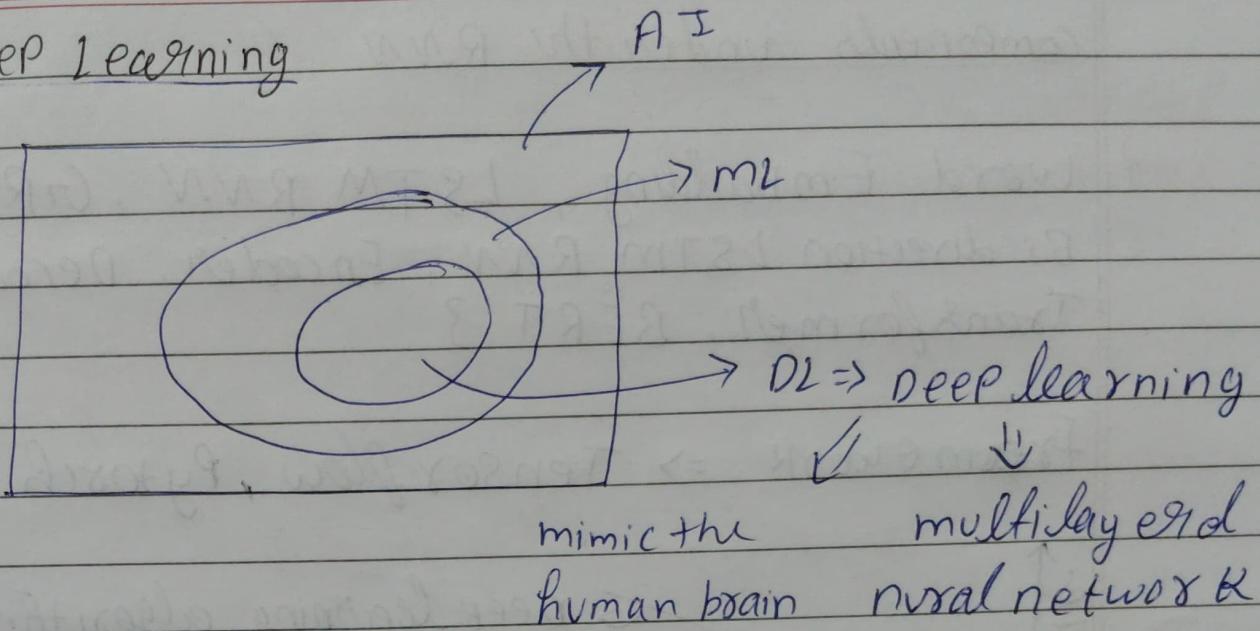


Deep Learning

(1) ANN \Rightarrow Artificial Neural Network \Rightarrow classification
 \Rightarrow Regression.

(2) CNN \Rightarrow convolutional Neural Network \Rightarrow

\downarrow
 I/P : Image / video frames \rightarrow RCNN,
 masked RCNN, YOLO

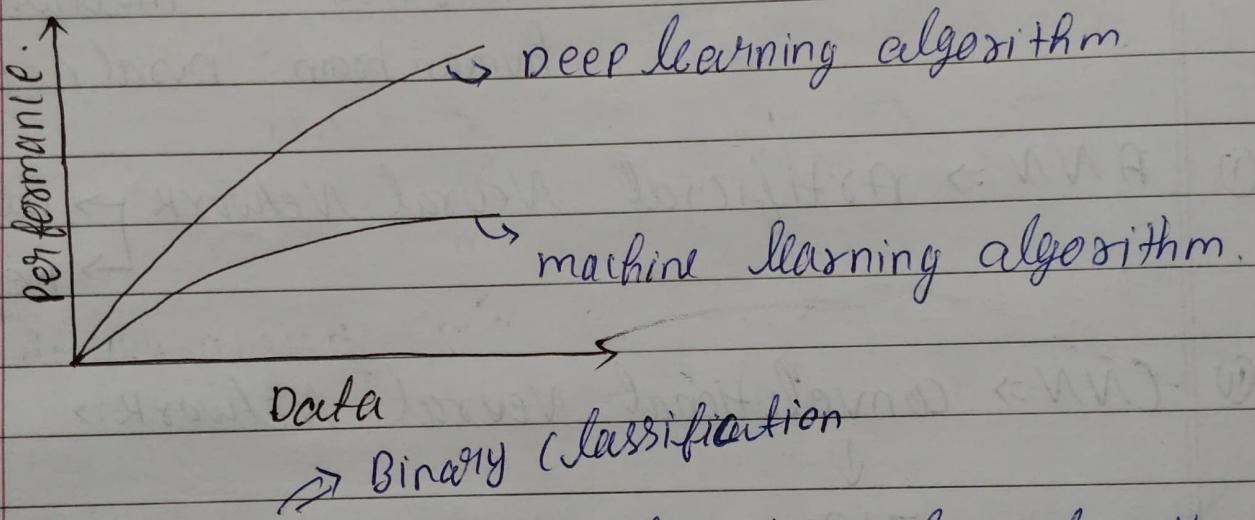
\downarrow
 Computer vision and object detection.

(3) RNN \Rightarrow Recurrent neural Network \Rightarrow NLP, Time-Series
 Text, Time series \leftarrow I/P

Components under the RNN

{ Word Embedding, LSTM RNN, GRU RNN
Bi direction LSTM RNN, Encoder, Decoder,
Transformer, BERT }

Framework \Rightarrow Tensorflow, Pytorch



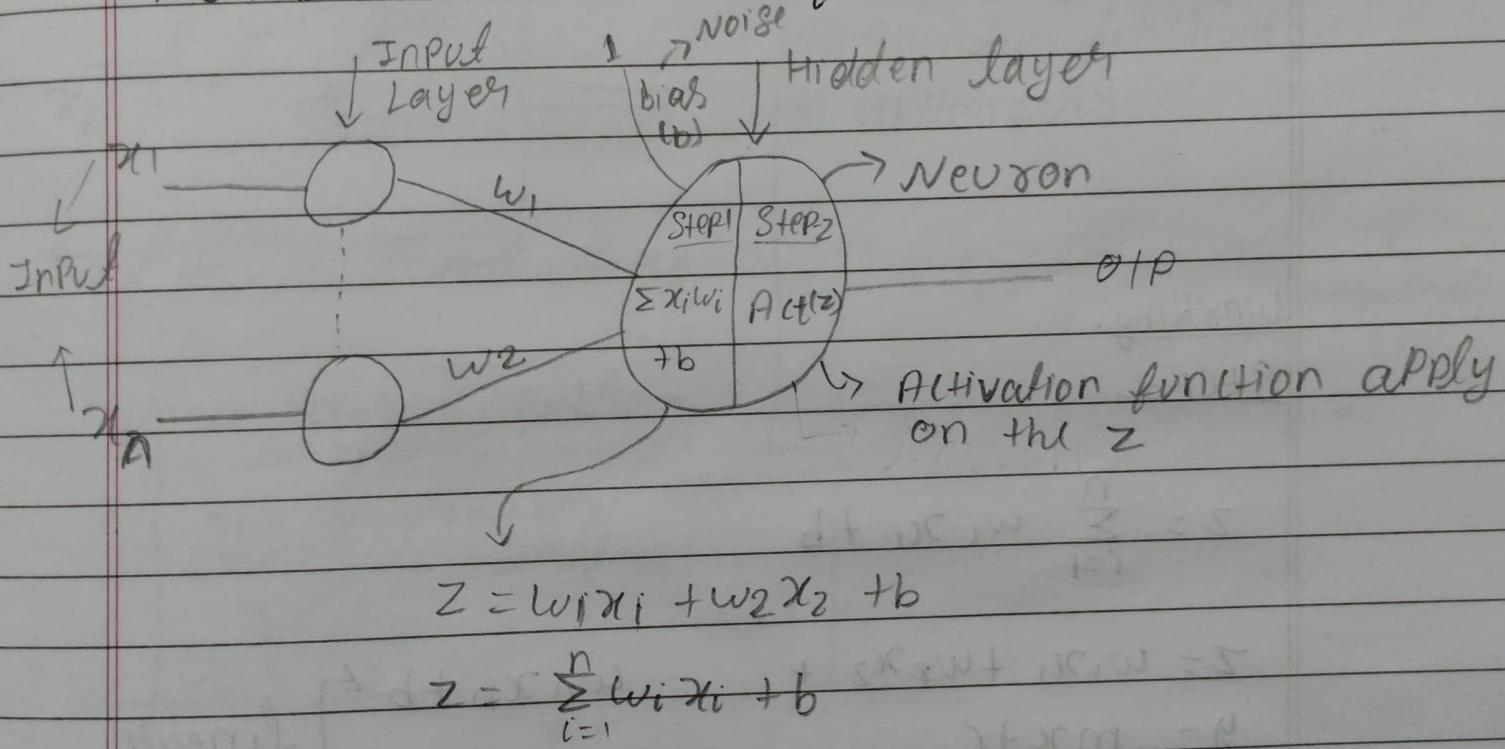
Perceptron: \Rightarrow A Single Neural Network or a neural network unit.

it consist 4 things:

- ① Input layer
- ② Hidden layer
- ③ weights
- ④ Activation function.

Dataset:

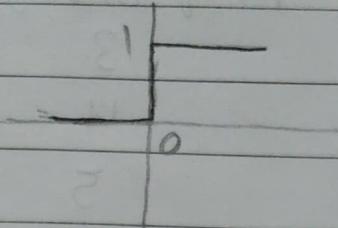
$IQ \rightarrow x_1$	no of Study hrs $\rightarrow x_2$	O/P: Pass/fail
95	3	$0 \rightarrow \text{fail}$
110	4	$1 \rightarrow \text{Pass}$
100	5	

A simple architecture of Perceptron.

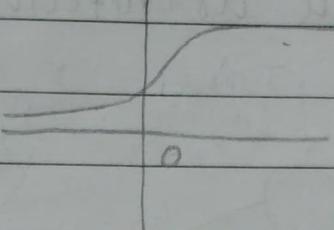
Activation function: Activation function is basically transform the output b/w some value like 0 to 1, -1 to 1

Activation function.

(i) Step function:



(ii) Sigmoid function:

Working.Step-1

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

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

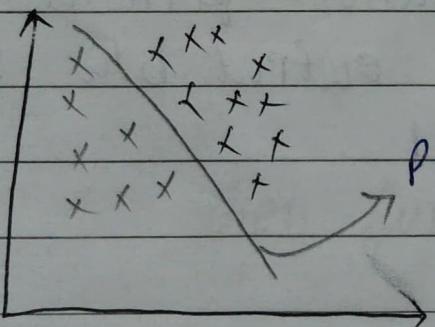
$$y = mx + c$$

$$y = \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \beta_0$$

linear

problem

statement.

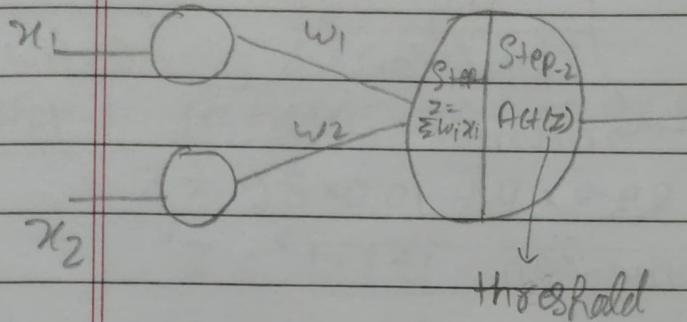


perception is a linear classifier.

[ANN]

Perceptron model

↓
Single layered
Perceptron



↓
multi layered
perceptron

① forward propagation

② backward propagation

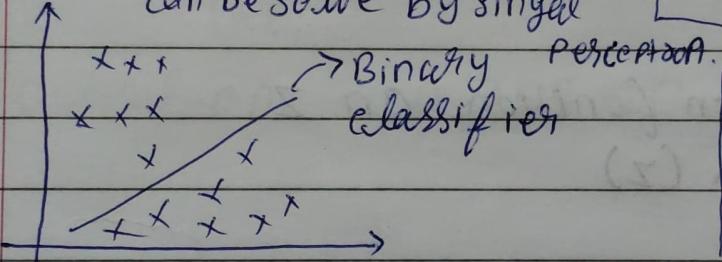
③ loss functions

④ activation functions.

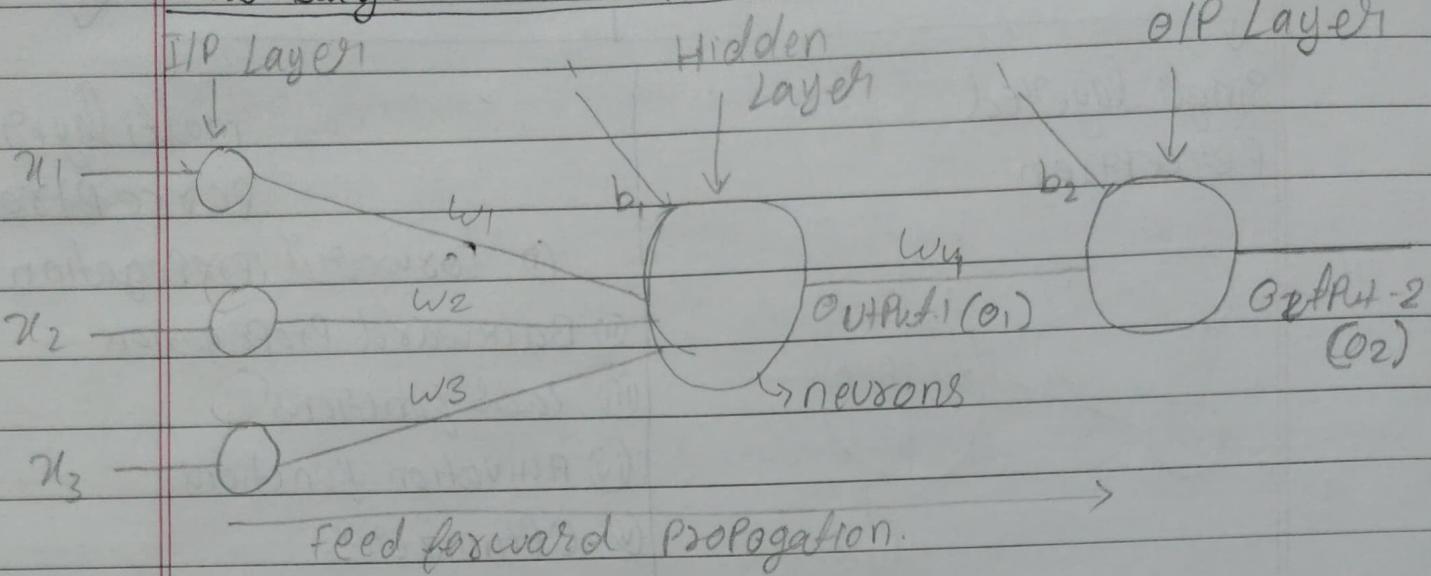
⑤ optimizers.

