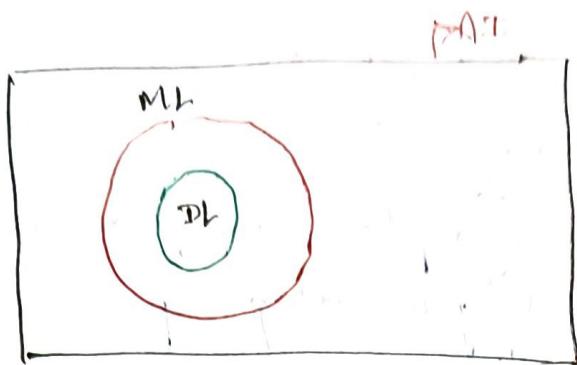


14 - Jan - 2023

## Deep Learning



- Can we make machine learn similarly like how we human being learning?
- Yes, scientist found a technique called deep learning which mimics the human brain with the help of multilayer neural network.

### Three Neural Network in Deep learning!

① ANN  $\rightarrow$  Artificial Neural Network  $\begin{matrix} \xrightarrow{\text{classification}} \\ \xrightarrow{\text{Regression}} \end{matrix}$

- Image classification can be solved
- used with tabular network

② CNN  $\rightarrow$  Convolutional Neural Network

Input - I/P will be Images, Videos frames

Eg: RCNN, masked RCNN, Detection, Yolo, ViT

- Object detection
- Computer vision

③ RNN  $\rightarrow$  Recurrent Neural Network

I/P will be in the form of time series or text data

Eg: LSTM RNN, RNN, GRU, Bidirectional LSTM RNN, Encoder Decoder, Transformers, BERT, Attention Model

## Two libraries:

① Pytorch

② Tensorflow (Keras is a part of Tensorflow)

Q why DL becoming very popular?

Ans In 2005, Facebook, Orkut (Social Media Platform)

People started uploading photos, text etc.  
Data started generating exponentially

Later on many platforms came

Problem: How so they can store these data efficiently?

In 2011 → Big Data is becoming very famous.

Platform like Hadoop, HDFS etc able to understand unstructured data.

In 2011-2016 → Amazon, Google, Facebook etc started storing data

S3 Bucket → AWS & stored data in my format

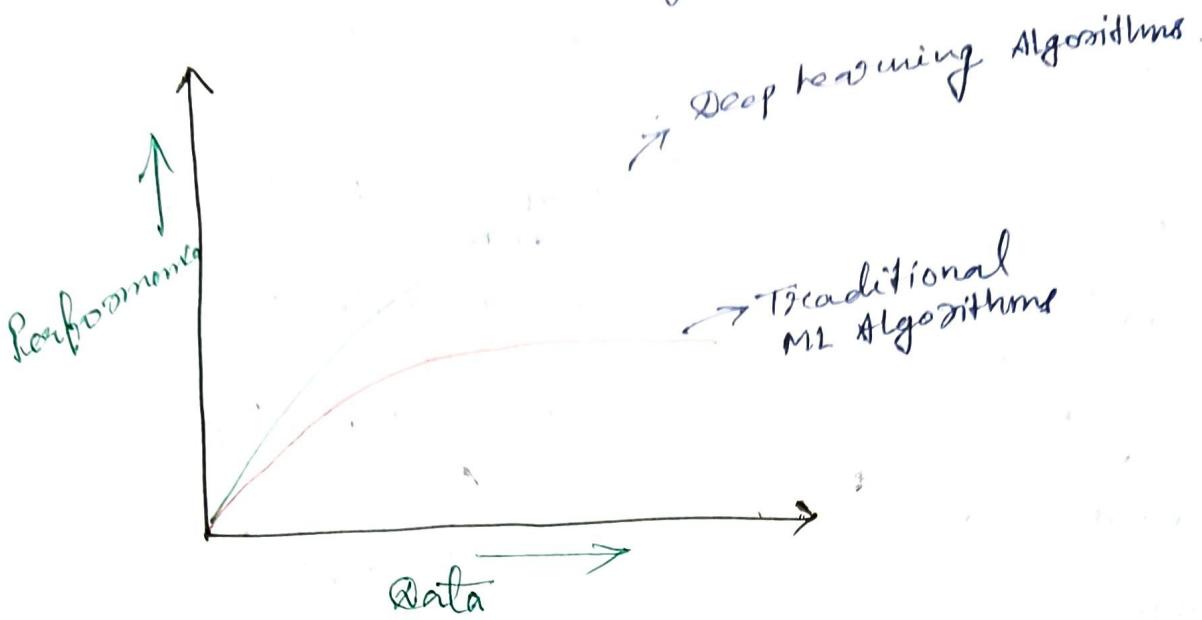
In 2015 → Problem: what will they do with the stored data?

Can we use the data to make the product better?

→ Lot of research more happening  
Deep learning comes into picture

→ For training huge data → Huge no. of GPUs required.  
Hardware requirement: GPUs  
→ GPU price started decreasing.

> Huge amount of data is getting generated and it is suitable for deep learning model.



- > Traditional ML algorithms get stagnate after sometime, but this is not the case in DL.
- > DL is been used in many domains.
- ① Medical : Prediction of disease, X-rays, Bone break, brain scan - 3D image. (detect tumor or not?)
- These points shows depression or not.
- ② E-commerce
- ③ Retail
- ④ Logistic

> Complex problem can be solved using DL.

## ① Perception

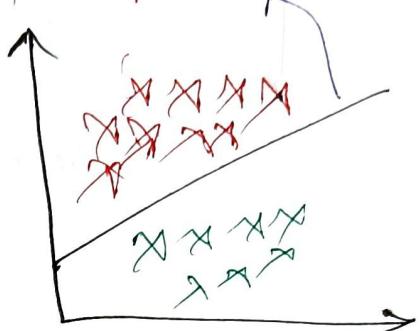
[Artificial Neuron or Neural Network Model]

single layer Neural Network

→ only need for binary classification.

Binary classification

Linear separable



① I/P layer

② hidden layer

③ weights

④ activation function

⑤ Bias

⑥ O/P layer

> Perception plays an important role.

> Neuron sends and process signal and send to the brain.

> Eyes are I/P layer for neuron. (Brain)

> When signal passes to eyes, lot of neurons are interconnected and each neuron process and transmit the signal to the brain.

> Our body has already gone through training since born but in case of a machine, we have to train.

