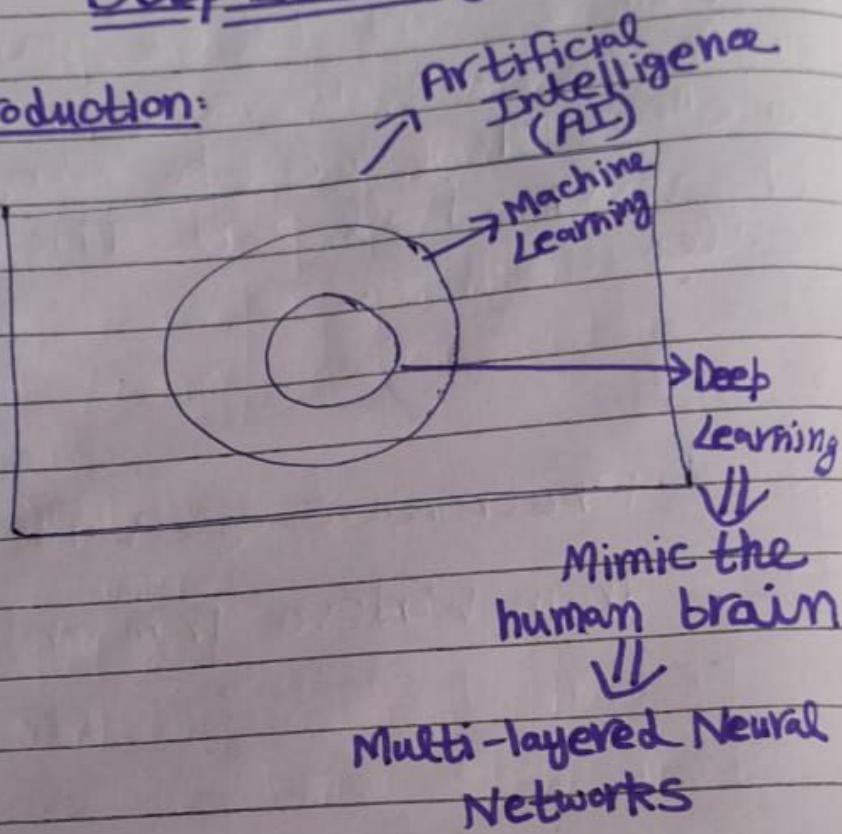


Deep Learning

Introduction:



Deep Learning:

① ANN → Artificial Neural Network → Regression
→ Classification

② CNN → Convolutional Neural Network

↳ Input data is in form of images and video frames

↳ RCNN, YOLO V5 etc.

used for object Detection in CV.

③ RNN → Recurrent Neural Network
↳ NLP → Input: Text, Time Series
↳ { Word Embedding, LSTM RNN, GRU RNN,
 Bidirectional LSTM RNN, Encoder-
 Decoder, Transformers, BERT }

Framework to be used: TensorFlow

↓
Developed by Google

* Why Deep Learning is becoming
popular?

2005 → Facebook, YouTube, WhatsApp, LinkedIn, Twitter

↓

2011 - 2012 →

Data was generated exponentially

(Images, Text, Documents, Videos)

① Hardware Requirement

↓

GPU's price decreasing

(Easy to train DL models) (HADOOP)

↳ Big Data

(Efficient storage and retrieval of data)

② Huge amount of data is getting
generated → Deep learning performs well

③ As the data getting increase
the traditional Machine Learning
Algorithms improve performance upto
a limit but Deep Learning models
keep on increasing performance

④ Deep Learning is being used in many
domains:

- (i) Medical
- (ii) Ecommerce
- (iii) Retail
- (iv) Marketing and many others

⑤ There are famous open source
frameworks for Deep Learning:

- (i) Tensorflow (by Google)
- (ii) Pytorch (by Facebook)

(Artificial Neuron or Neural Network unit)

⇒ Perceptron Intuition:

- Perceptron is a single layered Neural Network.

Some Important things to be discussed here are:

- ① Input Layer
- ② Hidden Layer
- ③ Weights
- ④ Activation Function

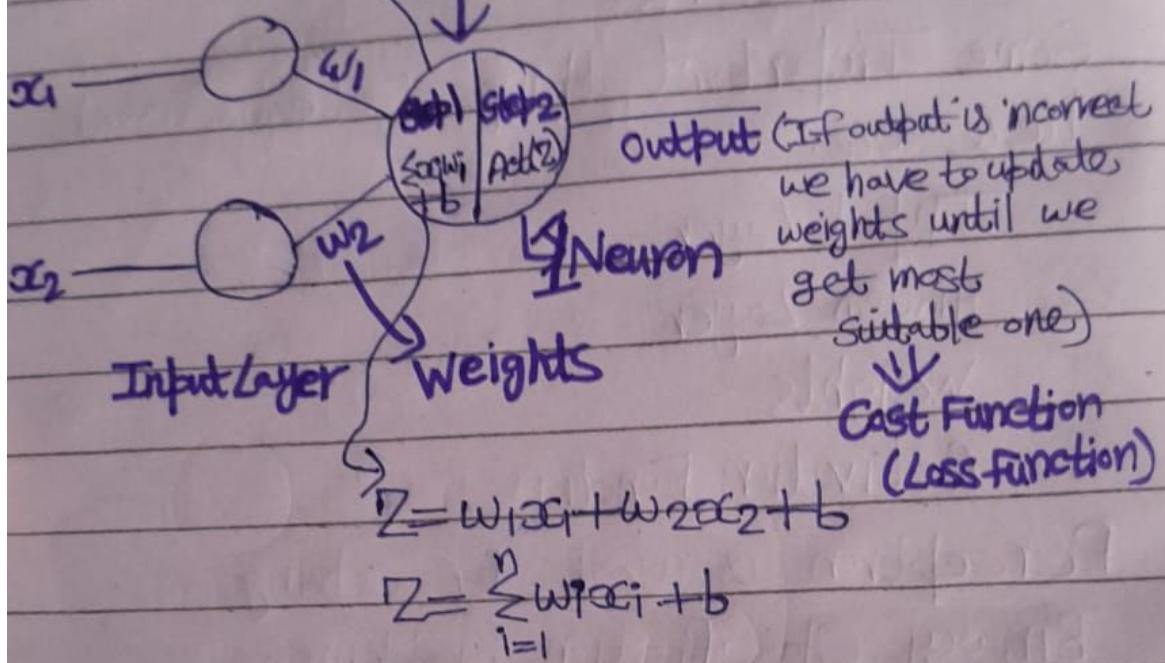
- Perceptron is used for solving Binary Classification problems.

Let we have Dataset:

IG	OC ₁	OC ₂	No. of study hours	Output (Pass/Fail)
95		3		0
110		4		1
100		5		1

As we have 2 input feature so in our Input Layer we will have 2 inputs

Noise \leftarrow bias (b) Hidden Layer

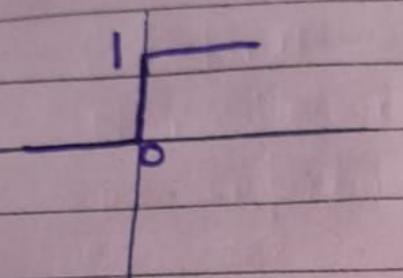


Activation Function:

In step 2, we will apply activation function on Σ . It is used to transform the output i.e from 0 to 1 or between 1 to -1.

Some activation functions are:

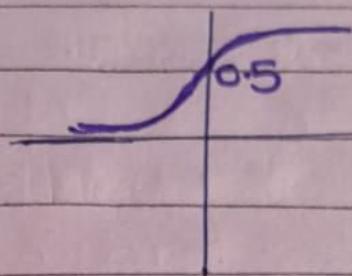
Step Function



(threshold value = 0)

$$\begin{cases} 0 \text{ if } z \leq 0 \\ 1 \text{ if } z > 0 \end{cases}$$

Sigmoid Function



(threshold value = 0.5)

$$\begin{cases} 1 \text{ if } z > 0.5 \\ 0 \text{ if } z \leq 0.5 \end{cases}$$

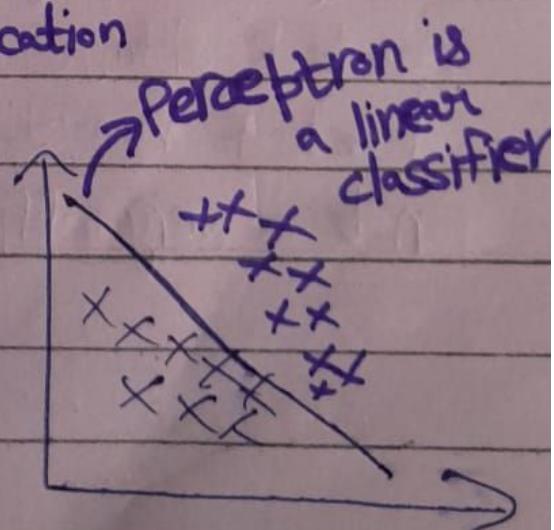
• Step 1:

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

$$z = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$$

↓ It is similar to linear classification

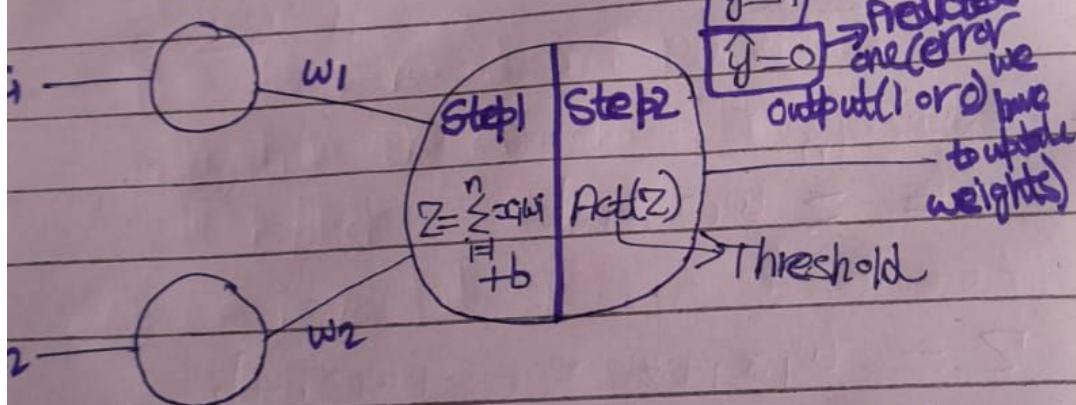
$$y = mx + c$$



Perception Models

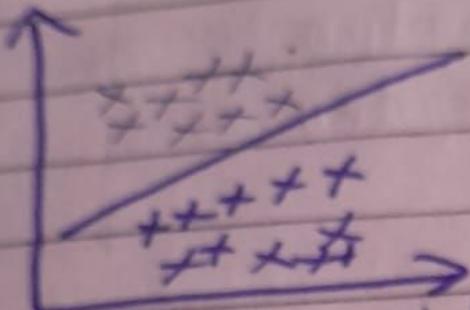
↓
Single Layer perception Model Multi-Layered Perception Model
↓
ANN

Feed Forward Neural Network



* Every time we get error we have to go back change weights and reperform Feed Forward Neural Network. It's not a very efficient technique.

But in Perceptron (Single Layered) it is
only good for Linearly Separable data



* It is not suitable for Non-Linear.
Multi-layered Neural Network
offers:

- ① Forward Propagation
- ② Backward Propagation
- ③ Loss Function
- ④ Activation Functions
- ⑤ optimizers

So Multi-layered Neural Network can be used for more complex problems and perform better.

