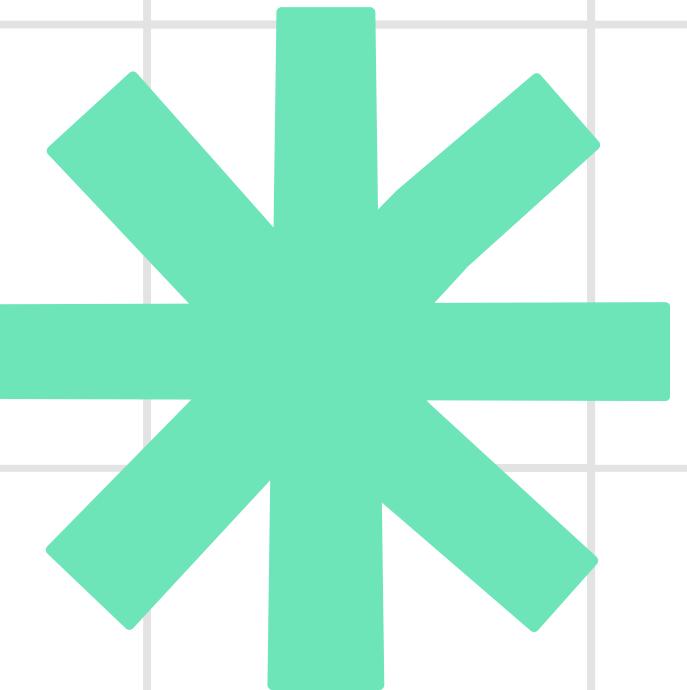


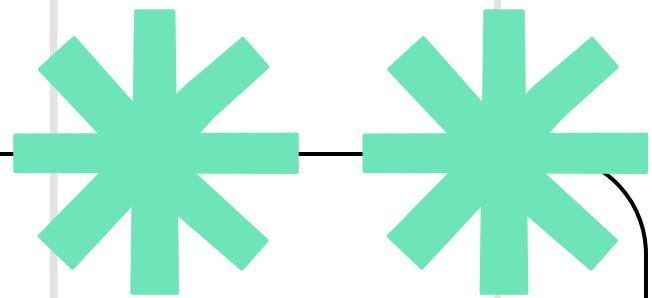
101

# Deep Learning 101

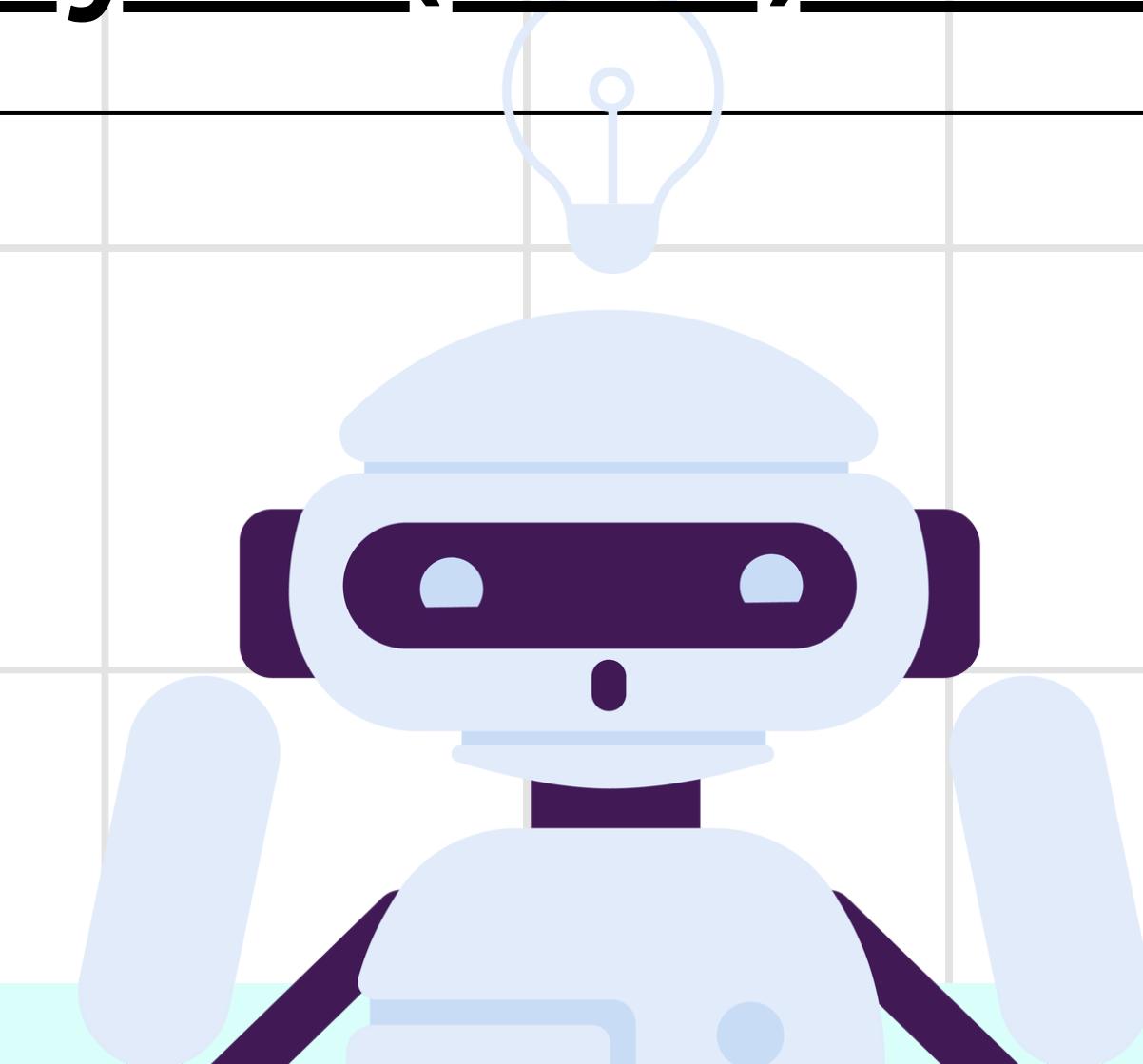
Selma MANI



For this session we will be using illustrations from  
the book

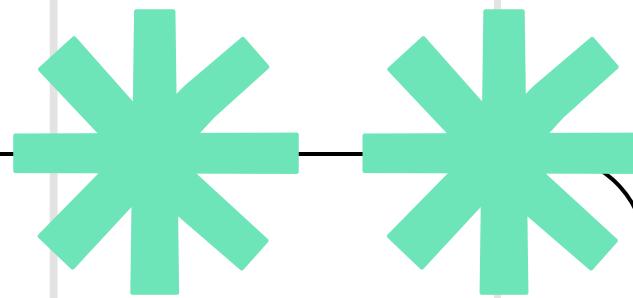
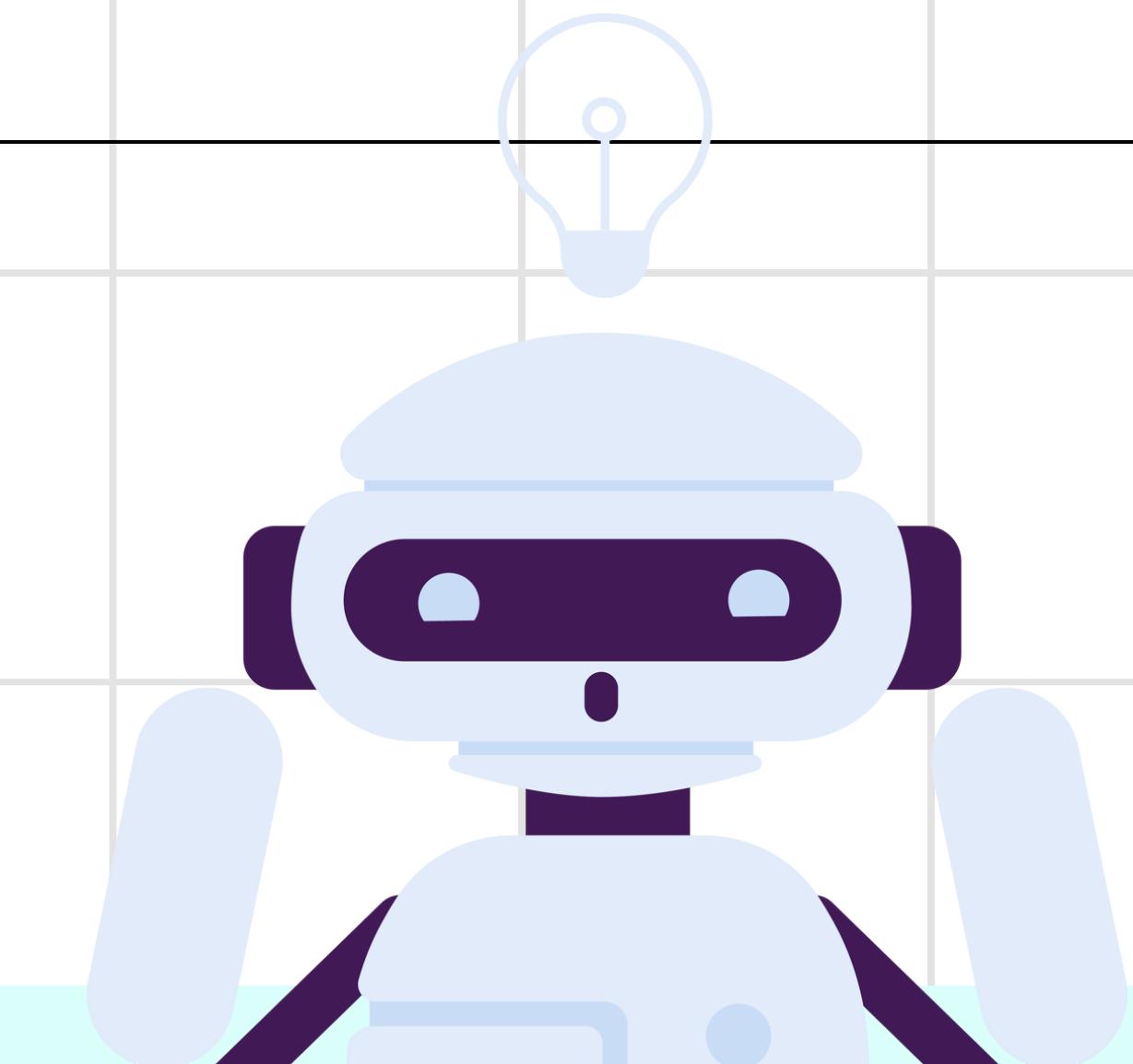