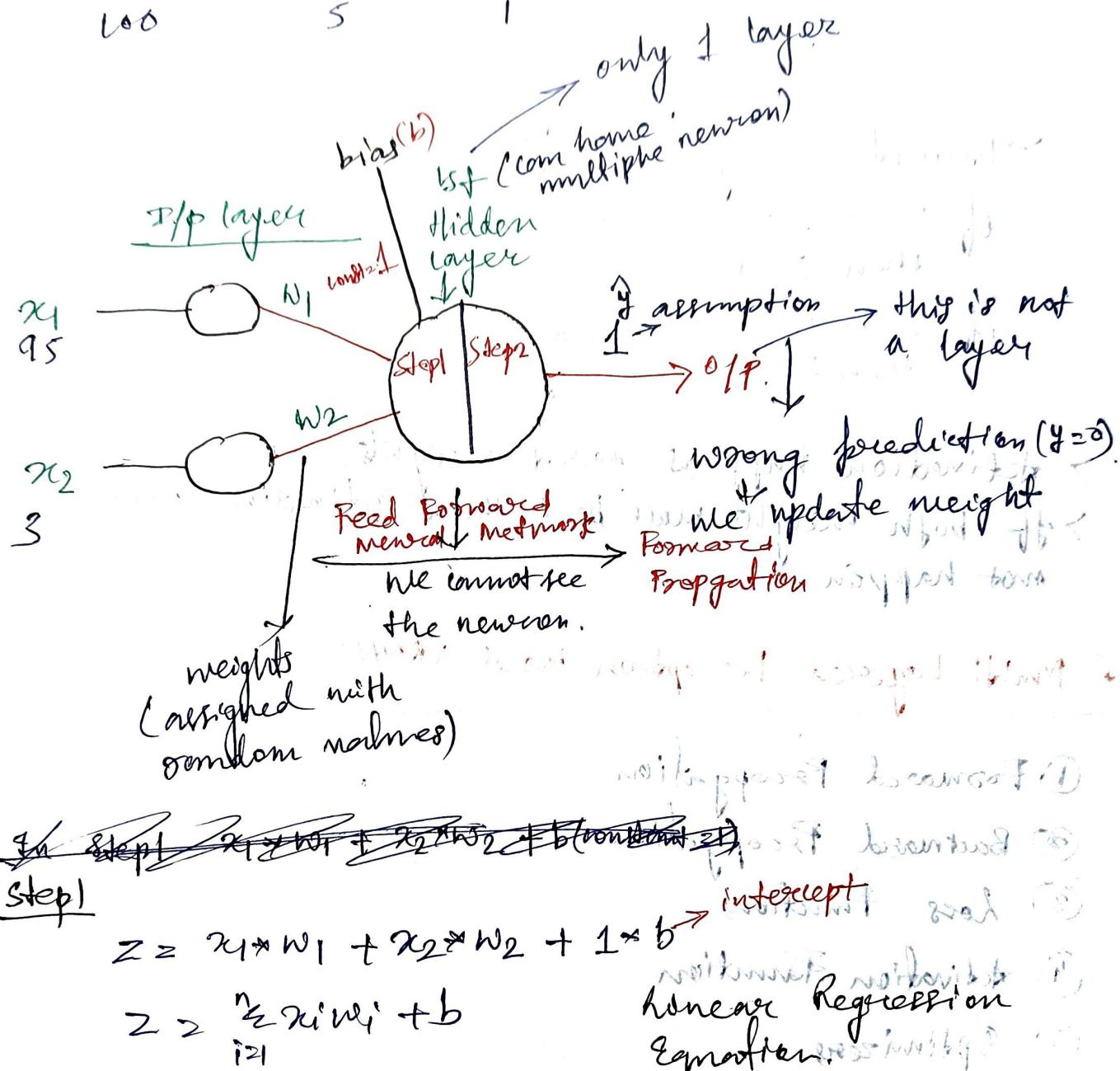
> Hidden layer are interconnected neuron.

> Output layer will be one reaction.

Eg.

Dataset

$x_1$ $\Sigma x$	No. of Study hours	O/P Pass/Fail
95	3	101
110	4	1
100	5	1



~~> In Step1  $z = w_1 x_1 + w_2 x_2 + b$~~

Step1

$$z = x_1 * w_1 + x_2 * w_2 + 1 * b$$

$$z = \sum_{i=1}^n x_i w_i + b$$

Intercept

Linear Regression  
Equation

> this creates a linear line so we can solve linear separable problem.

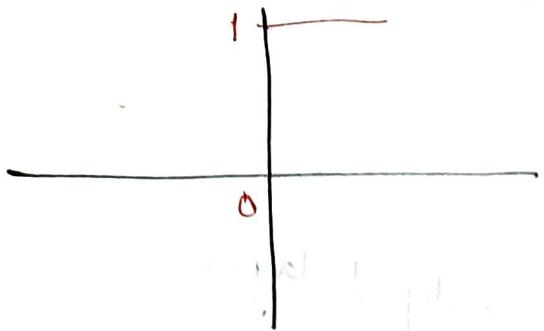
> 'b' will be intercept

Step2

> we apply activation function, called step function.

## Step Function

> Threshold = 0.



> Transforming  $Z_i$  to  $[0, 00, 1]$ .

if  $Z = 5, \geq Z \geq 1$   
then  $Z = 1$

if  $Z = -2, \leq Z < 0$   
then  $Z = 0$

> Activation happens w.r.t. weights.

> If both weight will be 0 then activation will not happen..

