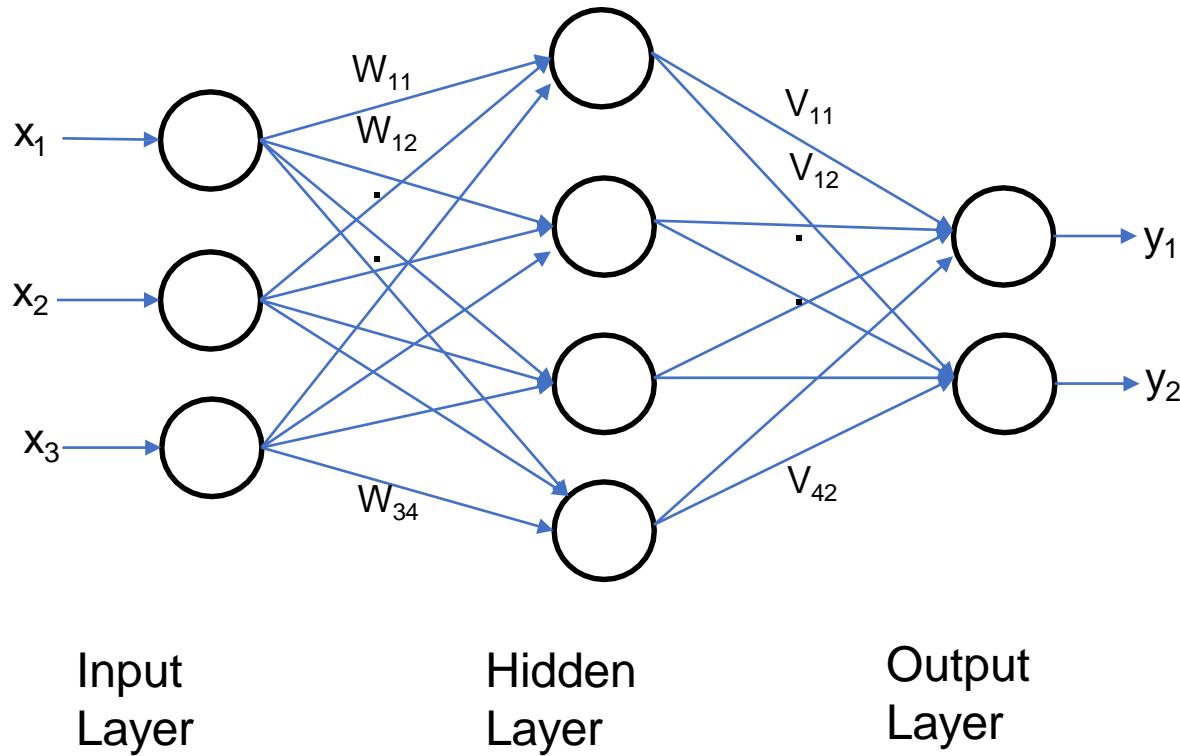


Computer Vision

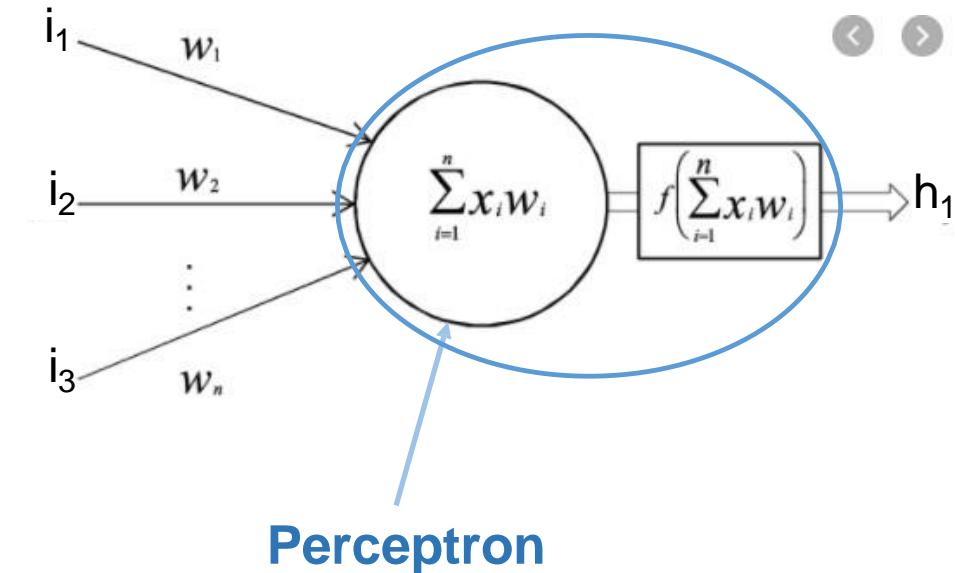
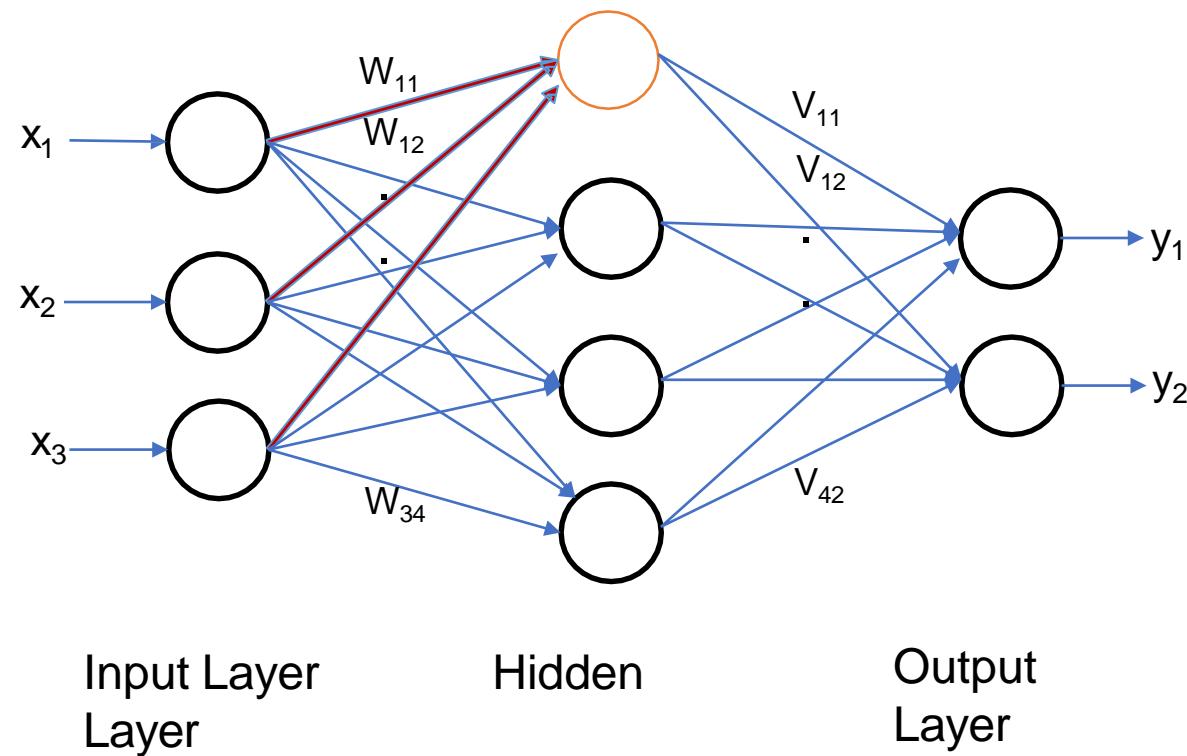
Chapter 5-3:

Multilayer Perceptron

What is Multilayer Perceptron NN?



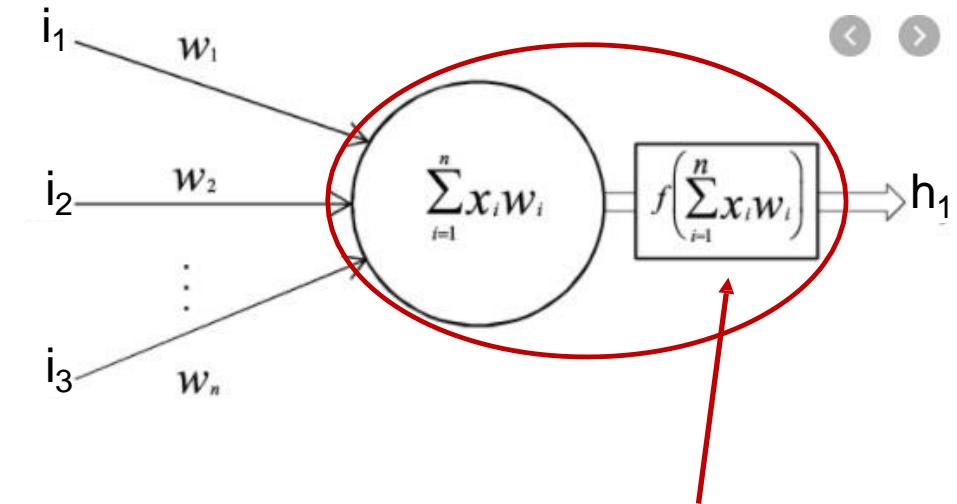
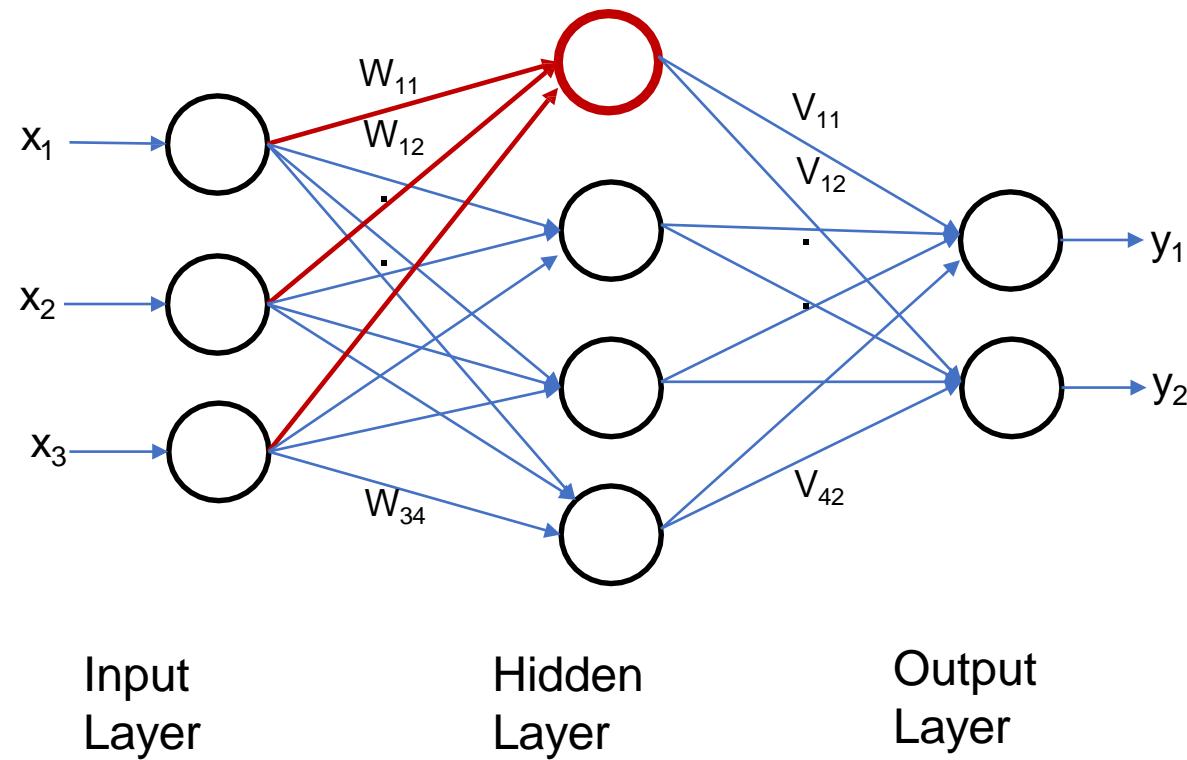
Multilayer Perceptron



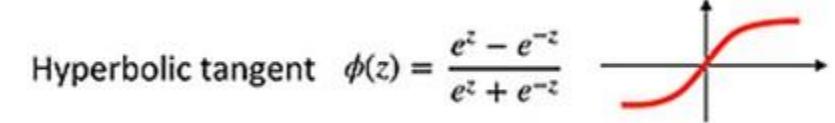
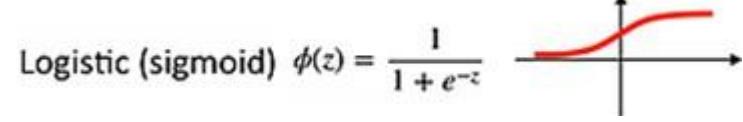
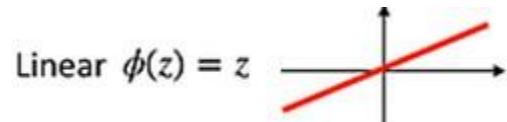
For a given node i , the perceptron is defined as weighted summation of the incoming data from the nodes of the previous layer.

