

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

# Lecture 9: Intro to Neural Networks

25 March 2025

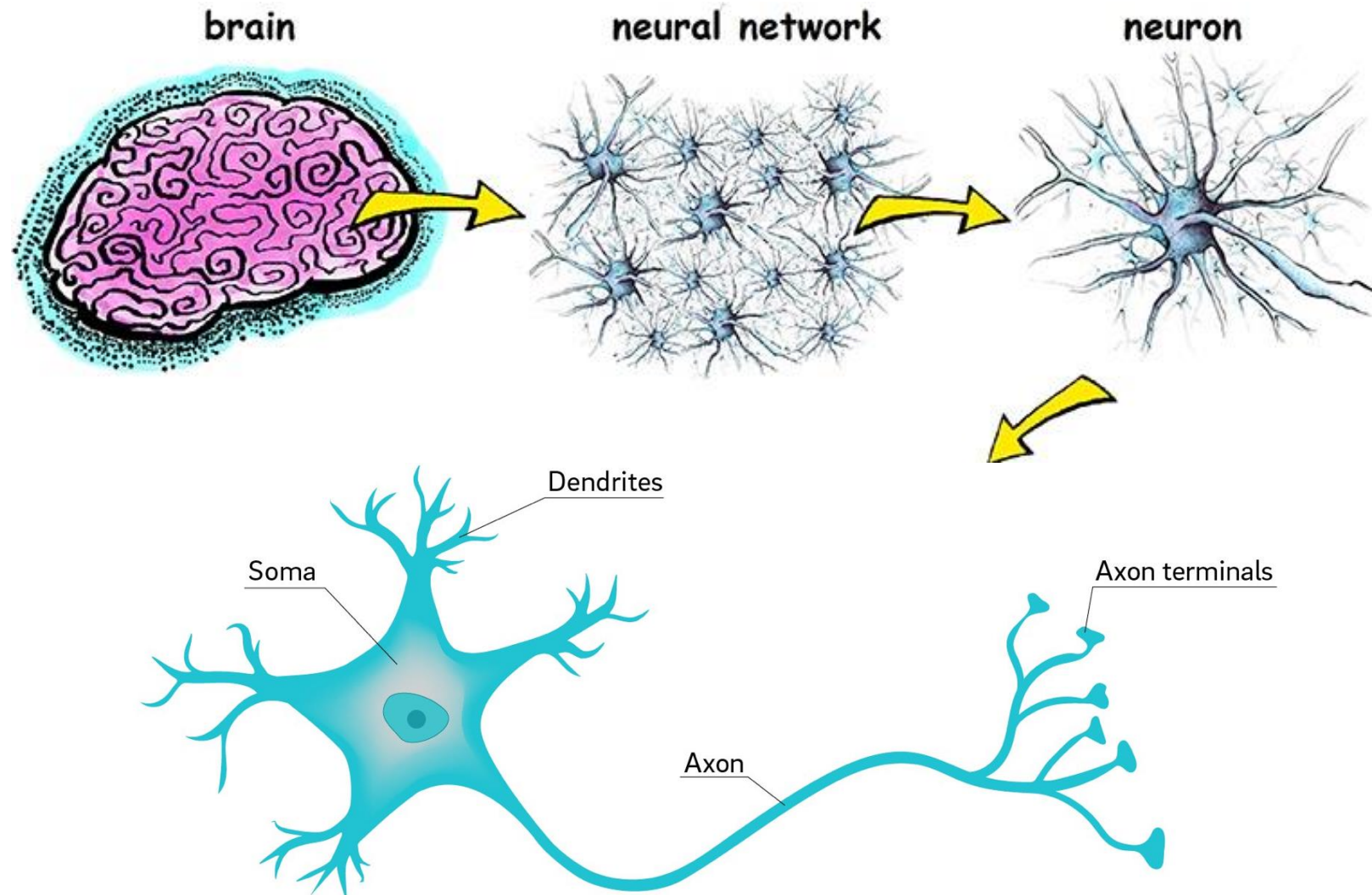
# Outline

- Perceptron
  - Biological inspiration
  - Perceptron Learning Algorithm
- Neural Network
  - Neuron
  - AND Gate Modelling
  - XNOR Gate Modelling
  - Single-layer and Multi-layer Neural Networks
- Multi-class Classification

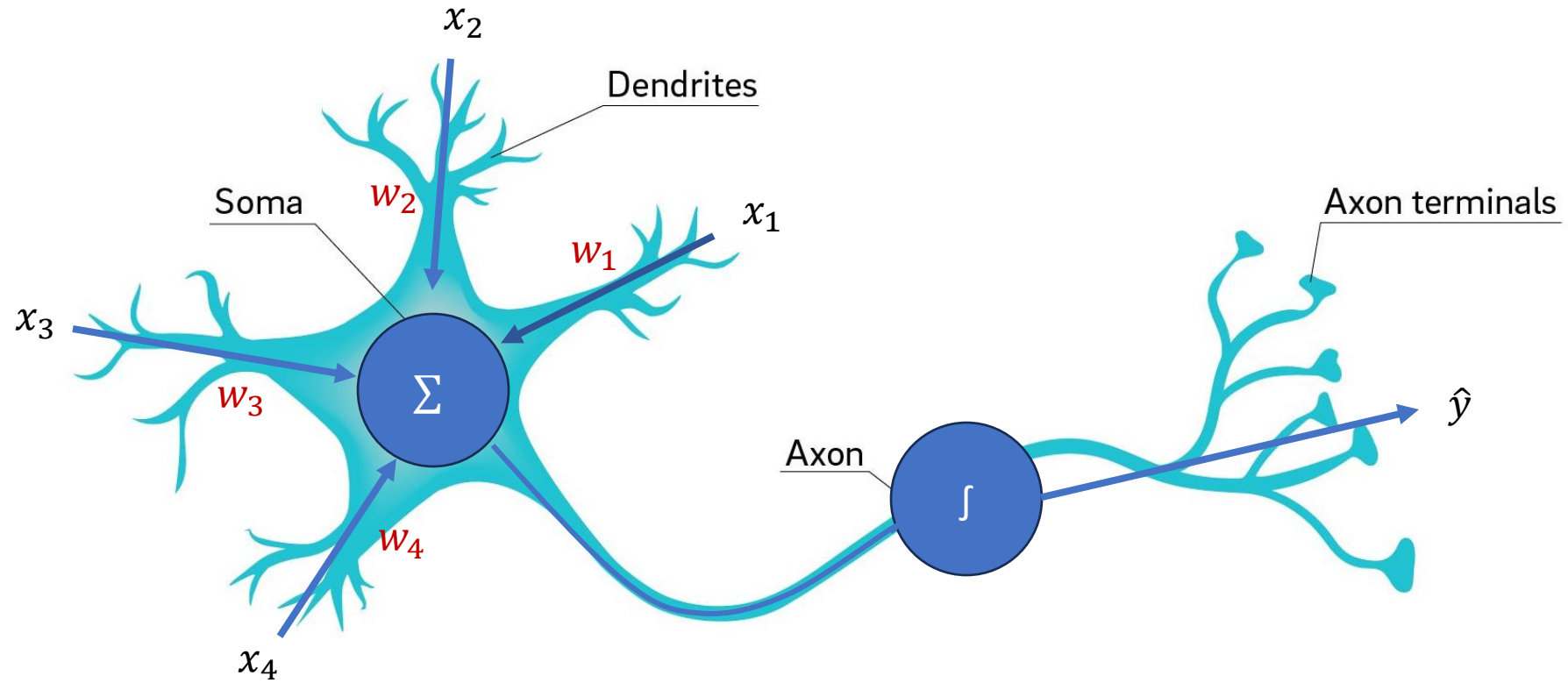
# Outline

- **Perceptron**
  - Biological inspiration
  - Perceptron Learning Algorithm
- Neural Network
  - Neuron
  - AND Gate Modelling
  - XNOR Gate Modelling
  - Single-layer and Multi-layer Neural Networks
- Multi-class Classification

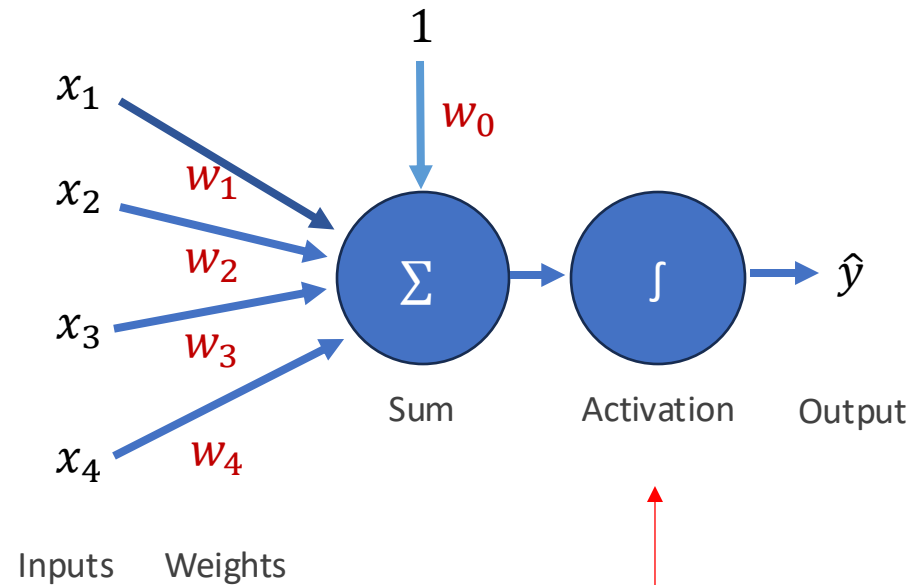
# Brain and Neuron



# Neuron



# Perceptron



Sign function

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

# Perceptron: Data

Suppose:

- We are given  $N$  data points.
- Each data point consists of **features** and a **target** variable.
- The features are described by a vector of **real numbers** in dimension  $d$ .
- The target is  $\{-1, 1\}$ , where  $-1$  is “negative” class and  $1$  is “positive” class.

# Perceptron: Data – Math

Suppose:

$$D = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\},$$

where for all  $i \in \{1, \dots, N\}$

Features:  $\mathbf{x}^{(i)} \in \mathbb{R}^d$

Targets:  $y^{(i)} \in \{-1, 1\}$



# Perceptron

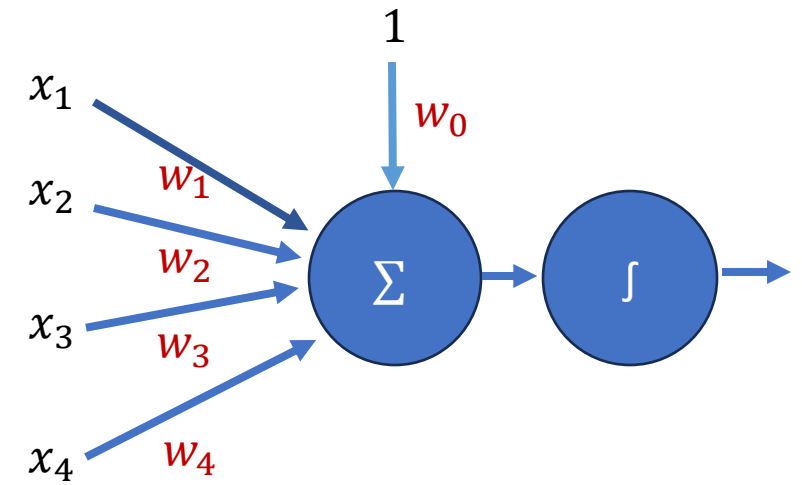
Given an input vector  $\mathbf{x}$  of dimension  $d$ , the Perceptron model is defined as:

$$h_{\mathbf{w}}(\mathbf{x}) = g(\mathbf{w}^T \mathbf{x}) = g\left(\sum_{j=0}^d w_j x_j\right) = g(w_0 x_0 + w_1 x_1 + w_2 x_2 + \cdots + w_d x_d)$$

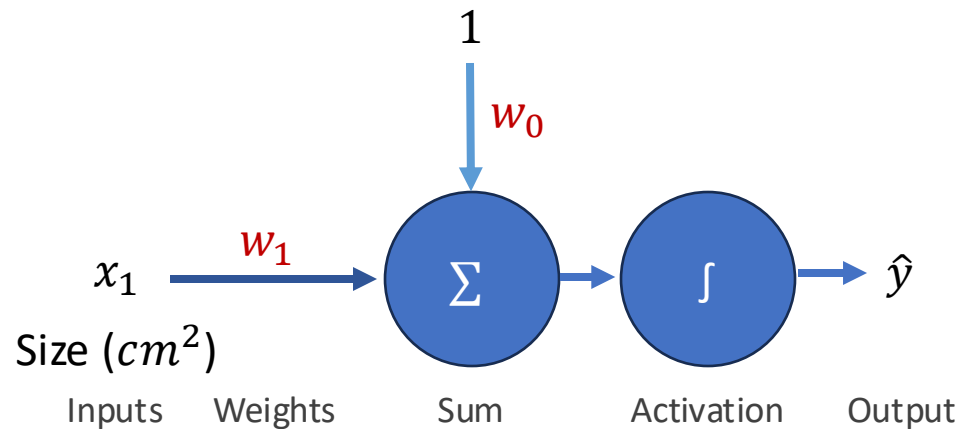
where  $w_0, \dots, w_d$  are **parameters/weights**,  $x_0 = 1$  is a dummy variable.

$g$  is sign function:

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



# Perceptron: An Example



Sign function

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

$$\hat{y} = h_{\mathbf{w}}(\mathbf{x}) = g\left(\sum_{j=0}^d w_j x_j\right)$$

$x_0 = 1$  (dummy variable)

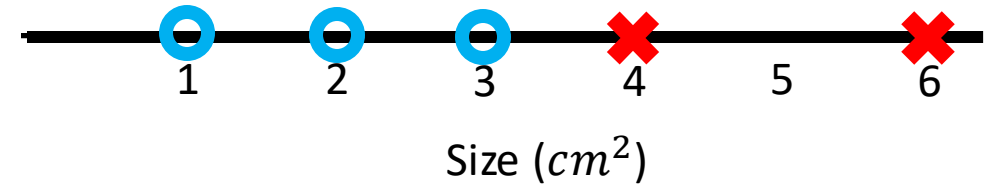
$$\hat{y} = \begin{cases} +1 & \text{if } w_0 x_0 + w_1 x_1 \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

$$w_0 = -1.5$$

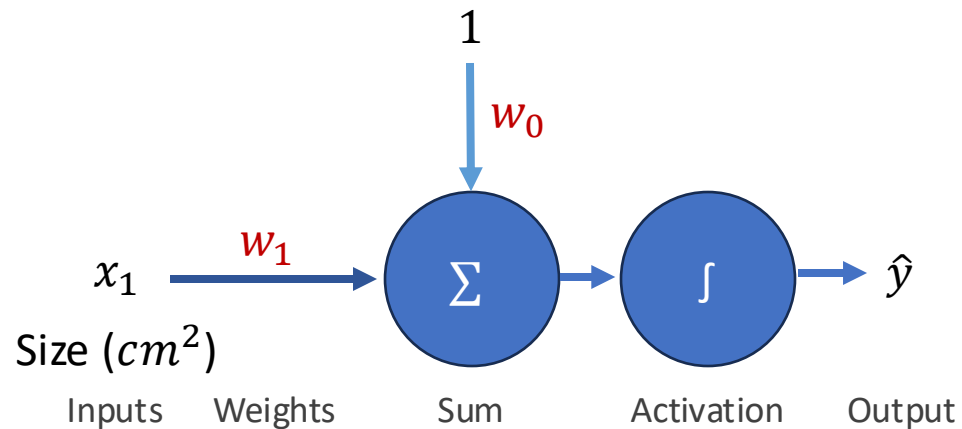
$$w_1 = 1$$

$$\hat{y} = \begin{cases} +1 & \text{if } -1.5 + 1x_1 \geq 0 \longrightarrow x_1 \geq 1.5 \\ -1 & \text{otherwise} \end{cases}$$

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



# Perceptron: An Example



Sign function

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

$$\hat{y} = h_{\mathbf{w}}(\mathbf{x}) = g\left(\sum_{j=0}^d w_j x_j\right)$$

$x_0 = 1$  (dummy variable)

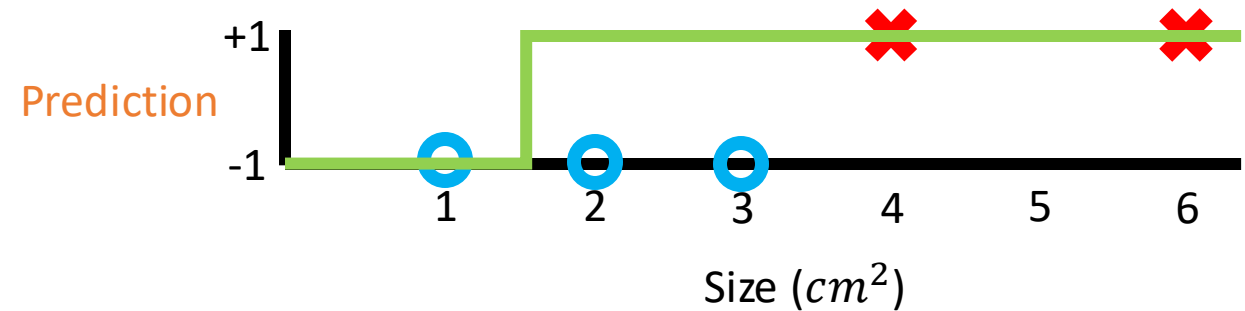
$$\hat{y} = \begin{cases} +1 & \text{if } w_0 x_0 + w_1 x_1 \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

$$w_0 = -1.5$$

$$w_1 = 1$$


$$\hat{y} = \begin{cases} +1 & \text{if } -1.5 + 1x_1 \geq 0 \longrightarrow x_1 \geq 1.5 \\ -1 & \text{otherwise} \end{cases}$$

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



Current model cannot predict all the labels correctly  
 -> Need to update  $\mathbf{w}$

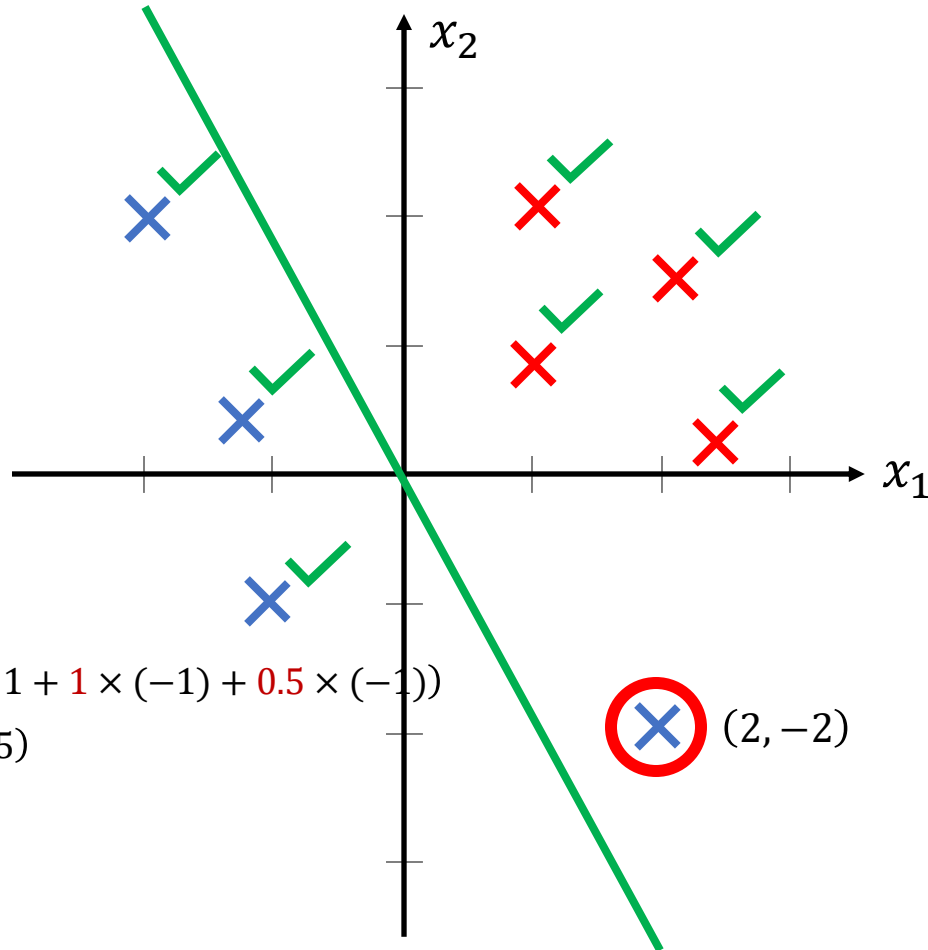
# Perceptron Learning Algorithm

- Initialize  $\mathbf{w}$
- Loop (until convergence or max steps reached)
  - For each instance  $(\mathbf{x}^{(i)}, y^{(i)})$ , classify  $\hat{y}^{(i)} = h_{\mathbf{w}}(\mathbf{x}^{(i)})$
  - Select one misclassified instance  $(\mathbf{x}^{(k)}, y^{(k)})$
  - Update weights:  $\mathbf{w} \leftarrow \mathbf{w} + \gamma(y^{(k)} - \hat{y}^{(k)})\mathbf{x}^{(k)}$   
  
Learning rate

Use  $y^{(i)}$  during training,  
supervised Learning!