Feed forward Neural Network

Linear Separable problem.
can be solve by singal



Multilayered Perceptron model [ANN]



Step-1: $Z = w_1x_1 + w_2x_2 + w_3x_3 + b_1$

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

Step-2: Activation function on Z
activation (Z)

Data set

I/O	Study hr	Play hr	O/P
95	4	4	1 \rightarrow Pass
100	5	2	1
95	2	7	0 \rightarrow fail

(1) Forward Propagation:

$$\text{Let } w = [w_1, w_2, w_3] \text{, } b = [0.001]$$

#

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

Step-1 in Hidden neuron or hidden Layer-1

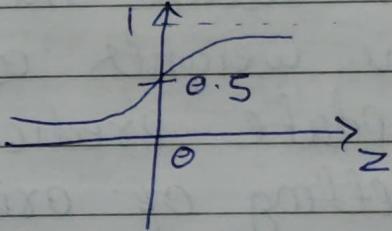
$$Z = 9.5 \times 0.01 + 4 \times 0.02 + 4 \times 0.03 + 1 \times 0.01$$

$$Z = 1.151$$

Step-2: Activation(Z)

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

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



$$\text{Output-one} = \underline{0.759} = o_1$$

Hidden layer-2

$$w_4 = 0.02$$

$$b_2 = 0.03$$

Step-1

$$\begin{aligned} Z &= 0_1 * w_4 + b_2 \\ &= 0.759 \times 0.02 + 0.03 \\ &= \underline{0.04518} \end{aligned}$$

Step-2 Activation(Z) = $\frac{1}{1+e^{-(0.04518)}}$ = $\underline{\underline{0.51129}} = o_2$

O_2 or Predicted o/p (\hat{y}) = 0.51129)

Loss function = error = (original - predicted)

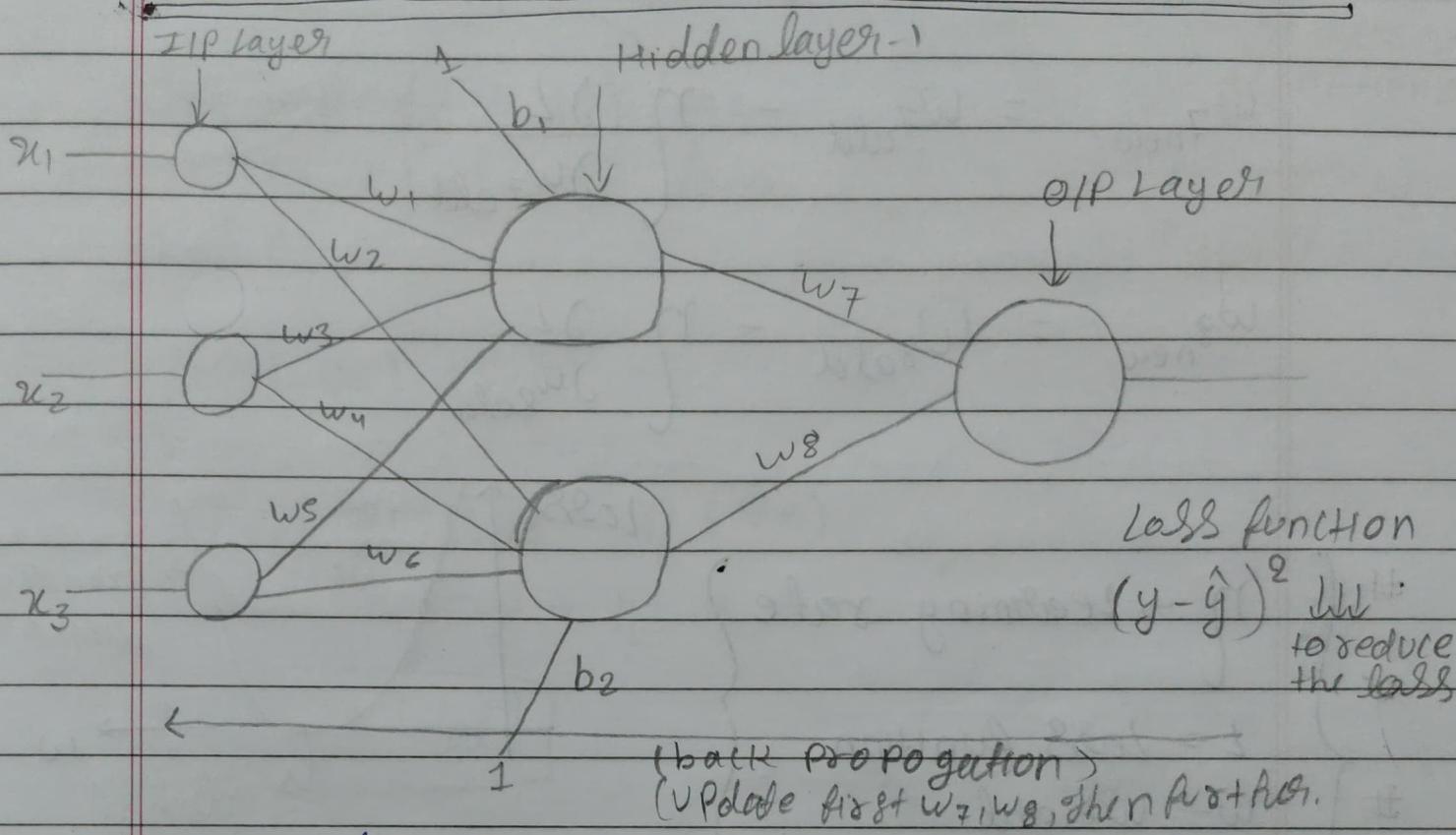
$$= (1 - 0.51129) \\ = \approx 0.49$$

→ So this error is very high so we have to reduce it.

So the way to reduce the loss function we have to update the weights so this is called the backward propagation. because the weights are ~~use~~ predict and again update predict and so on until the getting of original ~~out~~ output.

After the first input the weights are rearranged and then the second input is comes to the pictures and then if it is again going the forward propagation and calculate the loss function if the value is not minimum the back propagation and then the repetition is going on.

Back propagation and weight updation formula.



Loss function

$$(y - \hat{y})^2$$

to reduce the loss

This looks as 3×2 matrix 3 inputs and 2 Hidden layers.

$$\begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} \begin{bmatrix} w_1 & w_2 \\ w_3 & w_4 \\ w_5 & w_6 \end{bmatrix}$$

Some loss functions:

Regression

- (i) MSE
- (ii) MAE
- (iii) Huber loss

Classification

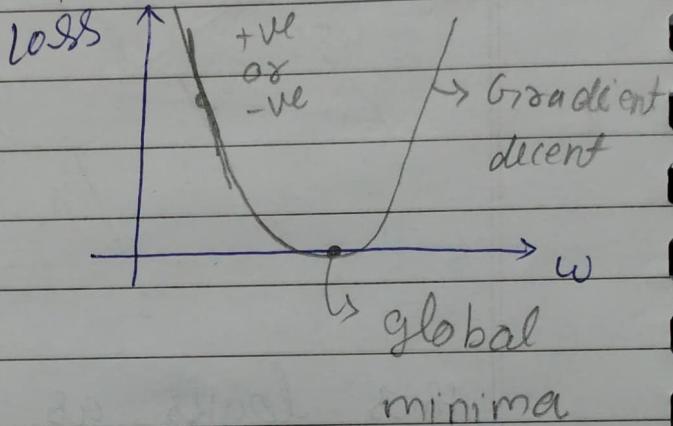
- (i) Binary cross entropy
- (ii) Categorical cross entropy

weight updation formula \rightarrow Learning rate

$$w_7_{\text{new}} = w_7_{\text{old}} - \eta \left(\frac{\partial L}{\partial w_7_{\text{old}}} \right) \rightarrow \text{Slope}$$

$$w_8_{\text{new}} = w_8_{\text{old}} - \eta \left(\frac{\partial L}{\partial w_8_{\text{old}}} \right)$$

$\eta \rightarrow \text{Learning rate}$
 # $L = \text{Loss function}$
 # $w = \text{weight}$

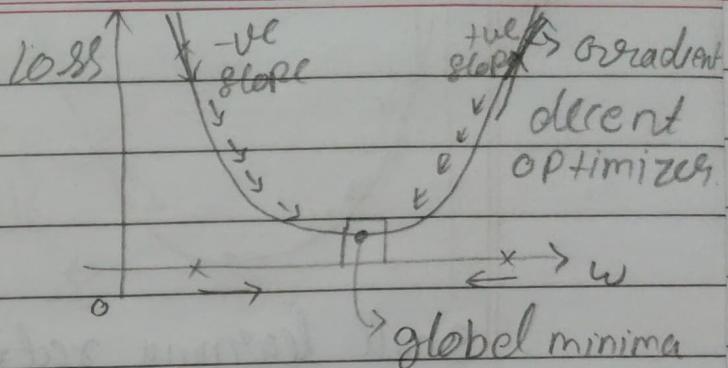


$$w_{\text{new}} = w_{\text{old}} - \eta \left(\frac{\partial L}{\partial w_{\text{old}}} \right)$$

weight updation formula.

Optimizers:

The role of optimizers is to reduce the loss value.

for (-ve) slope

$$\begin{aligned} w_{\text{new}} &= w_{\text{old}} - \eta (-\text{ve}) \\ &= w_{\text{old}} + \eta (+\text{ve}) \end{aligned}$$

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

for (+ve) slope

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

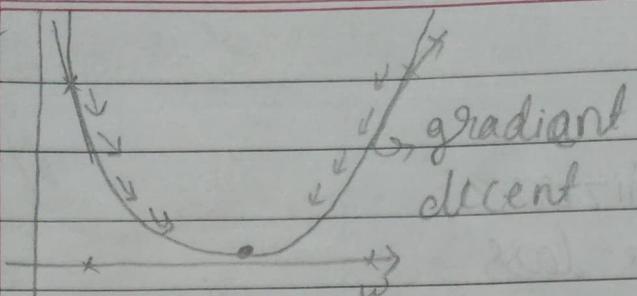
$$w_{\text{new}} \ll w_{\text{old}}$$

Learning rate: So basically learning rate is decide the step size how fast it is achieve the global minima.

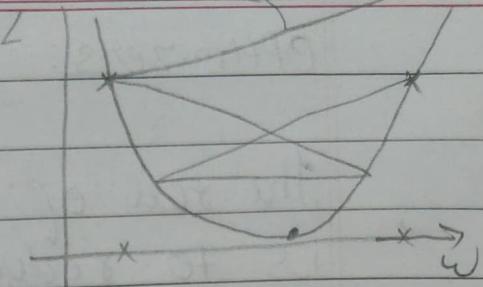
∴

So always take the small learning rate.

Never archive
the minima



small learning rate
 $\eta = 0.01$



high learning rate
 $\eta = \text{large value}$

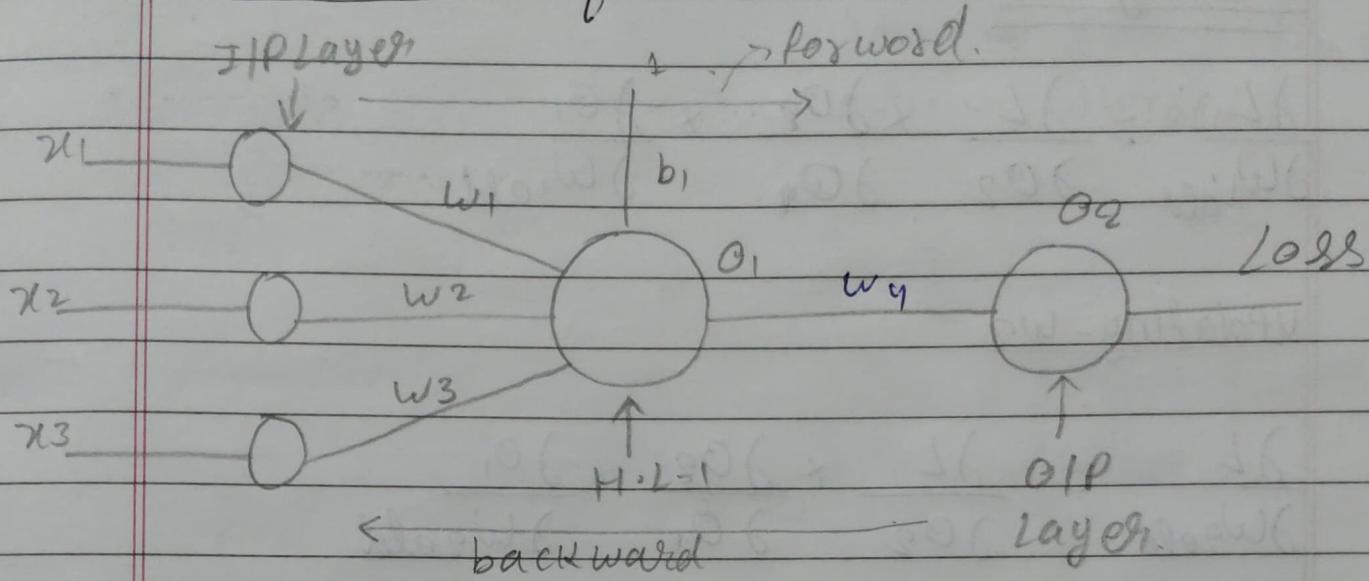
when W reaches the
global minima

} it will just never
converge to get the
minima and reduce
the loss.

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

no need to update
the weight.

Chain Rule of Derivative.



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

Updating w_4 :

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

$$\frac{\partial \text{loss}}{\partial w_{4\text{old}}} = \frac{\partial L}{\partial O_2} \times \frac{\partial O_2}{\partial w_{4\text{old}}}$$

\Rightarrow chain rule of derivation.

Explanation: Loss is depend on the O_2 (Output-2)

and the O_2 (Output-2) is depend on the w_4 .

Updating - w_1

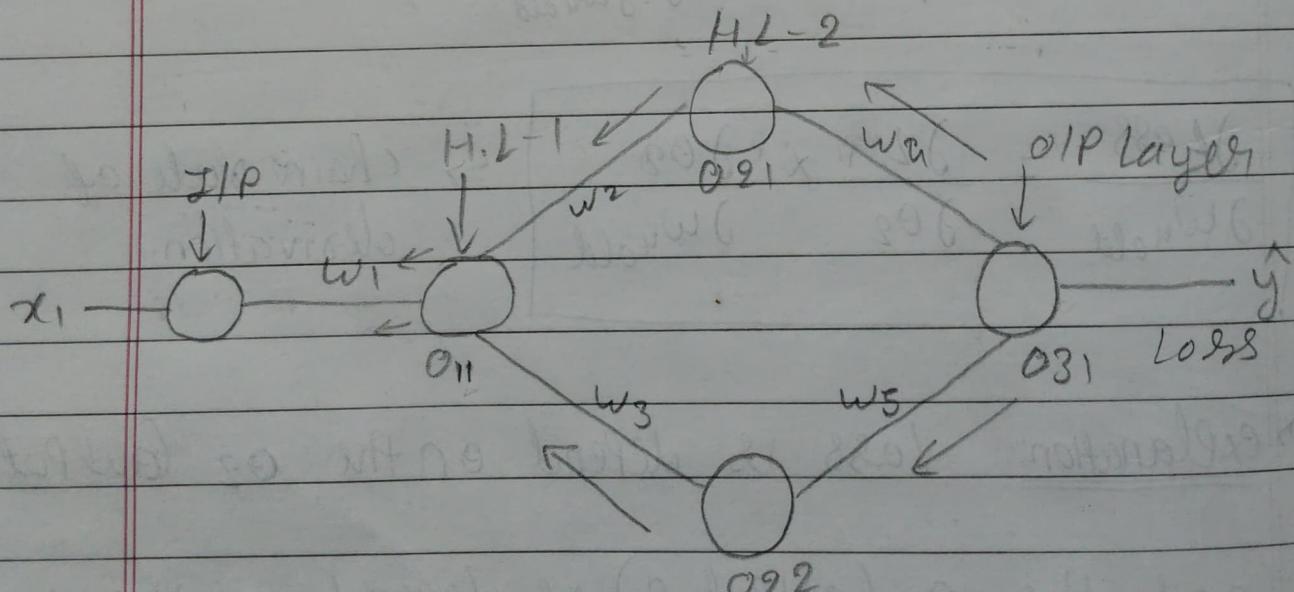
$$\frac{\partial L}{\partial w_{1, \text{old}}} = \frac{\partial L}{\partial o_2} \times \frac{\partial o_2}{\partial o_1} \times \frac{\partial o_1}{\partial w_{1, \text{old}}}$$

Updating - w_2

$$\frac{\partial L}{\partial w_{2, \text{old}}} = \frac{\partial L}{\partial o_2} \times \frac{\partial o_2}{\partial o_1} \times \frac{\partial o_1}{\partial w_{2, \text{old}}}$$

Updating - w_3

$$\frac{\partial L}{\partial w_{3, \text{old}}} = \frac{\partial L}{\partial o_2} \times \frac{\partial o_2}{\partial o_1} \times \frac{\partial o_1}{\partial w_{3, \text{old}}}$$