## Multi Layered Perception Model (ANN)

① Forward Propagation

② Backward Propagation

③ Loss Function

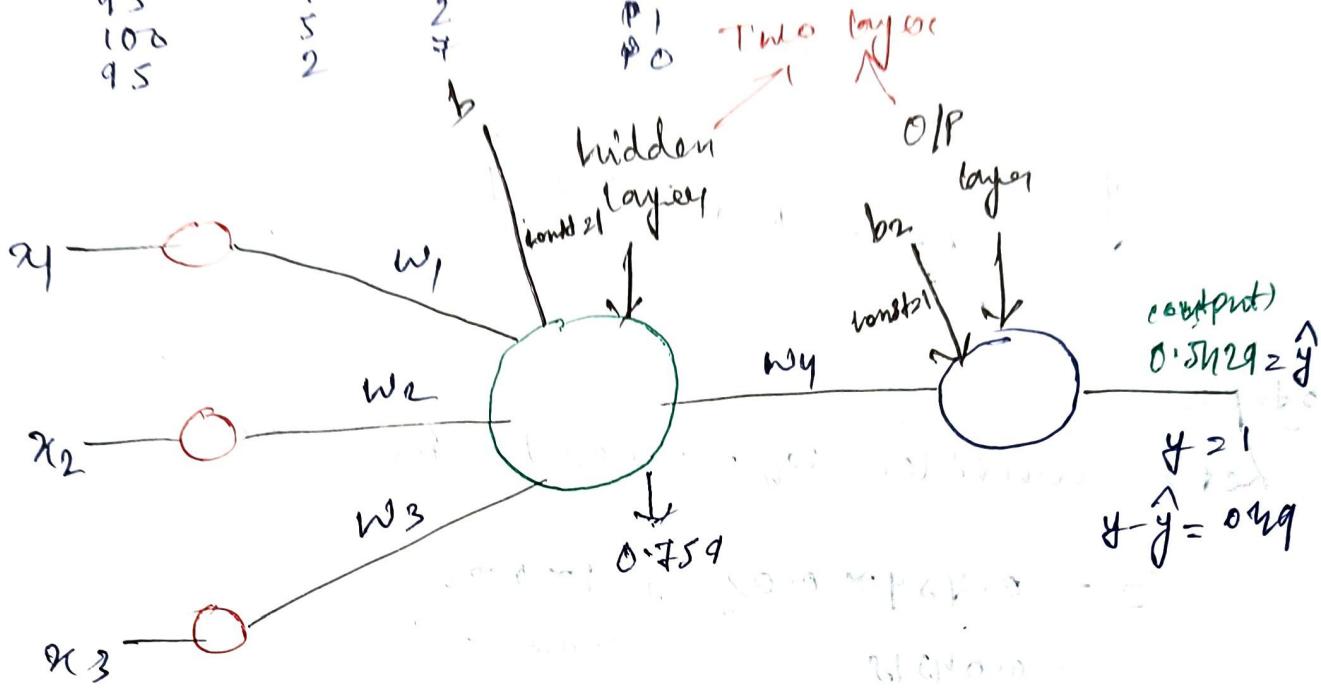
④ Activation Function

⑤ Optimizers

## 2 layer neural network

Advantages - same in 1 layer vs neural network.

(x <sub>1</sub> ) 2 or .	(x <sub>2</sub> ) Study hrs	(x <sub>3</sub> ) Play hrs	O/P (Pass/Fail)
95	4	4	P
100	5	2	P
95	2	7	P



let's consider  $\begin{bmatrix} w_1 & w_2 & w_3 \\ 0.01 & 0.02 & 0.03 \end{bmatrix}$  and  $b = 0.01$

Forward Propagation

Step 1

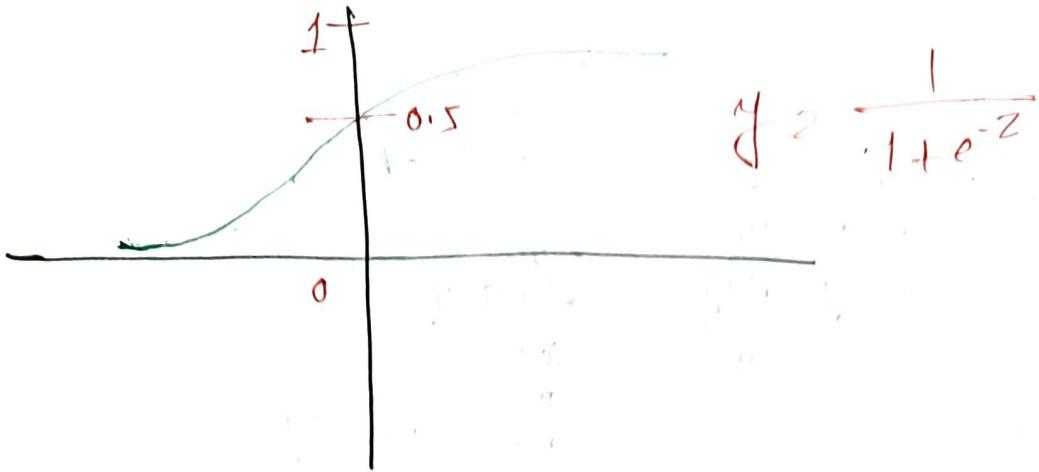
$$z = 95 \times 0.01 + 4 \times 0.02 + 4 \times 0.03 + 0.01 = 1.151$$

Step 2

Apply activation function.

Sigmoid activation function - Transformed  $z$  into 0, 1.

$$\text{sigmoid}(z) = \frac{1}{1 + e^{-z}} = \frac{1}{1 + e^{-1.151}} = 0.759$$



$$y = \frac{1}{1+e^{-0.759}} = 0.759$$

Step 3

Let's consider  $w_1 = [0.02]$   $b_2 = 0.03$

$$z = 0.759 \approx 0.02 + 1 \approx 0.03$$

↓  
constant

Step 4

Apply Activation function to z addition that

$$y = \frac{1}{1+e^{-z}}$$

$$z = \frac{1}{1+e^{-0.759+0.03}}$$

$$= 0.5429$$

Step 5:

$$y = 0.5429$$

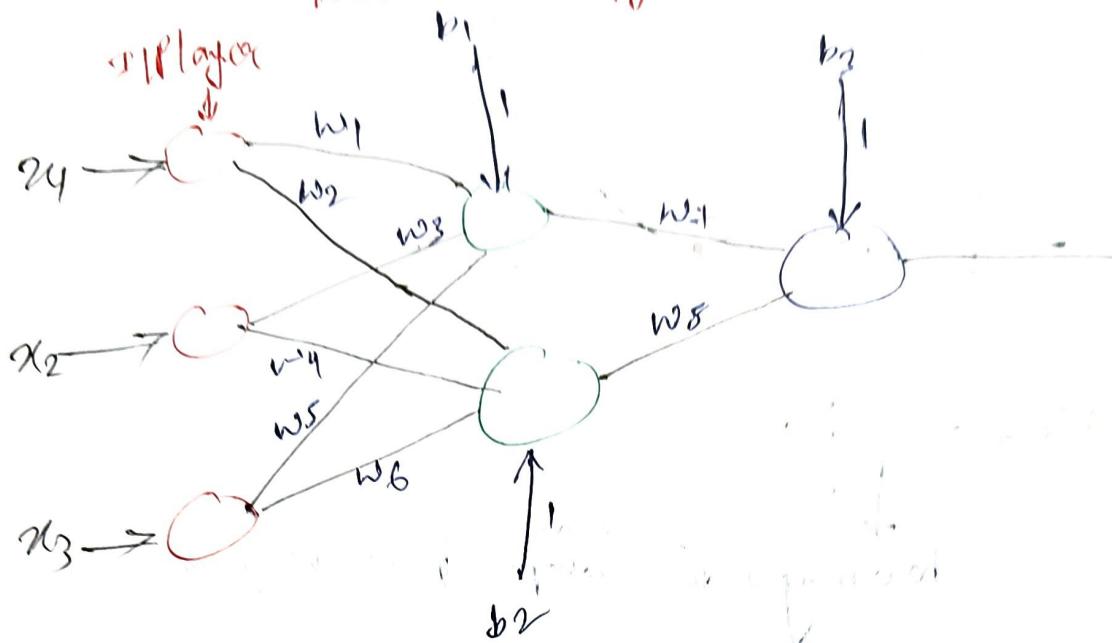
$$y = 1$$

$$\boxed{\text{loss}} = y - \hat{y} = 1 - 0.5429 = 0.459 \quad (\text{let } \hat{y} \text{ on higher side})$$

> Since loss is high, we will perform backward propagation by changing all the weights

### Backward Propagation and updation formula

#### Forward Propagation



#### Backward Propagation

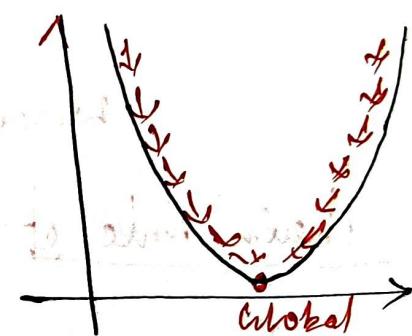
$$\text{loss function} = (y - \hat{y})^2$$