# Perceptron Learning Algorithm: An Example

$$0x_0 + 1x_1 + 0.5x_2 = 0 \quad \text{red } \times : +1 \quad \text{blue } \times : -1$$



$$\begin{aligned} &g(0 \times 1 + 1 \times (-1) + 0.5 \times (-1)) \\ &= g(-1.5) \\ &= -1 \end{aligned}$$

$$\hat{y} = h_{\mathbf{w}}(\mathbf{x}) = g\left(\sum_{j=0}^d w_j x_j\right)$$

$x_0 = 1$  (dummy variable)

Sign function

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

$$\begin{aligned} \hat{y} = h_{\mathbf{w}}(\mathbf{x}) &= g(w_0 x_0 + w_1 x_1 + w_2 x_2) \\ &= g(0x_0 + 1x_1 + 0.5x_2) \end{aligned}$$

weights  
initialization

$$g(0 \times 1 + 1 \times 2 + 0.5 \times (-2)) = g(1) = +1$$

misclassified  
instance

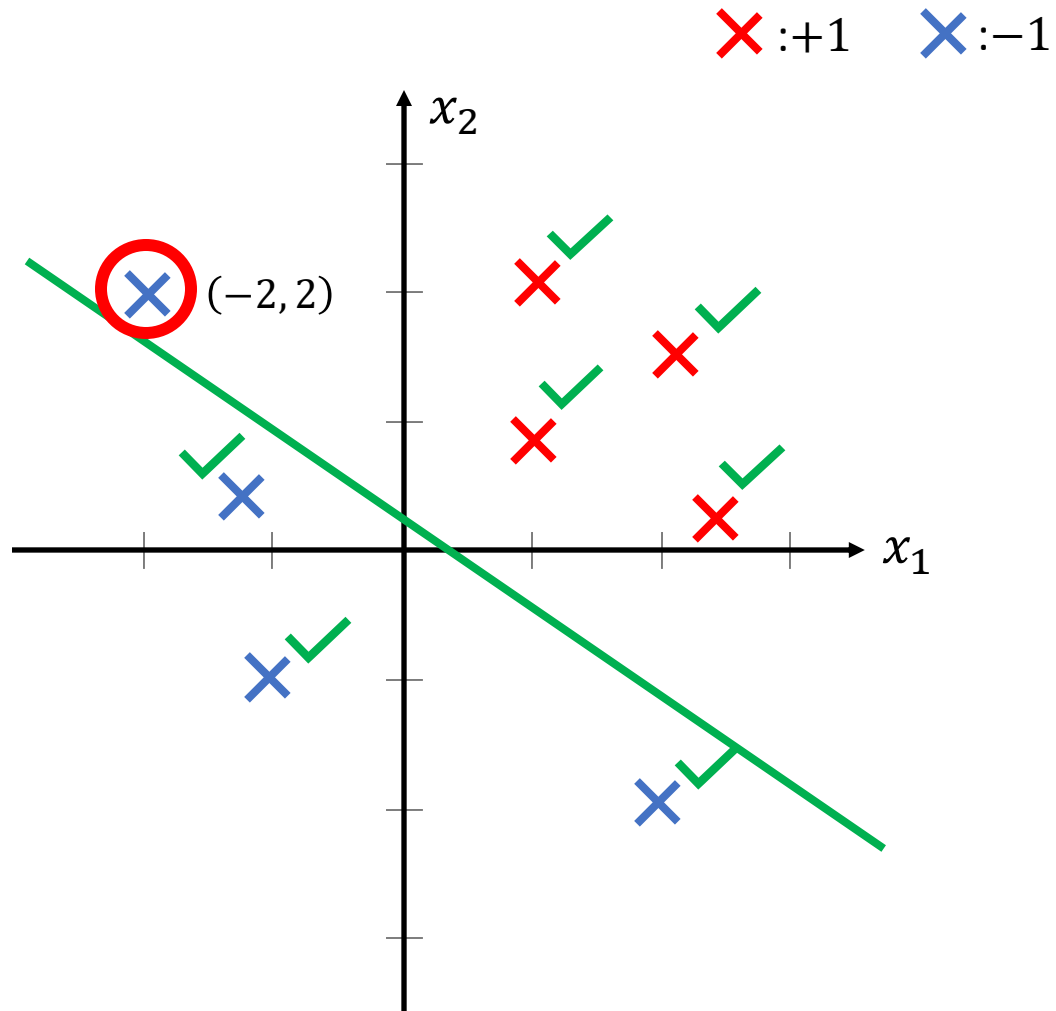
$$\mathbf{w} \leftarrow \mathbf{w} + \gamma(y^{(k)} - \hat{y}^{(k)}) \mathbf{x}^{(k)}$$

$$\begin{bmatrix} 0 \\ 1 \\ 0.5 \end{bmatrix} + 0.1(-1 - 1) \begin{bmatrix} 1 \\ 2 \\ -2 \end{bmatrix} = \begin{bmatrix} -0.2 \\ 0.6 \\ 0.9 \end{bmatrix}$$

$$\text{New } h_{\mathbf{w}}(\mathbf{x}) = g(-0.2x_0 + 0.6x_1 + 0.9x_2)$$

update  
weights

# Perceptron Learning Algorithm: An Example



$$\hat{y} = h_{\mathbf{w}}(\mathbf{x}) = g\left(\sum_{j=0}^d w_j x_j\right)$$

$x_0 = 1$  (dummy variable)

Sign function

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

$$\begin{aligned} \hat{y} = h_{\mathbf{w}}(\mathbf{x}) &= g(w_0 x_0 + w_1 x_1 + w_2 x_2) \\ &= g(-0.2x_0 + 0.6x_1 + 0.9x_2) \end{aligned}$$

$$g(-0.2 \times 1 + 0.6 \times (-2) + 0.9 \times 2) = g(0.4)$$

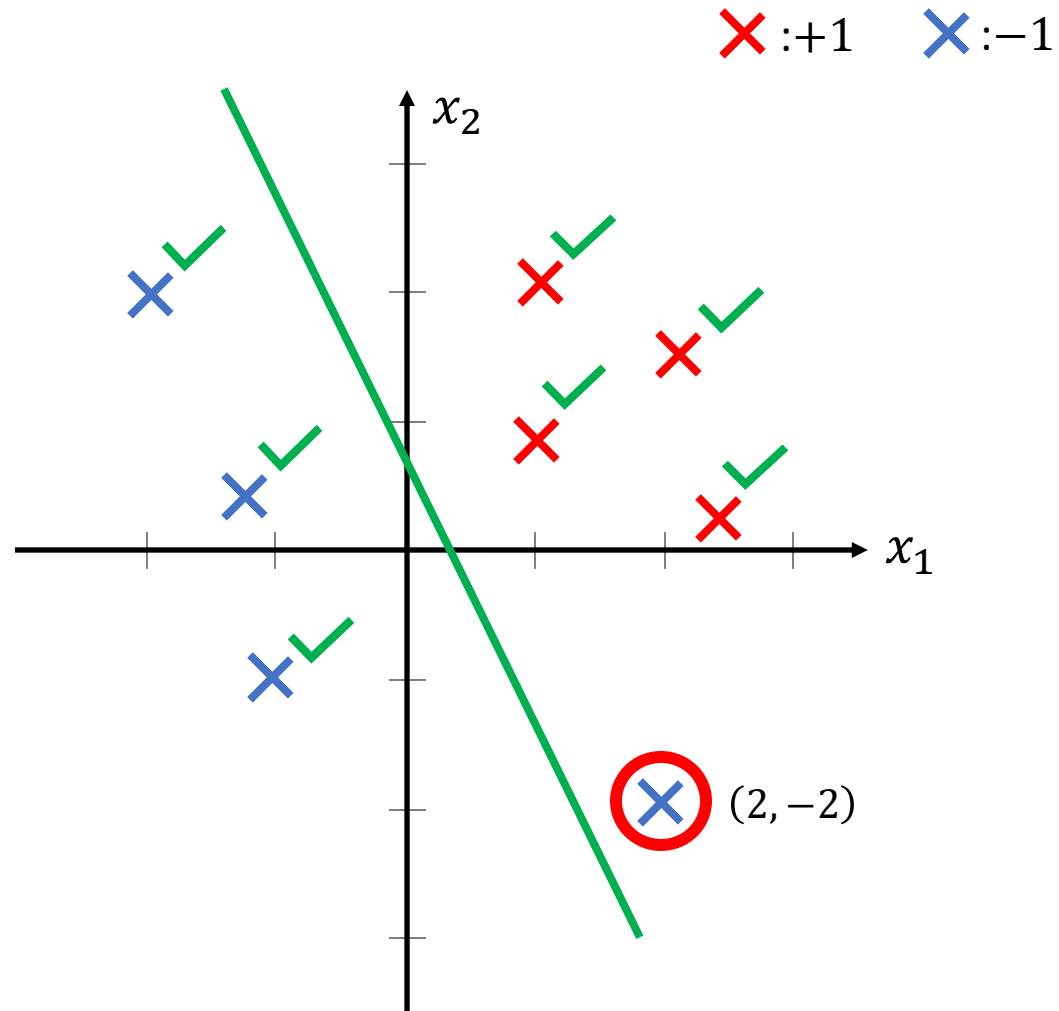
= +1

$$\mathbf{w} \leftarrow \mathbf{w} + \gamma(y^{(k)} - \hat{y}^{(k)}) \mathbf{x}^{(k)}$$

$$\begin{bmatrix} -0.2 \\ 0.6 \\ 0.9 \end{bmatrix} + 0.1(-1 - 1) \begin{bmatrix} 1 \\ -2 \\ 2 \end{bmatrix} = \begin{bmatrix} -0.4 \\ 1 \\ 0.5 \end{bmatrix}$$

$$\text{New } h_{\mathbf{w}}(\mathbf{x}) = g(-0.4x_0 + 1x_1 + 0.5x_2)$$

# Perceptron Learning Algorithm: An Example



$$\hat{y} = h_{\mathbf{w}}(\mathbf{x}) = g\left(\sum_{j=0}^d w_j x_j\right)$$

$x_0 = 1$  (dummy variable)

Sign function

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

$$\begin{aligned} \hat{y} = h_{\mathbf{w}}(\mathbf{x}) &= g(w_0 x_0 + w_1 x_1 + w_2 x_2) \\ &= g(-0.4x_0 + 1x_1 + 0.5x_2) \end{aligned}$$

$$g(-0.4 \times 1 + 1 \times 2 + 0.5 \times (-2)) = g(0.6)$$

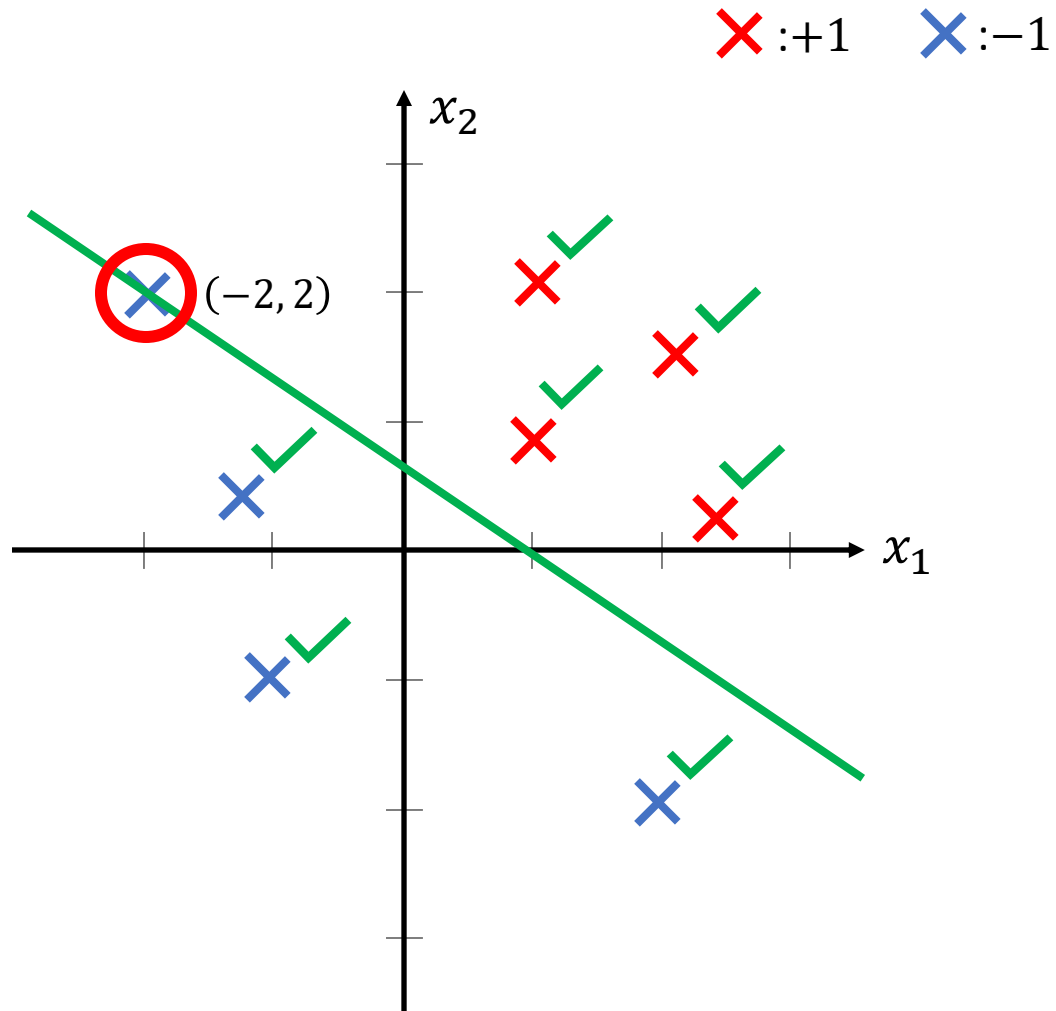
= +1

$$\mathbf{w} \leftarrow \mathbf{w} + \gamma(y^{(k)} - \hat{y}^{(k)}) \mathbf{x}^{(k)}$$

$$\begin{bmatrix} -0.4 \\ 1 \\ 0.5 \end{bmatrix} + 0.1(-1 - 1) \begin{bmatrix} 1 \\ 2 \\ -2 \end{bmatrix} = \begin{bmatrix} -0.6 \\ 0.6 \\ 0.9 \end{bmatrix}$$

New  $h_{\mathbf{w}}(\mathbf{x}) = g(-0.6x_0 + 0.6x_1 + 0.9x_2)$

# Perceptron Learning Algorithm: An Example



$$\hat{y} = h_{\mathbf{w}}(\mathbf{x}) = g\left(\sum_{j=0}^d w_j x_j\right)$$

$x_0 = 1$  (dummy variable)

Sign function

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

$$\begin{aligned} \hat{y} = h_{\mathbf{w}}(\mathbf{x}) &= g(w_0 x_0 + w_1 x_1 + w_2 x_2) \\ &= g(-0.6x_0 + 0.6x_1 + 0.9x_2) \end{aligned}$$

$$g(-0.6 \times 1 + 0.6 \times (-2) + 0.9 \times 2) = g(0)$$

$= +1$

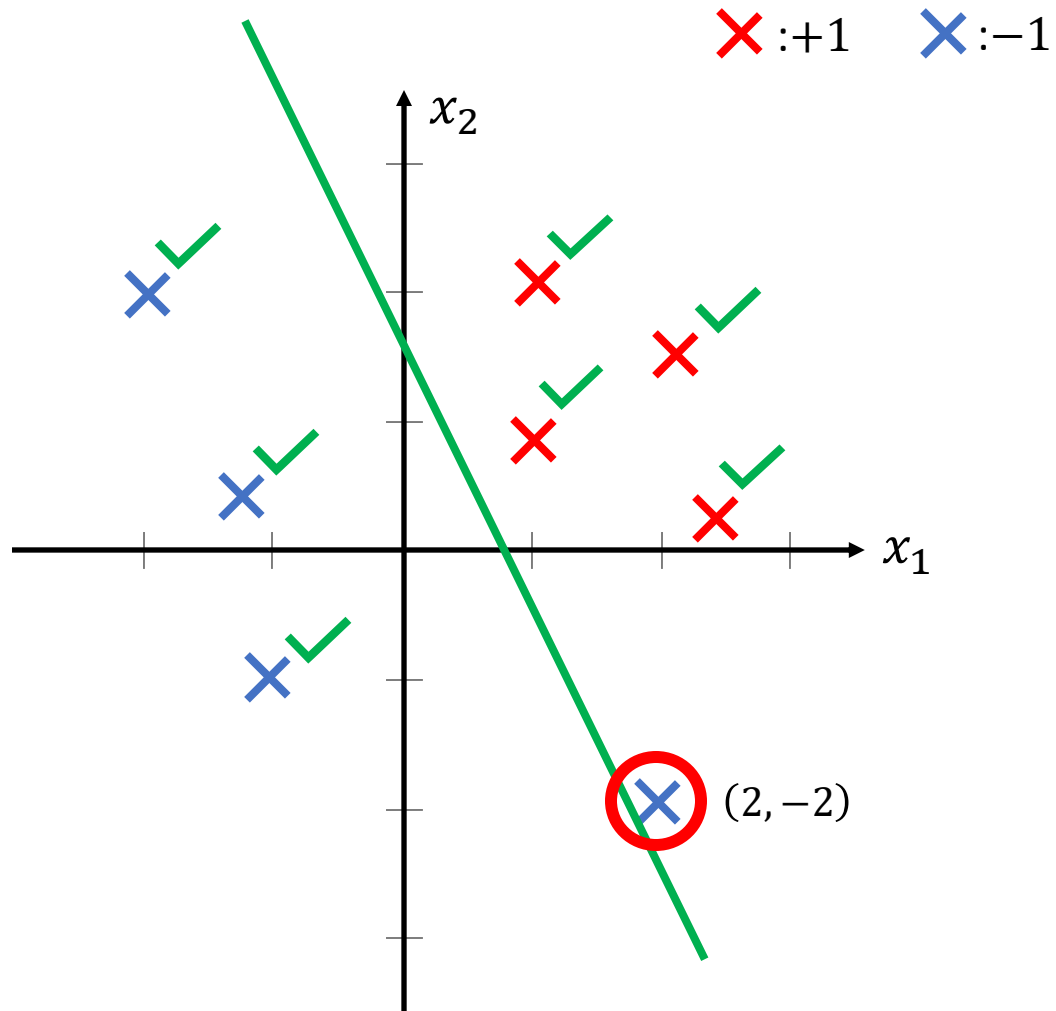
$$\mathbf{w} \leftarrow \mathbf{w} + \gamma(y^{(k)} - \hat{y}^{(k)}) \mathbf{x}^{(k)}$$

$$\begin{bmatrix} -0.6 \\ 0.6 \\ 0.9 \end{bmatrix} + 0.1(-1 - 1) \begin{bmatrix} 1 \\ -2 \\ 2 \end{bmatrix} = \begin{bmatrix} -0.8 \\ 1 \\ 0.5 \end{bmatrix}$$

New  $h_{\mathbf{w}}(\mathbf{x}) = g(-0.8x_0 + 1x_1 + 0.5x_2)$



# Perceptron Learning Algorithm: An Example



$$\hat{y} = h_{\mathbf{w}}(\mathbf{x}) = g\left(\sum_{j=0}^d w_j x_j\right)$$

$x_0 = 1$  (dummy variable)

