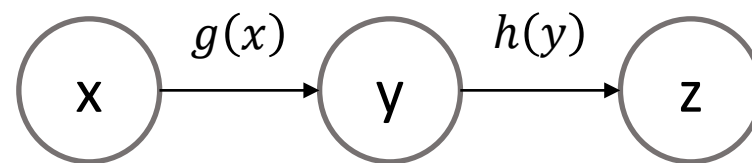
- Single layer might be enough, but it requires ‘enough’ neurons.
- Informally, ‘shallower and wider’ networks require exponentially more hidden units to compute ‘narrower and deeper’ neural networks
 - [Lecture 2 | The Universal Approximation Theorem - YouTube](#)

The Chain Rule

The Chain Rule

- A single variable chain rule

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



$$f: h \circ g$$

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

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

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

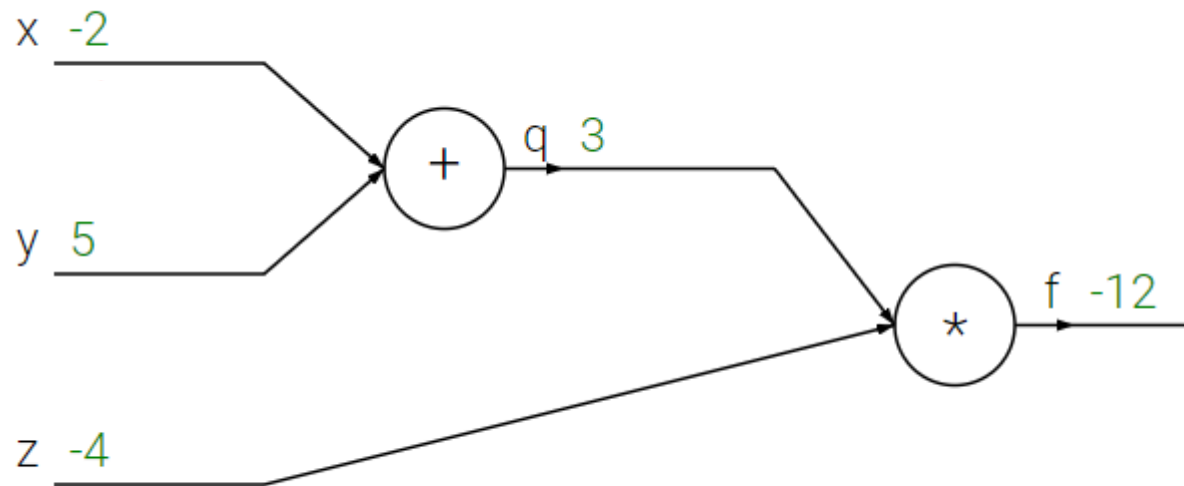
Simple Example

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

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

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

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



Simple Example

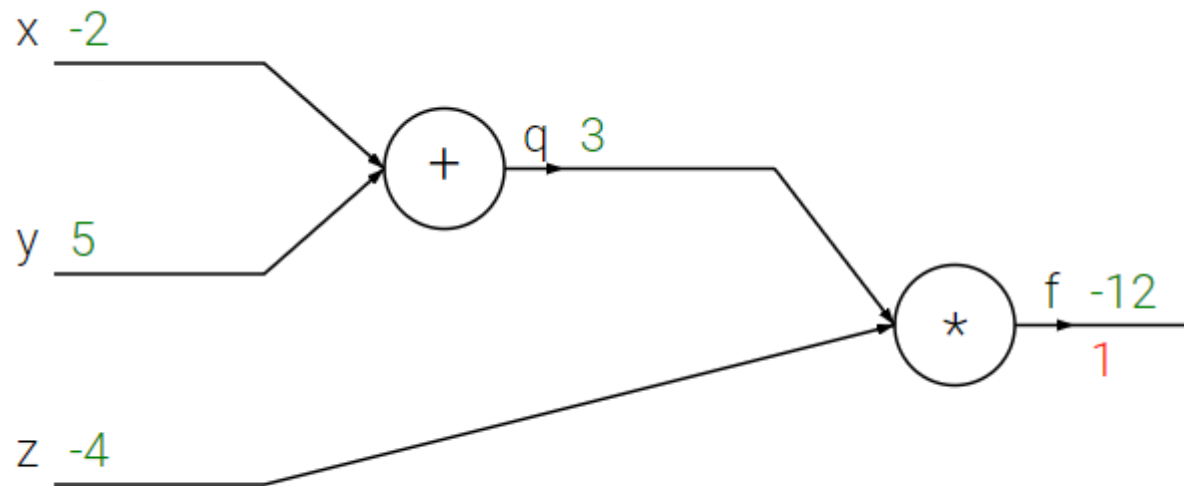
$$f(x, y, z) = (x + y)z$$

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

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