#### Weight Updation Formula

$$w_{\text{new}} = w_{\text{old}} - \eta \left[ \frac{\partial L}{\partial w_{\text{old}}} \right]$$

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\frac{\partial \text{loss}}{\partial w_{\text{old}}}}{\frac{\partial \text{loss}}{\partial w_{\text{old}}}}$$

learning rate  $\eta$   $\rightarrow$  gradient



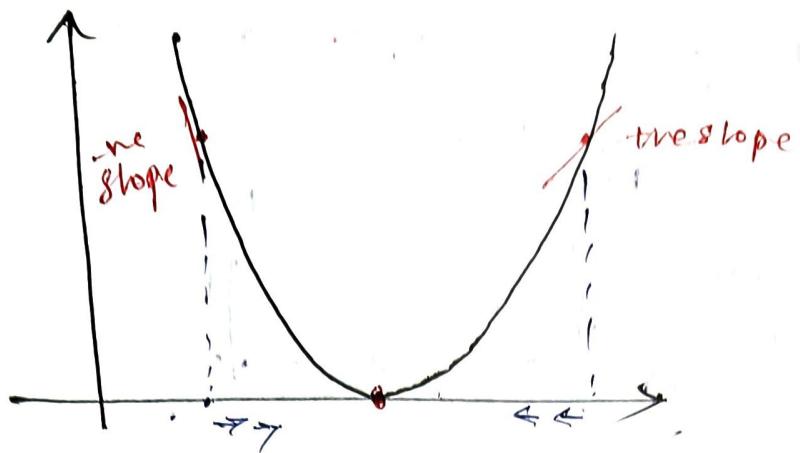
$$b_{\text{new}} = b_{\text{old}} - \eta \frac{\partial \text{loss}}{\partial b_{\text{old}}}$$

Bias updatation formula.

> Weight updation can be done with the help of optimizers like gradient descent.

## Optimizers

• To reduce the loss value.



### -ve slope

$$W_{\text{new}} = W_{\text{old}} - \eta \frac{\partial \text{loss}}{\partial W_{\text{old}}}$$

learning rate  $\eta$  > slope of conergence

### -ve slope

$$W_{\text{new}} > W_{\text{old}}$$

### the slope

$$W_{\text{new}} = W_{\text{old}} - \eta \frac{\partial \text{log}}{\partial W_{\text{old}}}$$

$$W_{\text{new}} < W_{\text{old}}$$

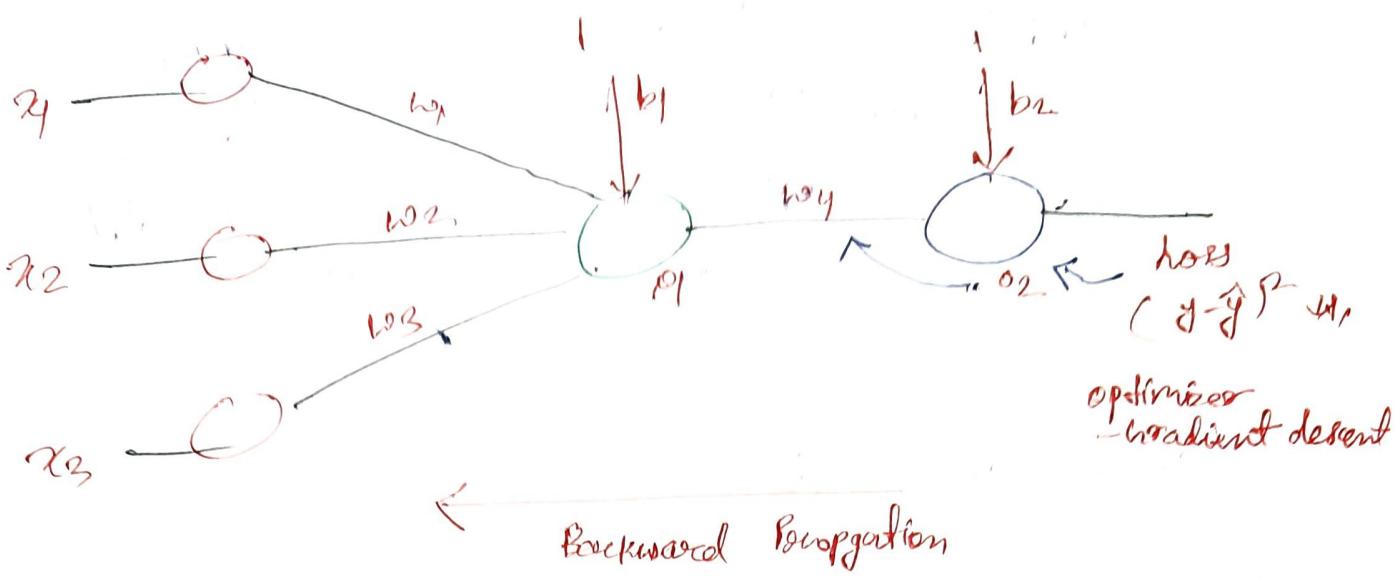
### chain rule of coordinate

$$\frac{\partial \text{loss}}{\partial W_{\text{old}}} = \frac{\partial \text{loss}}{\partial W_{\text{new}}} \cdot \frac{\partial W_{\text{new}}}{\partial W_{\text{old}}}$$

gradient working with  $\left\{ \begin{array}{l} \text{loss} \\ \text{data} \end{array} \right\}$  in each iteration

repeat until loss with new weight till point of  
desired threshold with reasonable

## chain Rule of Backpropagation



$$w_{\text{new}} = w_{\text{old}} - \gamma \left( \frac{\partial \text{loss}}{\partial w} \right) \rightarrow \text{slope}$$

> loss is dependent on  $\theta_2$ ,  $\theta_2$  is dependent on  $w_4$  etc.

$$\frac{\partial \text{loss}}{\partial w_{\text{old}}} = \frac{\partial \text{loss}}{\partial \theta_2} \times \frac{\partial \theta_2}{\partial w_{\text{old}}}$$