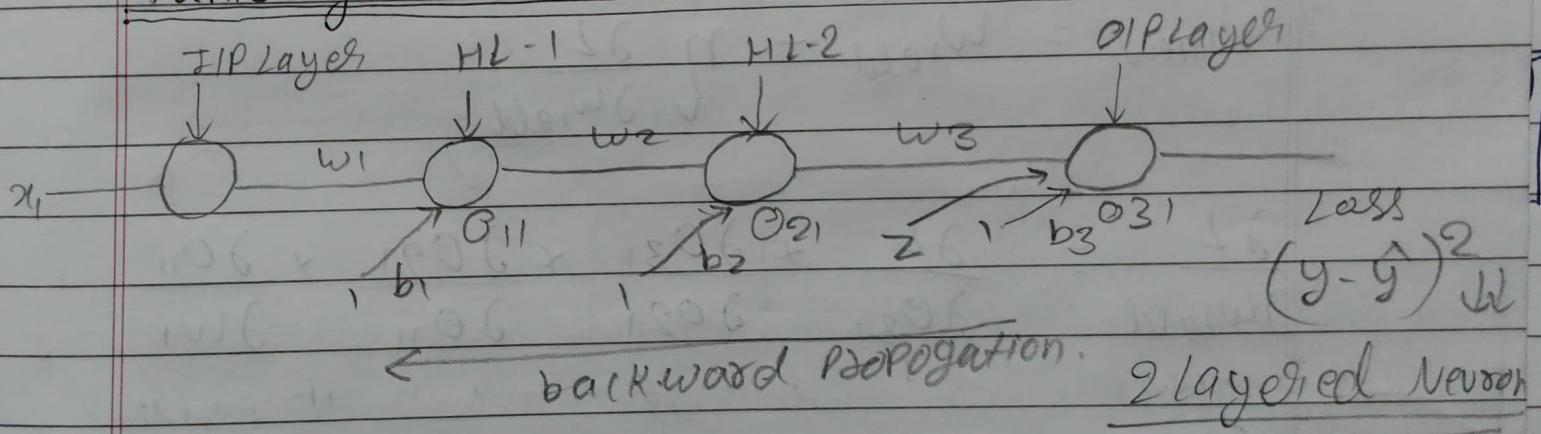
Updating - w_1

$$\frac{\partial L}{\partial w_{1,old}} = \left[\frac{\partial L}{\partial \theta_{31}} \times \frac{\partial \theta_{31}}{\partial \theta_{21}} \times \frac{\partial \theta_{21}}{\partial \theta_{11}} \times \frac{\partial \theta_{11}}{\partial w_{1,old}} \right] \rightarrow 1^{\text{st}} \text{ side}$$

+

$$\left[\frac{\partial L}{\partial \theta_{31}} \times \frac{\partial \theta_{31}}{\partial \theta_{22}} \times \frac{\partial \theta_{22}}{\partial \theta_{11}} \times \frac{\partial \theta_{11}}{\partial w_{1,old}} \right] \rightarrow 2^{\text{nd}} \text{ side}$$

Vanishing Gradient Problem And Activation Function.



first step in every neuron is

$$\{ z = w_3 \times \theta_{21} + b_3 \}$$

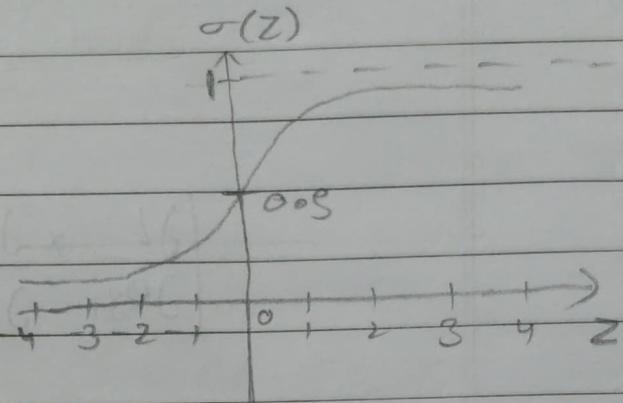
$$z = \sum_{i=1}^n w_i x_i + b$$

and second step to apply the activation function.

Sigmoid activation function (σ)

$\sigma(z)$ to convert the value between 0 to 1

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



Updating the w_1

values 0.01

$$w_{1\text{new}} = w_{1\text{old}} - \eta \left[\frac{\partial L}{\partial w_{1\text{old}}} \right]$$

$$\frac{\partial L}{\partial w_{1\text{old}}} = \frac{\partial L}{\partial o_{31}} \times \boxed{\frac{\partial o_{31}}{\partial o_{21}}} \times \boxed{\frac{\partial o_{21}}{\partial o_{11}}} \times \boxed{\frac{\partial o_{11}}{\partial w_1}}$$

for $o_{31} \Rightarrow$

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

$\underbrace{w_3 \times o_{21} + b_3}_{z}$

$$z = o_{21} \times w_3 + b_3$$

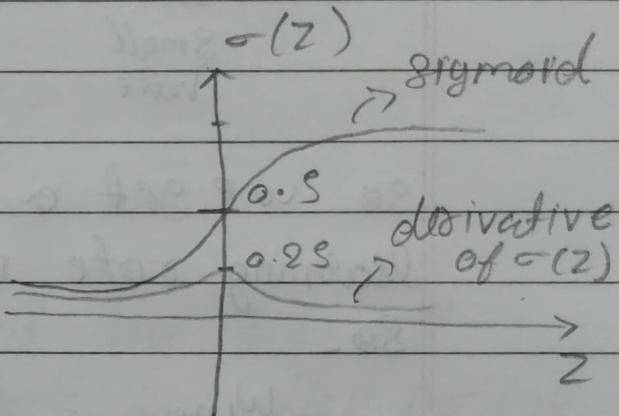
$$o_{31} = \sigma(z)$$

$$\frac{\partial O_3}{\partial O_2} = \frac{\partial(\sigma(z))}{\partial z} \times \frac{\partial z}{\partial O_2} \quad \left. \right\} \text{chain rule.}$$

$0 \leq \sigma(z) \leq 1 \rightarrow$ when we applied the sigmoid

derivative of $\sigma(z)$

$$0 \leq \sigma(z) \leq 0.25$$



$$\frac{\partial(\sigma(z))}{\partial z} = 0 \leq \sigma(z) \leq 0.25 \times \frac{\partial((w_3 \times O_2) + b_3)}{\partial O_2}$$

$$\Rightarrow \frac{\partial O_3}{\partial O_2} = 0 \leq \sigma(z) \leq 0.25 \times \frac{\partial((w_3 \times O_2) + b_3)}{\partial(O_2)}$$

$$\Rightarrow \frac{\partial O_3}{\partial O_2} = 0 \leq \sigma(z) \leq 0.25 \times w_{3, \text{old}}$$

so we get the $\frac{\partial O_3}{\partial O_2}$ in the range of $(0 - 0.25)$

which is small.

So in case of updating the w_i weight

$$\frac{\partial L}{\partial w_{i,\text{old}}} = \frac{\partial L}{\partial \theta_{31}} \times \frac{\partial \theta_{31}}{\partial \theta_{21}} \times \frac{\partial \theta_{21}}{\partial \theta_{11}} + \frac{\partial \theta_{11}}{\partial w_{i,\text{old}}}$$

↓ ↓ ↓ ↓
 same (0 - 0.25) (0 - 0.25) (0 - 0.25)
 small
value

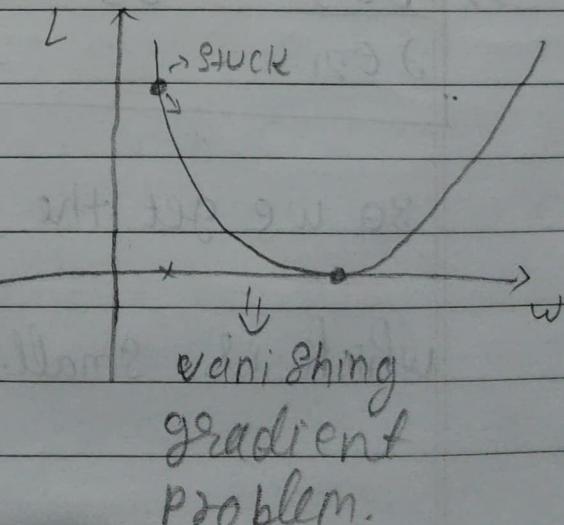
so we get a smaller value and our learning rate is also a smaller value

So -

$$w_{i,\text{new}} = w_{i,\text{old}} - \left(\eta \frac{\partial L}{\partial w_{i,\text{old}}} \right) \rightarrow \text{very small.}$$

$$\Rightarrow \{ w_{i,\text{new}} \approx w_{i,\text{old}} \}$$

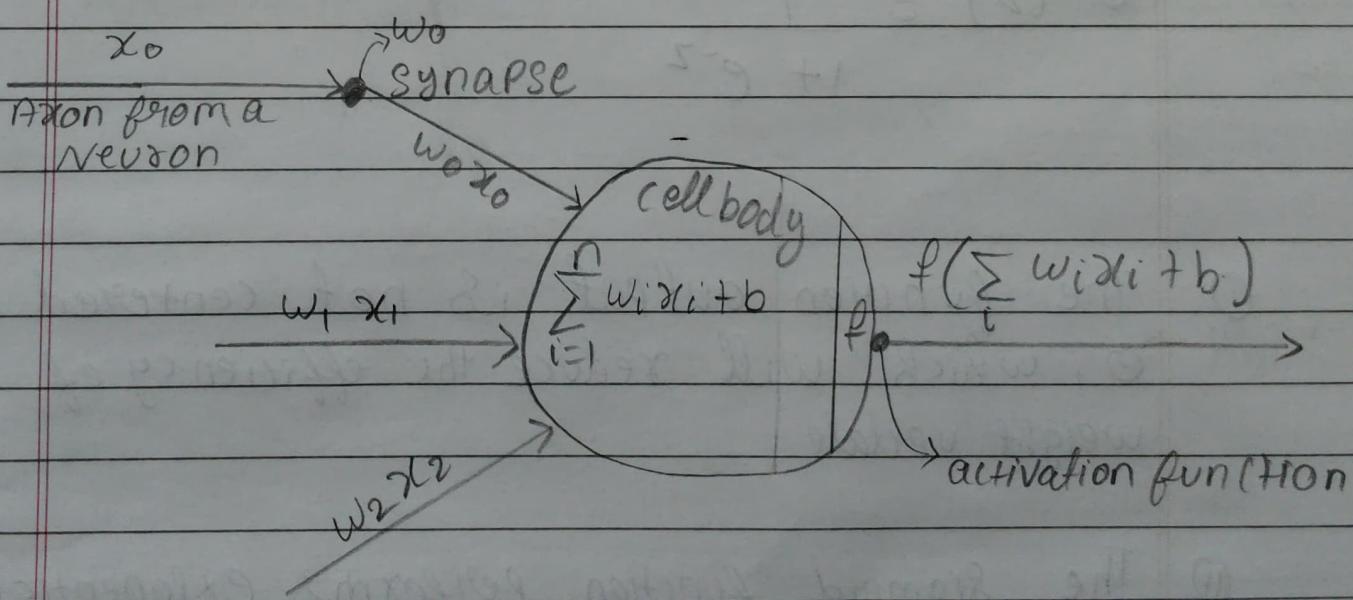
So in this case $w_{i,\text{old}}$ is like never updated. and stuck to achieve the global minima this is called the gradient vanishing problem.



To fix the Vanishing gradient problem the Researchers started to exploring other activation function.

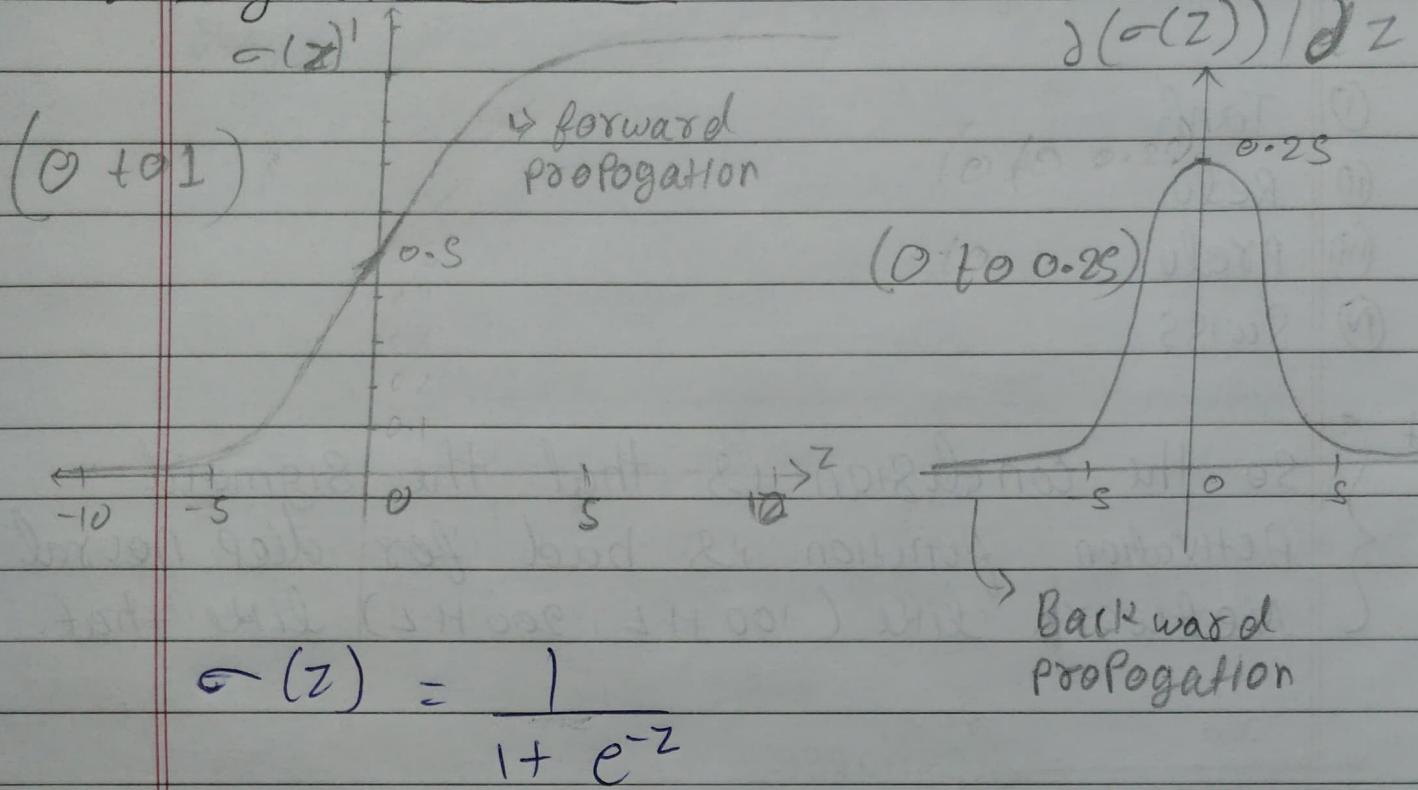
- (I) Tanh
- (II) Relu
- (III) Prelu
- (IV) Swiss

{ So the conclusion is that the Sigmoid - Activation function is bad for deep neural Network like (100 H-L, 200 H-L) like that.



Activation Functions.

