

**CS2109S: Introduction to AI and Machine Learning**

# Lecture 9: Backpropagation

27 October 2023

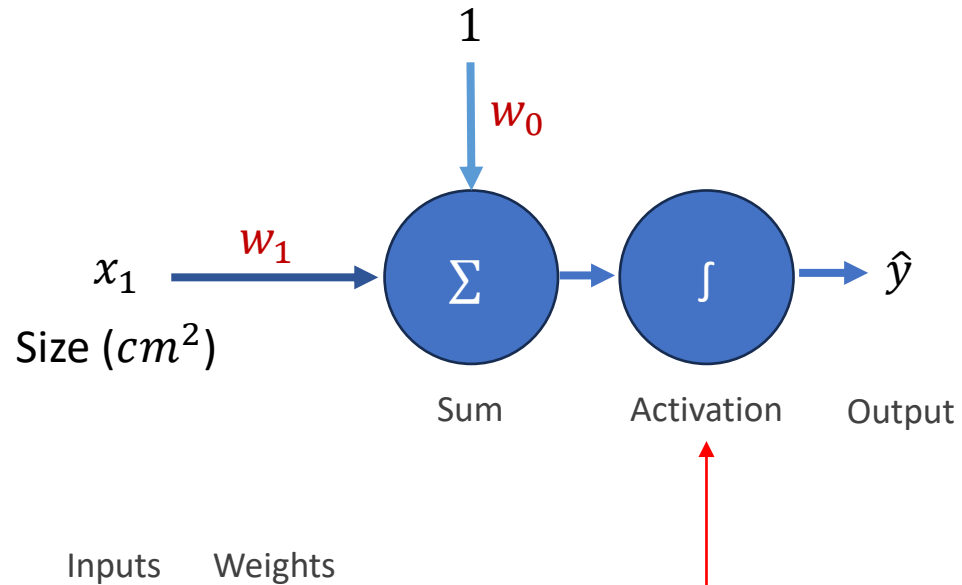
# Announcement

- Midterm grading is still on progress, sorry :(

# Recap

- Perceptron
  - Biological inspiration: brain, neural network, neuron
  - Perceptron Learning Algorithm:
    - $w \leftarrow w + \gamma(y^{(j)} - \hat{y}^{(j)})x^{(j)}$  on a misclassified instance
- Neural Networks
  - Single-layer Neural Networks: AND, OR, NOR
  - Multi-layer Neural Networks: XNOR, Universal Approximation Theorem
  - Regression and Classification
- Neural Networks with Gradient Descent
- Neural Networks vs Other Models: learned feature mapping!

# Perceptron: An Example



Sign function

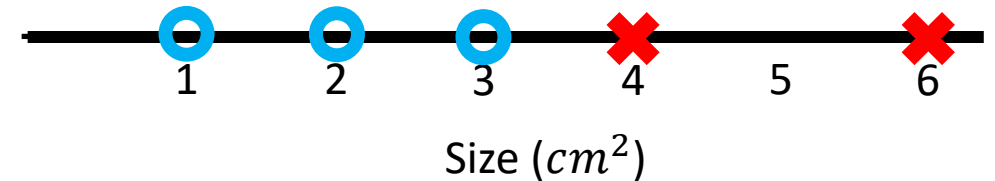
$$g(z) = \begin{cases} +1 & \text{if } z \geq 0 \\ -1 & \text{if } z < 0 \end{cases}$$

$$\hat{y} = h_w(x) = g\left(\sum_{i=0}^n w_i x_i\right)$$

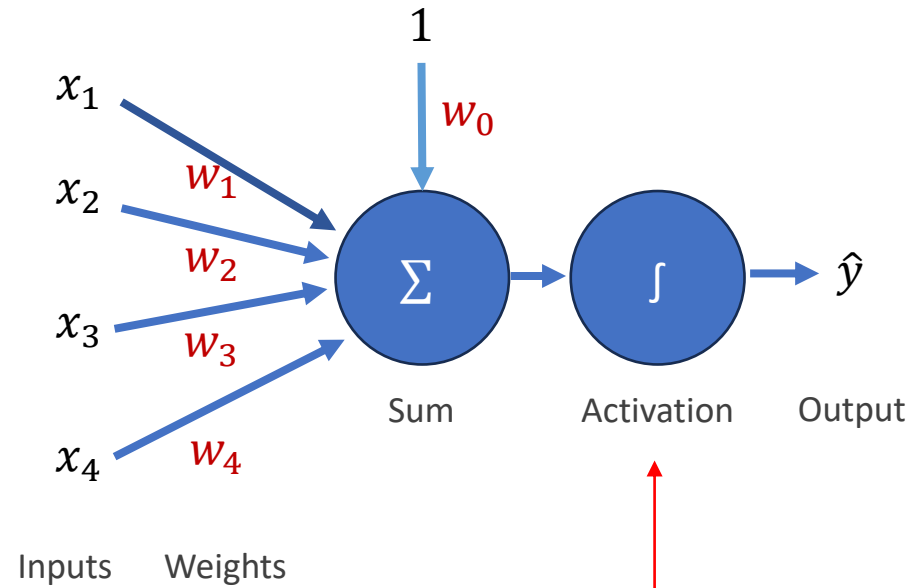
$$h_w(x) = \begin{cases} +1 & \text{if } w_0 + w_1 x \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

$$h_w(x) = \begin{cases} +1 & \text{if } -3.5 + 1x \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

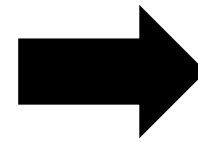
Cancer prediction: benign (-1), malignant (1)



# Single-layer Neural Networks



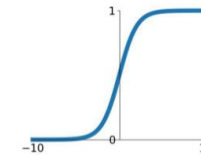
**Activation Function**  
(usually non-linear)



$$\hat{y} = h_w(x) = g\left(\sum_{i=0}^n w_i x_i\right)$$

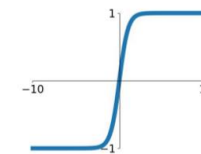
**Sigmoid**

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



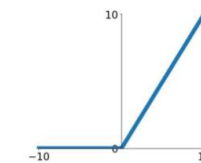
**tanh**

$$\tanh(x)$$



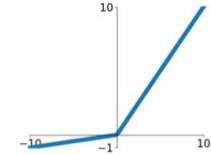
**ReLU**

$$\max(0, x)$$



**Leaky ReLU**

$$\max(0.1x, x)$$

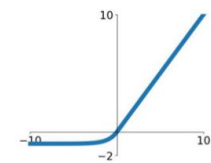


**Maxout**

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

**ELU**

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



# Background: Chain Rule

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

$$\frac{dz}{dx} = \frac{dz}{dh} \frac{dh}{dg} \frac{dg}{df} \frac{df}{dx}$$

# Aside: Derivative of Sigmoid Function

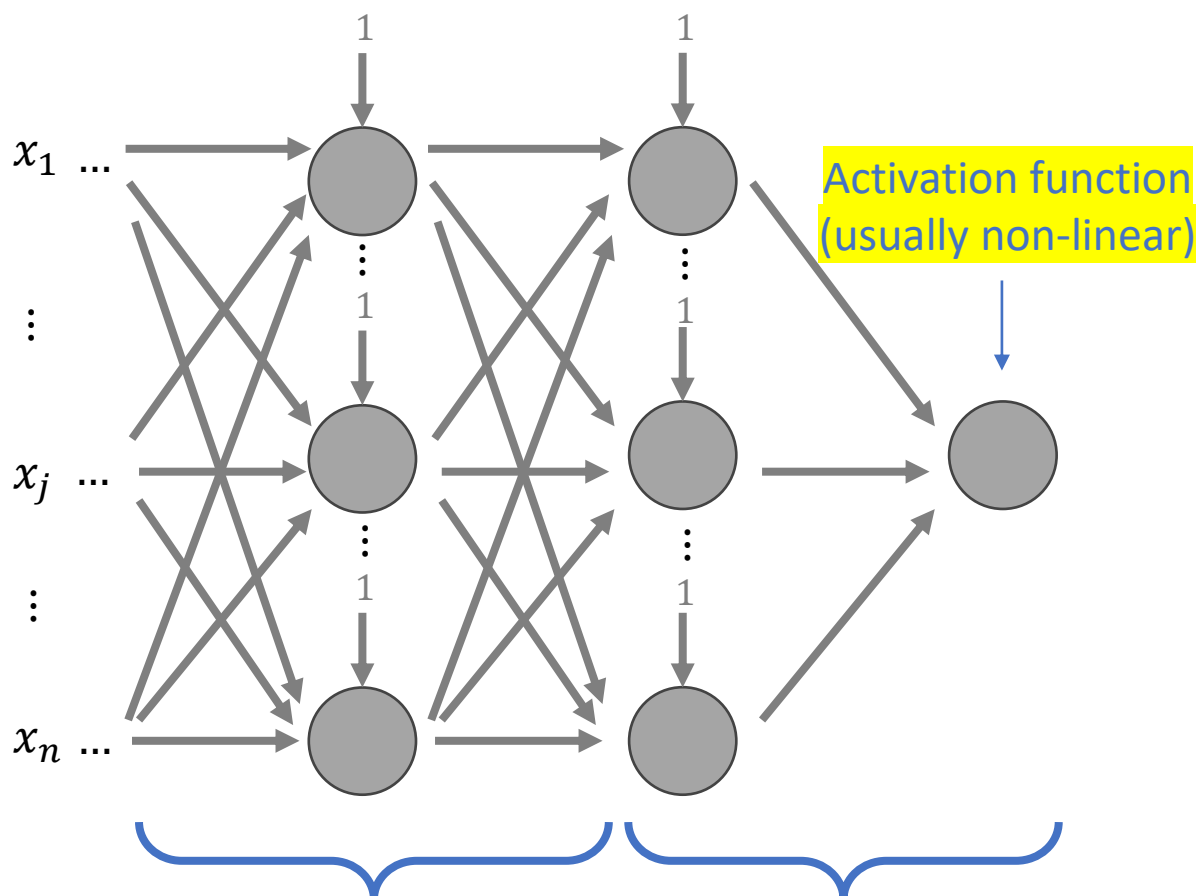
$$\begin{aligned}\sigma(x) &= \frac{1}{1 + e^{-x}} & \sigma'(x) &= \frac{d\left(\frac{1}{1 + e^{-x}}\right)}{dx} \\ & & &= \frac{d((1 + e^{-x})^{-1})}{dx} \\ & & &= -(1 + e^{-x})^{-2} (-e^{-x}) \\ & & &= \frac{e^{-x}}{(1 + e^{-x})^2} \\ & & &= \frac{1}{1 + e^{-x}} \cdot \frac{e^{-x}}{1 + e^{-x}} \\ & & &= \frac{1}{1 + e^{-x}} \left(1 - \frac{1}{1 + e^{-x}}\right) \\ & & &= \boxed{\sigma(x)(1 - \sigma(x))}\end{aligned}$$

