

# Deep Learning with Deep NLP

## Artificial Intelligence

### Machine Learning

### Deep Learning

Today we are diving into the basics of deep learning.

- Q - what is deep learning? why is there so much hype about deep learning. what is the history of deep learning. How did it evolve and why company spend so much money on deep learning.

Deep learning  $\rightarrow$  which is a subset of machine learning which itself is a part of much greater Artificial intelligence.

Deep learning has algorithms that are inspired by the structure and the function of the brain.



# Deep Learning with Deep NLP

## Artificial Intelligence

### Machine Learning

### Deep Learning

Today we are die into the basics of deep learning

- Q- what is deep learning? why is ~~some~~ much hype about deep learning. what is the history of deep learning. How did it evolve and why company spend so much money on ~~deep~~ deep learning.

Deep learning → which is a subset of machine learning which itself is a part of much greater Artificial intelligence.

Deep learnings has algorithms that are inspired by the structure and the function of the brain.

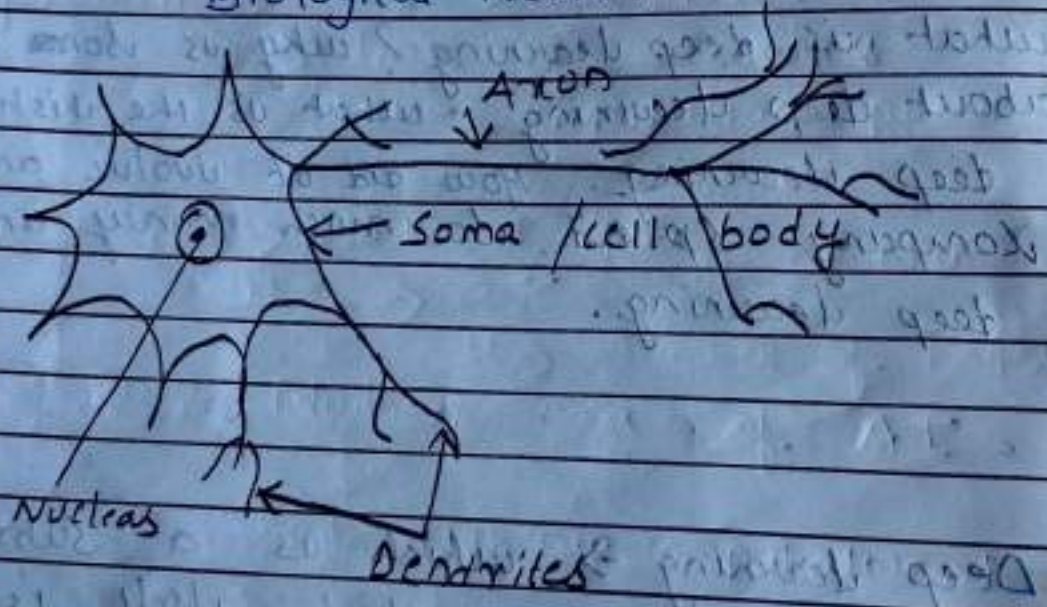


# History of Deep Learning

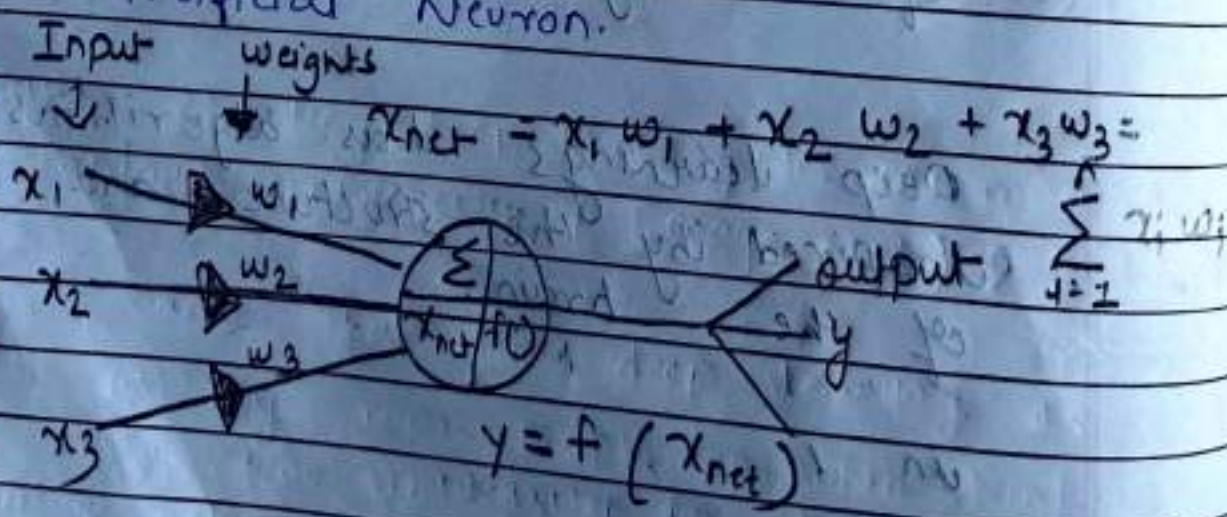
→ it started in 1958 where we get first idea of perceptron, the perceptron is inspired by the biological neuron.