1) Sigmoid function:



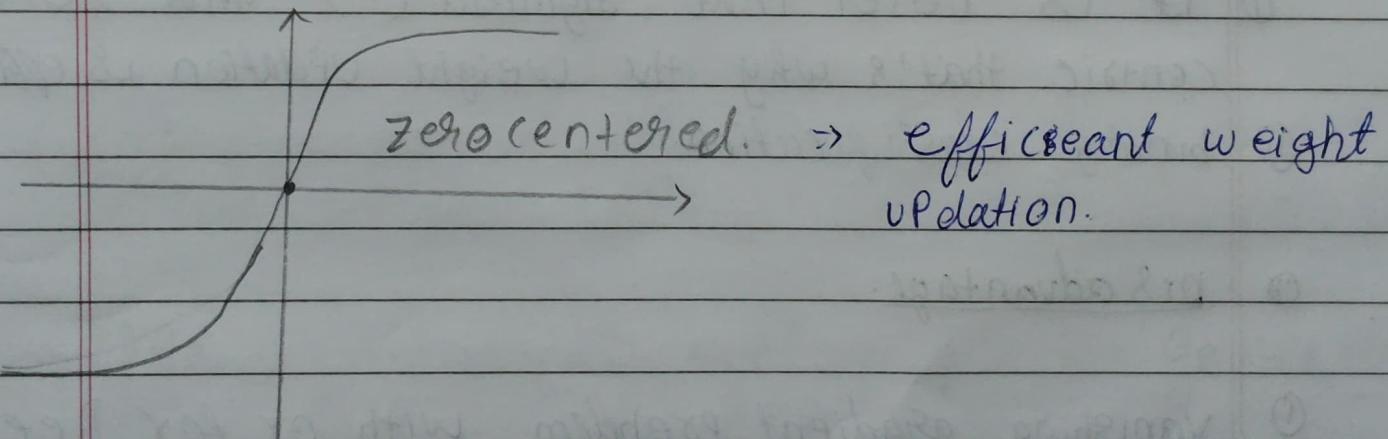
- ① The function output is not centered on 0, which will reduce the efficiency of weight update.
- ② The sigmoid function performs exponential operations, which ~~will~~ is slower the computer.

Advantage of Sigmoid function:-

- ① Smooth gradient, preventing "Jumps" in output
- ② Output values bound between 0 and 1, normalizing the output of each neuron.
- ③ Clear Predictions, i.e. very close to 1 or 0.
- ④ Suitable for binary classification.

Disadvantages:

- ① prone the Gradient vanishing.
- ② function output is not zero-centered.
- ③ Power operations are relatively time consuming.

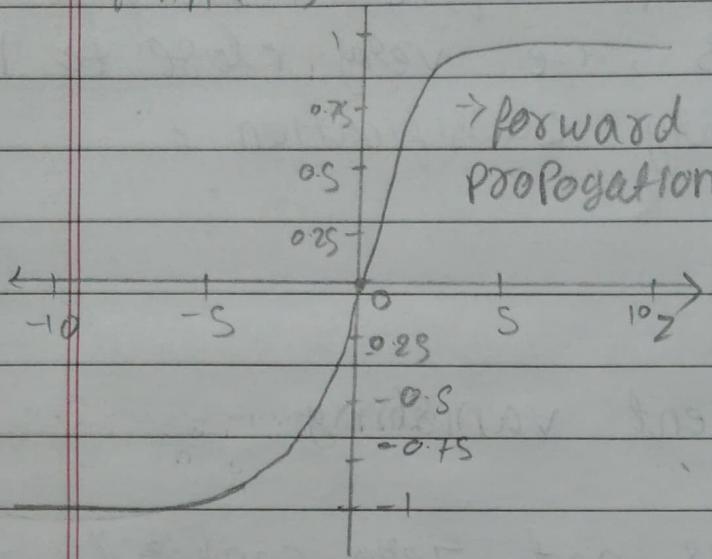


② tanh function:

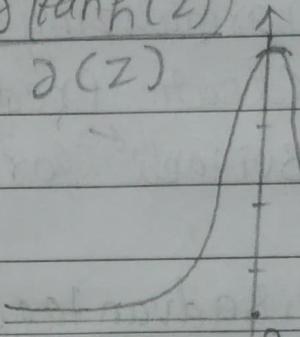
$$\{ z = \sum_{i=1}^n w_i x_i + b \}$$

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

$\tanh(z)$



$$\frac{\partial \tanh(z)}{\partial z}$$



Advantage:

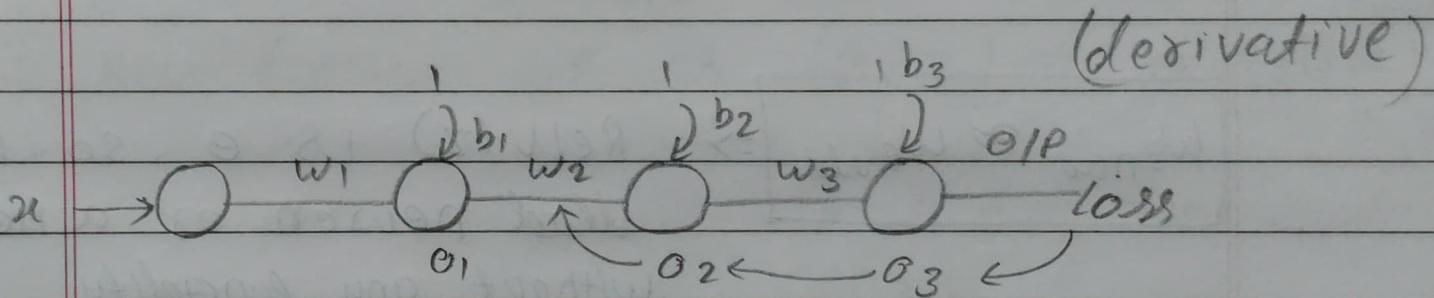
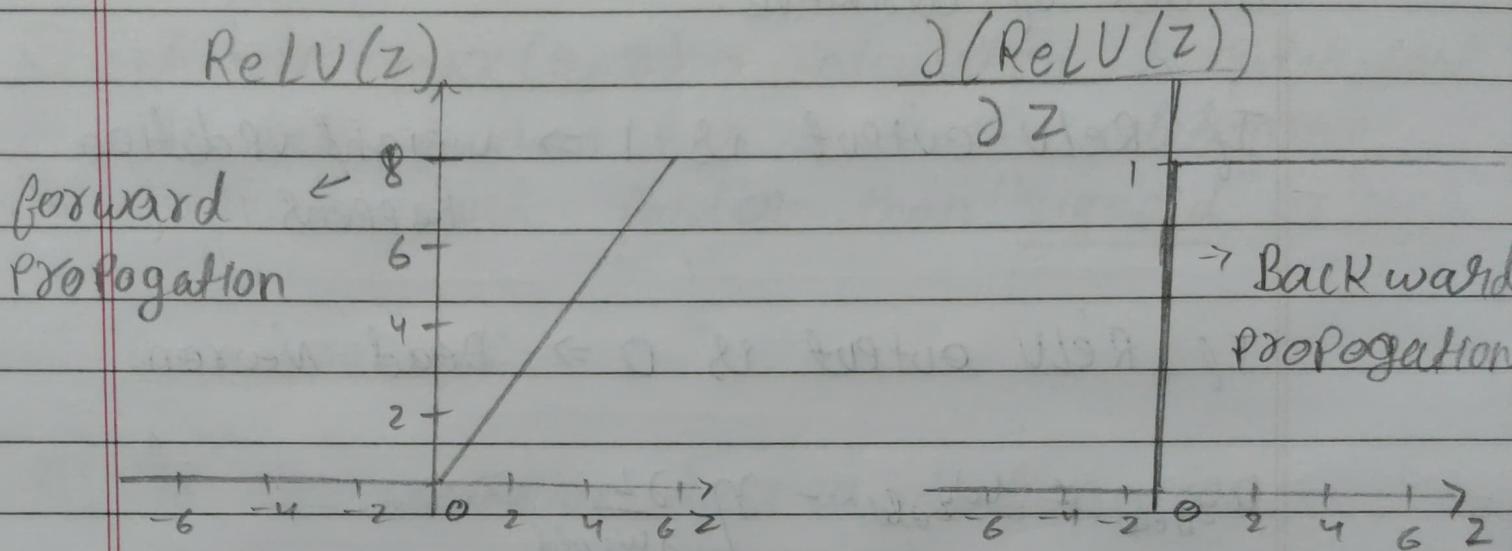
- ① It is better than sigmoid (σ) and zero-centric that's why the weight updation is efficient.
- ② binary classification.

Disadvantage:

- ① Vanishing gradient problem with or for Deep neural network.
- ② It take too much time for computation because of exponential power.

(2) ReLU Activation function:-

$$\text{ReLU}(z) = \max(0, z) \rightarrow \text{formula.}$$



$$\Rightarrow \frac{\partial L}{\partial w_{2, \text{old}}} = \frac{\partial L}{\partial o_3} \times \left| \frac{\partial o_3}{\partial o_2} \right| \times \frac{\partial o_2}{\partial w_2} \quad \left\{ w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w_{\text{old}}} \right\}$$

$$\Rightarrow \frac{\partial o_3}{\partial o_2} = \left| \frac{\partial (\text{ReLU}(z))}{\partial z} \right| \times \frac{\partial z}{\partial o_2} \quad \left\{ \because z = o_2 \times w_3 + b_3 \right\}$$

$\downarrow \downarrow \downarrow$

$$\left[\begin{matrix} 0 & 0 & 1 \end{matrix} \right] \times \frac{\partial [(o_2 \times w_3) + b_3]}{\partial o_2}$$

$$\Rightarrow [0 \ 0 \text{ or } 1] \times w_3 \Rightarrow (\text{small value get.})$$

In case of derivative.

If ReLU output is 1 \Rightarrow weight updation happens.

If ReLU output is 0 \Rightarrow Dead Neuron.

$$w_{2\text{new}} = w_{2\text{old}} - \eta \left[\frac{\partial L}{\partial w_{2\text{old}}} \right] > 0$$

$w_{2\text{new}} \approx w_{2\text{old}}$ \Rightarrow $\text{ReLU}(z)$ is 0, so it is dead neuron or a neuron without any functionality.

① if $z = +ve$ $\frac{\partial(\text{ReLU}(z))}{\partial z} = 1$

② if $z = -ve$ $\frac{\partial(\text{ReLU}(z))}{\partial z} = 0$

\hookrightarrow the dead neuron is just pass the input as it is to the next neuron.

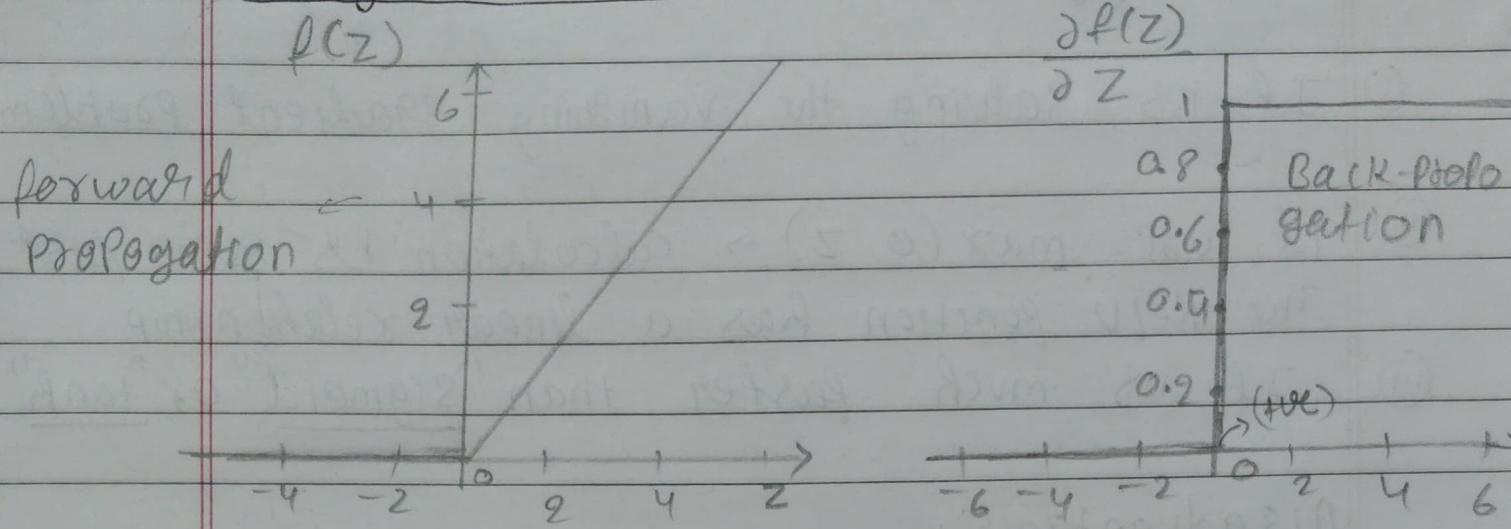
Advantages:

- ① It is solving the vanishing gradient problem
- ② It use $\max(0, z) \rightarrow$ calculation is super-fast
The ReLU function has a linear relationship
- ③ It is much faster than "Sigmoid" or "tanh"

Disadvantages:

- ① It is create or make a dead neuron.
- ② ReLU function O/P is $(0, z) \Rightarrow$ ○ ○ & (+ve)
number
 \Downarrow
If is not zero
centric.

(4). Leaky ReLU and parametric ReLU:-



$$\begin{array}{l} (\text{Leaky}) \rightarrow \text{formula} = \\ (\text{ReLU}) \quad f(z) = \max(0.01z, z) \end{array}$$

$\alpha \rightarrow [0.01, 0.02, 0.03]$

$$\begin{array}{l} (\text{Parametric}) \rightarrow \text{formula} \Rightarrow f(z) = \max(\underbrace{\alpha z}_{\text{hyper parameter}}, z) \\ (\text{ReLU}) \end{array}$$

So in this Leaky ReLU activation function have some (+ve) value instead of 0 so it is never give the output 0 0

So it remove the problem of Dead neuron.

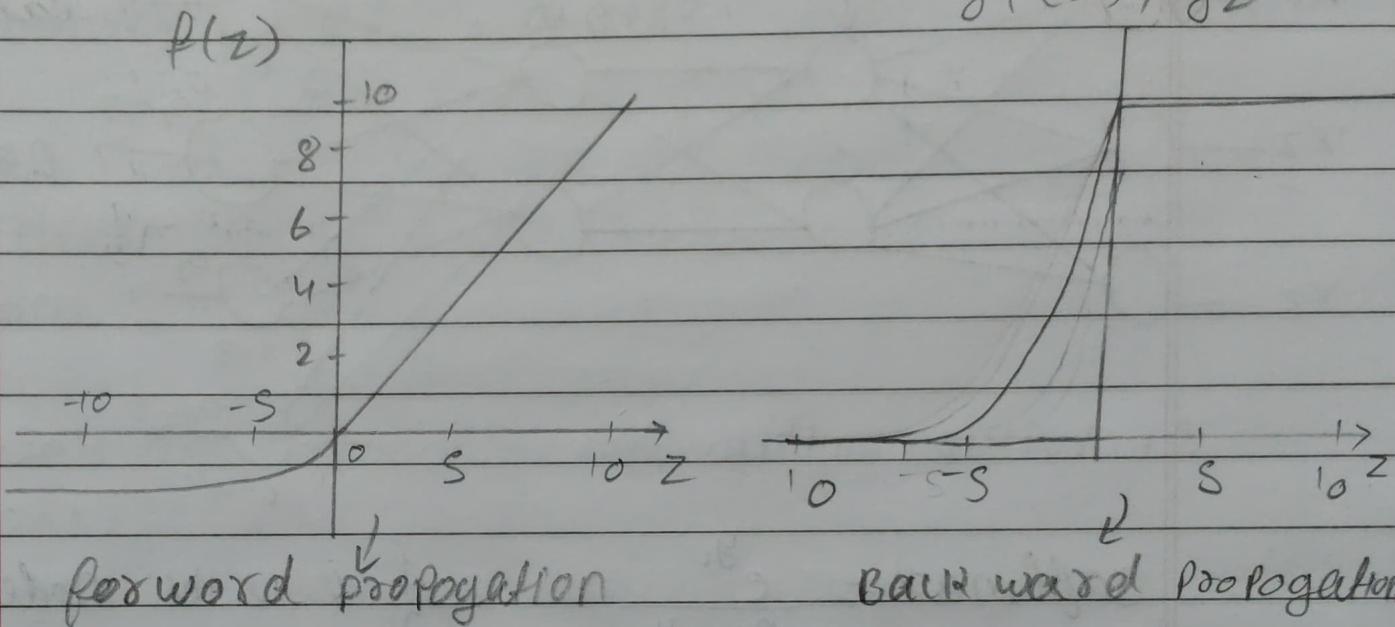
ReLU \Rightarrow Rectified Linear Unit.