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

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



Simple Example

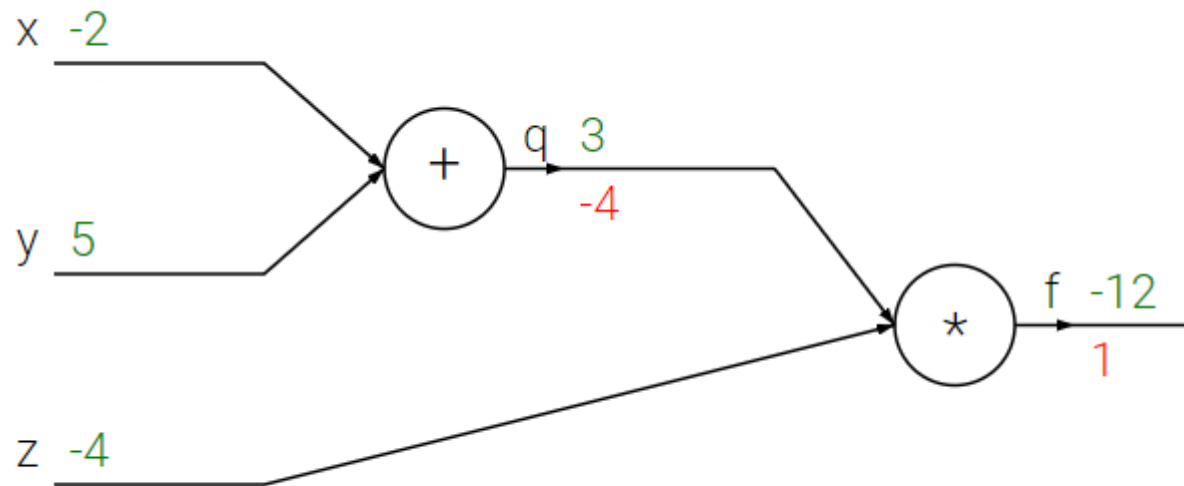
$$f(x, y, z) = (x + y)z$$

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

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

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

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



Simple Example

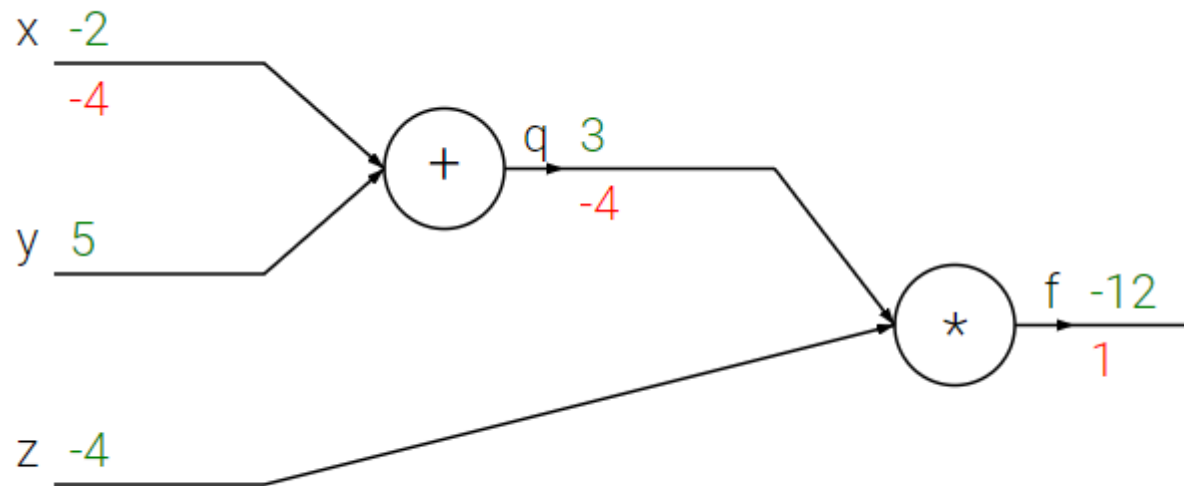
$$f(x, y, z) = (x + y)z$$

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

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

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

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



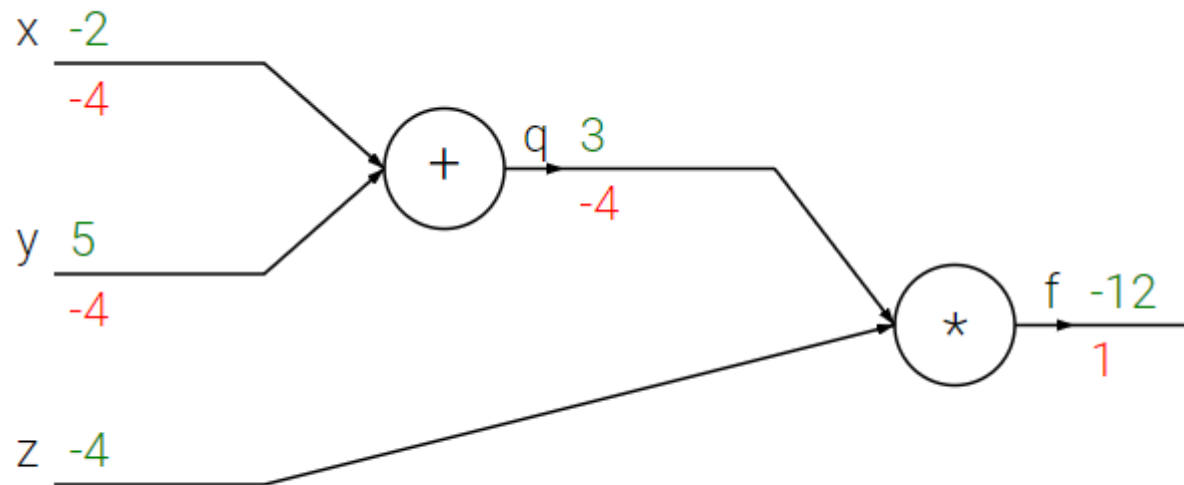
Simple Example

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

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

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

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



Simple Example

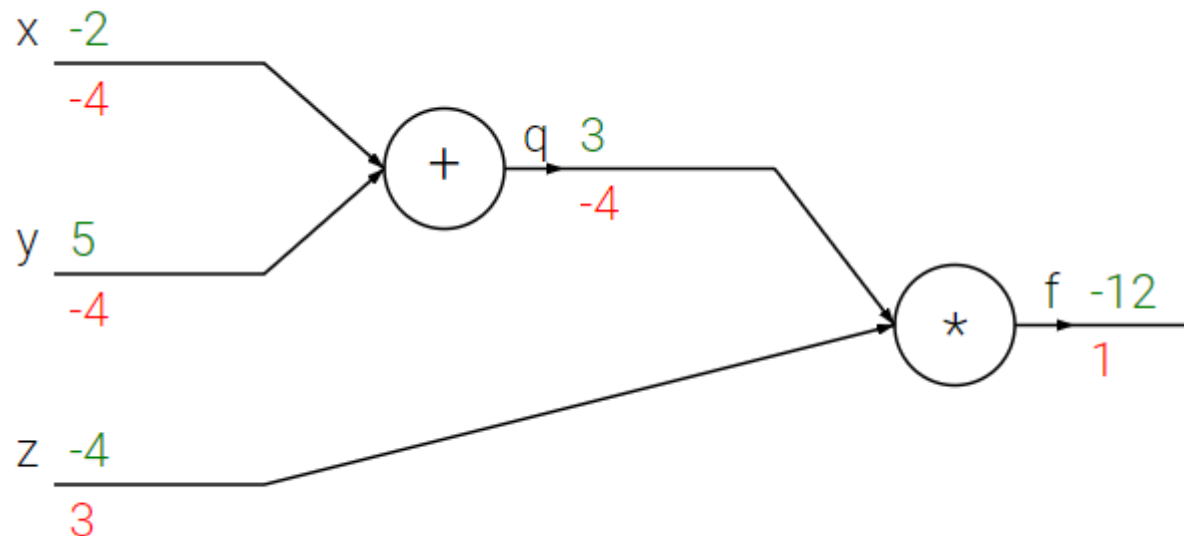
$$f(x, y, z) = (x + y)z$$