neuron is a basic building block of the brain.

## Biological Neuron



## Artificial Neuron





$$x_{net} = x_1 w_1 + x_2 w_2 + x_3 w_3 = \sum_{i=1}^n x_i w_i$$

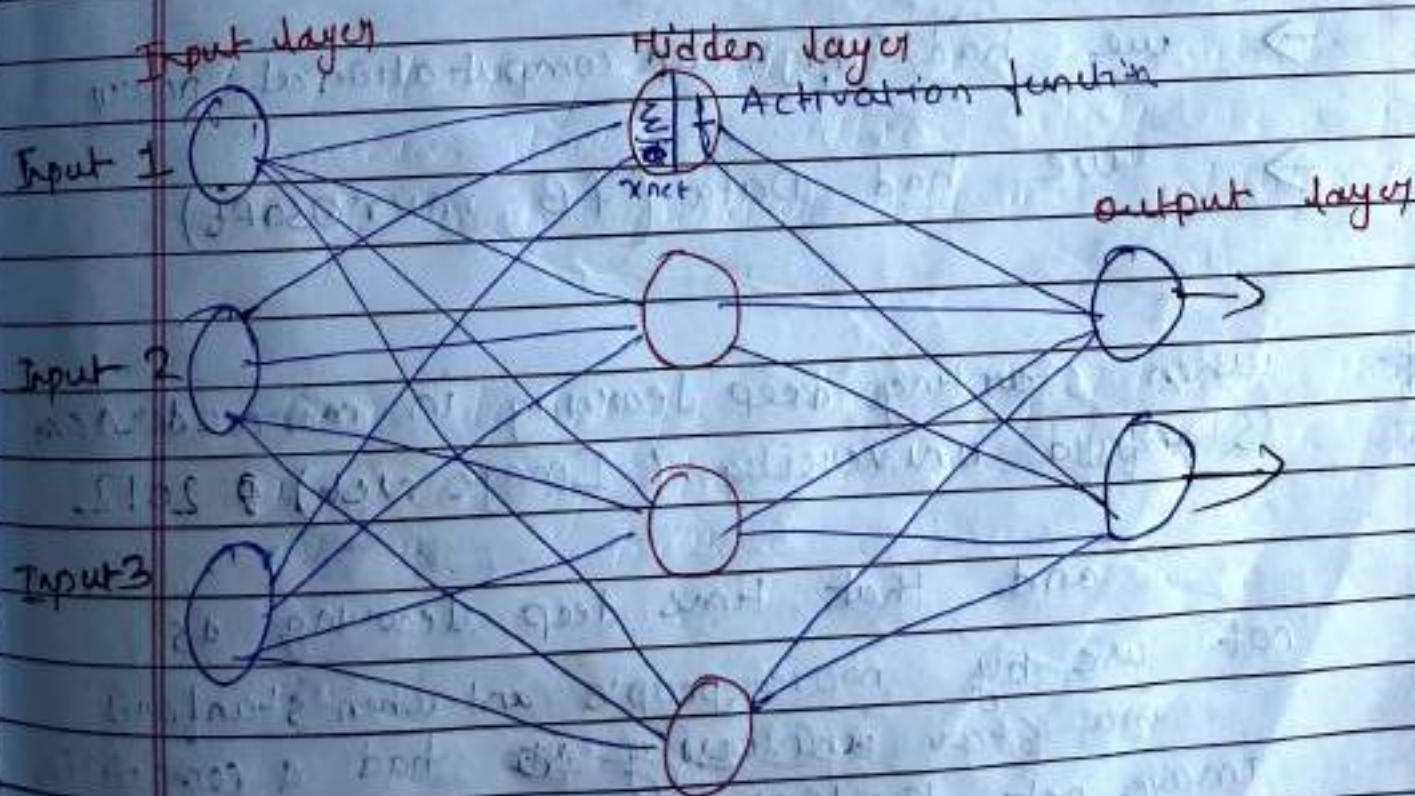
$f(x) \rightarrow$  passing the function. and this function is known as Activation function