It is not a zero-centric.

⑤ ELU (Exponential Linear Units)

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

$$\frac{\partial f(z)}{\partial z}$$



Advantages

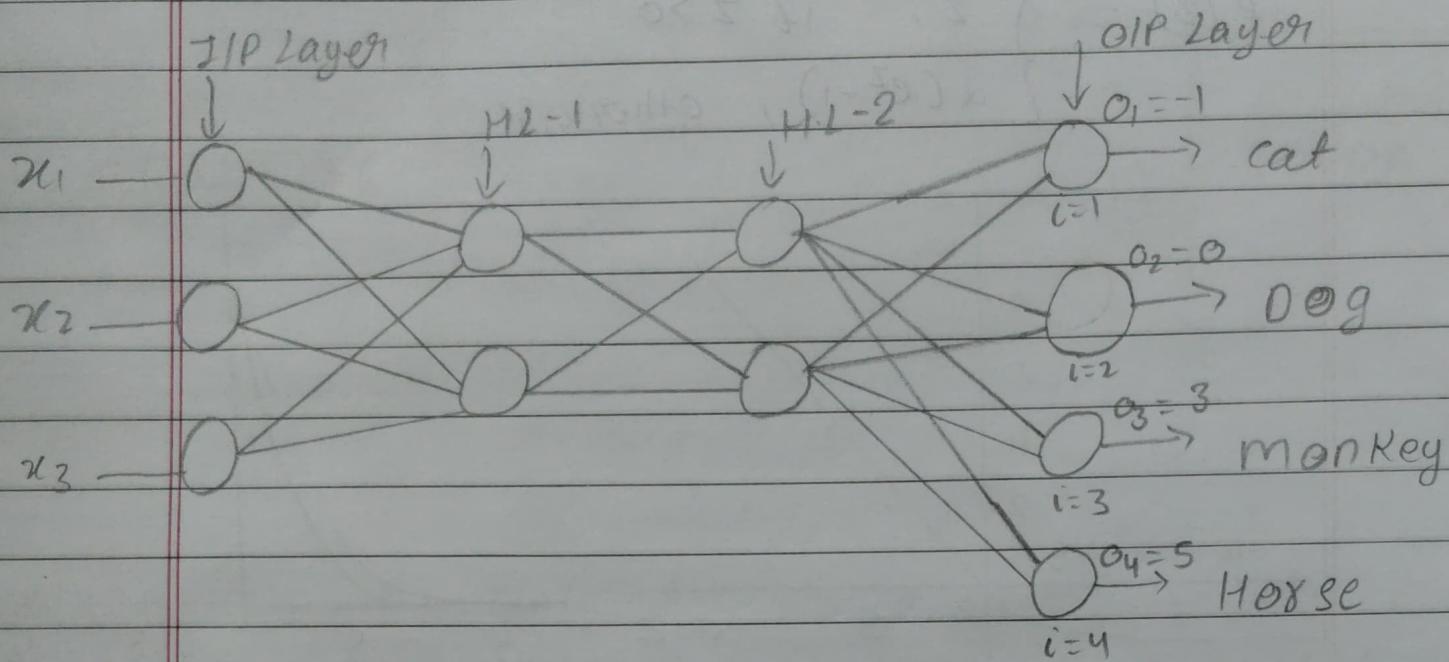
- ① No dead ReLU issues
- ② It is zero centric

Disadvantages.

- ① Slightly more computationally intensive.

If it is solve the ReLU problem. }

⑥ Softmax activation function \Rightarrow [multi class classification problem].



$$\text{softmax} = \frac{e^{y_i}}{\sum_{k=1}^n e^{y_k}}$$

{ i is the node / neuron no.

y_i = output of that neuron.

$y_i = \text{output} \times w + b$

apply the softmax:

$$\text{Cat} = \frac{e^{-1}}{e^{-1} + e^0 + e^3 + e^5} = \frac{e^{-1}}{e^{-1+0+3+5}} = \frac{e^{-1}}{e^7} = 0.00033$$

if $i=1$ $e^{y_i} = e^{-1}$ because y_1 of output is -1

$y_2 = 0, y_3 = 3, y_4 = 5$.

$$\text{Dog} = \frac{e^0}{e^{-1+0+3+5}} = \frac{e^0}{e^7} = 0.0024$$

$$\text{monkey} = \frac{e^3}{e^7} = 0.0183$$

$$\text{Horse} = \frac{e^5}{e^7} = 0.1353$$

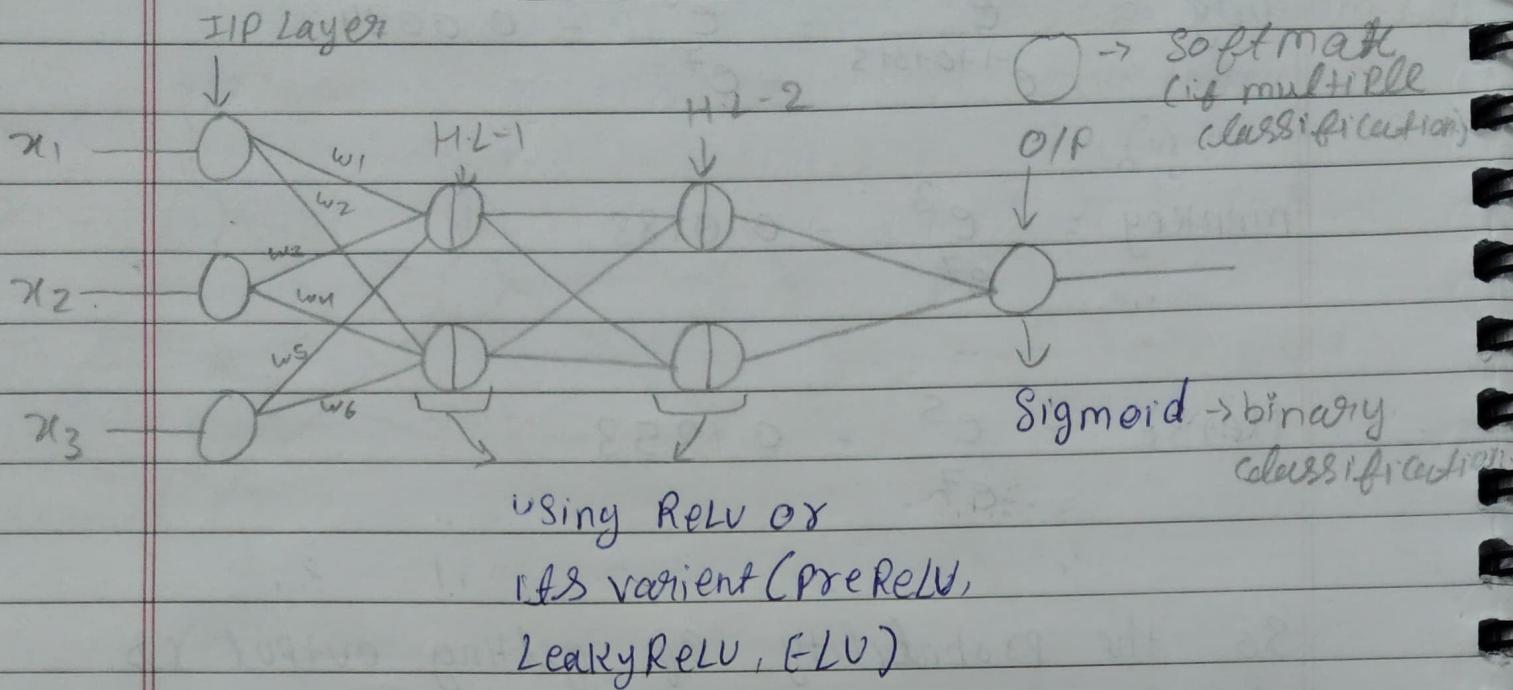
So the probability of getting output is horse is

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

$$= 0.86 \approx 86\%$$

So there is 86% chance to the output is horse at the given INPUT, so likewise we can get or find out the others probability.

⑦ Which activation function To ~~use~~ / apply when?



So if your neural network is very deep then use the ReLU activation function.

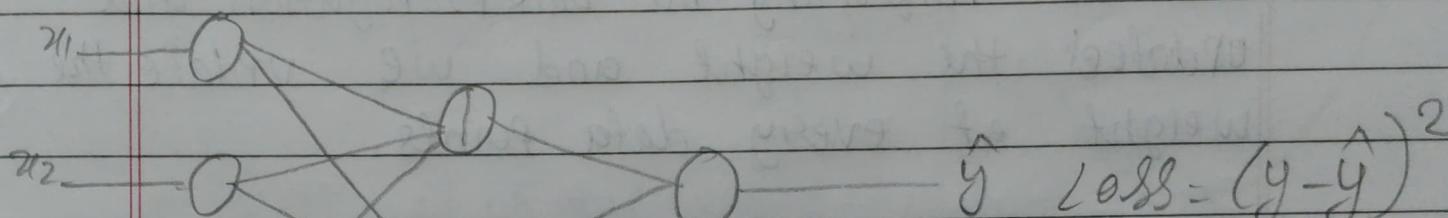
① if we use the sigmoid activation function or tanh, it leads the vanishing gradient problem So to handle these type of problem we use the ReLU activation function or its variants.

② but in the output layer if we require the binary output then inside the O/P layer we are using the Sigmoid activation function.

=> If we require the multiple classification like dog, cat, horse etc. we are using the Softmax activation function in the O/P layer.

Loss function And Cost function

forward



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

backward

	x_1	x_2	x_3	O/P
-	-	-	-	0
-	-	-	-	1
-	-	-	-	0
-	-	-	-	1

loss function

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

cost function:

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

update at every point.

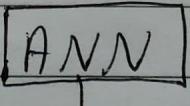
→ take the batch and update once

Loss function:

In case of loss function we give the input and forward ~~for~~ propagation happens then it is give the \hat{y} the predicted value and we calculate the error and to reduce that error we are performing the back propagation and updated the weight and we update the weight at every data points.

Cost function:

In case of cost function we give the inputs at once and after performing the operations. and get the \hat{y} value after this it is calculate the errors and ~~is~~ update the weight once only. in the back propagation.



- ① mean square Error (MSE)
- ② mean Absolute Error (MAE)
- ③ Huber Loss
- ④ RMSE (Root mean square Error)

Different Loss and cost function in the Regression Problem

classmate

Date _____

Page _____

① Mean Squared Error (MSE) \Rightarrow

$$\text{Loss function} = (y - \hat{y})^2 \rightarrow \text{Quadratic equation}$$

Cost function:

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

Loss

or

Cost
fun

gradient
gradient of

outliers: it is that type of values that deviate significantly from other values in a data set.

global minima

Advantage:

(i) MSE is differentiable at any point.

(ii) It has only 1 local minima or global minima.

(iii) It converges faster because of smooth gradient descent

Disadvantage:

• Not robust to outliers.

outlier
penalizing
the error
(\uparrow)

in case of outlier the best fit line shift much towards the outlier

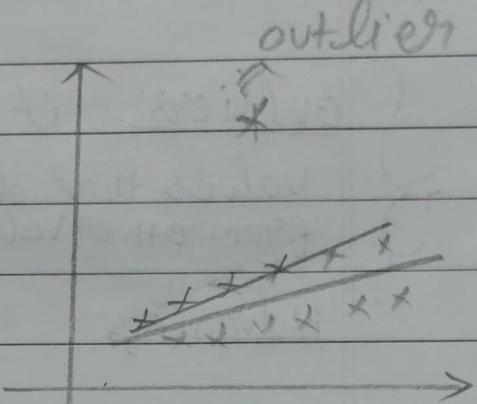
② mean Absolute Error (MAE)

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

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

Advantage:

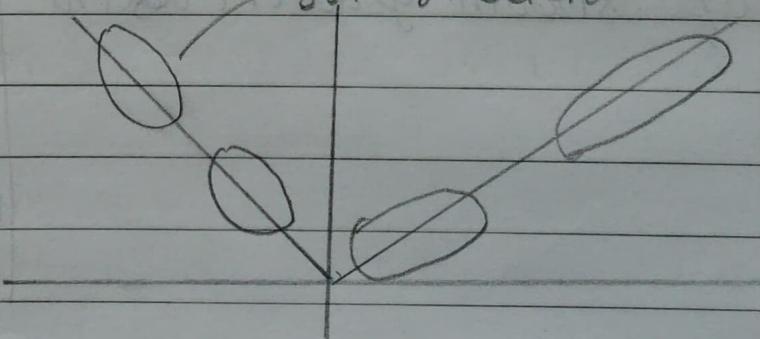
It is Robust to outlier.
in case of coming of outlier
the best fit line shift very
less toward the outlier.



Disadvantage:

the convergence take time ^{in MAE} because the curve
is not quadratic so it is create a
Sub-gradient descent and calculate the slope.

→ Sub-gradient



(3) Huber loss: It is a combination of MSE and MAE with the ~~some~~ outliers and Hyperparameter. (MSE) ^{no outliers present} ^{hyperparameter}

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

(MAE)

outlier present

(4) Root mean squared Error (RMSE)

$$\text{cost function} := \sqrt{\frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{N}}$$

N = No. of Data points.

Advantages:

Differentiable at any point

easy Interpretability of different models.

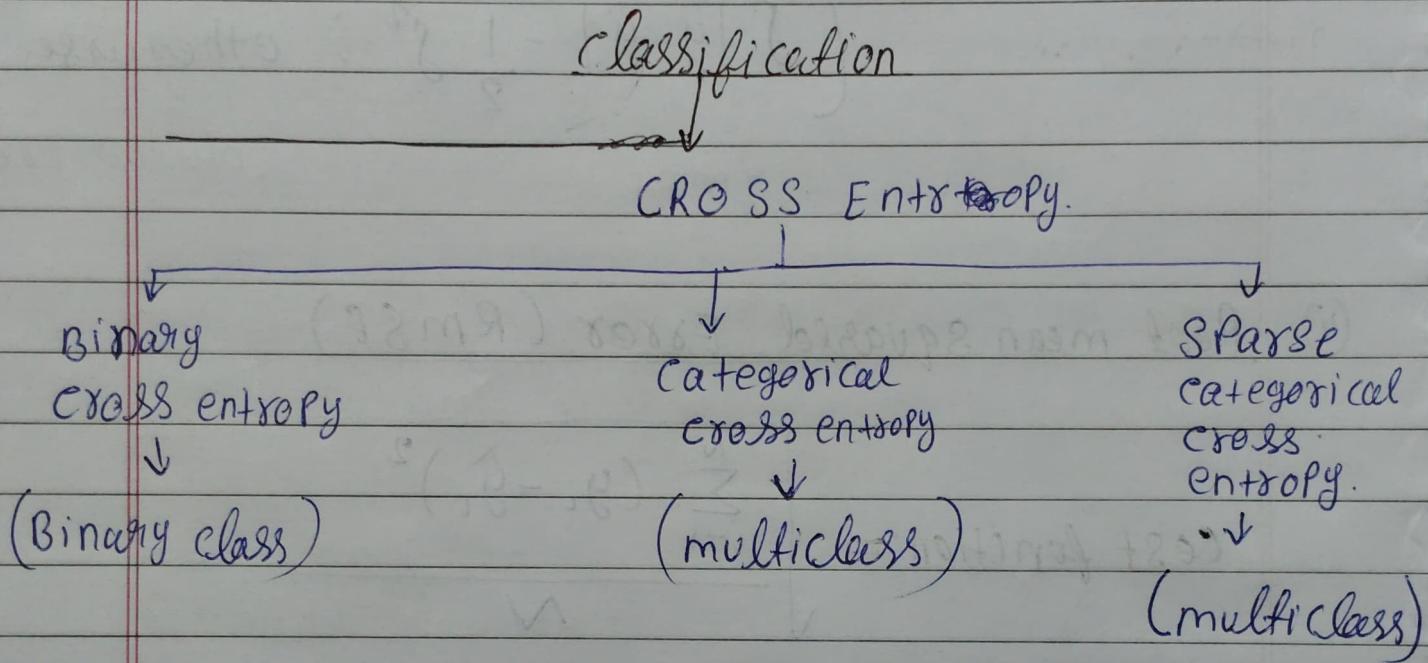
Disadvantages:

Not Robust to outliers.

scale dependence.