Sign function

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

$$\begin{aligned} \hat{y} = h_{\mathbf{w}}(\mathbf{x}) &= g(w_0 x_0 + w_1 x_1 + w_2 x_2) \\ &= g(-0.8 x_0 + 1 x_1 + 0.5 x_2) \end{aligned}$$

$$g(-0.8 \times 1 + 1 \times 2 + 0.5 \times (-2)) = g(0.2)$$

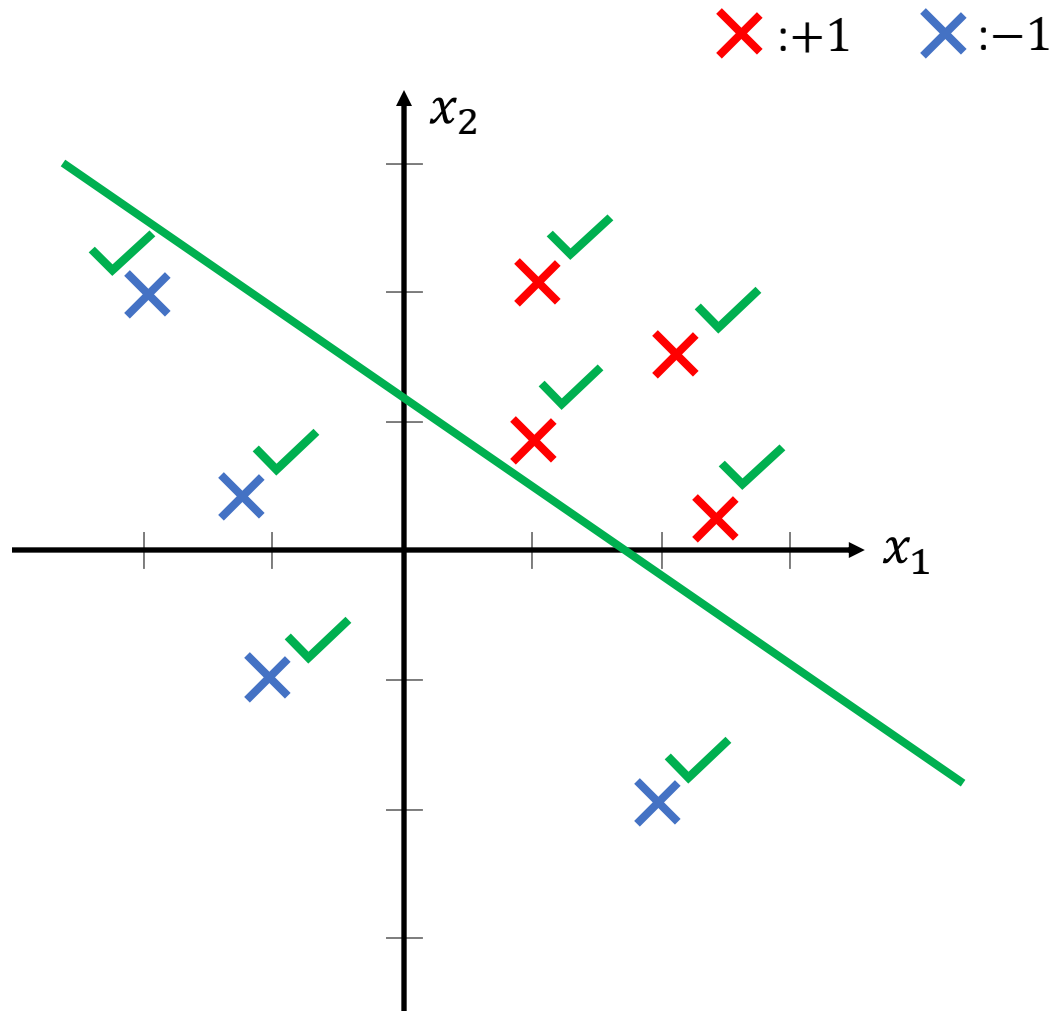
$= +1$

$$\mathbf{w} \leftarrow \mathbf{w} + \gamma(y^{(k)} - \hat{y}^{(k)}) \mathbf{x}^{(k)}$$

$$\begin{bmatrix} -0.8 \\ 1 \\ 0.5 \end{bmatrix} + 0.1(-1 - 1) \begin{bmatrix} 1 \\ 2 \\ -2 \end{bmatrix} = \begin{bmatrix} -1 \\ 0.6 \\ 0.9 \end{bmatrix}$$

$$\text{New } h_{\mathbf{w}}(\mathbf{x}) = g(-1 x_0 + 0.6 x_1 + 0.9 x_2)$$

# Perceptron Learning Algorithm: An Example



$$\hat{y} = h_{\mathbf{w}}(\mathbf{x}) = g\left(\sum_{j=0}^d w_j x_j\right)$$

$x_0 = 1$  (dummy variable)

Sign function

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

$$\begin{aligned}\hat{y} = h_{\mathbf{w}}(\mathbf{x}) &= g(w_0 x_0 + w_1 x_1 + w_2 x_2) \\ &= g(-1x_0 + 0.6x_1 + 0.9x_2)\end{aligned}$$

**No misclassifications! Converged!**

**Different from SVM, we do not maximize margin here.**

**What if it's not linearly separable?**

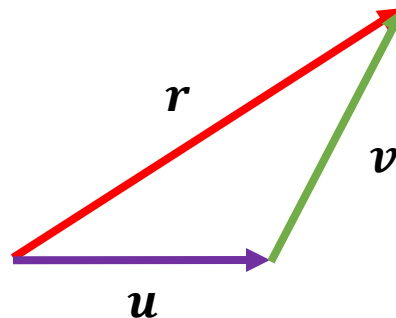
The algorithm will not converge

# Background: Vector Addition

Suppose that we have two vectors,  $\mathbf{u}$  and  $\mathbf{v}$ ,

$$\mathbf{r} = \mathbf{u} + \mathbf{v}$$

Vector addition can be visualized in a vector diagram by drawing the vectors to be added tip-to-tail:



# Background: Dot Product and Angle

Suppose that we have two vectors,  $\mathbf{u}$  and  $\mathbf{v}$ , dot product:

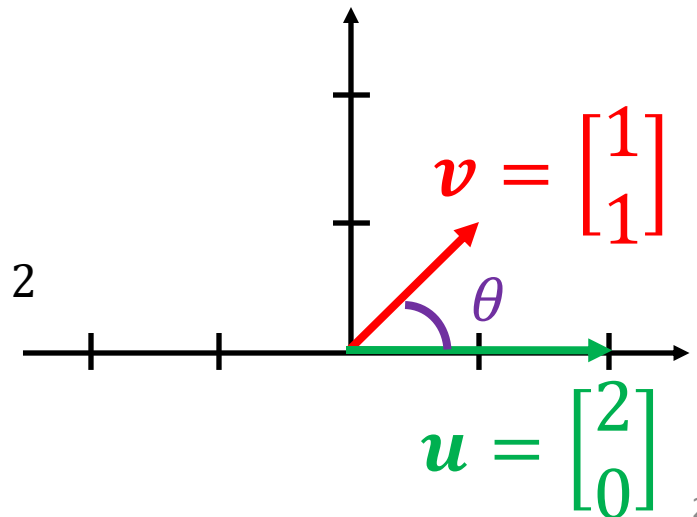
$$\mathbf{u} \cdot \mathbf{v} = \mathbf{u}^T \mathbf{v} = \|\mathbf{u}\| \|\mathbf{v}\| \cos \theta$$

where  $\theta$  is the angle between two vectors and  $0 \leq \theta \leq \pi$ .

If  $\mathbf{u} = \begin{bmatrix} 2 \\ 0 \end{bmatrix}$ ,  $\mathbf{v} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ , we have  $\|\mathbf{u}\| = 2$ ,  $\|\mathbf{v}\| = \sqrt{2}$ ,  $\theta = \frac{\pi}{4}$ ,  $\cos \theta = \frac{\sqrt{2}}{2}$

$$\mathbf{u} \cdot \mathbf{v} = \mathbf{u}^T \mathbf{v} = \begin{bmatrix} 2 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 2$$

$$\|\mathbf{u}\| \|\mathbf{v}\| \cos \theta = 2 \times \sqrt{2} \times \frac{\sqrt{2}}{2} = 2$$



# Perceptron Learning Algorithm: Why?

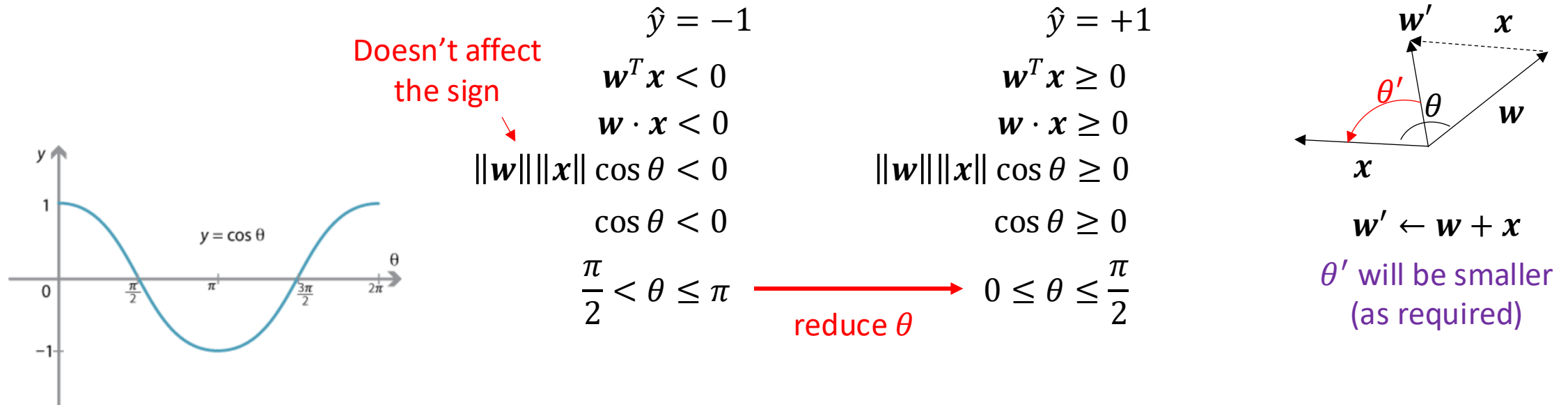
$$\hat{y} = g\left(\sum_{j=0}^d w_j x_j\right) = g(\mathbf{w}^T \mathbf{x}) = \begin{cases} +1, & \text{if } \mathbf{w}^T \mathbf{x} \geq 0 \\ -1, & \text{if } \mathbf{w}^T \mathbf{x} < 0 \end{cases}$$

When there is a misclassification:

- **Case 1:**  $y = +1, \hat{y} = -1$

What we have:

What we want:



# Perceptron Learning Algorithm: Why?

$$\hat{y} = g\left(\sum_{j=0}^d w_j x_j\right) = g(\mathbf{w}^T \mathbf{x}) = \begin{cases} +1, & \text{if } \mathbf{w}^T \mathbf{x} \geq 0 \\ -1, & \text{if } \mathbf{w}^T \mathbf{x} < 0 \end{cases}$$

When there is a misclassification:

- **Case 2:**  $y = -1, \hat{y} = +1$

What we have:

$$\hat{y} = +1$$

$$\mathbf{w}^T \mathbf{x} \geq 0$$

$$\mathbf{w} \cdot \mathbf{x} \geq 0$$

$$\|\mathbf{w}\| \|\mathbf{x}\| \cos \theta \geq 0$$

$$\cos \theta \geq 0$$

$$0 \leq \theta \leq \frac{\pi}{2}$$

What we want:

$$\hat{y} = -1$$

$$\mathbf{w}^T \mathbf{x} < 0$$

$$\mathbf{w} \cdot \mathbf{x} < 0$$

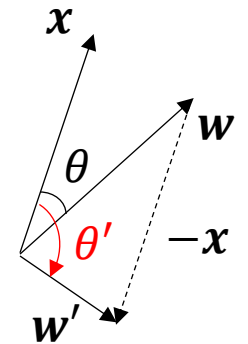
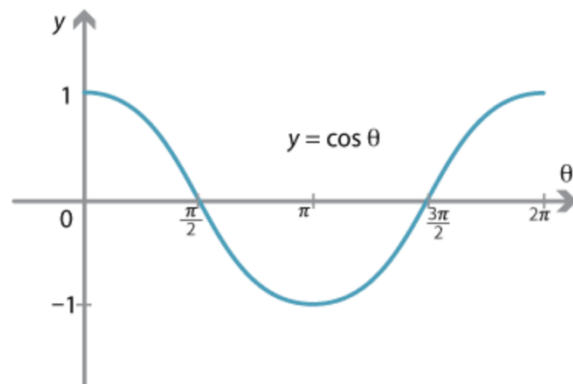
$$\|\mathbf{w}\| \|\mathbf{x}\| \cos \theta < 0$$

$$\cos \theta < 0$$

$$\frac{\pi}{2} < \theta \leq \pi$$

increase  $\theta$

Doesn't affect  
the sign



$$\mathbf{w}' \leftarrow \mathbf{w} - \mathbf{x}$$

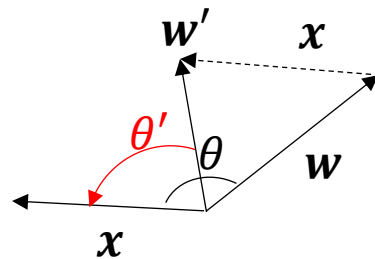
$\theta'$  will be larger  
(as required)

# Perceptron Learning Algorithm: Why?

$$\hat{y} = g\left(\sum_{j=0}^d w_j x_j\right) = g(\mathbf{w}^T \mathbf{x}) = \begin{cases} +1, & \text{if } \mathbf{w}^T \mathbf{x} \geq 0 \\ -1, & \text{if } \mathbf{w}^T \mathbf{x} < 0 \end{cases}$$

When there is a misclassification:

$$\begin{aligned} \hat{y} &= -1 \\ y &= +1 \end{aligned}$$



$$\mathbf{w}' \leftarrow \mathbf{w} + \mathbf{x}$$

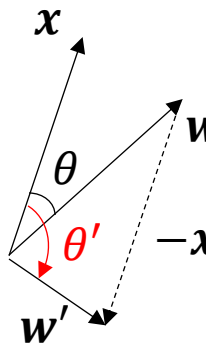
$\theta'$  will be smaller  
(as required)

$$\mathbf{w}' \leftarrow \mathbf{w} + 2\gamma \mathbf{x} \quad \gamma \text{ is constant}$$

$$\mathbf{w} + \gamma(+1 - (-1))\mathbf{x}$$

$$\mathbf{w} + \gamma(y - \hat{y})\mathbf{x}$$

Weight updating rule in PLA



$$\begin{aligned} \hat{y} &= +1 \\ y &= -1 \end{aligned}$$

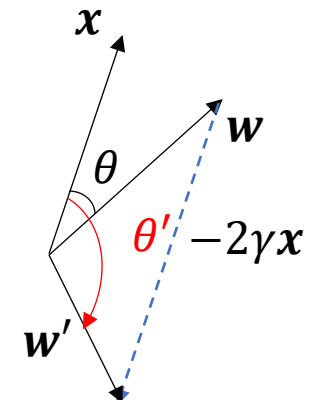
$$\mathbf{w}' \leftarrow \mathbf{w} - \mathbf{x}$$

$\theta'$  will be larger  
(as required)

$$\mathbf{w}' \leftarrow \mathbf{w} - 2\gamma \mathbf{x}$$

$$\mathbf{w} + \gamma(-1 - (+1))\mathbf{x}$$

$$\mathbf{w} + \gamma(y - \hat{y})\mathbf{x}$$



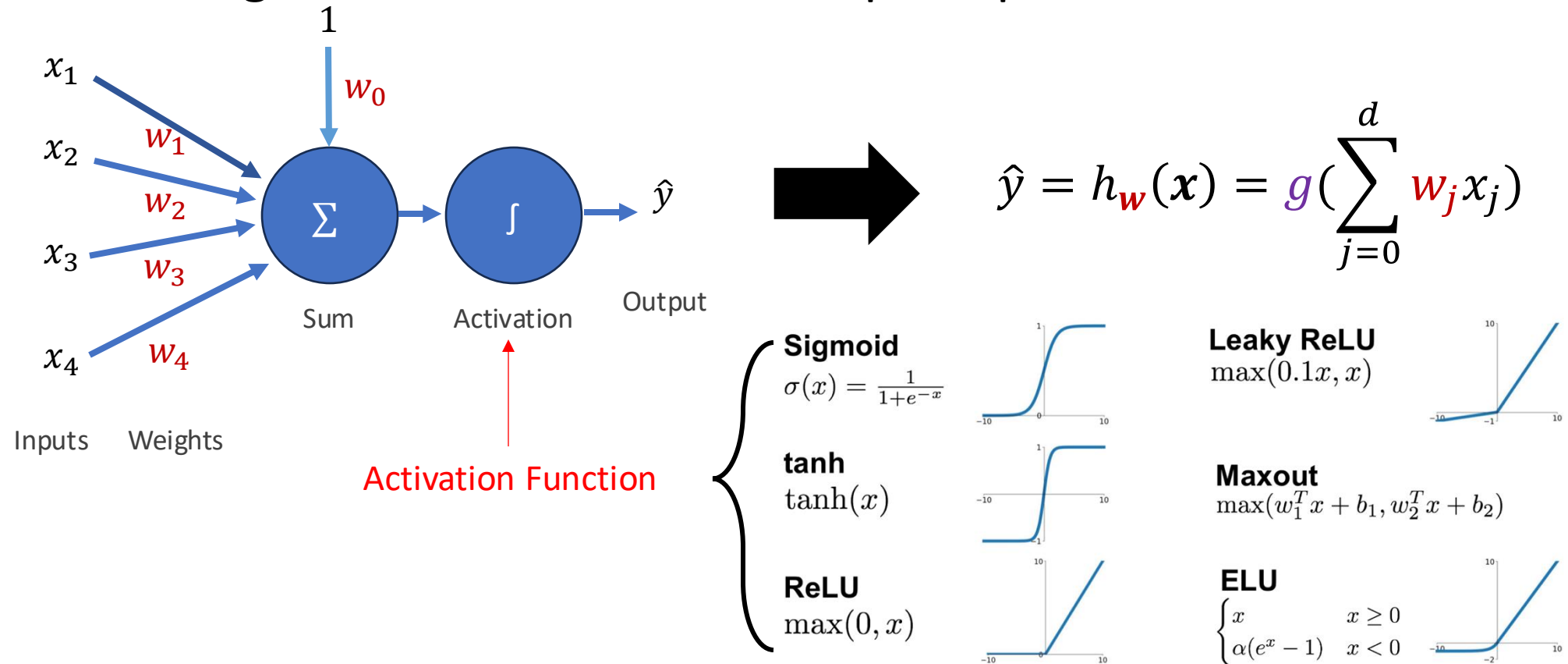
# Outline

- Perceptron
  - Biological inspiration
  - Perceptron Learning Algorithm
- **Neural Network**
  - Neuron
  - AND Gate Modelling
  - XNOR Gate Modelling
  - Single-layer and Multi-layer Neural Networks
- Multi-class Classification



# Neuron

- Neuron is the building block of neural networks. An artificial neuron is a more generalized version of the perceptron.



# Linear Regression Model

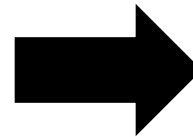
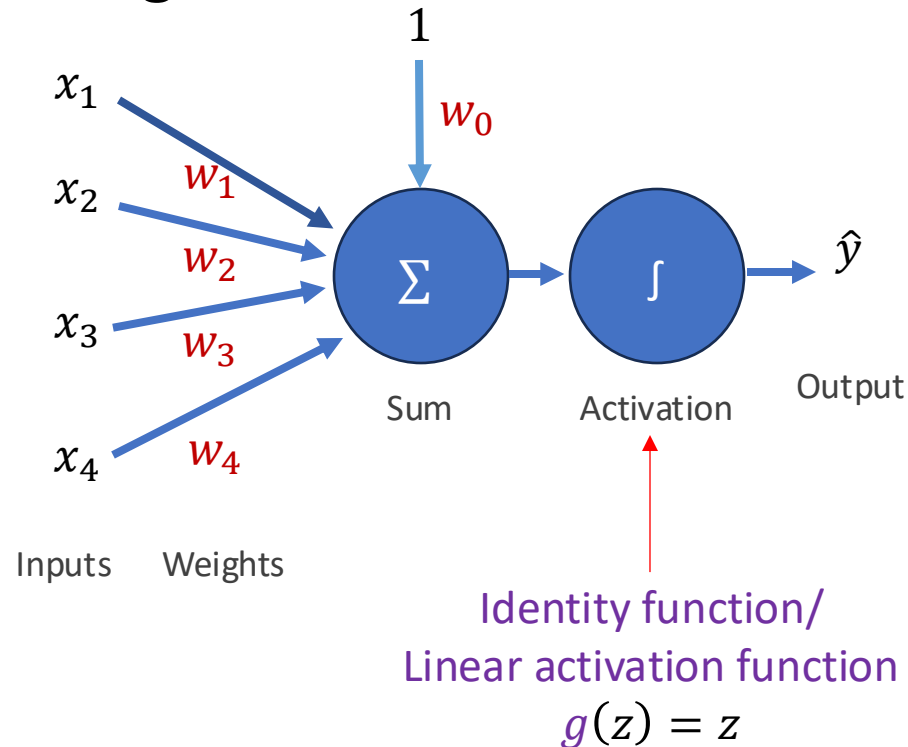
Given an input vector  $x$  of dimension  $d$ , the hypothesis class of linear models is defined as the set of functions:

$$h_{\mathbf{w}}(x) = \sum_{j=0}^d \mathbf{w}_j x_j$$

Where  $\mathbf{w}_0, \dots, \mathbf{w}_d$  are **parameters/weights**.