## **Fundamentals of ML for Predictive Data** **Analytics (2020) 2nd edition**

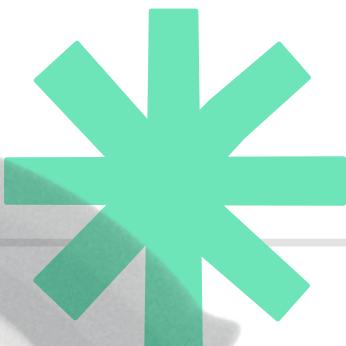


The notebook to this session

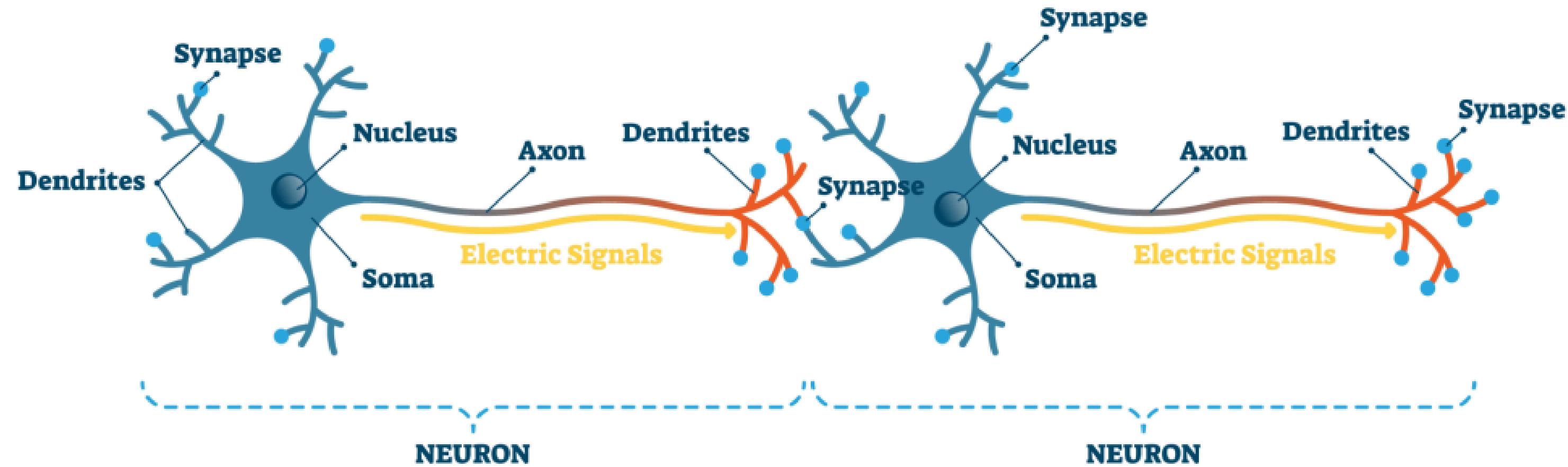
## **Deep Learning 101 Notebook**



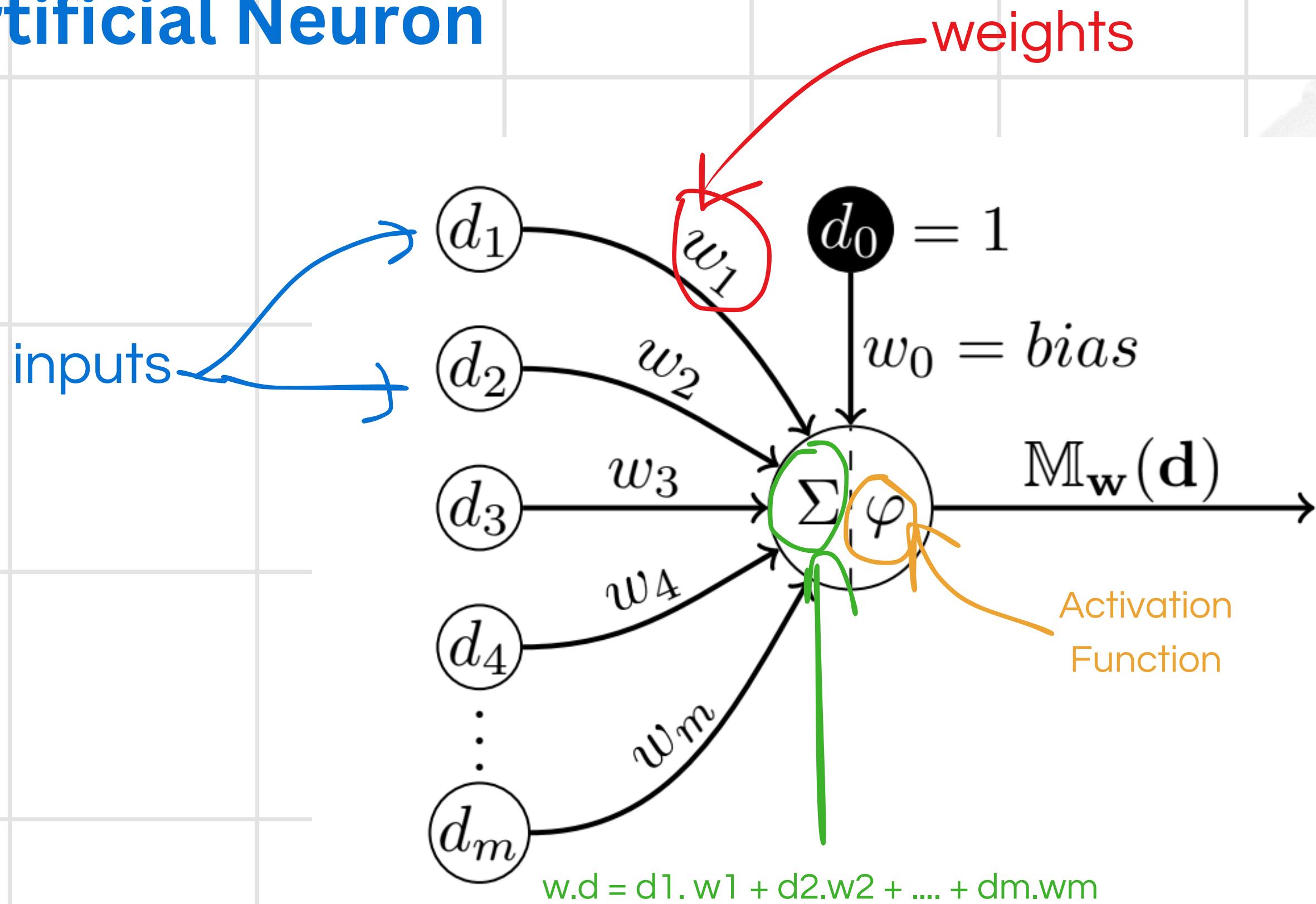
# How does our brain work?



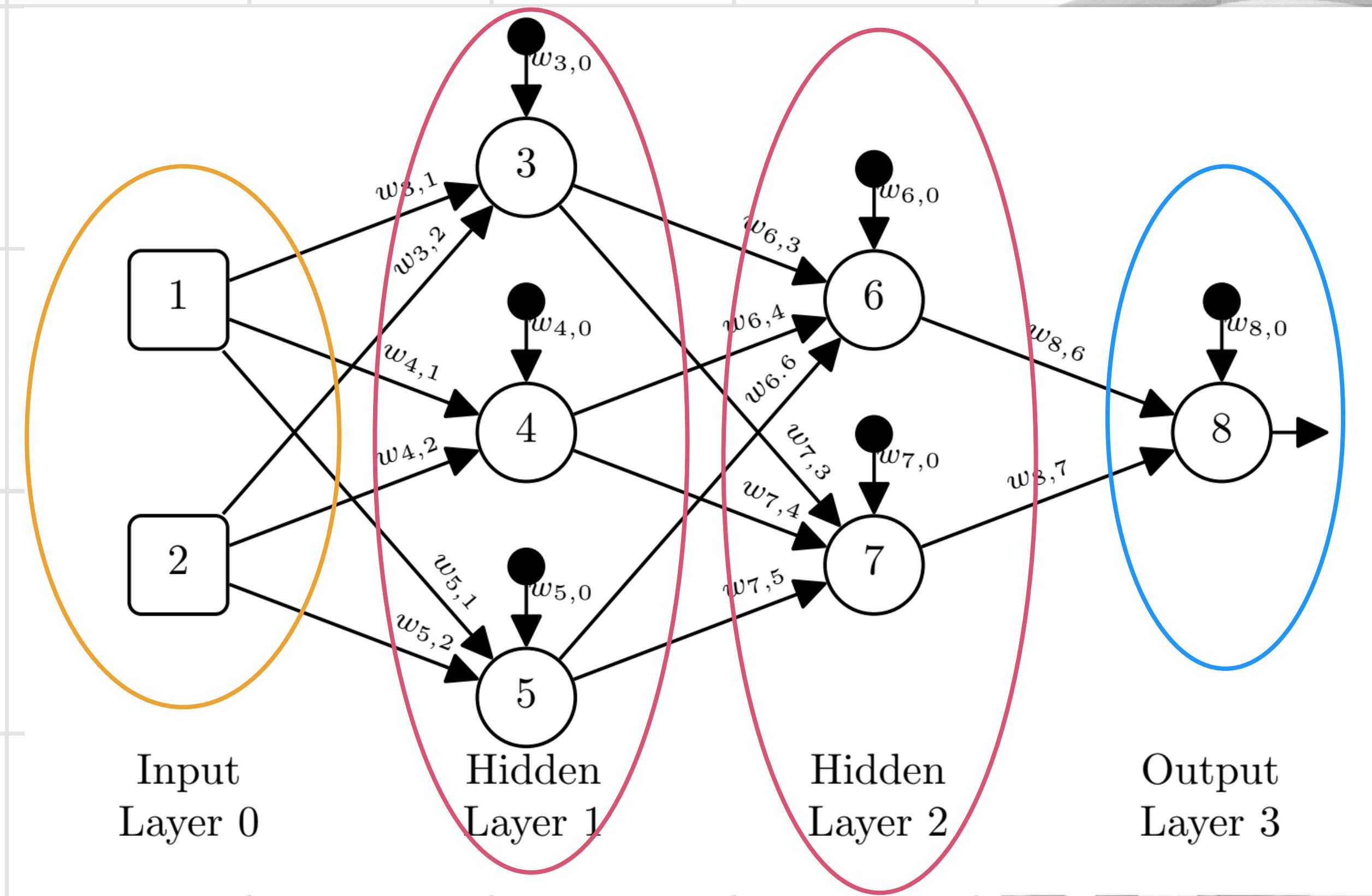
# How does our brain work?