Outliers: It is a data points that behaves differently from all other data points in the data set.

Loss or cost function for classification problems:



① Binary cross Entropy

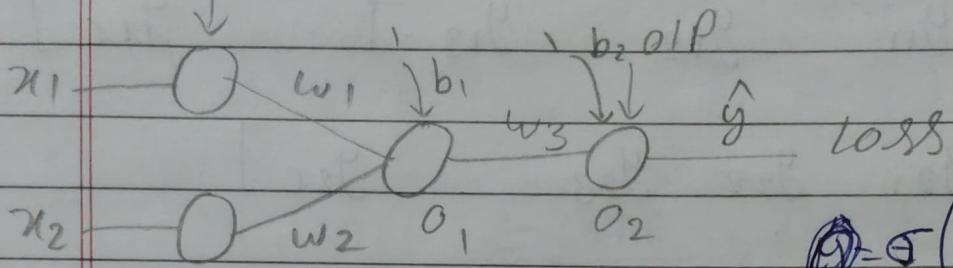
Log loss function :-

$$\text{loss} = -y \times \log(\hat{y}) - (1-y) \times \log(1-\hat{y})$$

$$\left. \begin{array}{l} y \Rightarrow \text{actual value} \\ \hat{y} \Rightarrow \text{predicted value.} \end{array} \right\}$$

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

I/O



$$\hat{y}_2 = \sigma(w_3 \times o_2 + b_2)$$

$$\boxed{\hat{y} = \frac{1}{1+e^{-z}}} \rightarrow \text{Sigmoid.}$$

② Categorical Cross Entropy (multiclassification)

f_1	f_2	f_3	O/I/P	\Rightarrow OHE \Rightarrow (One Hot encoding)
2	3	4	Good	$C = \text{No. of categories.}$
5	6	7	Bad	
8	9	10	Neutral	

on the output feature.

$j=1$	$j=2$	$j=3$	\Rightarrow When output is Good the Good is 1 and others are 0.
Good	Bad	Neutral	
1	0	0	
0	1	0	
0	0	1	

$$\text{Loss}(x_i, y_i) = - \sum_{j=1}^c y_{ij} * \ln(\hat{y}_{ij})$$

$$1^{\text{st}} \text{ data point} \rightarrow y_{1j} = [y_{11} \quad y_{12} \quad y_{13} \quad \dots \quad y_{1c}]$$

$$2^{\text{nd}} \text{ data point} \rightarrow = [y_{21} \ y_{22} \ y_{23} \ \dots \ y_{2C}]$$

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

Prediction $\rightarrow \hat{y}_{ij} \Rightarrow$ apply softmax activation function
on the oIP layer.

$$\Rightarrow S_{\text{eff}}(z) = \frac{e^z}{\sum_{j=1}^K e^{z_j}}$$

② OIP $\Rightarrow \hat{y}_{ij}$ is in form of Probabilities,

$\Rightarrow [0.2, 0.5, 0.3] \Rightarrow$ if we have

↳ Categorical cross Entropy. three categories.

The Categorical cross entropy is also give the probability of other categories.

⑤ Sparse Categorical Cross Entropy:

\hat{y} = [0.2, 0.3, 0.5] $^{0^{\text{th}}} \text{ 1}^{\text{st}} \text{ 2}^{\text{nd}} \rightarrow \text{categories.}$

O/P $\leftarrow \begin{cases} \text{2}^{\text{nd}} \text{ index have highest probability} \\ \downarrow \\ \text{category} \end{cases}$

so it is not going to give the others probability.

Suppose you have 5 categories and their

\hat{y} = [0.2, 0.3, 0.1, 0.2, 0.2] $^{0^{\text{th}}} \text{ 1}^{\text{st}} \text{ 2}^{\text{nd}} \text{ 3}^{\text{rd}} \text{ 4}^{\text{th}} \rightarrow \text{category}$

so your final output is 1st index category directly.

Disadvantage \Rightarrow in case of sparse we

lose the info about the probability of other categories.

which loss function is use when.

	Hidden layer activation function	O/P layer	Problem	Loss function	① Grad
①	ReLU or (var)	Sigmoid	Binary classification	Binary cross entropy	② Sto
②	ReLU ..	softmax	multi class classification	categorical / sparse cross Entropy.	③ mini SGD
③	ReLU	Linear	Regression	MSE, MAE, Huber loss, RMSE	④ Adag

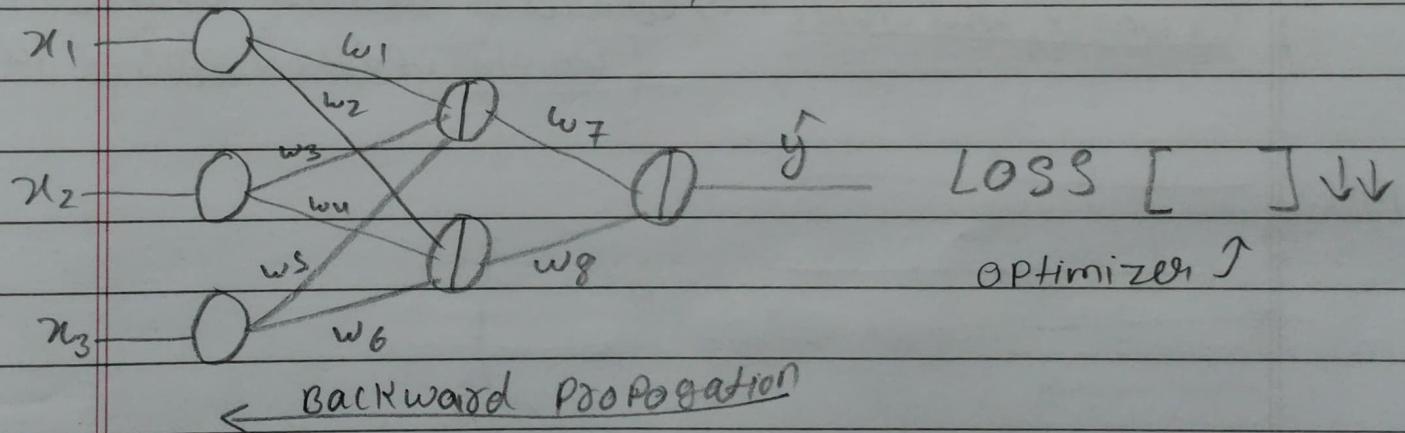
Optimizers.

- ① Gradient Descent
- ② Stochastic Gradient Descent (SGD)
- ③ mini batch SGD
- ④ SGD with momentum
- ⑤ Adagard and RMSprop
- ⑥ Adam optimizer.

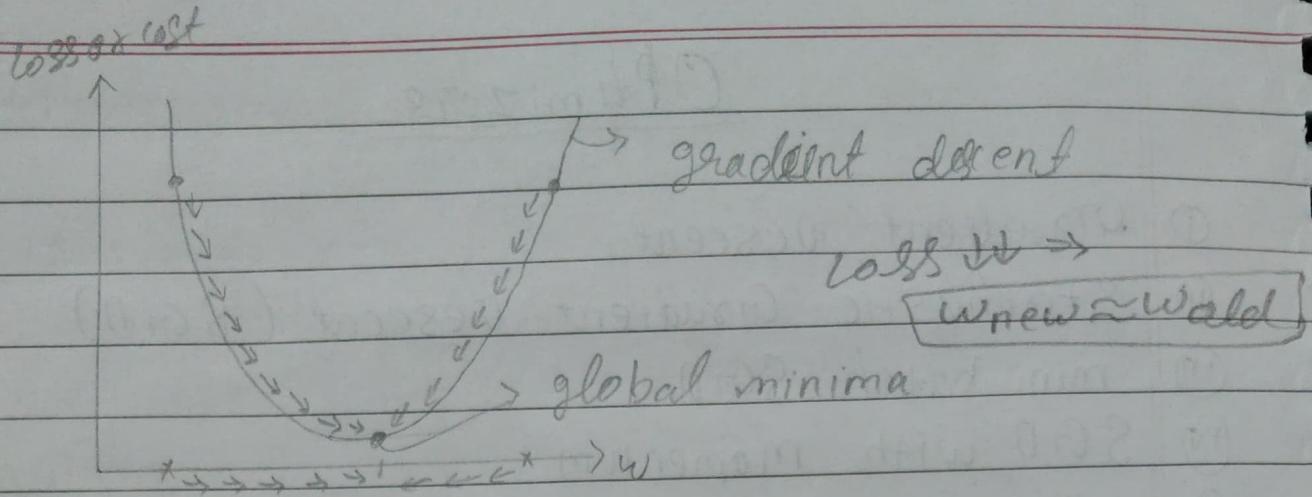
Optimizers are responsible for reducing the loss function in the back propagation.

(i) Gradient Descent Optimizer:-

forward propagation →



$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w_{\text{old}}} \quad \begin{cases} \rightarrow \text{weight updation} \\ \rightarrow \text{Learning rate} \end{cases}$$

MSE

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

cost function

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

Epochs, iterationData set \Rightarrow 1000 Data Point.

1-Epoch

1000 data point $\rightarrow \hat{y} \rightarrow$ cost function J(w)

\leftarrow weight will update \rightarrow gradient descent optimizer
 [weight updation applied]

2-Epoch

\rightarrow
 \leftarrow

100 Epoch

\rightarrow
 \leftarrow

Advantages:

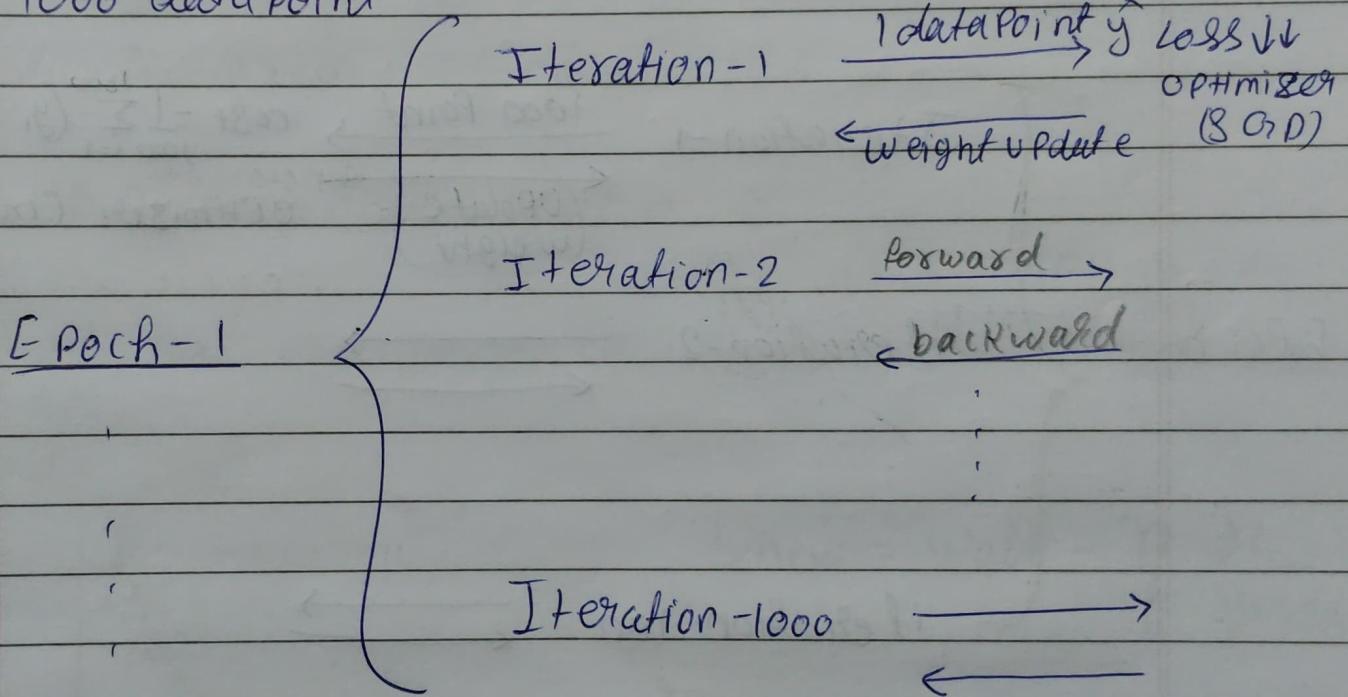
① convergence will happen.

Disadvantage:

① It \oplus require huge resource like RAM, GPU
 \downarrow
 it is Resource intensive.

② Stochastic Gradient Descent (SGD)

1000 data point



\Rightarrow Advantage:

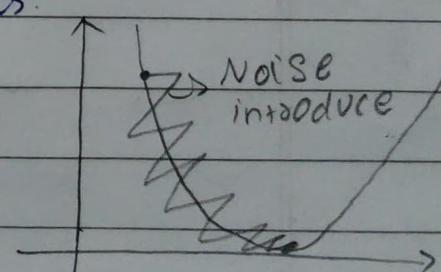
\Rightarrow it is solve the resource issues.

\Rightarrow Disadvantages:

\Rightarrow Huge Time complexity ($\uparrow\uparrow$)

\Rightarrow Convergence take more time

\Rightarrow Noise get introduce.



③ Mini batch SGD:

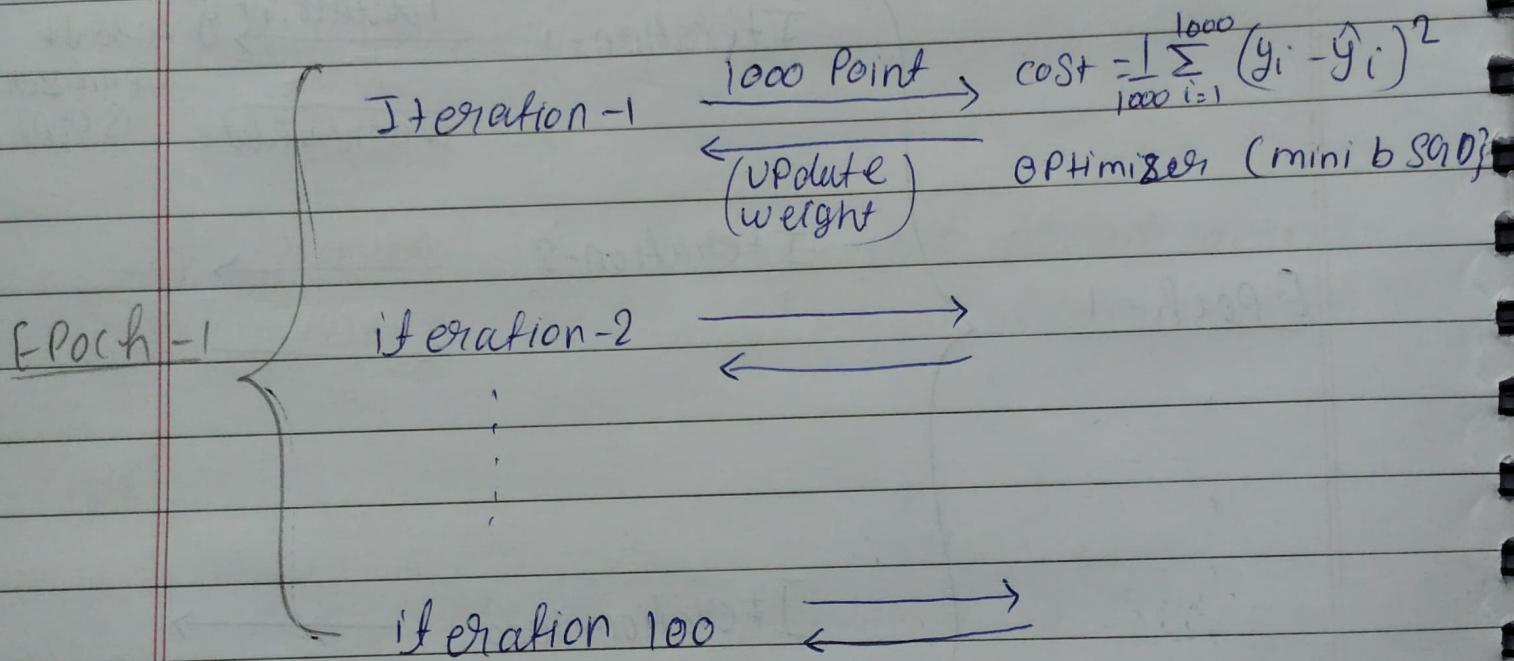
In mini batch SGD we introduce the Batch-size.