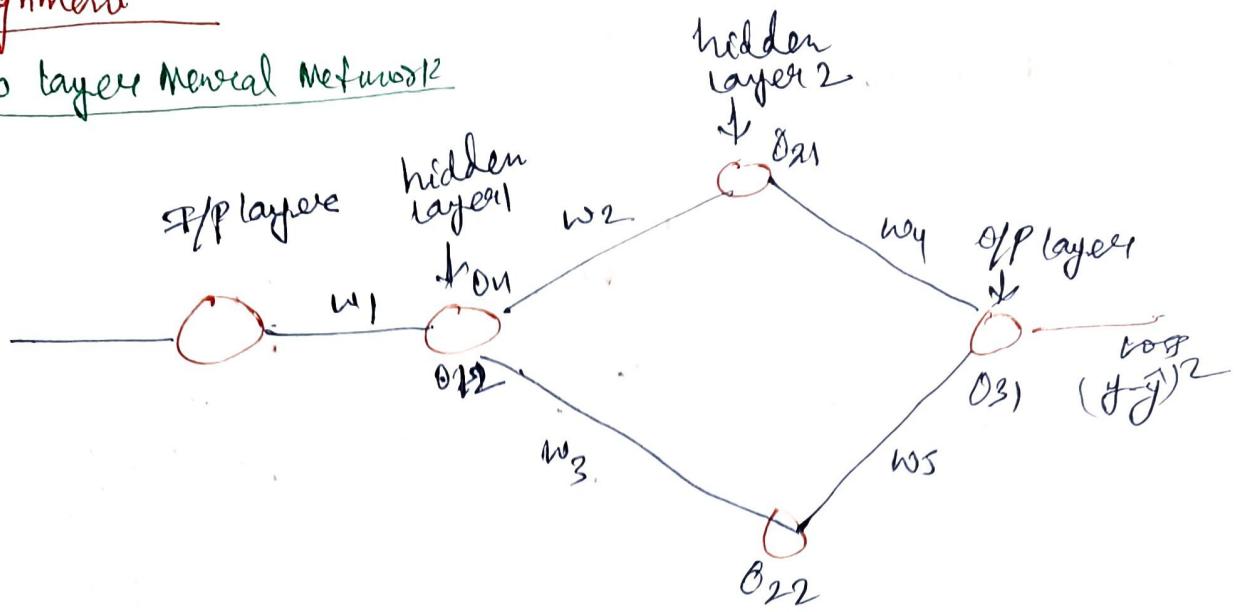
> chain Rule of Backpropagation

$$w_{\text{new}} = w_{\text{old}} - \gamma \frac{\partial \text{loss}}{\partial w_{\text{old}}}$$

$$w_{\text{new}} = w_{\text{old}} - \gamma \frac{\partial \text{loss}}{\partial \theta_2} \times \frac{\partial \theta_2}{\partial w_{\text{old}}} \times \frac{\partial \theta_1}{\partial w_{\text{old}}}$$

## Assignment

### Two layer Neural Network<sup>12</sup>



$$\text{When } \frac{\partial \text{loss}}{\partial \text{word}} = \frac{\partial \text{loss}}{\partial \text{word}}$$

> there will be two path

#### 1st Path

$$\left[ \begin{array}{c} \frac{\partial \text{loss}}{\partial \text{word}} \\ \end{array} \right]_1 = \left[ \begin{array}{c} \frac{\partial \text{loss}}{\partial o_{31}} \xrightarrow{\frac{\partial o_{31}}{\partial o_{21}}} \frac{\partial o_{21}}{\partial o_{11}} \xrightarrow{\frac{\partial o_{11}}{\partial w_1}} \frac{\partial \text{loss}}{\partial w_1} \\ \end{array} \right]$$

#### 2nd Path

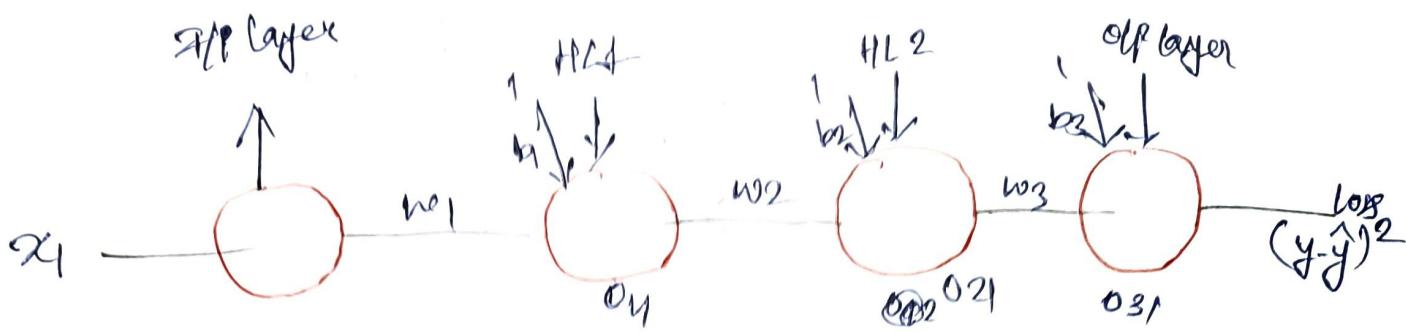
$$\left[ \begin{array}{c} \frac{\partial \text{loss}}{\partial \text{word}} \\ \end{array} \right]_2 = \left[ \begin{array}{c} \frac{\partial \text{loss}}{\partial o_{31}} \xrightarrow{\frac{\partial o_{31}}{\partial o_{22}}} \frac{\partial o_{22}}{\partial o_{12}} \xrightarrow{\frac{\partial o_{12}}{\partial w_3}} \frac{\partial \text{loss}}{\partial w_3} \\ \end{array} \right]$$

$$\left[ \begin{array}{c} \frac{\partial \text{loss}}{\partial \text{word}} \\ \end{array} \right] = \left[ \begin{array}{c} \frac{\partial \text{loss}}{\partial \text{word}} \\ \end{array} \right]_1 + \left[ \begin{array}{c} \frac{\partial \text{loss}}{\partial \text{word}} \\ \end{array} \right]_2$$

>  $w_4$  is dependent on both the path. so they have been added.