# Neuron v.s. Linear Regression Model

- A neuron with linear activation function is equivalent to linear regression model.



$$\hat{y} = h_{\mathbf{w}}(\mathbf{x}) = g\left(\sum_{j=0}^d w_j x_j\right) = \sum_{j=0}^d w_j x_j$$

# Logistic Regression Model

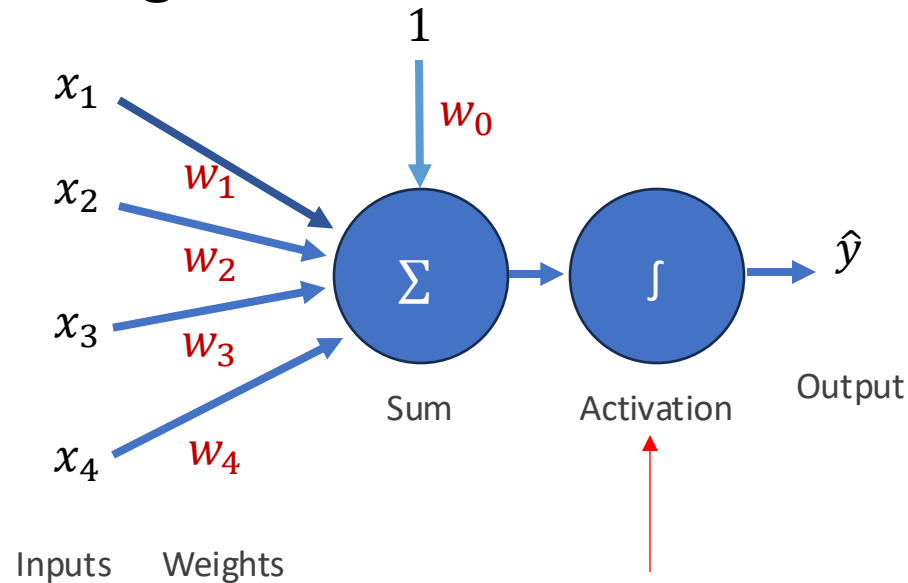
Given an input vector  $x$  of dimension  $d$ , the hypothesis class of linear models is defined as the set of functions:

$$h_{\mathbf{w}}(x) = \sigma\left(\sum_{j=0}^d w_j x_j\right)$$

Where  $w_0, \dots, w_d$  are **parameters/weights** and  $\sigma$  is the sigmoid function.

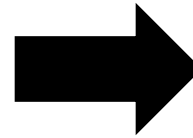
# Neuron v.s. Logistic Regression Model

- A neuron with sigmoid activation function is equivalent to logistic regression model.



Sigmoid function

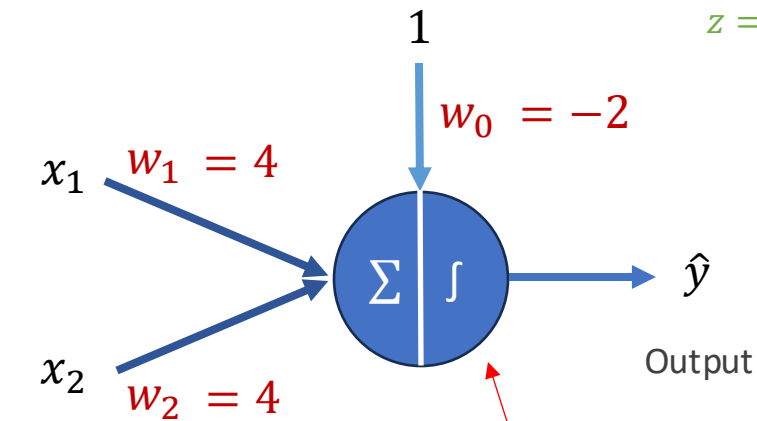
$$g(z) = \frac{1}{1 + e^{-z}}$$



$$\hat{y} = h_{\mathbf{w}}(\mathbf{x}) = g\left(\sum_{j=0}^d w_j x_j\right) = \sigma\left(\sum_{j=0}^d w_j x_j\right)$$

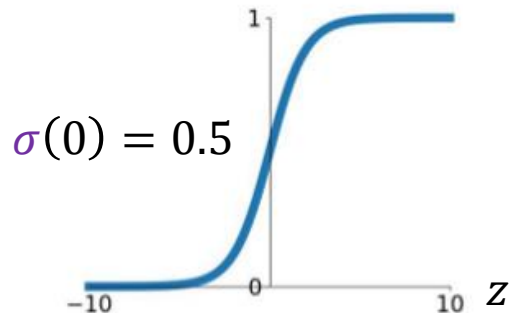
# OR Gate Modelling

With Logistic Regression Model



Sigmoid function

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



$$\hat{y} = h_{\mathbf{w}}(\mathbf{x}) = \sigma\left(\sum_{j=0}^d w_j x_j\right) \quad x_0 = 1 \text{ (dummy variable)}$$

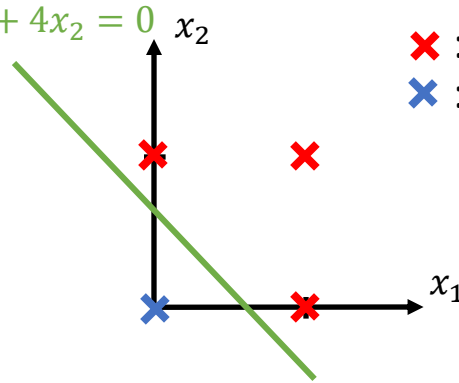
$$\text{Decision threshold: } \hat{y} = \sigma\left(\sum_{j=0}^d w_j x_j\right) = 0.5 = \sigma(0)$$

$$\text{Decision boundary: } z = \sum_{j=0}^d w_j x_j = 0$$

Consider  $x_1, x_2 \in \{1, 0\}$

$x_0$	$x_1$	$x_2$	$y$
1	0	0	0
1	0	1	1
1	1	0	1
1	1	1	1

$\times$  : 1  
 $\times$  : 0



$$\hat{y} = \sigma(z), z = w_0 x_0 + w_1 x_1 + w_2 x_2 \quad \begin{array}{l} z \geq 0, \text{ predicted as 1} \\ z < 0, \text{ predicted as 0} \end{array}$$

$$\mathbf{x}^{(1)} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, y^{(1)} = 0: \quad \mathbf{x}^{(2)} = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}, y^{(2)} = 1:$$

$$z^{(1)} = w_0 < 0$$

$$z^{(2)} = w_0 + w_2 \geq 0$$

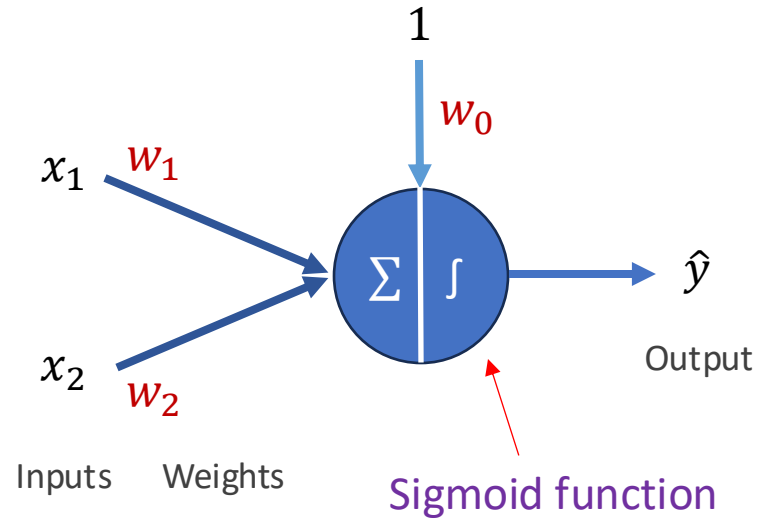
$$\mathbf{x}^{(3)} = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}, y^{(3)} = 1: \quad \mathbf{x}^{(4)} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}, y^{(4)} = 1$$

$$z^{(3)} = w_0 + w_1 \geq 0$$

$$z^{(4)} = w_0 + w_1 + w_2 \geq 0$$

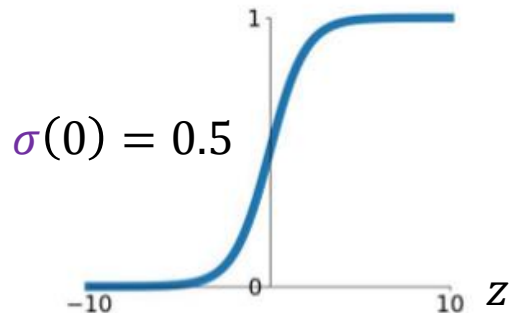
# AND Gate Modelling

With Logistic Regression Model



Sigmoid function

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



$$\hat{y} = h_{\mathbf{w}}(\mathbf{x}) = \sigma\left(\sum_{j=0}^d \mathbf{w}_j x_j\right) \quad x_0 = 1 \text{ (dummy variable)}$$

Decision threshold:  $\hat{y} = \sigma\left(\sum_{j=0}^d \mathbf{w}_j x_j\right) = 0.5 = \sigma(0)$

Decision boundary:  $z = \sum_{j=0}^d \mathbf{w}_j x_j = 0$

Consider  $x_1, x_2 \in \{1, 0\}$

$x_0$	$x_1$	$x_2$	$y$
1	0	0	0
1	0	1	0
1	1	0	0
1	1	1	1

$\hat{y} = \sigma(z), z = \mathbf{w}_0 x_0 + \mathbf{w}_1 x_1 + \mathbf{w}_2 x_2$   
 $z \geq 0$ , predicted as 1  
 $z < 0$ , predicted as 0

# Poll Everywhere

To correctly model the AND gate, which of the following options can be used to set the parameters in the logistic regression model?

- A:  $w_0 = -6, w_1 = 4, w_2 = 4$
- B:  $w_0 = -5, w_1 = 2, w_2 = 2$
- C:  $w_0 = 3, w_1 = -2, w_2 = -2$
- D:  $w_0 = -4, w_1 = -2, w_2 = -2$

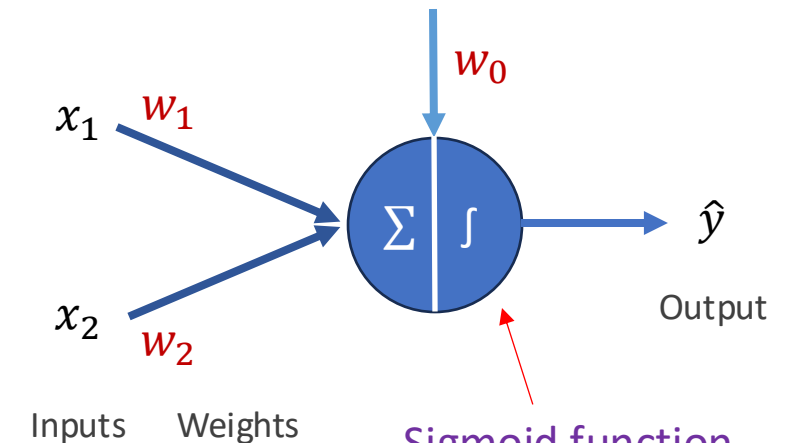
Consider  $x_1, x_2 \in \{1,0\}$

$x_0$	$x_1$	$x_2$	$y$
1	0	0	0
1	0	1	0
1	1	0	0
1	1	1	1

$$\hat{y} = \sigma(z), z = w_0 x_0 + w_1 x_1 + w_2 x_2$$

$z \geq 0$ , predicted as 1

$z < 0$ , predicted as 0



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



# Poll Everywhere

To correctly model the AND gate, which of the following options can be used to set the parameters in the logistic regression model?

**A:**  $w_0 = -6, w_1 = 4, w_2 = 4$

**B:**  $w_0 = -5, w_1 = 2, w_2 = 2$

**C:**  $w_0 = 3, w_1 = -2, w_2 = -2$

**D:**  $w_0 = -4, w_1 = -2, w_2 = -2$

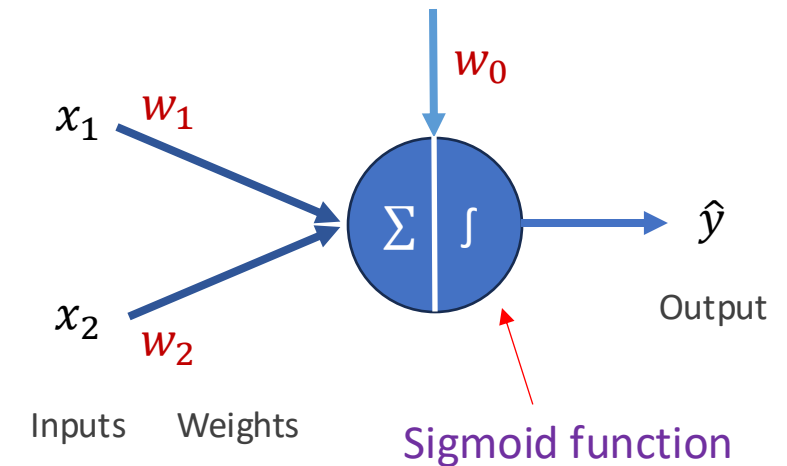
Consider  $x_1, x_2 \in \{1,0\}$

$x_0$	$x_1$	$x_2$	$y$
1	0	0	0
1	0	1	0
1	1	0	0
1	1	1	1

$$\hat{y} = \sigma(z), z = w_0 x_0 + w_1 x_1 + w_2 x_2$$

$z \geq 0$ , predicted as 1

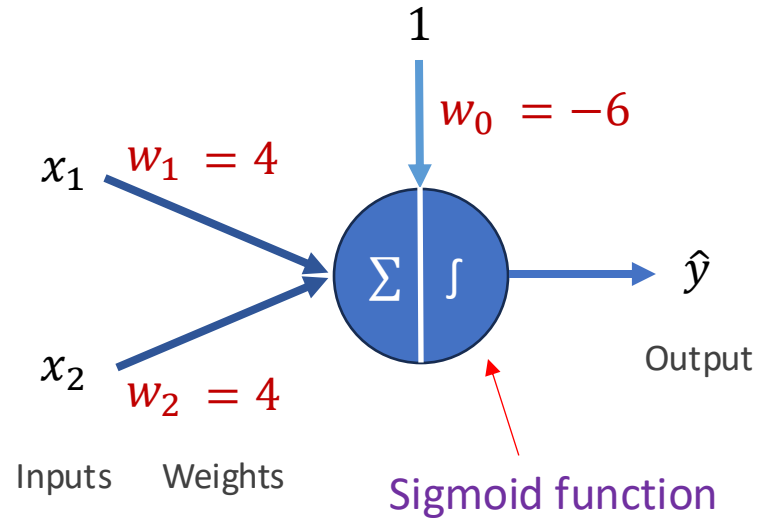
$z < 0$ , predicted as 0



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

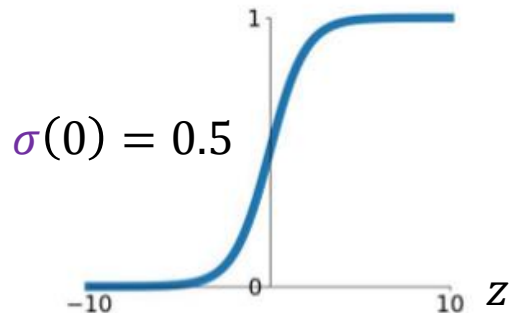
# AND Gate Modelling

With Logistic Regression Model



Sigmoid function

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

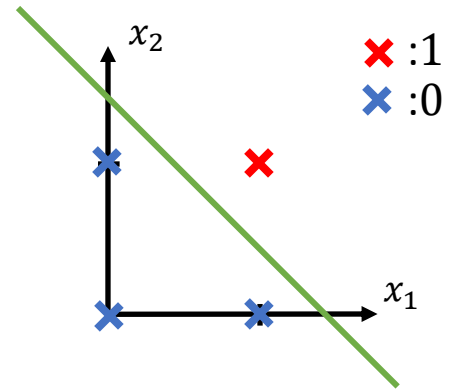


$$\hat{y} = h_{\mathbf{w}}(\mathbf{x}) = \sigma\left(\sum_{j=0}^d \mathbf{w}_j x_j\right) \quad x_0 = 1 \text{ (dummy variable)}$$

$$\text{Decision threshold: } \hat{y} = \sigma\left(\sum_{j=0}^d \mathbf{w}_j x_j\right) = 0.5 = \sigma(0)$$

$$\text{Decision boundary: } z = \sum_{j=0}^d \mathbf{w}_j x_j = 0$$

$$z = -6 + 4x_1 + 4x_2 = 0$$



Consider  $x_1, x_2 \in \{1, 0\}$

$x_0$	$x_1$	$x_2$	$y$
1	0	0	0
1	0	1	0
1	1	0	0
1	1	1	1

$$\hat{y} = \sigma(z), \quad z = \mathbf{w}_0 x_0 + \mathbf{w}_1 x_1 + \mathbf{w}_2 x_2$$

$z \geq 0$ , predicted as 1  
 $z < 0$ , predicted as 0

$$\mathbf{x}^{(1)} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, y^{(1)} = 0: \quad \mathbf{x}^{(2)} = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}, y^{(2)} = 0:$$

$$z^{(1)} = \mathbf{w}_0 < 0$$

$$z^{(2)} = \mathbf{w}_0 + \mathbf{w}_2 < 0$$

$$\mathbf{x}^{(3)} = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}, y^{(3)} = 0$$

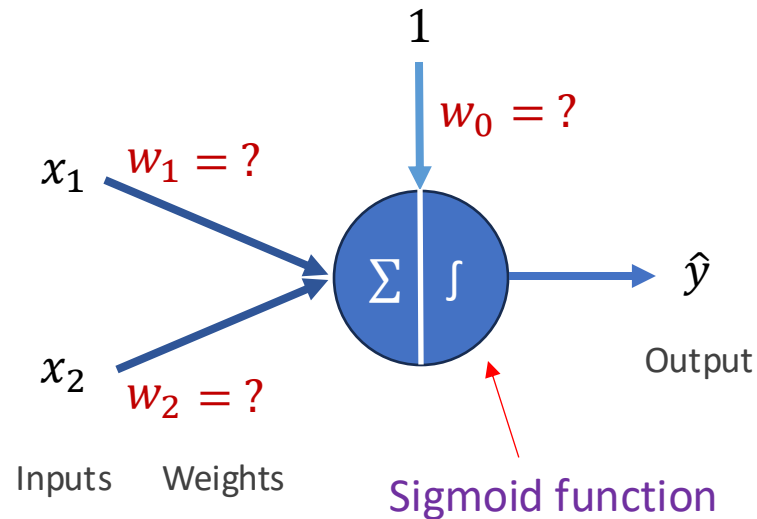
$$z^{(3)} = \mathbf{w}_0 + \mathbf{w}_1 < 0$$

$$\mathbf{x}^{(4)} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}, y^{(4)} = 1$$

$$z^{(4)} = \mathbf{w}_0 + \mathbf{w}_1 + \mathbf{w}_2 \geq 0$$

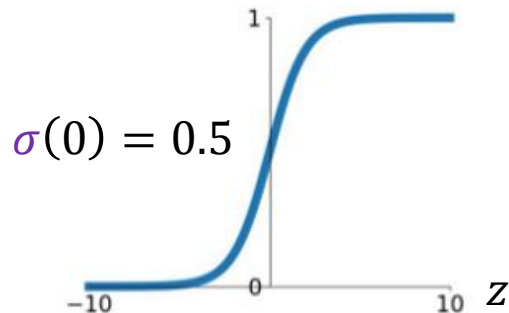
# XNOR Gate Modelling

With Logistic Regression Model



Sigmoid function

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



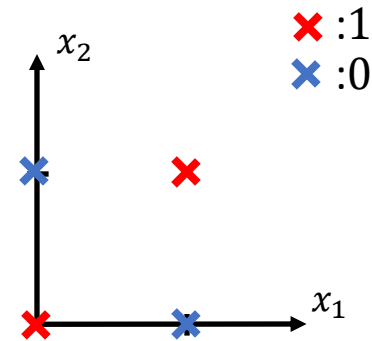
$$\hat{y} = h_{\mathbf{w}}(\mathbf{x}) = \sigma\left(\sum_{j=0}^d \mathbf{w}_j x_j\right) \quad x_0 = 1 \text{ (dummy variable)}$$

$$\text{Decision threshold: } \hat{y} = \sigma\left(\sum_{j=0}^d \mathbf{w}_j x_j\right) = 0.5 = \sigma(0)$$

$$\text{Decision boundary: } z = \sum_{j=0}^d \mathbf{w}_j x_j = 0$$

Consider  $x_1, x_2 \in \{1, 0\}$

$x_0$	$x_1$	$x_2$	$y$
1	0	0	1
1	0	1	0
1	1	0	0
1	1	1	1



$$\hat{y} = \sigma(z), z = \mathbf{w}_0 x_0 + \mathbf{w}_1 x_1 + \mathbf{w}_2 x_2$$

Not linearly separable

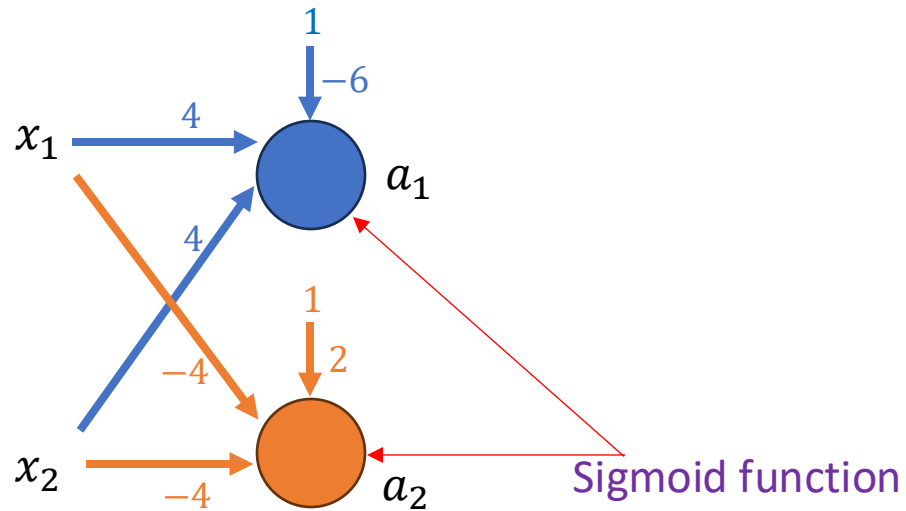
Feature engineering is needed to transform existing features:

- Manually create new features, e.g.,  $x_1^2$ ,  $e^{x_2}$ ,  $x_1 x_2$  ...
- Use neurons to transform features

# XNOR Gate Modelling

With More Neurons

$a_1$  and  $a_2$  are new features created based on existing features  $x_1$  and  $x_2$ .