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

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

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

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



Sigmoid Example

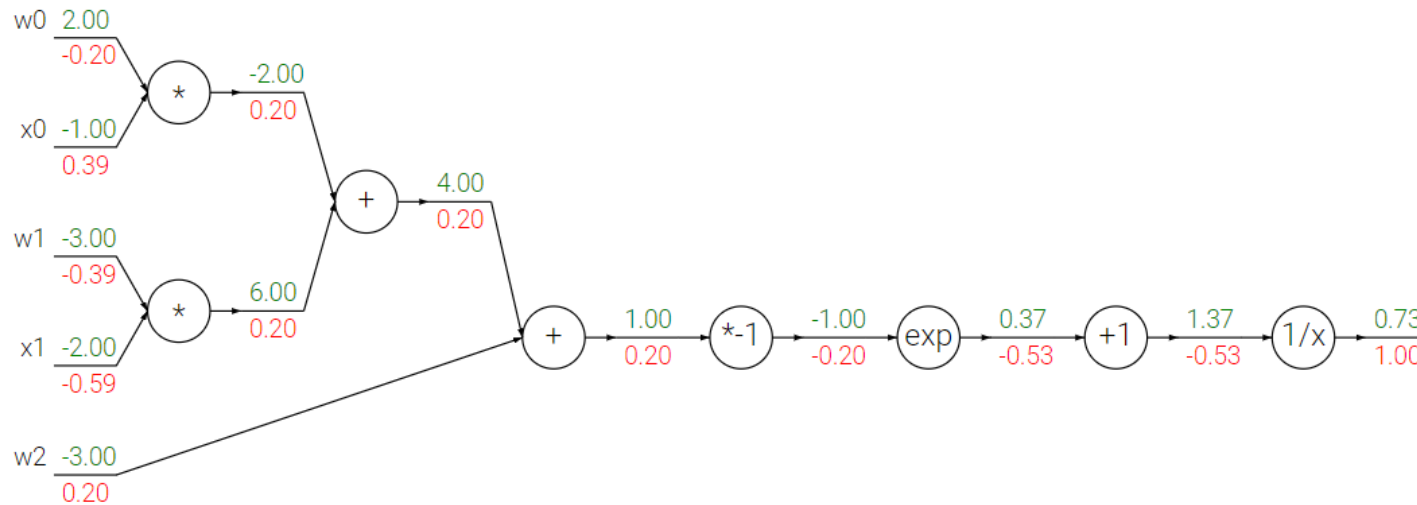
$$\sigma(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

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

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

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

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

$$y = f(x)$$

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

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

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

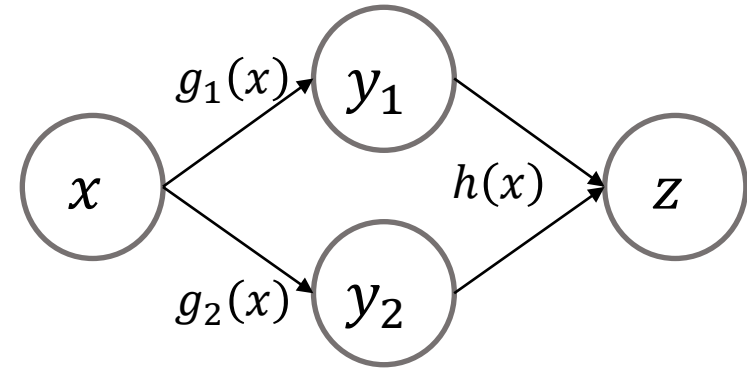
- Multi-variable chain rule

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

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

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

$$\frac{dz}{dx} = \frac{dz}{dy_1} \frac{dy_1}{dx} + \frac{dz}{dy_2} \frac{dy_2}{dx} \quad \text{(Total derivative)}$$