ERRATA

# Neural Networks vs Other Methods

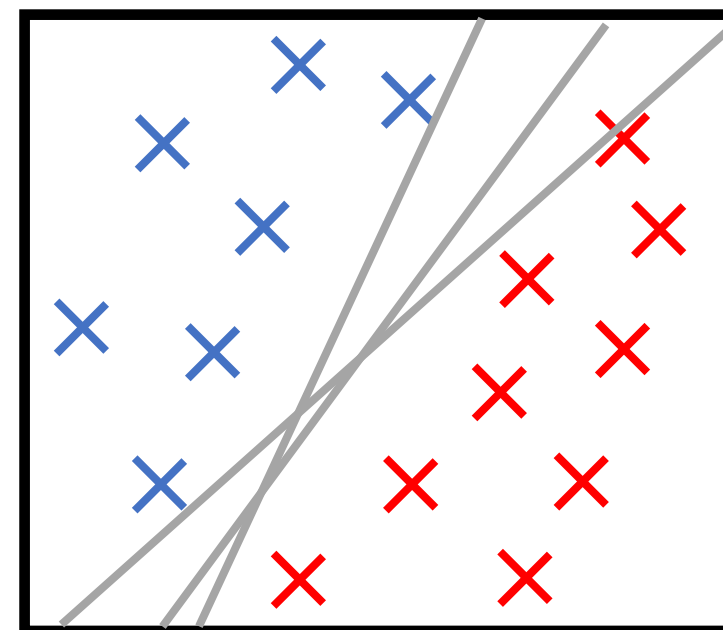


Credit: Internet meme, original source unknown



Learned Feature Mapping

Linear Model

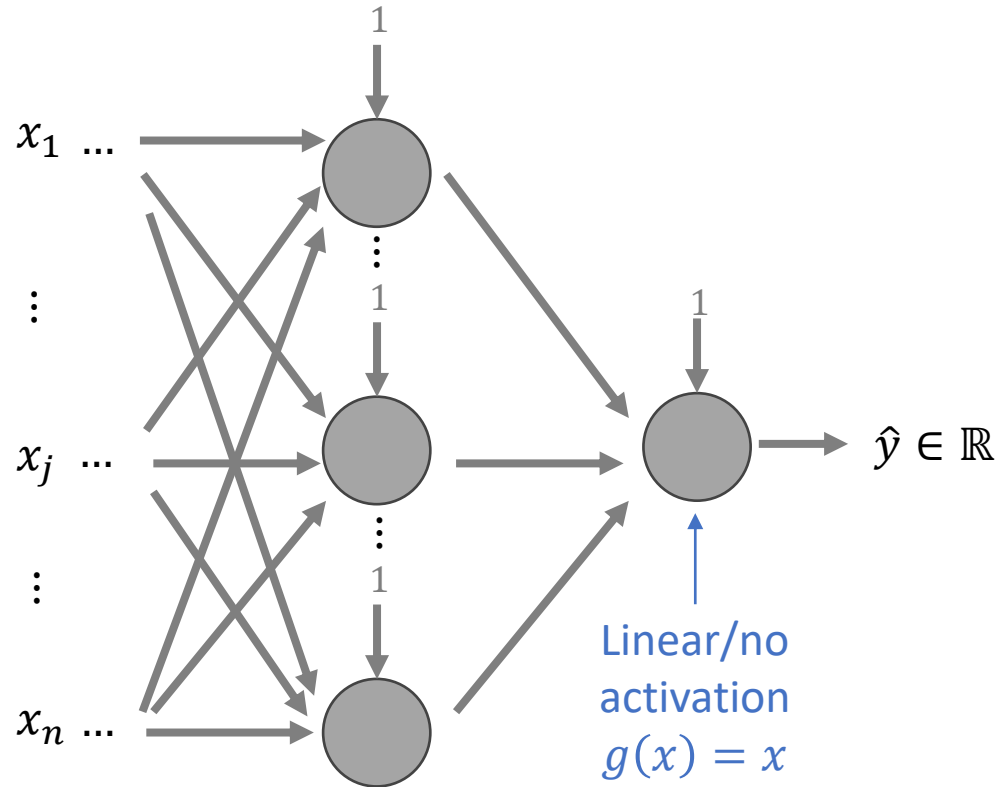


Non-linear, non-robust decision boundary

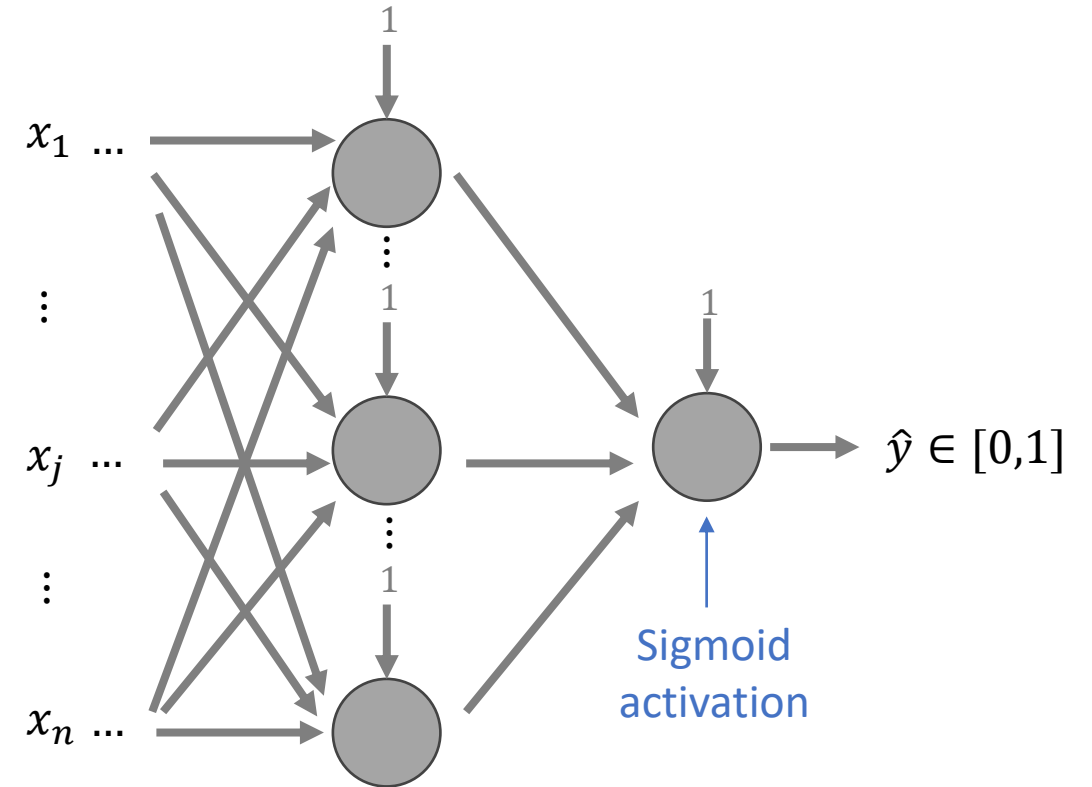
Prone to misclassification since the decision boundary can be too close to data points



# Regression and Classification

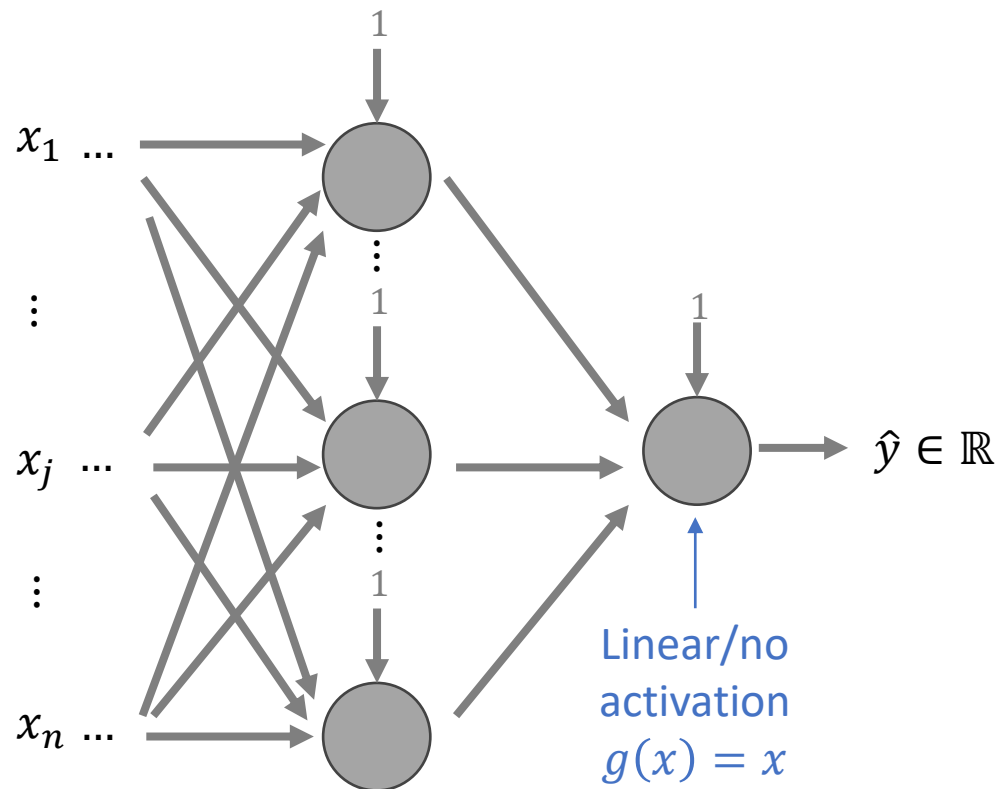


Regression

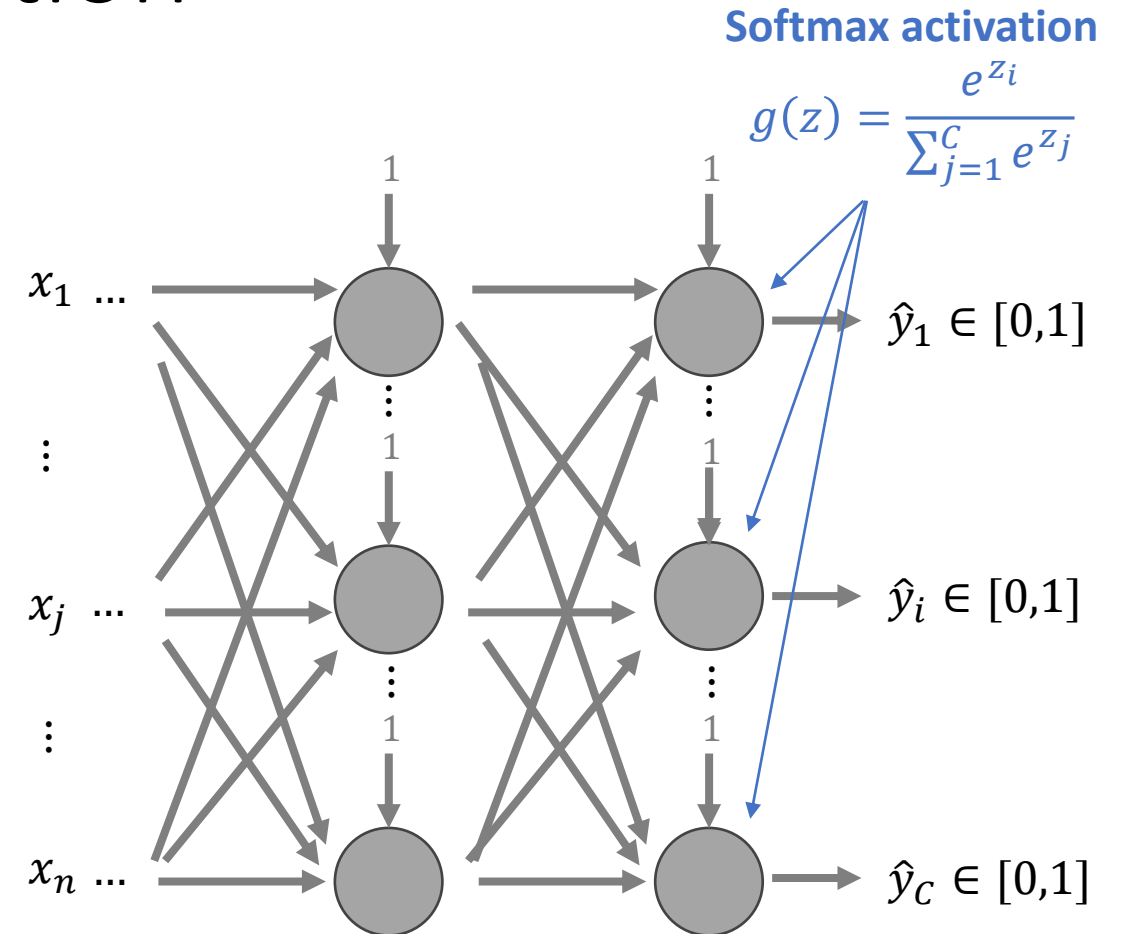


Binary Classification

# Regression and Classification



Regression

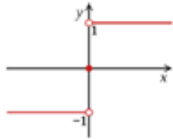
Multi-class Classification  
with  $C$  classes