Decision threshold:  $\hat{y} = 0.5 = \sigma(0)$

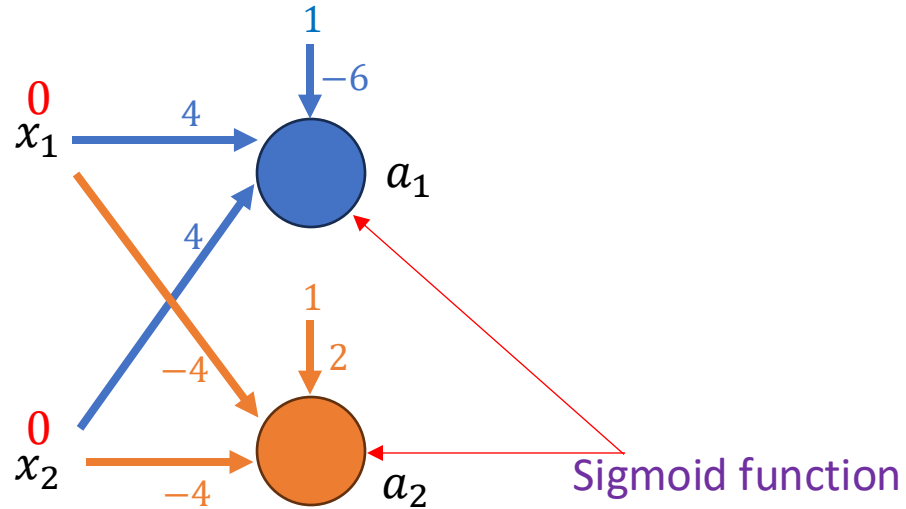
$$a_1 = \sigma(-6 + 4x_1 + 4x_2)$$

$$a_2 = \sigma(2 + (-4)x_1 + (-4)x_2)$$

# XNOR Gate Modelling

With More Neurons

$a_1$  and  $a_2$  are new features created based on existing features  $x_1$  and  $x_2$ .



Decision threshold:  $\hat{y} = 0.5 = \sigma(0)$

$$\begin{aligned} a_1 &= \sigma(-6 + 4x_1 + 4x_2) \\ &= \sigma(-6 + 4 \times 0 + 4 \times 0) = 0.002 \end{aligned}$$

$$\begin{aligned} a_2 &= \sigma(2 + (-4)x_1 + (-4)x_2) \\ &= \sigma(2 + (-4) \times 0 + (-4) \times 0) = 0.881 \end{aligned}$$

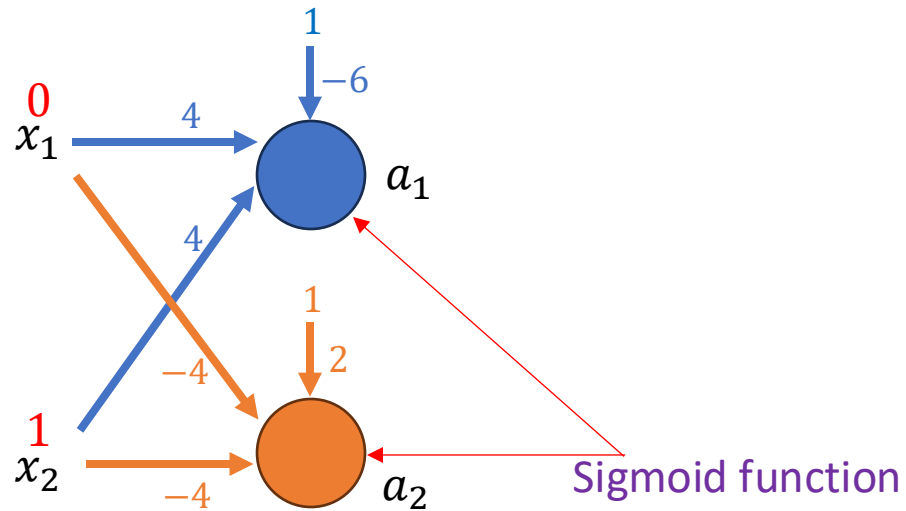
Consider  $x_1, x_2 \in \{1, 0\}$

$x_1$	$x_2$	$y$	$a_1$	$a_2$
0	0	1	0.002	0.881
0	1	0		
1	0	0		
1	1	1		

# XNOR Gate Modelling

With More Neurons

$a_1$  and  $a_2$  are new features created based on existing features  $x_1$  and  $x_2$ .



Decision threshold:  $\hat{y} = 0.5 = \sigma(0)$

$$\begin{aligned} a_1 &= \sigma(-6 + 4x_1 + 4x_2) \\ &= \sigma(-6 + 4 \times 0 + 4 \times 1) = 0.119 \end{aligned}$$

$$\begin{aligned} a_2 &= \sigma(2 + (-4)x_1 + (-4)x_2) \\ &= \sigma(2 + (-4) \times 0 + (-4) \times 1) = 0.119 \end{aligned}$$

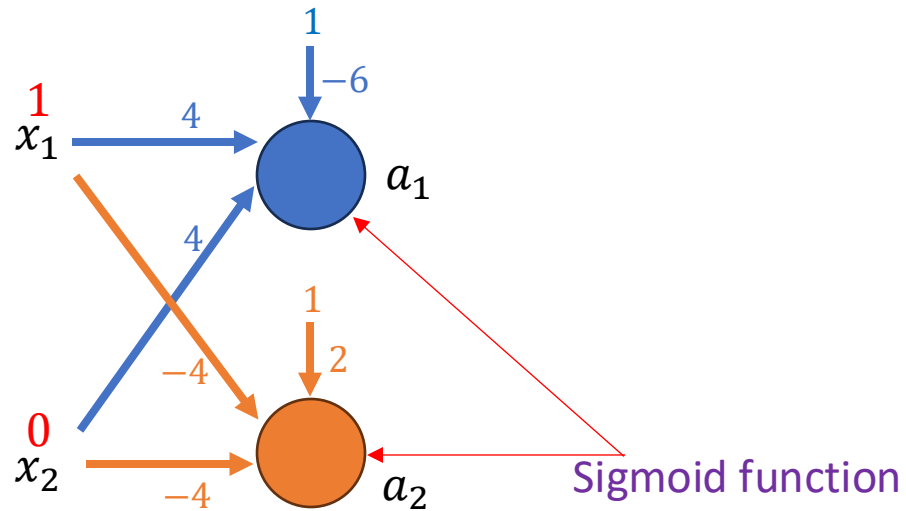
Consider  $x_1, x_2 \in \{1, 0\}$

$x_1$	$x_2$	$y$	$a_1$	$a_2$
0	0	1	0.002	0.881
0	1	0	0.119	0.119
1	0	0		
1	1	1		

# XNOR Gate Modelling

With More Neurons

$a_1$  and  $a_2$  are new features created based on existing features  $x_1$  and  $x_2$ .



Decision threshold:  $\hat{y} = 0.5 = \sigma(0)$

$$\begin{aligned} a_1 &= \sigma(-6 + 4x_1 + 4x_2) \\ &= \sigma(-6 + 4 \times 1 + 4 \times 0) = 0.119 \end{aligned}$$

$$\begin{aligned} a_2 &= \sigma(2 + (-4)x_1 + (-4)x_2) \\ &= \sigma(2 + (-4) \times 1 + (-4) \times 0) = 0.119 \end{aligned}$$

Consider  $x_1, x_2 \in \{1, 0\}$

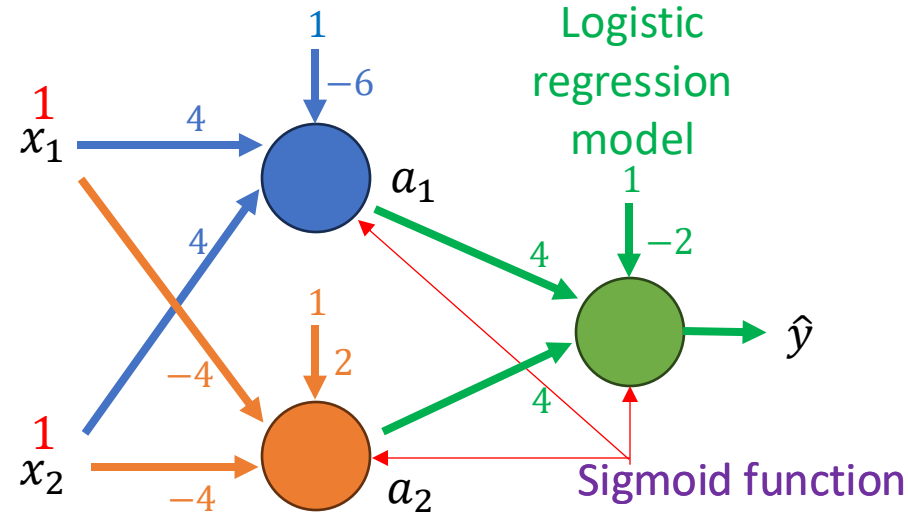
$x_1$	$x_2$	$y$	$a_1$	$a_2$
0	0	1	0.002	0.881
0	1	0	0.119	0.119
1	0	0	0.119	0.119
1	1	1		

Decision threshold:  $\hat{y} = 0.5 = \sigma(0)$

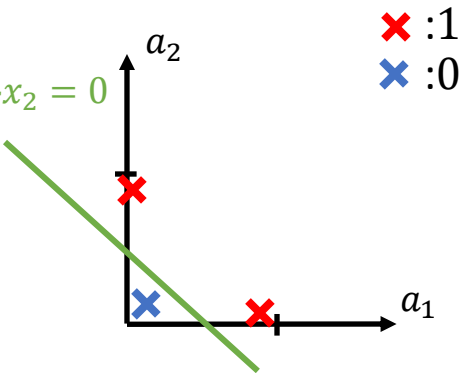
# XNOR Gate Modelling

With More Neurons

$a_1$  and  $a_2$  are new features created based on existing features  $x_1$  and  $x_2$ .



$$z = -2 + 4x_1 + 4x_2 = 0$$



✖ :1  
✖ :0

$$\begin{aligned} a_1 &= \sigma(-6 + 4x_1 + 4x_2) \\ &= \sigma(-6 + 4 \times 1 + 4 \times 1) = 0.881 \\ a_2 &= \sigma(2 + (-4)x_1 + (-4)x_2) \\ &= \sigma(2 + (-4) \times 1 + (-4) \times 1) = 0.002 \end{aligned}$$

Consider  $x_1, x_2 \in \{1,0\}$

$x_1$	$x_2$	$y$	$a_1$	$a_2$
0	0	1	0.002	0.881
0	1	0	0.119	0.119
1	0	0	0.119	0.119
1	1	1	0.881	0.002

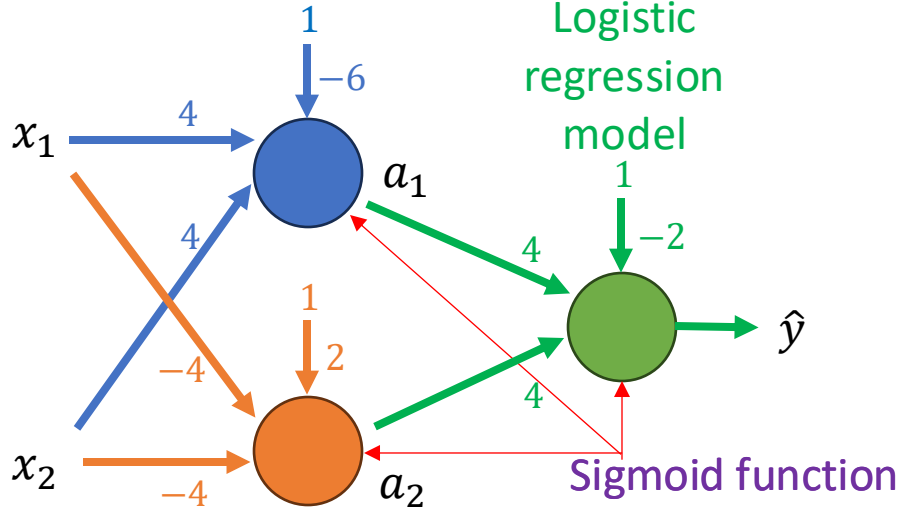


Decision threshold:  $\hat{y} = 0.5 = \sigma(0)$

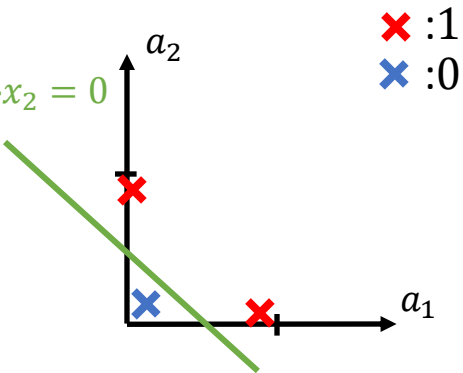
# XNOR Gate Modelling

With More Neurons

$a_1$  and  $a_2$  are new features created based on existing features  $x_1$  and  $x_2$ .



Multi-layer neural network



$$a_1 = \sigma(-6 + 4x_1 + 4x_2)$$

$$a_2 = \sigma(2 + (-4)x_1 + (-4)x_2)$$

$$\hat{y} = \sigma(-2 + 4a_1 + 4a_2)$$

Consider  $x_1, x_2 \in \{1,0\}$

$x_1$	$x_2$	$y$	$a_1$	$a_2$	$\hat{y}$	Pred. label
0	0	1	0.002	0.881	1.532	1
0	1	0	0.119	0.119	-1.048	0
1	0	0	0.119	0.119	-1.048	0
1	1	1	0.881	0.002	1.532	1

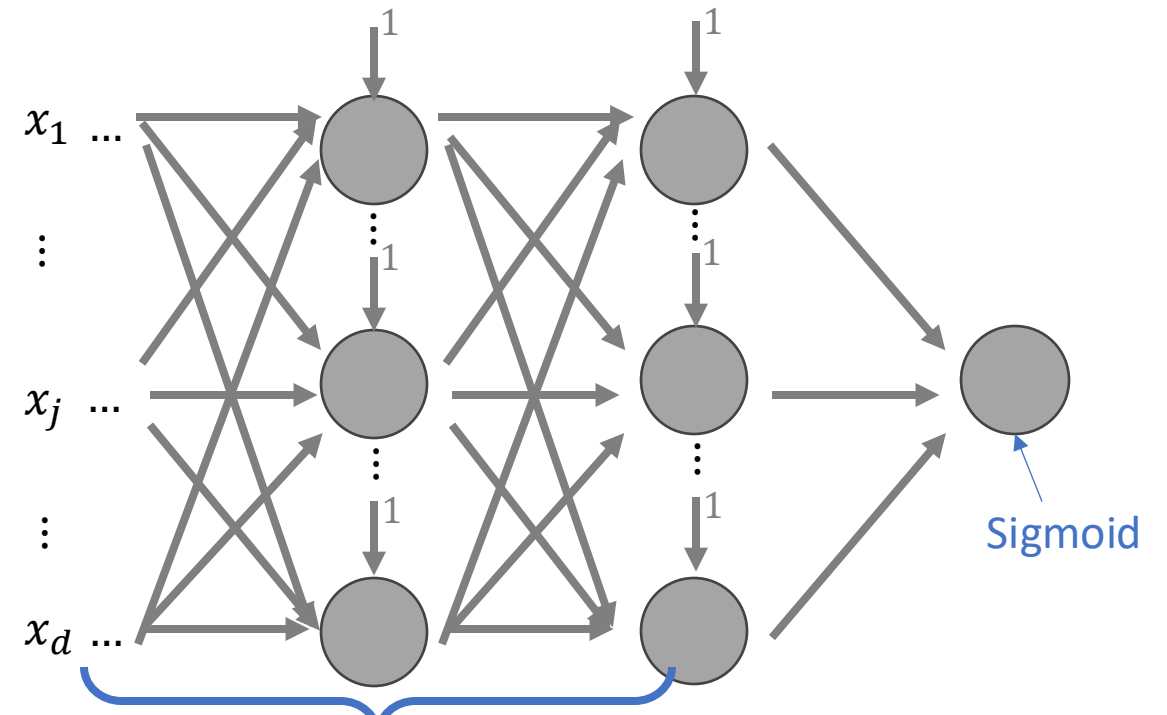
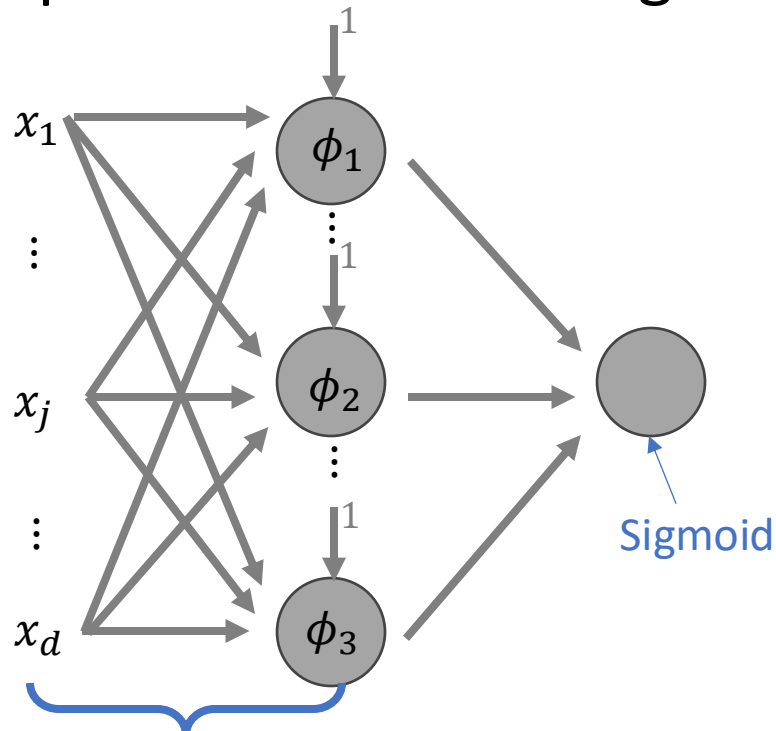


How to learn the weights for neurons will be covered in the next lecture.