$$p_i = x_0w_{0i} + x_1w_{1i} + \dots + x_nw_{ni} = \sum_{j=1}^n x_j w_j$$

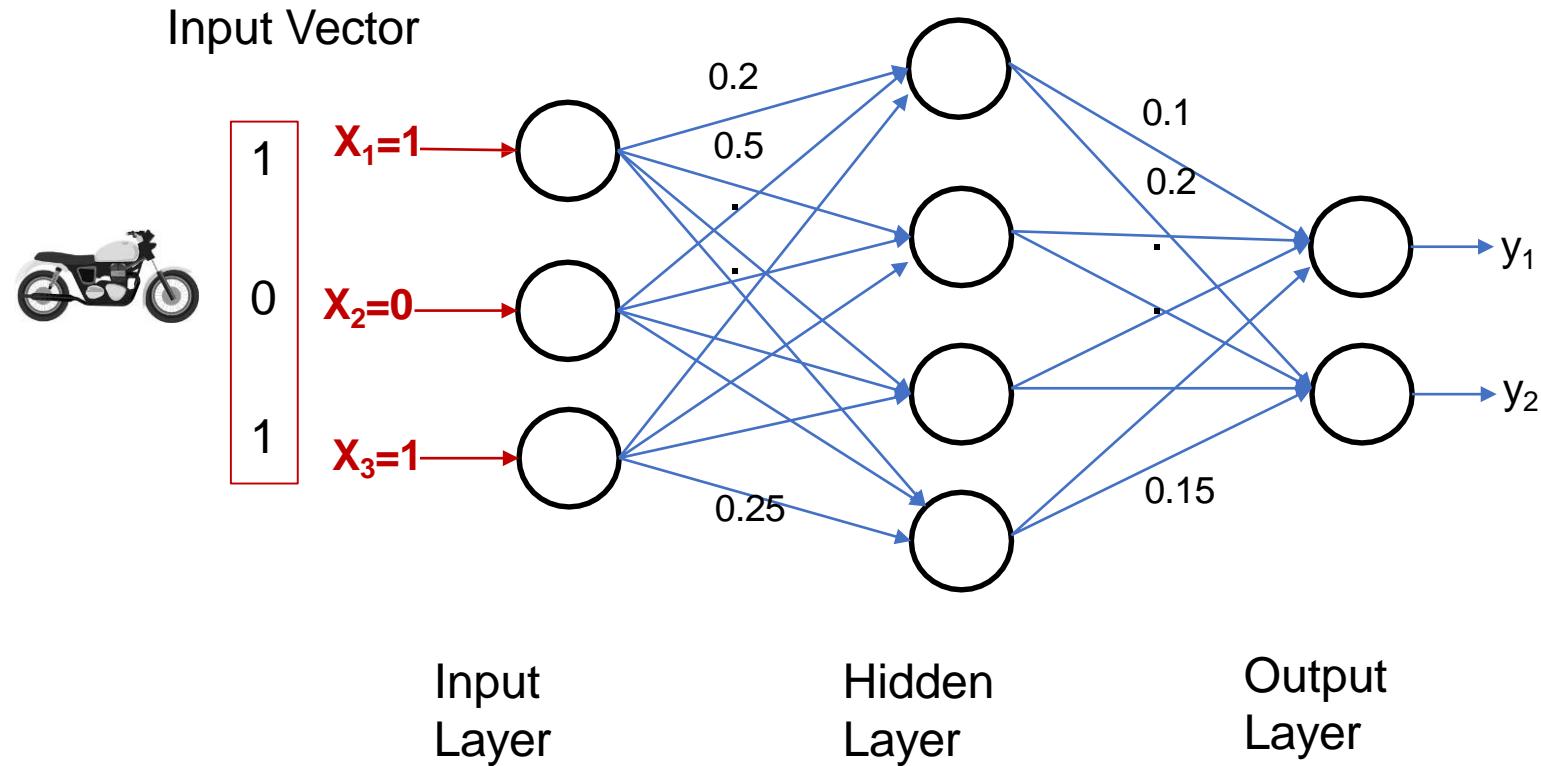
Multilayer Perceptron



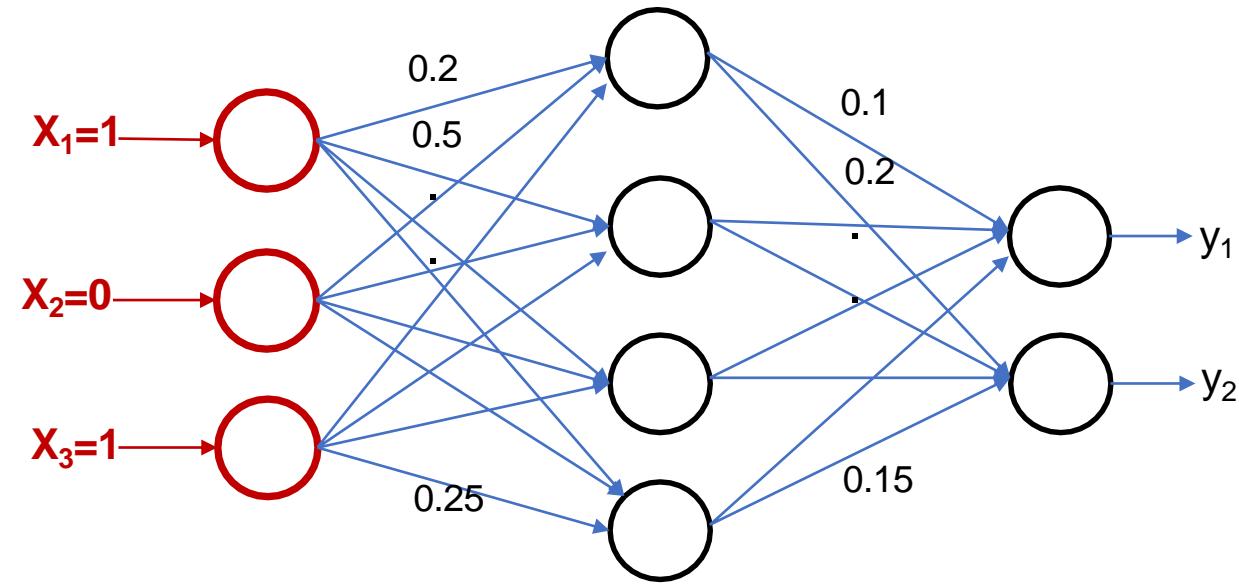
Activation Function



Forward Pass



Forward Pass



Input Layer

Hidden Layer

Output Layer

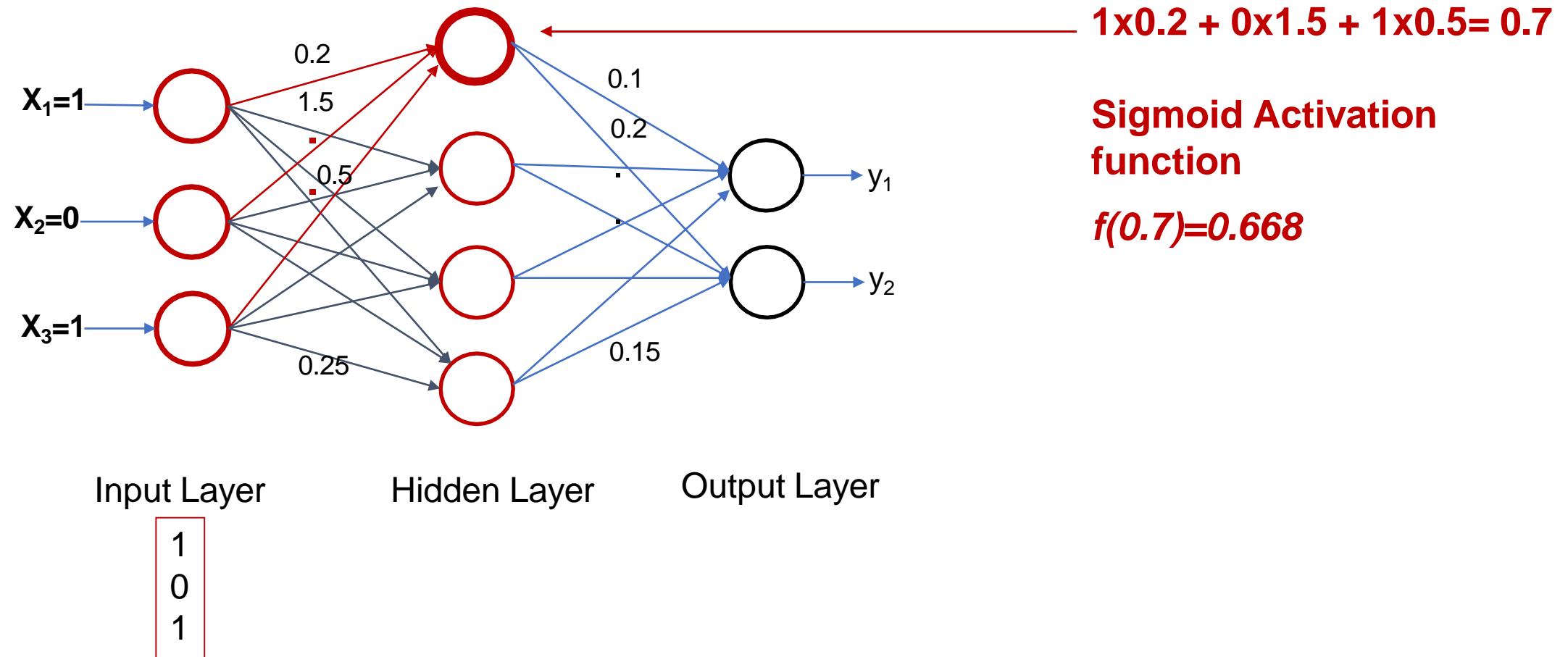
1
0
1

$x = \text{Perceptron}(x)$

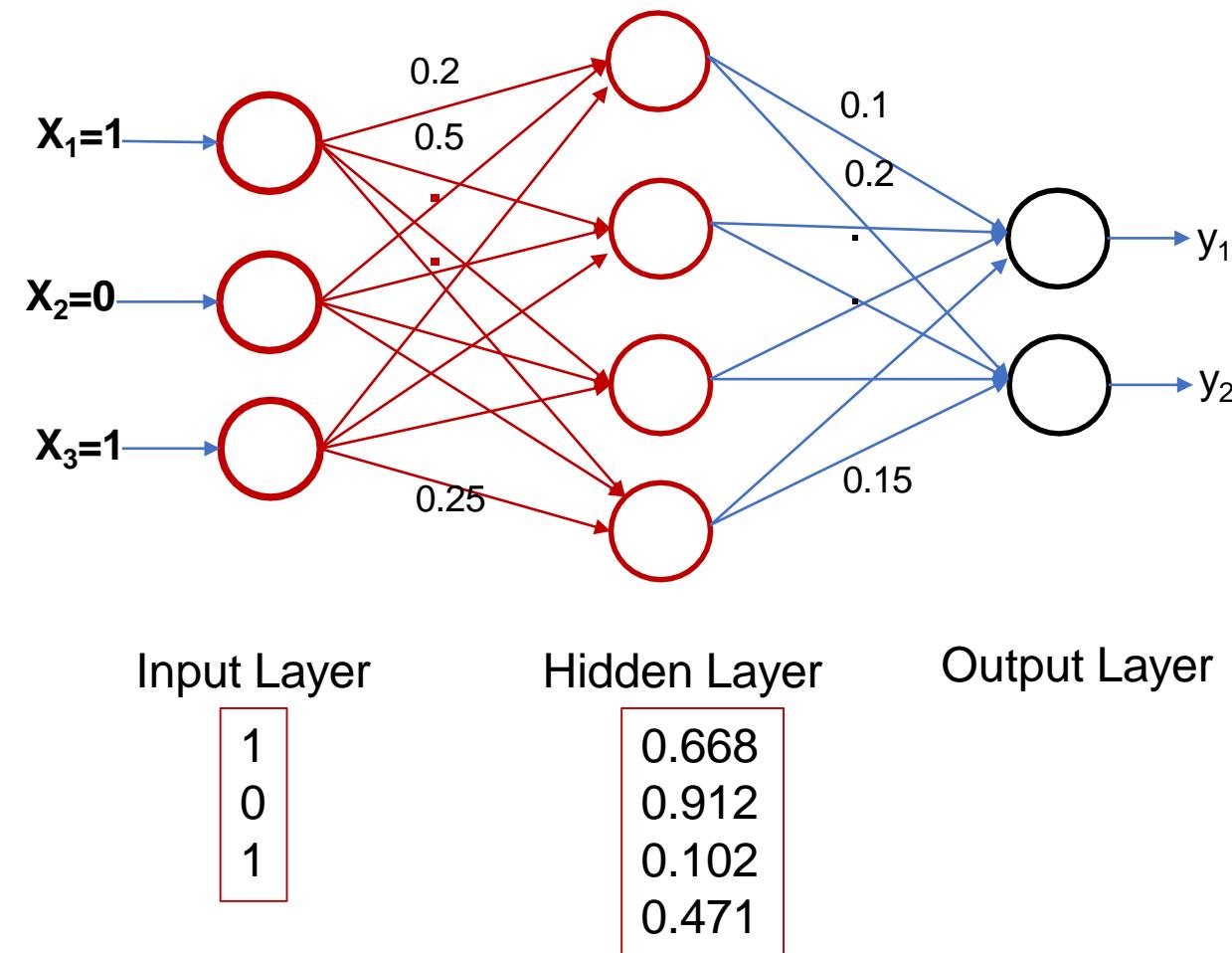
Linear Activation function

$x = f(x)$

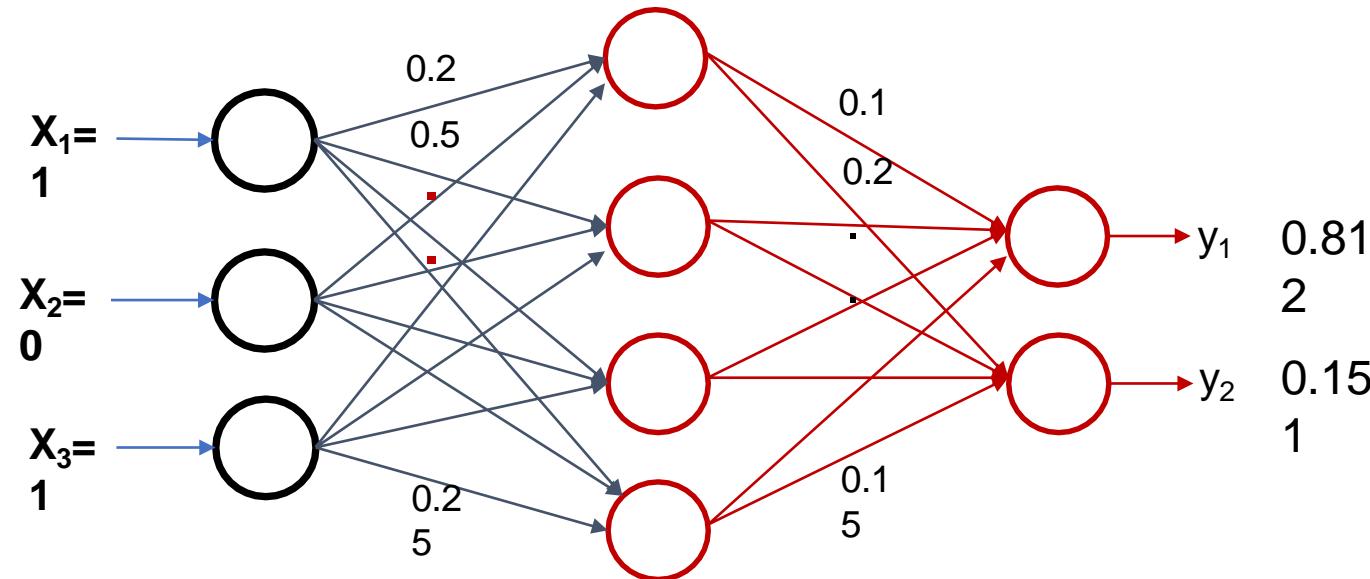
Forward Pass



Forward Pass



Forward Pass



Input
Layer

0
1

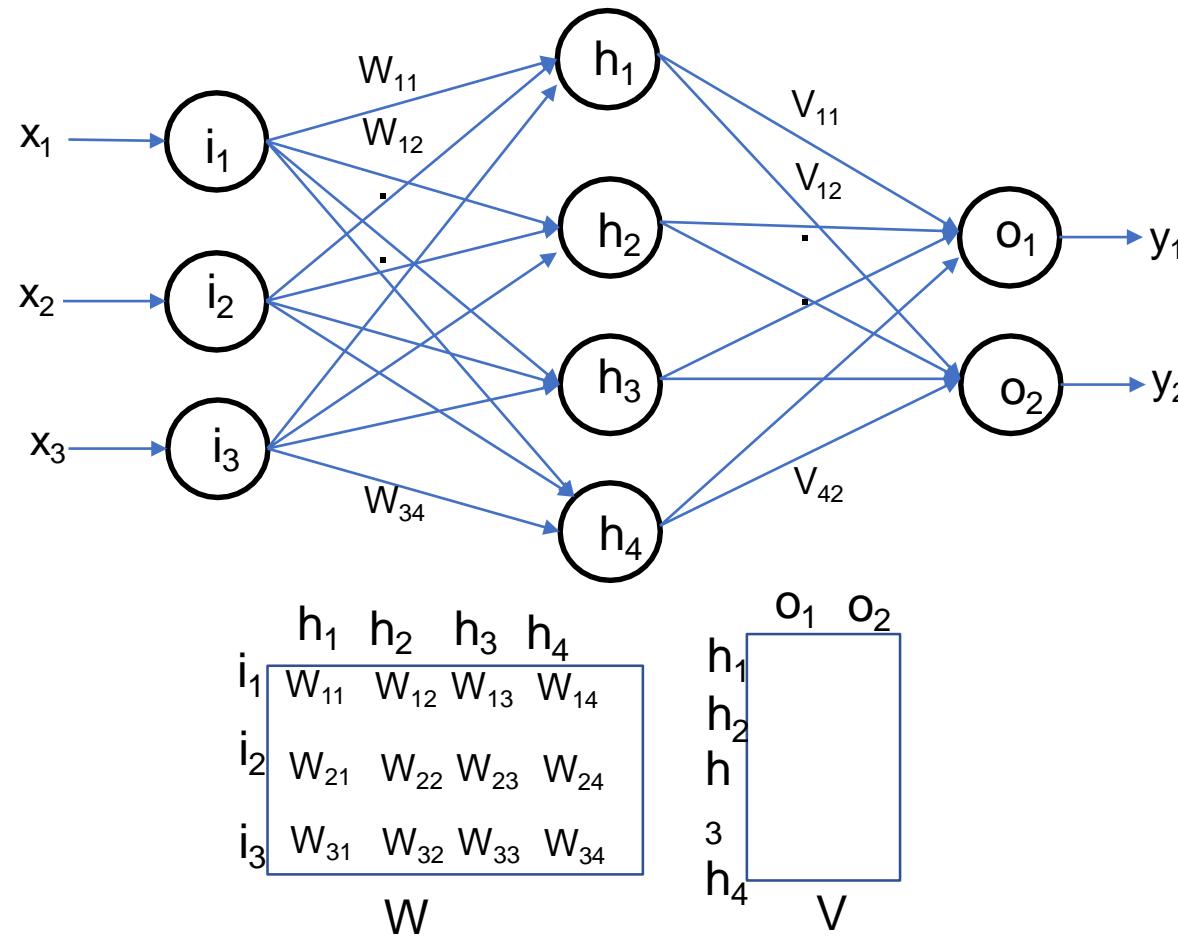
Hidden
Layer

0.668
0.912
0.102
0.471

Output
Layer

0.812
0.151

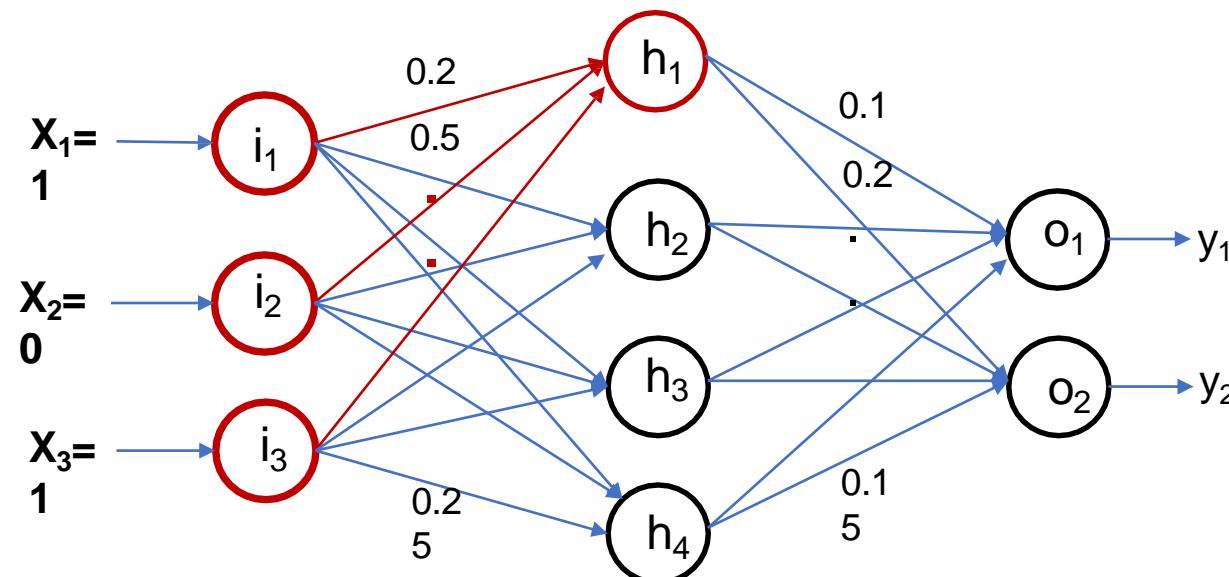
Weight Matrix



Forward Pass

$$p_1 = x_1 W_{11} + x_2 W_{21} + x_3 W_{31}$$

$$h_1 = \text{Sigmoid}(p_1)$$

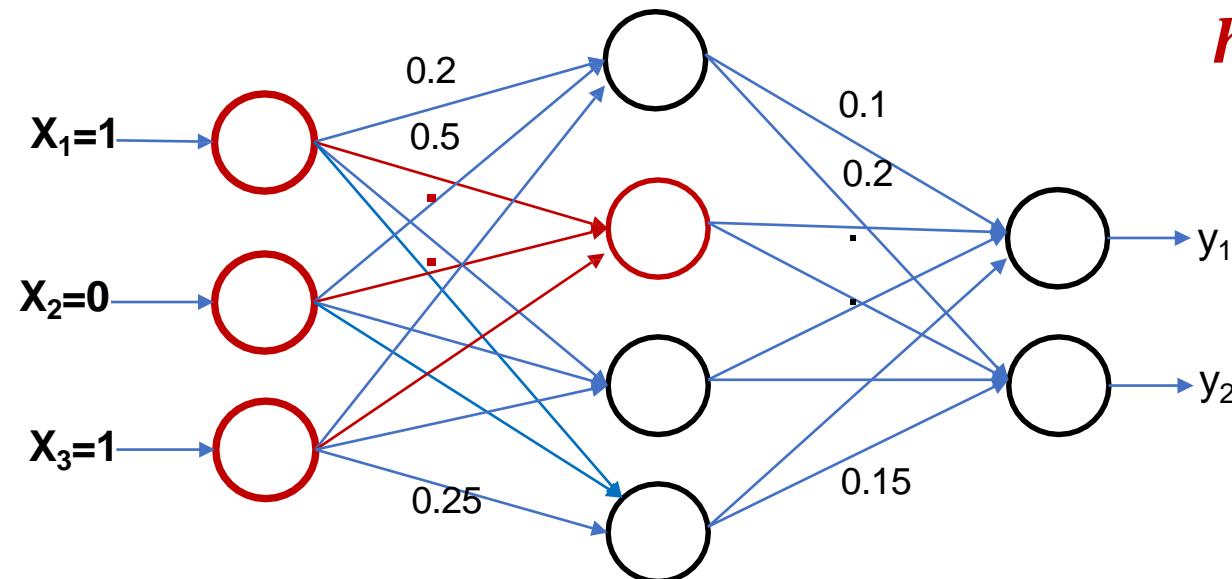


$$\begin{bmatrix} 1 & 0 & 1 \end{bmatrix} \times \begin{bmatrix} h_1 & h_2 & h_3 & h_4 \\ i_1 & & & \\ i_2 & & & \\ i_3 & & & \end{bmatrix} = \boxed{0.668}$$

Forward Pass

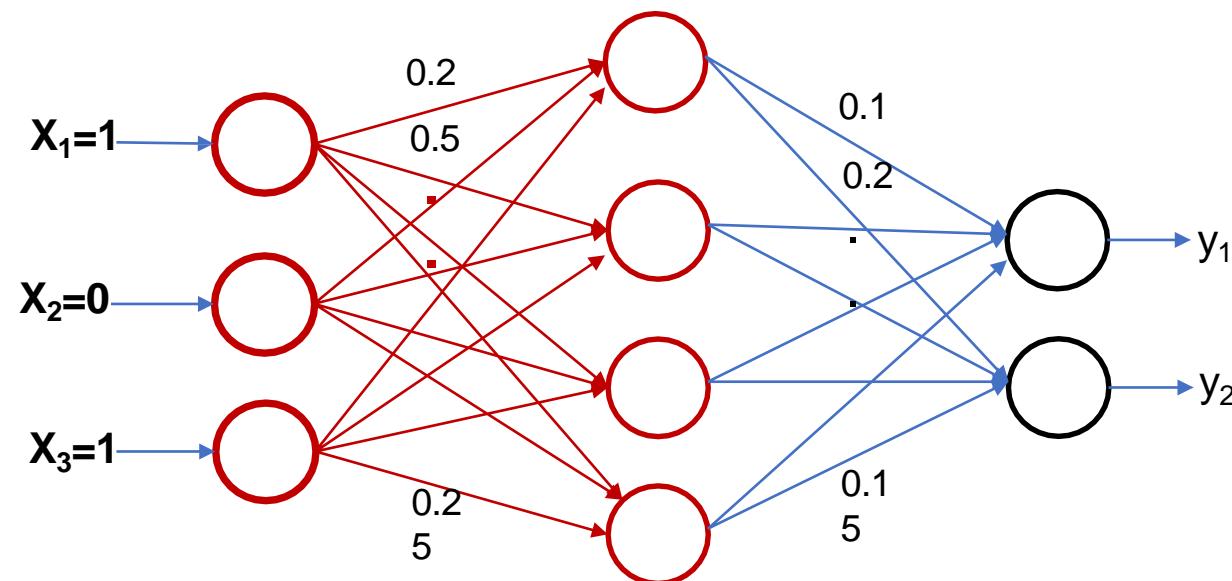
$$p_2 = x_1 W_{12} + x_2 W_{22} + x_3 W_{32}$$

$$h_2 = \text{Sigmoid}(p_2)$$



$$\begin{bmatrix} 1 & 0 & 1 \end{bmatrix} \times \begin{bmatrix} i_1 & h_1 & h_2 & h_3 & h_4 \\ i_2 & & & & \\ i_3 & & & & \end{bmatrix}_W = \begin{bmatrix} 0.668 & \boxed{0.912} \end{bmatrix}$$

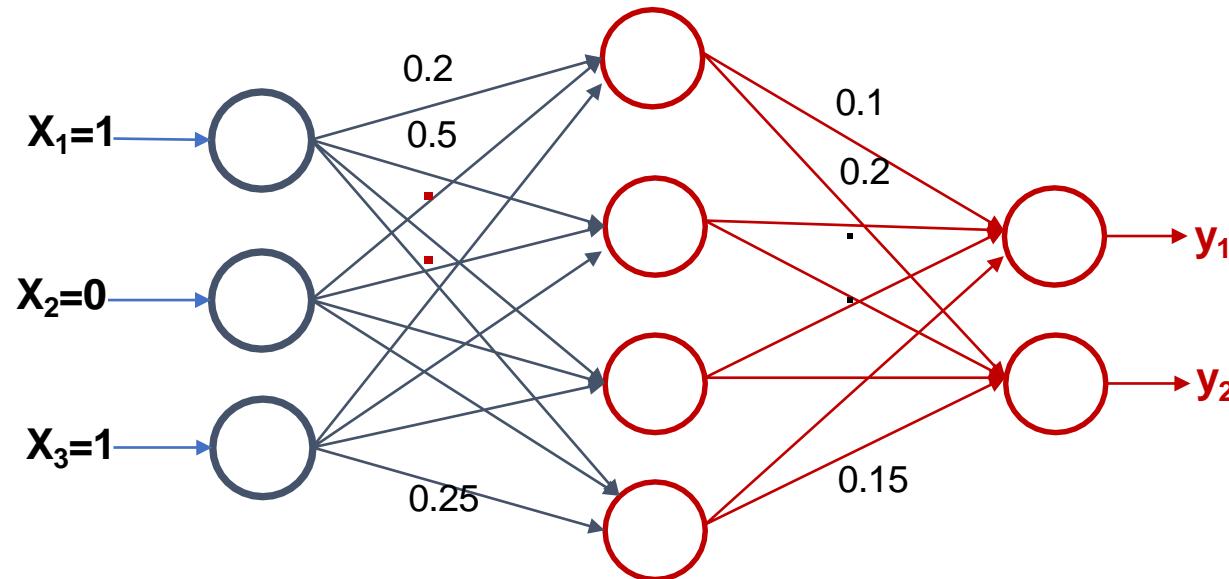
Forward Pass



$$\bar{h}^T = \text{Sigmoid}(\bar{p}^T)$$

$$\begin{bmatrix} 1 & 0 & 1 \\ \bar{x}^T \end{bmatrix} \times \begin{bmatrix} h_1 & h_2 & h_3 & h_4 \\ i_1 & & & \\ i_2 & & & \\ i_3 & & & \end{bmatrix} = \begin{bmatrix} h_1 & h_2 & h_3 & h_4 \\ 0.668 & 0.912 & 0.102 & 0.471 \\ \bar{h}^T \end{bmatrix}$$

Forward Pass

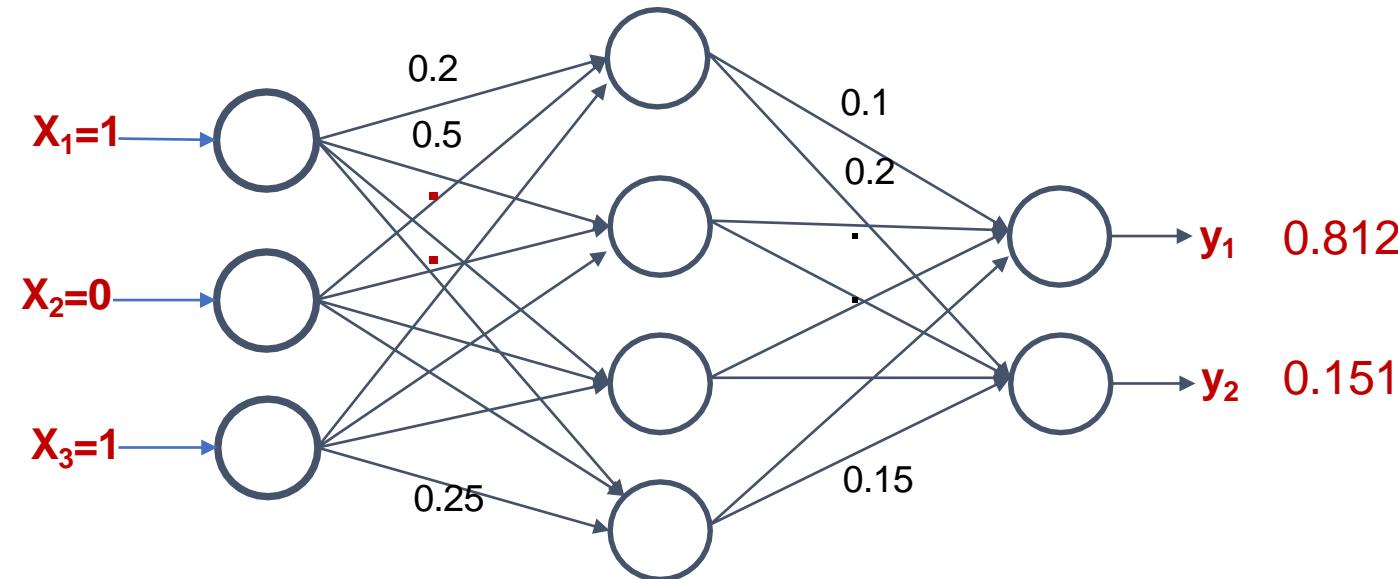


$$\bar{p}^T = \bar{h}^T \cdot V$$

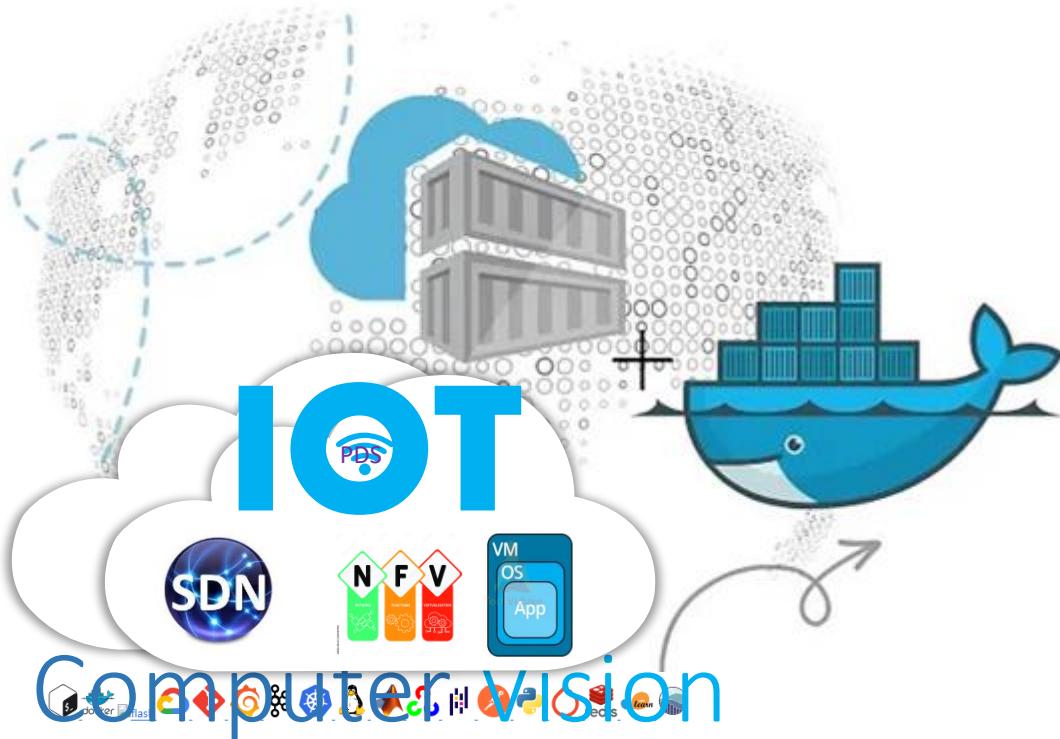
$$\bar{y}^T = \text{Sigmoid}(\bar{p}^T)$$

$$\begin{array}{cccc}
 0.668 & 0.912 & 0.102 & 0.471 \\
 \bar{h}^T
 \end{array}
 X
 \begin{array}{c}
 \begin{matrix} o_1 & o_2 \end{matrix} \\
 \begin{matrix} h_1 & \\ h_2 & \\ h_3 & \\ h_4 & \end{matrix}
 \end{array}
 V =
 \begin{array}{cc}
 0.812 & 0.151 \\
 \bar{y}^T
 \end{array}$$