$$\sum_{i=1}^C \hat{y}_i = 1$$

# Multi-layer Neural Networks: $|x - 1|$

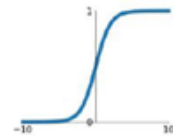
**Step**

$$\text{sgn}(x) = \begin{cases} +1 & x \geq 0 \\ -1 & x < 0 \end{cases}$$



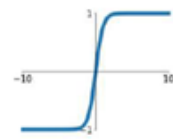
**Sigmoid**

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



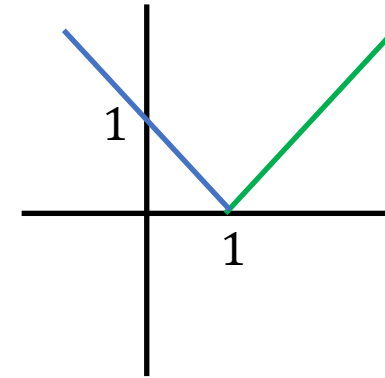
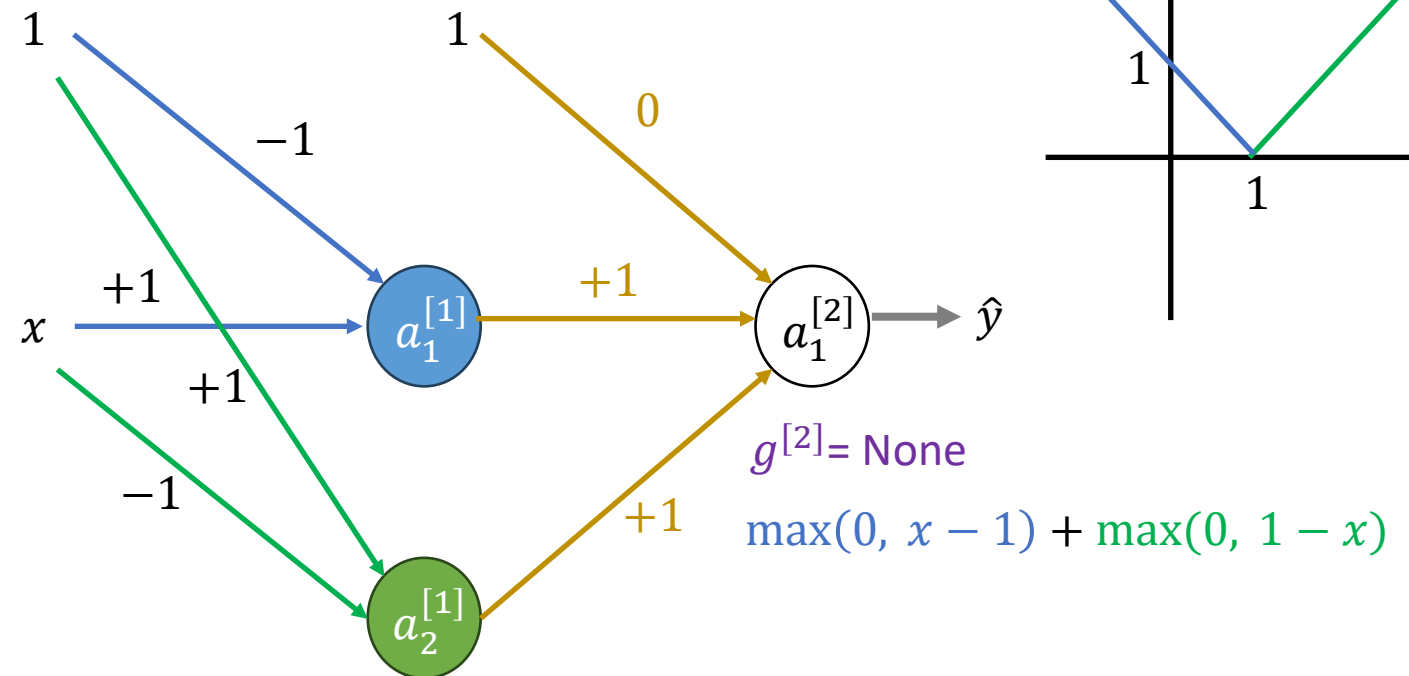
**tanh**

$$\tanh(x)$$



**ReLU**

$$\max(0, x)$$



$$g^{[2]} = \text{None}$$

$$\max(0, x - 1) + \max(0, 1 - x)$$

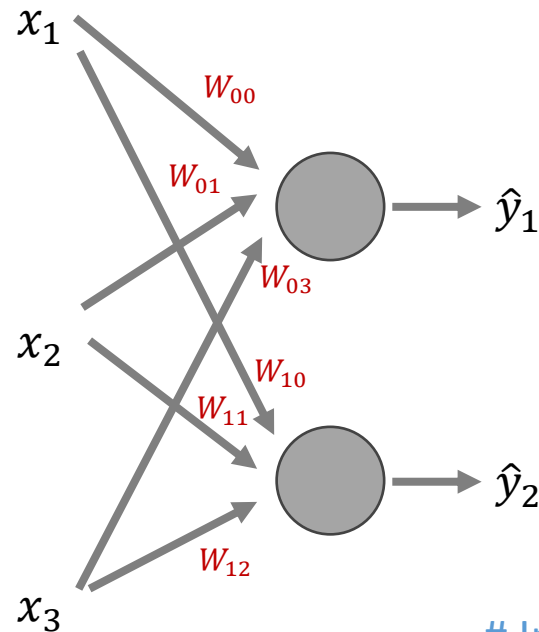
Which activation function(s)?

$$g^{[1]} = \text{ReLU } \max(0, x)$$

$$\max(0, x - 1)$$

$$\max(0, 1 - x)$$

# Neural Networks and Matrix Multiplication (1)

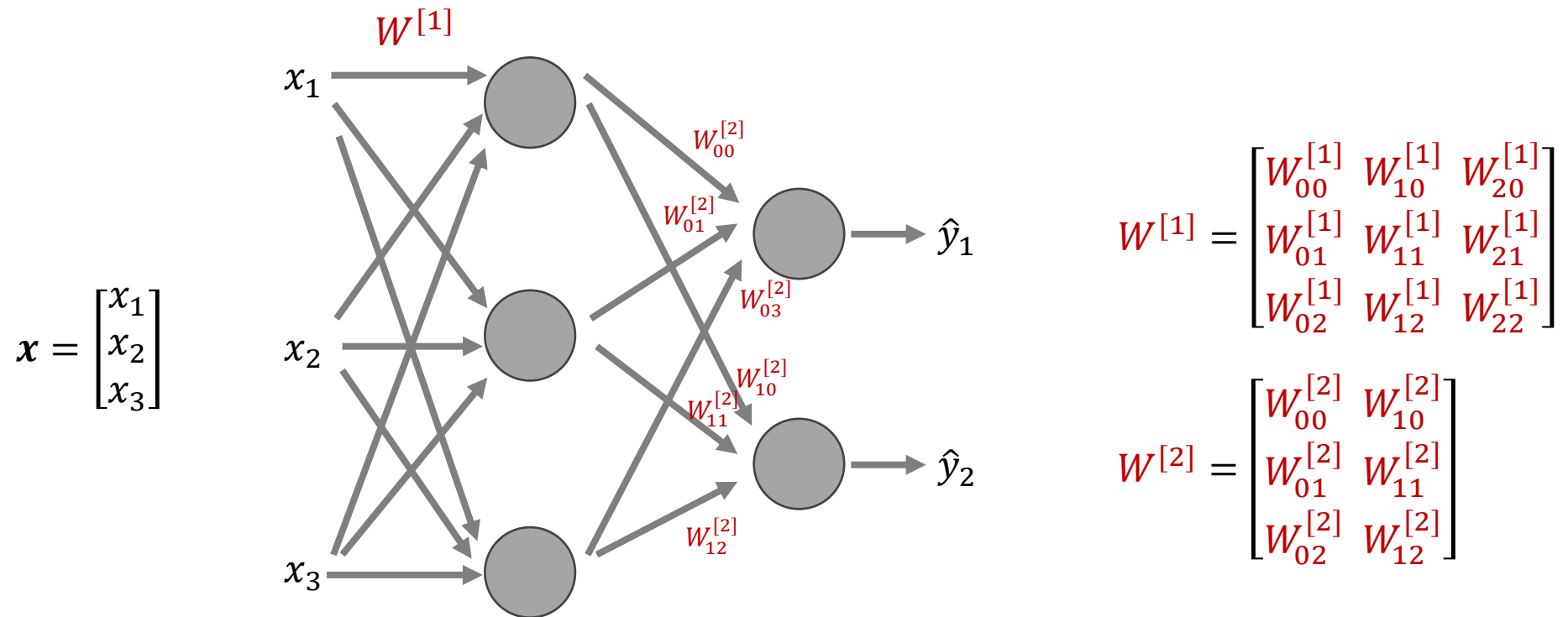


# Input (number of weights per neuron / input variables)

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \mathbf{W} = \begin{bmatrix} W_{00} & W_{10} \\ W_{01} & W_{11} \\ W_{02} & W_{12} \end{bmatrix} \quad \hat{\mathbf{y}} = g(\mathbf{W}^T \mathbf{x}) = g \left( \begin{bmatrix} W_{00} & W_{10} \\ W_{01} & W_{11} \\ W_{02} & W_{12} \end{bmatrix}^T \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \right) = g \left( \begin{bmatrix} W_{00} & W_{01} & W_{02} \\ W_{10} & W_{11} & W_{12} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \right) = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \end{bmatrix}$$