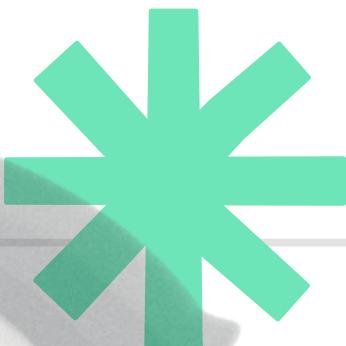
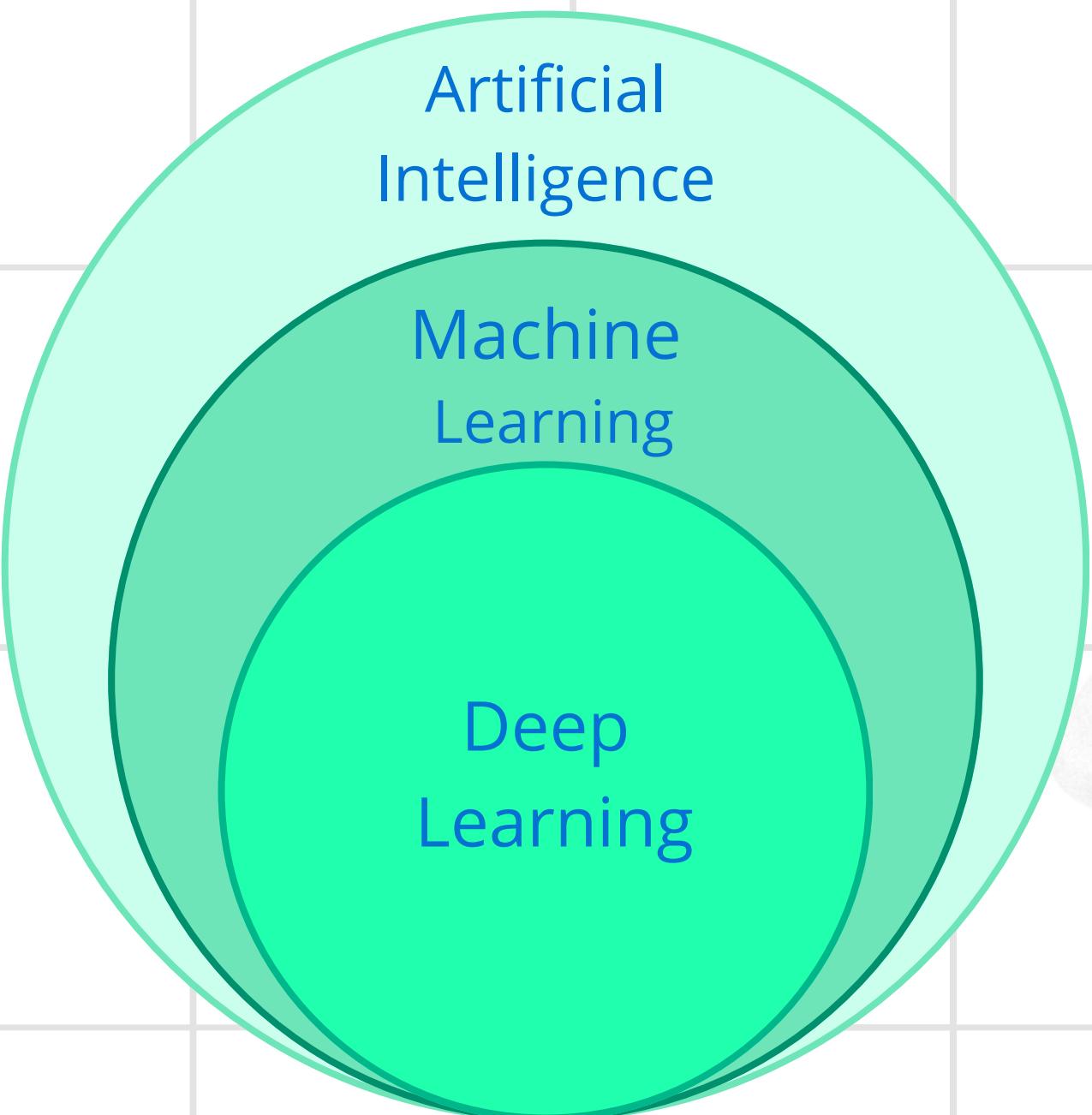
# Artificial Neuron



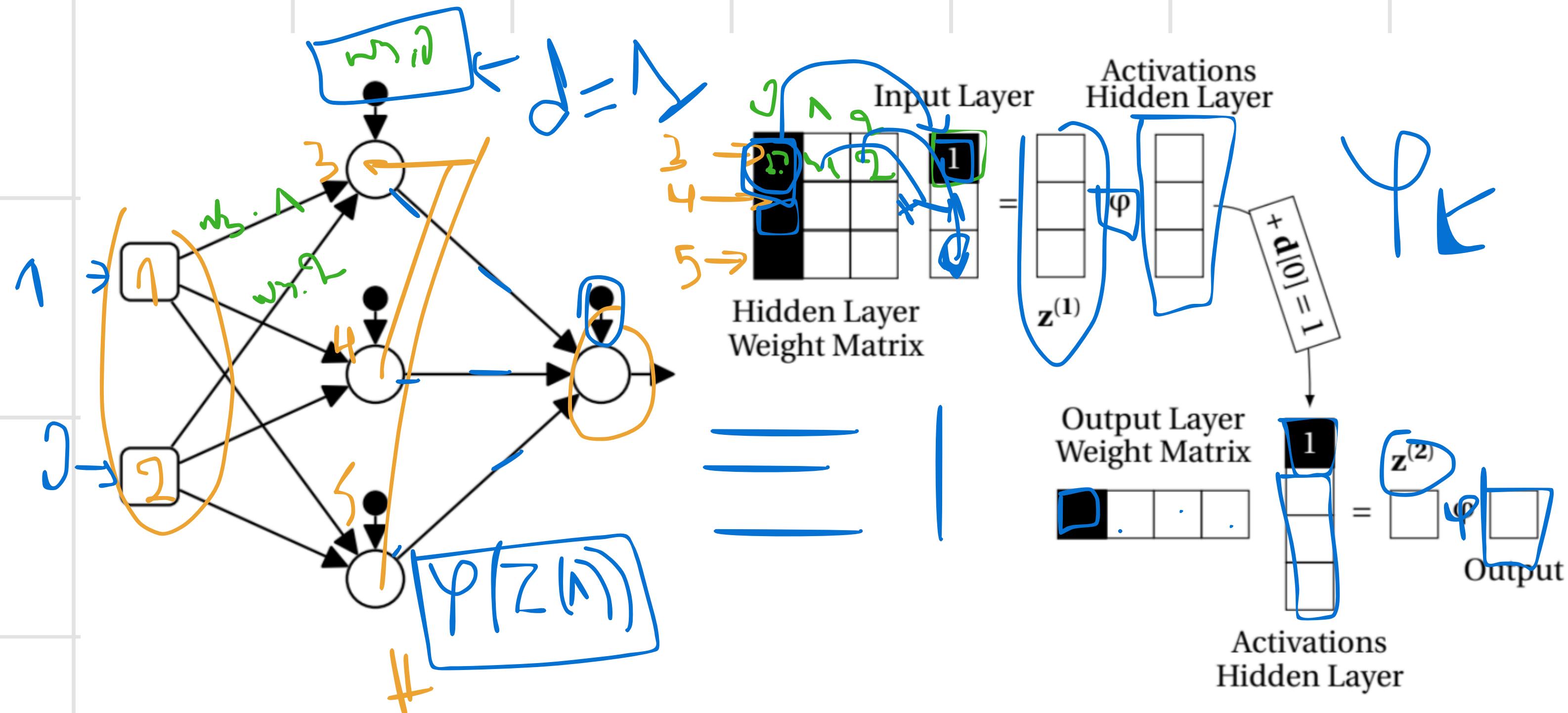
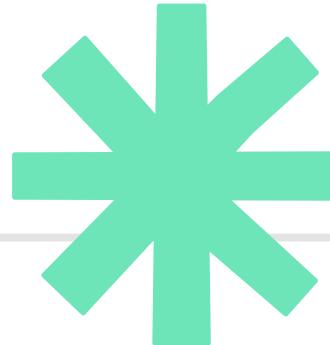
# Artificial Neuron Network