\* ~~pro~~ problems

$\rightarrow$  A single perceptrons/node

$\rightarrow$  No learning

2  $\rightarrow$  MLP (Multi layer perceptron) 1980





<Hidden layer> → That ~~mean~~ means all the layer that is between the input layer and the output layer are known as hidden layer.

these ~~are~~ layers are neither directly ~~connecte~~ connected to the input or are they directly connected to the output.

problem → This only proved ~~deep~~ theory:

2 → lack of resources

3 → less computational power

4 → lack of data

→ 2012 Modern era of AI

→ we had high computational power

→ we had Data (FB, Microsoft)

\* when starting deep learning in ~~modern~~ modern era  
 \* Stanford university (Image net) 2012

and that time deep learning is not use by more people and when Stanford ~~university~~ ~~had~~ had a competition image net ~~deep~~ and since ~~deep~~ deep



Learning is not popular. People used to stick with ~~machine~~ machine learning algo like Random Forest, SVM. But a team use basic deep learning algo to solve this image net problem.

image net  $\rightarrow$  To identify objects from images  
(very difficult)

Deep learning beats all the ML algo.

\* Now Days

$\rightarrow$  Now we have high resources

$\rightarrow$  Lots of data

$\rightarrow$  Because of globalization we keep getting new algorithms to solve new problems.

CV  $\rightarrow$  CNNs

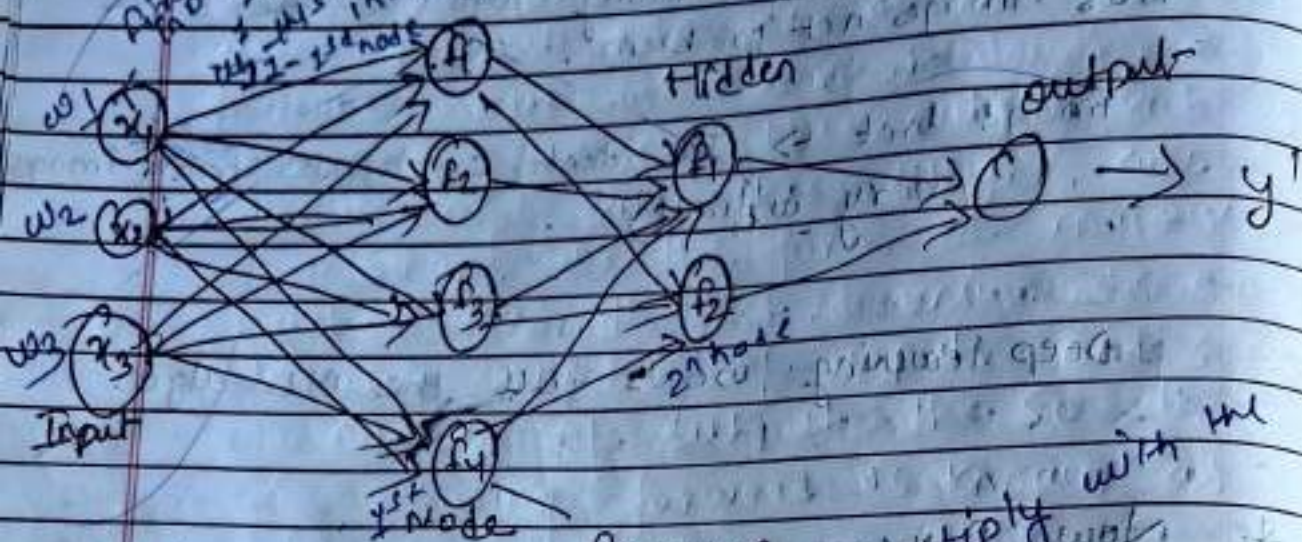
80s/90s

1) History of deep learning and we read 2 different concepts perceptrons and multilayer perceptron or (MLP). So let us quickly recap. we had cell body here, we inputs going in we performed first the summation for which is followed by an activation



# Optimisers & Gradient Descent

## Basic structure of neuron and MLP



> So, first inputs you will be multiply with the weights.