Multi-layered Neural Network

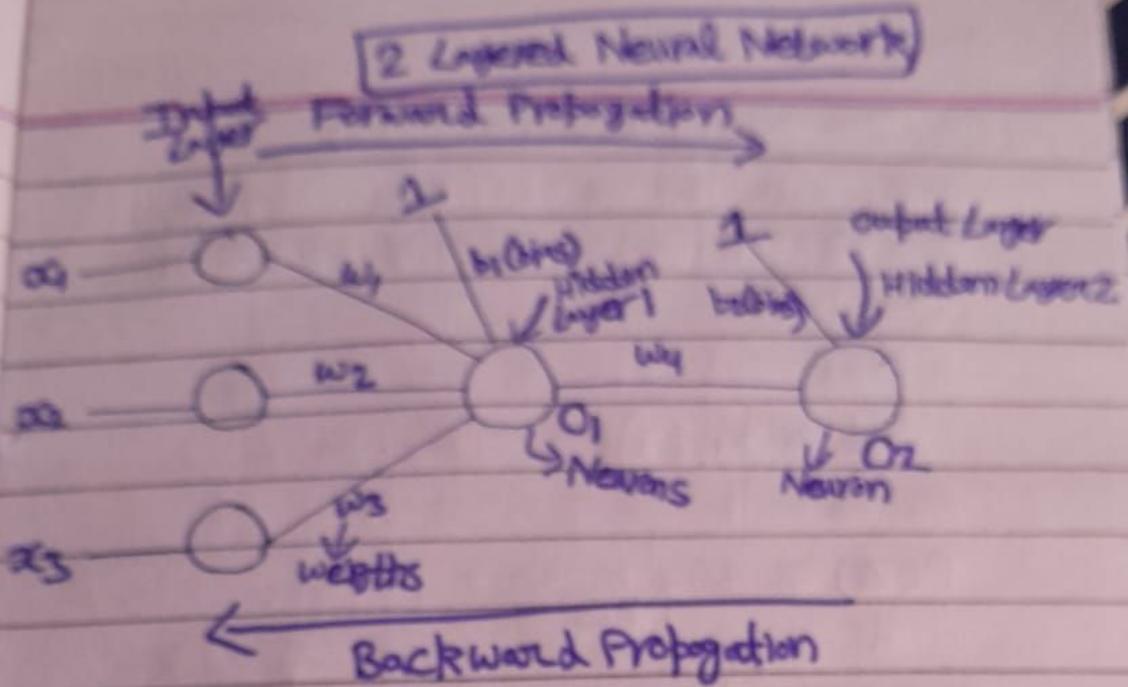
⇒ Artificial Neural Network
(ANN) (Multilayered Perceptron Model)

- ① Forward Propagation
- ② Backward Propagation
- ③ Loss Function
- ④ Optimizers
- ⑤ Activation Function

* Due to Forward and Backward propagation the efficient update of weights will happen

Let we have dataset

α_1	α_2	α_3	O/P (Pass/Fail)
IQ	Study hours	Play hours	
95	4	4	1
100	5	2	1
95	2	7	0



① Forward Propagation:

Hidden Layer 01:

Step 1:

$$Z = x_1 w_1 + x_2 w_2 + x_3 w_3 + b_1$$

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

Let's initialize our weights randomly

$$w_1 = 0.01, w_2 = 0.02, w_3 = 0.03$$

and we are picking 1st record of our dataset

$$x_1 = 95 \quad x_2 = 4 \quad x_3 = 4$$

$$\text{Let } b_1(\text{bias}) = 0.01$$

AS $Z = \alpha_1 w_1 + \alpha_2 w_2 + \alpha_3 w_3 + b_1$

So,

$$Z = 95 \times 0.01 + 4 \times 0.02 + 4 \times 0.03 + 0.01$$

$$\boxed{Z = 1.51}$$

Step 02:

Activation (Z)

Sigmoid activation

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

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

$$f(z) = \frac{1}{1+e^{-1.51}}$$

(Sigmoid activation
keeps the value
b/w 0 and 1)

$$\boxed{O_1 = 0.759}$$

↓
output of Hidden Layer 1

Hidden Layer 2:

Step 1:

Now let $w_4 = 0.02$

and $b_2(\text{bias}) = 0.03$

$$Z = O_1 w_4 + b_2$$

$$Z = (0.759)(0.02) + 0.03$$

$$\boxed{Z = 0.04518}$$

Step 2:

Activation (z)

$$f(z) = \frac{1}{1+e^{-z}} \quad (\text{Sigmoid activation})$$

$$f(z) = \frac{1}{1+e^{-0.04518}}$$

$$f(z) = 0.51129$$

$$\boxed{O_2 = 0.51129}$$

↓
Output of Output Layer/Hidden Layer 2

As we get output Now we will calculate

Loss Function:

$$\hat{y} = O_2 = 0.51129 \quad (\text{Predicted output})$$

$$y = 1 \quad (\text{Actual output of record we taken})$$

$$\text{Loss Function} = y - \hat{y} = 1 - 0.51129$$

$$\boxed{\text{Loss Function} = 0.49}$$

↓
Error

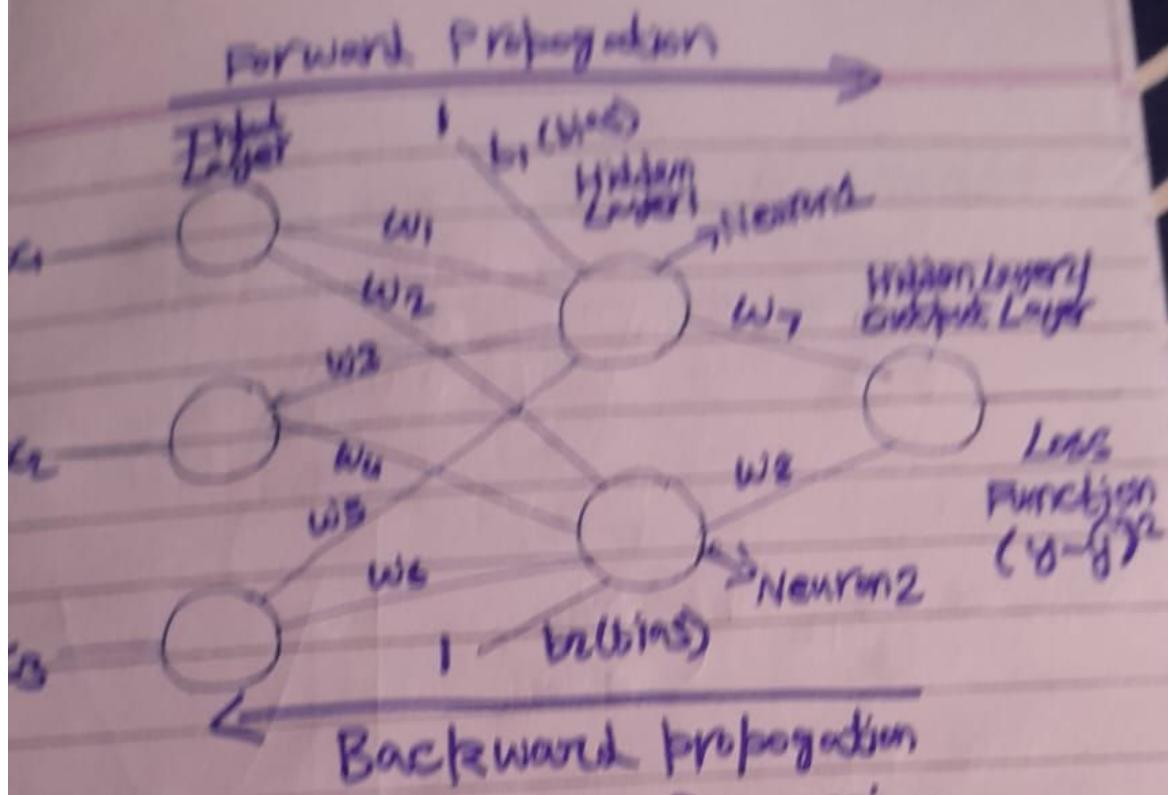
- Our main aim is to reduce error.

- We will reduce error by updating weights in the manner of **backward propagation**.
- After updation of weights by backward propagation, the next record will be passed from dataset and again forward propagation will happen. This forward and backward propagation will happen until we get minimized loss function

> Backward Propagation and Weight Updation Formulas:

Let we have dataset

Input Features			Pass/Fail Output
x_1	x_2	x_3	
IQ	Study hours	Play hours	
95	4	4	1
100	5	2	1
95	2	7	0



Loss Functions for Regression:

- ① MSE
- ② MAE
- ③ Huber Loss

Loss Functions for Classification:

- ① Binary Cross Entropy
- ② Categorical cross Entropy

When loss function is high we reduce it by updating weights in backward propagation.

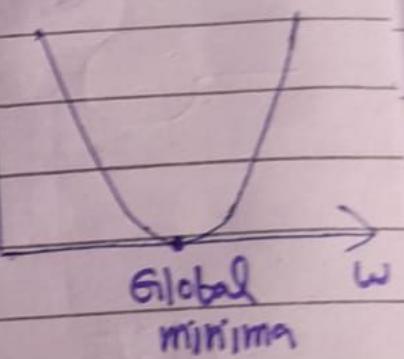
Weight Updation Formula:

$$w_7 \text{new} = w_7 \text{old} - \eta$$

$$\boxed{\frac{\partial L}{\partial w_7 \text{old}}} \rightarrow \begin{array}{l} \text{slope} \\ (\text{derivative} \\ \text{of loss with} \\ w_7) \end{array}$$

$$w_8 \text{new} = w_8 \text{old} - \eta \frac{\partial L}{\partial w_8 \text{old}}$$

Learning Loss
Rate



Generically,

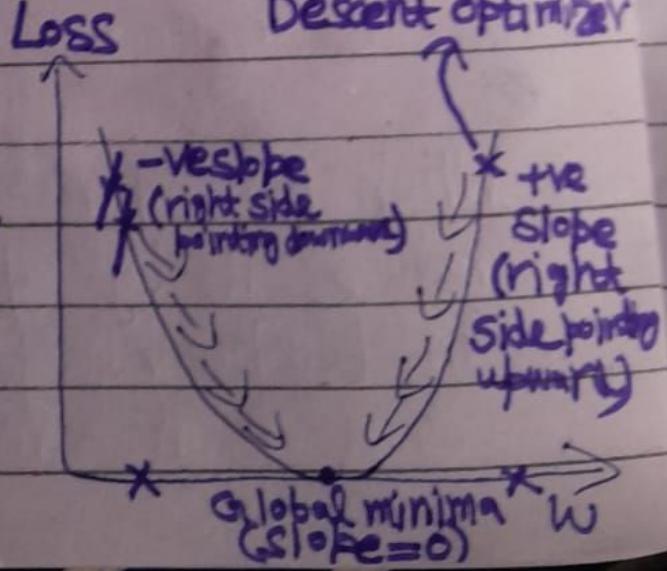
$$w \text{new} = w \text{old} - \eta \frac{\partial L}{\partial w \text{old}}$$

(our aim is to
reach global
minima as slope
is 0 there)

↓
Weight Updation
Formula

Optimizers:

To reduce the
loss value



For -ve slope

$$w_{\text{new}} = w_{\text{old}} - \eta \quad (-\text{ve})$$

$$w_{\text{new}} = w_{\text{old}} + \eta \quad (+\text{ve})$$

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

(Increase weight)

For +ve slope

$$w_{\text{new}} = w_{\text{old}} - \eta \quad (+\text{ve})$$

$$w_{\text{new}} = w_{\text{old}} - \eta \quad (+\text{ve})$$

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

(Decrease weight)