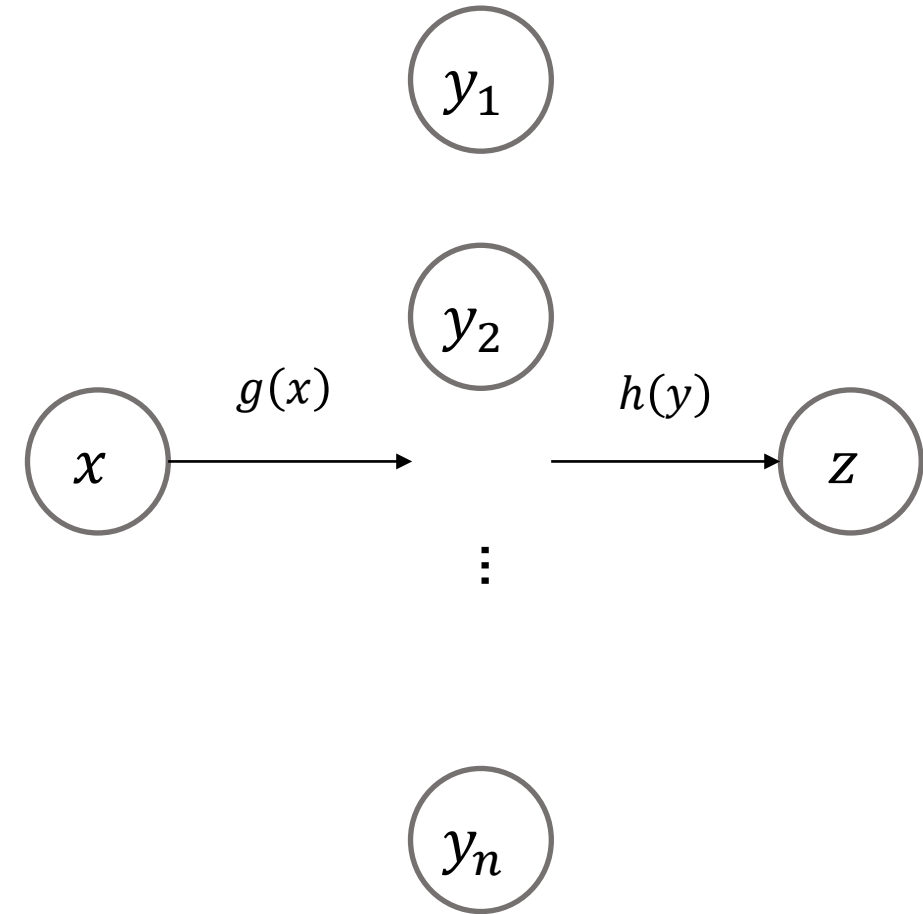
The Chain Rule

- Multi-variable chain rule

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

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

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



$$\frac{\partial z}{\partial x} = \sum_{i=1}^n \frac{\partial z}{\partial y_i} \frac{dy_i}{dx} = \underbrace{\frac{\partial z}{\partial y}}_{\in \mathbb{R}^{1 \times n}} \underbrace{\frac{\partial y}{\partial x}}_{\in \mathbb{R}^{n \times 1}}$$

The Chain Rule

- Multi-variable chain rule

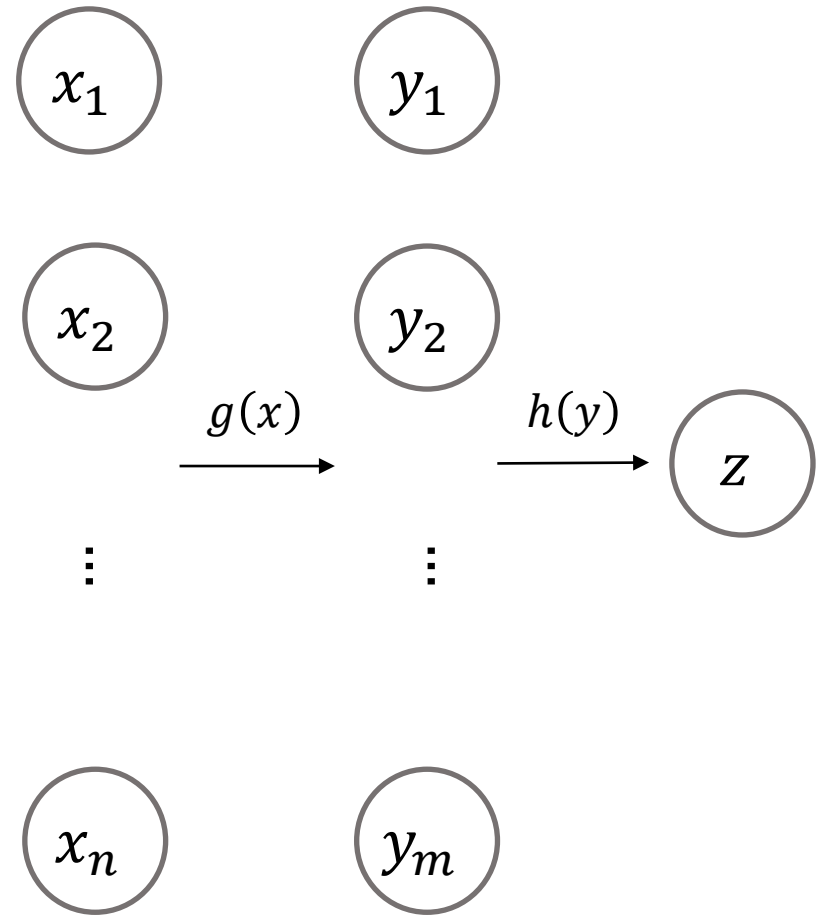
$$x \in \mathbb{R}^n, y \in \mathbb{R}^m, z \in \mathbb{R}$$