$w_{34} \rightarrow 3 \rightarrow$  stands of Number of inputs

$4 \rightarrow$  Number of nodes

$$Z_1 = w_{11}x_1 + w_{21}x_2 + w_{31}x_3 + \text{bias}$$

$$h_1 = f(Z_1)$$

$$h_2 = f(Z_2)$$

$$h_3 = f(Z_3)$$

$$h_4 = f(Z_4)$$

$$h_5 = f(Z_5)$$

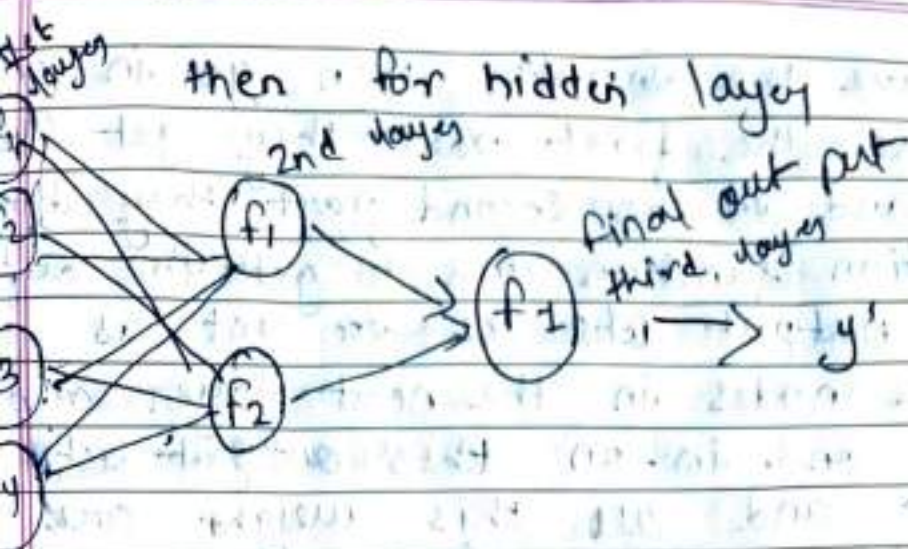
$$h_6 = f(Z_6)$$

$$h_7 = f(Z_7)$$

$$h_8 = f(Z_8)$$



$w$  = weight  
 $h$  = hidden



$$z_1^2 = w_{11}^2 h_1^1 + w_{21}^2 h_2^1 + w_{31}^2 h_3^1 + w_{41}^2 h_4^1 + \text{bias}$$

$$\rightarrow h_1^2 = f(z_1^2)$$

26 This whole process going forward, you have the inputs go and make a go forward first to the input layer then hidden layer then on and on they will reach the output layer that is predicted by the model ( $y'$ )

→ That moving forward feed (very important)

~~that's~~ ~~clears~~

⇒ this whole process is going forward you have

⇒ so, since in one way it is moving forward it's called forward feed, we



we have the inputs, then go in one by one. the first part they get summed up and in the second part they have activation function and to get the outputs this output ~~also~~ ~~also~~ act as input for the nodes in the next layer and so on. in on this output ~~act as~~ ~~input~~ and all this weight are shape in the form of matrix where, the first dimension is equal to the number of inputs and second number of output. So that moving forward is known as forward feed.

So, the next part we learning is ~~more~~ backword ~~in~~ ~~that~~ ~~learning~~ ~~op~~

\* Optimisers - Gradient Descent

So, in the last class we learned about forward feed

forward feed  $\rightarrow$  you have the input and each input get multiply with weights to add weightage to the input and get passed through a multilayer perceptron to given a output. in a multilayer perceptron we have nodes so the input



get multiply with the weights and each node they get all summed up together and pass through an function that is an activation function. This gives us an output which acts as a input for the nodes in some other layers and this goes on and on until you get the predicted value.