{Epoch, Iteration, Batch-size}

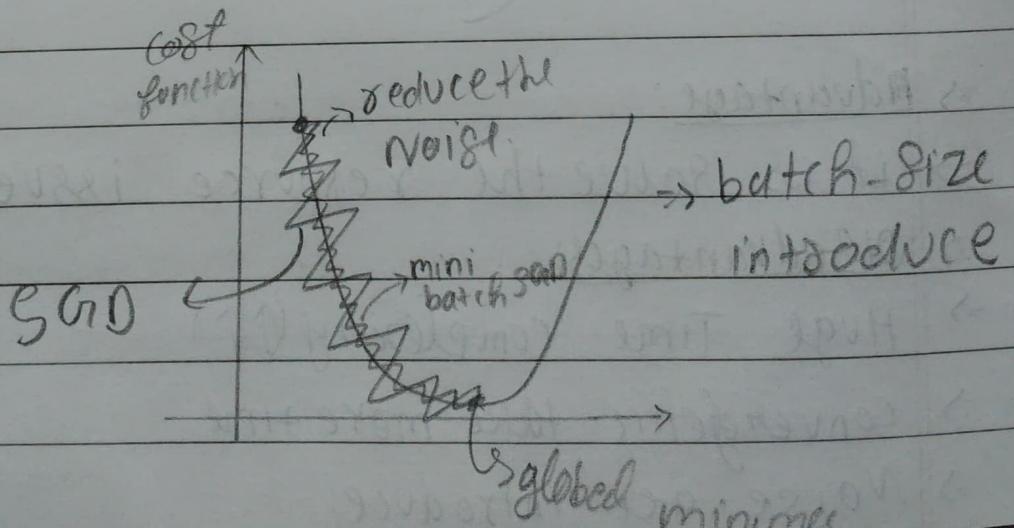
Data Points = 100K

⇒ batch-size = 1000 ,

$$\text{No of iteration} = \frac{100000}{1000} = 100 \text{ iteration}$$



①



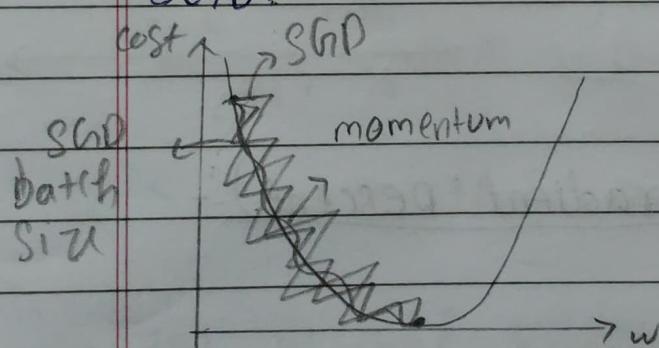
Advantage:

- ① Convergence speed will increase
- ② Noise will be less compared with SGD
- ③ Efficient resource usage (RAM, GPU)

Disadvantage:

- ④ Noise still exists.
- ⑤ SGD with momentum \Rightarrow

In this optimizer we ~~are~~ are going to smooth the noise of SGD.



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

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

update w.r.t time.

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

Exponential weight Average (Smoothing)

ARIMA, SARIMAX

Time series

time	t_1	t_2	t_3	t_4	- - -	t_n
value	a_1	a_2	a_3	a_4	- - -	a_n

$$\Rightarrow v_{t_1} = a$$

$$\Rightarrow v_{t_2} = \beta * v_{t_1} + (1-\beta) * a_2 \quad \left. \begin{array}{l} \{\beta = \text{control the} \\ \text{smoothing} \\ \text{function} \end{array} \right\}$$

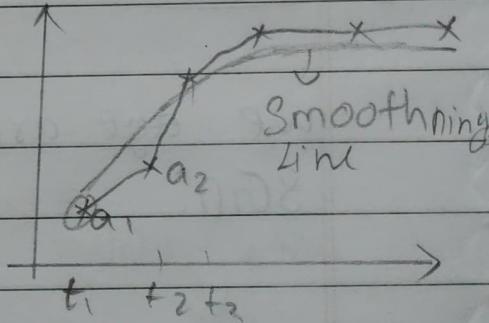
$$\Rightarrow v_{t_3} = \beta * v_{t_2} + (1-\beta) * a_3 \quad \left. \begin{array}{l} \{\beta \leq \beta \leq 1 \end{array} \right\}$$

Advantage:

- ① Reduce the noise
 - ② Quick convergence.
- ⑤ Adagrad: (Adaptive Gradient Descent):-

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

\Rightarrow We know that the learning rate is fixed
we are going to make it dynamic



Dynamic learning rate

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

$$\eta' = \eta$$

$$\sqrt{d_t + \epsilon} \rightarrow \text{very small value}$$

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

$$\xrightarrow{t=1}$$

$$\xrightarrow{t=2}$$

$$\xrightarrow{t=3}$$

$$\eta = 0.01$$

$$\eta = 0.003$$

$$\eta = 0.003$$

Disadvantage:-

- ① $\eta' \rightarrow$ Possibility to become a very small value ≈ 0

\Rightarrow so in this case $w_t \approx w_{t-1}$

the weight updation is not happen

⑥ Adadelta and RMSProp

$$\eta' = \frac{\eta}{\sqrt{Sd_{w_t} + \epsilon_0}}$$

We use the Exponential weight average to restrict the η' to reduce too much.

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

[Dynamic learning rate + Smoothening EWA]

$$w_t = w_{t-1} - \eta' \left(\frac{\partial L}{\partial w_{t-1}} \right)$$

↓ ↓
New old

⑦ Adam optimizer. $\rightarrow \{\text{Best Optimizer}\}$

\hookrightarrow SGD with momentum + RMSprop [Dynamic LR + Smoothening]

$$w_t = w_{t-1} - \eta' v_{dw} \quad \rightarrow \text{weight updation}$$

$$b_t = b_{t-1} - \eta v_{db} \quad \rightarrow \text{bias updation.}$$

$$\eta' = \frac{\eta}{\sqrt{S_{dw_t} + \epsilon_0}}$$

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

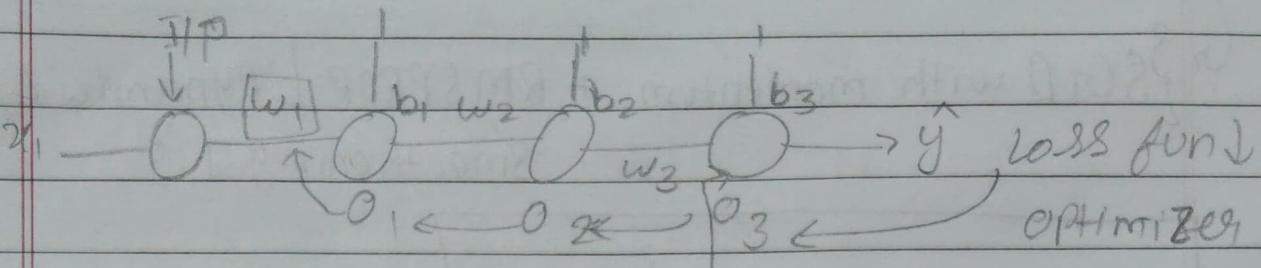
initially $\Rightarrow S_{dw_t} = 0$

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

$$V_{db_t} = \beta * V_{db_{t-1}} + (1-\beta) \left(\frac{\partial L}{\partial b_{t-1}} \right)$$

$\hookrightarrow \{\text{momentum Smoothening}\}$

Exploding Gradient problem



Updating the $w_1 \rightarrow Z = (O_2 * w_3 + b_3) \Rightarrow$

$$w_{1\text{new}} = w_{1\text{old}} + \eta \left(\frac{\partial L}{\partial w_{1\text{old}}} \right)$$

$$\frac{\partial L}{\partial w_{1\text{old}}} = \frac{\partial L}{\partial O_3} \times \frac{\partial O_3}{\partial O_2} + \frac{\partial O_2}{\partial O_1} \times \frac{\partial O_1}{\partial w_{1\text{old}}}$$

↳ this part

$$\frac{\partial O_3}{\partial Z} = \frac{\partial G(Z)}{\partial Z} \times \frac{\partial Z}{\partial O_2}$$

$$\Rightarrow [0 - 0.28] \times \frac{\partial L(O_2 \times w_3 + b_3)}{\partial O_2}$$

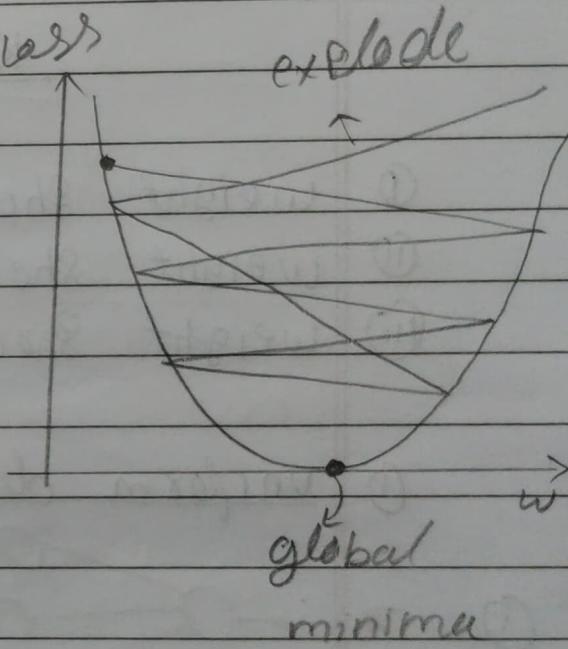
$$\left\{ \frac{\partial O_3}{\partial O_2} \rightarrow [0 - 0.28] \times w_3 \right\}$$

if weight initialization is very high value.

then the $\left(\frac{\partial L}{\partial w_{i, \text{old}}} \right)$ is very high it may be (+ve) or (-ve).

in this case

$$\left\{ \begin{array}{l} w_{i, \text{new}} \ggg w_{i, \text{old}} \\ \text{or} \\ w_{i, \text{new}} \lll w_{i, \text{old}} \end{array} \right\} \Rightarrow$$



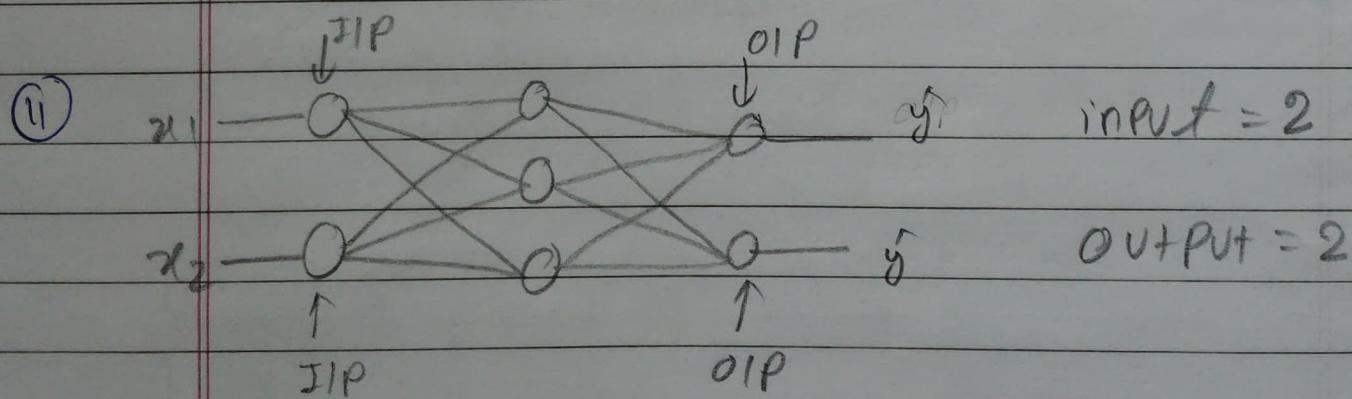
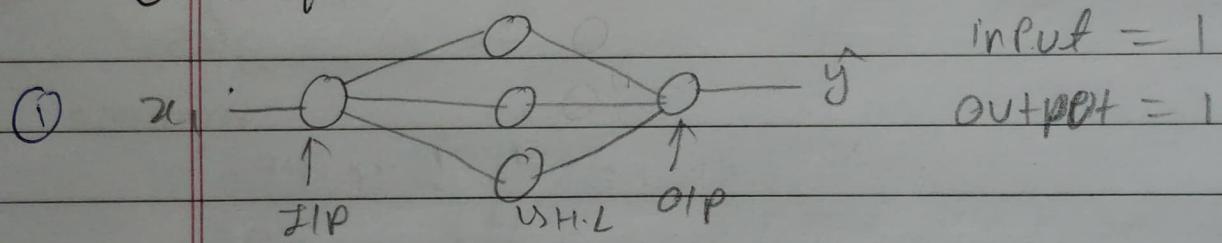
Weight Initializing Techniques.

- ① Uniform distribution
- ② Xavier / Glorot Initialization
- ③ Kaiming He Initialization.

Key Points

- ① weight should be small.
- ② weight should not be same.
- ③ weight should have good variance.

① Uniform distribution:-



[Neural Networks]

weights

$$w_{ij} \approx \text{Uniform distribution} \left[\frac{-1}{\sqrt{\text{input}}}, \frac{1}{\sqrt{\text{input}}} \right]$$

$$\text{in case of 2nd} \Rightarrow \left[\frac{-1}{\sqrt{2}}, \frac{1}{\sqrt{2}} \right]$$

② Xavier / Glorot Initialization:-

two types :-

- ① Xavier normal initialization.
- ② Xavier uniform initialization.

① Xavier Normal

$$w_{ij} \approx N(0, \sigma^2) \quad \text{mean} \rightarrow (\text{standard deviation})$$

$$\sigma = \sqrt{\frac{2}{(\text{input} + \text{output})}}$$

② Xavier Uniform:

$$w_{ij} \approx \text{Uniform distribution} \left[\frac{-\sqrt{6}}{\sqrt{(\text{input} + \text{output})}}, \frac{\sqrt{6}}{\sqrt{(\text{input} + \text{output})}} \right]$$

③ Haiming He Initialization:

① He Normal

$$\text{weights } w_{ij} \approx N(0, \sigma)$$

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

② He uniform

$$w_{ij} \text{ uniform } \left[-\sqrt{\frac{6}{\text{IIP}}}, \sqrt{\frac{6}{\text{IIP}}} \right]$$

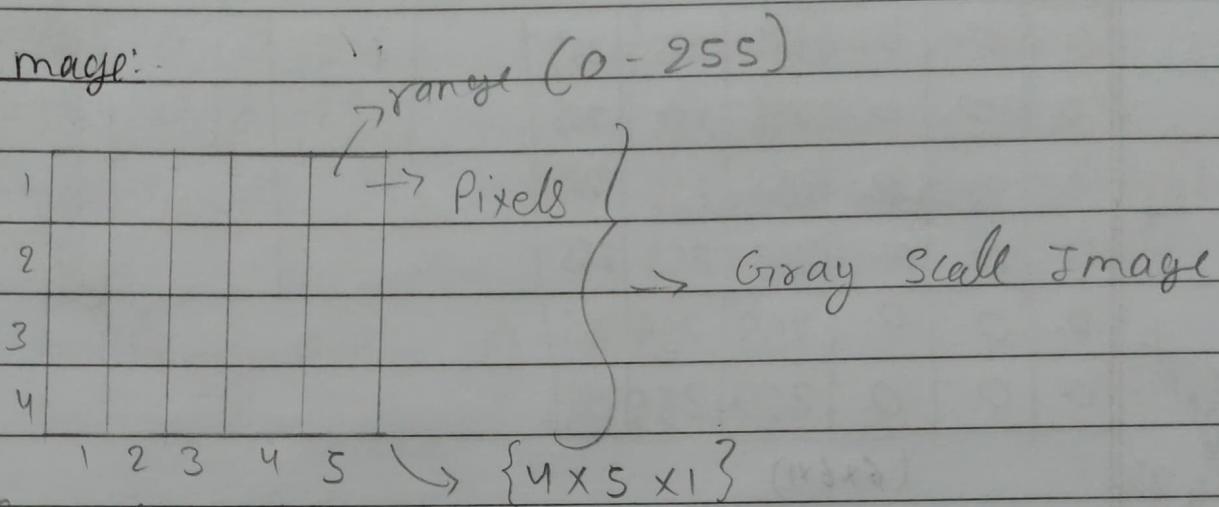
Dropout layer.

overfitting in neural network when they perform well on the data they have seen during training but fail to generalize to new data.