Forward Pass



$$\bar{y}^T = \text{Sigmoid}(\text{Sigmoid}(\bar{x}^T W)^T V)$$



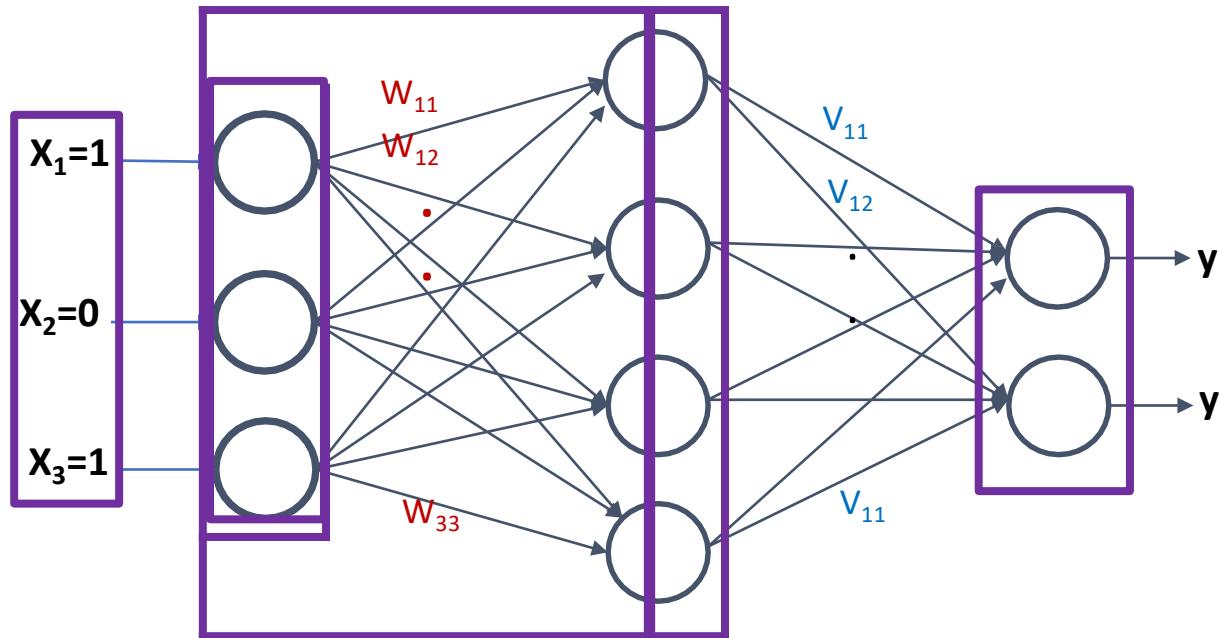
Computer Vision

Chapter 5-4:

Learning the parameters : Backpropagation through time

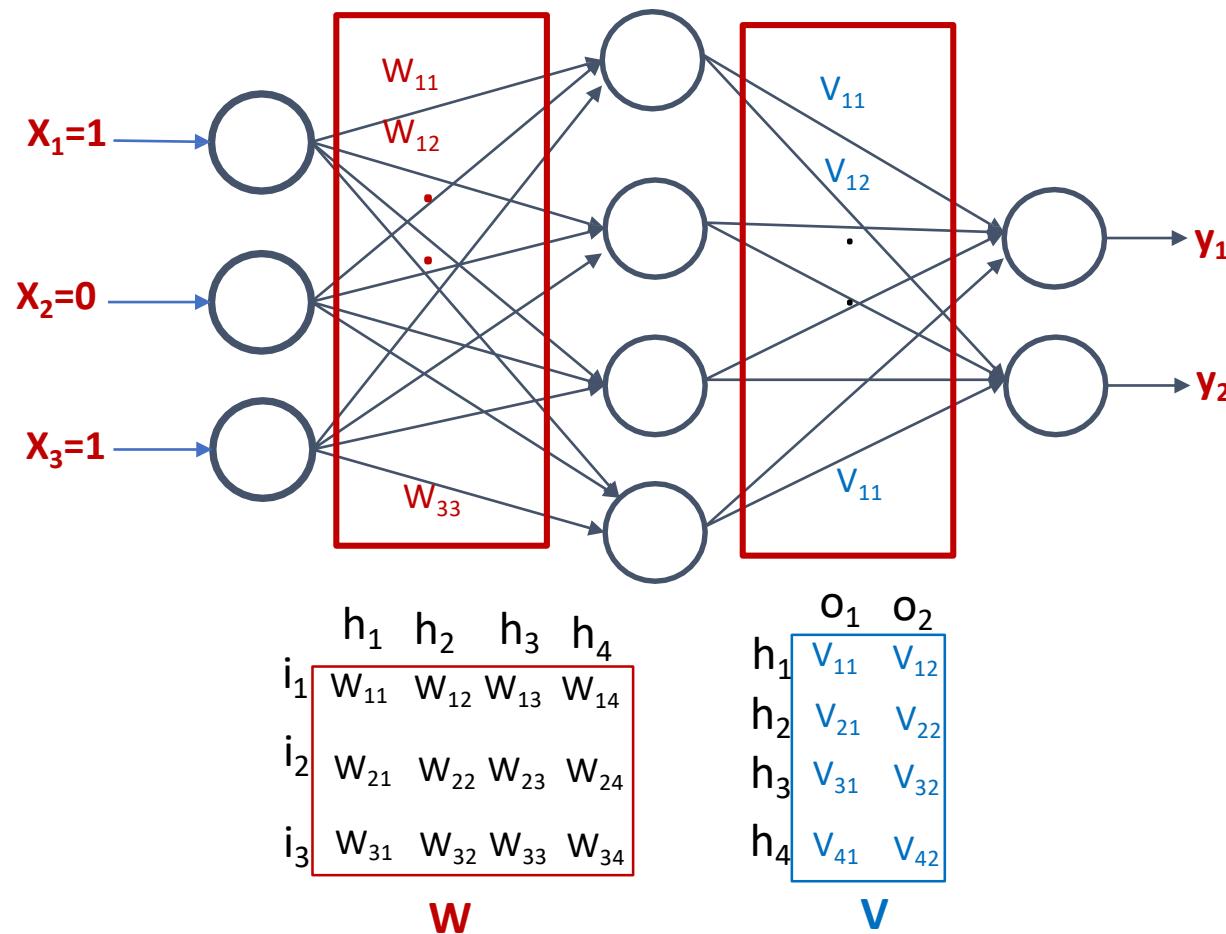
- What is loss function?
- What are the parameters of a multilayer perceptron neural network?
- How to estimate parameters using backpropagation through time?

What are the parameters?



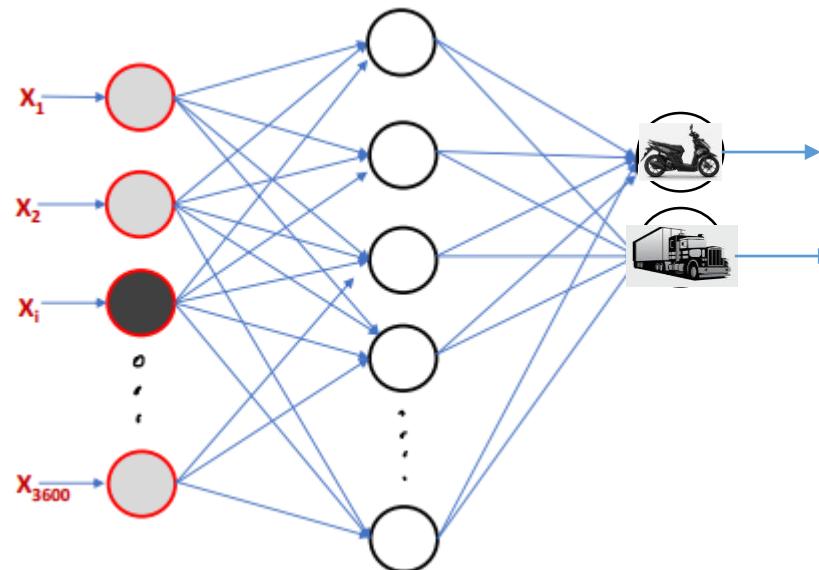
$$\bar{y}^T = f_{hidden}(f_{output}(\bar{x}^T \bar{W})^T \cdot \bar{V})$$

What are the parameters?



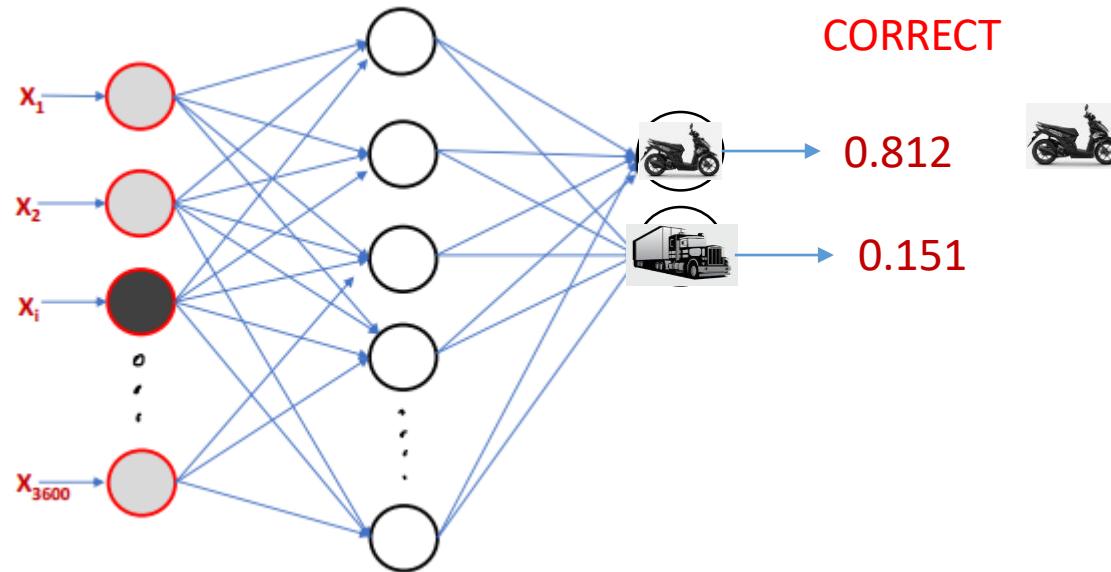


60x60=3600



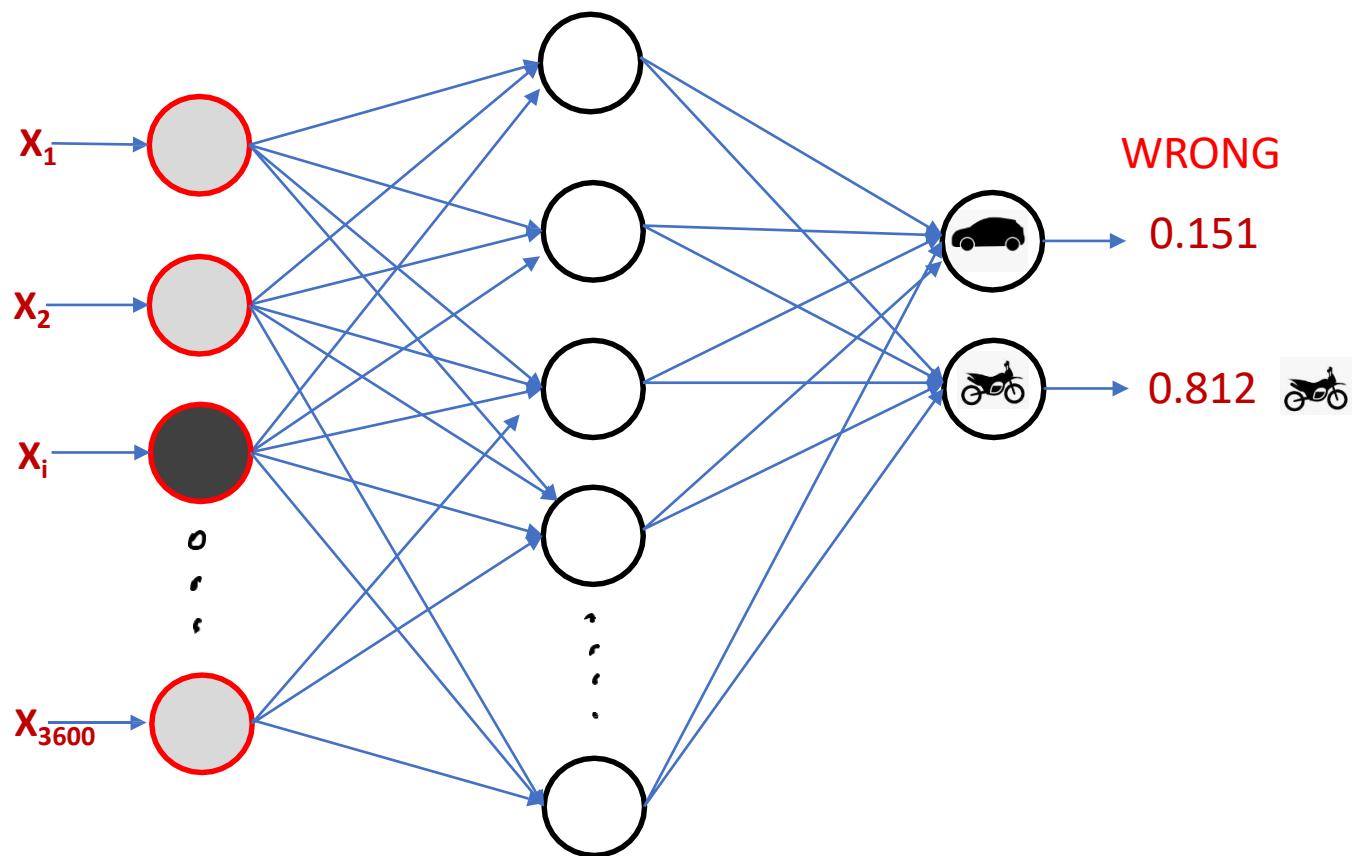


60x60=3600

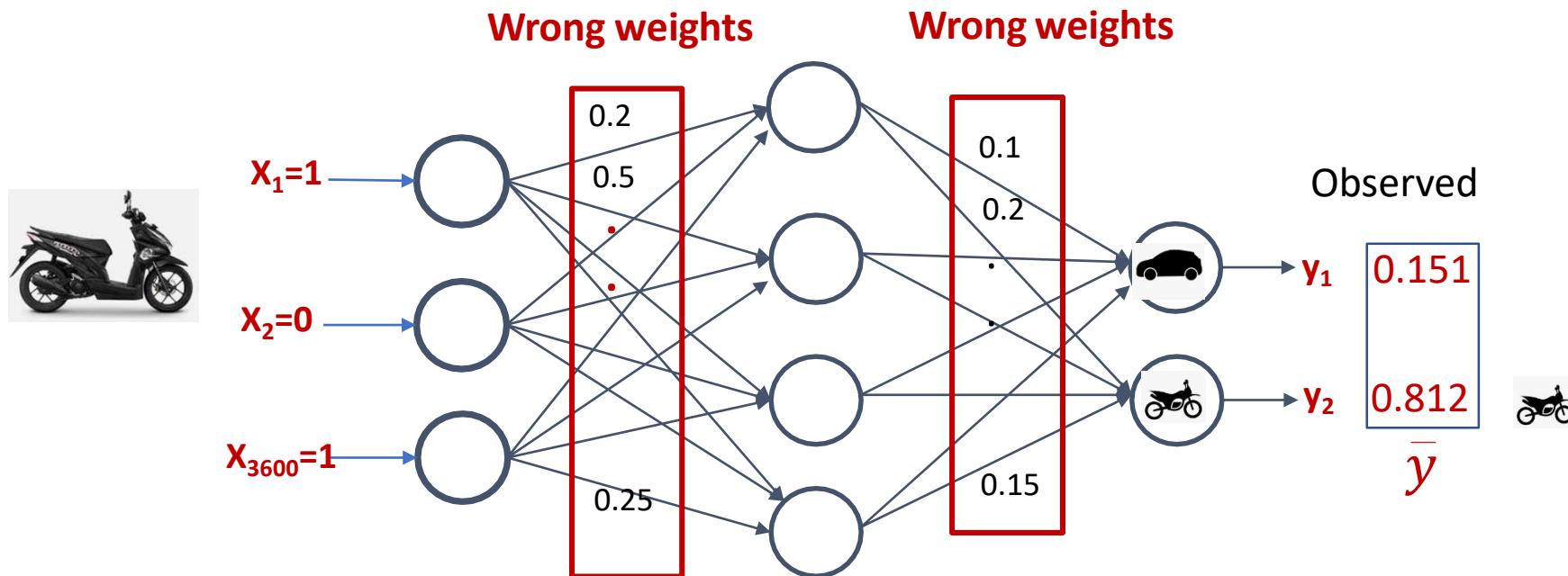




60x60=3600

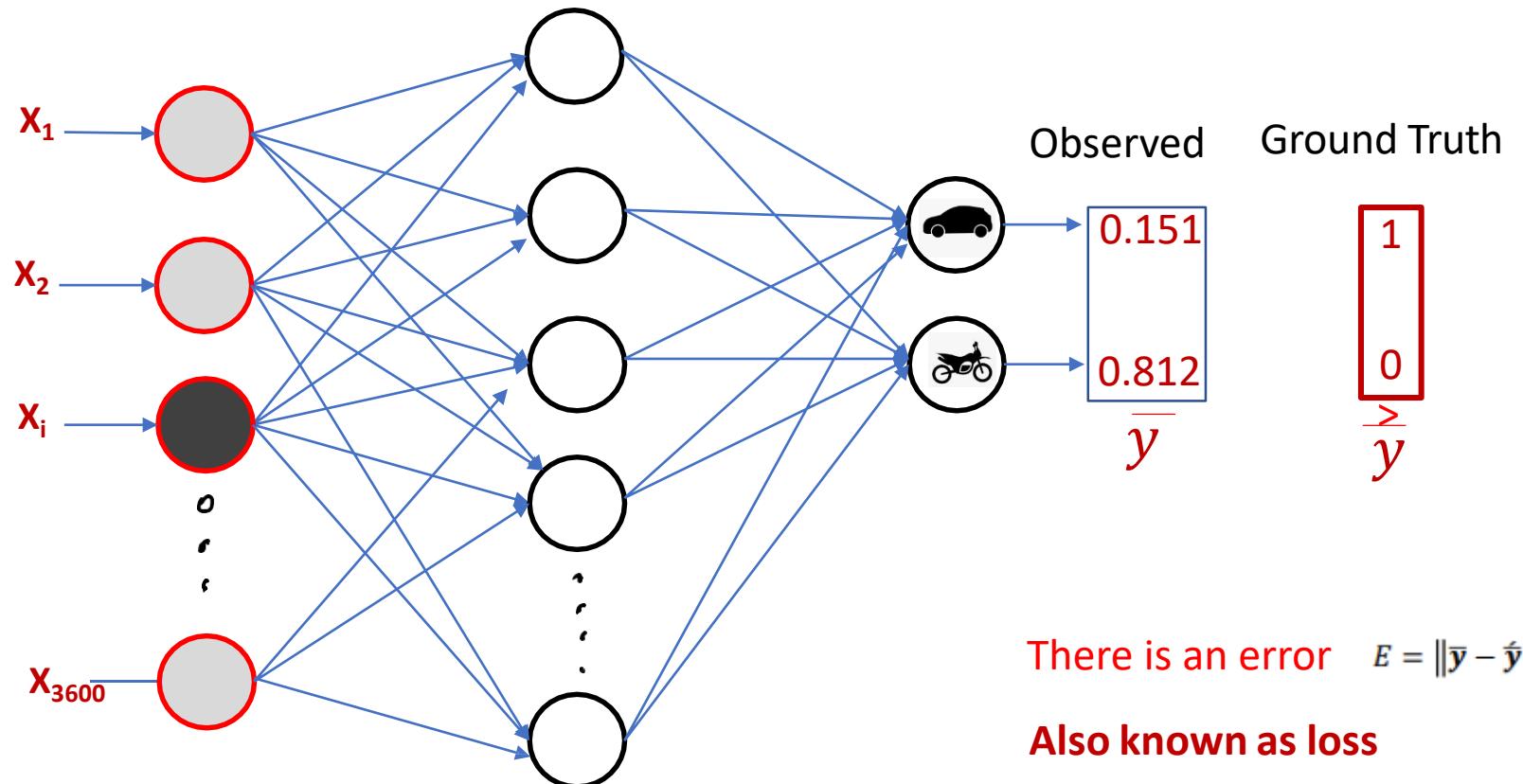


Why there is an error in prediction?