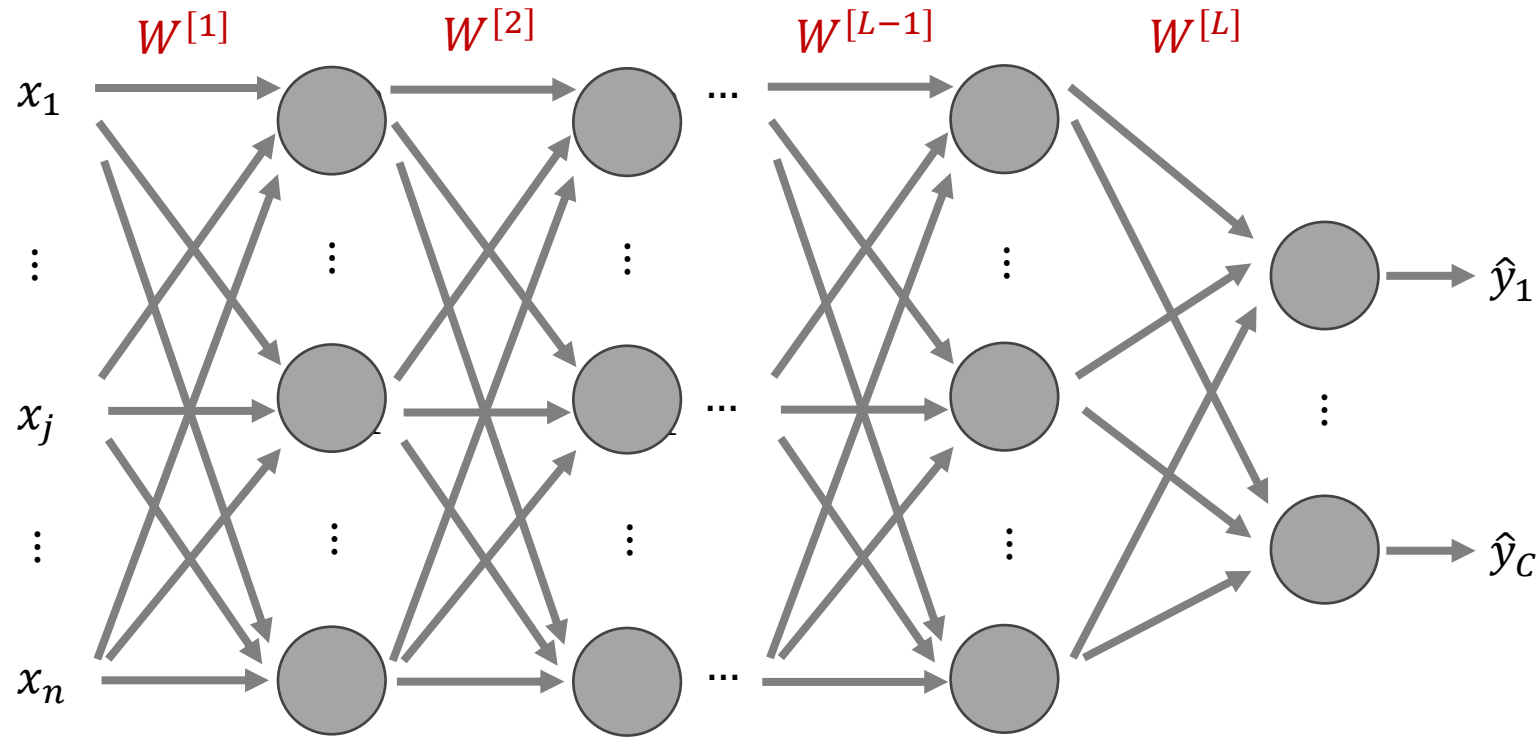
# Output (number of layer's neurons / output variables)

# Neural Networks and Matrix Multiplication (2)



$$\hat{\mathbf{y}} = g^{[2]} \left( W^{[2]T} g^{[1]} \left( W^{[1]T} \mathbf{x} \right) \right) = g^{[1]} \left( \begin{bmatrix} W_{00}^{[2]} & W_{10}^{[2]} \\ W_{01}^{[2]} & W_{11}^{[2]} \\ W_{02}^{[2]} & W_{12}^{[2]} \end{bmatrix}^T g^{[2]} \left( \begin{bmatrix} W_{00}^{[1]} & W_{10}^{[1]} & W_{20}^{[1]} \\ W_{01}^{[1]} & W_{11}^{[1]} & W_{21}^{[1]} \\ W_{02}^{[1]} & W_{12}^{[1]} & W_{22}^{[1]} \end{bmatrix}^T \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \right) \right) = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \end{bmatrix}$$

# Neural Networks and Matrix Multiplication (3)



Forward Propagation

$$\hat{\mathbf{y}} = g^{[L]} \left( W^{[L]T} \dots g^{[L-1]} \left( W^{[L-1]T} \dots g^{[l]} \left( W^{[l]T} \dots g^{[2]} \left( W^{[2]T} g^{[1]} \left( W^{[1]T} \mathbf{x} \right) \right) \right) \right) \right) = \begin{bmatrix} \hat{y}_1 \\ \dots \\ \hat{y}_c \end{bmatrix}$$

# Outline

- Notation and Math Refresher
- Backpropagation
  - Backpropagation on different scenarios
  - Biological plausibility of backpropagation
- Automatic Differentiation
  - Reverse mode automatic differentiation
  - Comparison with other methods
- Introduction to PyTorch
  - Tensors
  - Modules & Functions
  - Loss function & Optimizers

# Outline

- **Notation and Math Refresher**
- Backpropagation
  - Backpropagation on different scenarios
  - Biological plausibility of backpropagation
- Automatic Differentiation
  - Reverse mode automatic differentiation
  - Comparison with other methods
- Introduction to PyTorch
  - Tensors
  - Modules & Functions
  - Loss function & Optimizers



# Notation

**Scalar:** not bolded, lower case

$x$

**Vector:** bolded, lower case

$$\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$

**Matrix:** bolded, upper case

$$\mathbf{X} = \begin{bmatrix} x_{11} & \cdots & x_{m1} \\ \vdots & \ddots & \vdots \\ x_{1n} & \cdots & x_{mn} \end{bmatrix}$$

$n$  = Number of features in  $\mathbf{x}$   
 $m$  = Number of instances in dataset

# Vector and Matrix Operations (1)

## Scalar-by-scalar:

- $y(x) = wx$  scale  $x$  by  $w$

## Scalar-by-vector:

- $y(x) = wx = w \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} wx_1 \\ wx_2 \end{bmatrix}$  scale  $x$  by  $w$

## Vector-by-vector:

- $y(x) = \mathbf{w} \cdot \mathbf{x} = \mathbf{w}^\top \mathbf{x} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = [w_1 \ w_2] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = w_1 x_1 + w_2 x_2$  weighted sum

# Vector and Matrix Operations (2)

**Vector-by-matrix:**

- $y(X) = Wx = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{pmatrix} w_{11}x_1 + w_{12}x_2 \\ w_{21}x_1 + w_{22}x_2 \end{pmatrix}$

**Matrix-by-matrix:**

- $Y(X) = WX = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{bmatrix} = \begin{bmatrix} w_{11}x_{11} + w_{12}x_{21} & w_{11}x_{12} + w_{12}x_{22} \\ w_{21}x_{11} + w_{22}x_{21} & w_{21}x_{12} + w_{22}x_{22} \end{bmatrix}$

**Hadamard product** ◦ (element-wise product):

- $Y(X) = W \circ X = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} \circ \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{bmatrix} = \begin{bmatrix} w_{11}x_{11} & w_{12}x_{12} \\ w_{21}x_{21} & w_{22}x_{22} \end{bmatrix}$

- For backpropagation (see later), element-wise multiplication between matrices

# Vector and Matrix Operations (1)

**Summation Series = Scalar**

$$\sum_{r=0}^n w_r x_r$$

$$w_1 x_1 + \cdots + w_r x_r + \cdots + w_n x_n$$

**Transposed Vector Multiplication = Scalar**

$$\mathbf{w}^T \mathbf{x} = [w_1 \quad \cdots \quad w_r \quad \cdots \quad w_n] \begin{bmatrix} x_1 \\ \vdots \\ x_r \\ \vdots \\ x_n \end{bmatrix}$$

**Vector Dot Product = Scalar**

$$\mathbf{w} \cdot \mathbf{x} = \begin{bmatrix} w_1 \\ \vdots \\ w_r \\ \vdots \\ w_n \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ \vdots \\ x_r \\ \vdots \\ x_n \end{bmatrix}$$

**Transposed Matrix Multiplication = Vector**

$$\mathbf{W}^T \mathbf{x} = \begin{bmatrix} w_{11} & \cdots & w_{1r} & \cdots & w_{1n} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ w_{r1} & \cdots & w_{rr} & \cdots & w_{rn} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ w_{n1} & \cdots & w_{nr} & \cdots & w_{nn} \end{bmatrix}^T \begin{bmatrix} x_1 \\ \vdots \\ x_r \\ \vdots \\ x_n \end{bmatrix}$$

Same

Multi-rows

# Derivatives

## Total Derivative: $d$

- $\frac{dy}{dx}$  is the derivative of  $y$  relative to  $x$

## Partial derivative: $\partial$

- $\frac{\partial y}{\partial x_1}$  is the derivative of  $y$  relative to  $x_1$
- But  $y$  also depends on other variables (e.g.,  $x_2$  so, we can also calculate  $\frac{\partial y}{\partial x_2}$ )

## Jacobian: $\nabla$

- To calculate the derivative relative to *all*  $x_1$  and  $x_2$  together
- $\nabla y(\mathbf{x})$  is the gradient of  $y$  relative to all variables  $\mathbf{x} = [x_1, \dots, x_n]^\top$

$$\nabla y(\mathbf{x}) = \frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \partial y / \partial x_1 \\ \vdots \\ \partial y / \partial x_n \end{bmatrix} = \begin{bmatrix} \frac{\partial y}{\partial x_1} & \dots & \frac{\partial y}{\partial x_n} \end{bmatrix}^\top$$

# Matrix Calculus

Scalar-by-Vector (= 1D Vector)

$$\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \vdots \\ \frac{\partial y}{\partial x_n} \end{bmatrix}$$

Scalar-by-Matrix (= 2D Matrix)

$$\frac{\partial y}{\partial \mathbf{X}} = \begin{bmatrix} \frac{\partial y}{\partial x_{11}} & \cdots & \frac{\partial y}{\partial x_{1m}} \\ \vdots & \ddots & \vdots \\ \frac{\partial y}{\partial x_{n1}} & \cdots & \frac{\partial y}{\partial x_{nm}} \end{bmatrix}$$

Vector-by-Vector (= 2D Matrix)

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_N}{\partial x_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_n} & \cdots & \frac{\partial y_N}{\partial x_n} \end{bmatrix}$$

Vector-by-Matrix (= 3D Matrix)

$$\frac{\partial \mathbf{y}}{\partial \mathbf{X}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_{11}} & \cdots & \frac{\partial y_1}{\partial x_{1m}} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_{n1}} & \cdots & \frac{\partial y_1}{\partial x_{nm}} \end{bmatrix} \cdots \begin{bmatrix} \frac{\partial y_N}{\partial x_{11}} & \cdots & \frac{\partial y_N}{\partial x_{1m}} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_N}{\partial x_{n1}} & \cdots & \frac{\partial y_N}{\partial x_{nm}} \end{bmatrix}$$

Along 3<sup>rd</sup> dimension

# Outline

- Notation and Math Refresher
- **Backpropagation**
  - Backpropagation on different scenarios
  - Biological plausibility of backpropagation
- Automatic Differentiation
  - Reverse mode automatic differentiation
  - Comparison with other methods
- Introduction to PyTorch
  - Tensors
  - Modules & Functions
  - Loss function & Optimizers