# Vanishing gradient Problem and Activation functions

Three layer neural network



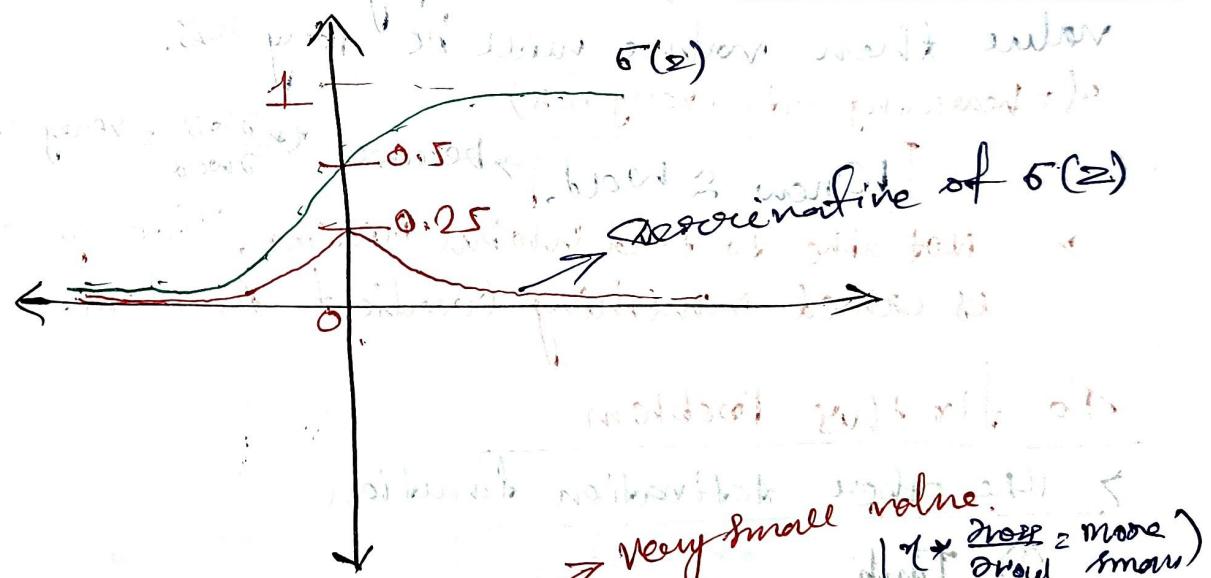
$o_i \rightarrow$  Outputs

Sigmoid Activation

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

$0 \rightarrow 1$

$0 \leq \delta(z) \leq 1$   
Domain of  $\delta(z)$   
 $0 \leq \delta(z) \leq 0.25$



$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial \text{loss}}{\partial w_{\text{old}}}$$

$$\frac{\partial \text{loss}}{\partial w_{\text{old}}} = \frac{\partial \text{loss}}{\partial o_3} \times \frac{\partial o_3}{\partial o_2} \times \frac{\partial o_2}{\partial w_1}$$

$$o_3 = \delta(w_3 \cdot o_2 + b_3) \quad z = w_3 \cdot o_2 + b_3$$

Small value

$$\frac{\partial O_{31}}{\partial O_{21}} = \frac{\frac{\partial (G(z))}{\partial (z)} \rightarrow \frac{\partial (z)}{\partial O_{21}}}{\{0-0.25\}} \quad \left\{ \text{chain rule} \right\}$$

↓

$$= \frac{\partial (w_3 * O_{21} + b_3)}{\partial O_{21}} \rightarrow w_3$$

$$\rightarrow \frac{\partial w_3 O_{21}}{\partial O_{21}} + \frac{\partial b_3}{\partial O_{21}}$$

$$\rightarrow w_3 + 0$$

$$\rightarrow w_3$$

$$\frac{\partial O_{31}}{\partial O_{21}} \rightarrow [0-0.25] * w_3$$

> As the chain rule is very big, derivative value becomes zero. If we multiply all the derivative value then value will be very less.

$\alpha$  = learning rate (very less)

- ( $\Rightarrow$ )  $\rightarrow$  When  $\approx 0$ ,  $\rightarrow$  become  $\frac{\partial z_{31}}{\partial w_{31}} = \text{very small}$
- > Not able to find global minima, This problem is called Vanishing Gradient Problem. -

### To fix this Problem

> Use other Activation function

① Tanh

② ReLU

③ Prelu

④ Elu

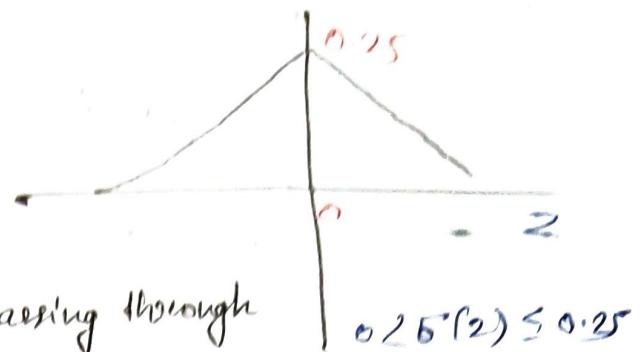
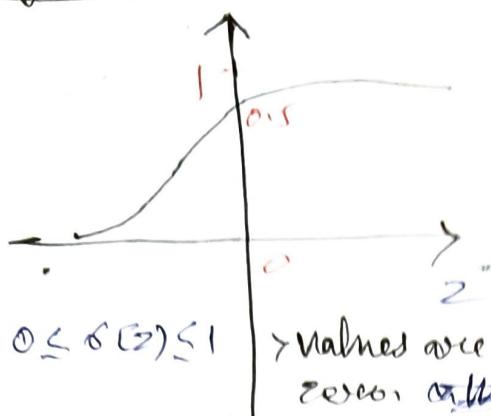
⑤ Swish

$$\text{sigmoid } (est + 180 * \text{act}) \rightarrow 180$$

value limit

## Activation Function

### ① Sigmoid Activation Function



### Advantages

- ① It gives clear prediction (0 to 1).

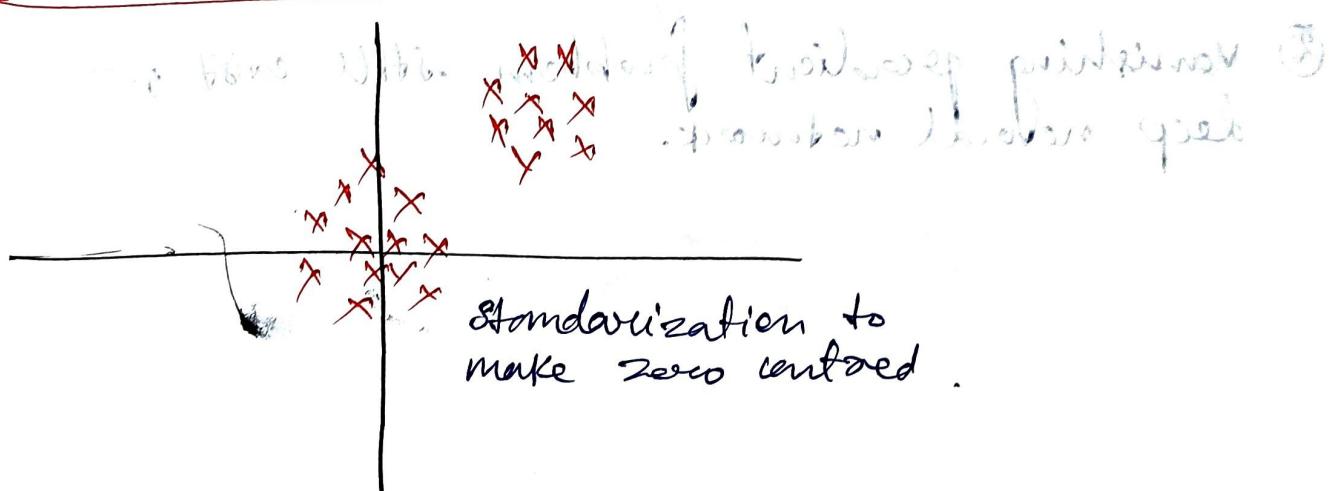
### Advantages

① It is prone to vanishing gradient problem.

② Function output is not zero centered.