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

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

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

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



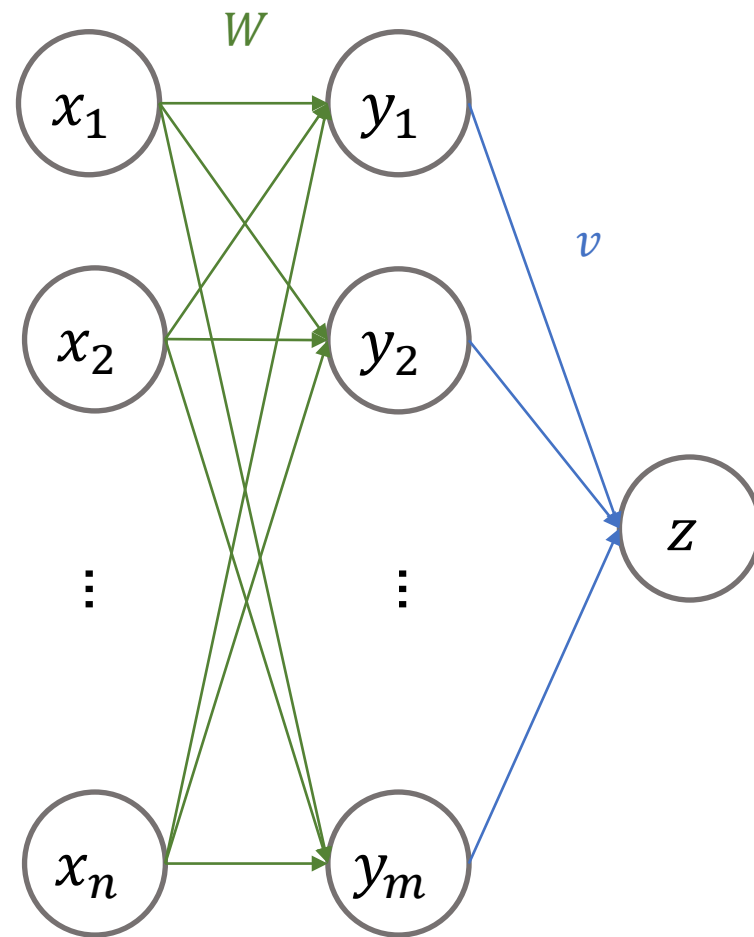
Two Layers MLP

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

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

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

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



Derivatives of Linear Layer

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

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

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

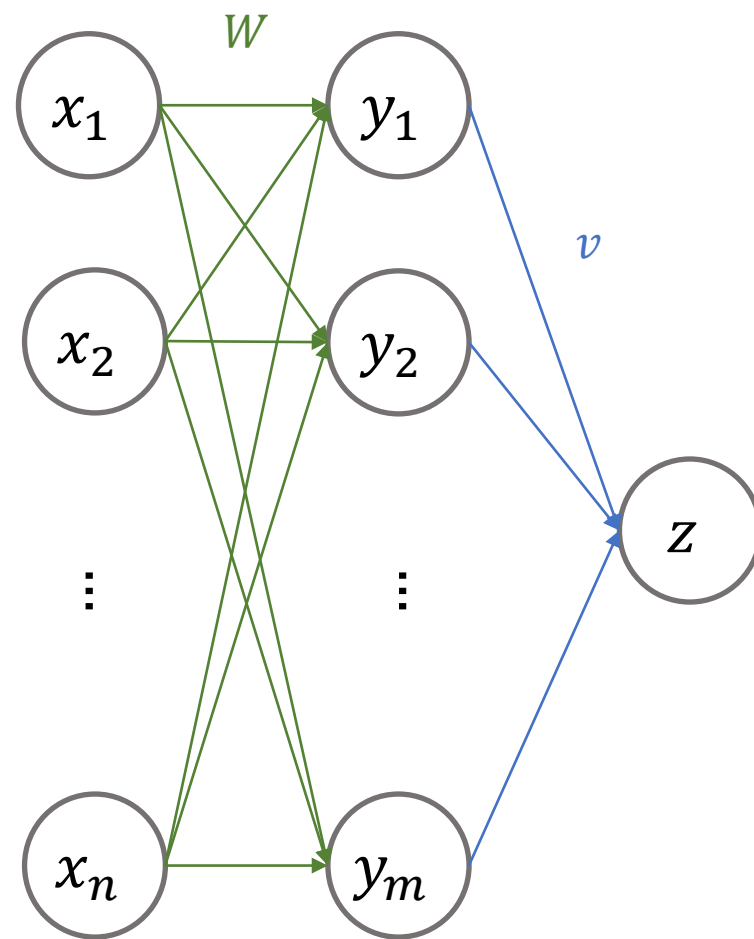
Two Layers MLP

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

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

$$\frac{\partial z}{\partial W} = \underbrace{\frac{\partial z}{\partial y}}_{\in \mathbb{R}^{1 \times m}} \underbrace{\frac{\partial y}{\partial W}}_{\in \mathbb{R}^{m \times m \times n}}$$

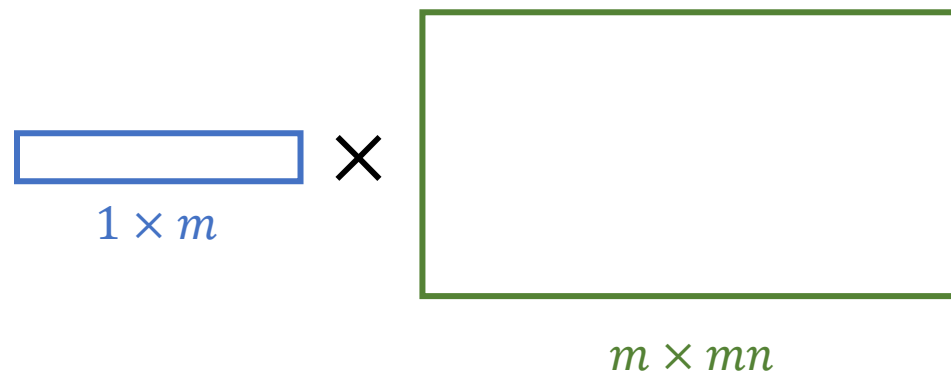
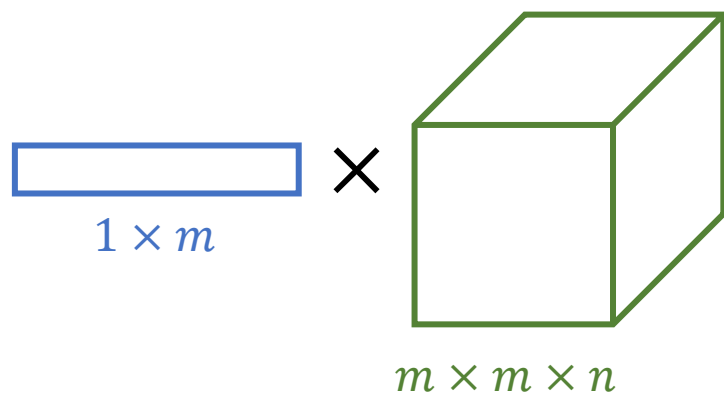
Tensor Product
(n-mode product)



Tensor Product

- N-mode product
 - Matricization -> matrix multiplication

$$\frac{\overset{\in \mathbb{R}^{m \times n}}{\partial z}}{\overset{\in \mathbb{R}^{1 \times m}}{\partial W}} = \frac{\overset{\in \mathbb{R}^{m \times m \times n}}{\partial y}}{\partial W}$$



Vector Jacobian Product (VJP)

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

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

$$y = Wx$$

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

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

Vector Jacobian Product (VJP)

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

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

$$y = Wx$$

[illegible]

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

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

Vector Jacobian Product (VJP)