# Gradient Descent

- Start at some  $w$
- Pick a nearby  $w$  that reduces  $J(w)$

$$w_j \leftarrow w_j - \gamma \frac{\partial J(w_0, w_1, \dots)}{\partial w_j}$$

- Repeat until minimum is reached

Learning Rate

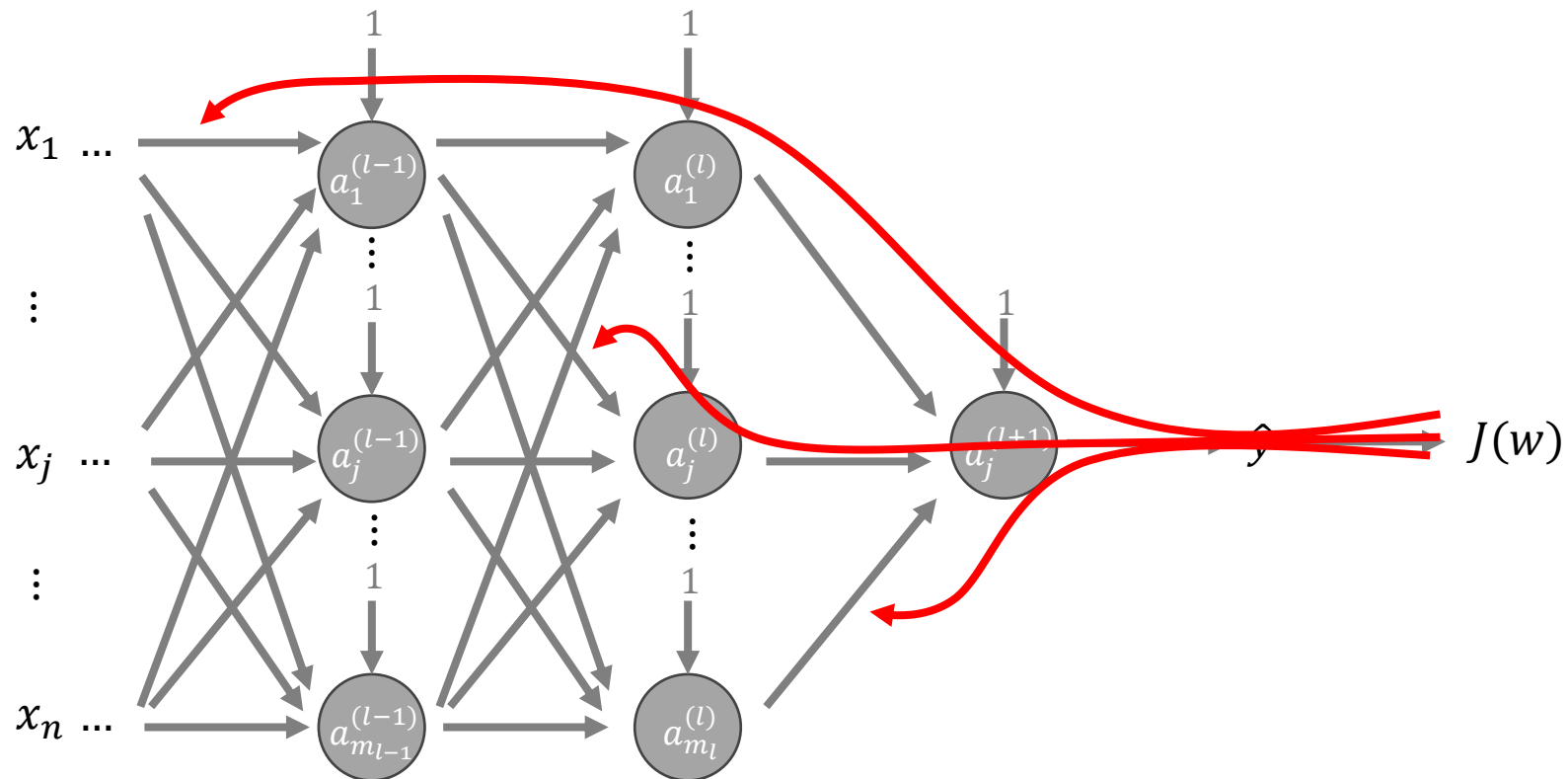
**Single-layer** Neural Networks with Sigmoid

$$w_i \leftarrow w_i - \gamma(\hat{y} - y)\hat{y}(1 - \hat{y})x_i$$

**Multi-layer** Neural Networks?



# Multi-layer Neural Networks



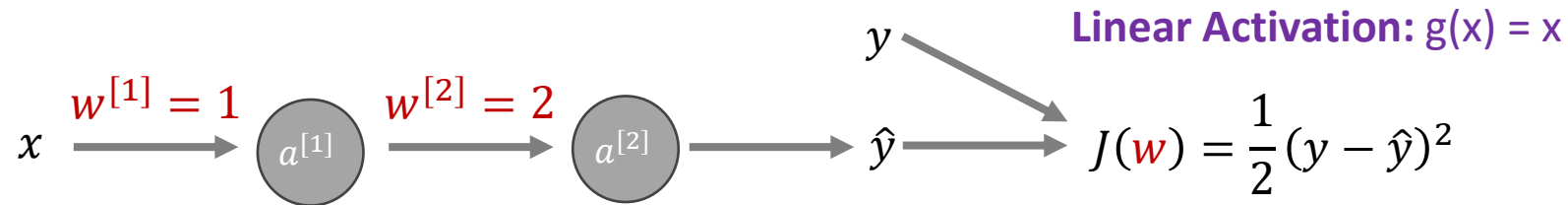
$$w_j \leftarrow w_j - \gamma \frac{\partial J(w_0, w_1, \dots)}{\partial w_j}$$

How to get  $\frac{\partial J(w_0, w_1, \dots)}{\partial w_j}$  for all  $w_j$ ?

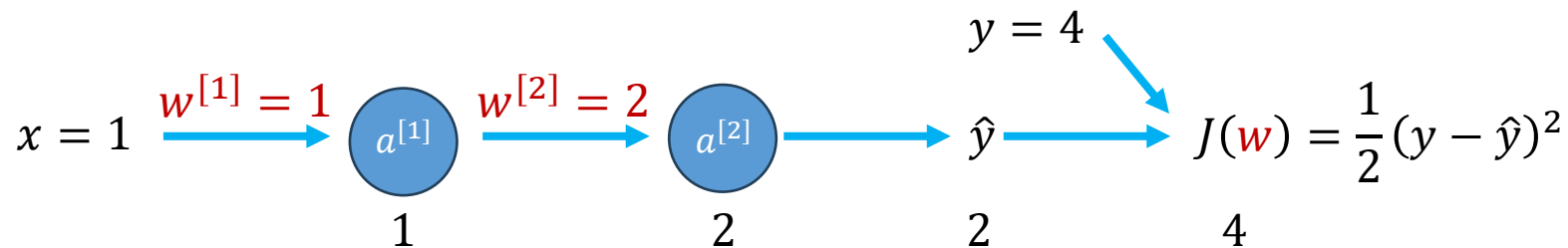
**Want:**

- Can be implemented as a program
- Efficient

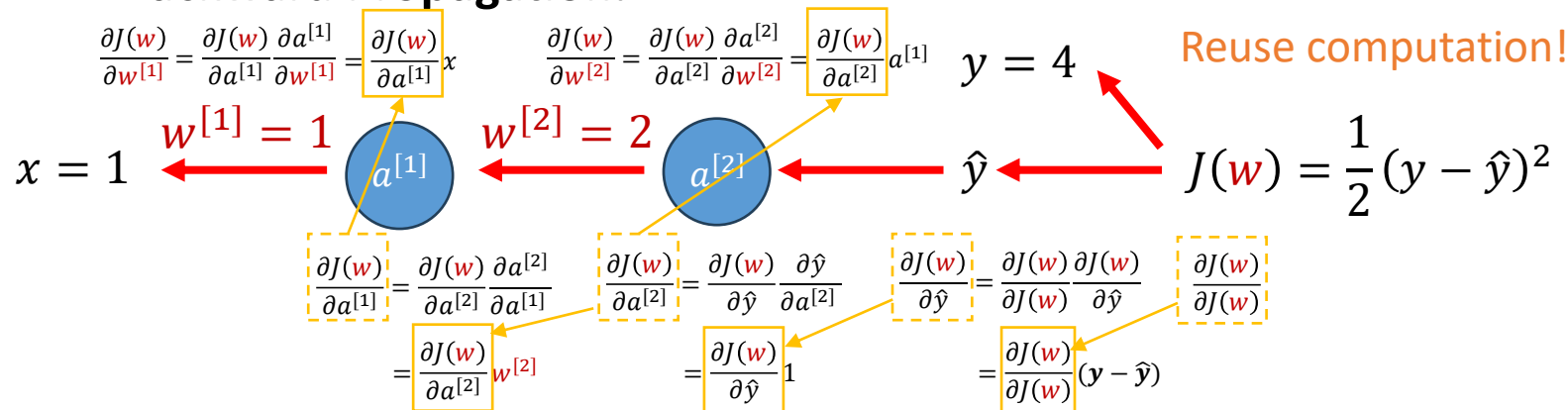
# Backpropagation



**Forward Propagation:**



**Backward Propagation:**



Let  $v_1, \dots, v_N$  be a *topological ordering* of the computation graph (i.e., parents comes before children).

$$v_N = J(\mathbf{w})$$

**Forward propagation**

For  $i$  in  $1, \dots, N$

Compute  $v_i$

**Backward propagation**

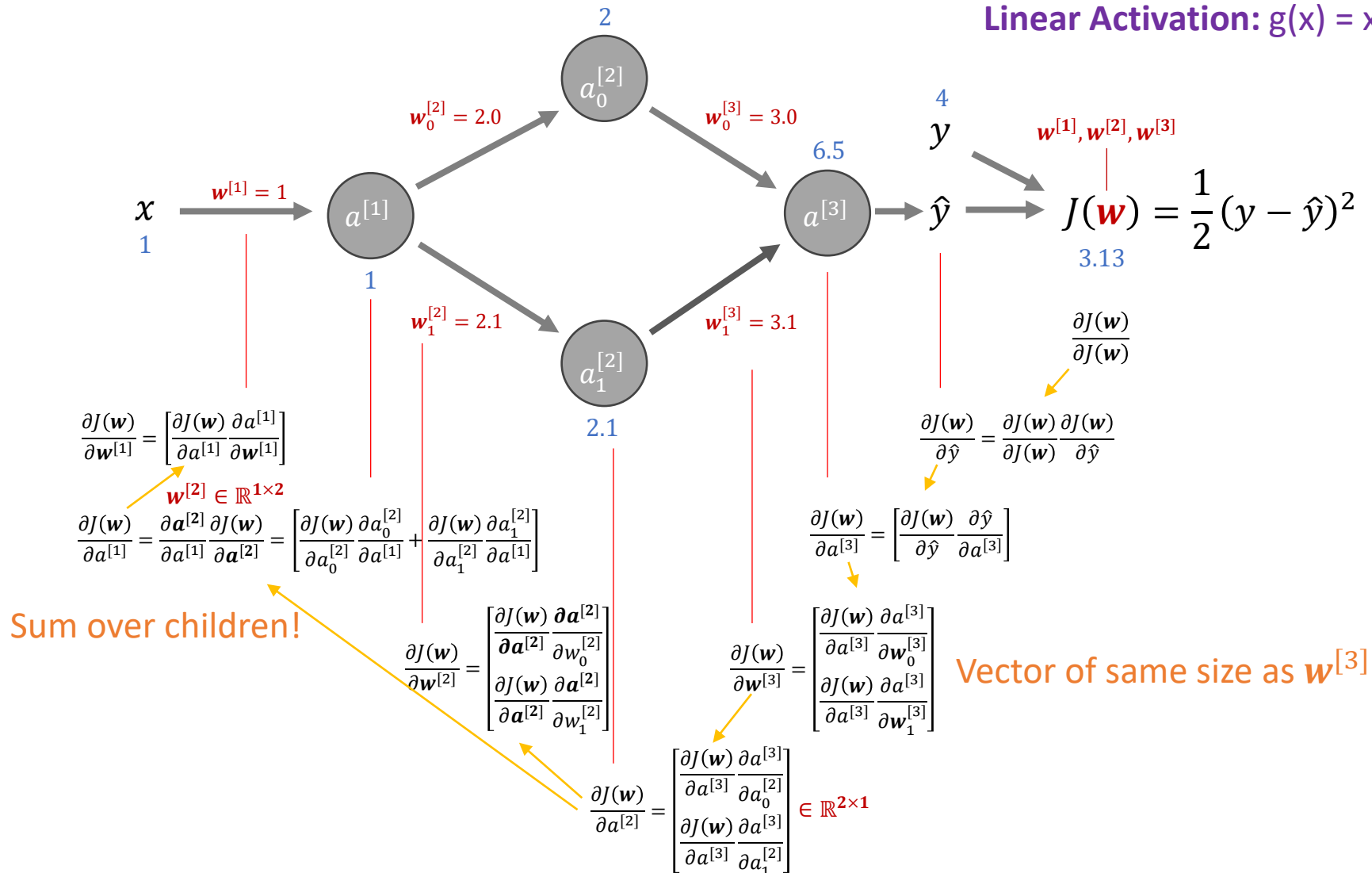
For  $i$  in  $N - 1, \dots, 1$

$$\frac{\partial J(\mathbf{w})}{\partial v_i} = \sum_{v_j \in Ch(v_i)} \frac{\partial J(\mathbf{w})}{\partial v_j} \frac{\partial v_j}{\partial v_i}$$