# Neural Network v.s. Logistic Regression Model

With Feature Engineering

Logistic regression relies on **manual feature engineering** to capture complex patterns, while a multi-layer neural network **learns** its own feature representations through its layers with activation functions.



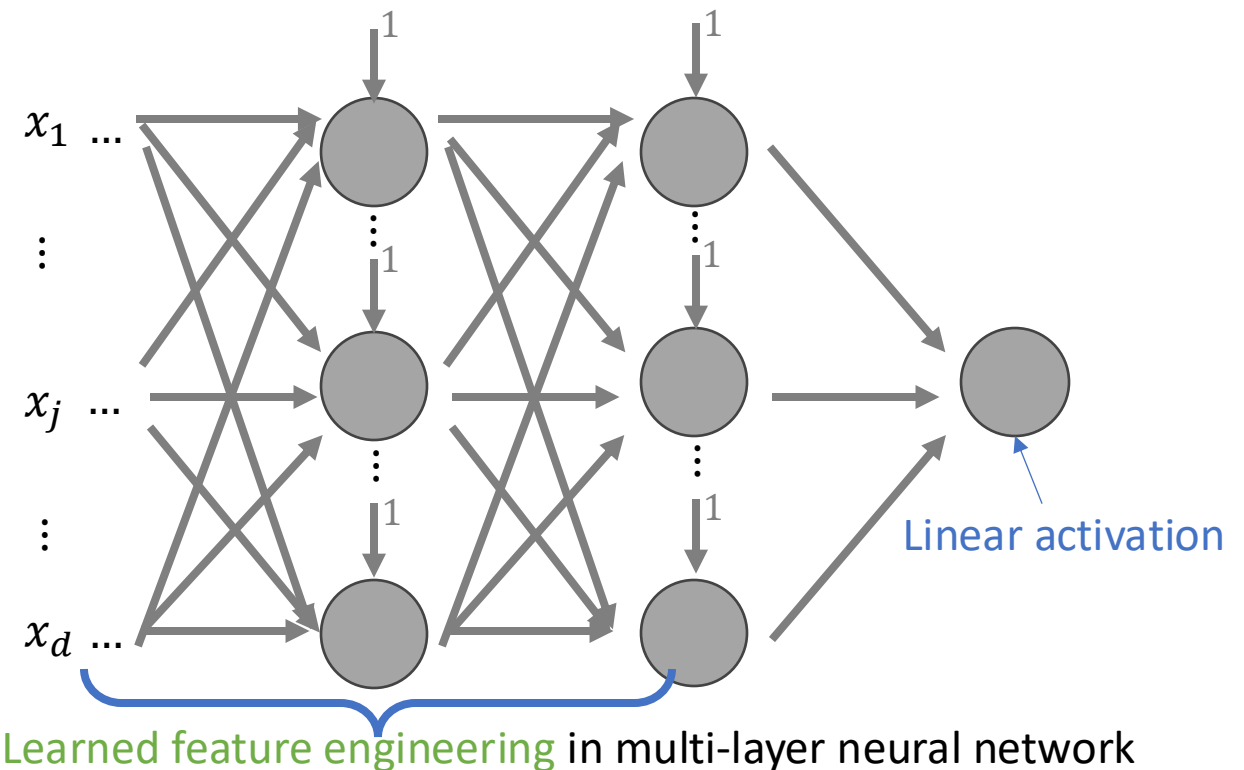
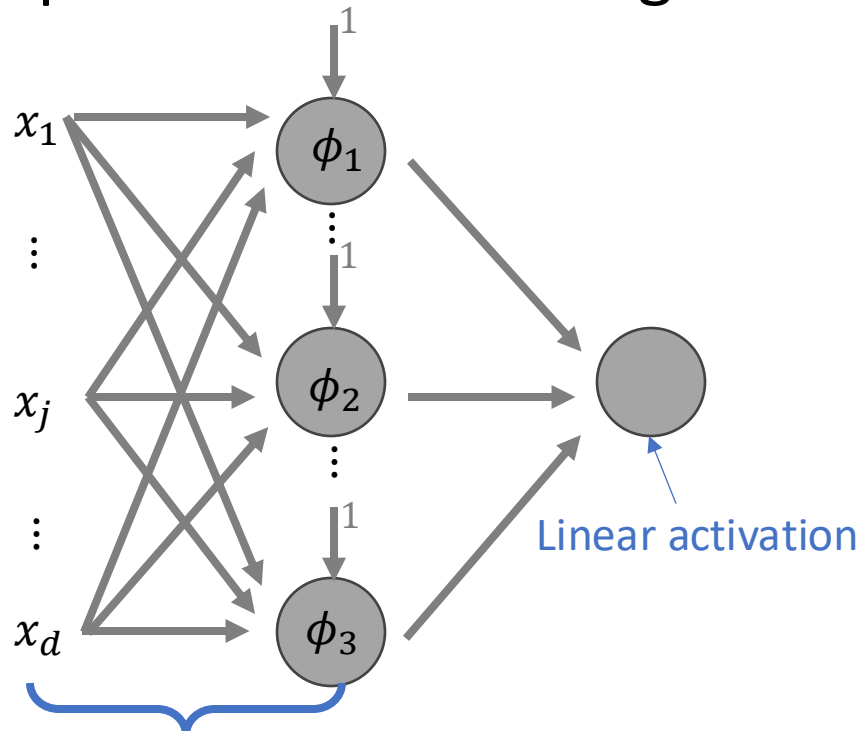
Manual feature engineering in logistic regression

Learned feature engineering in multi-layer neural network

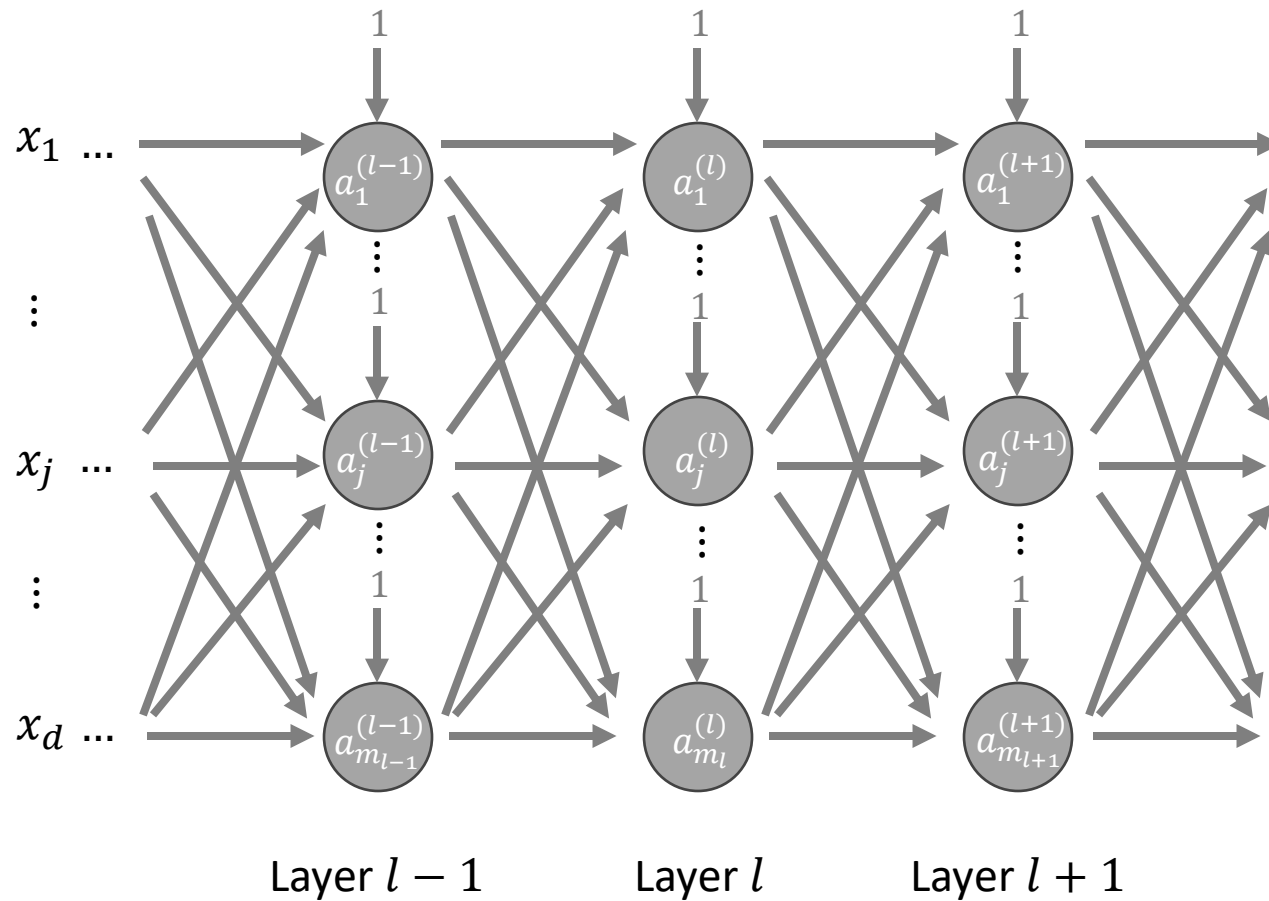
# Neural Network v.s. Linear Regression Model

With Feature Engineering

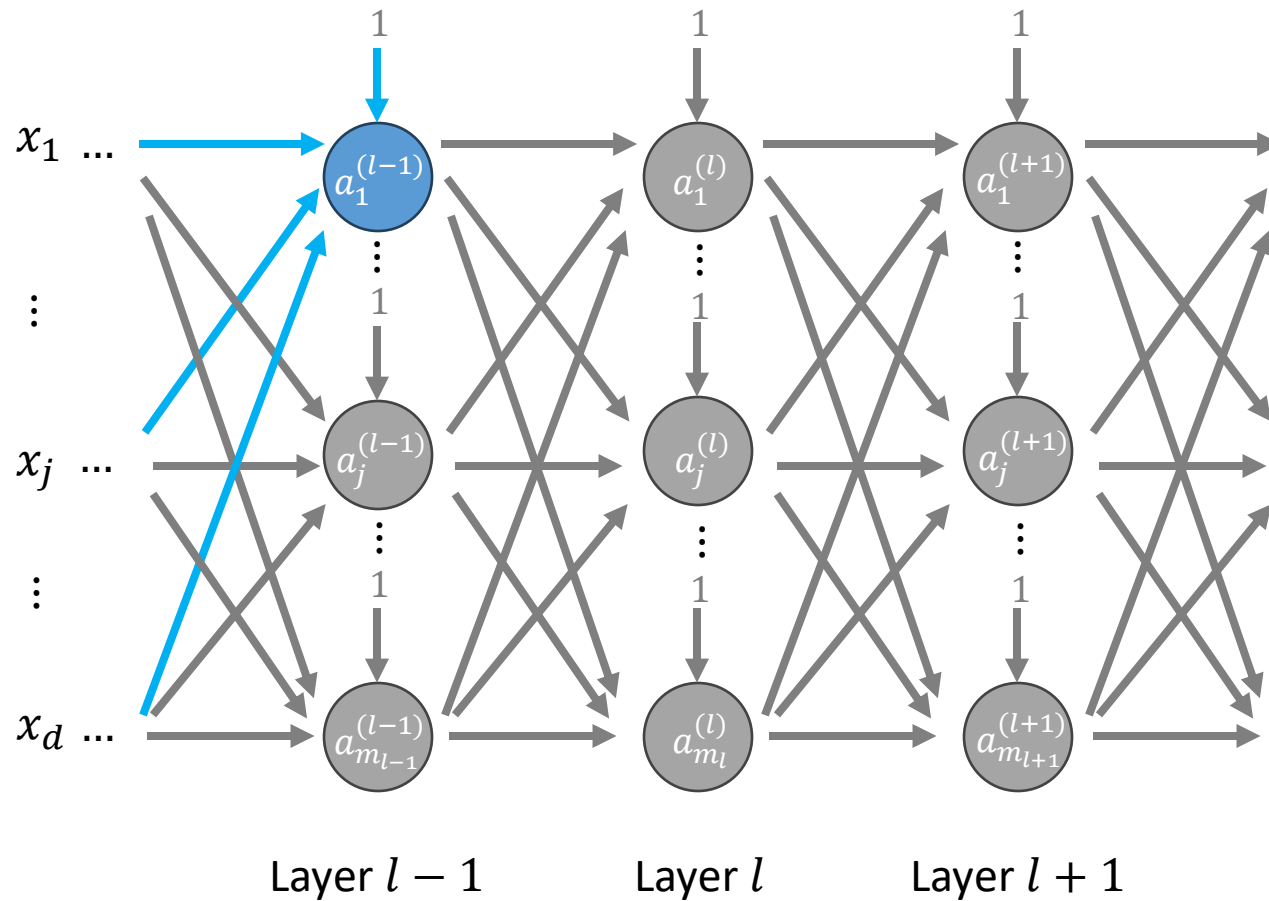
Linear regression relies on **manual feature engineering** to capture complex patterns, while a multi-layer neural network **learns** its own feature representations through its layers with activation functions.