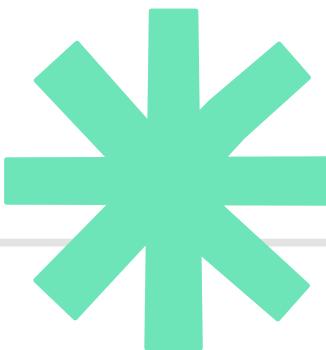
# Machine Learning VS Deep Learning



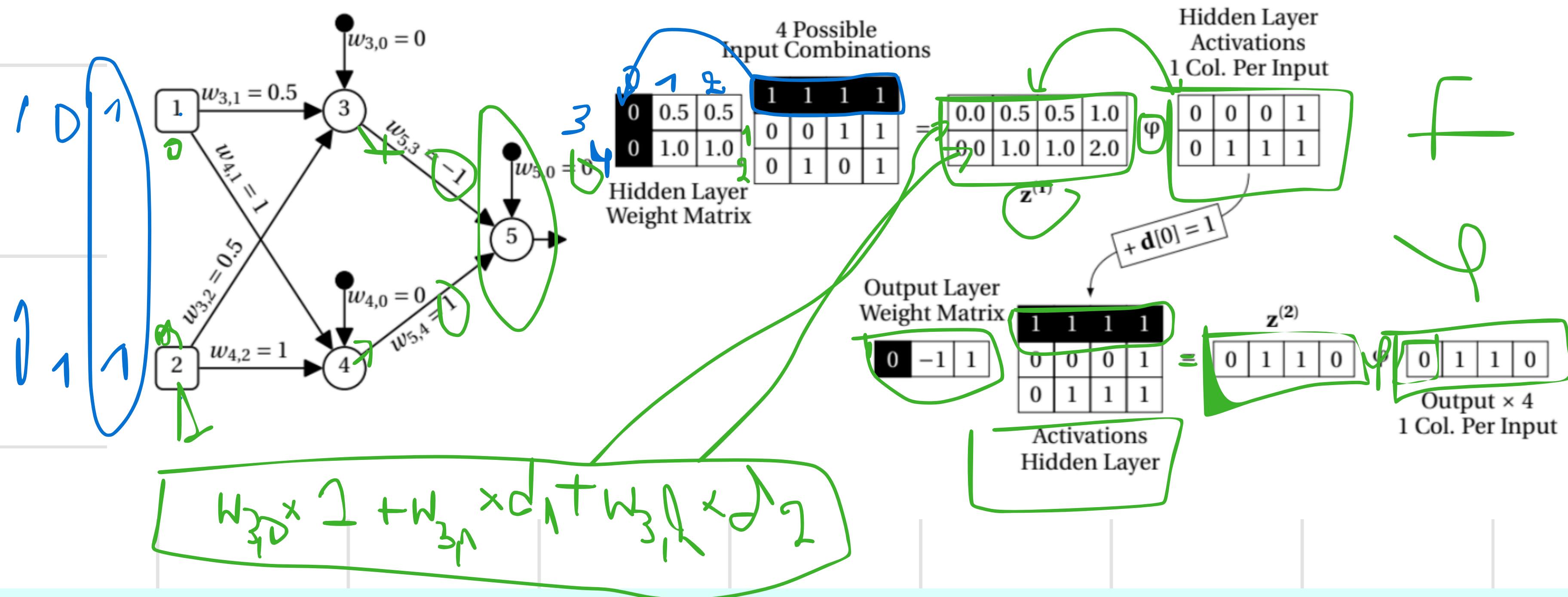
# Matrices in Deep Learning



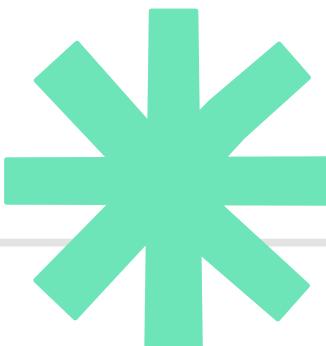
# Activation Function



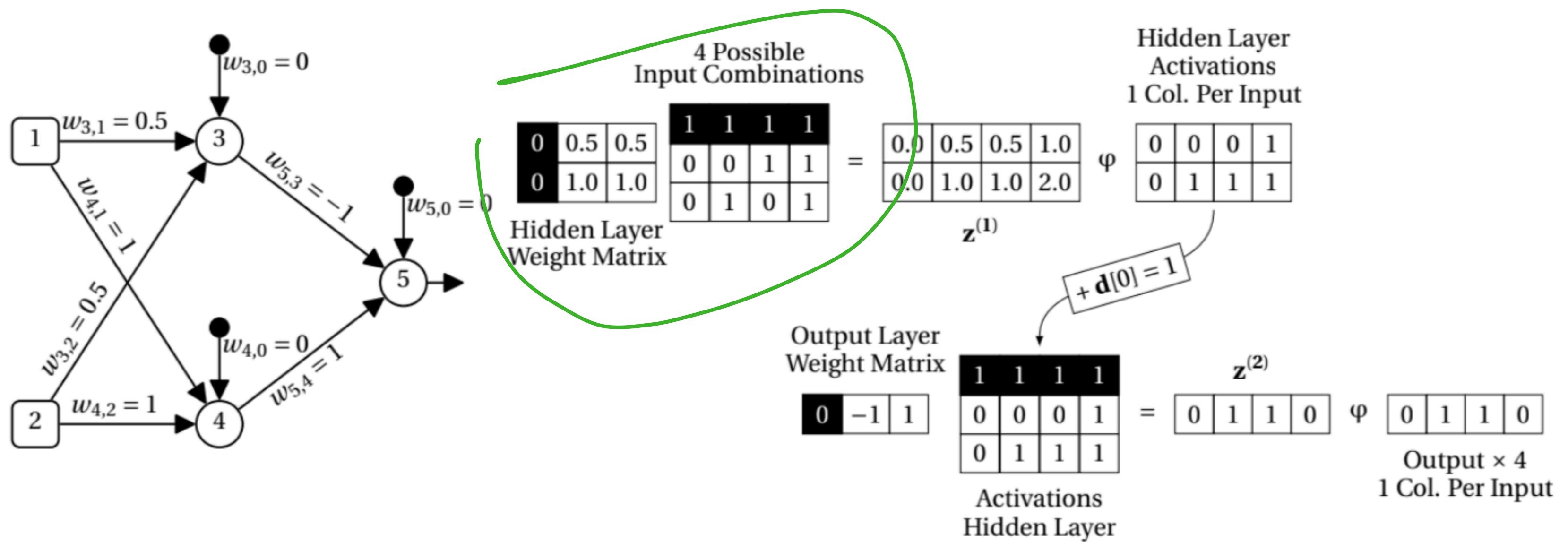
What is the activation function in this example?



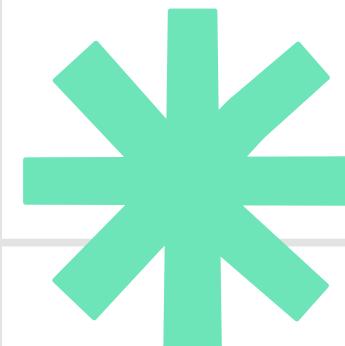
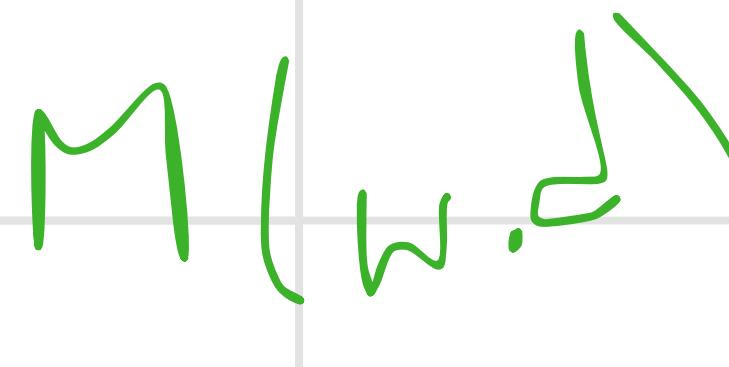
# Activation Function



$M(w.d) = 0 \text{ if } z < 1 \text{ else } 1$



# Activation Function



$$M_w(\mathbf{d}) = \varphi (\mathbf{w}[0] \times \mathbf{d}[0] + \mathbf{w}[1] \times \mathbf{d}[1] + \dots + \mathbf{w}[m] \times \mathbf{d}[m])$$

$$= \varphi \left( \sum_{i=0}^m w_i \times d_i \right) = \varphi \left( \underbrace{\mathbf{w} \cdot \mathbf{d}}_{dot\ product} \right)$$

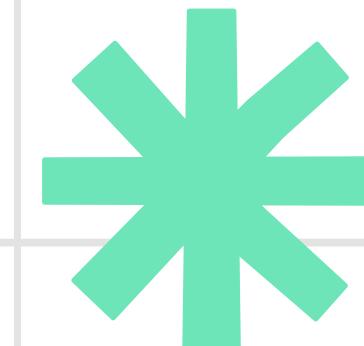
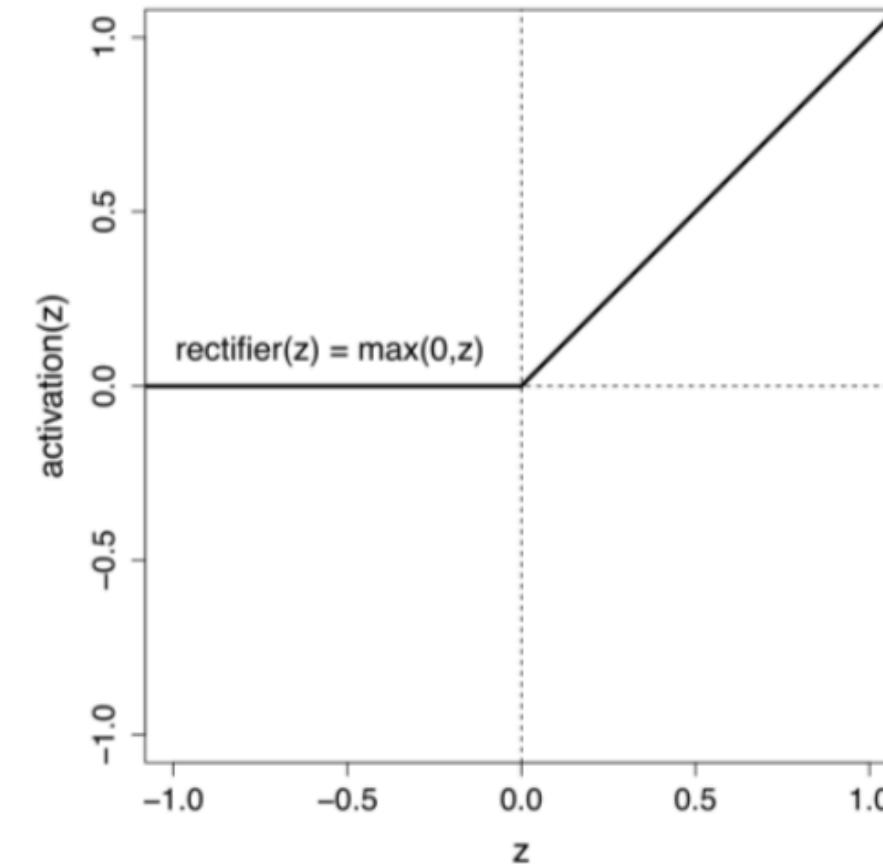
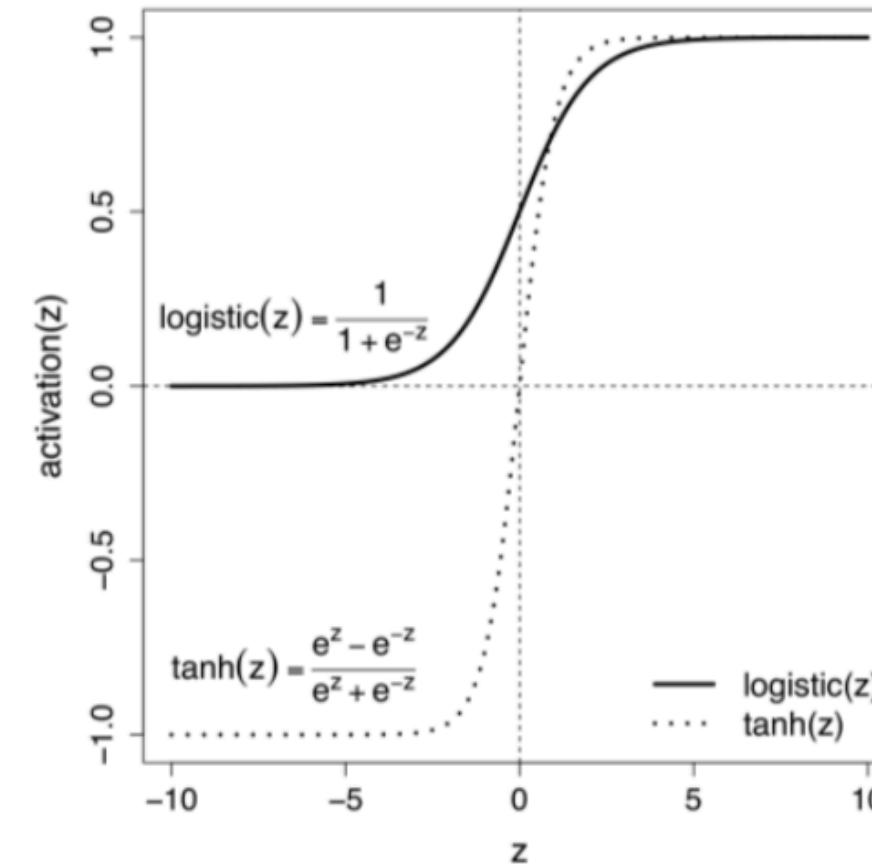
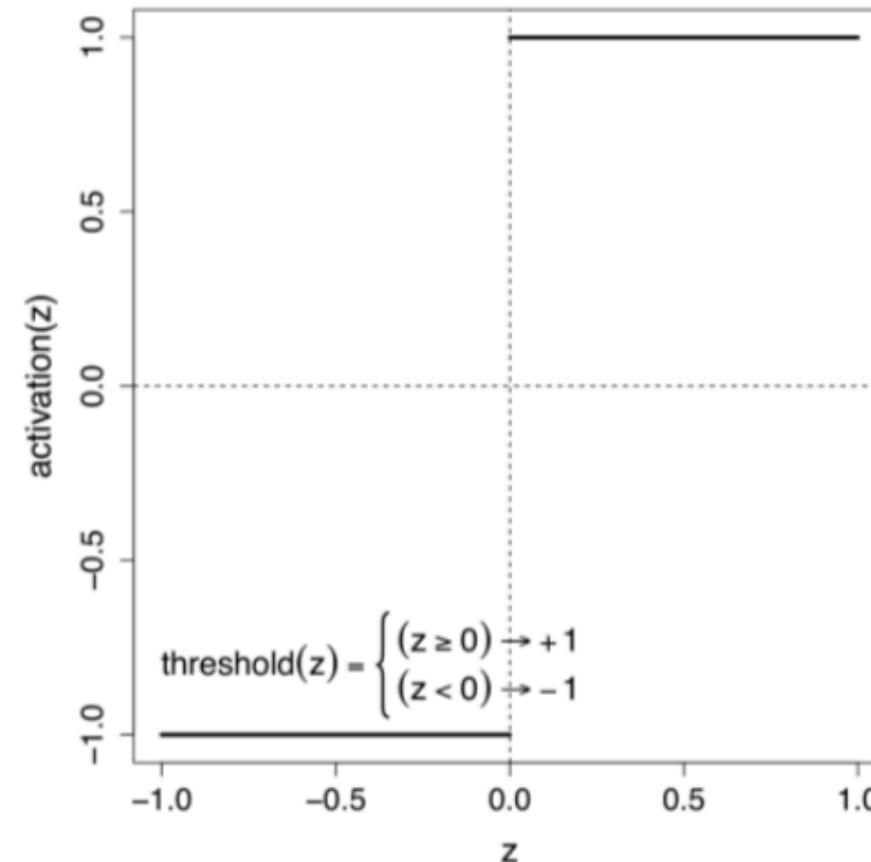
$$= \varphi \left( \underbrace{\mathbf{w}^T \mathbf{d}}_{matrix\ product} \right) = \varphi \left( [w_0, w_1, \dots, w_m] \begin{bmatrix} d_0 \\ d_1 \\ \vdots \\ d_m \end{bmatrix} \right)$$