⇒ Dropout is a technique that randomly disable the neurons during the training iteration that prevent the dependent of learning pattern and then it perform well on new / unseen data.

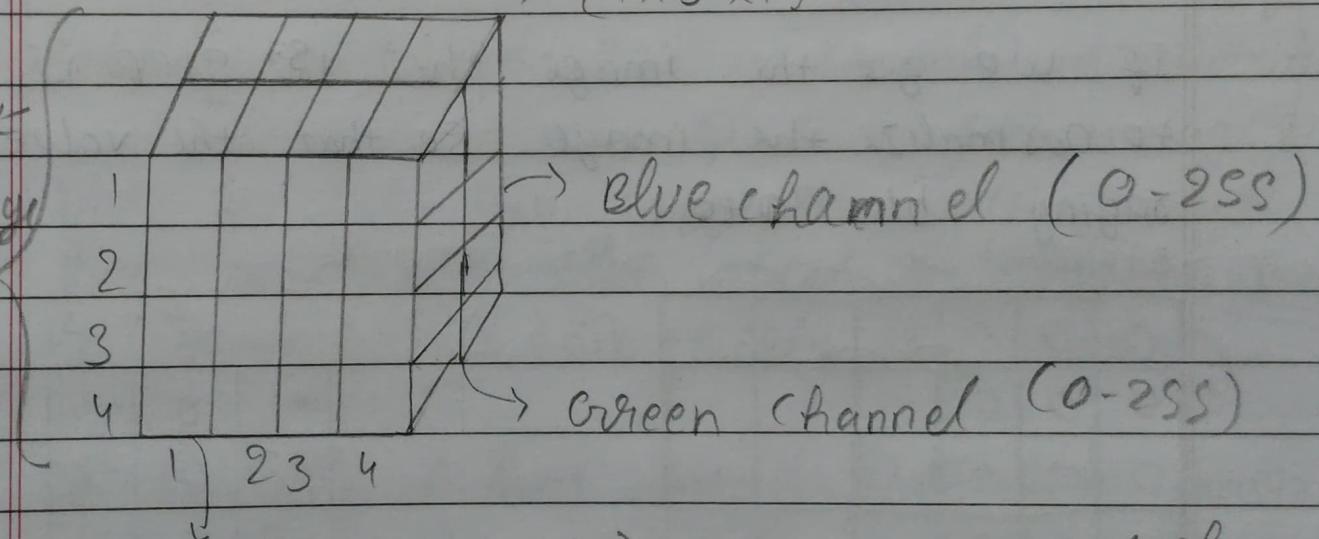
$\text{CNN} \rightarrow \text{Convolutional Neural Network}$.

Image:-



color⁴

Image



No of channels.

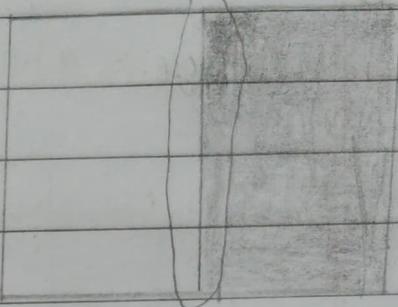
Red channel (0-255)

$$\text{three channel} = \underbrace{\{4 \times 4 \times 3\}}_{\text{No of Pixels}}$$

No of Pixels

Convolution operation In CNN:

0	0	0	255	255	255
0	0	0	255	255	255
0	0	0	255	255	255
0	0	0	255	255	255
0	0	0	255	255	255
0	0	0	255	255	255



(6x6x1)

If we get the image the 1st step is to normalize the image so that the value is ranging b/w 0 to 1.

slide 4 $s + \sigma d = 1$ \rightarrow (3+3) move

0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1

$n=6$

$d = 3$

(3+3)
filters.

(6x6x1)

0	-4	-4	0
0	-4	-4	0
0	-4	-4	0
0	-4	-4	0

OP = 4

\Rightarrow denormalize

4x4

2SS	0	0	2SS
2SS	0	0	2SS
2SS	0	0	2SS
2SS	0	0	2SS

 \Rightarrow

edge detected \rightarrow vertical.

$$\left\{ \text{formula} = O/P = n - f + 1 \right\}$$

Padding in CNN \Rightarrow So in convolution operation we give a 6×6 image size but we get 4×4 image size O/P. So we lost some information about the original image. To prevent this we apply the padding.

So basically Padding means here we construct or add the lines around the 6×6 matrix and it will become the 8×8 matrix and after apply the filter of 3×3 matrix we get the 6×6 in the out-put.

So formula for finding the Padding. is

$$\left\{ O/P \text{ size required} = \text{original size} - \text{filter size} + 2P + 1 \right\}$$

$$n=6, \quad f=3 \quad \Rightarrow \quad \{n-f+2P+1 = 0 \mid P\}$$

$$6 - 3 + 2P + 1 = 6$$

$$3 + 2P = 5$$

$$2P = 2$$

$$P = 212$$

$$P = 1$$

so we use only 1 padding on the original given image.

paddany
layer

0	0	0	1	1	1				
0	0	0	1	1	1	+1	0	-1	
0	0	0	1	1	1	X	+2	0	-2
0	0	0	1	1	1		+1	0	-1
0	0	0	1	1	1				
0	0	0	1	1	1				
0	0	0	1	1	1				
0	0	0	1	1	1				
0	0	0	1	1	1				

$$6 \times 6 \rightarrow 8 \times 8$$

This is called
the convolutional
layer.

0	-4	-4	0	
0	-4	-4	0	
0	-4	-4	0	
0	-4	-4	0	91

$$4 \times 4 \rightarrow 6 \times 6$$

Q. If you have 7×7 image size and you apply 3×3 filter on the original image and if you have require 7×7 output. Then find How much Padding you require.

$$\begin{array}{c|c|c|c} 7 \times 7 & \times & 3 \times 3 & = 7 \times 7 \\ n=7 & & f=3 & O/P=7 \end{array}$$

Padding require \Rightarrow

$$\Rightarrow n - f + 2P + 1 = O/P$$

$$\Rightarrow 7 - 3 + 2P + 1 = 7$$

$$\Rightarrow 4 + 2P + 1 = 7$$

$$\Rightarrow 2P = 7 - 5$$

$$\Rightarrow P = \frac{2}{2} = 1$$

\Rightarrow So padding is two types to initialize.

① Zero Padding \Rightarrow fill all the value = 0

② Neighbouring Padding \Rightarrow fill the value of as the value of Neighbouring Point.

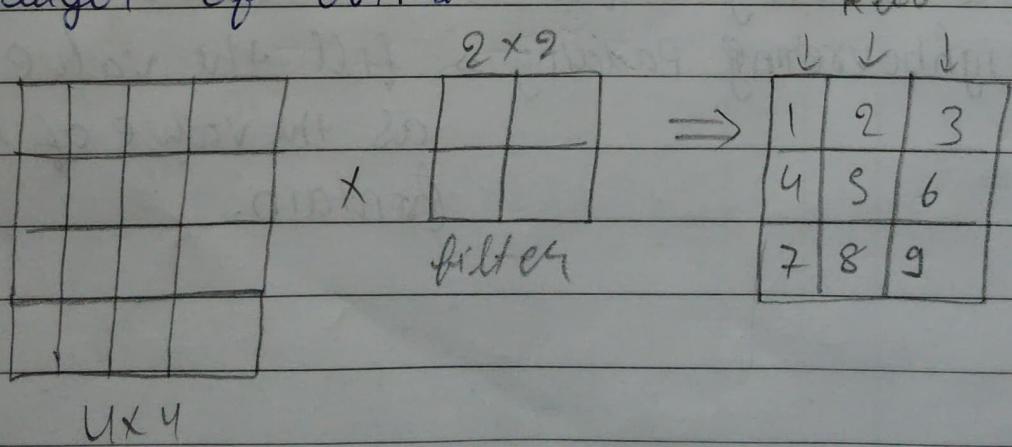
#⇒ In CNN the filters are going to be updated so in the output matrix we use the ReLU activation function for the back propagation.

The need of ReLU or it's variant:-

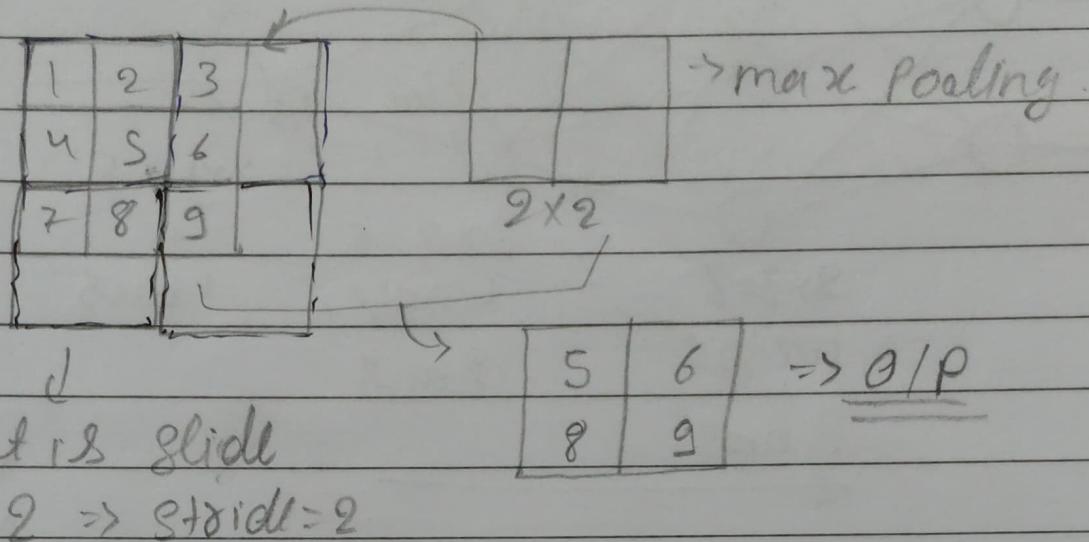
We use the ReLU for finding the derivative and using that derivative with applying the learning rate, we are going to update the values of filters in case of CNN for backward propagation.

Max Pooling, Min Pooling and Mean Pooling.

So basically we use the max pooling to extract the imp features of the given image and it is apply on the convolutional layer of output.



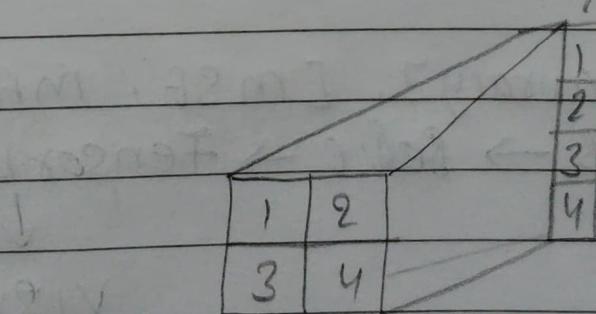
APPLY THE MAX POOLING ON O/P.



- ⇒ In max pooling we extract the max value
- ⇒ In min pooling we extract the min value
- ⇒ In mean pooling we extract the mean value.

Flatten layer:

flatten layer:

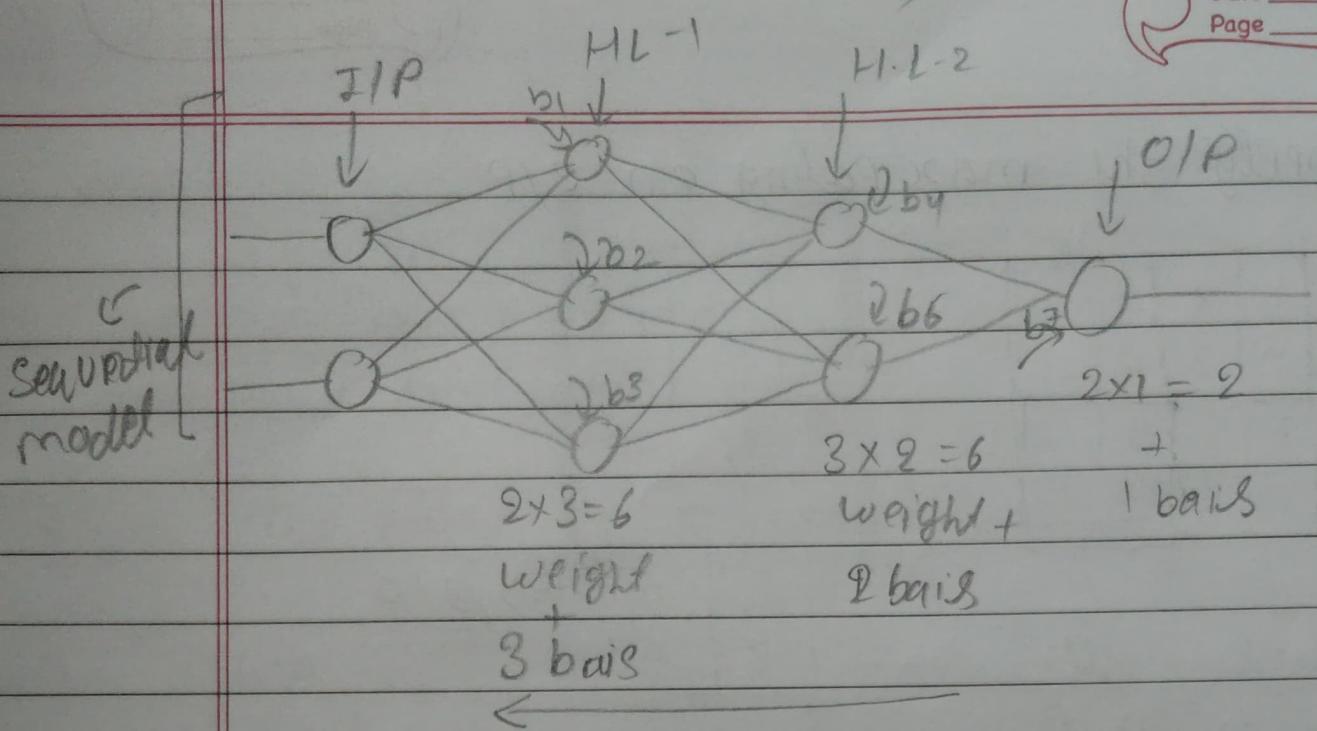


Training of ANN

classmate

Date _____

Page _____



$$\text{total trainable parameter} = (6+3) + (6+2) + (2+1) = \underline{\underline{20}}$$

Steps

- (I) created a sequential network
- (II) Dense = 64 \Rightarrow $\begin{bmatrix} 0 \\ 1 \\ 2 \\ 3 \end{bmatrix} \Rightarrow 64 \text{ Neurons}$
- (III) Apply the activation function on each node.
- (IV) Optimizer
- (V) Loss function.
- (VI) matrix \rightarrow [Accuracy], [mSE, mAE].
- (VII) Training \rightarrow Logs \rightarrow save \rightarrow Tensorboard \downarrow Visualization.