Children

# Backpropagation with Branches

Linear Activation:  $g(x) = x$



Let  $v_1, \dots, v_N$  be a *topological ordering* of the computation graph (i.e., parents comes before children).

$$v_N = J(\mathbf{w})$$

## Forward propagation

For  $i$  in  $1, \dots, N$

Compute  $v_i$

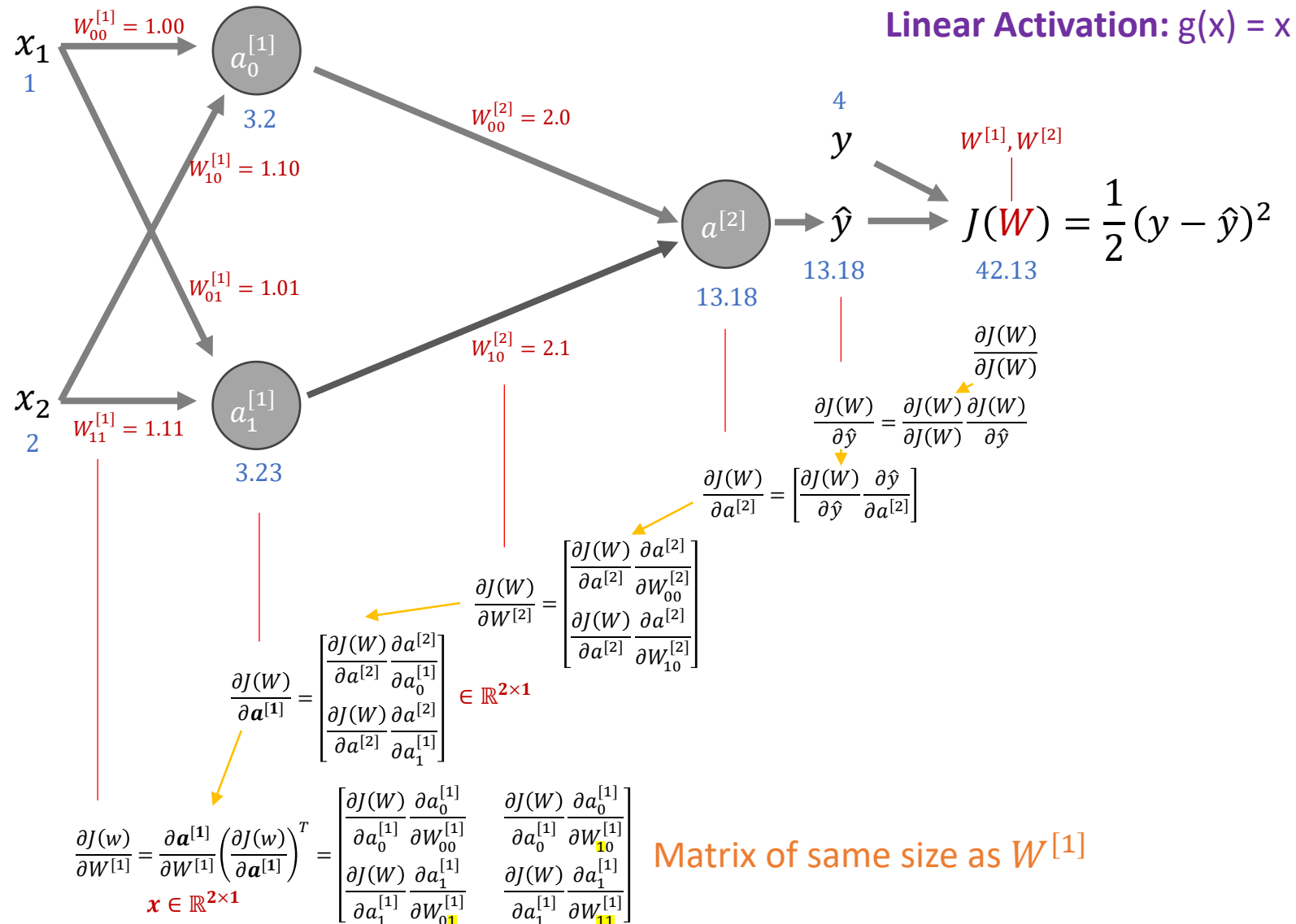
## Backward propagation

For  $i$  in  $N - 1, \dots, 1$

$$\frac{\partial J(\mathbf{w})}{\partial v_i} = \sum_{v_j \in \text{Children}(v_i)} \frac{\partial J(\mathbf{w})}{\partial v_j} \frac{\partial v_j}{\partial v_i}$$

Children

# Backpropagation with Many Features



Let  $v_1, \dots, v_N$  be a *topological ordering* of the computation graph (i.e., parents comes before children).

$$v_N = J(W)$$

## Forward propagation

For  $i$  in  $1, \dots, N$

Compute  $v_i$

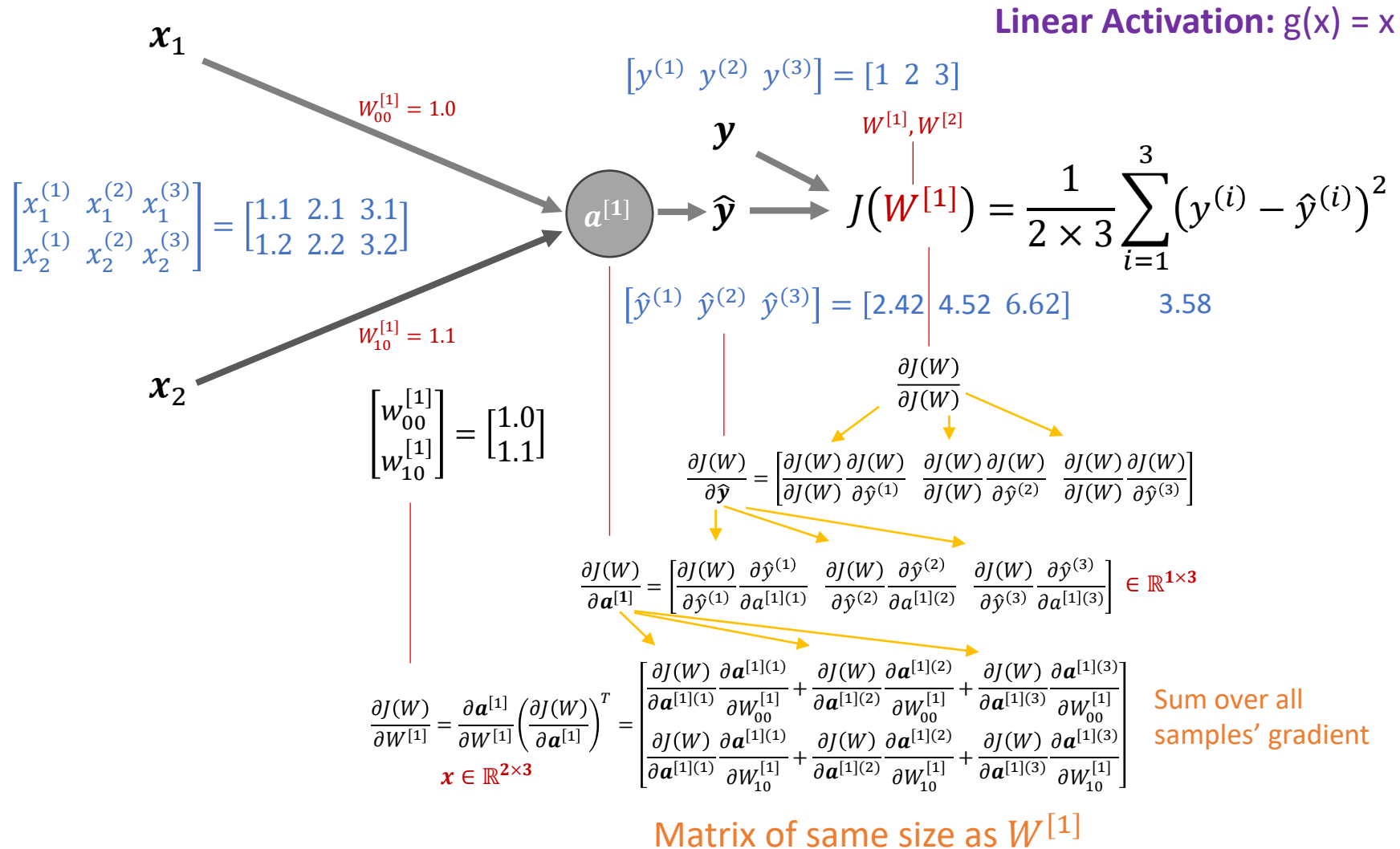
## Backward propagation

For  $i$  in  $N - 1, \dots, 1$

$$\frac{\partial J(W)}{\partial v_i} = \sum_{v_j \in \text{Children}(v_i)} \frac{\partial J(W)}{\partial v_j} \frac{\partial v_j}{\partial v_i}$$

Children

# Backpropagation with Many Samples



Let  $v_1, \dots, v_N$  be a *topological ordering* of the computation graph (i.e., parents comes before children).

$$v_N = J(W)$$

## Forward propagation

For  $i$  in  $1, \dots, N$

Compute  $v_i$

## Backward propagation

For  $i$  in  $N - 1, \dots, 1$

$$\frac{\partial J(W)}{\partial v_i} = \sum_{v_j \in \text{Children}(v_i)} \frac{\partial J(W)}{\partial v_j} \frac{\partial v_j}{\partial v_i}$$

Children

# Backpropagation in (Generalized) Matrix Form

$$\hat{y} = h(x) = g^{[L]} \left( f^{[L]} \left( g^{[L-1]} \left( \dots \left( g^{[l]} \left( f^{[l]} \left( g^{[l-1]} \left( \dots \left( g^{[1]} \left( f^{[1]}(x^{[0]}) \right) \right) \right) \right) \right) \right) \right) \right) \right) \right) \quad \begin{array}{l} \text{Activation function} \\ \text{Weighted sum} \end{array}$$

$$\mathbf{a}^{[l]} = g^{[l]}(f^{[l]})$$

$$f^{[l]} = (\mathbf{W}^{[l]})^T \mathbf{a}^{[l-1]}$$

$$\frac{\partial g^{[L]}}{\partial \mathbf{W}^{[L]}} = \frac{\partial f^{[L]}}{\partial \mathbf{W}^{[L]}} \left( \frac{\partial g^{[L]}}{\partial f^{[L]}} \right)^T \quad \delta^{[L]}$$

$$\frac{\partial g^{[L]}}{\partial \mathbf{W}^{[l+1]}} = \frac{\partial f^{[l+1]}}{\partial \mathbf{W}^{[l+1]}} \left( \frac{\partial g^{[l+1]}}{\partial f^{[l+1]}} \cdots \frac{\partial f^{[L]}}{\partial g^{[L-1]}} \frac{\partial g^{[L]}}{\partial f^{[L]}} \right)^T \quad \delta^{[l+1]}$$

$$\frac{\partial g^{[L]}}{\partial \mathbf{W}^{[l]}} = \frac{\partial f^{[l]}}{\partial \mathbf{W}^{[l]}} \left( \frac{\partial g^{[l]}}{\partial f^{[l]}} \frac{\partial f^{[l+1]}}{\partial g^{[l]}} \frac{\partial g^{[l+1]}}{\partial f^{[l+1]}} \cdots \frac{\partial f^{[L]}}{\partial g^{[L-1]}} \frac{\partial g^{[L]}}{\partial f^{[L]}} \right)^T$$

$$= \frac{\partial f^{[l]}}{\partial \mathbf{W}^{[l]}} \left( \frac{\partial g^{[l]}}{\partial f^{[l]}} \frac{\partial f^{[l+1]}}{\partial g^{[l]}} \delta^{[l+1]} \right)^T \quad \delta^{[l]}$$