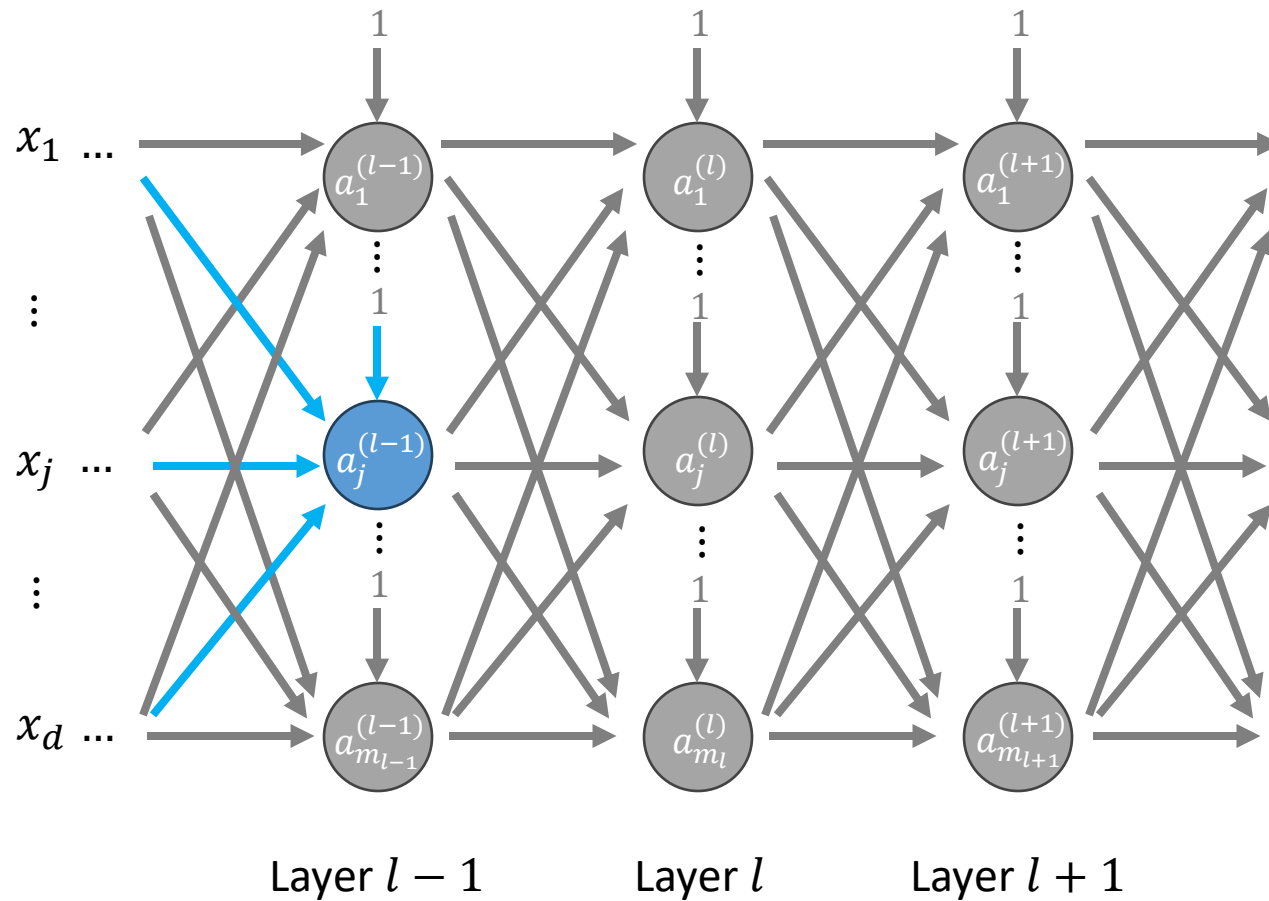
# Multi-layer Neural Networks



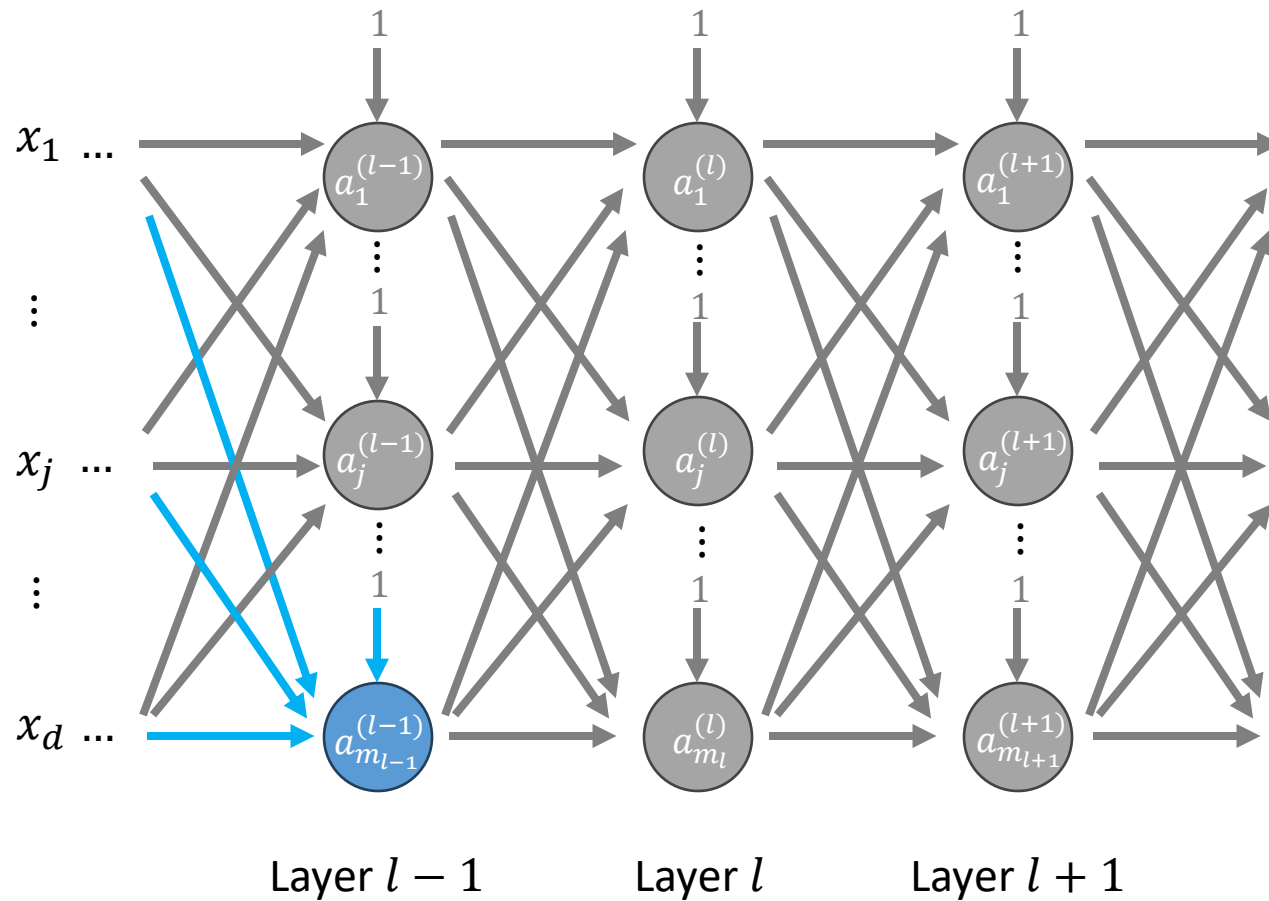
# Multi-layer Neural Networks



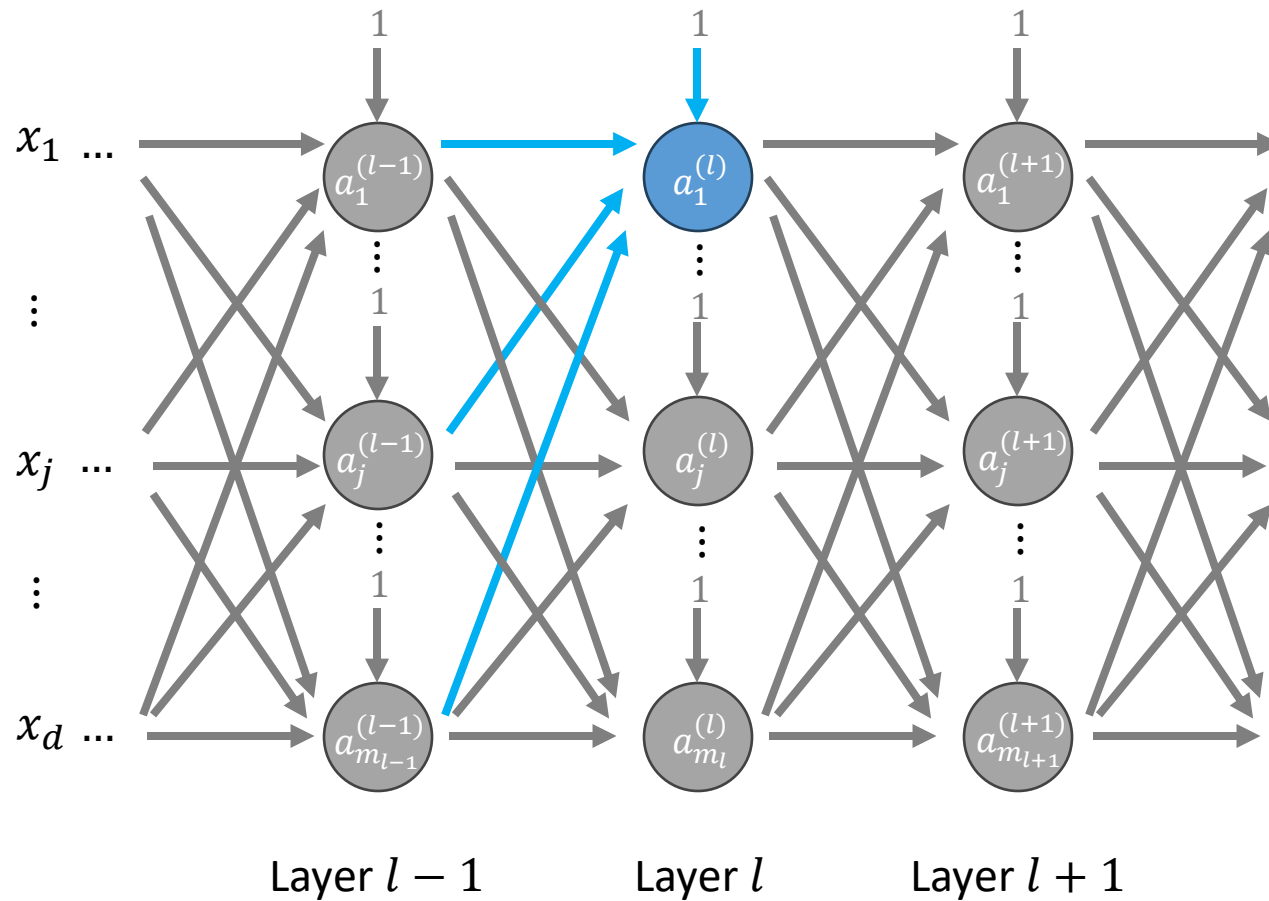
# Multi-layer Neural Networks



# Multi-layer Neural Networks

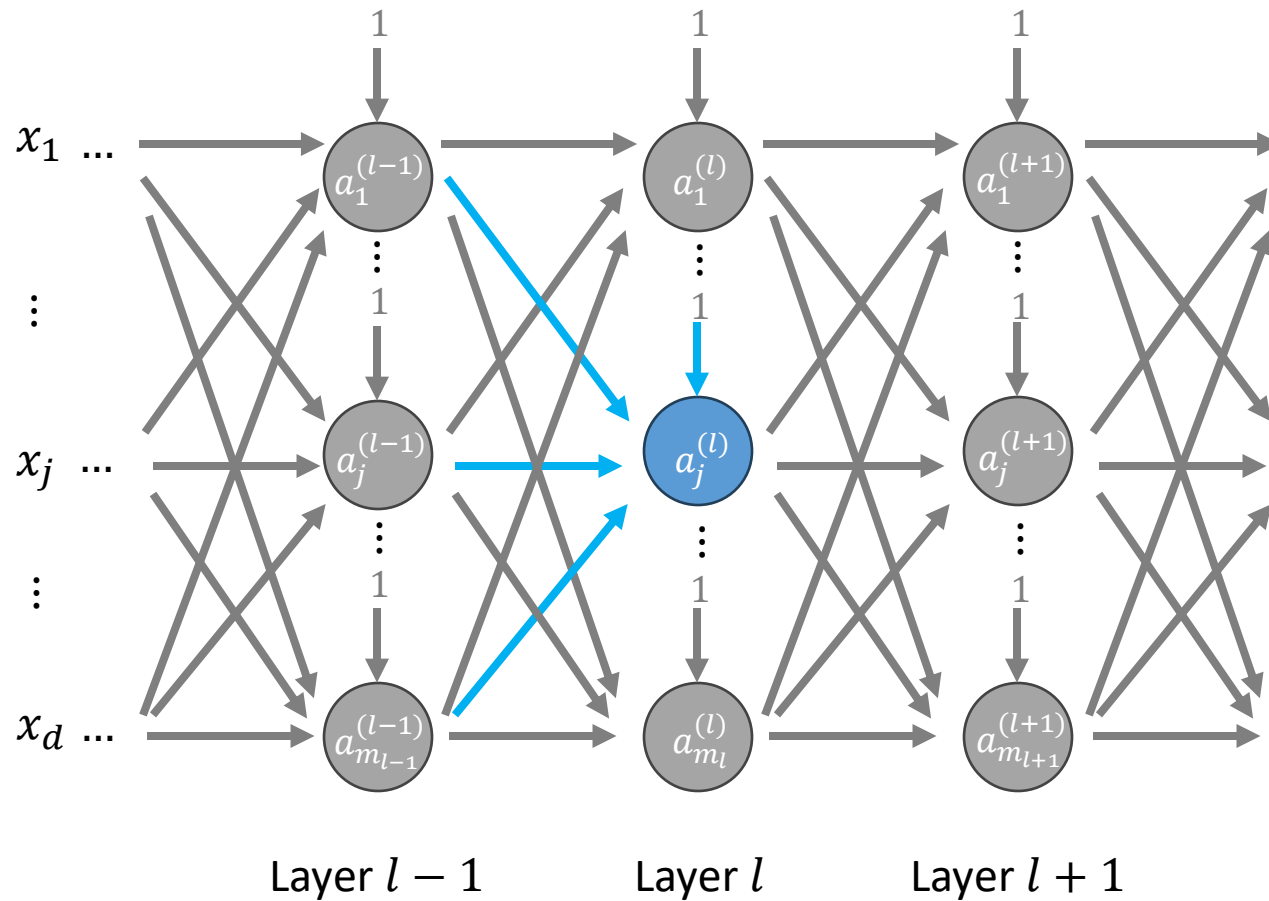


# Multi-layer Neural Networks

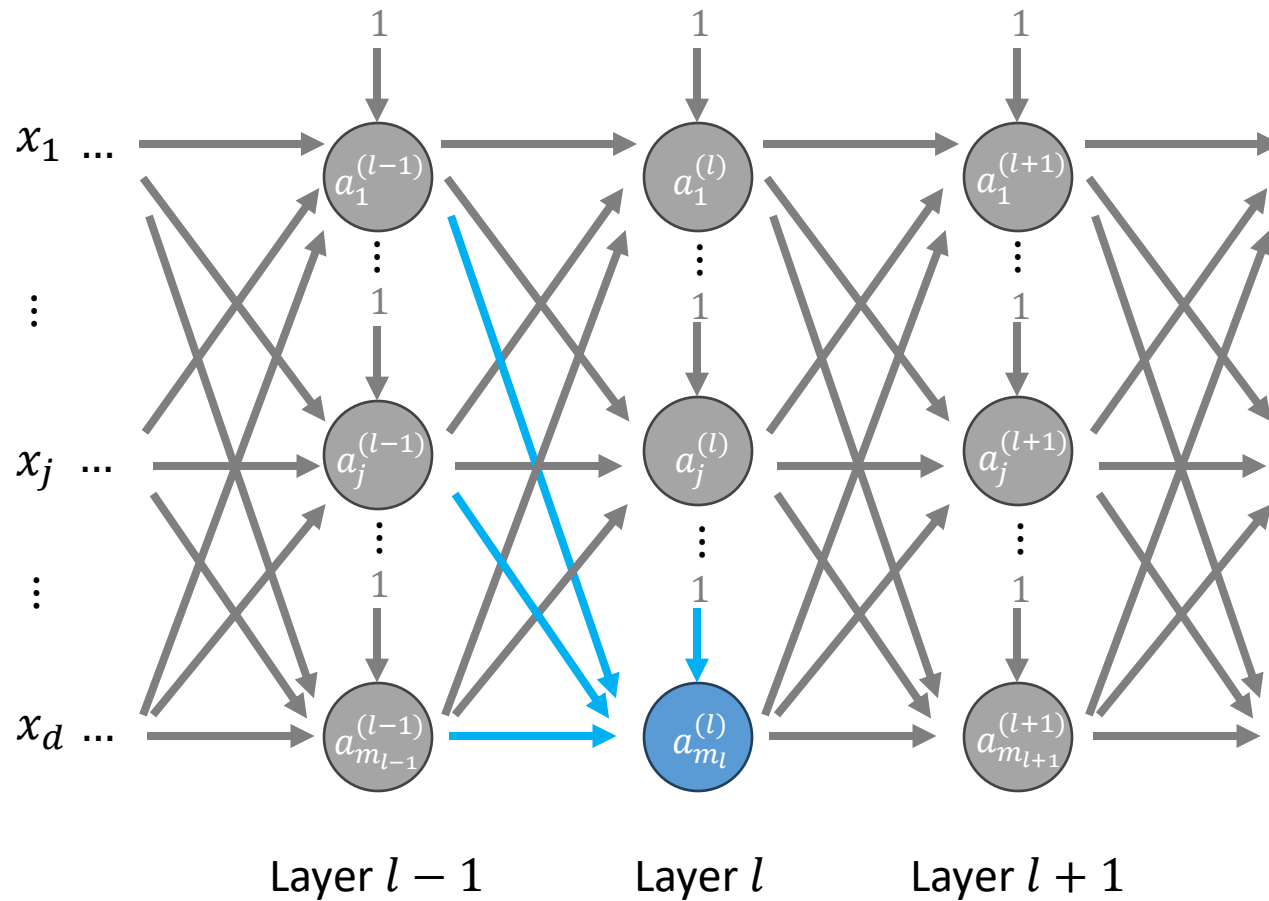




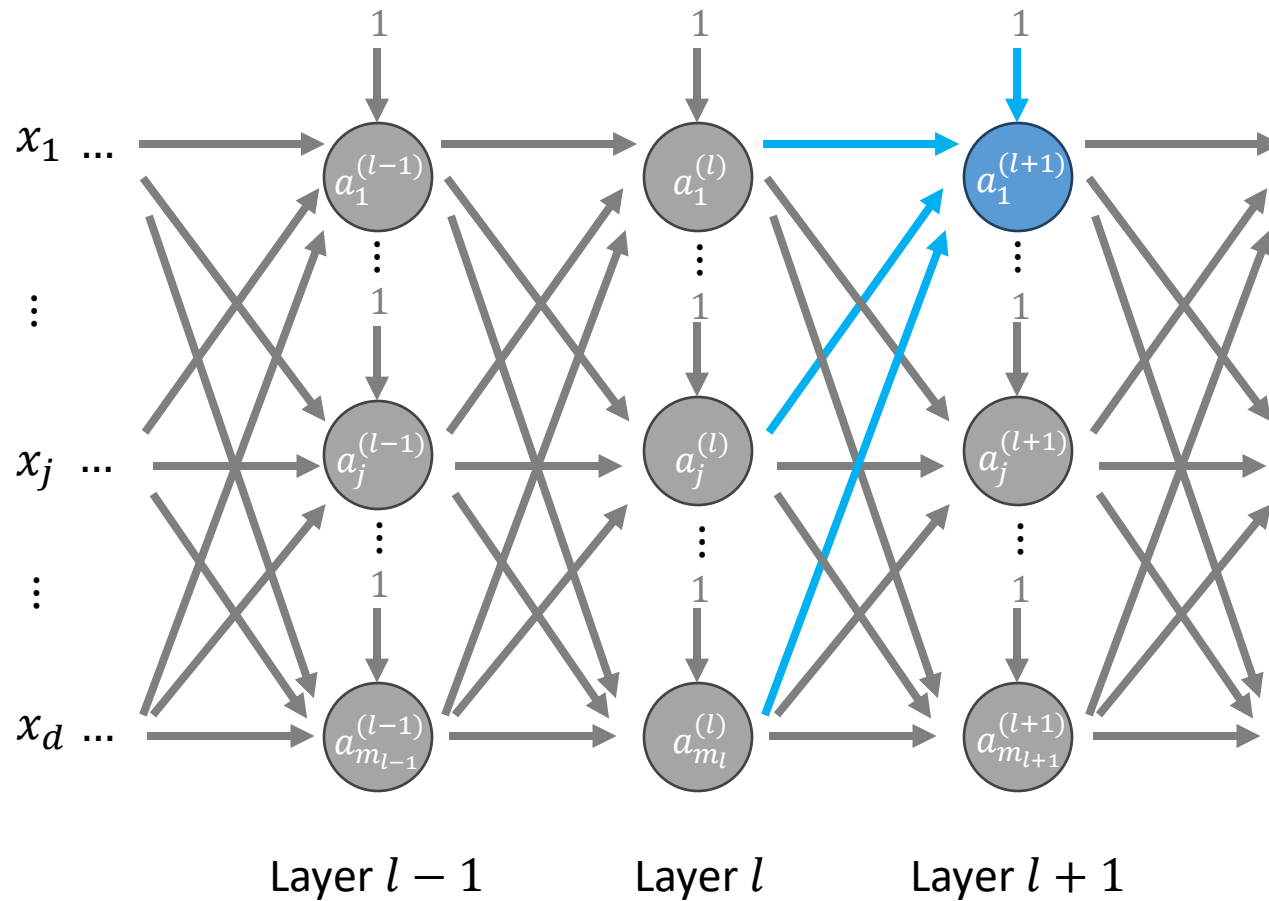
# Multi-layer Neural Networks



# Multi-layer Neural Networks

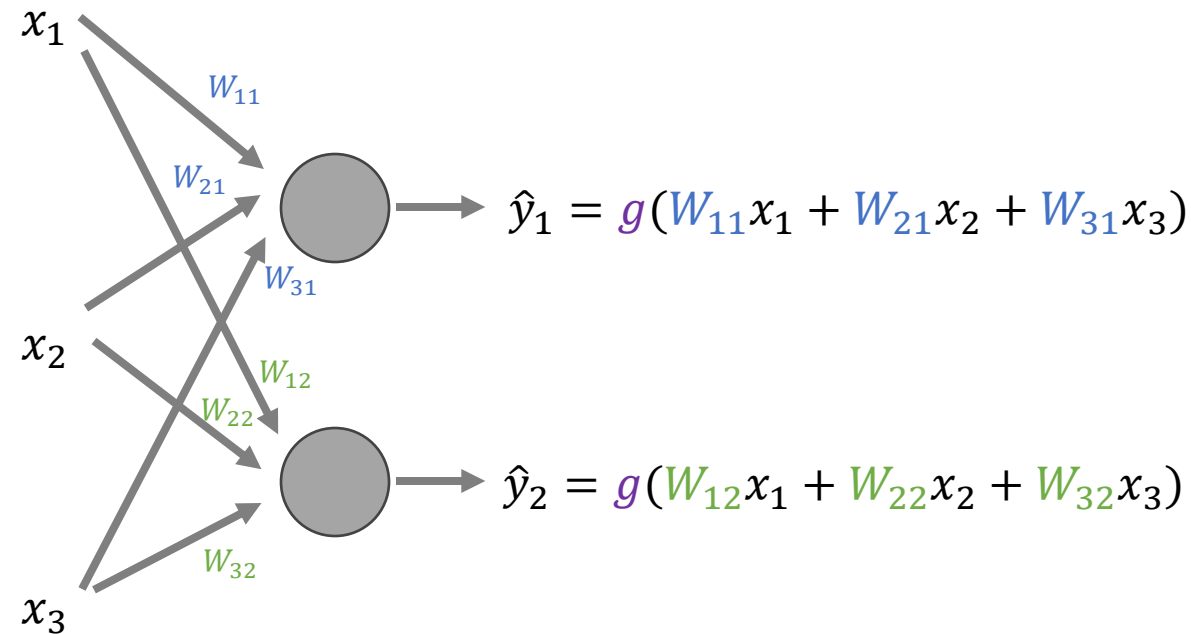


# Multi-layer Neural Networks



# Neural Networks and Matrix Multiplication

Single-layer



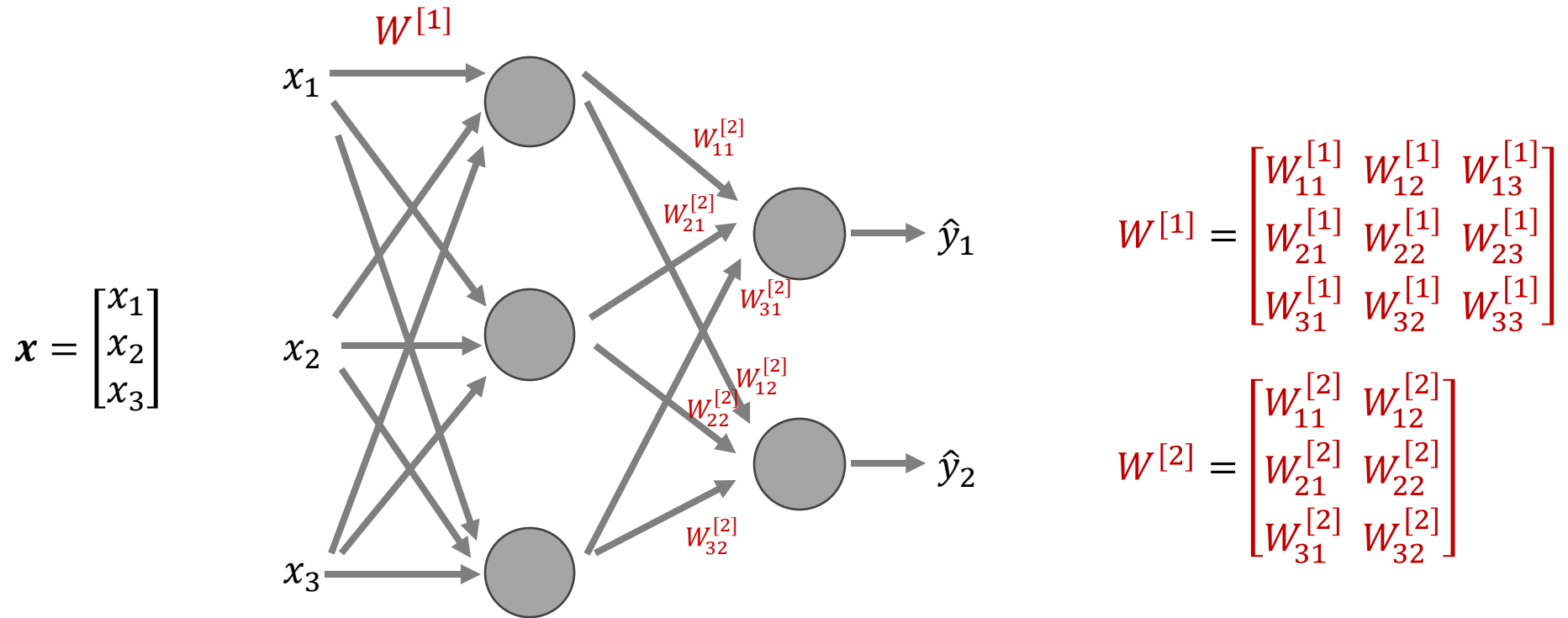
# Input (number of weights per neuron / input variables)