- Explicit formation of Jacobian is too expensive

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

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

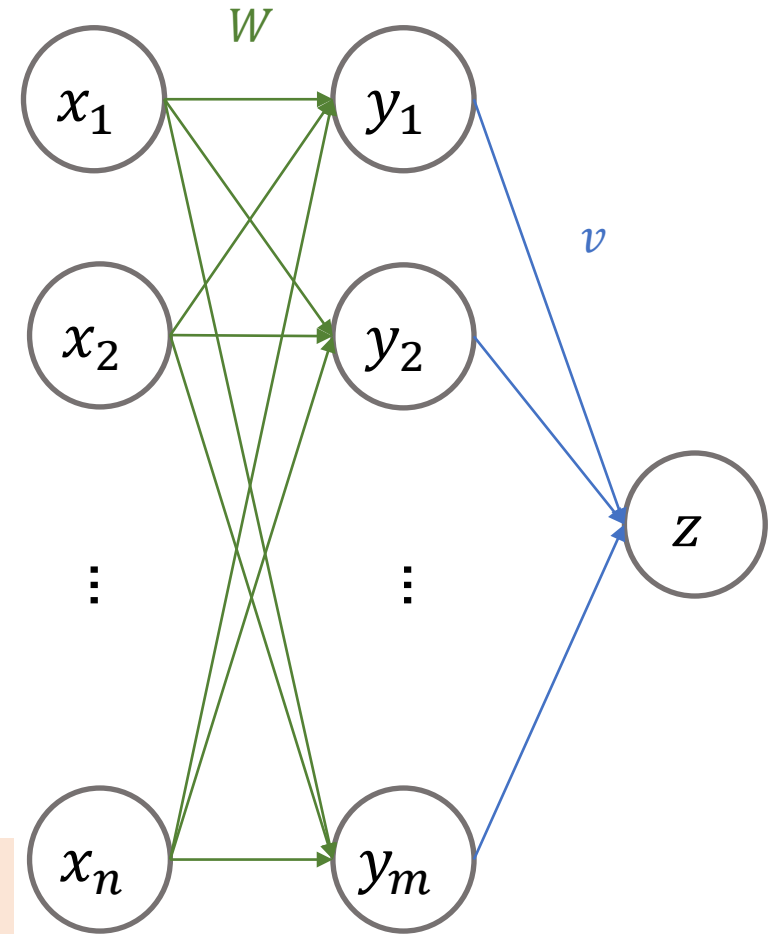
$$\frac{\partial z}{\partial W} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial W}$$

$\frac{\partial z}{\partial W} \in \mathbb{R}^{1 \times m}$
 $\frac{\partial y}{\partial W} \in \mathbb{R}^{m \times m \times n}$

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

$\left(\frac{\partial z}{\partial y} \right)^\top \in \mathbb{R}^{m \times 1}$
 $x^\top \in \mathbb{R}^{1 \times n}$

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



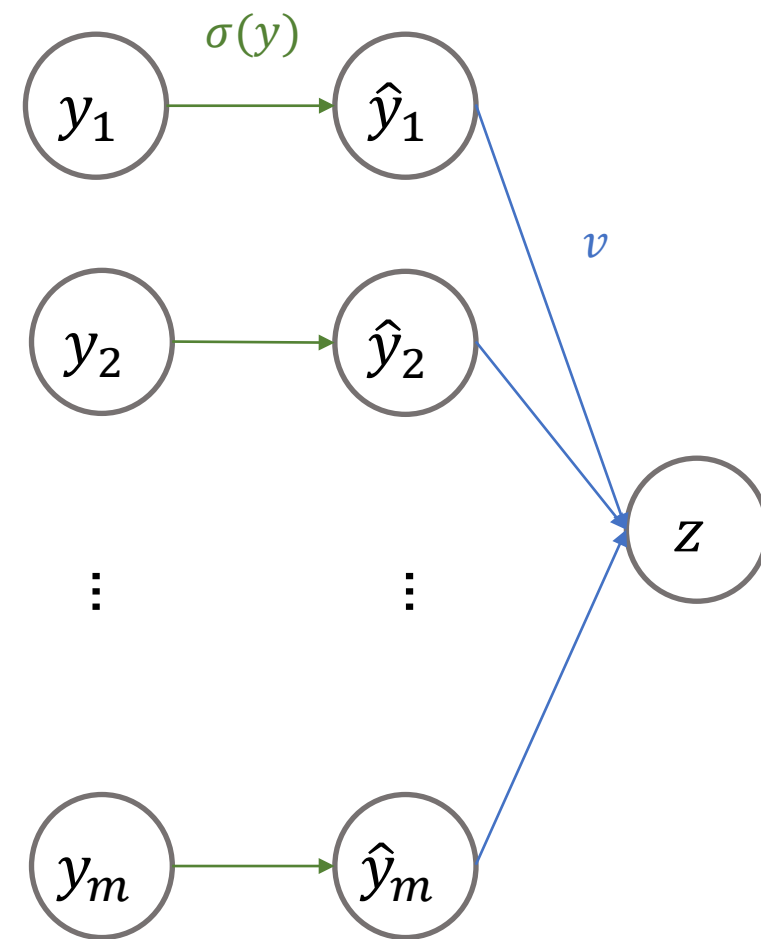
Vector Jacobian Product (VJP)

- Elementwise activation functions

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

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

$$\begin{aligned} \frac{\partial z}{\partial y} &= \frac{\partial z}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial y} \\ &\in \mathbb{R}^{1 \times m} \quad \in \mathbb{R}^{m \times m} \quad \in \mathbb{R}^{1 \times m} \\ \frac{\partial \hat{y}}{\partial y} &= \begin{bmatrix} \frac{\partial \hat{y}_1}{\partial y_1} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \frac{\partial \hat{y}_m}{\partial y_m} \end{bmatrix} \\ &= \begin{bmatrix} \sigma(y_1)(1 - \sigma(y_1)) & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \sigma(y_m)(1 - \sigma(y_m)) \end{bmatrix} \end{aligned}$$



Vector Jacobian Product (VJP)

- Elementwise activation functions

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

Element-wise product

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

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)

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

$$b = f(a)$$

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

$$c = g(b)$$

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

$$d = h(c)$$

$$e = i(d)$$

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

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

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)

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

$$b = f(a)$$

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

$$c = g(b)$$

$$d = h(c)$$

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

$$e = i(d)$$

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

Loss function:
scalar function

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

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

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

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

$$b = f(a)$$

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

$$c = g(b)$$

$$d = h(c)$$

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

$$e = i(d)$$

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

Loss function:
scalar function

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

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

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

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

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

Vector-Jacobian
Product (VJP)

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

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

$$\frac{\partial e}{\partial a} = \overset{\in \mathbb{R}^{1 \times n_2}}{\frac{\partial e}{\partial e}} \overset{\in \mathbb{R}^{n_2 \times n_1}}{\frac{\partial e}{\partial d} \frac{\partial d}{\partial c} \frac{\partial c}{\partial b} \frac{\partial b}{\partial a}} = \frac{\partial e}{\partial b} \frac{\partial b}{\partial a}$$

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?