$$\delta^{[l]} = g'^{[l]}(f^{[l]}) \mathbf{W}^{[l+1]} \delta^{[l+1]} = g'^{[l]}(f^{[l]}) \circ (\mathbf{W}^{[l+1]} \delta^{[l+1]})$$

$$\frac{\partial f^{[l]}}{\partial \mathbf{W}^{[l]}} = \mathbf{a}^{[l-1]}$$

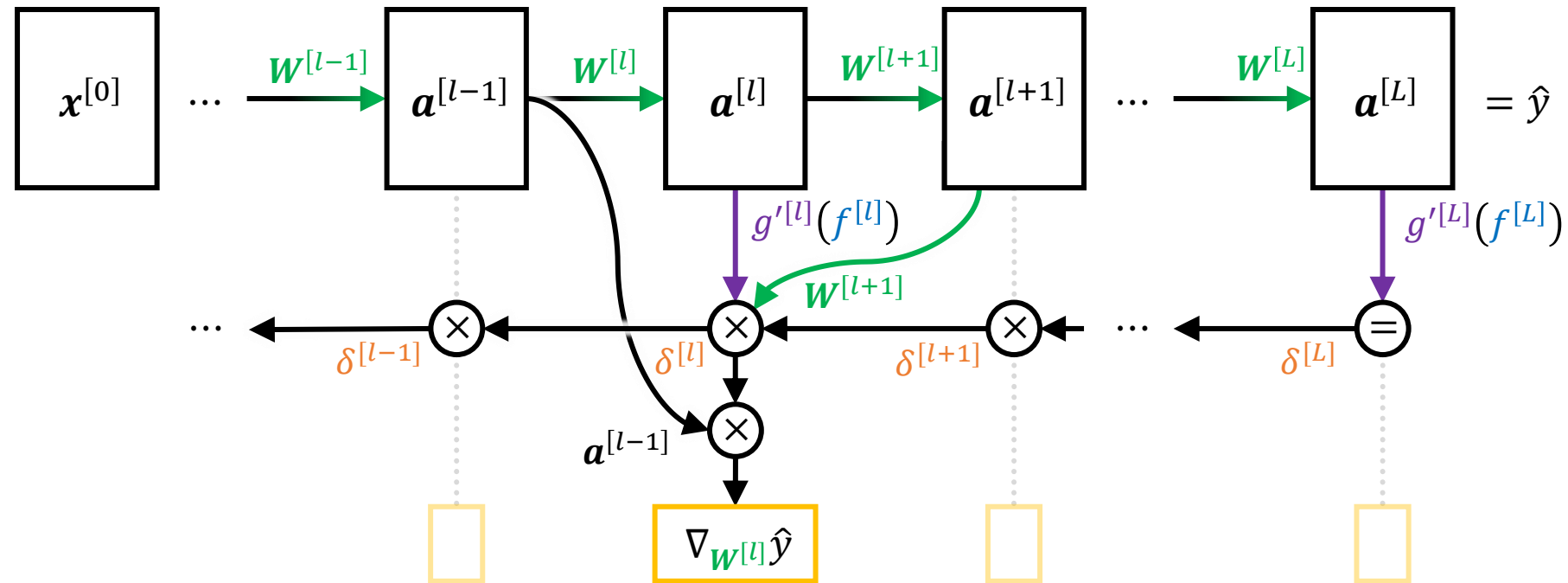
$$= \mathbf{a}^{[l-1]} (\delta^{[l]})^T$$

$$\frac{\partial g^{[l]}}{\partial f^{[l]}} = g'^{[l]}(f^{[l]})$$

$$\frac{\partial f^{[l+1]}}{\partial g^{[l]}} = \frac{\partial f^{[l+1]}}{\partial \mathbf{a}^{[l]}} = \mathbf{W}^{[l+1]}$$

$$\frac{\partial g^{[L]}}{\partial \mathbf{W}^{[l]}} = \mathbf{a}^{[l-1]} (\delta^{[l]})^T$$

# Backpropagation in (Generalized) Matrix Form



$$\nabla_{\mathbf{W}^{[l]}} \hat{\mathbf{y}} = \mathbf{a}^{[l-1]} (\delta^{[l]})^\top$$

$$\delta^{[l]} = [g'^{[l]}(f^{[l]})] \circ (\mathbf{W}^{[l+1]} \delta^{[l+1]})$$

Matrix form is only for **notational convenience**. When in doubts, refer to the element wise derivatives

# Notes on Backpropagation

- Backpropagation is used to train most neural networks
- Despite its success, backprop is believed to be biologically implausible
  - **Synapses are unidirectional**, backprop runs two passes: forward and backward
  - How **derivative** is calculated in each **neuron**, like backprop, is unclear
  - Backpropagation path is **linear**, neurons aren't linear
  - Backpropagation signal is **instantaneous**, brain is not



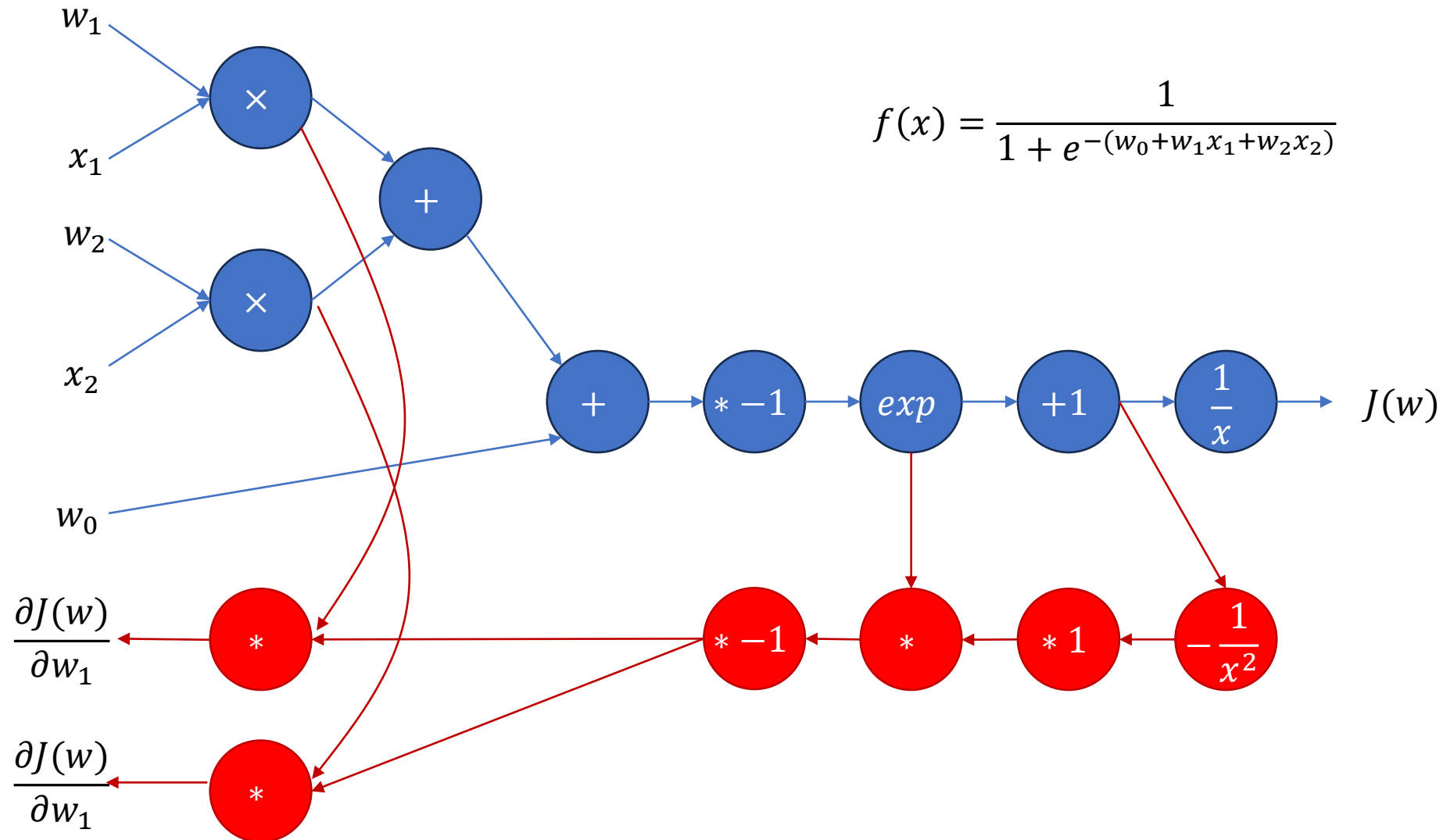
# Outline

- Notation and Math Refresher
- Backpropagation
  - Backpropagation on different scenarios
  - Biological plausibility of backpropagation
- **Automatic Differentiation**
  - Reverse mode automatic differentiation
  - Comparison with other methods
- Introduction to PyTorch
  - Tensors
  - Modules & Functions
  - Loss function & Optimizers

# Automatic Differentiation (AD)

- Backpropagation is a **special case** of automatic differentiation
  - Applied to neural networks, i.e.,  $\mathbb{R}^N \rightarrow \mathbb{R}$  (many features to one loss)
  - In practice, we use AD for neural networks
- AD has two modes: **forward** mode and **reverse** mode
  - Backpropagation is a special case of reverse mode AD
- Example library:
  - **PyTorch autograd**
    - Just need to implement the forward pass (e.g., layers), the backward pass is done automatically
  - **Jax grad**
  - **Tensorflow autodiff**

# Reverse Mode Differentiation



# AD vs Symbolic Differentiation

What we have done in the previous lectures!

$n$	$l_n$	$\frac{d}{dx} l_n$	$\frac{d}{dx} l_n$ (Simplified form)
1	$x$	1	1
2	$4x(1-x)$	$4(1-x) - 4x$	$4 - 8x$
3	$16x(1-x)(1-2x)^2$	$16(1-x)(1-2x)^2 - 16x(1-2x)^2 - 64x(1-x)(1-2x)$	$16(1 - 10x + 24x^2 - 16x^3)$
4	$64x(1-x)(1-2x)^2(1-8x+8x^2)^2$	$128x(1-x)(-8+16x)(1-2x)^2(1-8x+8x^2) + 64(1-x)(1-2x)^2(1-8x+8x^2)^2 - 64x(1-2x)^2(1-8x+8x^2)^2 - 256x(1-x)(1-2x)(1-8x+8x^2)^2$	$64(1 - 42x + 504x^2 - 2640x^3 + 7040x^4 - 9984x^5 + 7168x^6 - 2048x^7)$

- Disadvantages?
  - Symbolic differentiation results in **complex** and redundant expressions
  - Can **only handle math**, not procedure
- The goal is to **compute derivatives**, not finding a derivative formula

# AD vs Numerical Differentiation

$$\frac{\partial f(x_1, \dots, x_N)}{\partial x_i} \approx \frac{f(x_1, \dots, x_i + \epsilon, \dots, x_N) - f(x_1, \dots, x_i - \epsilon, \dots, x_N)}{2\epsilon}$$