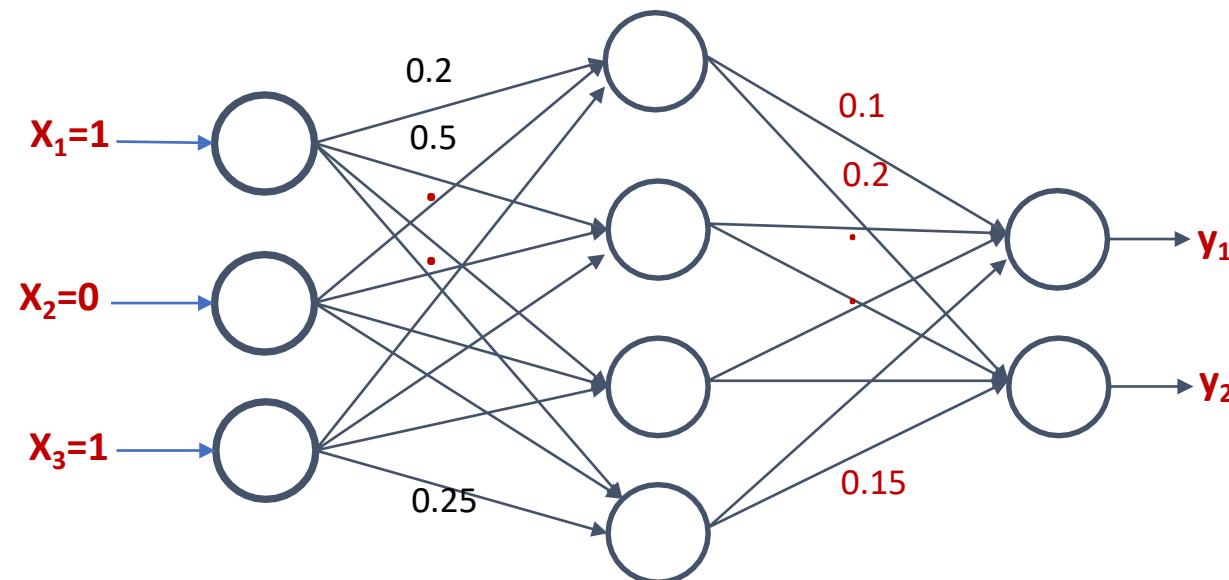
$60 \times 60 = 3600$



There is an error $E = \|\bar{y} - \hat{y}\|$

Also known as loss

Backpropagation Through Time



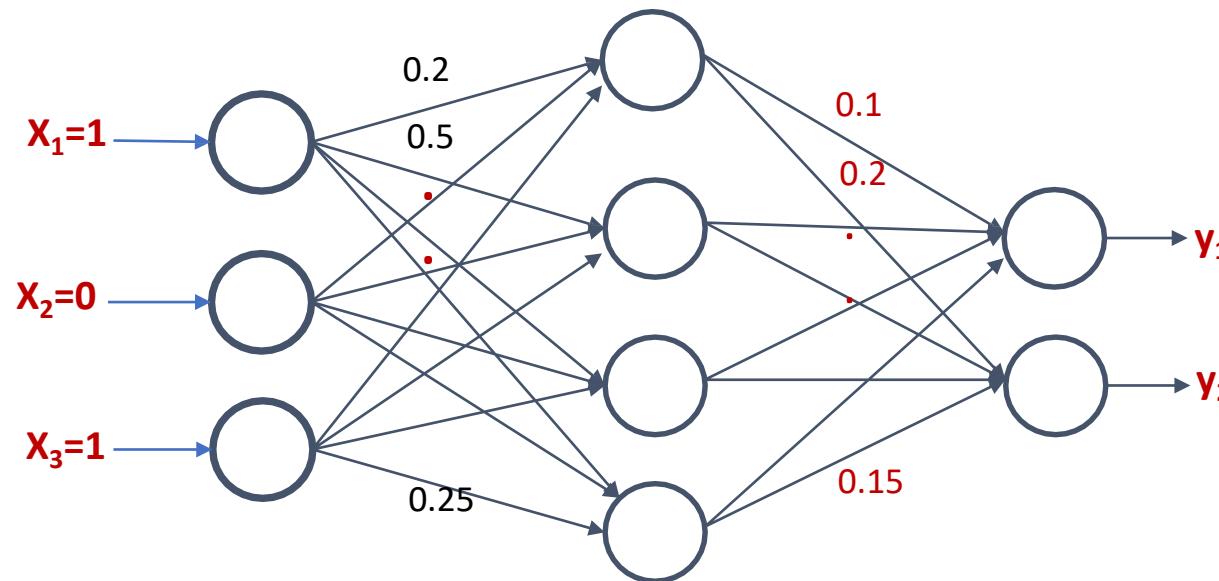
Minimize the Loss function

$$E = \|\bar{y} - \hat{y}\|$$

Backpropagation Through Time

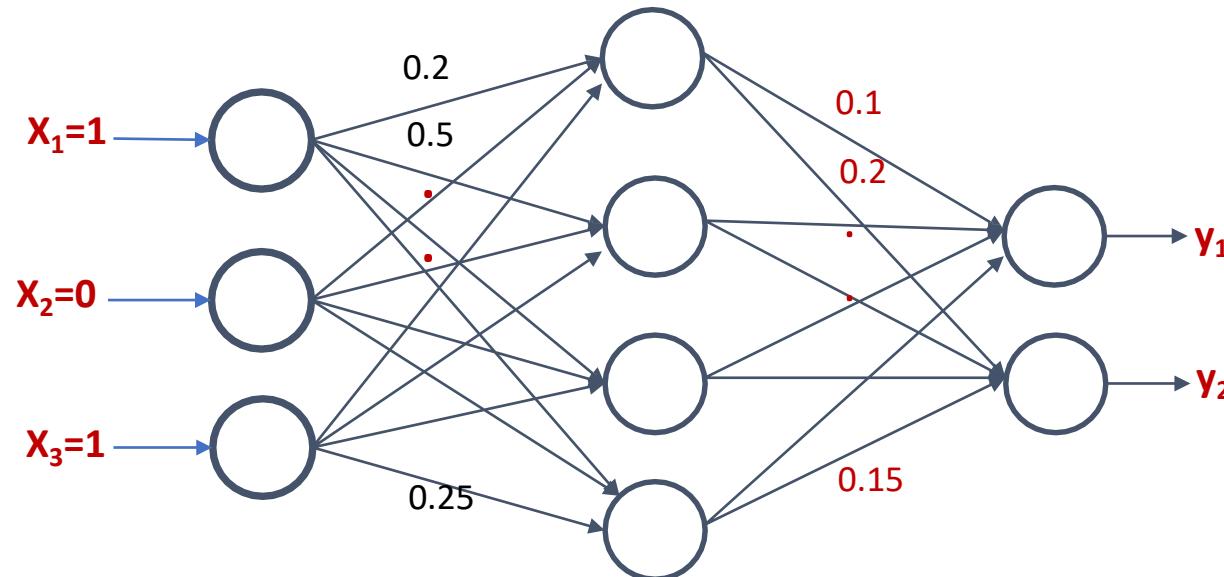
Minimize the Loss function

$$E = \|\bar{y} - \hat{y}\|$$



Backpropagation Through Time

Minimize the Loss function $E = \|\bar{y} - \hat{y}\|$



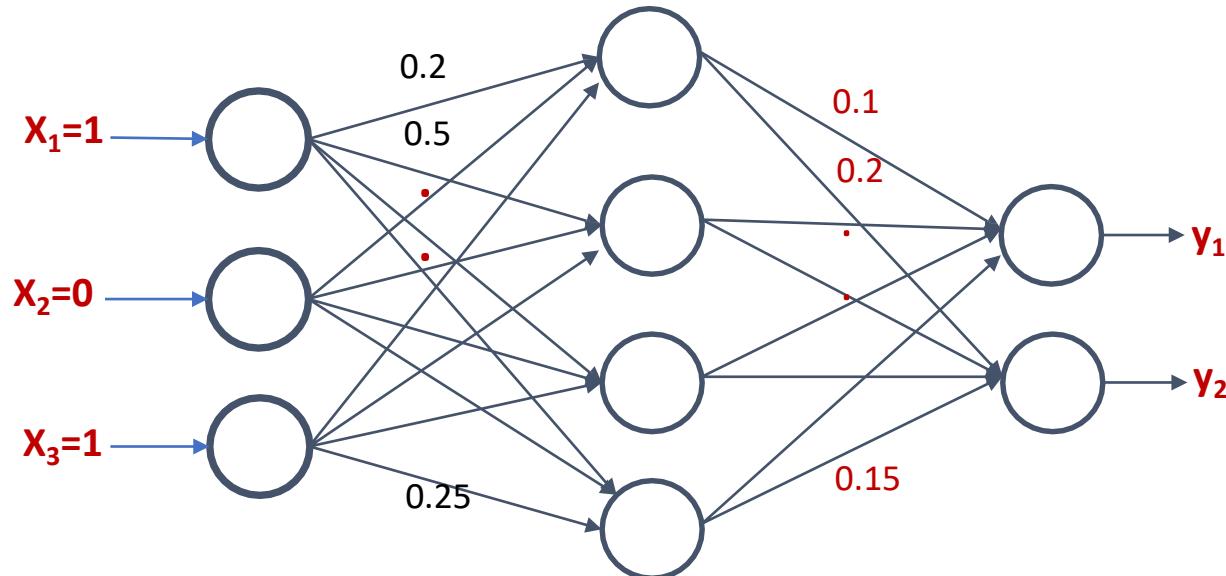
$$\nabla = \frac{\delta E}{\delta V} = 0$$

$$\nabla_{ij} = \frac{\delta E}{\delta V_{ij}} = \frac{\delta \|\bar{y} - \hat{y}\|}{\delta V_{ij}} = 0$$

Backpropagation Through Time

Minimize the Loss function

$$E = \|\bar{y} - \hat{y}\|$$



$$\nabla = \frac{\delta E}{\delta V} = 0$$

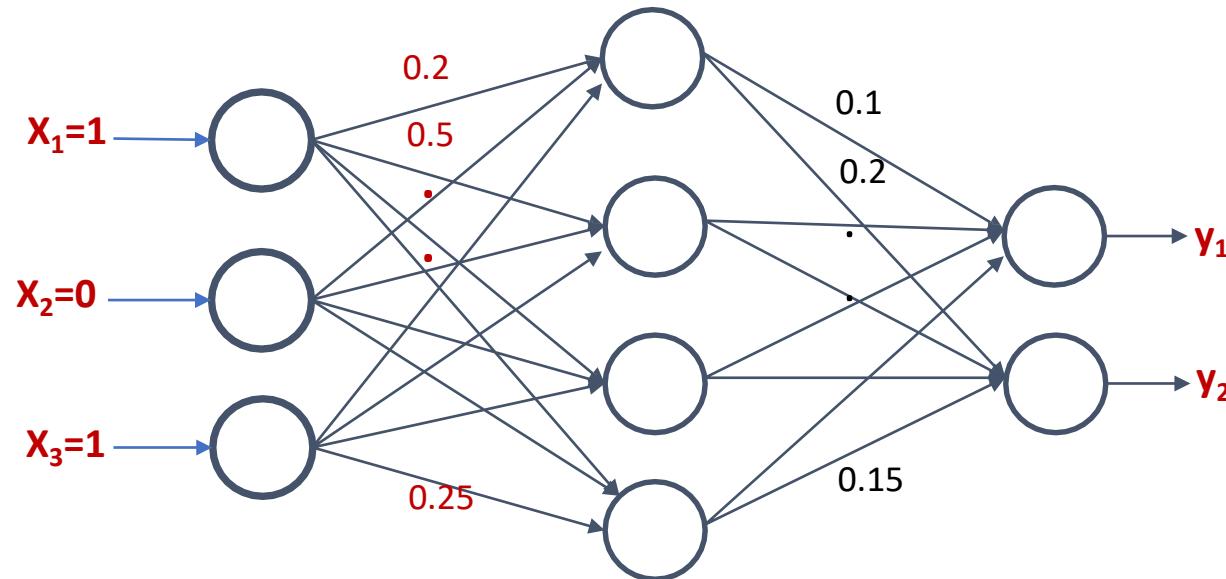
$$\nabla_{ij} = \frac{\delta E}{\delta V_{ij}} = \frac{\delta \|\bar{y} - \hat{y}\|}{\delta V_{ij}} = 0$$

$$\nabla = \frac{\delta E}{\delta V} = 0$$

$$V_{ij}^t = V_{ij}^{t-1} + \eta \nabla_{ij}^t$$

$$\eta = [0, 1]$$

Backpropagation Through Time

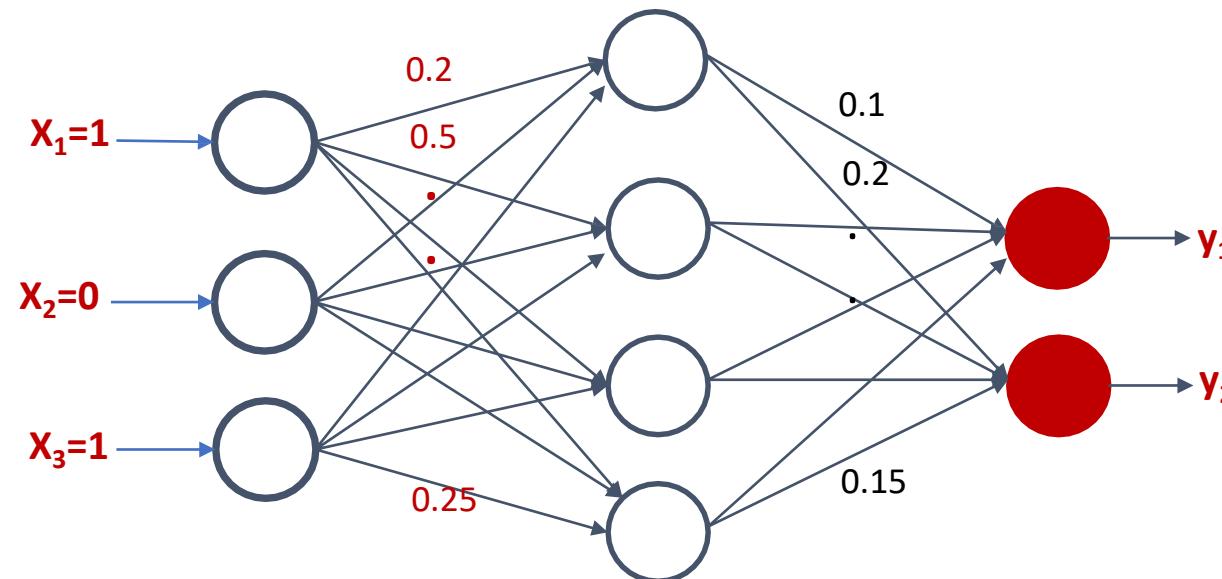


$$E = \|\bar{y} - \hat{y}\| = 0$$

or

$$E = \|\bar{y} - \hat{y}\| \leq \epsilon$$

Backpropagation Through Time

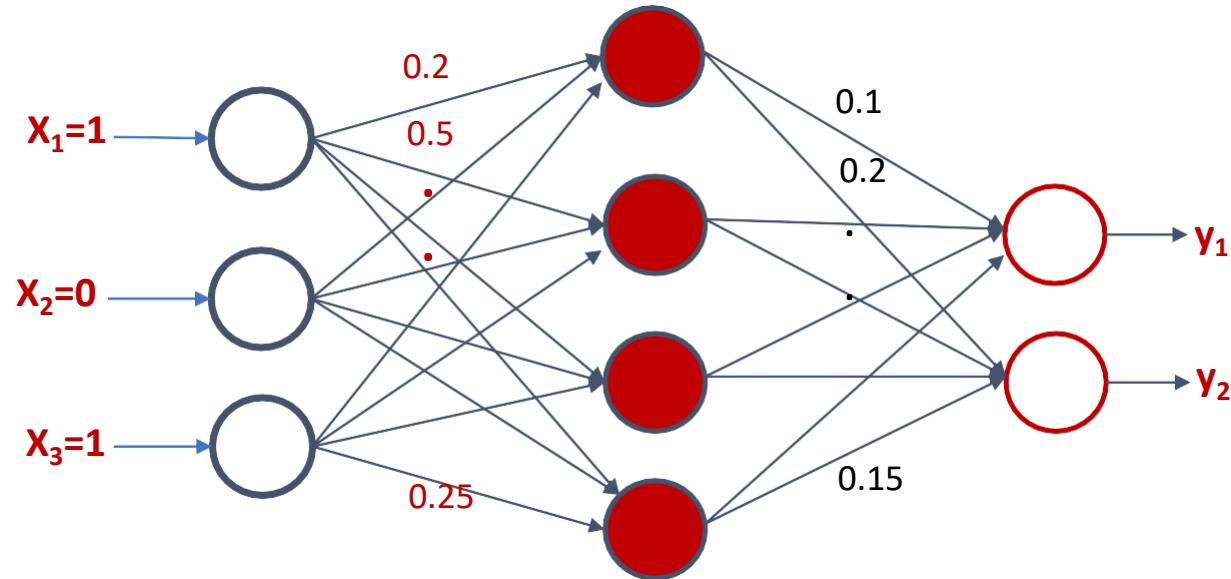


$$E = \|\bar{y} - \hat{y}\| = 0$$

or

$$E = \|\bar{y} - \hat{y}\| \leq \epsilon$$

Backpropagation Through Time

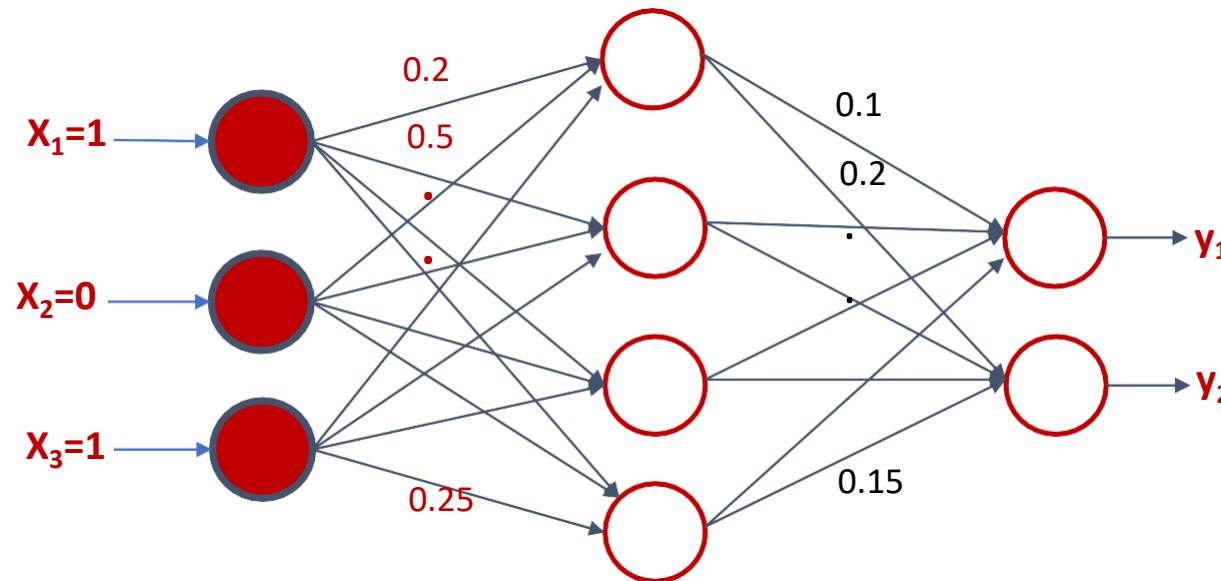


$$E = \|\bar{y} - \hat{y}\| = 0$$

or

$$E = \|\bar{y} - \hat{y}\| \leq \epsilon$$

Backpropagation Through Time

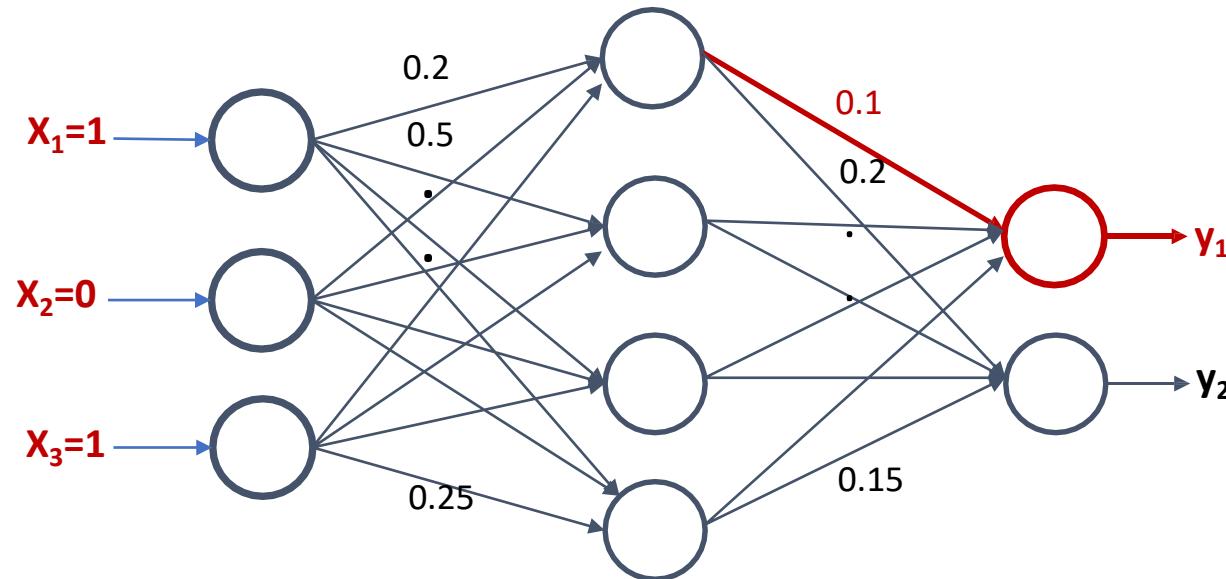


$$E = \|\bar{y} - \hat{y}\| = 0$$

or

$$E = \|\bar{y} - \hat{y}\| \leq \epsilon$$

How are the Derivatives performed



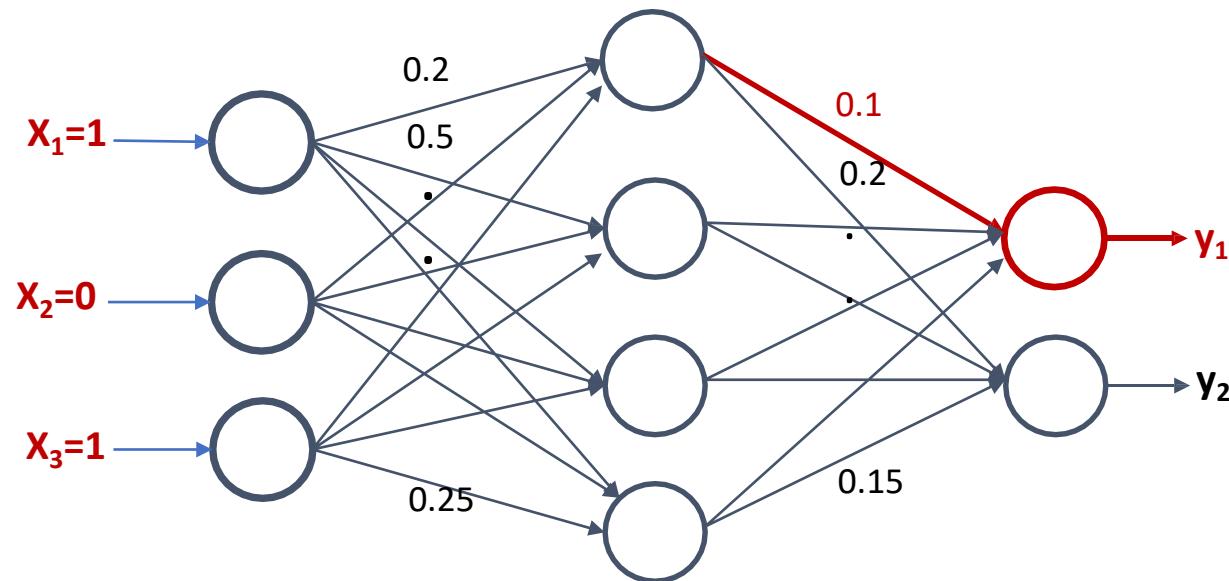
Loss function $E = \|\bar{y} - \hat{y}\|$

$$\nabla = \frac{\delta E}{\delta V} = 0$$

$$\nabla_{11} = \frac{\delta E}{\delta V_{11}} = 0$$

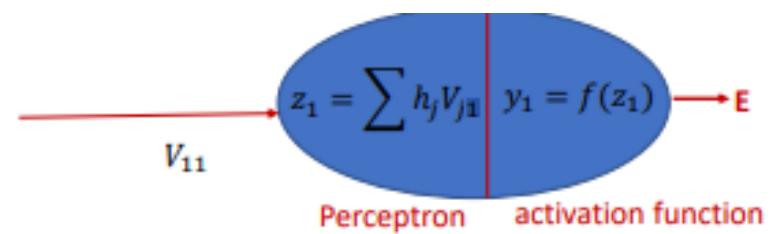
How are the Derivatives performed

Loss function $E = \|\bar{y} - \hat{y}\|$



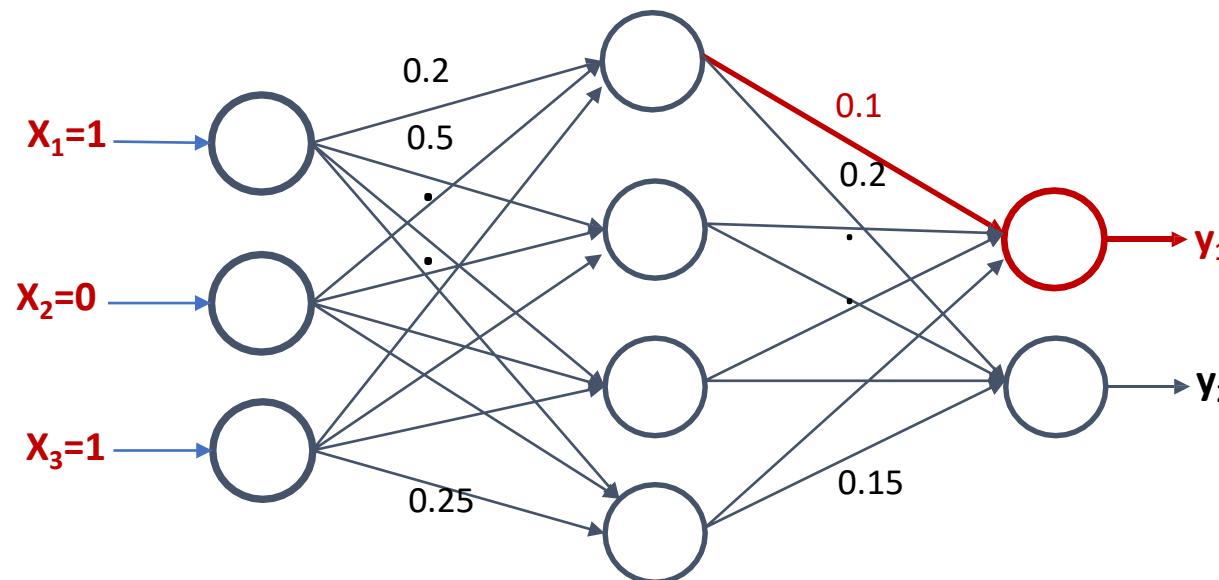
$$\nabla = \frac{\delta E}{\delta V} = 0$$