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

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

3 X
Computation

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

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

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

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

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

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

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

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

$$b = f(a)$$

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

$$c = g(b)$$

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

$$d = h(c)$$

$$e = L(d)$$

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

$$\frac{\partial b}{\partial a} = \overset{\in \mathbb{R}^{n_1 \times 1}}{\frac{\partial b}{\partial a}} \overset{\in \mathbb{R}^{1 \times 1}}{\frac{\partial a}{\partial a}} = \frac{\partial b}{\partial a} \mathbf{1}$$

Forward-Mode AD

single variable

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

$$b = f(a)$$

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

$$c = g(b)$$

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

$$d = h(c)$$

$$e = L(d)$$

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

$$\frac{\partial b}{\partial a} = \frac{\partial b}{\partial a} \frac{\partial a}{\partial a} = \frac{\partial b}{\partial a} \overset{\in \mathbb{R}^{n_1 \times 1}}{\underset{\in \mathbb{R}^{1 \times 1}}{1}}$$

$$\frac{\partial c}{\partial a} = \frac{\partial c}{\partial b} \frac{\partial b}{\partial a} \frac{\partial a}{\partial a} = \frac{\partial c}{\partial b} \frac{\partial b}{\partial a} \overset{\in \mathbb{R}^{n_2 \times n_1}}{\underset{\in \mathbb{R}^{n_1 \times 1}}{1}}$$

Forward-Mode AD

single variable

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

$$b = f(a)$$

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

$$c = g(b)$$

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

$$d = h(c)$$

$$e = L(d)$$

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

$$\frac{\partial b}{\partial a} = \frac{\overset{\in \mathbb{R}^{n_1 \times 1}}{\partial b} \overset{\in \mathbb{R}^{1 \times 1}}{\partial a}}{\partial a} = \frac{\partial b}{\partial a} 1$$

$$\frac{\partial c}{\partial a} = \frac{\overset{\in \mathbb{R}^{n_2 \times n_1}}{\partial c} \overset{\in \mathbb{R}^{n_1 \times 1}}{\partial b} \partial a}{\partial b \partial a} = \frac{\partial c}{\partial b} \frac{\partial b}{\partial a}$$

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

Forward-Mode AD

single variable

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

$$b = f(a)$$

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

$$c = g(b)$$

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

$$d = h(c)$$

$$e = L(d)$$

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

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

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

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?

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

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

3 X
Computation

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

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

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

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

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

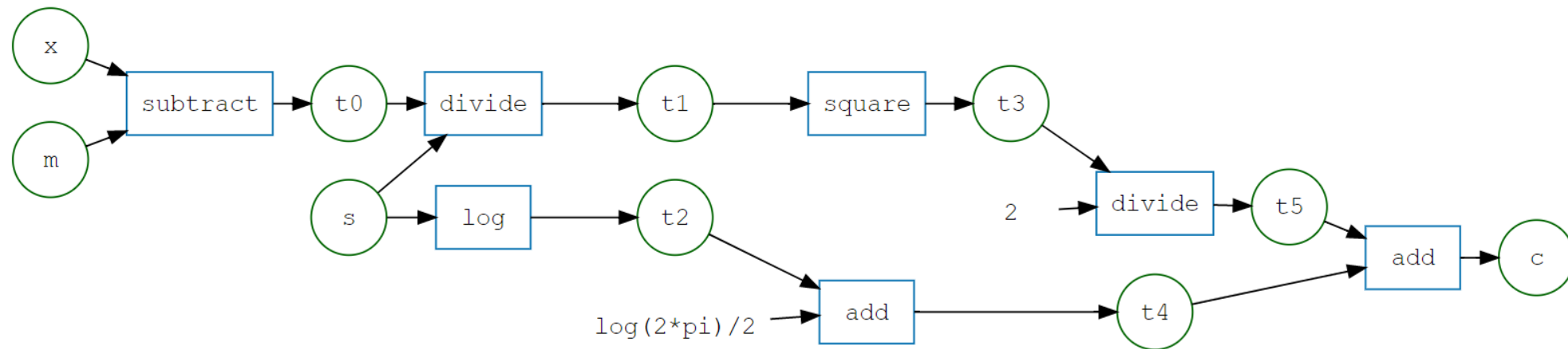
$\in \mathbb{R}^{n_4 \times n_3}$ $\in \mathbb{R}^{n_3 \times 3}$

Automatic Differentiation (AD)