*Efficient weight updates.*

### Zero centered



→ Zero centered smoothes the activation function

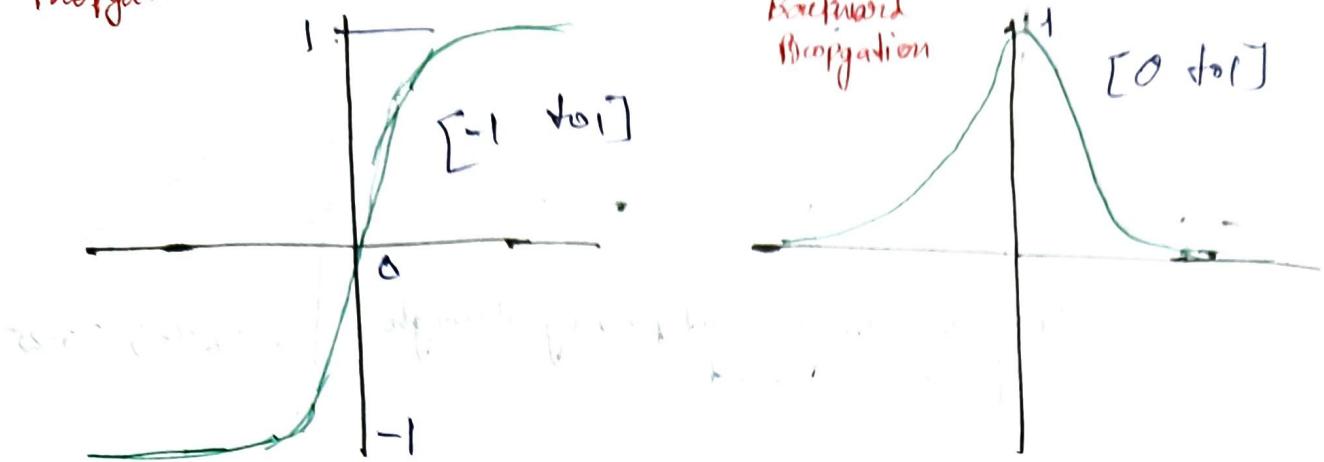
## ② Tanh Activation Function

Forward Propagation

$$\text{Tanh } z = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

(-1 to 1)

Backward  
Propagation



### Advantages

- ① It is zero centered.

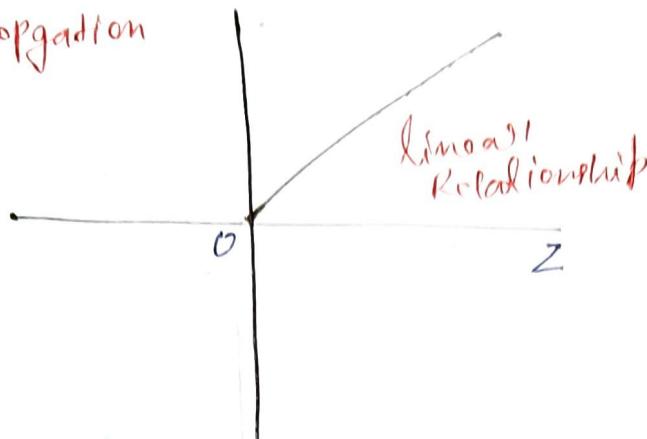
### Disadvantages

- ② More time complexity due to deactivating.
- ③ Vanishing gradient problem still exist for deep neural network.

### ③ ReLU Activation Function

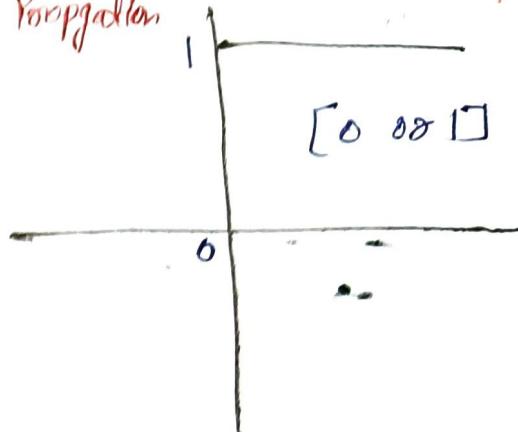
Forward  
Propagation

ReLU: max(0, z)



Backward  
Propagation

$d\text{ReLU}(z)/dz$



#### Advantages

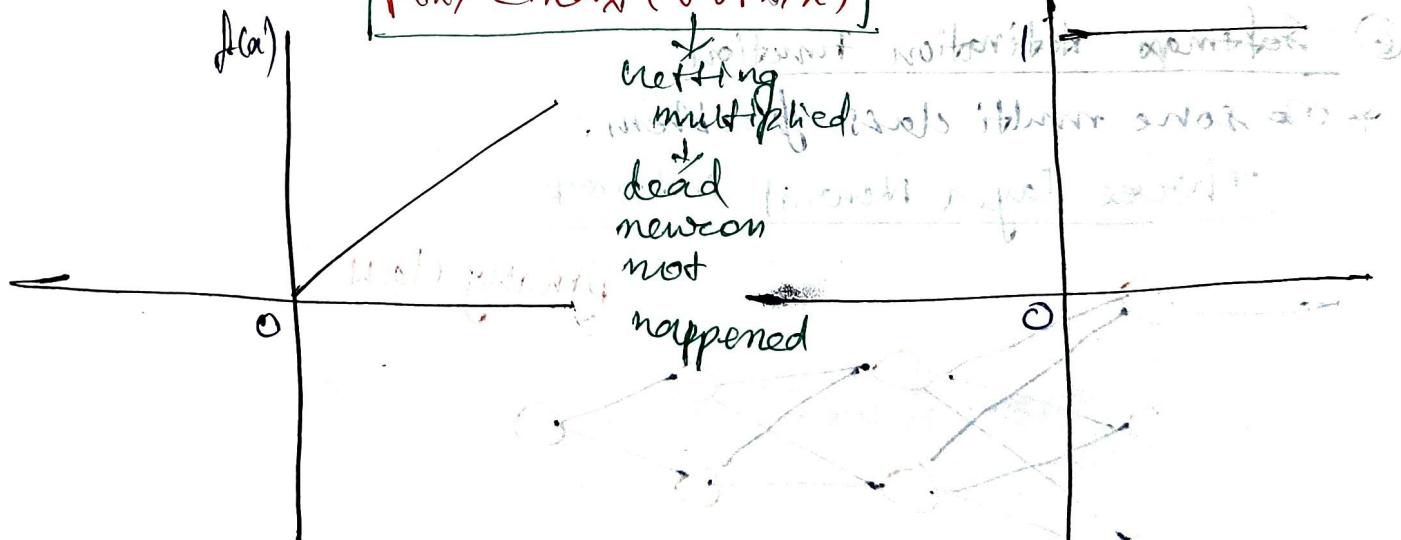
- ① Solved vanishing gradient problem.
- ② It is faster than tanh and sigmoid.

#### Disadvantages

- ① If derivative is zero, then it becomes dead neuron which will skip the layer.
- ② It is not zero centric.

### ④ Leaky ReLU or Parameteric ReLU

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



#### Advantages

- ① Dead neuron problem is solved.