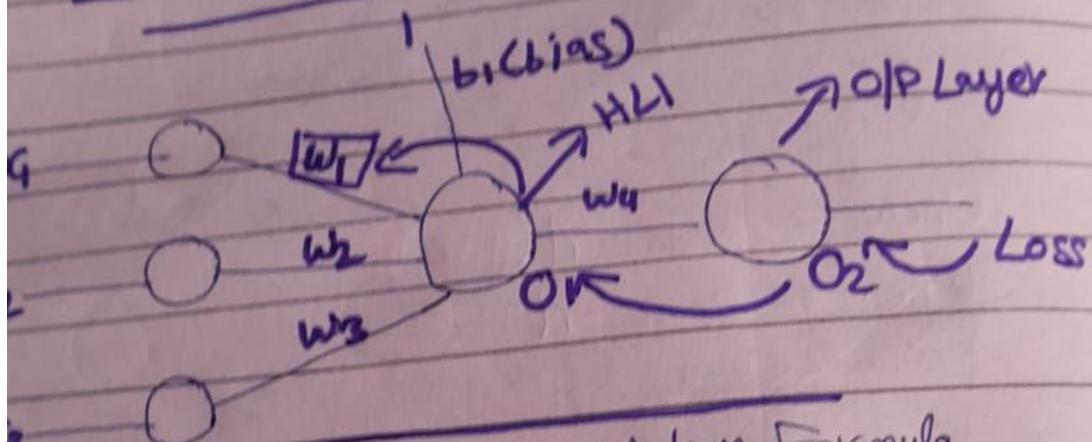
→ Learning rate (η) decides the step size how soon we have to converge so better to have smaller learning rate i.e 0.001

When w reaches global minima

$$\text{slope} \Rightarrow \frac{\partial L}{\partial w_{\text{old}}} = 0$$

$$w_{\text{new}} = w_{\text{old}}$$

Chain Rule of Derivative:



We know Weight update Formula

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

We start backward so first update

w4

$$w_{4\text{new}} = w_{4\text{old}} - \eta \boxed{\frac{\partial L}{\partial w_{4\text{old}}}}$$

↙ Expand

$$\frac{\partial L}{\partial w_{4\text{old}}} = \frac{\partial L}{\partial O_2} \times \frac{\partial O_2}{\partial w_4} \Rightarrow \text{Chain Rule of Derivative}$$

Now update w1

$$w_{new} = w_{old} - \eta$$

$$\boxed{\frac{\partial L}{\partial w_{old}}}$$

$$\frac{\partial L}{\partial w_{old}} = \frac{\partial L}{\partial w_{CO_2}} \times \frac{\partial w_{CO_2}}{\partial w_{old}} \times \frac{\partial w_{O_1}}{\partial w_{old}}$$

$$w_{new} = w_{old} - \eta$$

$$\boxed{\frac{\partial L}{\partial w_{CO_2}}}$$

$$\frac{\partial L}{\partial w_{CO_2}} = \frac{\partial L}{\partial w_{CO_2}} \times \frac{\partial w_{CO_2}}{\partial w_{old}} \times \frac{\partial w_{O_1}}{\partial w_{old}}$$

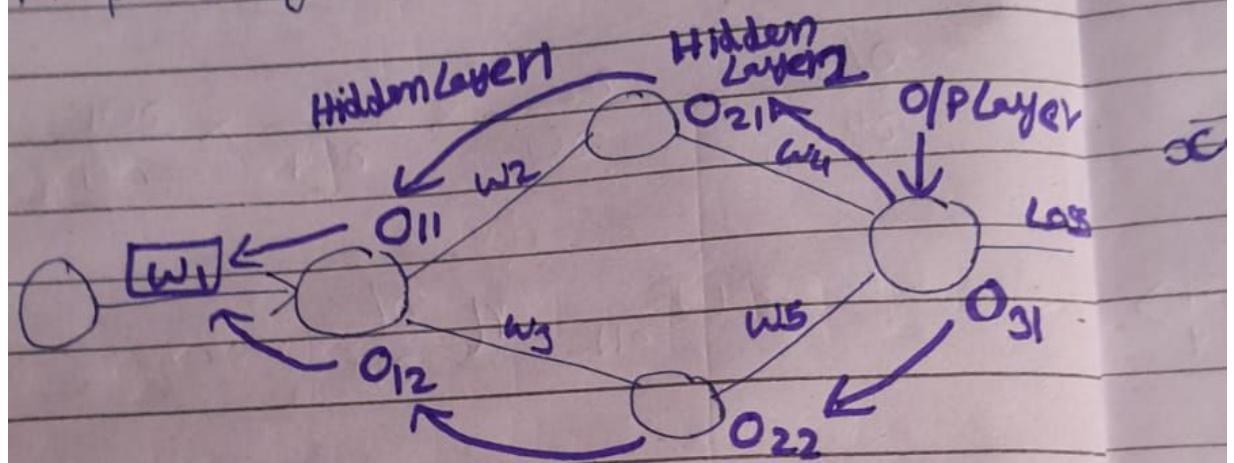
(Loss is dependent on O_2 , O_2 is dependent on O_1 and η is dependent on w_{old})

$$w_{new} = w_{old} - \eta$$

$$\boxed{\frac{\partial L}{\partial w_{O_1}}}$$

$$\frac{\partial L}{\partial w_{O_1}} = \frac{\partial L}{\partial w_{CO_2}} \times \frac{\partial w_{CO_2}}{\partial w_{old}} \times \frac{\partial w_{O_1}}{\partial w_{old}}$$

Now let's take another example
for practicing Chain Rule of Derivation

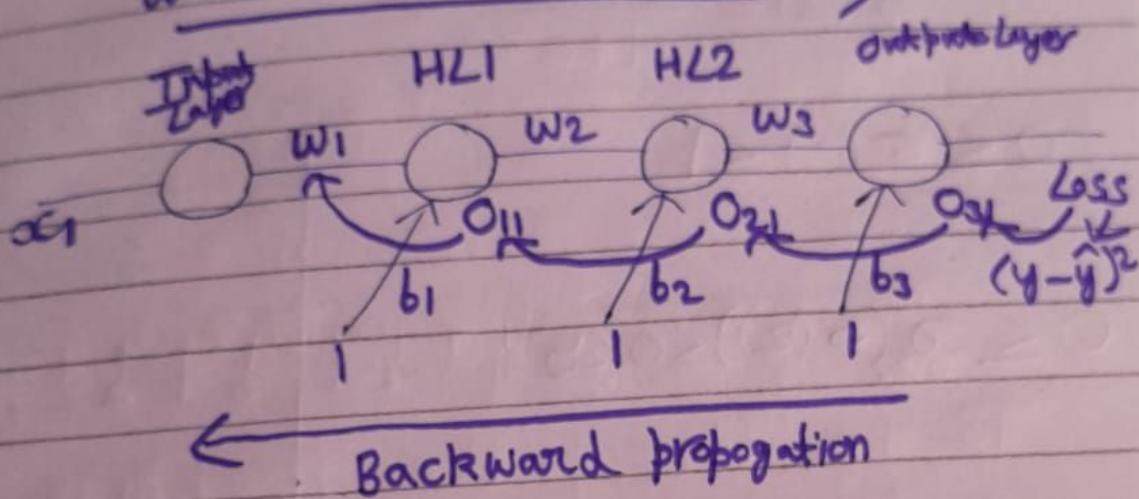


Now let we have to update w_1

$$w_{1\text{new}} = w_{1\text{old}} - \eta \frac{\partial L}{\partial w_{1\text{old}}}$$

$$\begin{aligned} \frac{\partial L}{\partial w_{1\text{old}}} &= \left(\frac{\partial L}{\partial O_{31}} \times \frac{\partial O_{31}}{\partial O_{21}} \times \frac{\partial O_{21}}{\partial O_{11}} \times \frac{\partial O_{11}}{\partial w_{1\text{old}}} \right) \\ &\quad + \left(\frac{\partial L}{\partial O_{31}} \times \frac{\partial O_{31}}{\partial O_{22}} \times \frac{\partial O_{22}}{\partial O_{12}} \times \frac{\partial O_{12}}{\partial w_{1\text{old}}} \right) \end{aligned}$$

⇒ Vanishing Gradient Problem
and Activation Functions:



We know

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

Then activation function is applied Let's use Sigmoid Activation

$$\sigma(z) \Rightarrow 0 \text{ to } 1 \text{ (transform)}$$

↓
Sigmoid

⇒ Activation Functions:

1. Sigmoid Activation Function:

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

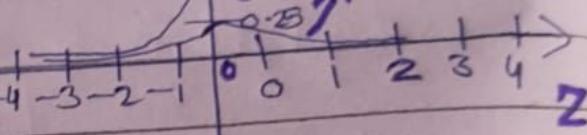
$$0 < \sigma(z) < 1$$

Derivative ($\sigma(z)$)

$$0 \leq \frac{\partial}{\partial z}(\sigma(z)) \leq 0.25$$

$$\sigma(z)$$

Derivative



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

↙ Expand (Chain Rule of Derivation)

$$\frac{\partial L}{\partial w_{i_{\text{old}}}} = \frac{\partial L}{\partial o_{31}} \times \boxed{\frac{\partial o_{31}}{\partial o_{21}}} \times \frac{\partial o_{21}}{\partial o_{11}} \times \frac{\partial o_{11}}{\partial w_{i_{\text{old}}}}$$

Now let's try to solve this

$$\frac{\partial o_{31}}{\partial o_{21}}$$

$$O_{31} = \sigma(w_3 \times O_{21} + b_3)$$

$$\text{Let } Z = w_3 \times O_{21} + b_3$$

$$\text{So, } O_{31} = \sigma(Z)$$

Now

$$\frac{\partial O_{31}}{\partial O_{21}} = \frac{\partial(\sigma(Z))}{\partial Z} \times \frac{\partial Z}{\partial O_{21}} \quad \{\text{Chain Rule}\}$$

We know

$$0 \leq \frac{\partial(\sigma(Z))}{\partial Z} \leq 0.25$$

↓
Derivative of
Sigmoid

So,

$$\frac{\partial O_{31}}{\partial O_{21}} = 0 \leq \sigma(Z) \leq 0.25 \times \frac{\partial(w_3 \times O_{21} + b_3)}{\partial(O_{21})}$$

$$\boxed{\frac{\partial O_{31}}{\partial O_{21}} = 0 \leq \sigma(Z) \leq 0.25 \times w_3}_{\text{old}}$$

- So we are getting smaller values b/w 0 to 0.25 so multiplying it we also get small value in $\frac{\partial L}{\partial w_{1012}}$

the change or update in w_1
will be very small. So weights
are actually never get updated
convergence not happening
this is called vanishing gradient
problem ($w_{\text{new}} \approx w_{\text{old}}$) and this
is happening due to Sigmoid activation.