Since all the weight and biases are randomly initialize there is no way that the predicted value is going to be close to the actual value.

→ we can ~~only~~ ~~change~~ have inputs and we can't change the input we can only change the weight and bias.

→ weight, bias need to be change so that predicted value very close the actual value.

~~Predicted~~  $\hat{y} \sim y \rightarrow$  Actual value

→ weight and biases are known as trainable or learnable parameters.

\* **LOSS** → loss basically means how <sup>diff</sup> ~~much~~ <sup>does</sup> ~~value~~ predicted value differ from the actual value.

we heard about ~~1~~ ~~2~~  $L1$  loss,  $L2$  loss.



Q → we need to change weights and biases so that we have very low loss.

For that we use optimiser gradient

we have to optimise the weights and biases so that we have a very low loss minimal loss.

Q → Let's talk about the basic optimiser Gradient Descent.



So, we slowly reach the minimum loss. You have to slowly descent downwards.

$$w_{\text{new}} = w_{\text{old}} - \frac{dL}{dw_{\text{old}}}$$

$$\frac{dL}{dw_{\text{old}}}$$

$$\frac{dy}{dx}$$

$$y = \text{loss}$$

$$x = \text{old weight}$$



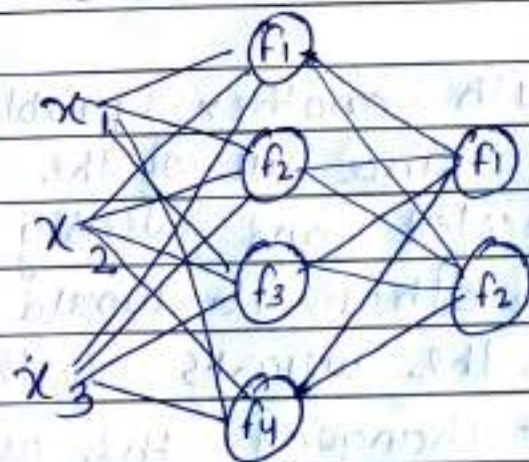
you have a chance of going 1 step at a time you have to chance to go 0.5 at a time. So, this magnitude is known as the learning rate. This is denoted -  $\eta$  - eta.

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{dL}{dw_{\text{old}}}$$

we can update the loss and differentiate the loss with respect to the weights in 1st node. But <sup>we need</sup> some how ~~we have~~ find the way to do that. then we <sup>came</sup> across and <sup>that is</sup> known as

\* Back propagation. This is a lot like chain rule.

\* ~~Process~~



in the last few ~~pages~~ <sup>clases</sup> we have learned about forward ~~feed~~ propagation. So you have the inputs they get multiply with the weights. The ~~and~~ get all summed together which is followed by an Activation function and



this on an on until you get a predicted value. No this value can not be equal to the actual value initially. So that's why you have the idea of loss which shows the how much different the predicted value is the actual value. Now you have the loss calculated so you do some how change the weight and biases to minimize the loss. to do this we use optimisers. to optimise the weight and function we talked about one optimiser Gradient descent.