$$\begin{bmatrix} d_0 \\ d_1 \\ \vdots \\ d_m \end{bmatrix}$$

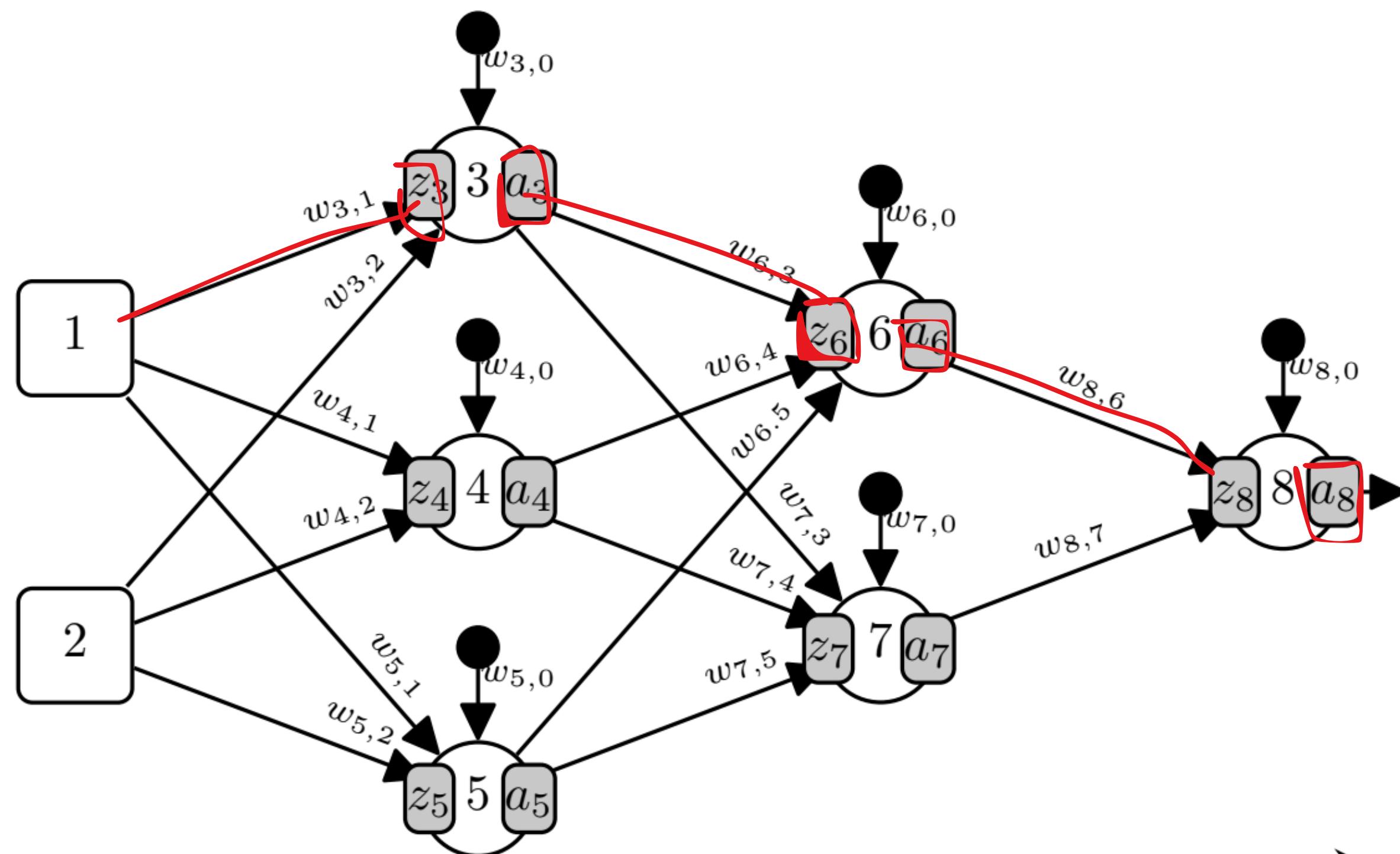
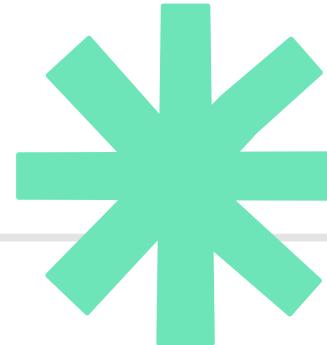
# Activation Function

$$M_w(d) = \begin{cases} 1 & \text{if } z \geq \theta \\ 0 & \text{otherwise} \end{cases}$$

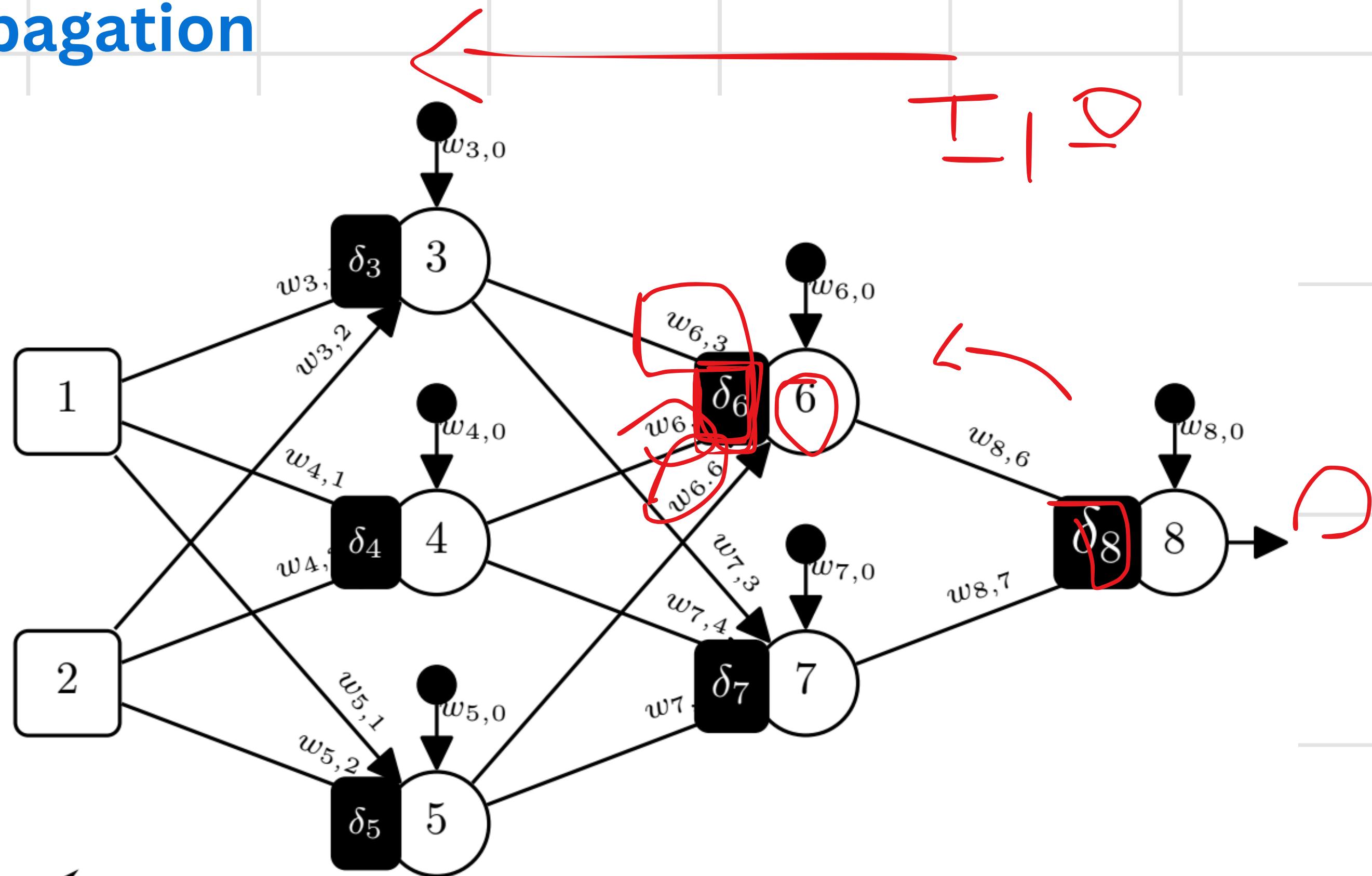
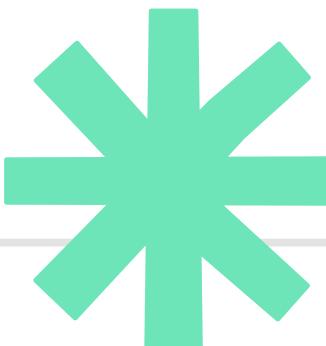
$$\text{rectifier}(z) = \max(0, z)$$