To fix this problem, researchers
started exploring other activation
functions like

tanh ② ReLU ③ LeReLU ④ SWISS

If neural network is very deep,
then Sigmoid Activation cause
Vanishing Gradient Problem

antages:

Binary Classification Suitable
Clear prediction i.e. 0 or 1

Disadvantages:

- ① Cause Vanishing Gradient Problem
- ② Function output is not zero-centered \Rightarrow weight update is not efficient

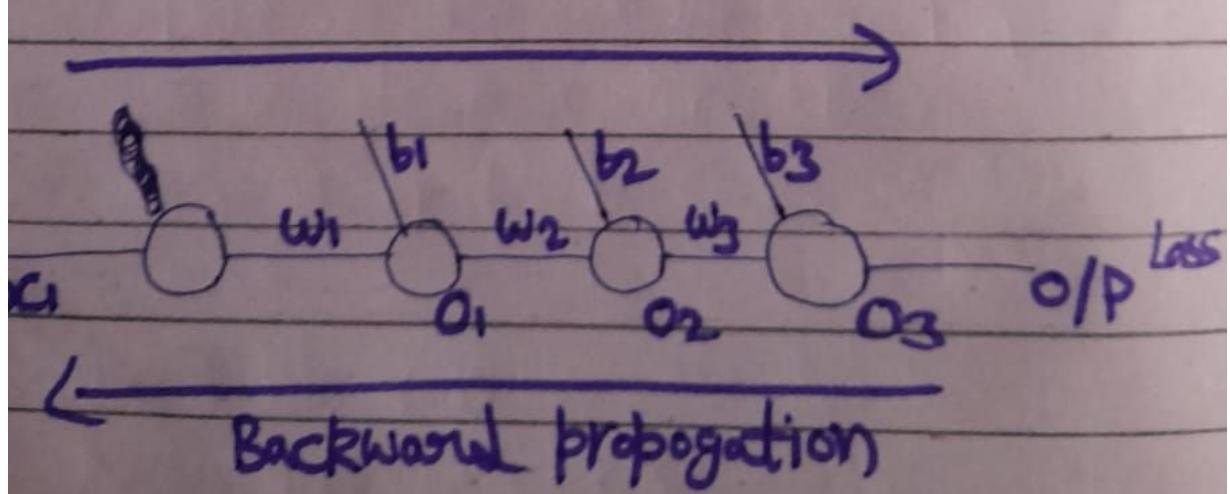
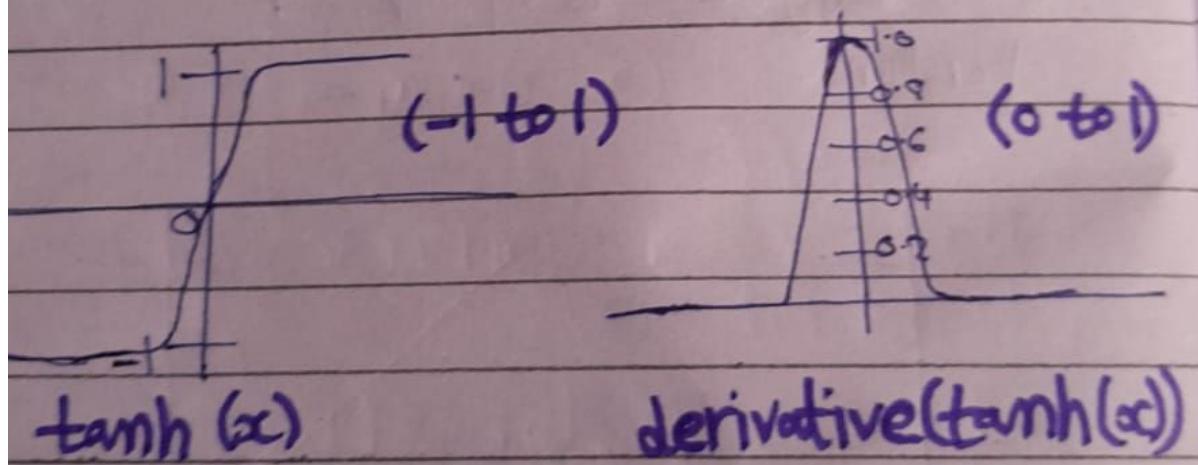
Mathematical operations
are relatively time-consuming.
 $(\text{e.g. } e^P)$ power
operations.

zero-centered
(Mean=0)

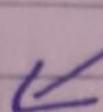
② tanh Activation Function:

$$\tanh x = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

- Its output is between -1 to 1
- The derivative of tanh ranges between 0 to 1



$$w_{2\text{new}} = w_{2\text{old}} - \eta \frac{\partial L}{\partial w_{2\text{old}}}$$



$$\frac{\partial L}{\partial w_{2\text{old}}} = \frac{\partial L}{\partial o_3} \times \frac{\partial o_3}{\partial o_2} \times \frac{\partial o_2}{\partial w_{2\text{old}}}$$

$$\frac{\partial o_3}{\partial o_2} = \frac{\partial (\tanh(z))}{\partial z} \times \frac{\partial z}{\partial o_2}$$

$$z = (o_2 \times w_3) + b_3$$

$$\frac{\partial o_3}{\partial o_2} = (o_2 \neq 1) \left(\frac{\partial ((o_2 \times w_3) + b_3)}{\partial o_2} \right)$$

$$\frac{\partial o_3}{\partial o_2} = (o_2 \neq 1) (w_{3\text{old}}) \Rightarrow \text{still small value in case of deep neural network and cause vanishing gradient problem}$$

vantages:

The output is zero-centric so efficient weight updation.

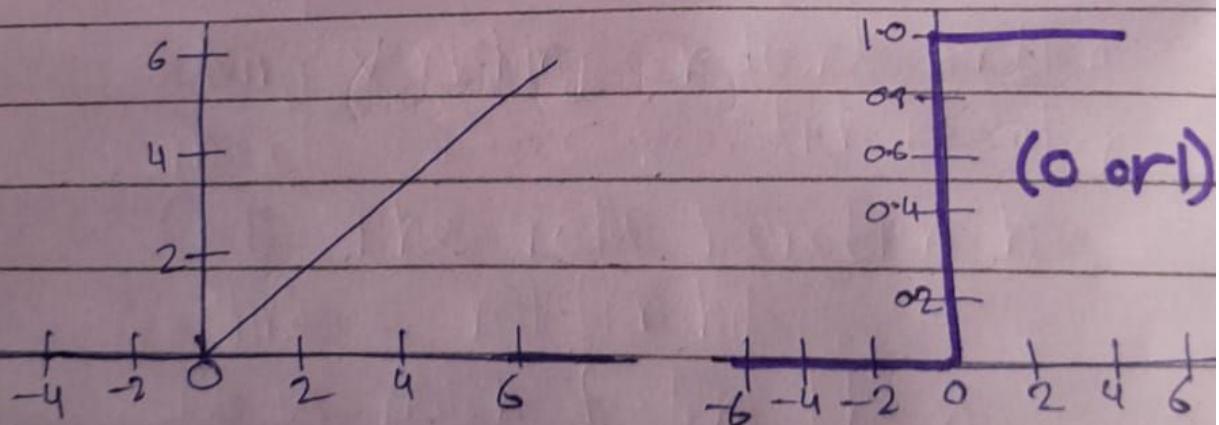
Disadvantages:

- ① Prone to vanishing gradient problem
- ② Still exponential mathematical operations are there which are time consuming.

→ Rectified linear unit

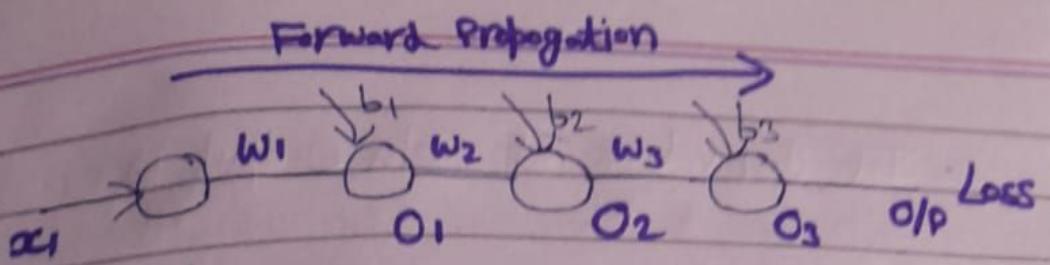
③ ReLU Activation Function:

$$\text{ReLU}(x) = \max(0, x)$$



$$\text{ReLU}(x) = \max(0, x)$$

derivative $\text{ReLU}'(x)$



Backward Propogation

$$w_{2 \text{ new}} = w_{2 \text{ old}} - \eta \left(\frac{\partial L}{\partial w_{2 \text{ old}}} \right)$$

$$\frac{\partial L}{\partial w_{2 \text{ old}}} = \frac{\partial L}{\partial O_3} \times \frac{\partial O_3}{\partial O_2} \times \frac{\partial O_2}{\partial w_{2 \text{ old}}}$$

$$\frac{\partial O_3}{\partial O_2} = \frac{\partial (\text{relu}(z))}{\partial z} \times \frac{\partial z}{\partial O_2}$$

$$\therefore (\text{relu}(z)) \Rightarrow 0 \text{ or } 1$$

$$\therefore z = (w_3 \times O_2) + b_3$$

$$\frac{\partial O_3}{\partial O_2} = [0 \text{ or } 1] \times \frac{\partial (w_3 \times O_2 + b_3)}{\partial O_2}$$

$$\frac{\partial O_3}{\partial O_2} = [0 \text{ or } 1] \times w_3$$

⇒ If Derivative of Relu output is 1
 ⇒ Weight Updation will happen

\Rightarrow If derivative of Relu output is 0
 \Rightarrow Dead Neuron (no weight update)
 $w_{new} \approx w_{old} \Rightarrow$ Dead Neuron

If $z=+ve$ then $\frac{\partial(\text{relu}(z))}{\partial z} = 1$

If $z=-ve$ then $\frac{\partial(\text{relu}(z))}{\partial z} = 0$

Advantages:

- ① Solving Vanishing Gradient Problem
- ② $\max(0, x) \rightarrow$ Calculation is super fast. The Relu Function has a linear relationship.
- ③ It is much faster than Sigmoid or tanh.

Disadvantages:

- ① Dead Neuron
- ② Not zero-centric

④ Leaky Relu and Parametric Relu

Relu:

$$f(x) = \max(\alpha x, x) \quad \xrightarrow{\text{Hyperparameter}} \alpha = 0.01, 0.02, 0.03 \text{ etc.}$$

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

- So the derivative of relu leaky will never be zero something greater than 0 even in case of -ve Z so no Dead Neurons.

Advantages:

1. Leaky Relu has all the advantages of Relu
2. It removes the Dead Relu problem

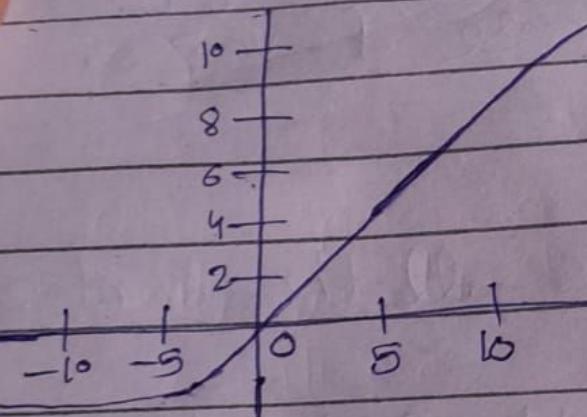
Disadvantage:

1. It is not zero-centric

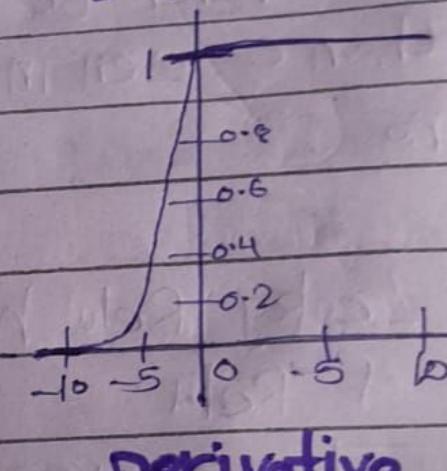
⑤ ELU (Exponential Linear Unit):

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

Forward Propogation



Backward propogation



Derivative
f'(x)

- Advantages:

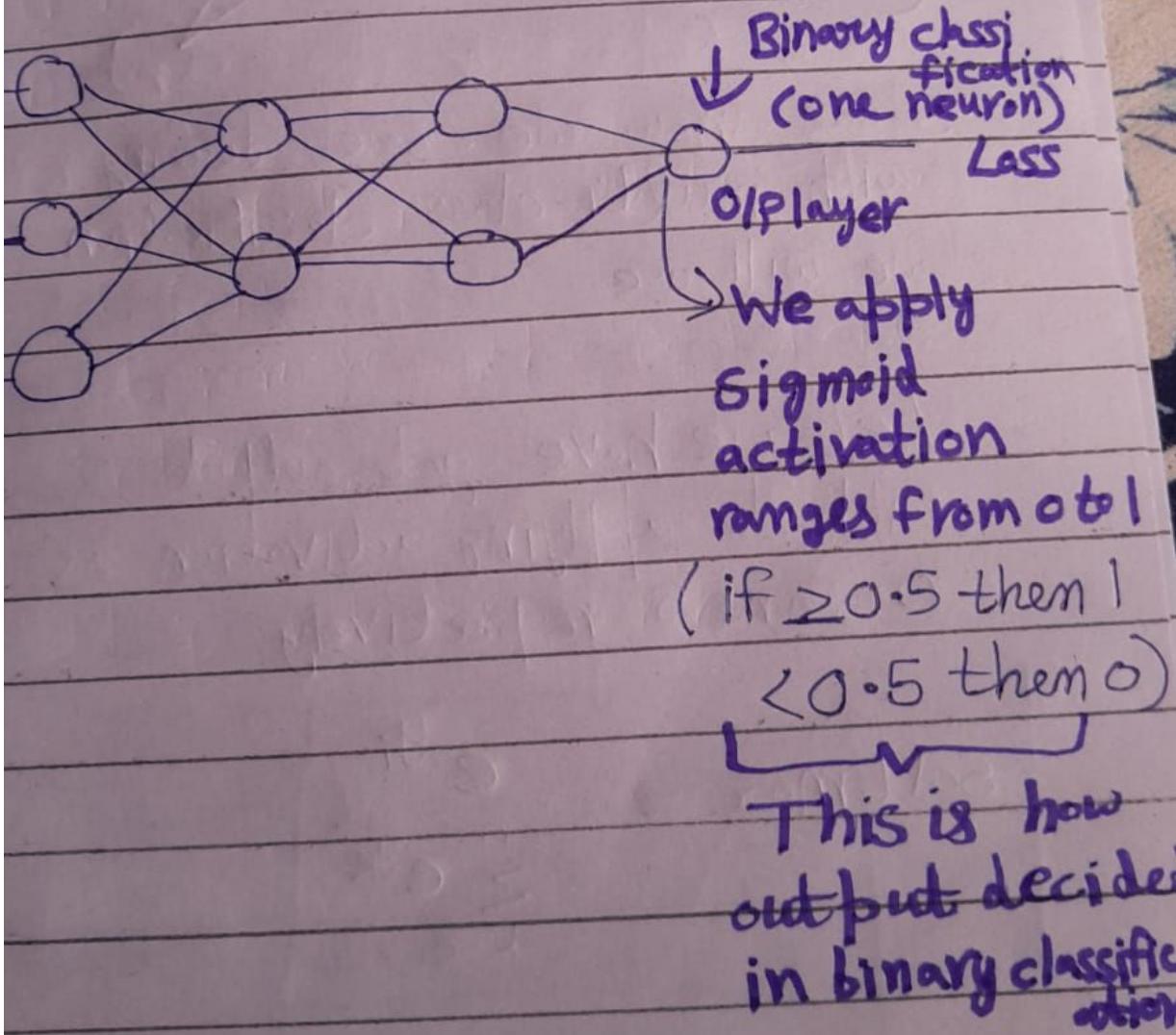
- ① No Dead Relu issues
- ② It is zero centered

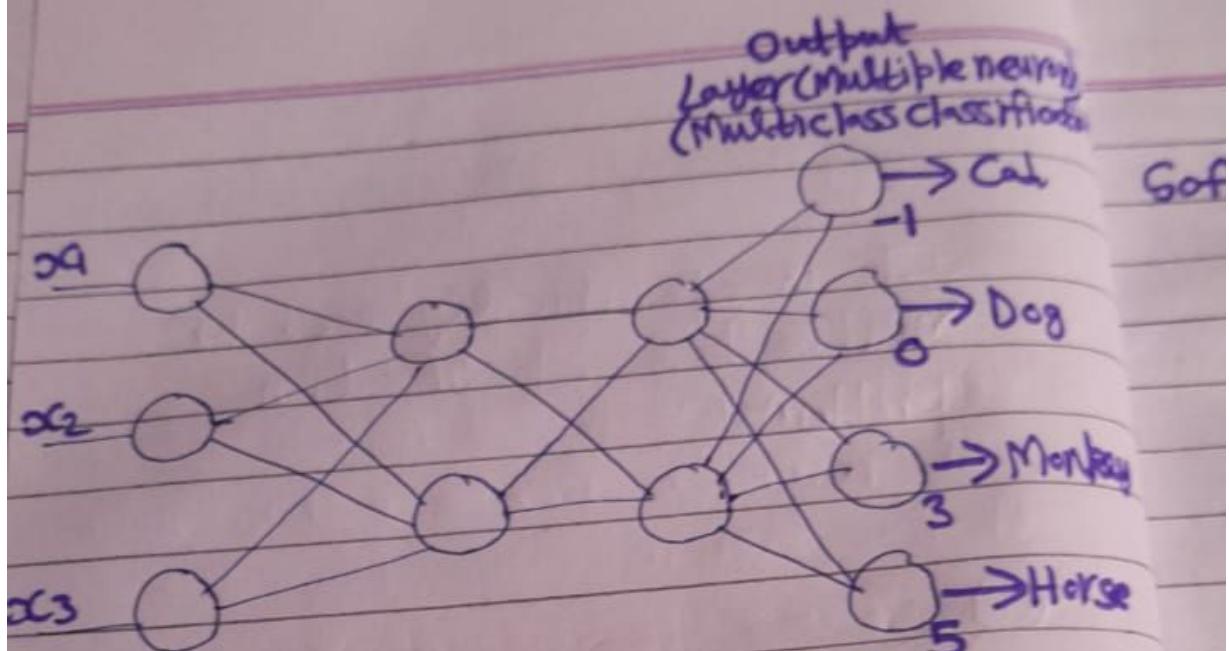
- Disadvantages:

- ① Slightly more computationally intensive because of exponential operation

Softmax Activation Function:

- It is specifically used for multi-class classification problem





- So as now we are dealing with multi-class classification we will use **Softmax activation function**.

- Let we have got outputs without applying activation as -1, 0, 3 and 5 respectively

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

$$f_{\text{max}} \Rightarrow \text{act} = \frac{e^{-1}}{(e^{-1+0+3+5})} = 0.00033$$

$\therefore y$ is output without apply activation

$$\text{Dog} \Rightarrow \frac{e^0}{e^{-1+0+3+5}} = 0.0024$$

$$\text{Monkey} \Rightarrow \frac{e^3}{e^{-1+0+3+5}} = 0.0183$$

$$\text{Horse} \Rightarrow \frac{e^5}{e^{-1+0+3+5}} = 0.1353$$

$$\Pr(\text{Horse}) = \frac{0.1353}{0.00033 + 0.0024 + 0.0183}$$

$$\boxed{\Pr(\text{Horse}) = 86\%}$$

⑦ Which Activation Function
to use when?

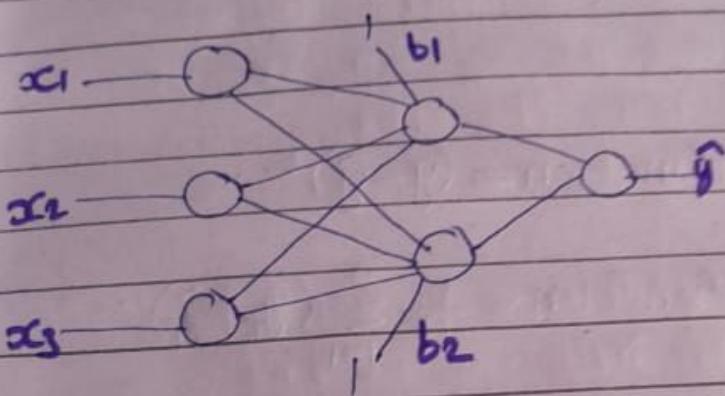
In Output Layer:

Binary Classification \rightarrow Sigmoid
Multi-class Classification \rightarrow Softmax

In Hidden Layers:

Use ReLU or its variants
like Pre-ReLU, Leaky-ReLU or ELU

Loss Function and Cost Function:



$x_1 \quad x_2 \quad x_3 \quad \text{O/P}$

-	-	-	0
-	-	-	1
-	-	-	0
-	-	-	1

Loss Function

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

Cost Function

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

- In loss function we pass one datapoint from dataset at a time.

- In cost function we pass all datapoints of dataset all at once.

\Rightarrow Regression Loss and Cost Functions:

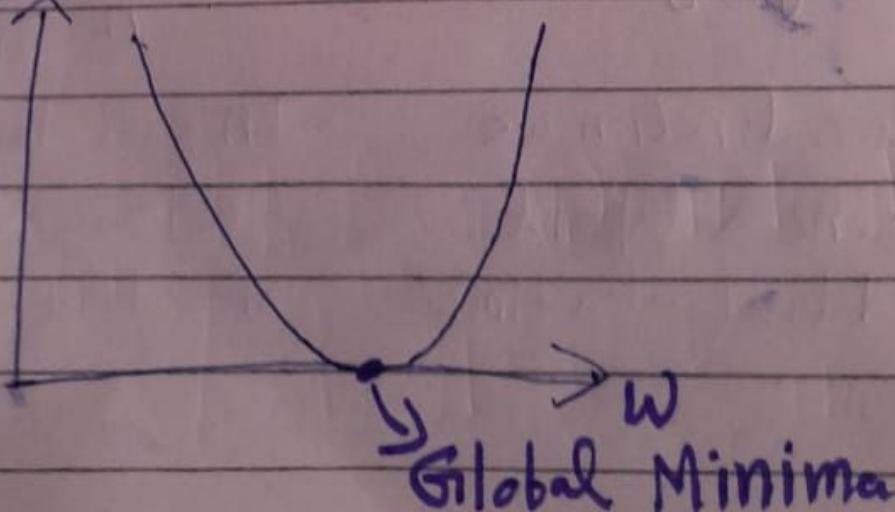
① Mean Squared Error (MSE):

$$\text{Loss Function} = (y - \hat{y})^2$$

$$\text{Cost Function} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

As the equation of Loss and cost is quadratic so it's graph with weight gives us gradient descent (parabola) which has one global minima and we have to reach that.

Loss



⇒ Advantages:

- ① MSE is differentiable.
- ② It has one local or global minima so we don't get stuck.
- ③ It converges faster.

⇒ Disadvantages:

- ① MSE is not robust to outlier as we square up the error and we are penalizing the cost function so line is changing much just due to one outlier.

Mean Absolute Error (MAE):

Loss Function = $|y - \hat{y}|$

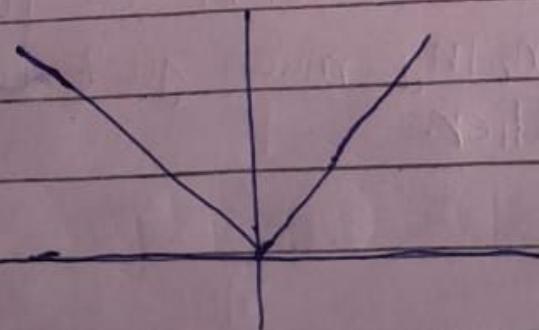
Cost Function = $\frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$

\Rightarrow Advantages:

- ① Robust to outliers because we are not squaring or penalizing the error.

\Rightarrow Disadvantages:

- ① Convergence usually takes time in MAE, as we don't get parabola curve and we have to use sub gradient concept.



③ Huber Loss:

$$\text{Cost Function} = \begin{cases} \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2, & \text{if } |y - \hat{y}| \leq \delta \\ \delta |y - \hat{y}| - \frac{1}{2} \delta^2, & \text{otherwise} \end{cases}$$

δ is a specific threshold treated as Hyperparameter.

- Huber Loss is actually combination of both MSE and MAE if there is no outlier than it use MSE equation otherwise adjusted MAE.