$$\nabla_{11} = \frac{\delta E}{\delta V_{11}} = 0$$



How are the Derivatives performed

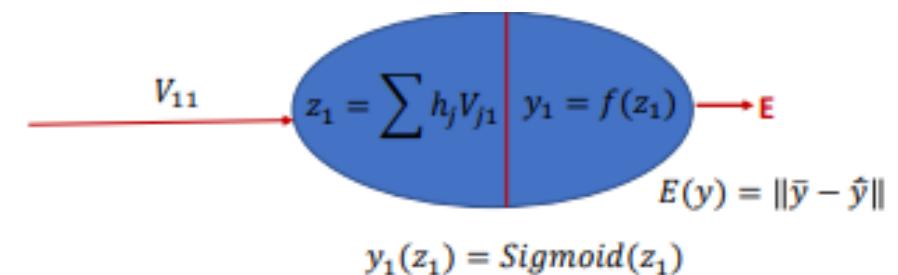
Loss function



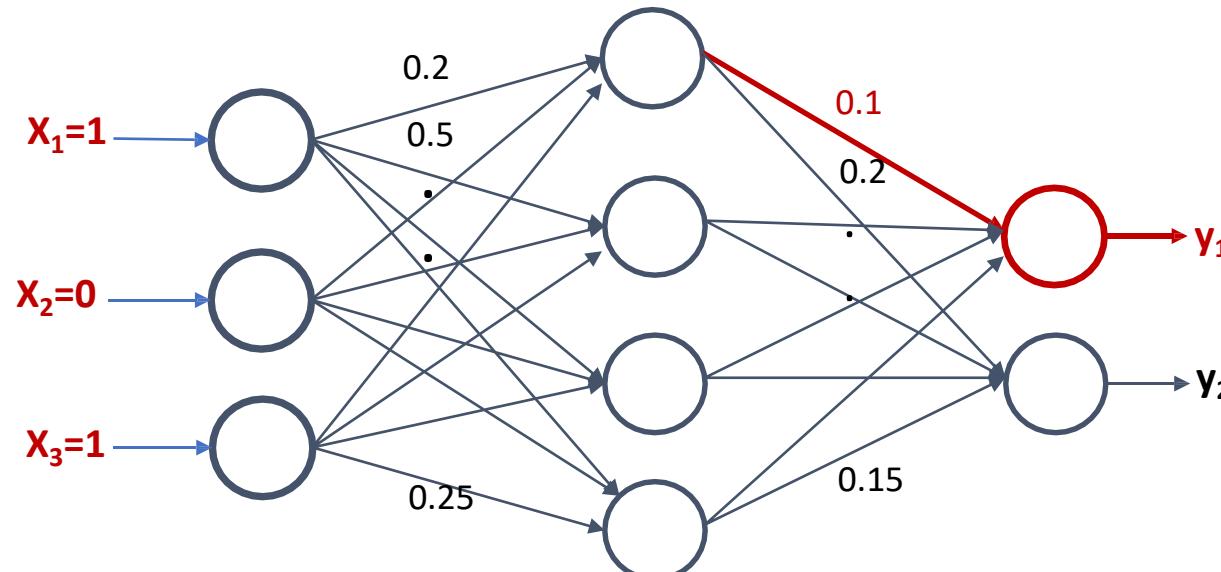
$$\nabla = \frac{\delta E}{\delta V} = 0$$

$$\nabla_{11} = \frac{\delta E}{\delta V_{11}} = 0$$

$$z_1(V_{11}) = h_1V_{11} + h_2V_{21} + h_3V_{31} + h_4V_{41}$$



How are the Derivatives performed



Loss function

$$E = \|\bar{y} - \hat{y}\|$$

$$\nabla = \frac{\delta E}{\delta V} = 0$$

$$\nabla_{ij} = \frac{\delta E}{\delta V_{ij}} = 0$$

$$\nabla_{11} = \frac{\delta E}{\delta V_{11}} = \frac{\delta z_1}{\delta V_{11}} \cdot \frac{\delta y_1}{\delta z_1} \cdot \frac{\delta E}{\delta y_1}$$

$$z_1 = h_1 V_{11} + h_2 V_{21} + h_3 V_{31} + h_4 V_{41}$$

$$y_1 = \text{Sigmoid}(z_1)$$

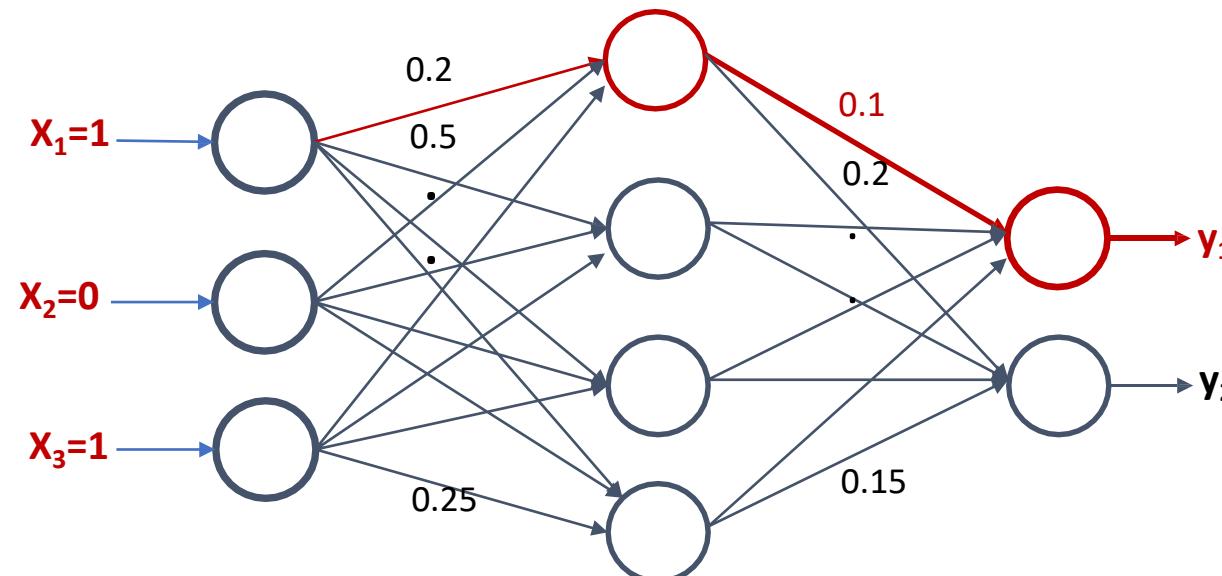
$$\frac{\delta E}{\delta V_{11}} = \frac{\delta z_1}{\delta V_{11}} \times \frac{\delta y_1}{\delta z_1} \times \frac{\delta E}{\delta y_1}$$

$$V_{11} \rightarrow z_1 = \sum h_j V_{j1} \mid y_1 = f(z_1) \rightarrow E$$

Backpropagation

Loss function

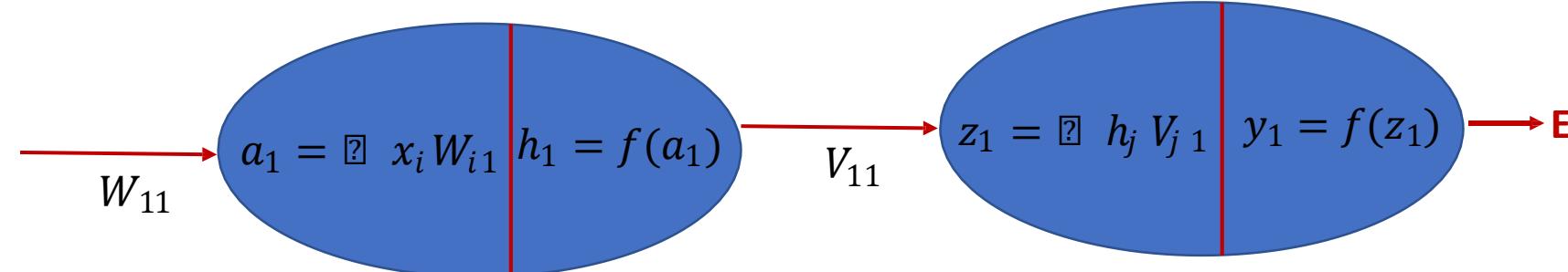
$$E = \|\bar{y} - \hat{y}\|$$



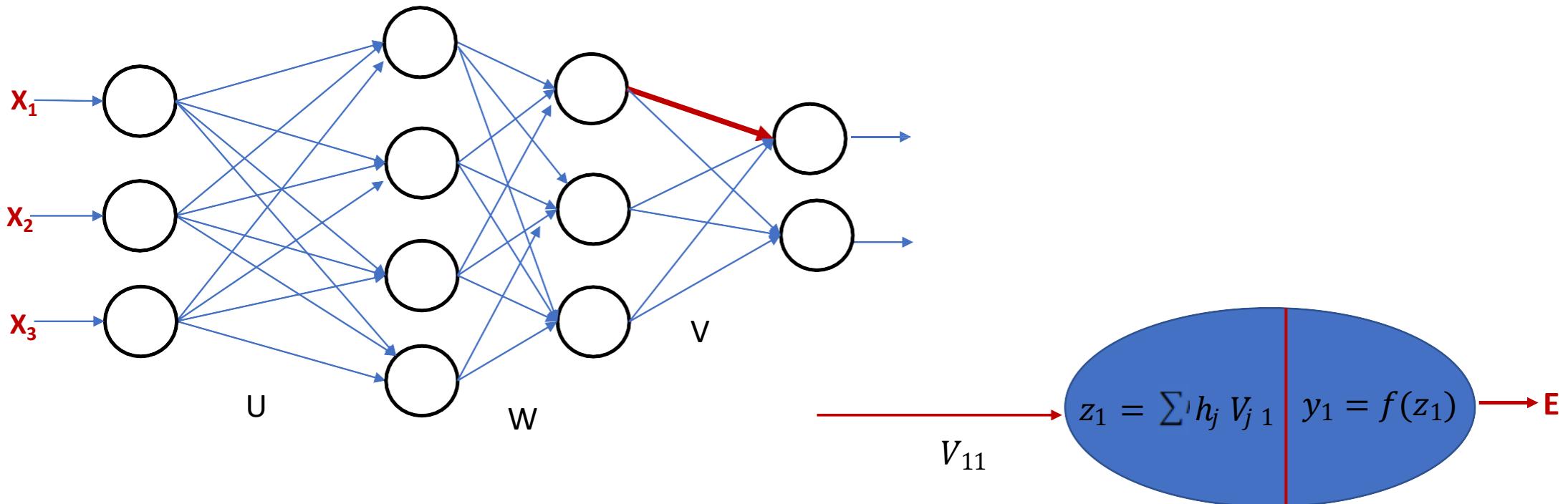
$$\nabla = \frac{\delta E}{\delta W} = 0$$

$$\nabla_i = \frac{\delta E}{\delta W_i} = 0$$

$$\frac{\delta E}{\delta W_{11}} = \frac{\delta a_1}{\delta W_{11}} \times \frac{\delta h_1}{\delta a_1} \times \frac{\delta z_1}{\delta h_1} \times \frac{\delta y_1}{\delta z_1} \times \frac{\delta E}{\delta y_1}$$

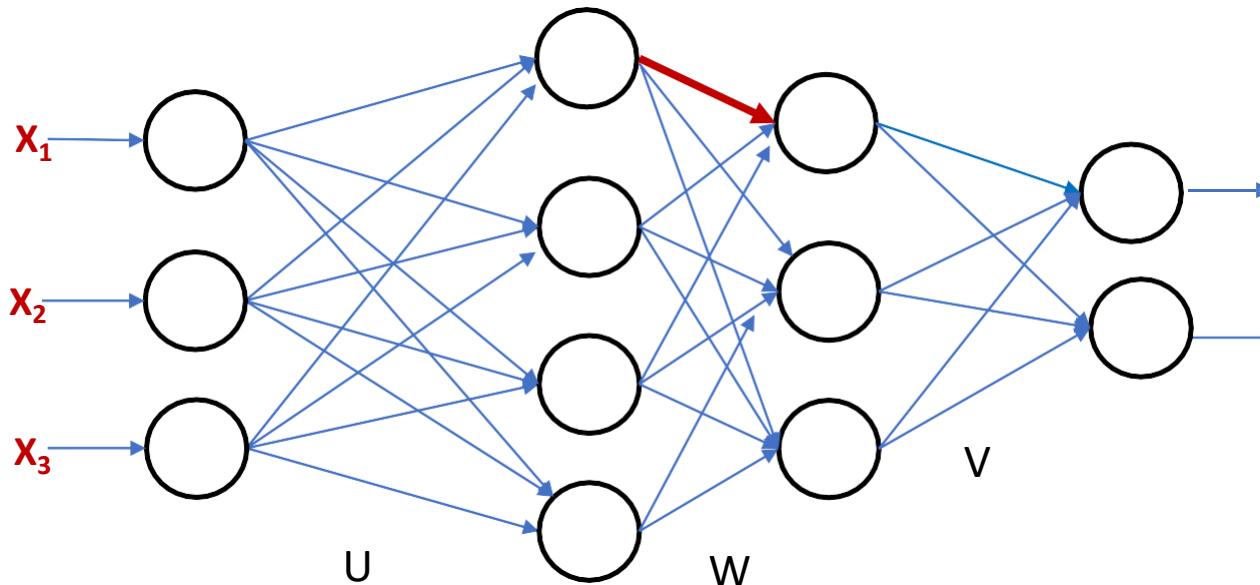


We may have multiple layers.

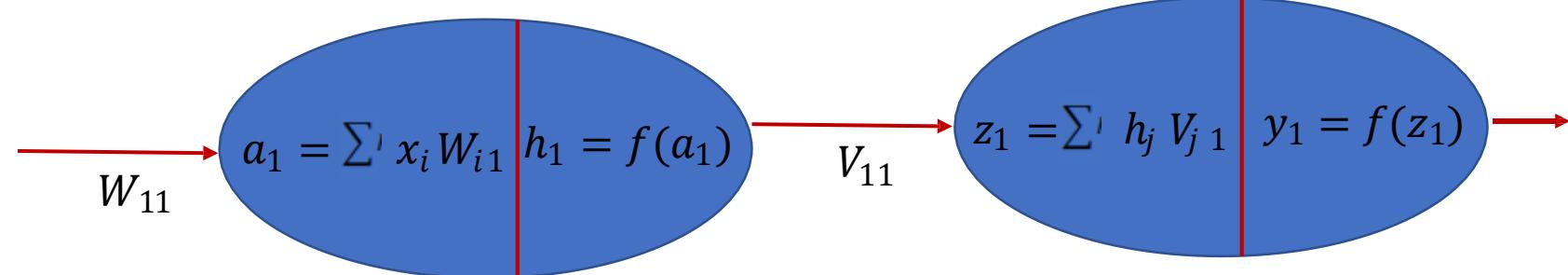


$$\frac{\delta E}{\delta V_{11}} = \frac{\delta z_1}{\delta V_{11}} \times \frac{\delta y_1}{\delta z_1} \times \frac{\delta E}{\delta y_1}$$

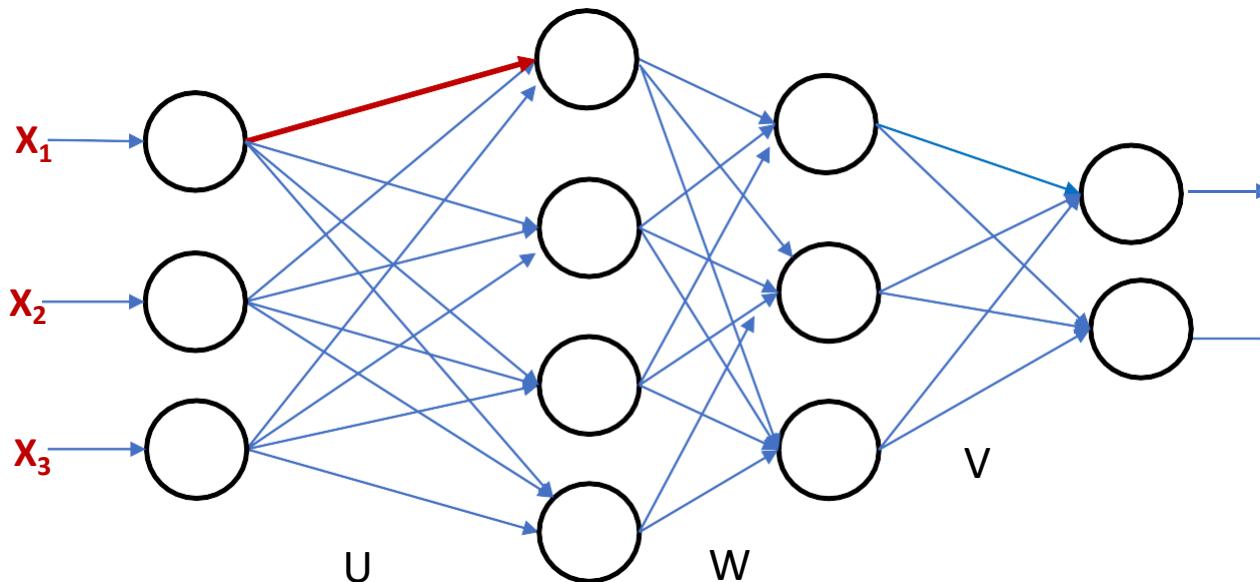
We may have multiple layers



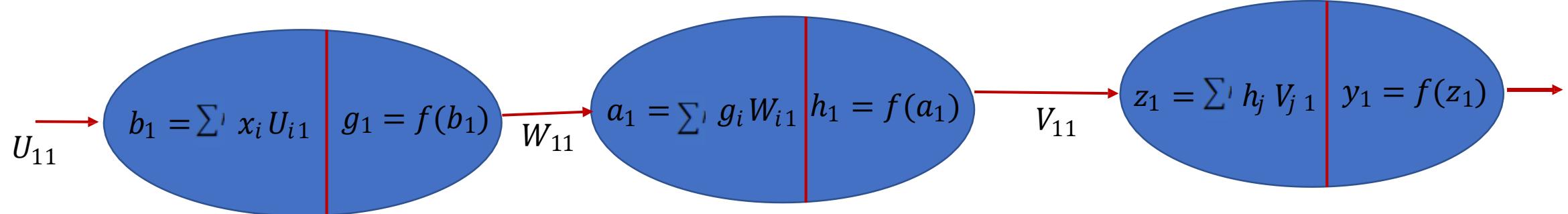
$$\frac{\delta E}{\delta W_{11}} = \frac{\delta a_1}{\delta W_{11}} \times \frac{\delta h_1}{\delta a_1} \times \frac{\delta z_1}{\delta h_1} \times \frac{\delta y_1}{\delta z_1} \times \frac{\delta E}{\delta y_1}$$

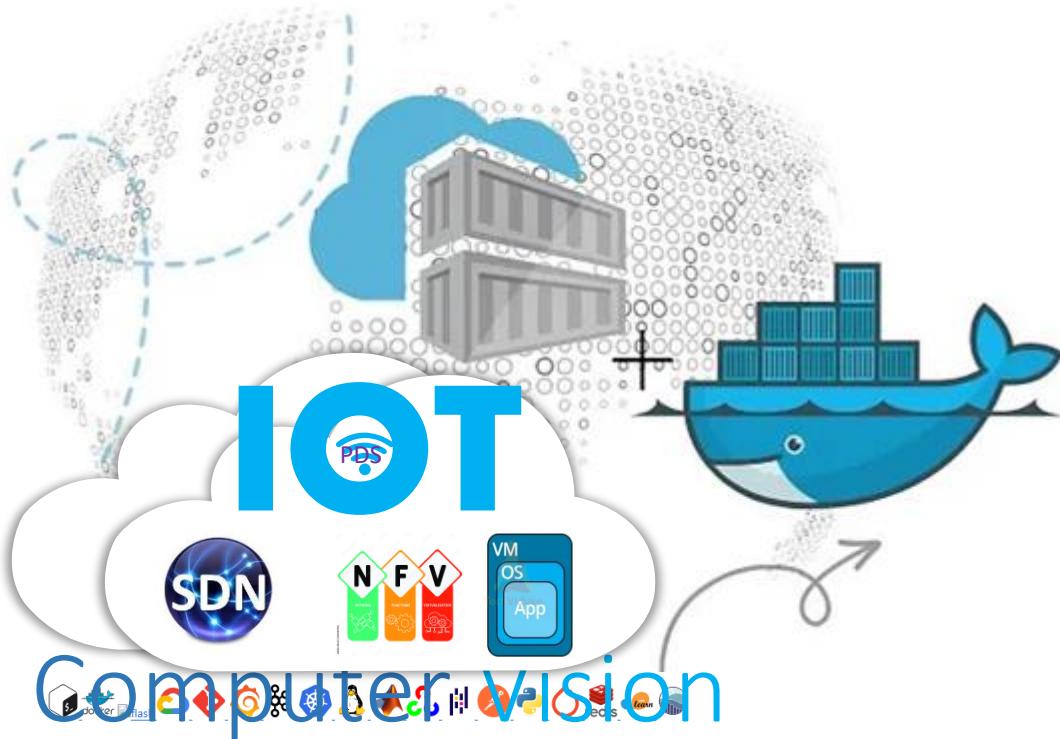


We may have multiple layers



$$\frac{\delta E}{\delta U_{11}} = \frac{\delta b_1}{\delta U_{11}} \times \frac{\delta g_1}{\delta b_1} \times \frac{\delta a_1}{\delta g_1} \times \frac{\delta h_1}{\delta a_1} \times \frac{\delta z_1}{\delta h_1} \times \frac{\delta y_1}{\delta z_1} \times \frac{\delta E}{\delta y_1}$$





Computer Vision

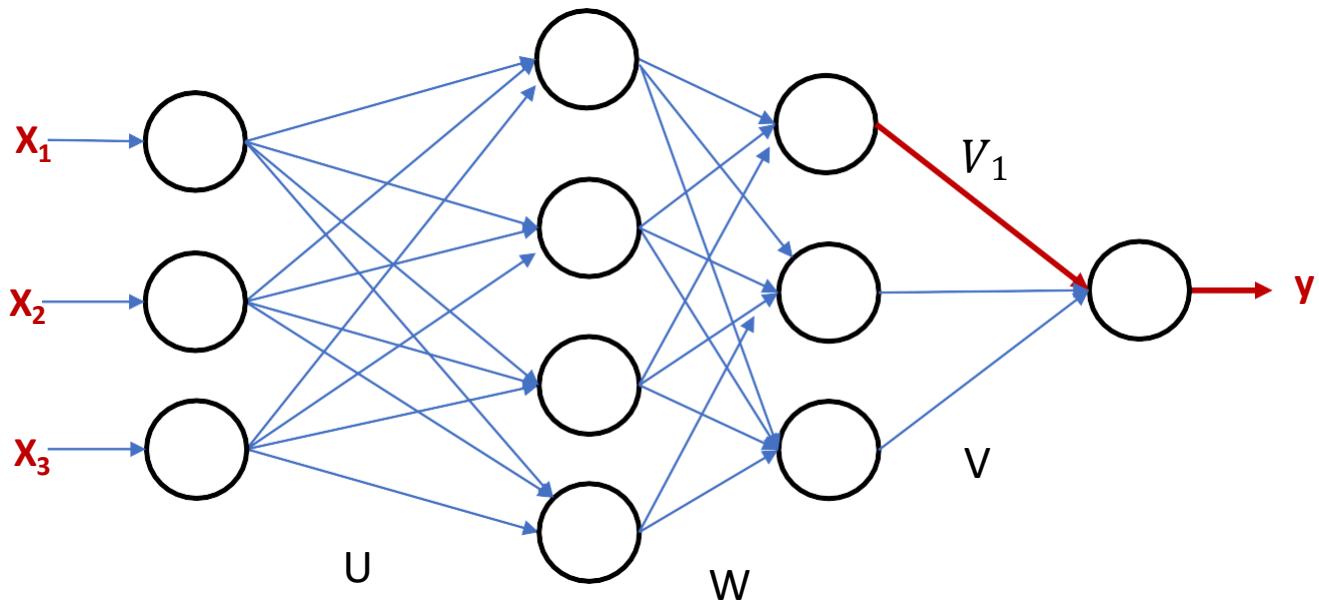
Chapter 5-5:

Learning with different Loss Functions and Their Derivatives

Two Commonly used Loss Functions are

- **Mean Square Error – Standard Loss Function for Regression**
- Cross Entropy Loss - Standard Loss Function for Classification

Mean Square Error (MSE)

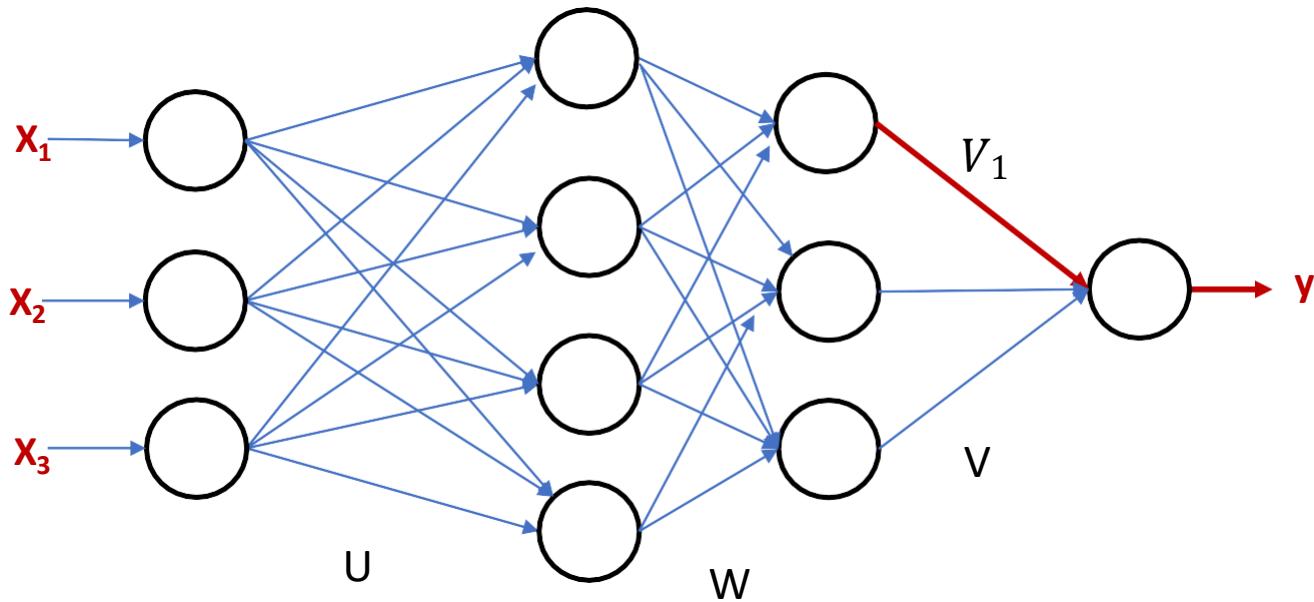


For the Single Sample

$$\text{MSE } E = (y - \hat{y})^2$$

Ground truth is
 y

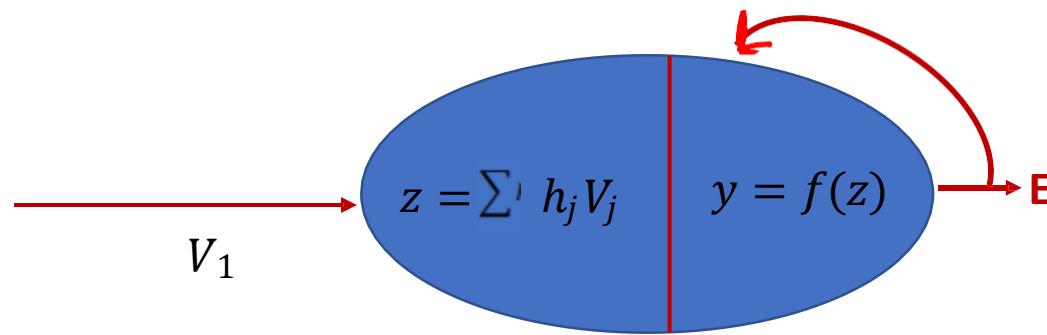
Mean Square Error (MSE)



For the Single Sample

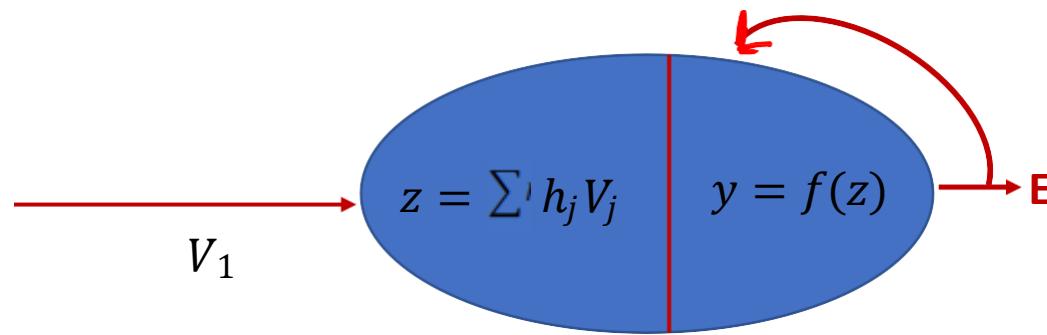
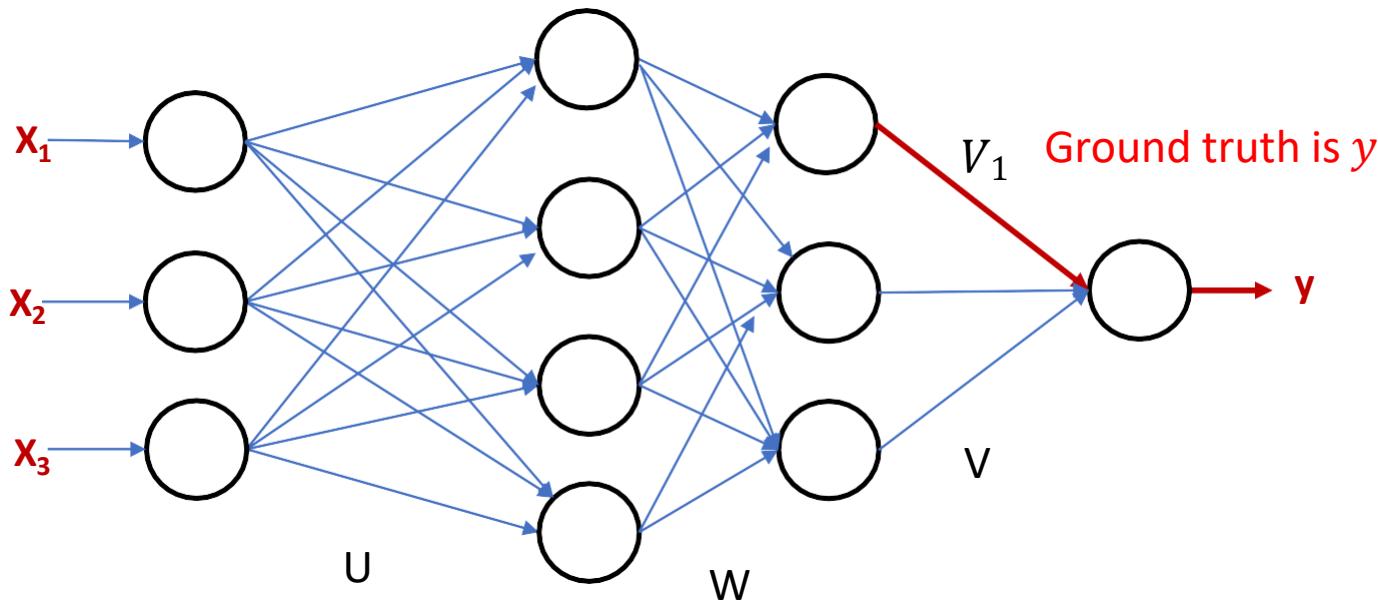
$$\text{MSE } E = (y - \hat{y})^2$$

Ground truth is
 y



$$\frac{\delta E}{\delta V_1} = \frac{\delta z}{\delta V_1} \times \frac{\delta y}{\delta z} \times \boxed{\frac{\delta E}{\delta y}}$$

Mean Square Errors



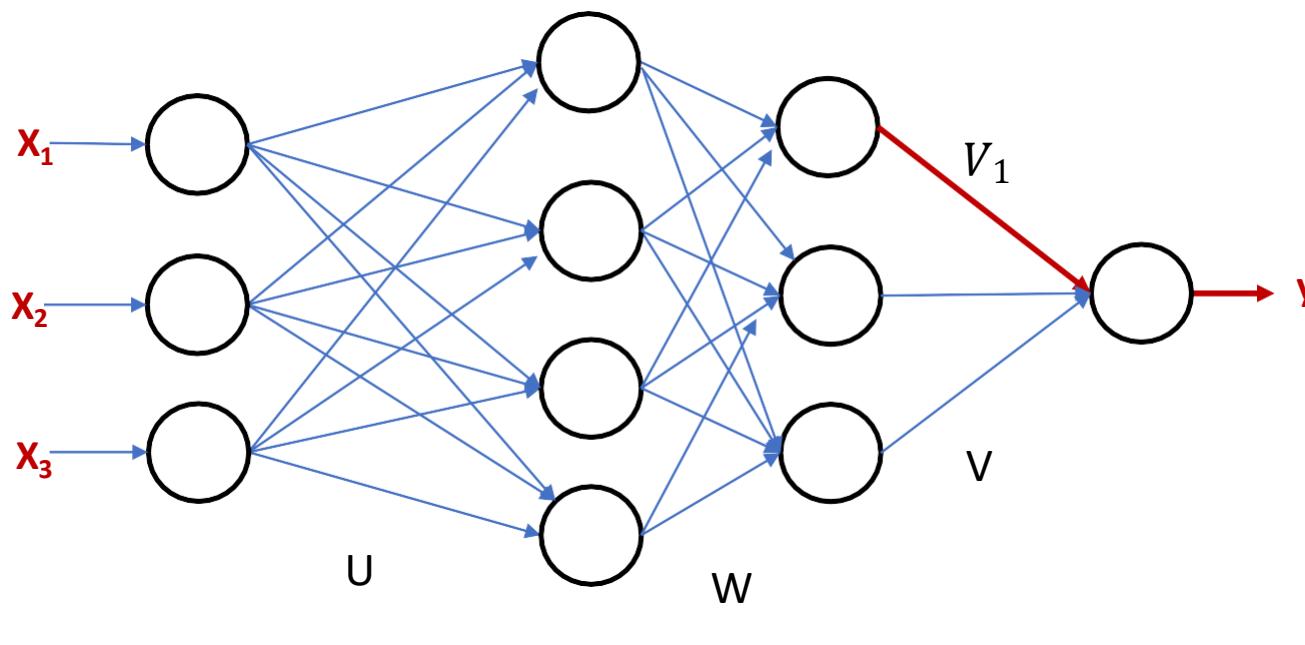
For the Single Sample

$$\text{MSE } E = (y - \hat{y})^2$$

$$\frac{\delta E}{\delta y} = \frac{\delta(y - \hat{y})^2}{\delta y} = 2(y - \hat{y})$$

$$\frac{\delta E}{\delta V_1} = \frac{\delta z}{\delta V_1} \times \frac{\delta y}{\delta z} \times \boxed{\frac{\delta E}{\delta y}}$$

Mean Square Errors

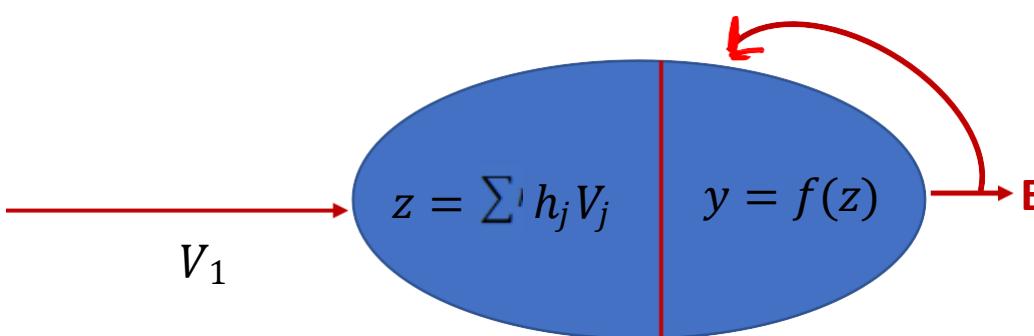


If the ground truth is y

For n Samples

$$E = \frac{1}{n} \sum_{i=1}^n (y^i - \hat{y})^2$$

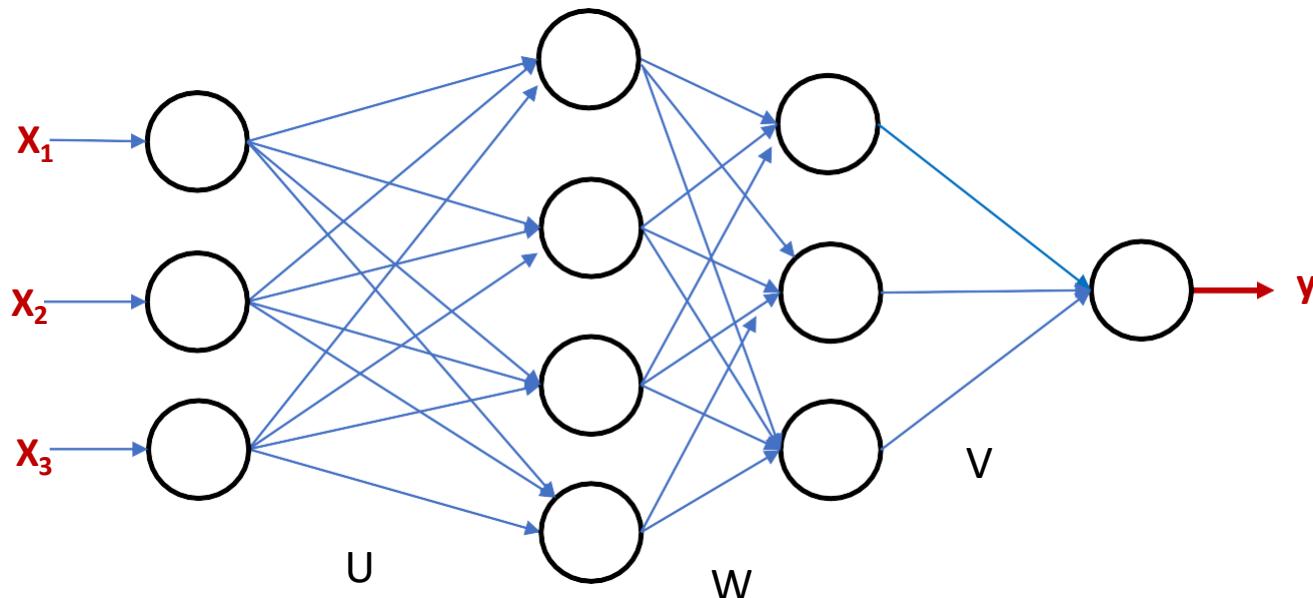
$$\frac{\delta E}{\delta y} = \frac{\delta \frac{1}{n} \sum_{i=1}^n (y^i - \hat{y})^2}{\delta y} = \frac{2}{n} \sum_{i=1}^n (y^i - \hat{y})$$



$$\frac{\delta E}{\delta V_1} = \frac{\delta z}{\delta V_1} \times \frac{\delta y}{\delta z} \times \boxed{\frac{\delta E}{\delta y}}$$

Σ

Mean Square Errors



If the ground truth is y

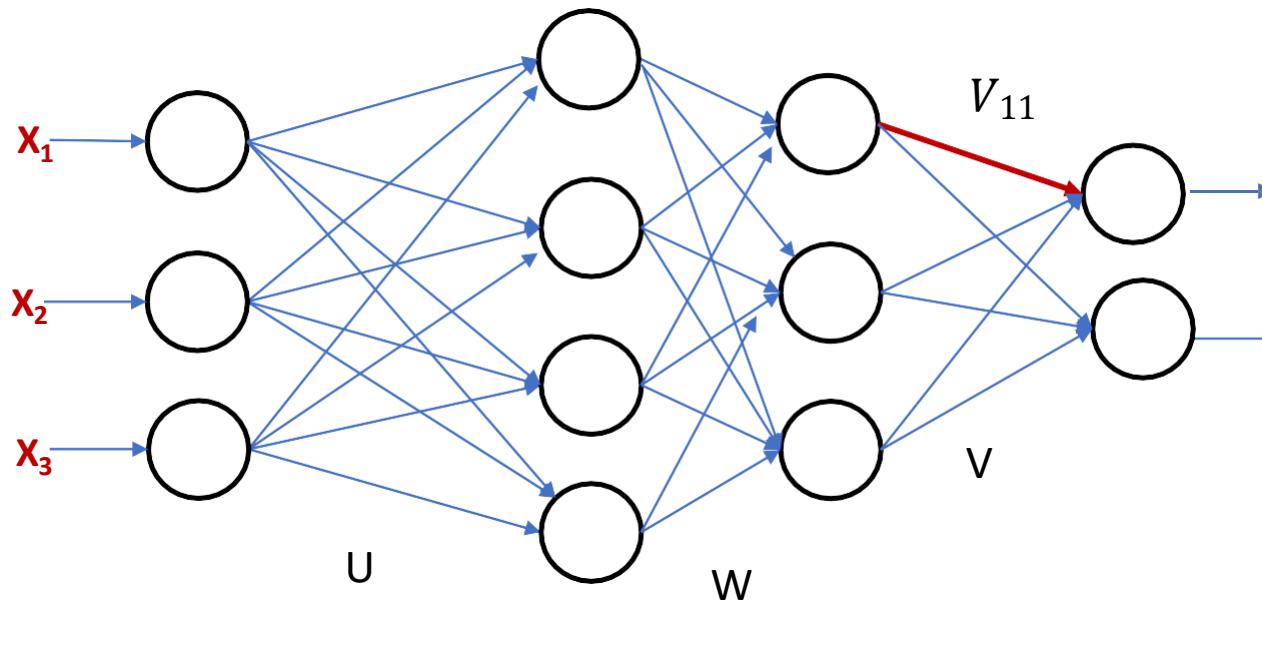
For n Samples

$$E = \frac{1}{n} \sum_{i=1}^n (y^i - \hat{y})^2$$

$$\frac{\delta E}{\delta y} = \frac{\delta \frac{1}{n} \sum_{i=1}^n (y^i - \hat{y})^2}{\delta y} = \frac{2}{n} \sum_{i=1}^n (y^i - \hat{y})$$

Backpropagation will be done after a batch of n Samples

Mean Square Errors



If the ground truth is y

For the n Sample

$$E = \frac{1}{n} \sum_{i=1}^n (y^i - \hat{y})^2$$

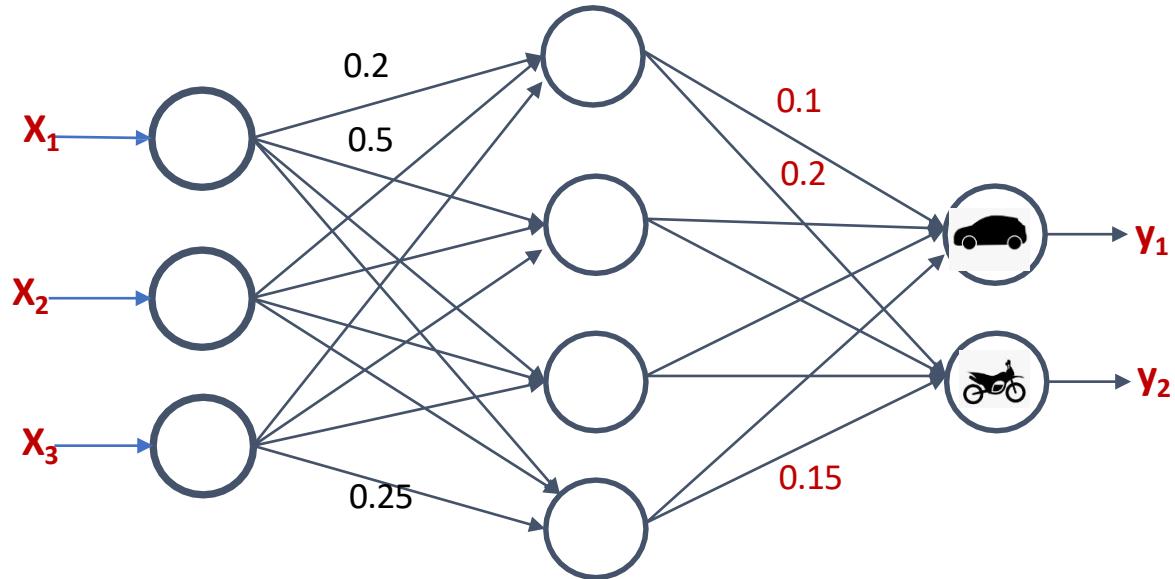
$$\frac{\delta E}{\delta y_1} = \frac{\delta \frac{1}{n} \sum_{i=1}^n (y^i - \hat{y})^2}{\delta y_1} = \frac{2}{n} \sum_{i=1}^n (y_1^i - \hat{y})$$

$$V_{11} \rightarrow z_1 = \sum h_j V_{j1} \quad | \quad y_1 = f(z_1) \rightarrow E$$

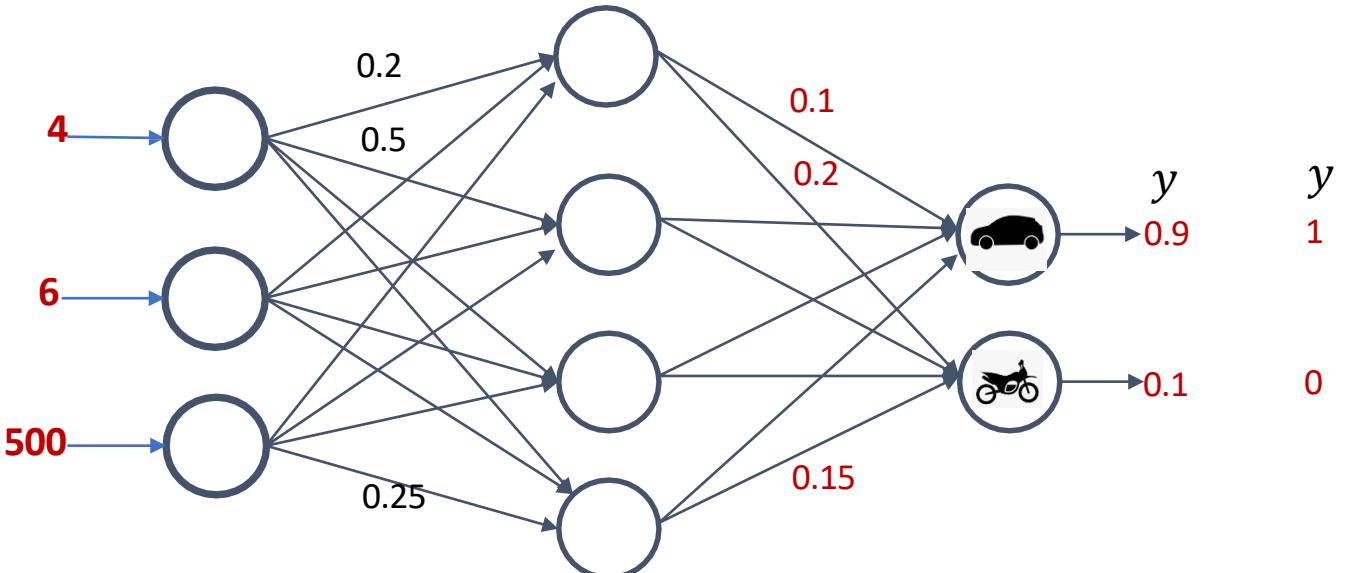
$$\frac{\delta E}{\delta V_{11}} = \frac{\delta z_1}{\delta V_{11}} \times \frac{\delta y_1}{\delta z_1} \times \frac{\delta E}{\delta y_1}$$