# Forward Propagation



# Backpropagation



# Backpropagation



```
1: Shuffle  $\mathcal{D}$  and create the mini-batches:  $[(\mathbf{X}^{(1)}, \mathbf{Y}^{(1)}), \dots, (\mathbf{X}^k, \mathbf{Y}^k)]$ 
2: Initialize the weight matrices for each layer:  $\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(L)}$ 
3: repeat
4:   for  $t=1$  to number of mini-batches do
5:      $\mathbf{A}^{(0)} \leftarrow \mathbf{X}^{(t)}$ 
6:     for  $l=1$  to  $L$  do
7:        $\mathbf{v} \leftarrow [1_0, \dots, 1_m]$ 
8:        $\mathbf{A}^{(l-1)} \leftarrow [\mathbf{v}; \mathbf{A}^{(l-1)}]$ 
9:        $\mathbf{Z}^{(l)} \leftarrow \mathbf{W}^l \mathbf{A}^{(l-1)}$ 
10:       $\mathbf{A}^{(l)} \leftarrow \varphi(\mathbf{Z}^{(l)})$ 
11:    end for
12:    for each weight  $w_{i,k}$  in the network do
13:       $\Delta w_{i,k} = 0$ 
14:    end for
15:    for each example in the mini-batch do
16:      for each neuron  $i$  in the output layer do
17:         $\delta_i = \frac{\partial \mathcal{E}}{\partial a_i} \times \frac{\partial a_i}{\partial z_i}$ 
18:      end for
19:      for  $l = L-1$  to  $1$  do
20:        for each neuron  $i$  in the layer  $l$  do
21:           $\delta_i = \frac{\partial \mathcal{E}}{\partial a_i} \times \frac{\partial a_i}{\partial z_i}$ 
22:        end for
23:      end for
24:      for each weight  $w_{i,k}$  in the network do
25:         $\Delta w_{i,k} = \Delta w_{i,k} + (\delta_i \times a_k)$ 
26:      end for
27:    end for
28:    for each weight  $w_{i,k}$  in the network do
29:       $w_{i,k} \leftarrow v_{i,k} - \alpha \times \Delta w_{i,k}$ 
30:    end for
31:  end for
32: shuffle $[(\mathbf{X}^{(1)}, \mathbf{Y}^{(1)}), \dots, (\mathbf{X}^k, \mathbf{Y}^k)]$ 
33: until convergence occurs
```

▷ Each repeat loop is one epoch  
▷ Each for loop is one iteration

▷ Create  $\mathbf{v}$  the vector of bias terms  
▷ Insert  $\mathbf{v}$  into the activation matrix

▷ Elementwise application of  $\varphi$  to  $\mathbf{Z}^{(l)}$

▷ Backpropagate the  $\delta$ s

▷ See Equation (9)<sup>[10]</sup>

▷ See Equation (11)<sup>[11]</sup>

▷ Equation (17)<sup>[15]</sup>

▷ Equation (18)<sup>[15]</sup>

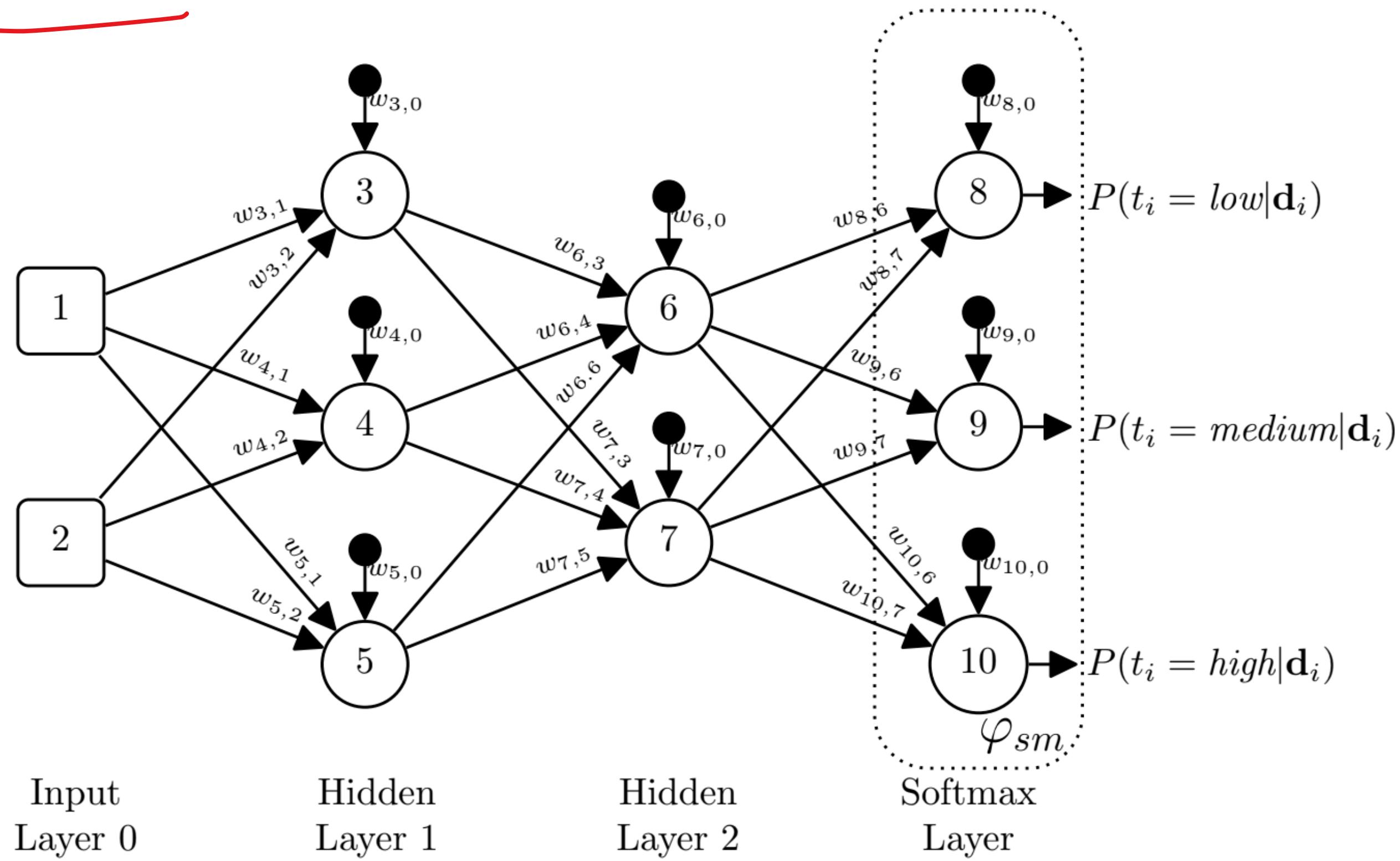
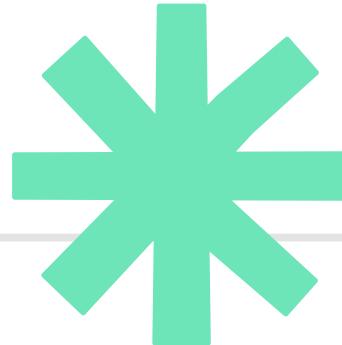


$w.$   $\delta$

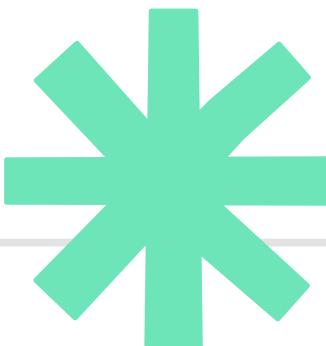
$b$   $y$

$l$   $(n-1)$

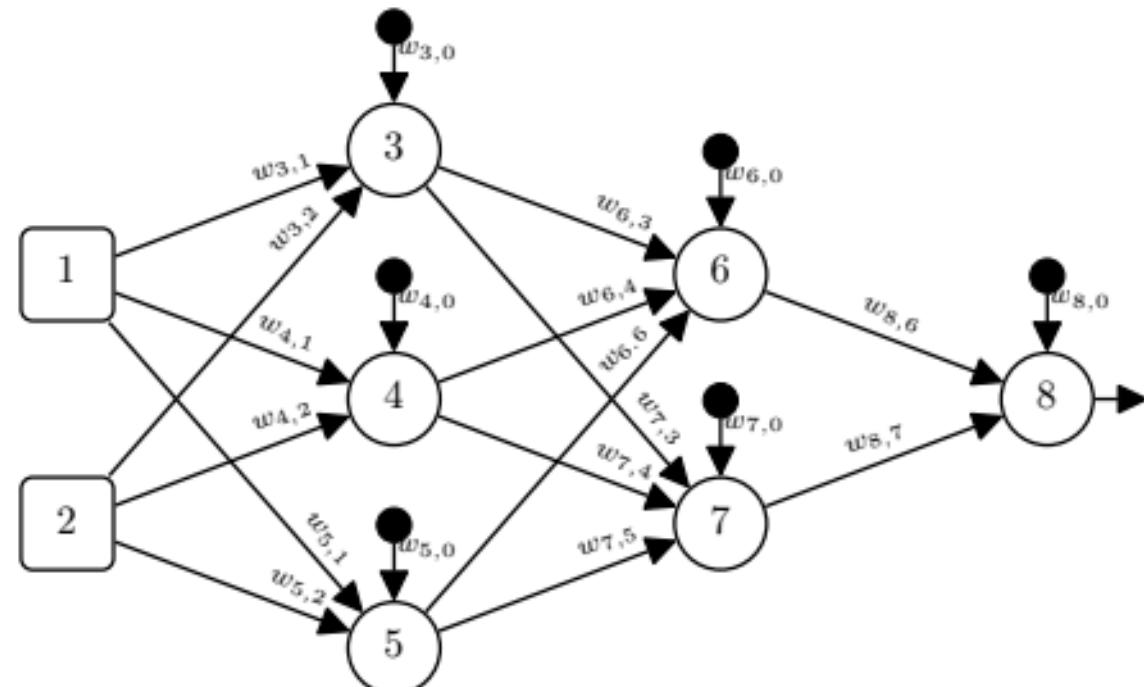
# Softmax Output Layer



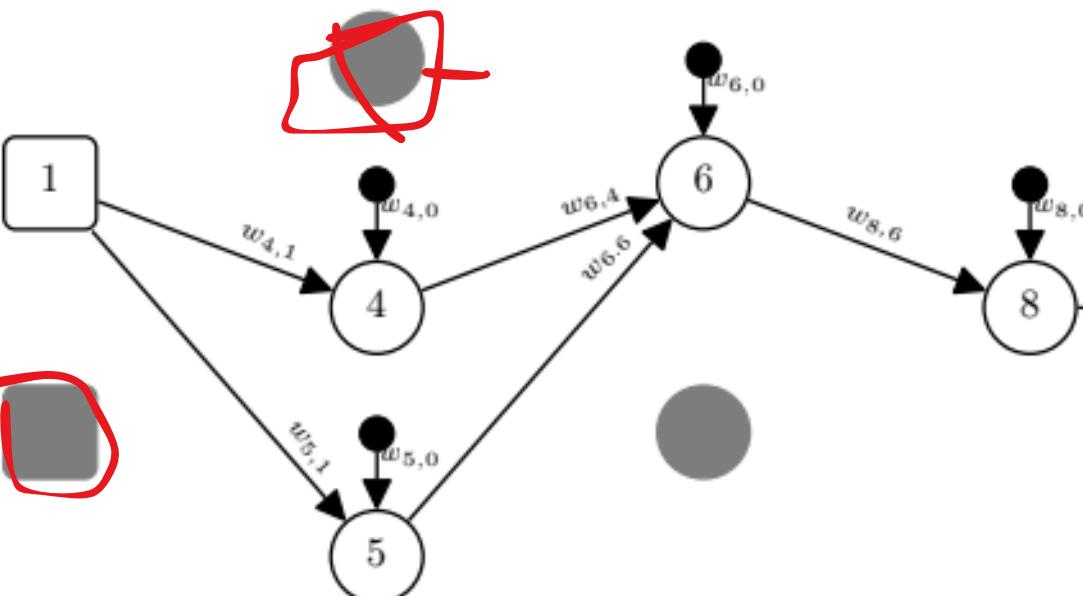
# How to prevent Overfitting?