④ Root Mean Squared Error (RMSE):

$$\text{Cost Function} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

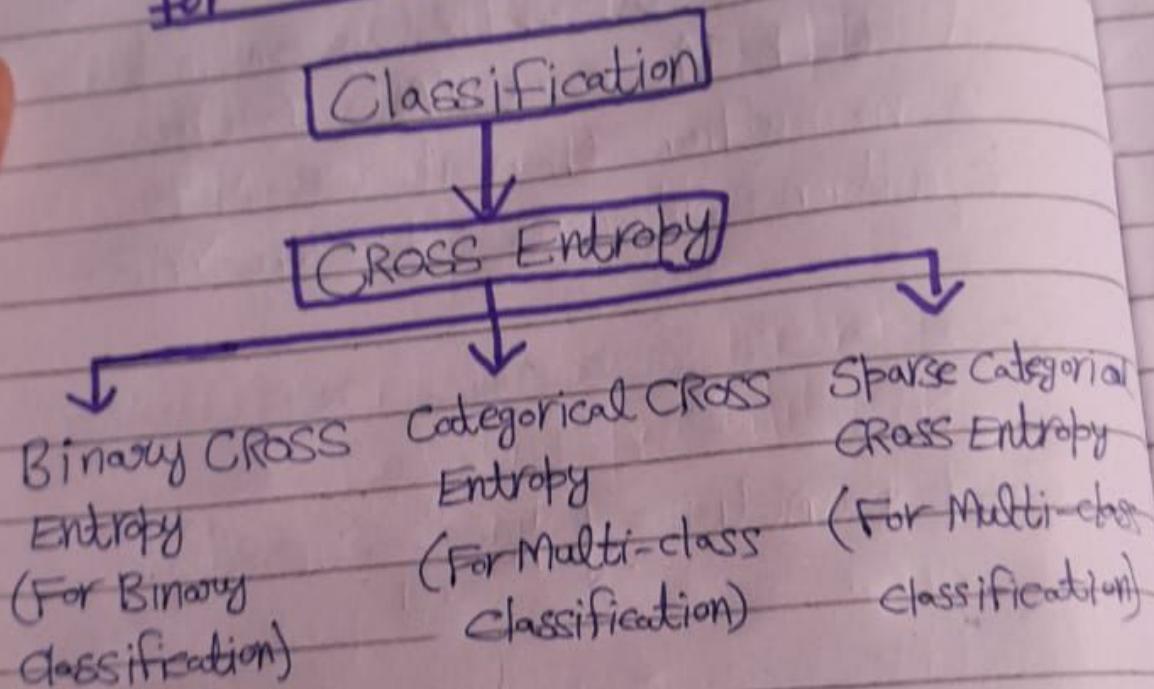
Advantages:

- Same unit
- Differentiable

Disadvantages:

- Not robust to outliers

⇒ Loss and Cost Functions
for Classification Problems.



① Binary CROSS Entropy:

↙ Log Loss

$$\text{Loss} = -y * \log(\hat{y}) - (1-y) * \log(1-\hat{y})$$

where

$y \Rightarrow$ Actual value, $\hat{y} \Rightarrow$ Predicted value

$$\text{Loss} = \begin{cases} -\log(1-\hat{y}) & \text{if } y=0 \\ -\log(\hat{y}) & \text{if } y=1 \end{cases}$$

- It is used for Binary Classification

② Categorical Cross Entropy.

- It is used for Multi-class classification problem

f_1	f_2	f_3	Output
2	3	4	Good
5	6	7	Bad
8	9	10	Neutral

↓ Step 1: Perform One Hot Encoding of output variable

$j=1$ $j=2$ $j=3$
Good Bad Neutral

1	0	0
0	1	0
0	0	1

$$L(x_i, y_i) = - \sum_{j=1}^C y_{ij} * \ln(g_{ij})$$

where $i = 1$ to n (For multiple records)

$$y_{ij} = \begin{cases} 1 & \text{if the element is in class} \\ 0 & \text{otherwise} \end{cases}$$

All \hat{y}_{ij} is calculated by getting output from neuron and apply softmax function.



gives us probabilities

- So Categorical CROSS Entropy also gives the probability of other categories

③ Sparse Categorical CROSS Entropy:

- Very similar to Categorical CROSS Entropy

Let we put some index with respect to categories and

mention their probabilities

$[0, 0.2, 0.3, 0.5]$ → categories



2nd index → output

Disadvantages:

- It don't give probability of other categories

→ Which Loss Function to use when?

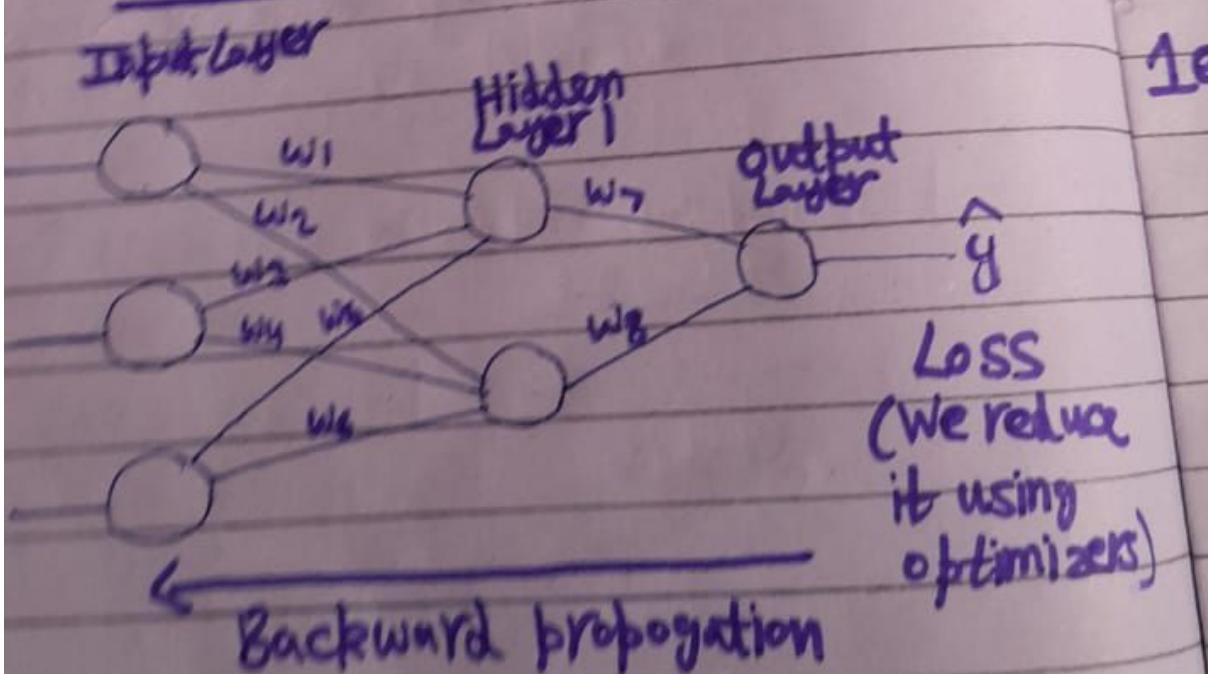
* Right combination

	Hidden Layers	O/P Layer	Problem	Loss Function
①	Relu or variants	Sigmoid	Binary Classification	BINARY CROSS ENTROPY
②	Relu or variants	Softmax	Multi-class	CATEGORICAL OR SPARSE
③	Relu or variants	Linear	Regression	MSE, MAE, Huber loss, RMS

Optimizers:

- ① Gradient Descent
- ② Stochastic Gradient Descent (SGD)
- ③ Mini Batch SGD
- ④ SGD with Momentum
- ⑤ Adagrad and RMSProp
- ⑥ Adam optimizers

① Gradient Descent Optimizer:



We know weight updation formula generally

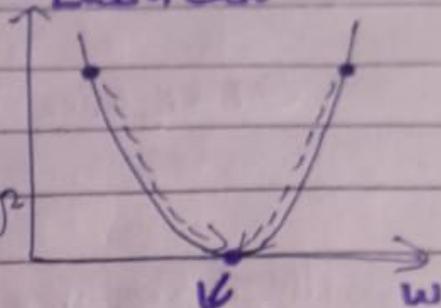
$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w_{\text{old}}} \quad \text{Learning rate}$$

MSE

$$\text{Loss Function} = (y - \hat{y})^2$$

$$\text{Cost Function} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Loss or cost



Epochs, Iteration

1 epoch {

1000 datapoint

$\hat{y}_i \Rightarrow$ Cost Function Our aim is to

update weights in backward propagation to reach global minima

weights will get updated by weight updation formula

(At global minima

slope = 0 so

$w_{\text{new}} \approx w_{\text{old}}$ no

weight updation.

- Mean in gradient descent optimizer we pass all datapoints at once not separately and use cost function

This is the point where loss is minimum)

- We can perform any number of epochs until we reach global minima

- In Gradient Descent optimizer,
1 epoch = 1 iteration
(as we pass all datapoints at once)

- Advantages:

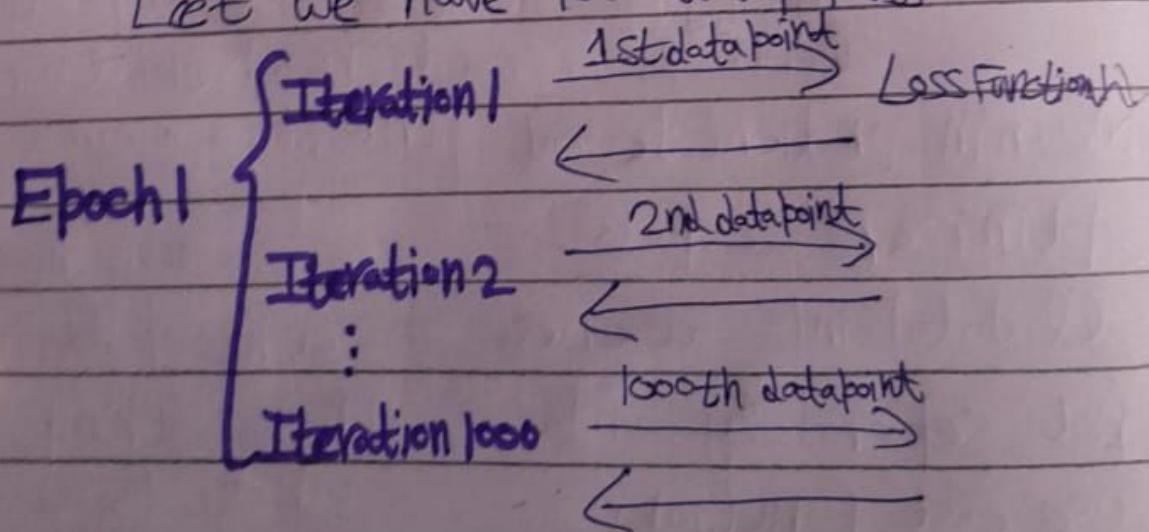
① Convergence will happen

- Disadvantage:

① Huge Resource Required i.e RAM & GPU
as we are passing all datapoints
So, Resource Intensive

② Stochastic Gradient Descent (SGD):

Let we have 1000 datapoints



in which we just one datapoint at a time and that is called learning with the first pass as many number of iteration as many datapoints we have

disadvantage:

the convergence issue

Advantage:

Time Complexity is huge

Convergence will also take more time

Noise gets introduced (as we are taking single datapoint)



③ Mini Batch SGD:

⇒ In Mini-Batch SGD Batch size is used means like SGD instead of passing one datapoint in each iteration we pass datapoints equals batch size

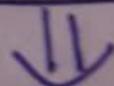
Let

$$\text{Datapoints} = 100000$$

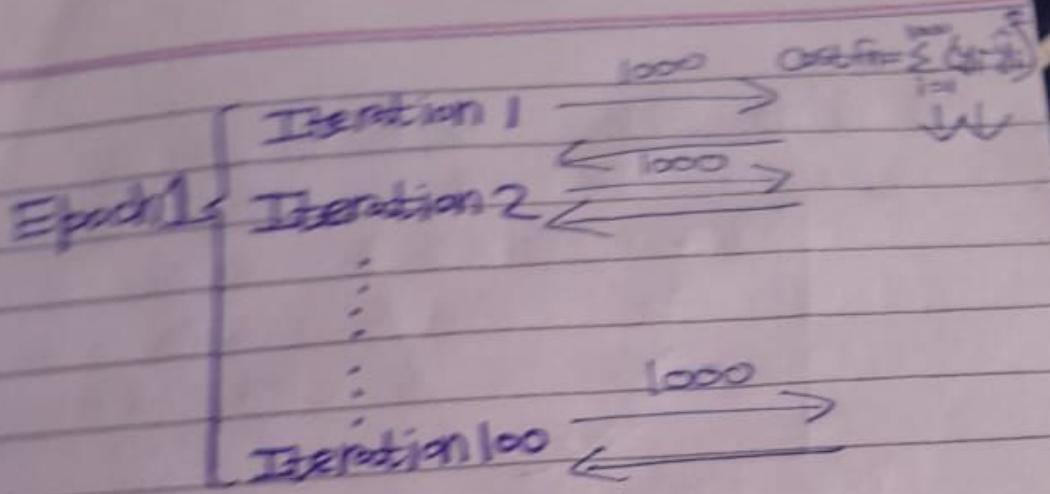
$$\text{batchsize} = 1000$$

$$\text{No. of iterations} = \frac{\text{Datapoints}}{\text{batchsize}} = \frac{100000}{1000}$$

$$\boxed{\text{No. of iterations} = 100}$$



Means there will be 100 iteration per epoch and we pass 1000 datapoints in one iteration.