Let us illustrate with a toy example

#Wheel	Height	Weight	
4	6	500	
4	5.5	600	
4	5	550	
2	3	200	
2	3.5	150	
2	4	250	

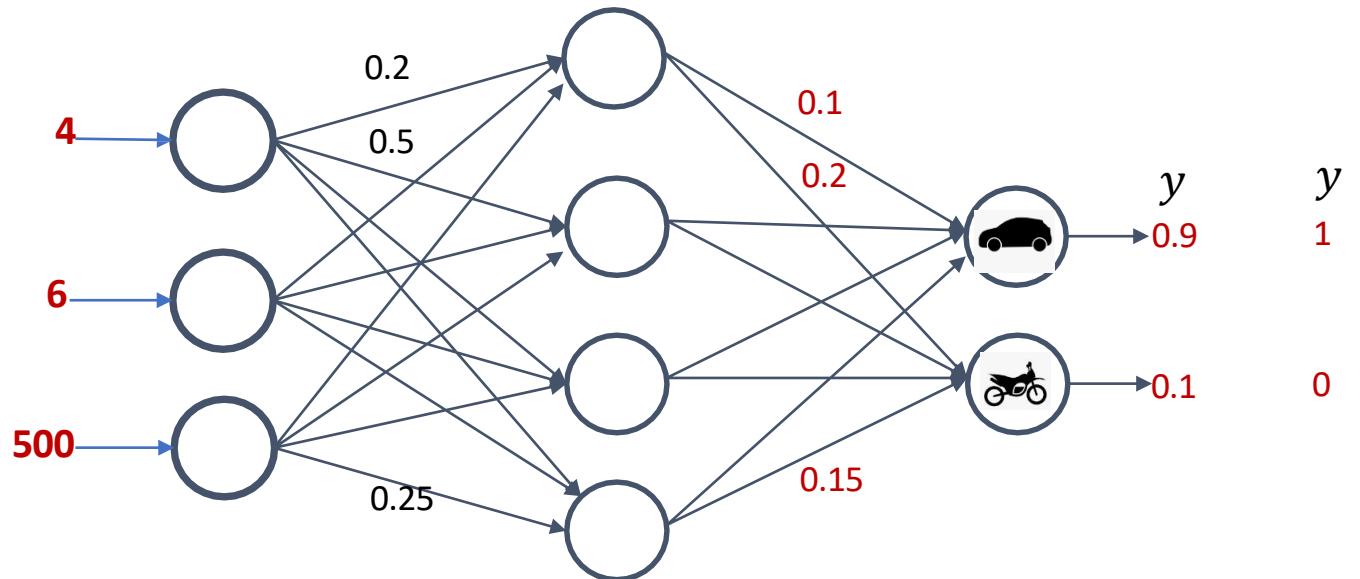


#Wheel	Height	Weight	
4	6	500	
4	5.5	600	
4	5	550	
2	3	200	
2	3.5	150	
2	4	250	



#Wheel Height Weight

			
4	6	500	
4	5.5	600	
4	5	550	
2	3	200	
2	3.5	150	
2	4	250	



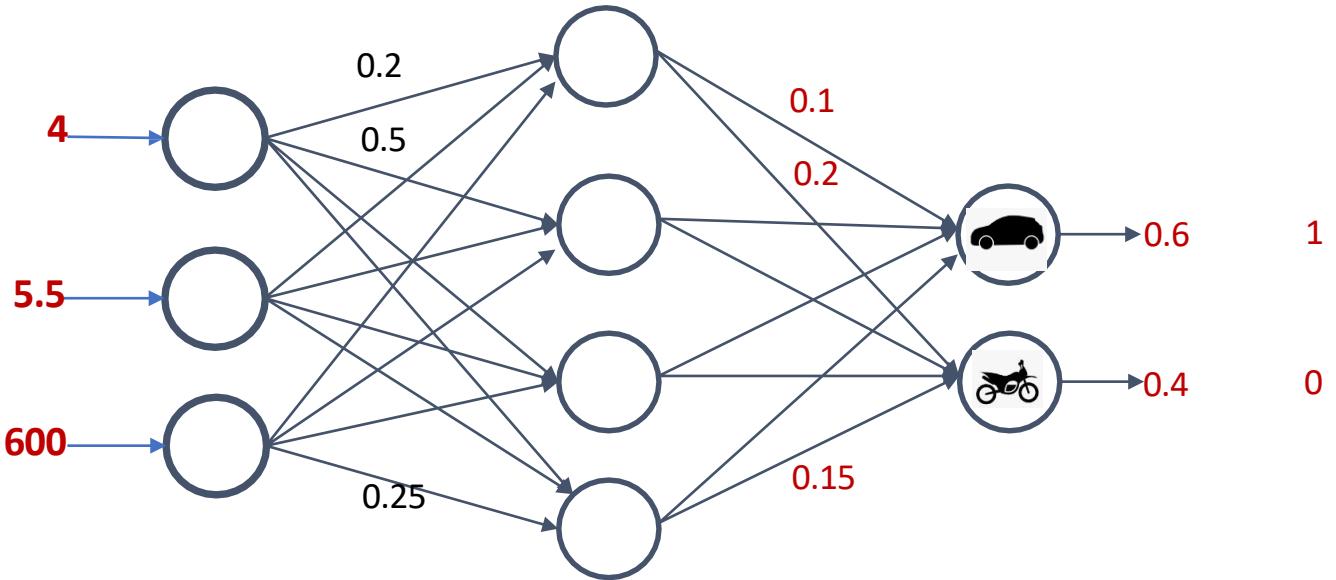
$$\text{Error } E = (y - \hat{y})^2$$

$$\text{Error } E_{y_1} = (0.9 - 1)^2 = 0.01$$

$$\text{Error } E_{y_2} = (0.1 - 0)^2 = 0.01$$

If these errors are not acceptable, then Backpropagate.

#Wheel	Height	Weight	
4	6	500	
4	5.5	600	
4	5	550	
2	3	200	
2	3.5	150	
2	4	250	



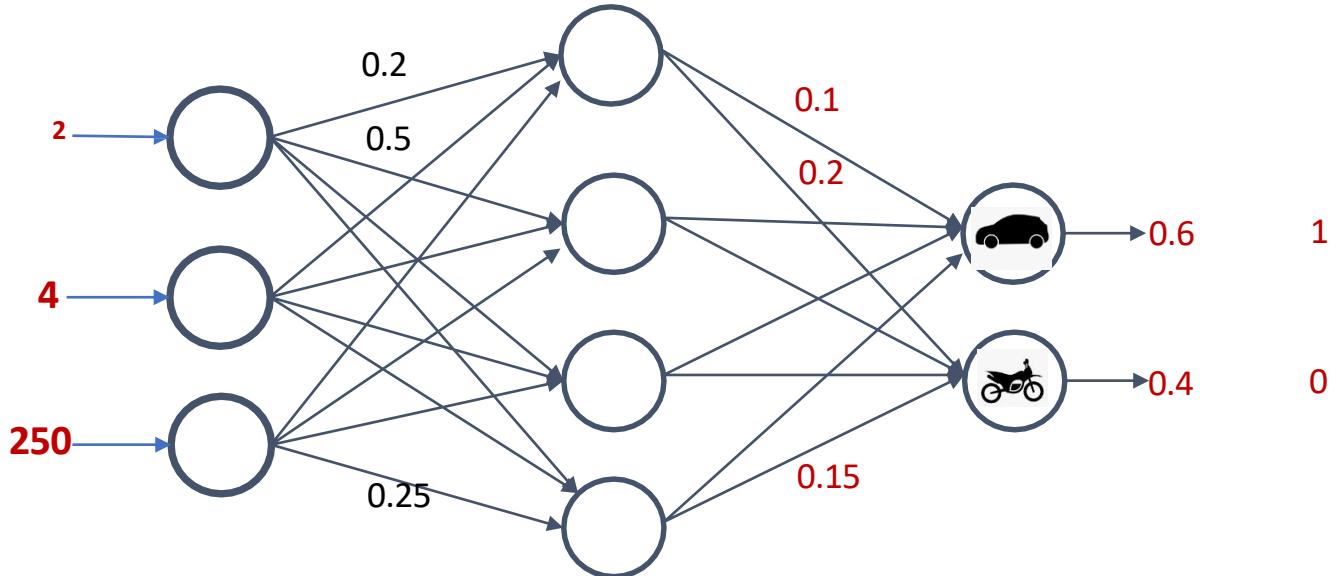
$$\text{Error } E = (y - \hat{y})^2$$

$$\text{Error } E_{y_1} = (0.6 - 1)^2 = 0.16$$

$$\text{Error } E_{y_2} = (0.4 - 0)^2 = 0.36$$

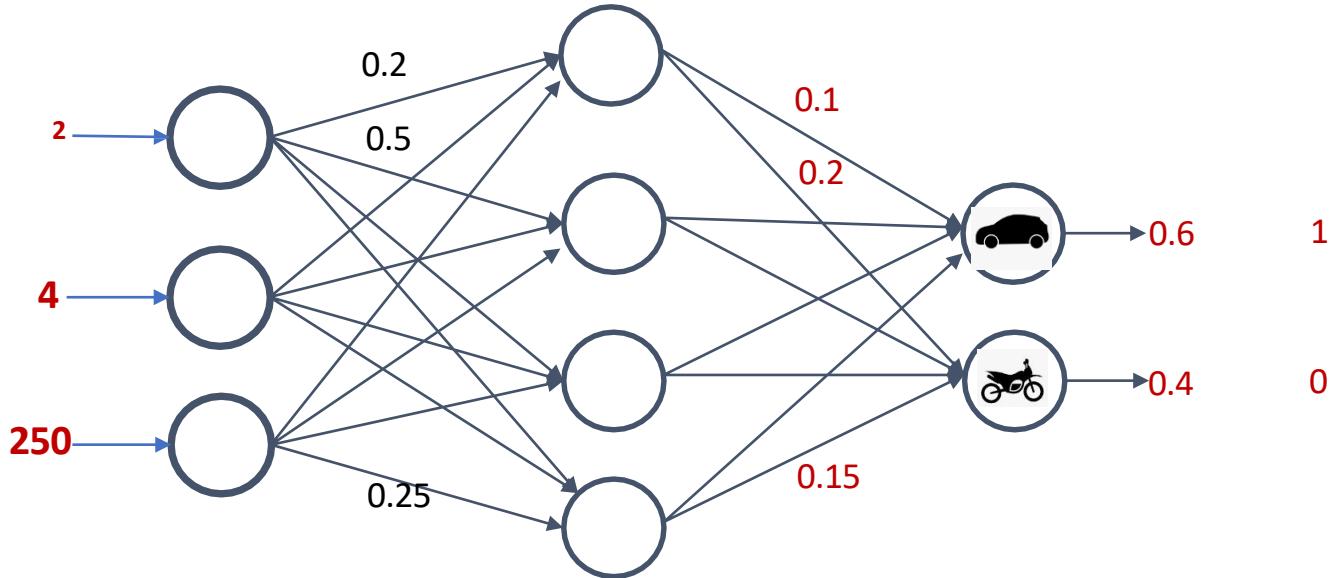
If these errors are not acceptable, then Backpropagate.

#Wheel	Height	Weight	
4	6	500	
4	5.5	600	
4	5	550	
2	3	200	
2	3.5	150	
2	4	250	



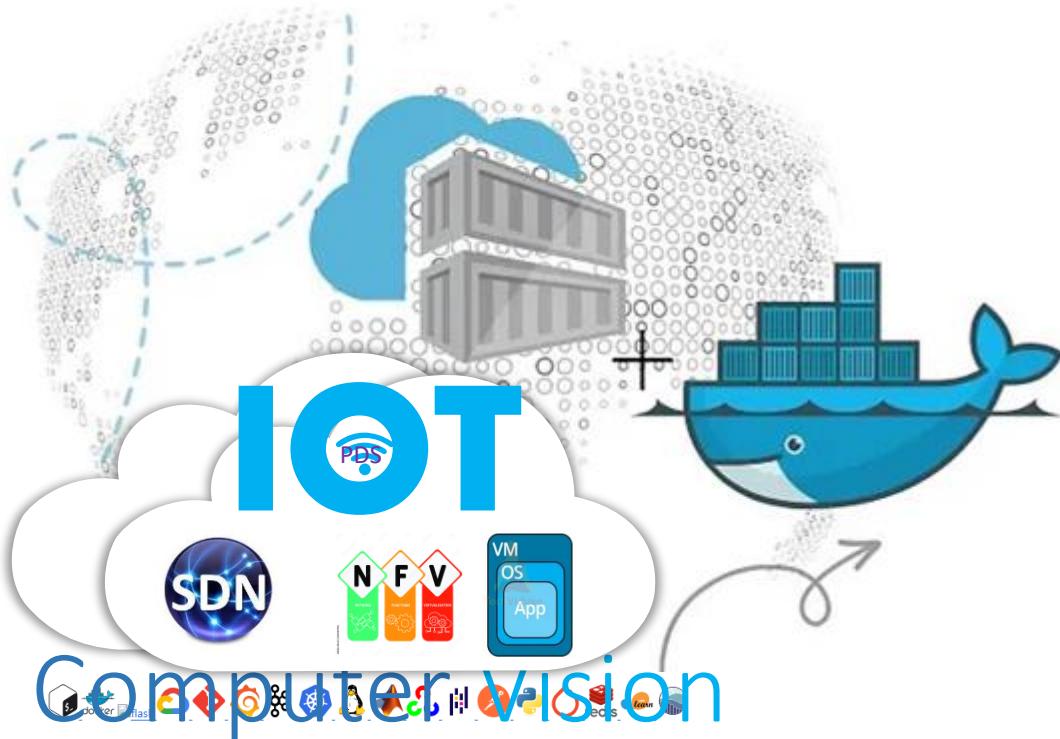
One complete cycle of training is called Epoch

#Wheel	Height	Weight	
4	6	500	
4	5.5	600	
4	5	550	
2	3	200	
2	3.5	150	
2	4	250	



Backpropagation after every sample is expensive.

Do it in batches.



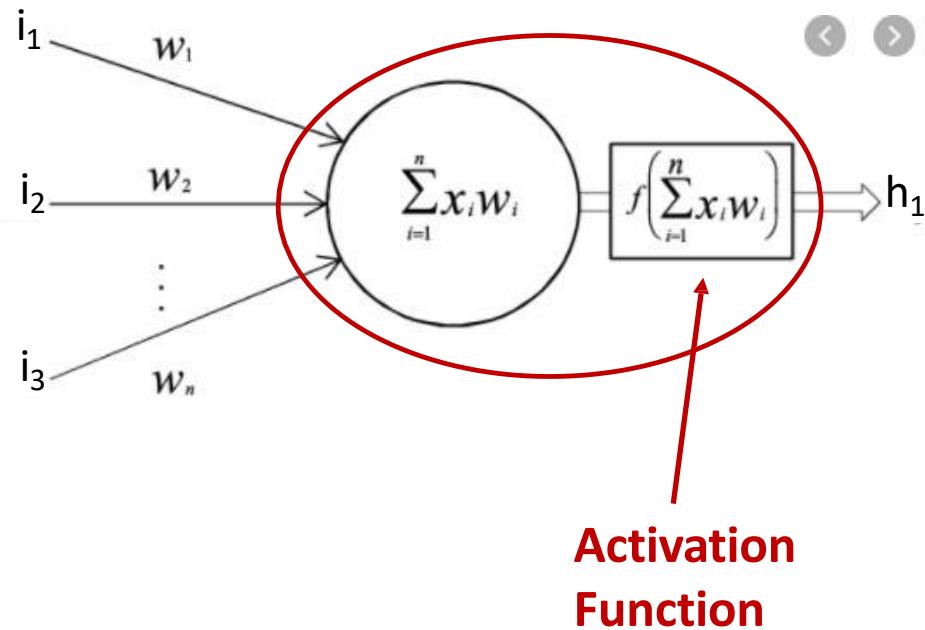
Computer Vision

Chapter 5-6:

Activation Functions and Their Derivatives

- Characteristics of different Activation Functions and their gradient
- The choice of activation function depend on the nature of the problem, nature of the target output and the deepness of the network.

Activation Functions



Activation Function is applied over the linear weighted summation of the incoming information to a node.

Convert linear input signals from perceptron to a linear/non-linear output signal.

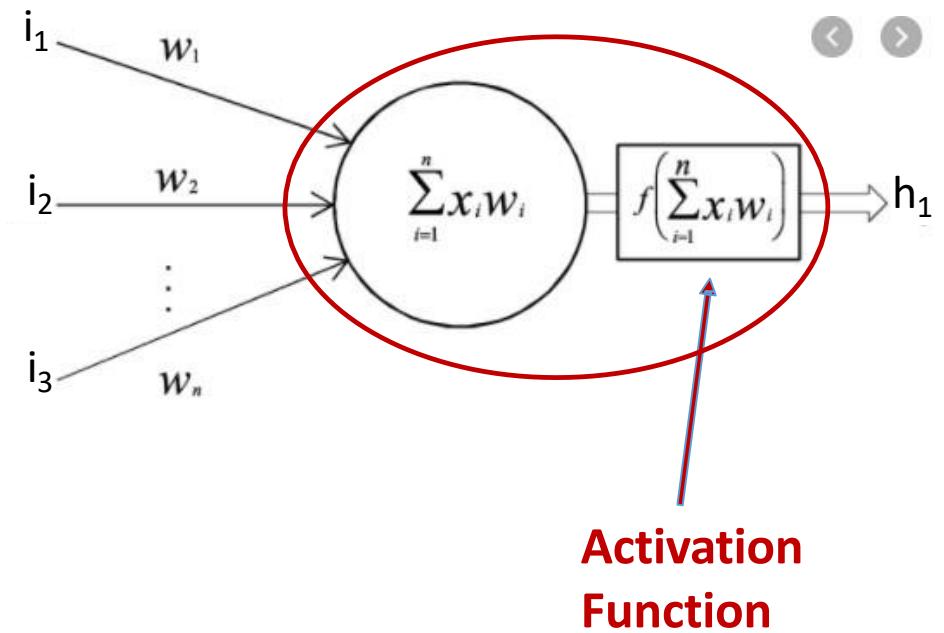
It decides whether to activate a node or not.

Activation Functions

Activation functions must be **monotonic**, **differentiable**, and **quickly converging**.

Types of Activation Functions:

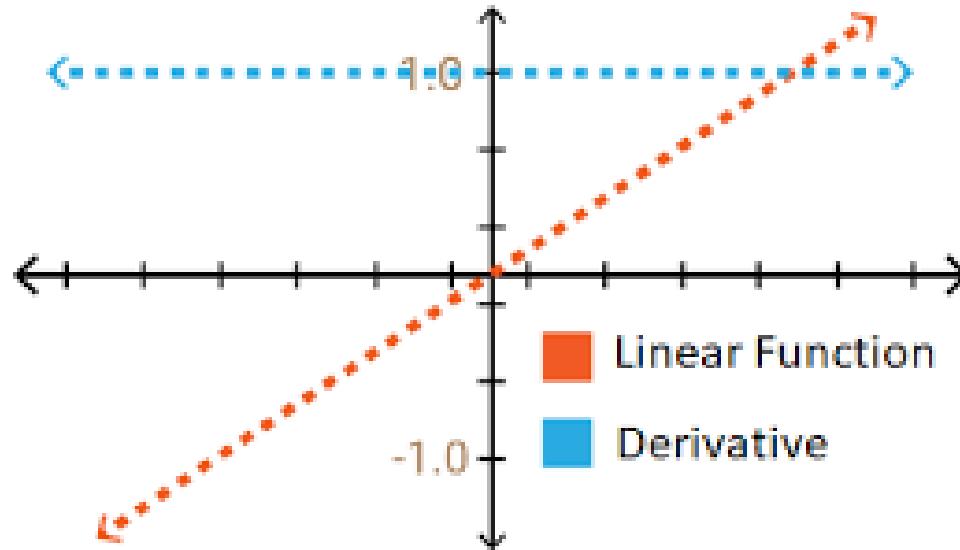
- Linear
- Non-Linear



Linear

$$f(x) = ax + b$$

$$\frac{df(x)}{dx} = a$$



Observations:

- Constant gradient
- Gradient does not depend on the change in the input

Linear

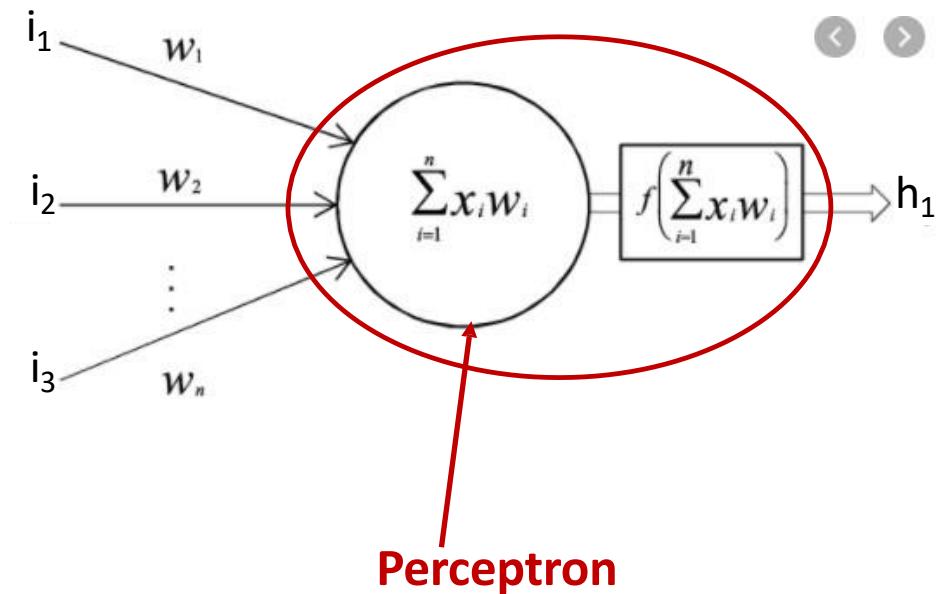
$$f(x) = ax + b$$

$$f(x) = a_1x_1 + a_2x_2 + a_3x_3 + \cdots + b$$

Linear

$$f(x) = ax + b$$

$$f(x) = a_1x_1 + a_2x_2 + a_3x_3 + \dots + b$$



Non-Linear

- Sigmoid (Logistic)
- Hyperbolic Tangent (Tanh)
- Rectified Linear Unit (ReLU)
 - *Leaky Relu*
 - *Parametric Relu*
- Exponential Linear Unit (ELU)