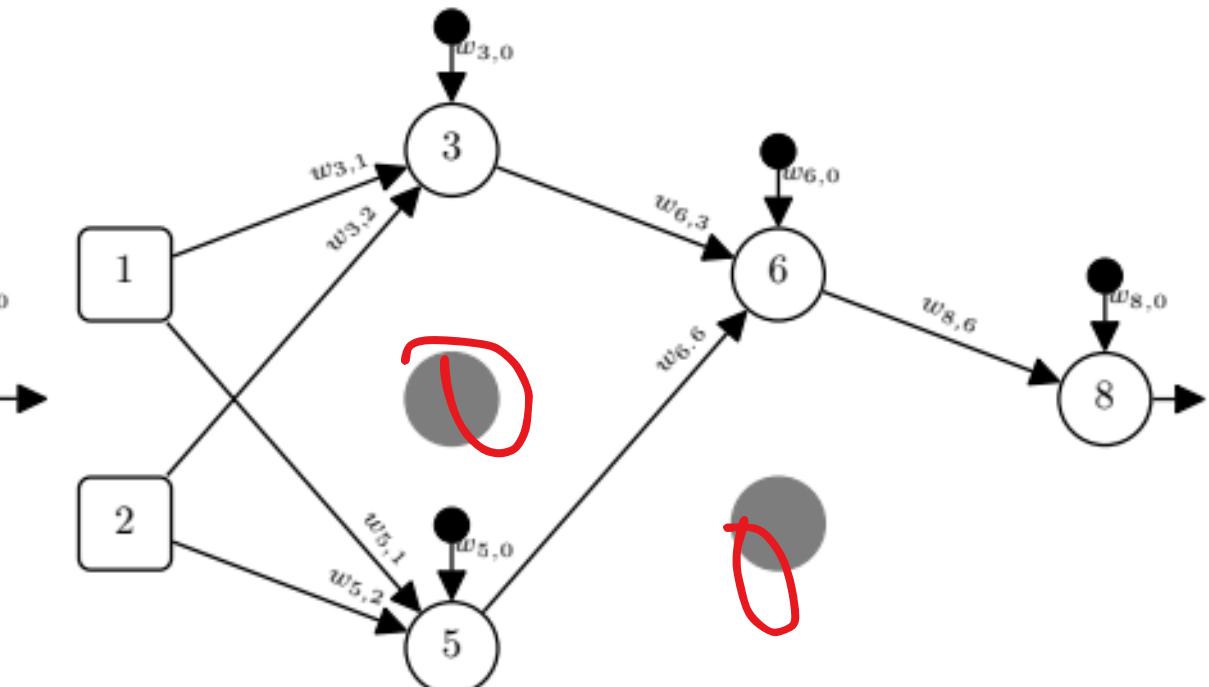
Dropout



No Dropout



Dropout for Example 1



Dropout for Example 2

# How to prevent Overfitting?



## Dropout

D | ↘

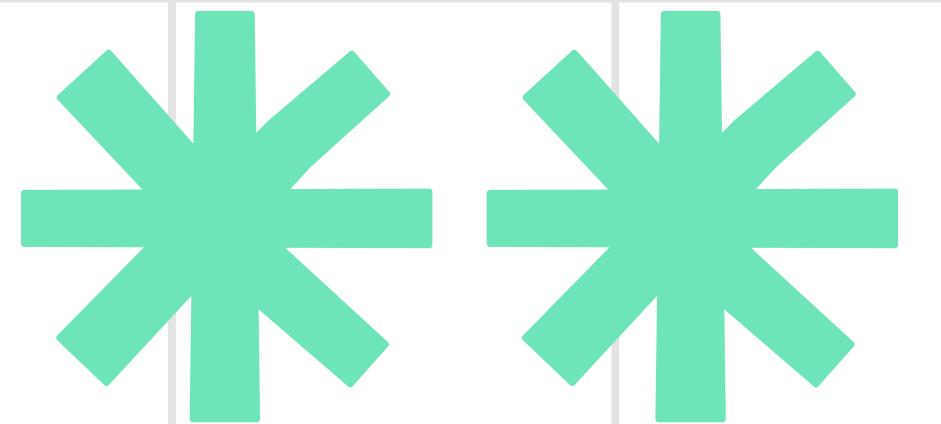
**Require:**  $\rho$  probability that a neuron in a layer will not be dropped

- 1: **for** each input or hidden layer  $l$  **do**
- 2:     $\text{DropMask}^{(l)} = \underbrace{(m_1, \dots, m_{\text{size}(l)})}_{\text{Bernoulli } (\rho)}$
- 3:     $\mathbf{a}^{(l)'} = \mathbf{a}^{(l)} \odot \text{DropMax}^{(l)}$
- 4:     $\mathbf{a}^{(l)''} = \underbrace{\frac{1}{\rho} \mathbf{a}^{(l)'}}_{\text{Forward Pass}}$
- 5: **end for**
- 6: **for** each layer  $l$  in backward pass **do**
- 7:     $\delta^{(l)} = \underbrace{\delta^{(l)} \odot \text{DropMax}^{(l)}}_{\text{Backward Pass}}$
- 8: **end for**