#### Disadvantages

- ① It is not zero centric

## ⑤ ELU (Exponential linear Unit)

Forward  
Propagation

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

hyperparameter

Back  
Propagation



gradient flow along with back propagation

Example: Loss function for ELU

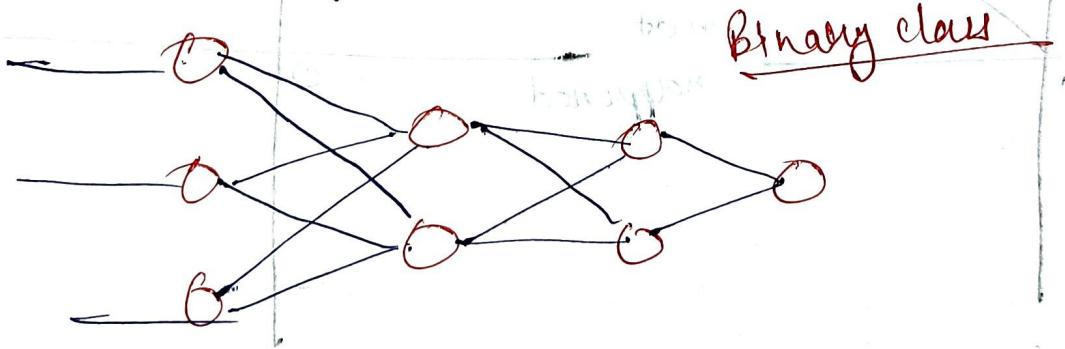
### Advantages

- ① It is zero centered
- ② Vanishing gradient problem got solved.

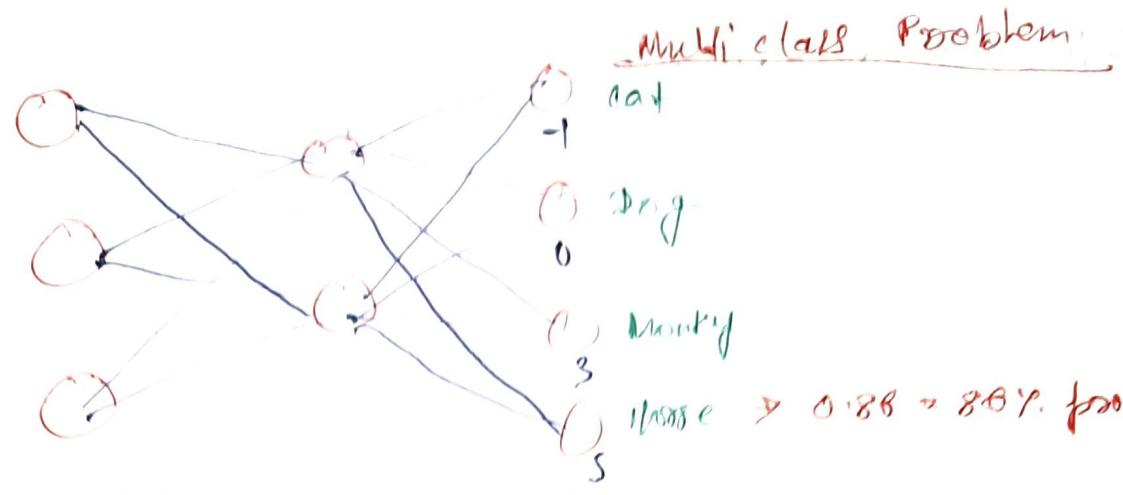
### Disadvantages

- ① More time complexity due to complex derivative
  - ② Softmax Activation Function
- To solve multi class problem.

### Three layer Neural Network



Network is mostly used for binary classification  
softmax error loss is used



$$\text{softmax} = \frac{e^{y_i}}{\sum_{k=0}^n e^{y_k}}$$

~~cat =  $\frac{e^{-1}}{e^{-1} + e^0 + e^3}$~~   
output after applying softmax function.

$$\text{cat} = \frac{e^{-1}}{e^{-1} + e^0 + e^3} = 0.00033$$

$$\text{Dog} = \frac{e^0}{e^{-1} + e^0 + e^3} = 0.0024$$

$$\text{Monkey} = \frac{e^3}{e^{-1} + e^0 + e^3} = 0.0183$$

$$\text{Horse} = \frac{e^5}{e^{-1} + e^0 + e^3} = 0.1853 \quad (\text{Highest})$$

### Probability

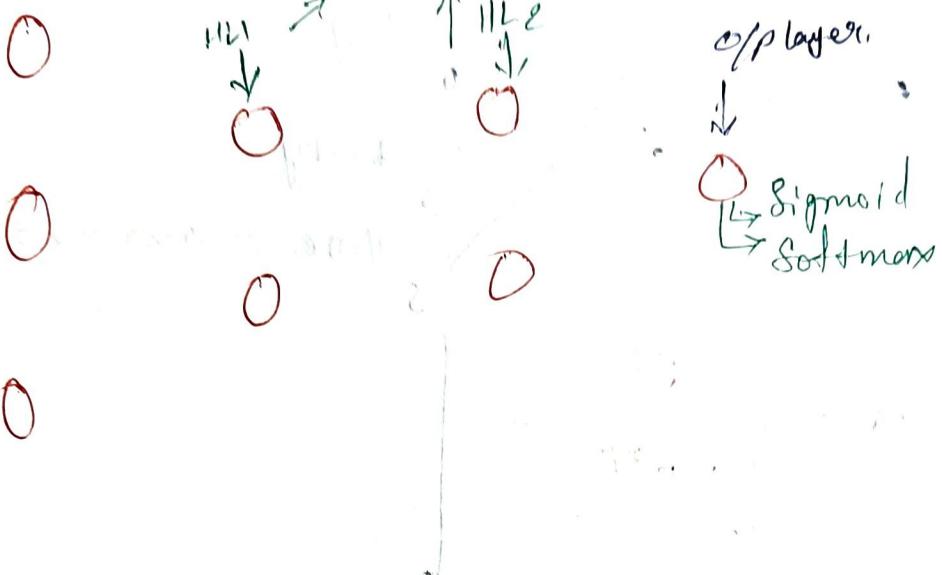
$$P(\text{Horse}) = \frac{0.1853}{0.00033 + 0.0024 + 0.0183 + 0.1853} \geq 0.86$$

→ In case of multi class problem, we need to use softmax activation function.

Q Which activation function to use when?

Ans

ReLU and its variation



Sigmoid  $\rightarrow$  binary classification

Softmax  $\rightarrow$  multi-class classification

Regression

① Linear Activation Function

$$\text{Output} = \frac{\theta_0 + \theta_1 x}{\theta_2 + \theta_3 x} = \text{postural}$$

$$\text{Output} = \frac{\theta_0 + \theta_1 x}{\theta_2 + \theta_3 x} = 0.999$$

$$0.999 \times 100 = 99.9 \text{ postural}$$

so what we want our model to do is to give us a pattern matching