⇒ In Mini-Batch SGD the noise gets reduce but still exist

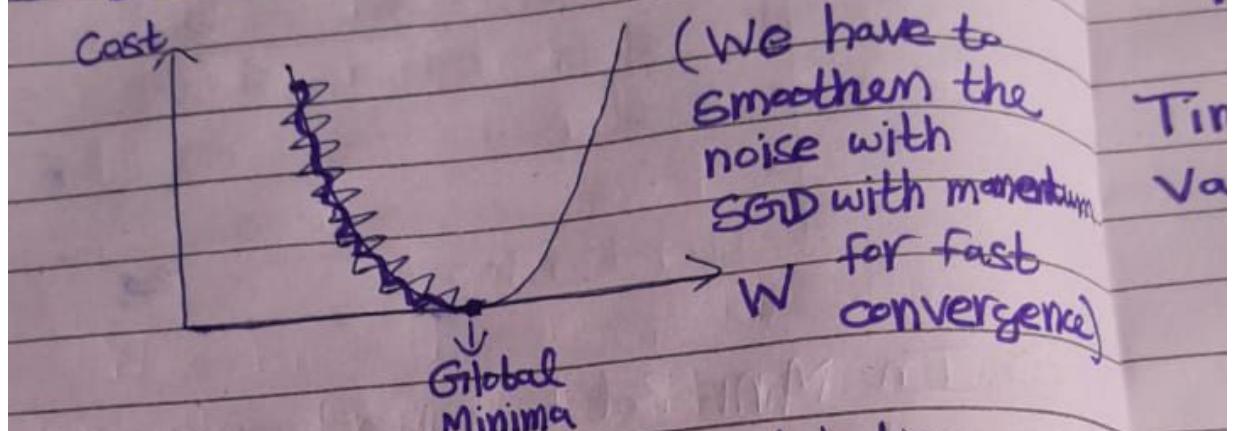
Advantages:

- ① Convergence speed will increase. ↑ compared to SGD
- ② Noise will be less when compared to SGD.
- ③ Efficient Resource Usage

Disadvantage:

- ① Noise still exist

④ SGD with Momentum:



We know the weight updation formula and bias updation formula

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

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

Now the updated weight updation formula for SGD with momentum is:

$$w_t = w_{t-1} - \eta \frac{\partial L}{\partial w_{t-1}}$$

↓
Weight Updation
Formula

Exponential Weighted
Average { Smoothening } \Rightarrow ARIMA, SARIMAX

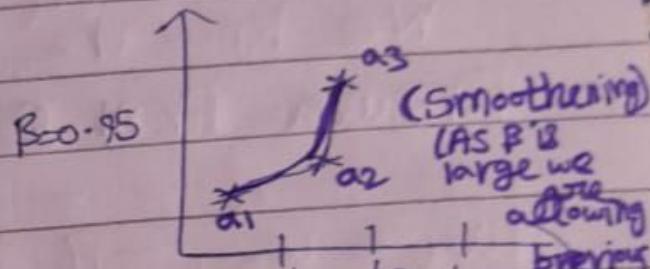
Time = $t_1 \ t_2 \ t_3 \ t_4 \ \dots \ t_n$
 Values = $a_1 \ a_2 \ a_3 \ a_4 \ \dots \ a_n$

Time Series

Now

$v_{t_1} = a_1$
 \downarrow
 Value at time t_1

$v_{t_2} = \beta * v_{t_1} + (1-\beta) * a_2$
 \downarrow
 $\beta = 0.95$ (parameter)



Let $\beta = 0.95$

$$v_{t_2} = (0.95) * a_1 + (0.05) a_2$$

$$v_{t_3} = \beta * v_{t_2} + (1-\beta) * a_3$$

$$v_{t_3} = 0.95 [0.95 * a_1 + 0.05 * a_2] + (0.05) * a_3$$

Advantage:

- ① Reduces the noise
- ② Quick convergence

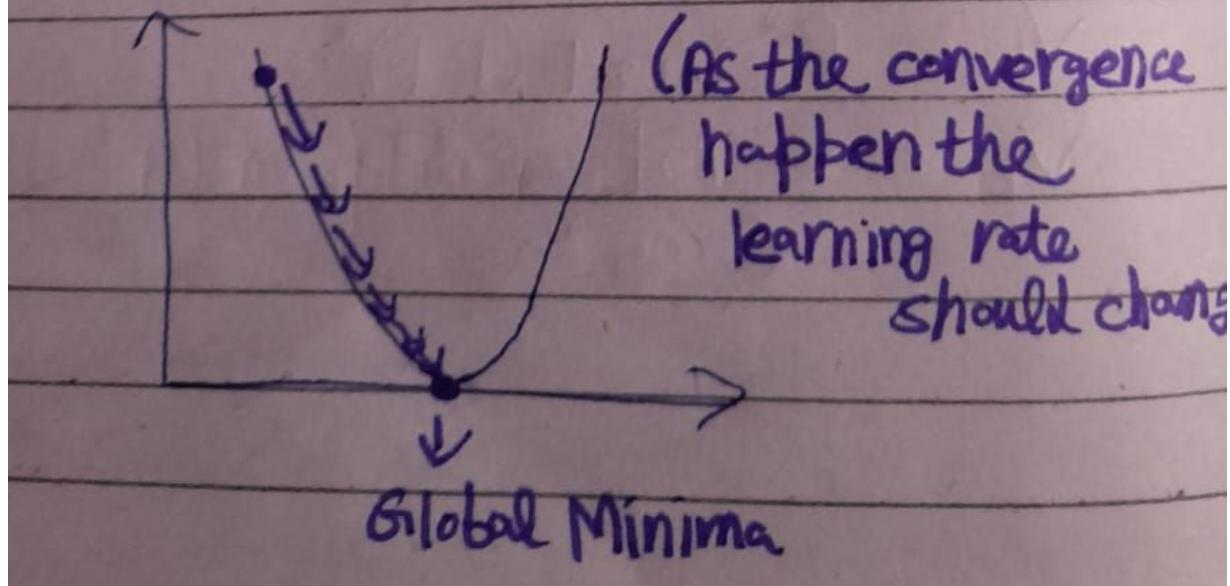
⑤ Adagrad (Adaptive Gradient Descent):

We know the new weight update formula from SGD with momentum

$$w_t = w_{t-1} - \eta \frac{\partial L}{\partial w_{t-1}}$$

→ Learning rate

- In previous optimizer η (learning rate) was fixed. Now we will make this η (learning rate) dynamic.



Now our new weight updation formula
for Adagrad is:

$$w_t = w_{t-1} - \eta' \frac{\partial L}{\partial w_{t-1}}$$

$$\eta' = \frac{\eta}{1 + \alpha_t + \epsilon}$$

Epsilon (small value)
(so that denominator
don't become 0
even when α_t
is 0)

$$\alpha_t = \sum_{i=1}^t \left(\frac{\partial L}{\partial w_i} \right)^2$$

As α_t keeps on increasing the learning rate keeps on decreasing means learning rate (η) will be large initially then reduce as we get closer to convergence

Disadvantage:

- ① $\eta' \rightarrow$ possibility to become a very small value ≈ 0 (in very deep neural network having large α_t)

⑥ Adadelta and RMSProp:

We know the weight update formula from Adagrad

$$w_t = w_{t-1} - \eta' \frac{\partial L}{\partial w_{t-1}}$$

But now,

$$\eta' = \frac{\eta}{\sqrt{S_{dw_t} + \epsilon}}$$

Exponential weighted average (smoothing)

$$S_{dw_t} = \beta * S_{dw_{t-1}} + (1-\beta) \left(\frac{\partial L}{\partial w_{t-1}} \right)^2$$

We initialize $S_{dw_t} = 0$

⑦ Adam optimizer:

⇒ Adam optimizer is the best optimizer

⇒ It combines:

SGD with momentum + RMSProp

$$w_t = w_{t-1} - \eta' Vd_{wt} \Rightarrow \text{Weight Update}$$

$$b_t = b_{t-1} - \eta' Vd_{bt} \Rightarrow \text{Bias Update}$$

We know

$$\eta' = \frac{\eta}{(Sd_{wt} + \epsilon)}$$

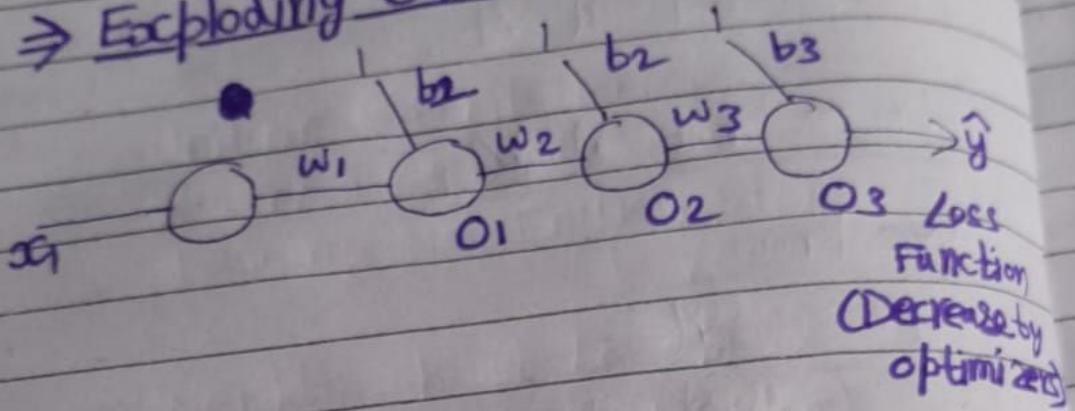
$$Sd_{wt} = \beta * Sd_{wt-1} + (1-\beta) \left(\frac{\partial L}{\partial w_{t-1}} \right)^2$$

Now

$$Vd_{wt} = \beta * Vd_{wt-1} + (1-\beta) \frac{\partial L}{\partial w_{t-1}}$$

$$Vd_{bt} = \beta * Vd_{bt-1} + (1-\beta) \frac{\partial L}{\partial b_{t-1}}$$

\Rightarrow Exploding Gradient Problem:



We know

$$w_{1\text{ new}} = w_{1\text{ old}} - \eta$$

$$\boxed{\frac{\partial L}{\partial w_{1\text{ old}}}}$$

$$\frac{\partial L}{\partial w_{1\text{ old}}} = \frac{\partial L}{\partial o_3} \times \boxed{\frac{\partial o_3}{\partial o_2}} \times \frac{\partial o_2}{\partial o_1} \times \frac{\partial o_1}{\partial w_1}$$

$$\frac{\partial o_3}{\partial o_2} = \frac{\partial (\sigma(z))}{\partial z} \times \frac{\partial z}{\partial o_2}$$

We know derivative of Sigmoid is from 0 to 0.25 and $z = o_2 \times w_3 + b_3$

$$\frac{\partial o_3}{\partial o_2} = [0 \text{ to } 0.25] * \frac{\partial (o_2 \times w_3 + b_3)}{\partial o_2}$$

$$\boxed{\frac{\partial o_3}{\partial o_2} = [0 \text{ to } 0.25] * w_3}$$

Now if the weight initialization is
with very large number like
 $w_{3012} = 500 \text{ to } 1000$

So when we multiply big value we
also get very big value and after
multiplying with learning rate we get
very big value now

$$w_{1, \text{new}} \gg w_{1, \text{old}} \quad \text{or} \quad w_{1, \text{new}} \ll w_{1, \text{old}}$$

($w_{1, \text{new}}$ is very bigger
than $w_{1, \text{old}}$)

($w_{1, \text{new}}$ is very
smaller than
 $w_{1, \text{old}}$)

⇒ As weight update is with very
large value so convergence become
difficult and this is called
exploding gradient problem which
happen due to wrong weight initiali

(ii) Xavier Uniform Initialization:

$w_{ij} \sim \text{Uniform Distribution } \left(-\frac{\sqrt{k}}{\sqrt{n_{\text{input}} + n_{\text{output}}}}, \frac{\sqrt{k}}{\sqrt{n_{\text{input}} + n_{\text{output}}}} \right)$

$w_{ij} \sim \text{Uniform Distribution } \left(-\frac{1}{\sqrt{n_{\text{input}}}}, \frac{1}{\sqrt{n_{\text{input}}}} \right)$

③ Kaiming He Initialization:

(i) He Normal

$w_{ij} \sim \text{Normal Distribution } (0, \sigma)$

$$\sigma = \sqrt{\frac{2}{n_{\text{input}}}}$$

In our case

$$\sigma = \sqrt{\frac{2}{3}}$$

$w_{ij} \sim \text{Normal Distribution } (0, \sqrt{\frac{2}{3}})$

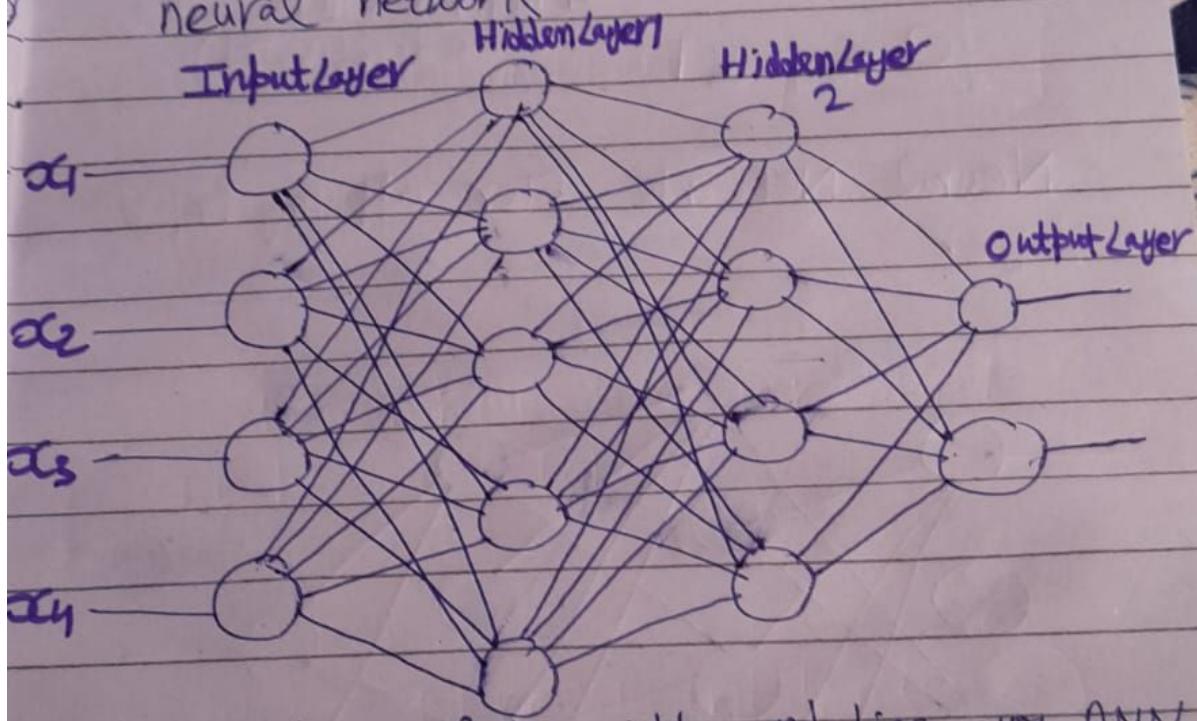
(ii) He Uniform:

$w_{ij} \sim \text{Uniform Distribution } \left[-\frac{\sqrt{6}}{\sqrt{n_{\text{input}}}}, \frac{\sqrt{6}}{\sqrt{n_{\text{input}}}} \right]$

$w_{ij} \sim \text{Uniform Distribution } \left[-\frac{\sqrt{6}}{\sqrt{3}}, \frac{\sqrt{6}}{\sqrt{3}} \right]$

⇒ Dropout Layer:

Let we have a multi-layered neural network



When we have multiple weights and bias our ANN can get overfitted

Overfitting:

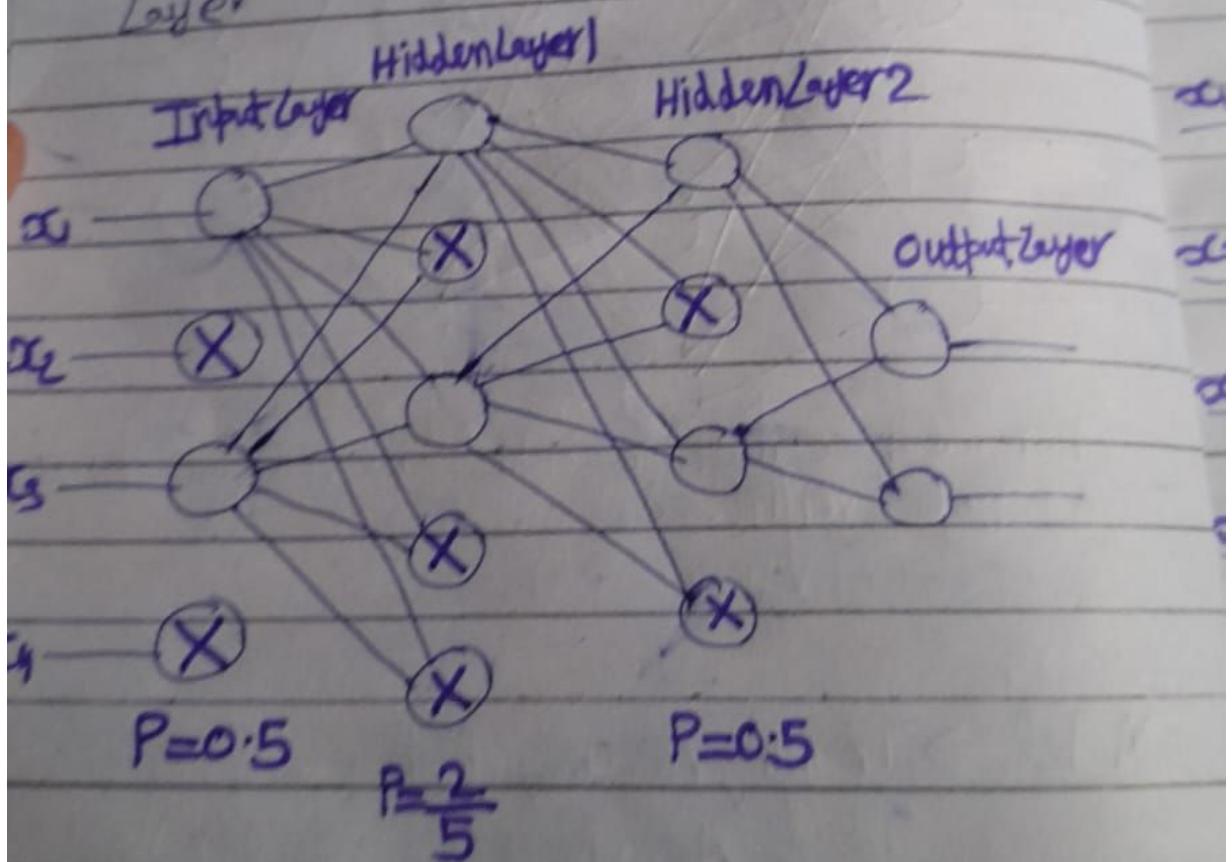
Model → Training Accuracy (high e.g 90%)

Test Accuracy (low e.g 60%)

⇒ But we have to get generalized model.

- To reduce overfitting in Neural network we use **Dropout Layer**.
Let $p=0.5$ (Means 50% of nodes from input will pass others will be deactivated.)
- OSPSI**

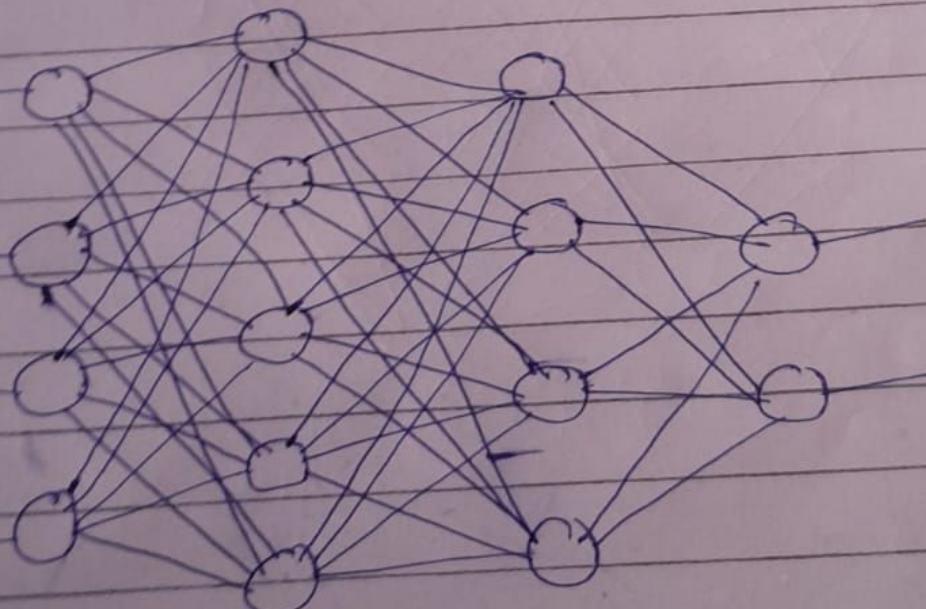
Neural Network after applying Dropout Layer



- In 1st epoch depending on P some neurons gets activated and others will be deactivated and they are

selected randomly and will not be same.

- The purpose is that we are performing dropout (deactivating some neurons) to avoid overfitting.
- P is the Hyperparameter $0 \leq P \leq 1$
- Now for test data



$$P=0.5 \quad P=\frac{2}{5}$$

- For test data we will not apply hidden layer all neurons will be activated:
Here we multiply the weights with the probability of hidden layer.