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \mathbf{W} = \begin{bmatrix} W_{11} & W_{12} \\ W_{21} & W_{22} \\ W_{31} & W_{32} \end{bmatrix} \quad \hat{\mathbf{y}} = g(\mathbf{W}^T \mathbf{x}) = g\left(\begin{bmatrix} W_{11} & W_{12} \\ W_{21} & W_{22} \\ W_{31} & W_{32} \end{bmatrix}^T \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}\right) = g\left(\begin{bmatrix} W_{11} & W_{21} & W_{31} \\ W_{12} & W_{22} & W_{32} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}\right) = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \end{bmatrix}$$

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

# Neural Networks and Matrix Multiplication

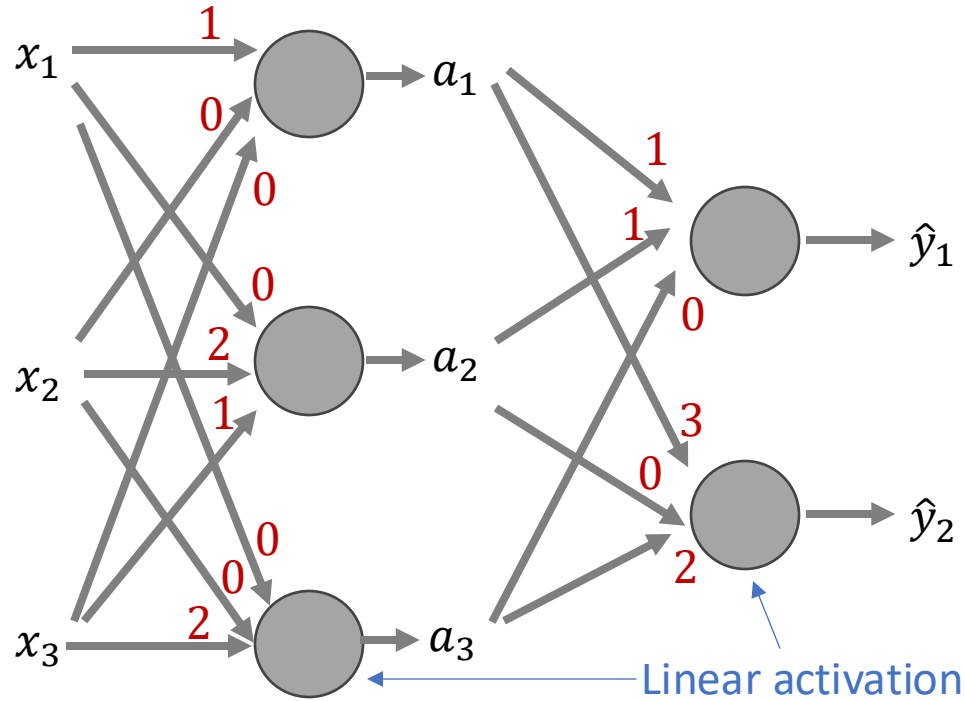
Multi-layer



$$\hat{\mathbf{y}} = g^{[2]} \left( W^{[2]T} g^{[1]} \left( W^{[1]T} \mathbf{x} \right) \right) = g^{[2]} \left( \begin{bmatrix} W_{11}^{[2]} & W_{12}^{[2]} \\ W_{21}^{[2]} & W_{22}^{[2]} \\ W_{31}^{[2]} & W_{32}^{[2]} \end{bmatrix}^T g^{[1]} \left( \begin{bmatrix} W_{11}^{[1]} & W_{12}^{[1]} & W_{13}^{[1]} \\ W_{21}^{[1]} & W_{22}^{[1]} & W_{23}^{[1]} \\ W_{31}^{[1]} & W_{32}^{[1]} & W_{33}^{[1]} \end{bmatrix}^T \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \right) \right) = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \end{bmatrix}$$

# Neural Networks and Matrix Multiplication

Multi-layer



$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} \quad \mathbf{W}^{[1]} = \begin{bmatrix} W_{11}^{[1]} & W_{12}^{[1]} & W_{13}^{[1]} \\ W_{21}^{[1]} & W_{22}^{[1]} & W_{23}^{[1]} \\ W_{31}^{[1]} & W_{32}^{[1]} & W_{33}^{[1]} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 1 & 2 \end{bmatrix}$$

$$\mathbf{a} = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} = \mathbf{W}^{[1]T} \mathbf{x} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 1 & 2 \end{bmatrix}^T \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2 & 1 \\ 0 & 0 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 5 \\ 2 \end{bmatrix}$$

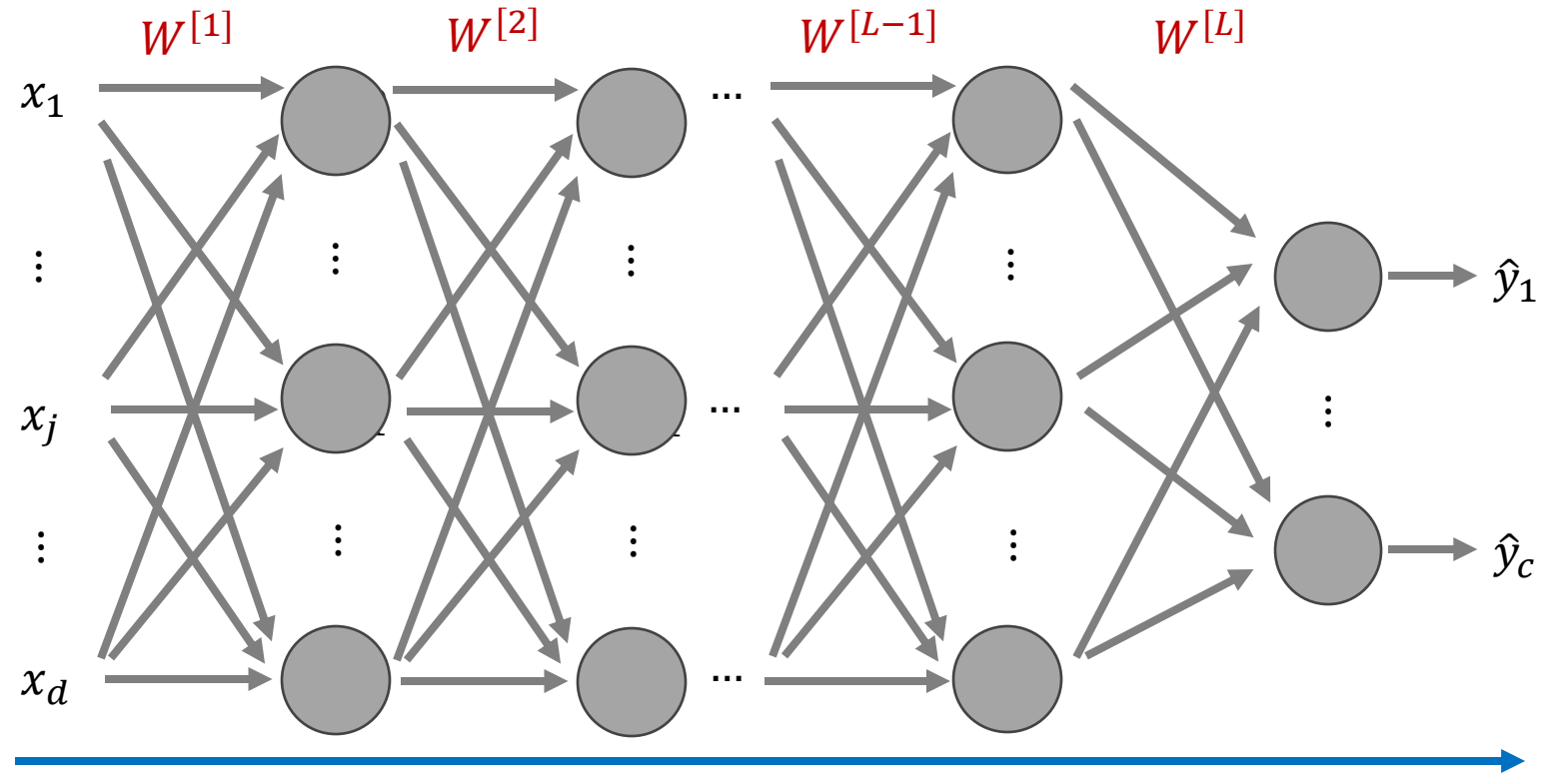
$$\mathbf{W}^{[2]} = \begin{bmatrix} W_{11}^{[2]} & W_{12}^{[2]} \\ W_{21}^{[2]} & W_{22}^{[2]} \\ W_{31}^{[2]} & W_{32}^{[2]} \end{bmatrix} = \begin{bmatrix} 1 & 3 \\ 1 & 0 \\ 0 & 2 \end{bmatrix}$$

$$\hat{\mathbf{y}} = \mathbf{W}^{[2]T} \mathbf{a} = \begin{bmatrix} 1 & 3 \\ 1 & 0 \\ 0 & 2 \end{bmatrix}^T \begin{bmatrix} 1 \\ 5 \\ 2 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 0 \\ 3 & 0 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 5 \\ 2 \end{bmatrix} = \begin{bmatrix} 6 \\ 7 \end{bmatrix}$$

$$\hat{\mathbf{y}} = g^{[2]} \left( \mathbf{W}^{[2]T} g^{[1]} \left( \mathbf{W}^{[1]T} \mathbf{x} \right) \right) = g^{[2]} \left( \begin{bmatrix} W_{11}^{[2]} & W_{12}^{[2]} \\ W_{21}^{[2]} & W_{22}^{[2]} \\ W_{31}^{[2]} & W_{32}^{[2]} \end{bmatrix}^T g^{[1]} \left( \begin{bmatrix} W_{11}^{[1]} & W_{12}^{[1]} & W_{13}^{[1]} \\ W_{21}^{[1]} & W_{22}^{[1]} & W_{23}^{[1]} \\ W_{31}^{[1]} & W_{32}^{[1]} & W_{33}^{[1]} \end{bmatrix}^T \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \right) \right) = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \end{bmatrix}$$

# Neural Networks and Matrix Multiplication

Multi-layer



The process by which input data passes through a neural network to generate an output prediction.

Forward Propagation

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

# Outline

- Perceptron
  - Biological inspiration
  - Perceptron Learning Algorithm
- Neural Network
  - Neuron
  - AND Gate Modelling
  - XNOR Gate Modelling
  - Single-layer and Multi-layer Neural Networks
- **Multi-class Classification**

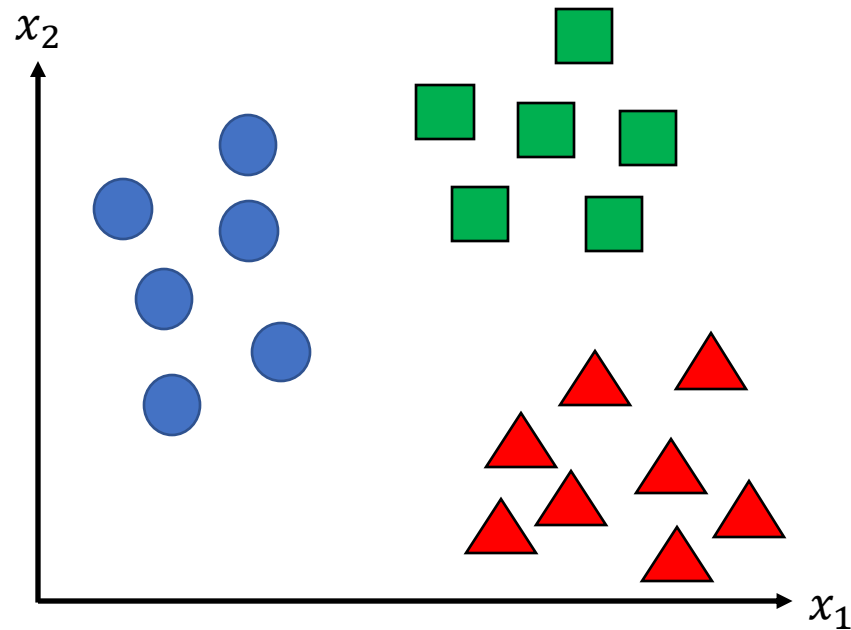


# Multi-class Classification

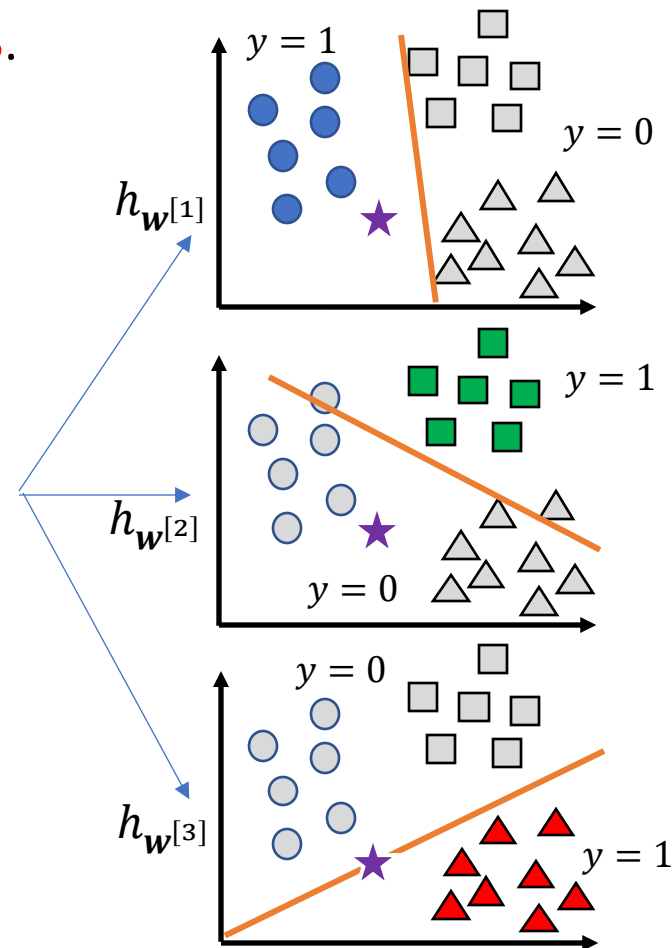
With Logistic Regression Model

Given a set of features (e.g.,  $x_1, x_2$ ),

Predict whether it belongs to **class 1**, **2**, or **3**.



- Fit one classifier per class,
- Fit against all other classes
- Pick **highest probability**



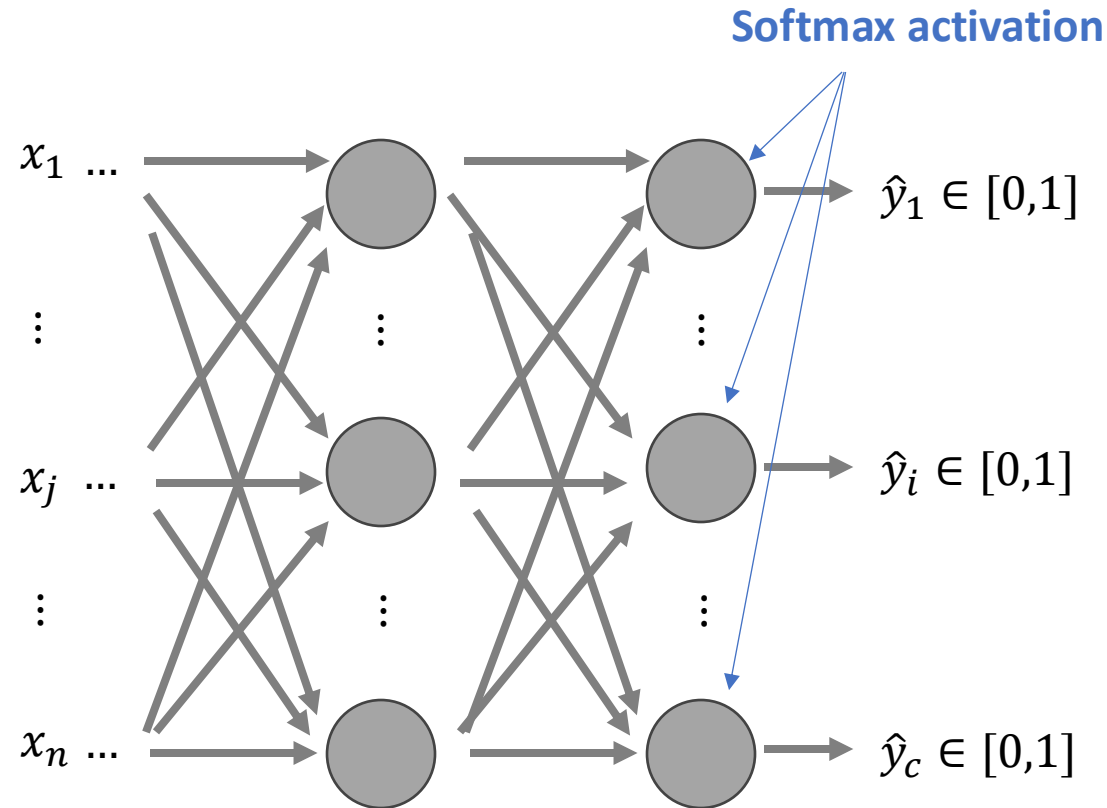
$$h_{\mathbf{w}^{[1]}}(\star) = 0.8 \quad \checkmark$$

$$h_{\mathbf{w}^{[2]}}(\star) = 0.1$$

$$h_{\mathbf{w}^{[3]}}(\star) = 0.5$$

# Multi-class Classification

With Neural Network



- Set the number of neurons in the last layer as  $c$
- Predict the probability of each class
- Pick **highest probability**

$c$ -class Classification

# Softmax Function

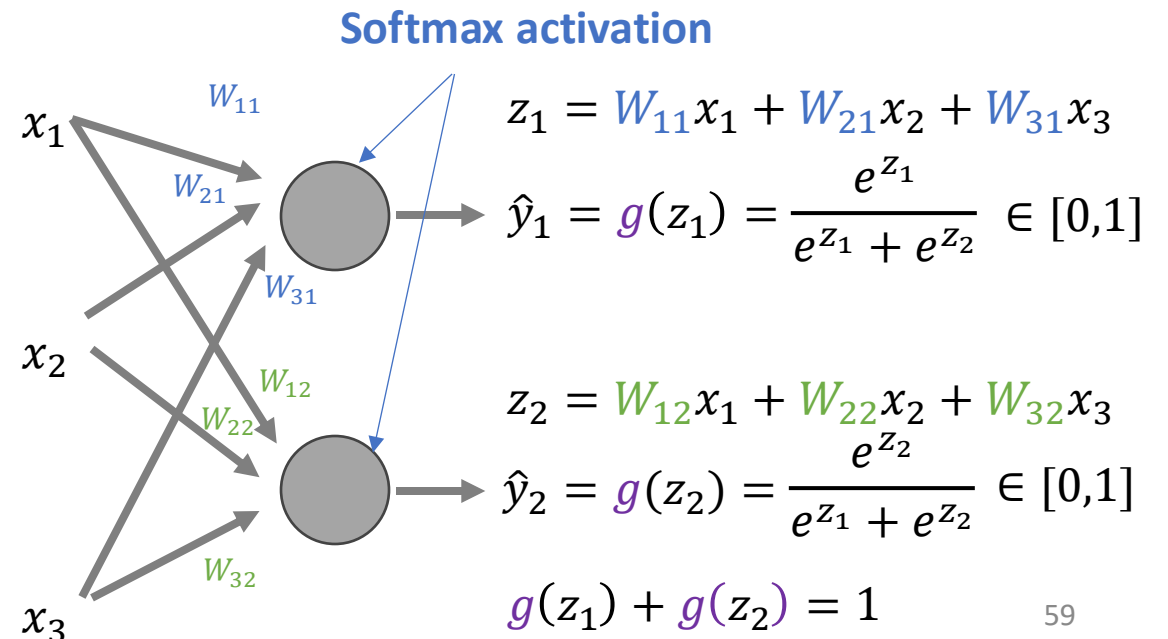
Given a vector  $\mathbf{z} = [z_1, z_2 \dots z_c]$ , the softmax function computes the output for each  $z_i$  as:

$$g(z_i) = \frac{e^{z_i}}{\sum_{j=1}^c e^{z_j}} \in [0,1]$$

And  $g(z_1) + g(z_2) \dots + g(z_c) = 1$ .

Softmax function can be used to

- Convert raw scores  $z_i$  into probabilities
- Ensure that the output values sum to 1.



# Summary

- 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
  - Linear regression model: A neuron with linear activation function
  - Logistic regression model: A neuron with sigmoid activation function
  - Neurons can be used to transform features
  - Multi-layer Neural Networks
- Multi-class Classification
  - Softmax activation function should be applied.

# Coming Up Next Week

- **More on Neural Networks**
  - Backpropagation
  - Convolution Neural Networks
  - ...

# To Do

- **Lecture Training 9**
  - +250 Free EXP
  - +100 Early bird bonus