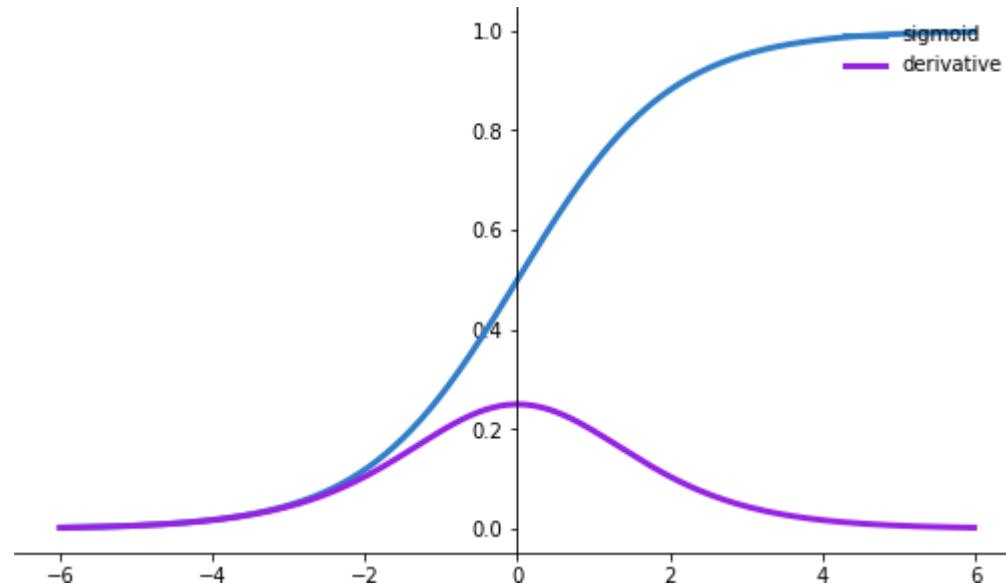
Sigmoid Activation Functions (Logistics)

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

$$\frac{df(x)}{dx} = f(x)(1 - f(x))$$

Observations:

- Output: 0 to 1
- Outputs are not zero-centered
- Can saturate and kill (vanish) gradients



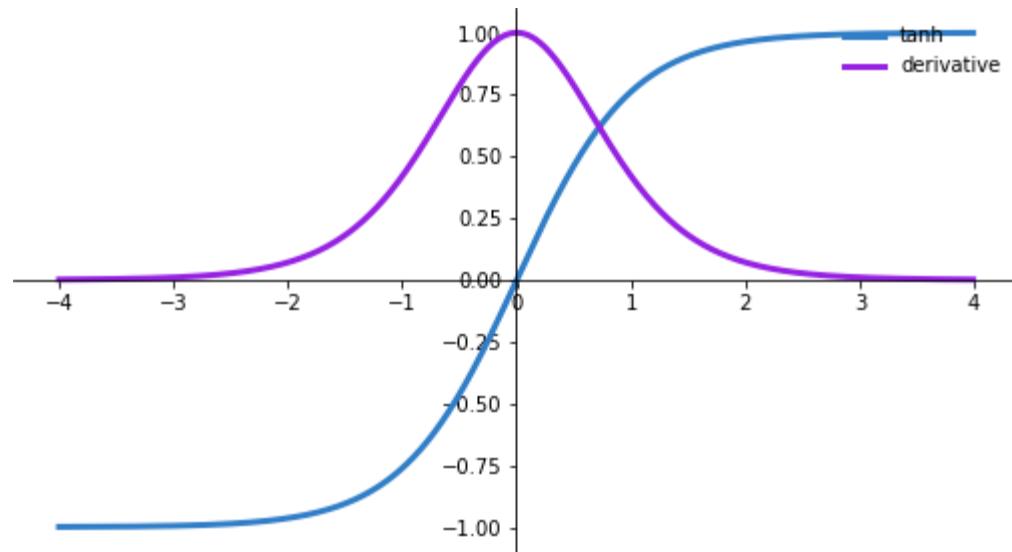
Tanh Activation Function

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

$$\frac{df(x)}{dx} = 1 - f(x)^2$$

Observations:

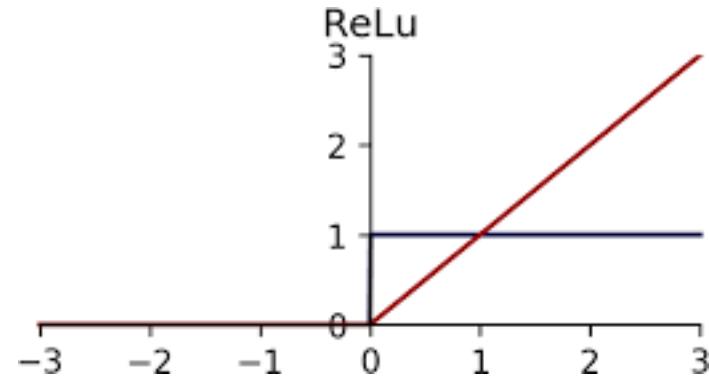
- Output: -1 to +1
- Outputs are zero-centered
- Can Saturate and kill (vanish) gradients
- Gradient is more steeped than Sigmoid, resulting in faster convergence



Rectified Linear Unit(ReLU)

$$f(x) = \max(0, x)$$

$$\frac{df(x)}{dx} = 1$$



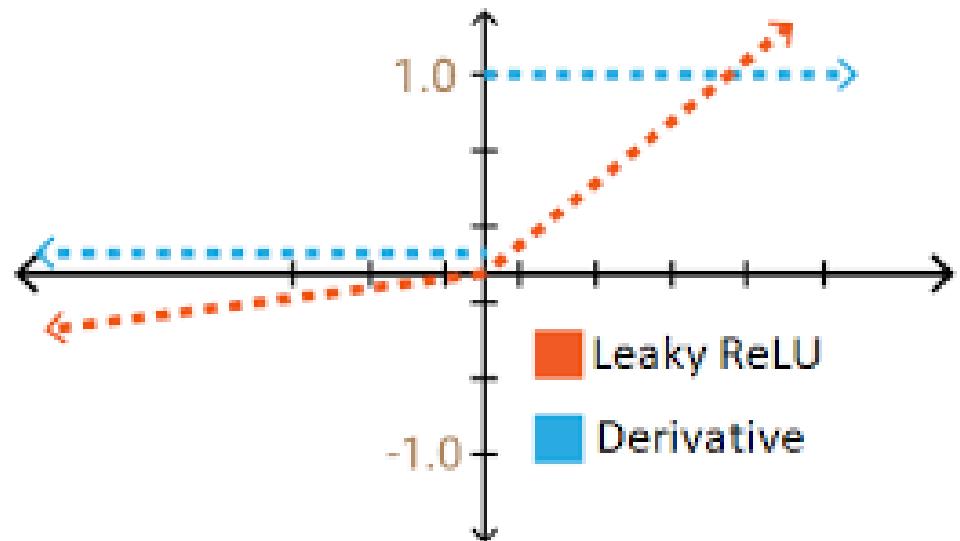
Observations:

- Greatly increase training speed compared to tanh and sigmoid
- Reduces likelihood of killing(vanishing) gradient
- It can blow up activation
- Dead nodes

Leaky-ReLU

$$f(x) = \max(0.01x, x)$$

$$\frac{df(x)}{dx} = \begin{cases} 0.01, & x < 0 \\ 1, & x \geq 0 \end{cases}$$



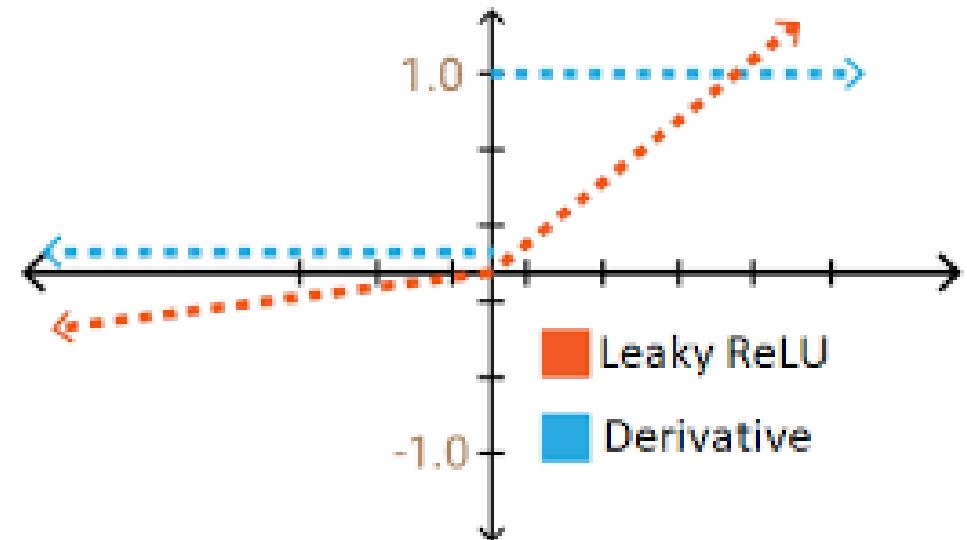
Observations:

- Fixed dying ReLU

Parameterized-ReLU

$$f(x) = \max(\alpha x, x)$$

$$\frac{df(x)}{dx} = \begin{cases} \alpha, & x < 0 \\ 1, & x \geq 0 \end{cases}$$



Observations:

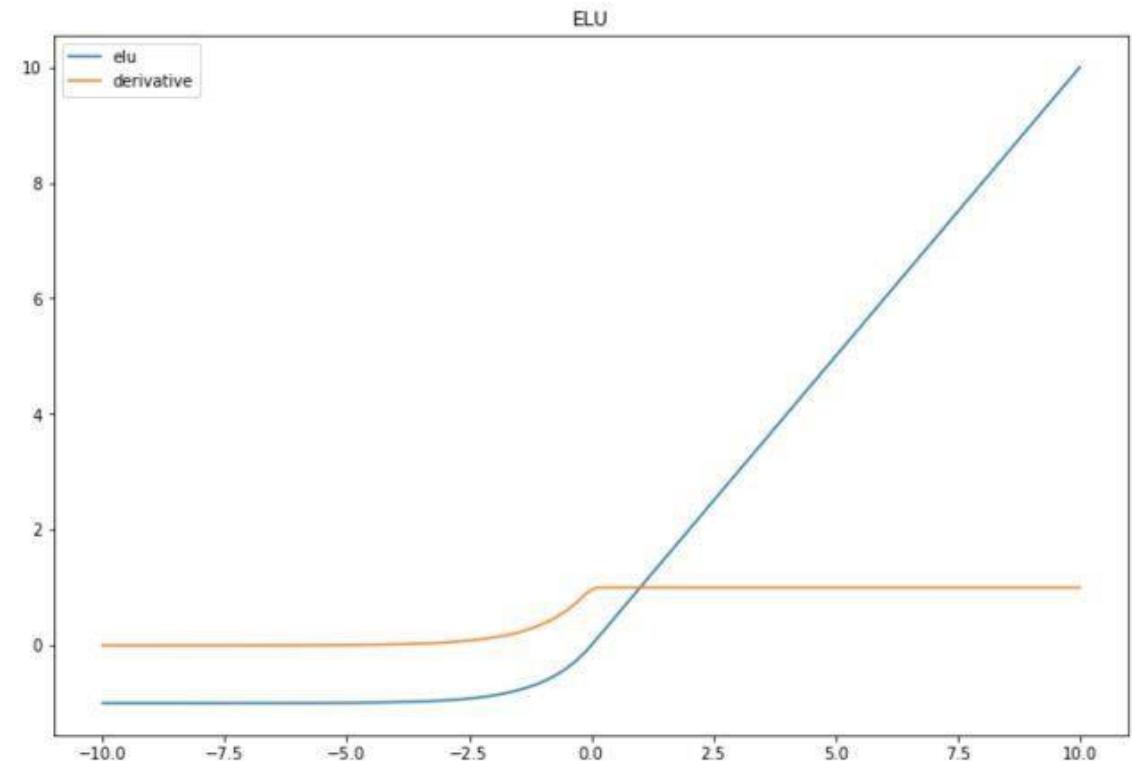
Exponential Linear Unit (ELU)

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

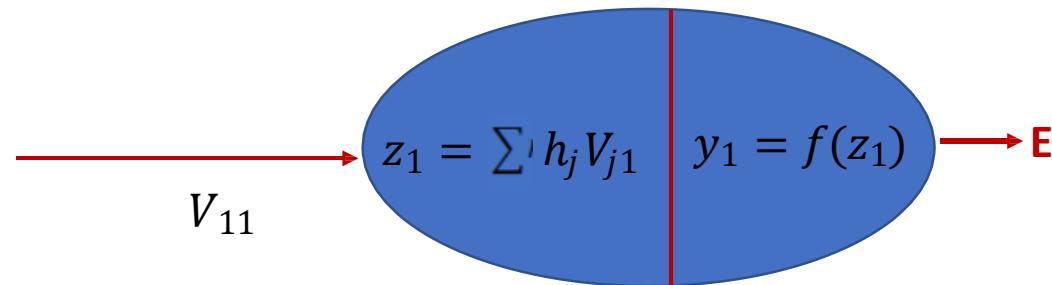
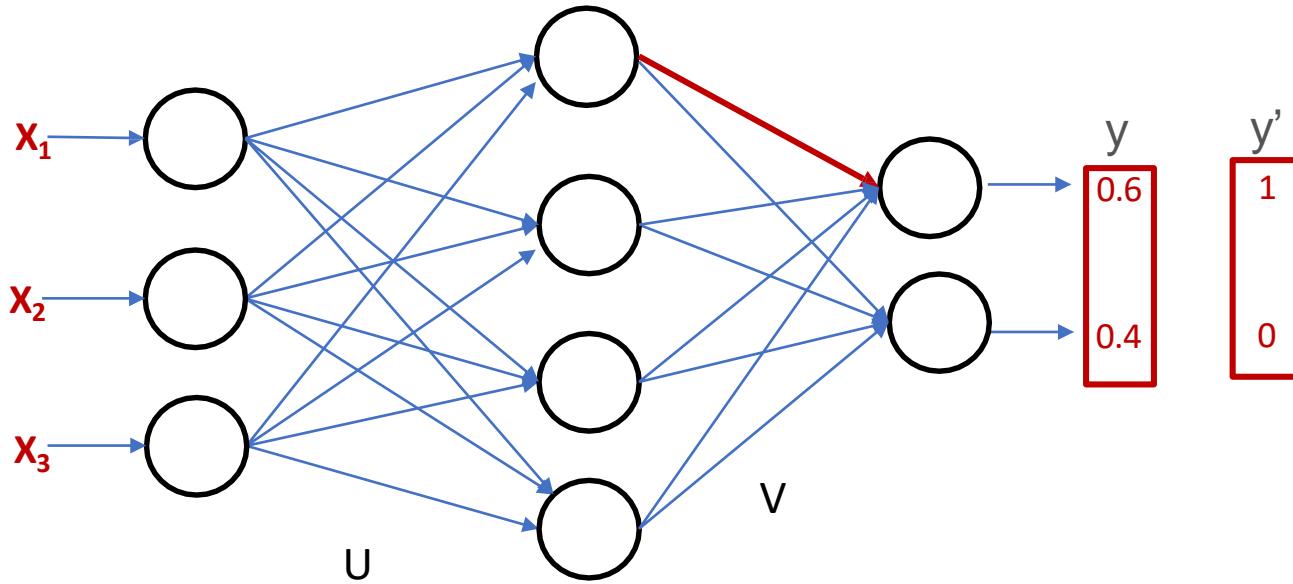
$$\frac{df(x)}{dx} = \begin{cases} \alpha e^x + \alpha, & x < 0 \\ 1, & x \geq 0 \end{cases}$$

Observations:

- It can produce –ve output
- It can blow up activation function

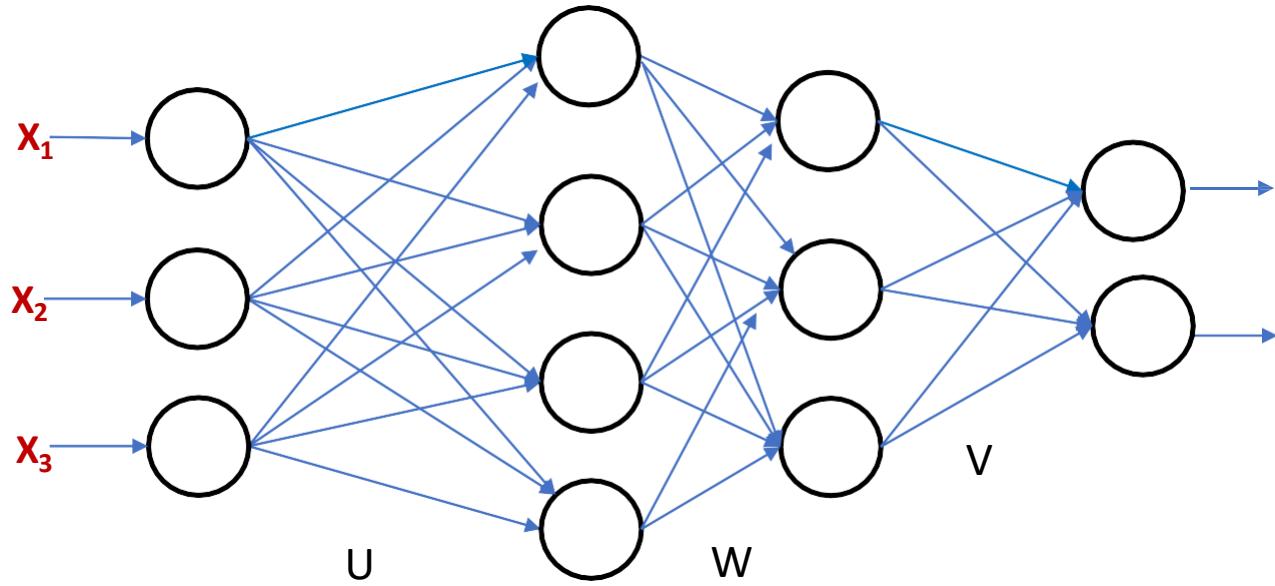


Complete Chain

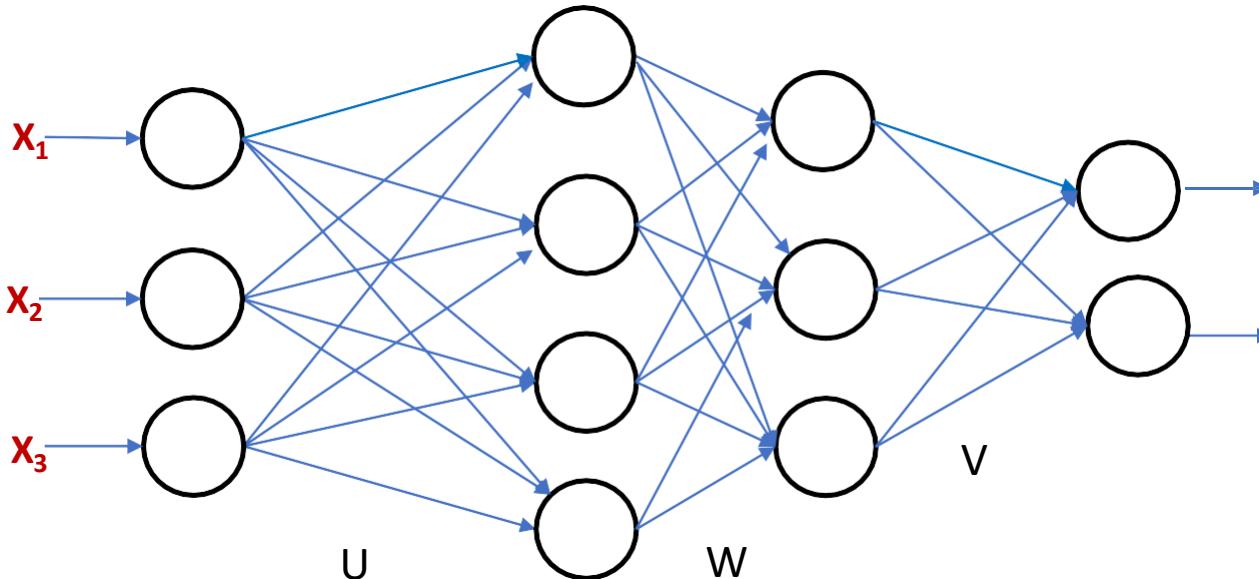


$$\frac{\delta E}{\delta V_{11}} = \frac{\delta z_1}{\delta V_{11}} \times \frac{\delta y_1}{\delta z_1} \times \frac{\delta E}{\delta y_1}$$

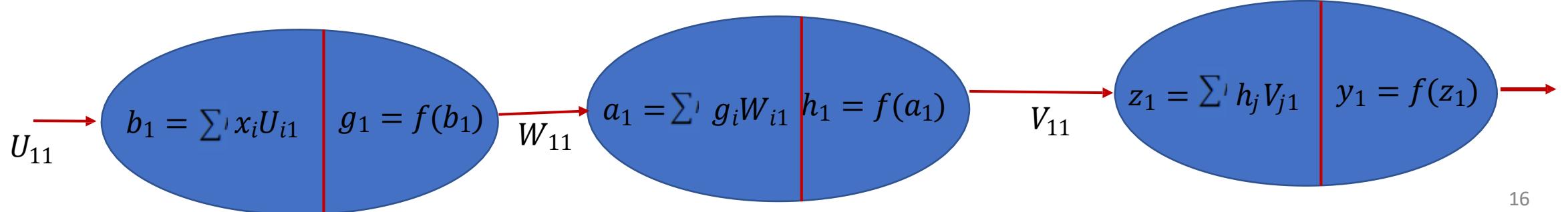
Deep Network

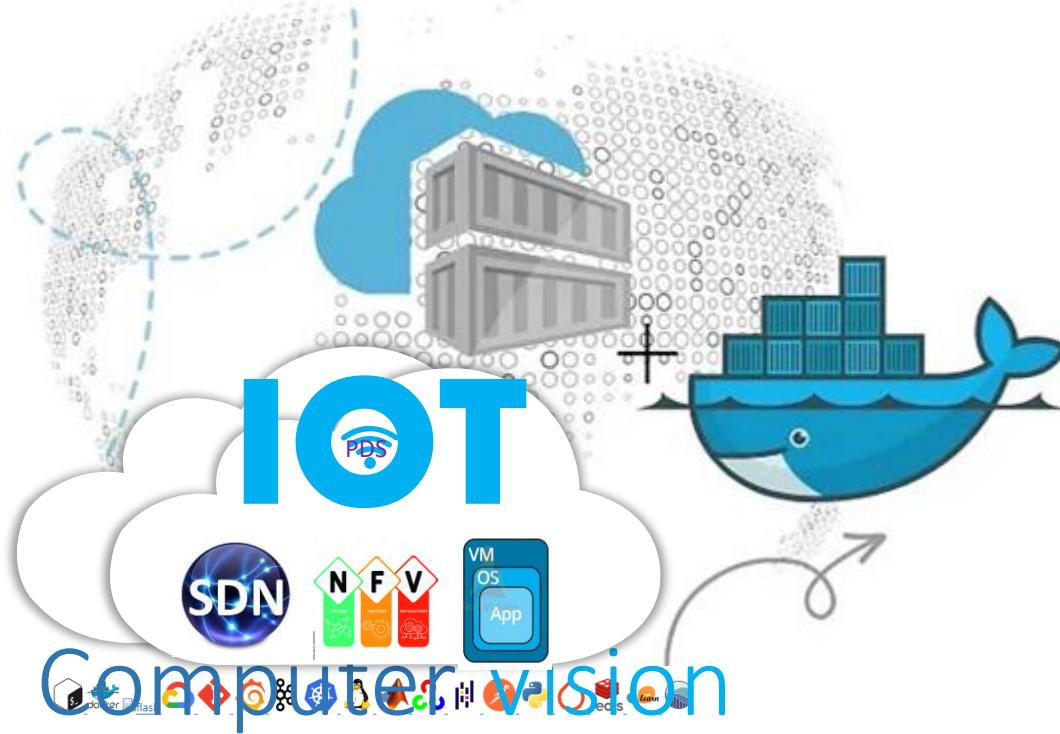


Deep Network - Vanishing/Exploding Gradient



$$\frac{\delta E}{\delta U_{11}} = \frac{\delta b_1}{\delta U_{11}} \times \frac{\delta g_1}{\delta b_1} \times \frac{\delta a_1}{\delta g_1} \times \frac{\delta h_1}{\delta a_1} \times \frac{\delta z_1}{\delta h_1} \times \frac{\delta y_1}{\delta z_1} \times \frac{\delta E}{\delta y_1}$$





Thank you!

Thank you For being a great class!

<https://github.com/siagianp>

https://github.com/amelcharolinesgn2/IoT_simulator-mqtt-NodeRed

<https://www.youtube.com/@AmelOline/videos>