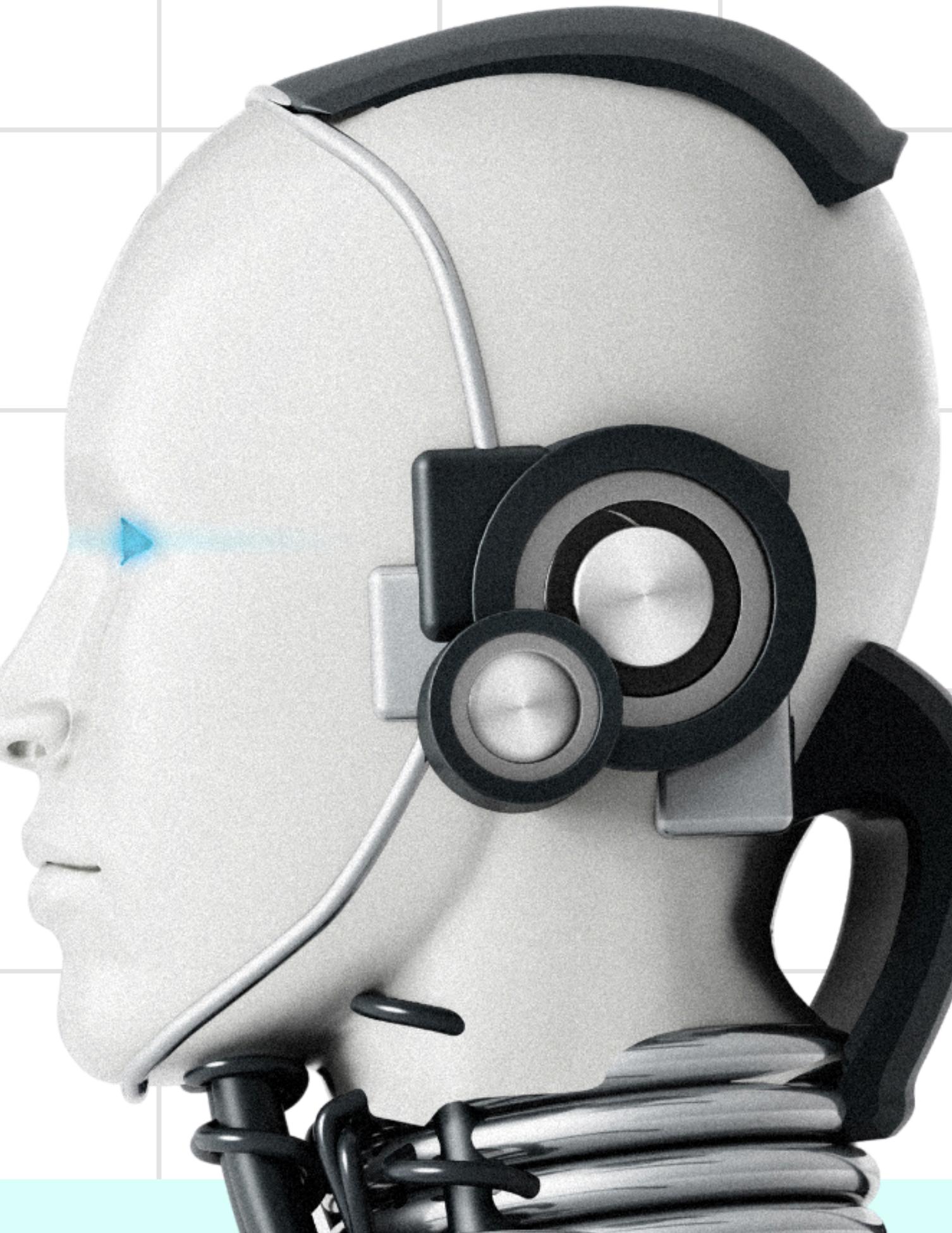
▷ Forward Pass

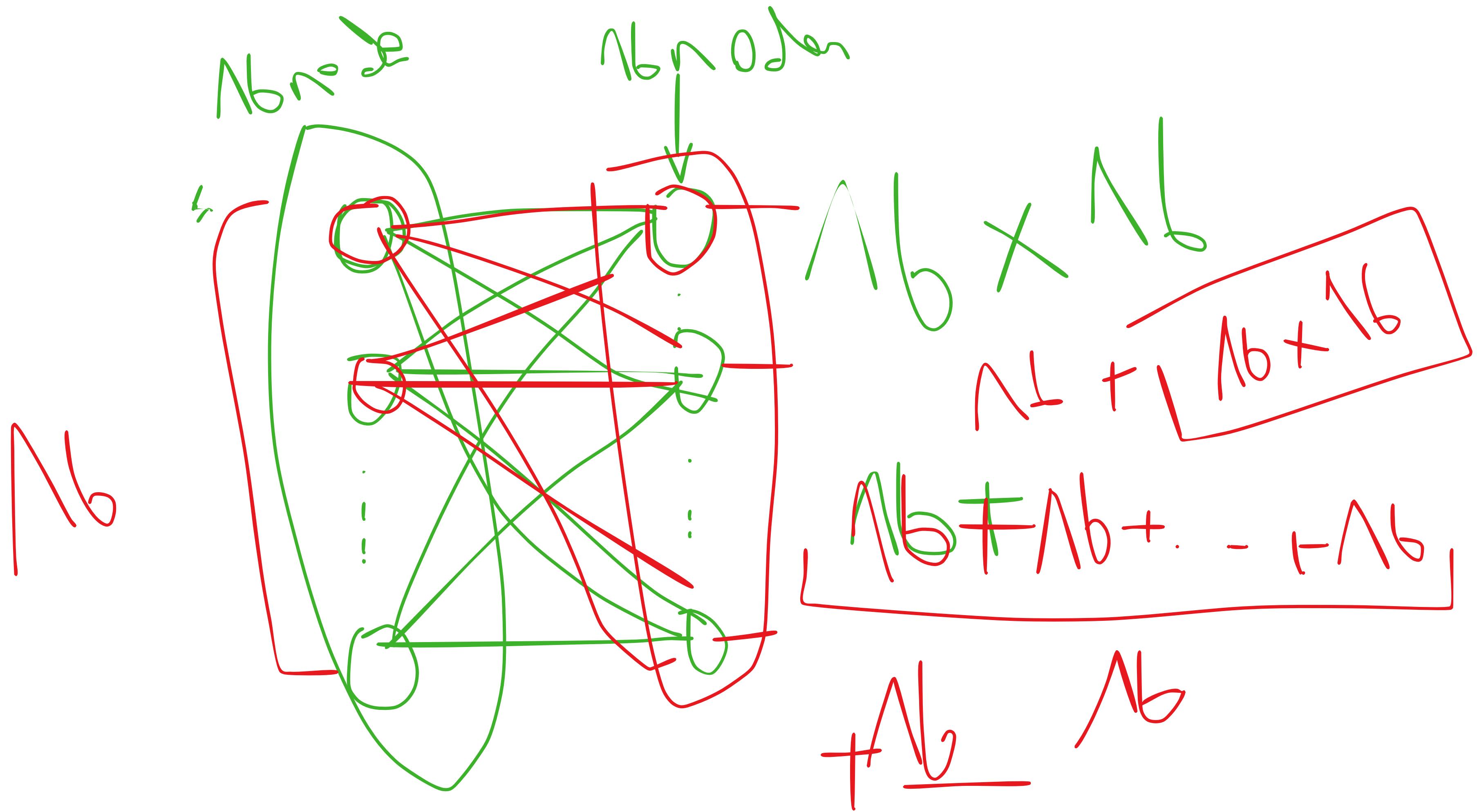
▷ Backward Pass

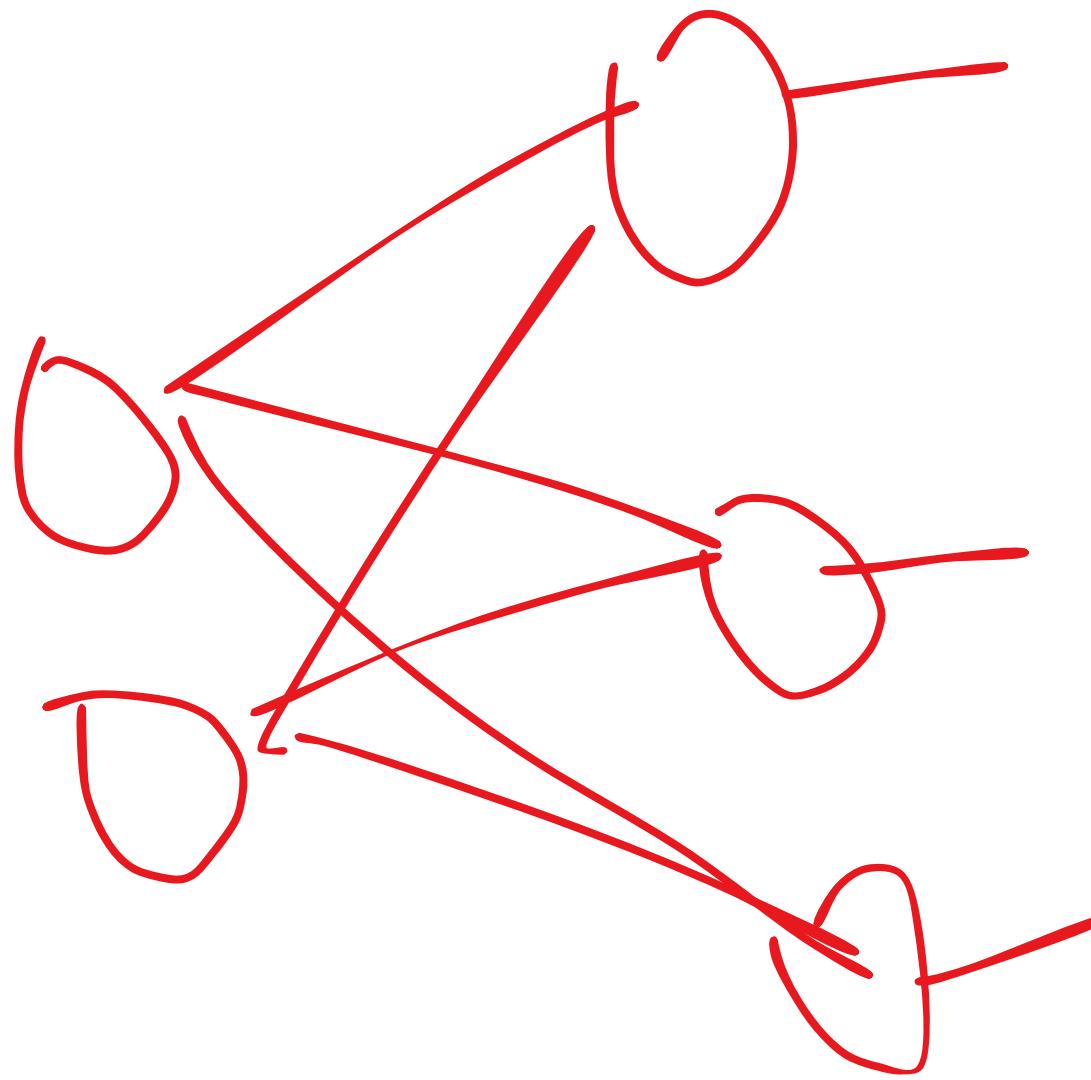


# Thank you!

I hope you  
enjoyed 

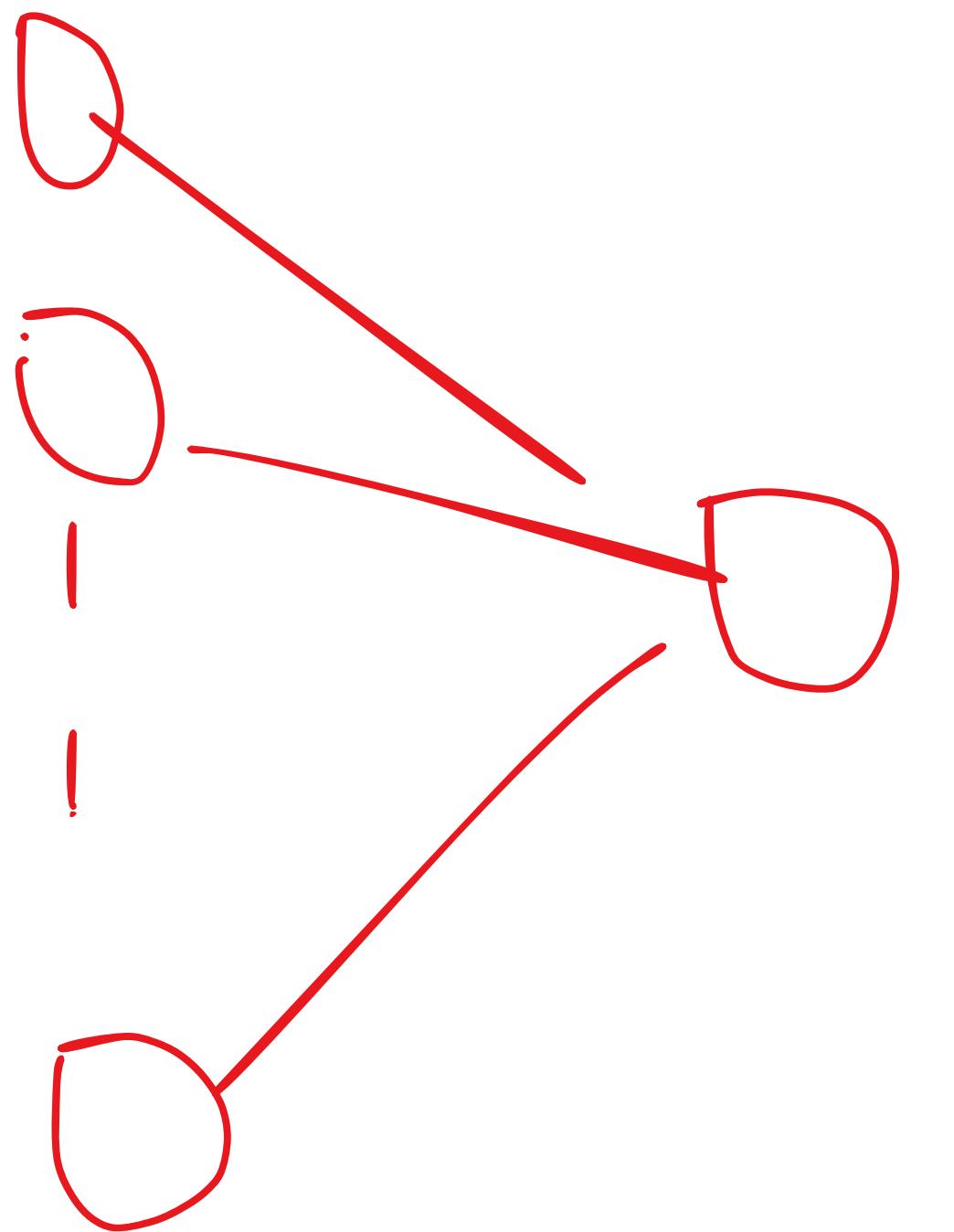






$$b + \overbrace{3} = \overbrace{c}$$

$$\overbrace{12} + \overbrace{3} = \overbrace{G}$$



Nb + -



$9 \times 9$

$9 \times 1$

$\sim P$

$O$

$\rightarrow D$