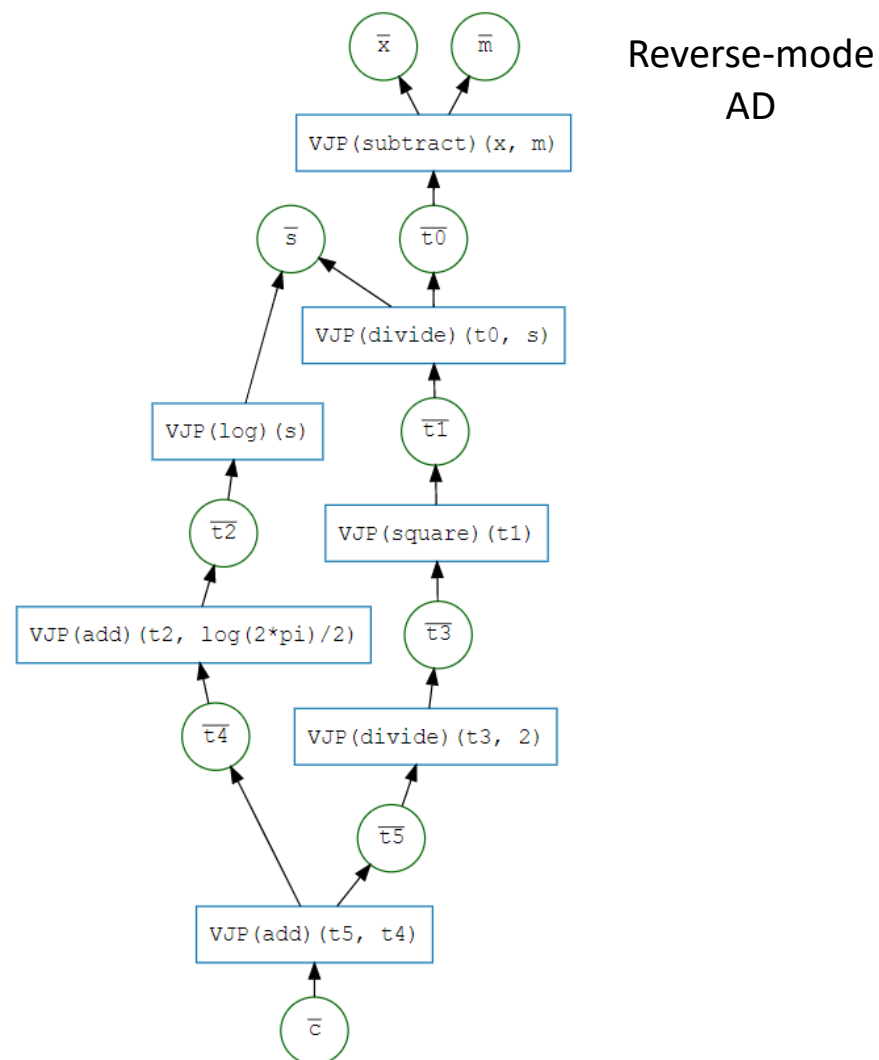
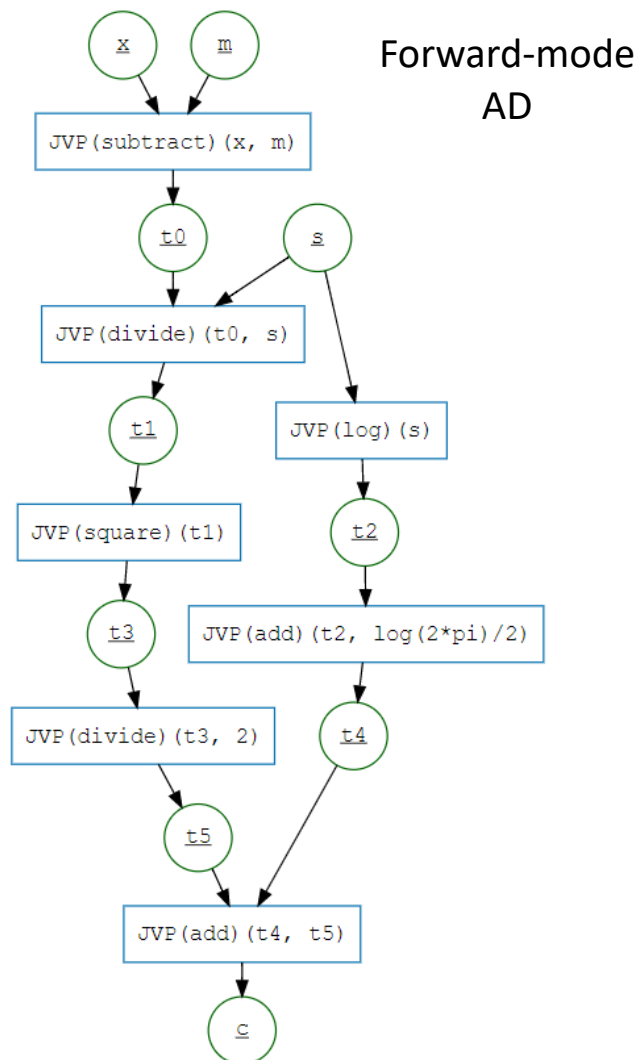
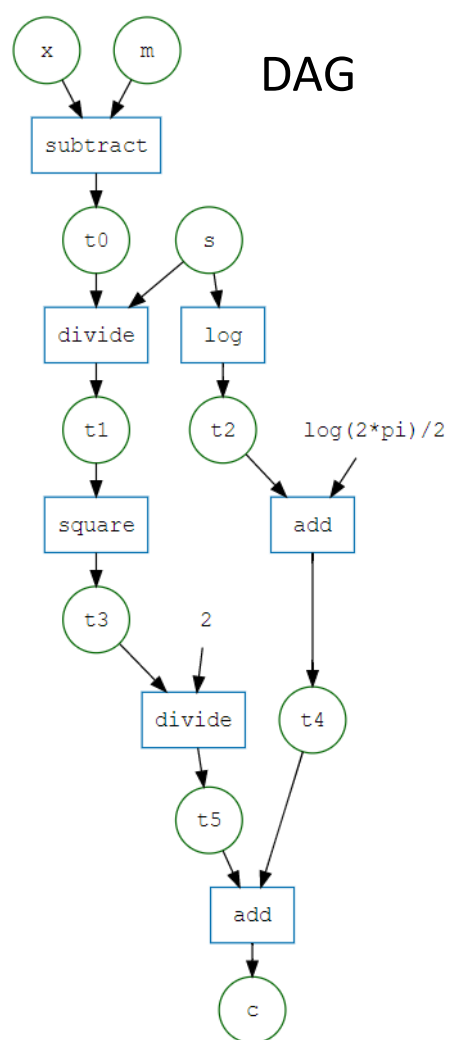
- For low dimensional outputs and high dimensional inputs
 - Objective function w/ deep neural networks
 - reverse-mode AD
- For high dimensional outputs and low dimensional inputs
 - Forward-mode AD

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 \(toronto.edu\)](#)