## Numerical Differentiation

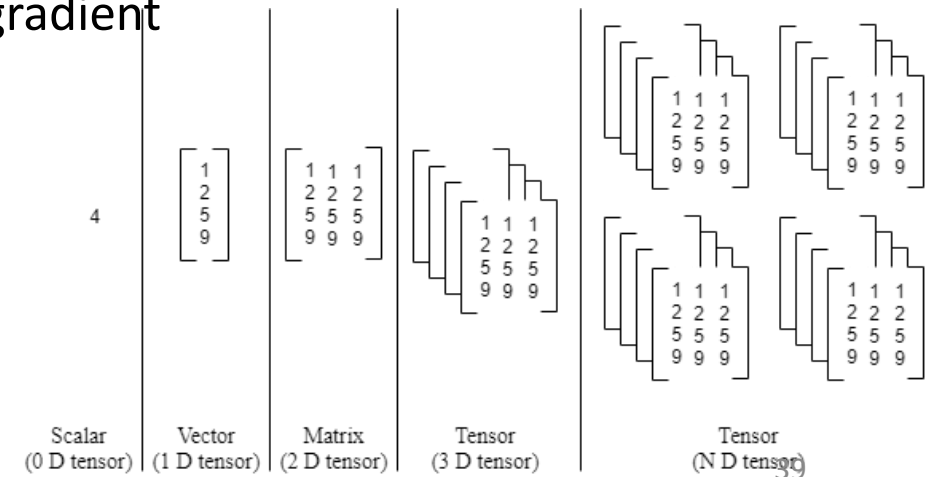
- **Expensive**, need to do forward pass for each derivative
- Introduce **numerical errors**
- Normally, only used for testing

# Outline

- Notation and Math Refresher
- Backpropagation
  - Backpropagation on different scenarios
  - Biological plausibility of backpropagation
- Automatic Differentiation
  - Reverse mode automatic differentiation
  - Comparison with other methods
- **Introduction to PyTorch**
  - Tensors
  - Modules & Functions
  - Loss function & Optimizers

# Tensors

- N-dimensional (e.g., 1D, 2D, 3D, ...) array representation
- **Similar to Numpy** arrays, with GPU support
- Contain information about computational graph
- Important functions:
  - `torch.tensor(..., requires_grad=True)`
    - Build tensor, and possibly set the tensor to require gradient
  - `backward()`
    - Performs backpropagation



# Tensors: Example

```
import torch  # Import PyTorch

x = torch.tensor([0.])
w1 = torch.tensor([0.], requires_grad=True)
w2 = torch.tensor([0.], requires_grad=True)

y = w2 * torch.sigmoid(w1 * x)
y.backward()

print(x, w1, w2)    tensor([0.]) tensor([0.], requires_grad=True) tensor([0.], requires_grad=True)
print(w1.grad)      tensor([0.])
print(w2.grad)      tensor([0.5000])
print(x.grad)       None
```



# Modules and Functions API

## Neural Networks Module (`torch.nn`)

- Containers

- `Module (torch.nn.Module)`
- `Sequential (torch.nn.Sequential)`

- Linear Layers

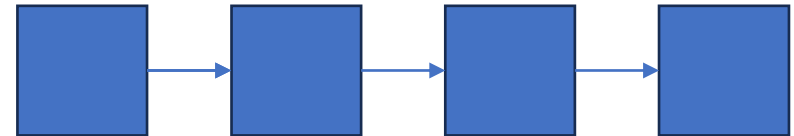
- Linear: Single Layer NN without activation (`torch.nn.Linear`)

- Non-linear activation functions

- `ReLU (torch.nn.ReLU)`
- `Sigmoid (torch.nn.Sigmoid)`
- `Softmax (torch.nn.Softmax)`

## Functional version (`torch.nn.functional`)

```
def __init__(self, ...):  
    ...  
def forward(self, ...):  
    ...
```

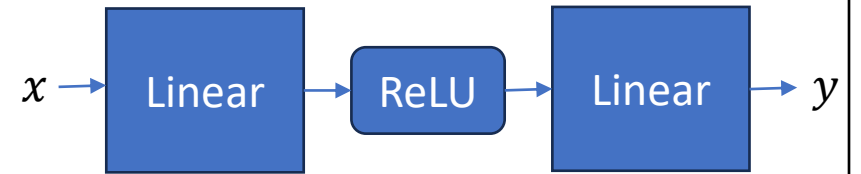


# Modules and Functions API: Example

```
class NeuralNetRegressor(torch.nn.Module):  
    def __init__(self, input_size, hidden_size):  
        super().__init__()  
        self.linear1 = torch.nn.Linear(input_size, hidden_size)  
        self.linear2 = torch.nn.Linear(hidden_size, 1)  
        self.relu = torch.nn.ReLU()  
  
    def forward(self, x):  
        f1 = self.linear1(x)  
        a1 = self.relu(f1)  
        f2 = self.linear2(a1)  
        return f2
```

```
model1 = NeuralNetRegressor(2, 8) # 2 features, 8 hidden neurons
```

```
model2 = torch.nn.Sequential(torch.nn.Linear(2, 8), torch.nn.ReLU(), torch.nn.Linear(8, 1)) # same
```



```
w1 = torch.tensor(8, 2, requires_grad=True)  
w2 = torch.tensor(2, 1, requires_grad=True)  
  
def neural_net_regressor(x): # also the same  
    f1 = torch.nn.functional.linear(x, w1)  
    a1 = torch.nn.functional.relu(f1)  
    return torch.nn.functional.linear(a1, w2)
```

# Loss Functions

- Mean Squared Error (`torch.nn.MSELoss`)
- Binary Cross Entropy (`torch.nn.BCELoss`)
- Cross Entropy (`torch.nn.CrossEntropyLoss`)

```
loss_function = torch.nn.MSELoss()  
loss_function = torch.nn.BCELoss()  
loss_function = torch.nn.CrossEntropyLoss()
```

# Optimizers

## Optimizers (`torch.optim`)

- Stochastic Gradient Descent (`torch.optim.SGD`)
- Adam (`torch.optim.Adam`)

## Important functions:

- `optimizer.zero_grad()`
  - Set all gradients to zero, before computing gradient
- `optimizer.step()`
  - Update the weights, and let the optimizer know that one step of optimization is done

```
optimizer = torch.optim.SGD([w_1, w_2], lr=0.01)
optimizer = torch.optim.Adam([w_1, w_2], lr=0.01)
```

# Example 1: Training NN using functional API

```
model = neural_net_regressor // w1 and w2 are declared before

optimizer = torch.optim.Adam([w1, w2], lr=0.01)

for epoch in range(num_epochs):
    optimizer.zero_grad() ← Zero the gradients in the weight tensor
    y_pred = model(x)
    loss = torch.nn.functional.mse_loss(y_pred, y)
    loss.backward() ← Do backpropagation
    optimizer.step() ← Update the weights
```

## Example 2: Training NN using modular API

```
model = NeuralNetRegressor()
```

Retrieve all the weights in the model



```
optimizer = torch.optim.Adam(model.parameters(), lr=0.01)
```

```
loss_function = torch.nn.MSELoss()
```

```
for epoch in range(num_epochs):
```

```
    optimizer.zero_grad() ← Zero the gradients in the weight tensor
```

```
    y_pred = model(x)
```

```
    loss = loss_function(y_pred, y)
```

```
    loss.backward() ← Do backpropagation
```

```
    optimizer.step() ← Update the weights
```

# Example 3: Building and Training a NN

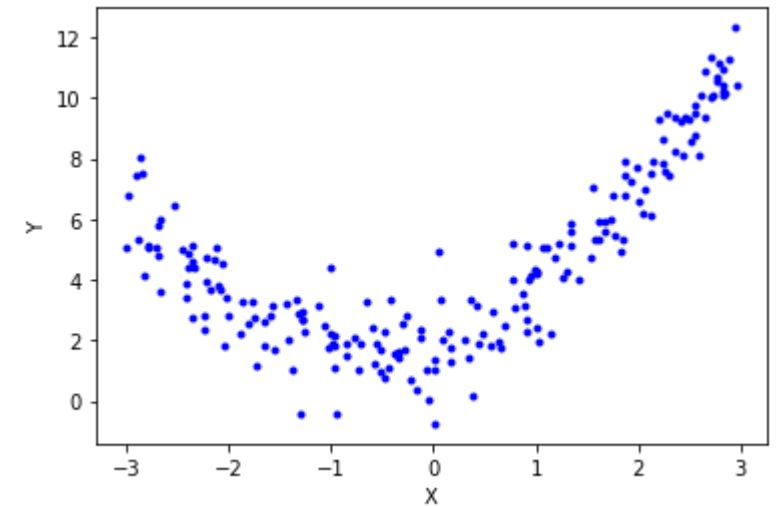
Step 0: Import the package

```
import torch
```

# Example 3: Building and Training a NN

## Step 1: Create a Dataset

```
# Generate 100 random data points  
# Normally distributed with 0 mean and 1 variance  
x = 5*torch.randn(100, 1)  
  
# Set true output y = f(x) = x^2  
y = torch.square(x)  
  
# Add noise to the true output y, noise is of the same shape as y  
# Noise is normally distributed with 0 mean and 1 variance  
y += torch.randn_like(y)
```



Credit: analyticsvidhya.com



# Example 3: Building and Training a NN

## Step 2: Prepare the Model and Optimizer

```
# Instantiate the model
model = NeuralNetRegressor(1, 8)

# Define the optimizer
optimizer = torch.optim.Adam(model.parameters(), lr=1e-1)

# Define loss using a predefined loss function
loss_function = torch.nn.MSELoss()
```

# Example 3: Building and Training a NN

## Step 3: Train the Model

```
n_epoch = 10000
for epoch in range(n_epoch):
    optimizer.zero_grad() # Set the gradients to 0

    y_pred = model(x) # Get the model predictions

    loss = loss_function(y_pred, y) # Get the loss

    if epoch%1000==0:
        print(f"Epoch {epoch}: traing loss: {loss}") # Print stats

    loss.backward() # Compute the gradients

    optimizer.step() # Take a step to optimize the weights
```

# Example 3: Building and Training a NN

## Step 4: Test the model

```
x2 = 5*torch.randn(5, 1) # Generate new data points
print(x2.tolist())
# [[5.367788791656494], [0.8833026885986328], [-0.6179562211036682], [4.864236831665039], [5.316009998321533]]

y2 = torch.square(x2) # Ground truth (true output) with no noise
print(y2.tolist())
# [[28.813156127929688], [0.7802236676216125], [0.3818698823451996], [23.660799026489258], [28.25996208190918]]

y_pred = model(x2)
print(y_pred.tolist())
# [[29.595108032226562], [0.8340979814529419], [0.5545129179954529], [23.755483627319336], [28.972129821777344]]
```

# Notes on Implementation

- If you want to implement custom layers, loss functions, etc, you have to make sure that they are **differentiable**.
- Non-differentiable functions can also be used, but it requires a gradient estimator to estimate the gradient for backpropagation.

# PyTorch Resources

## **Books on Deep Learning** (with PyTorch implementation)

- <https://d2l.ai>

## **Good online tutorial:**

- <https://pytorch.org/tutorials/>
- [https://web.stanford.edu/class/cs224n/materials/CS224N\\_PyTorch\\_Tutorial.html](https://web.stanford.edu/class/cs224n/materials/CS224N_PyTorch_Tutorial.html)

Feel free to try out the examples in the slides!

# Summary

- Backpropagation
  - Backpropagation on different scenarios:
    - Path, branches, many features, and many samples (sum the gradients)
  - Biological plausibility of backpropagation – believed to be **not feasible**
- Automatic Differentiation
  - Reverse mode automatic differentiation – **backprop is a special case**:  $\mathbb{R}^N \rightarrow \mathbb{R}$
  - Comparison with other methods: **symbolic** and **numerical** differentiation
- Introduction to PyTorch
  - Tensors
    - **n-dimensional array representation** with GPU support
    - Maintain **computational graph**
  - Modules & Functions: Linear (linear), ReLU (relu), etc – they are equivalent
  - Loss function & Optimizers

# Coming Up Next Week

- **Neural Networks (Deep Learning) Architectures**
  - Convolutional Neural Networks
  - Recurrent Neural Networks
  - Attention Neural Networks (the basis of **ChatGPT**) – maybe...
- **Applications of Deep learning**
  - Vision
  - Speech and Language
- **Issues with Deep Learning**
  - Overfitting
  - Exploding/vanishing gradients

# To Do

- **Lecture Training 9**
  - +100 Free EXP
  - +50 Early bird bonus