formulated

$$w_{\text{new}} = w_{\text{old}} - \frac{dL}{dw_{\text{old}}}$$

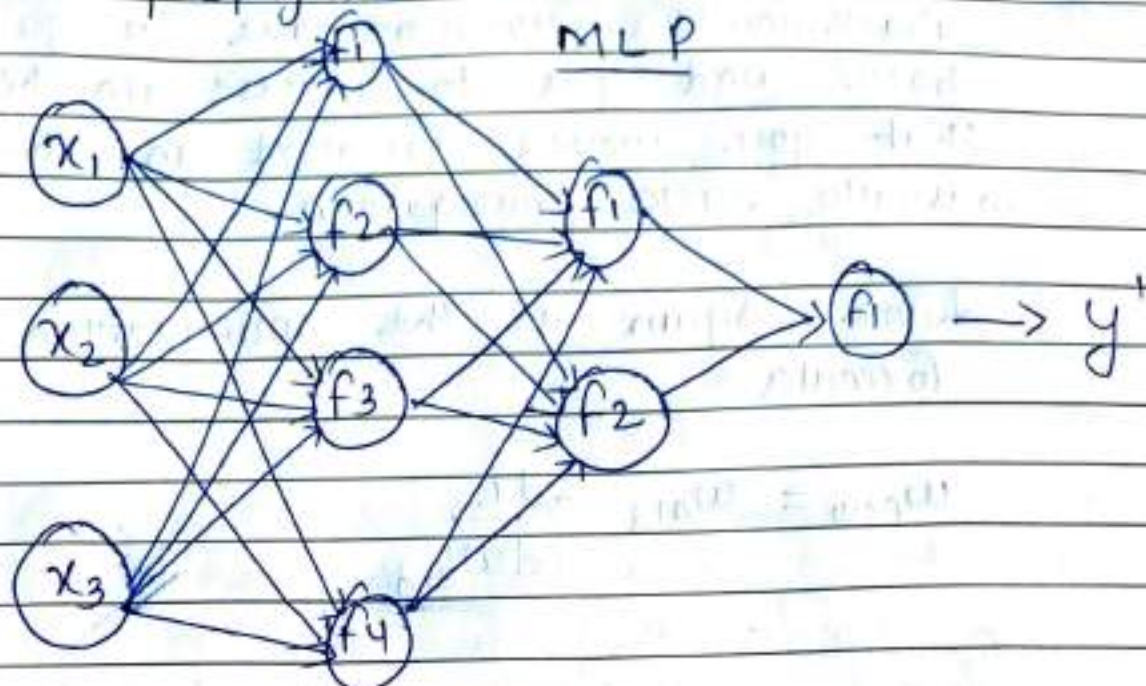
we came across with another problem that was happening. Since all of the nodes are interconnected and densely connected with each other. It would be impossible to update the weights of the inner layer without changing the weights of the outer layer. So this was a problem. Hence we came across what was known as the Back propagation.

Here we follow chain rule of differentiation.



Back propagation

MLP



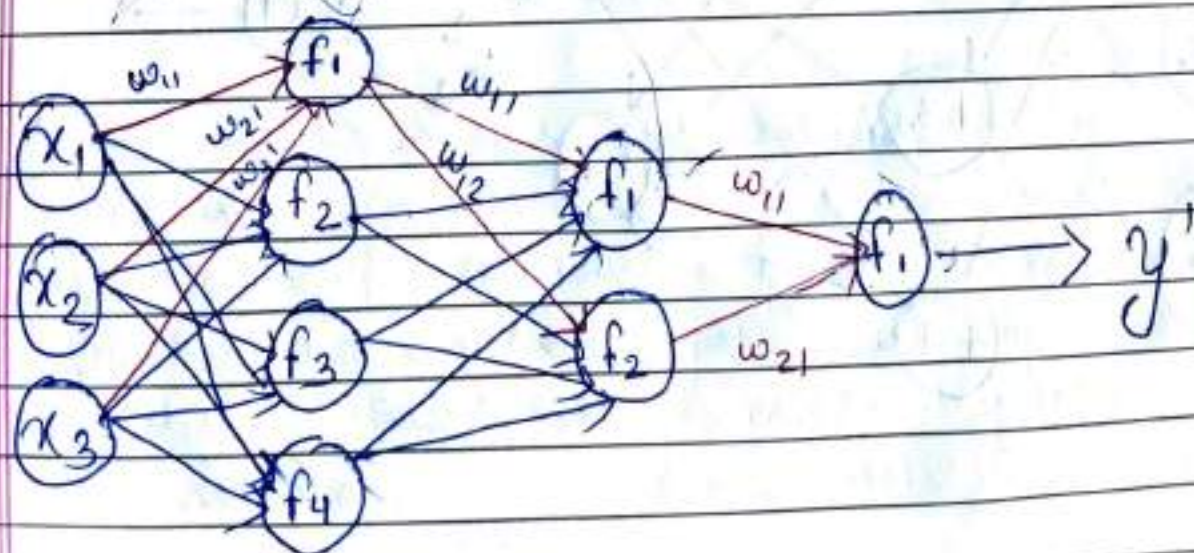
So, we had this output

$$w_{ij}^k$$

$k \rightarrow$  represent which layer

$i \rightarrow$  from which node

$j \rightarrow$  into which node





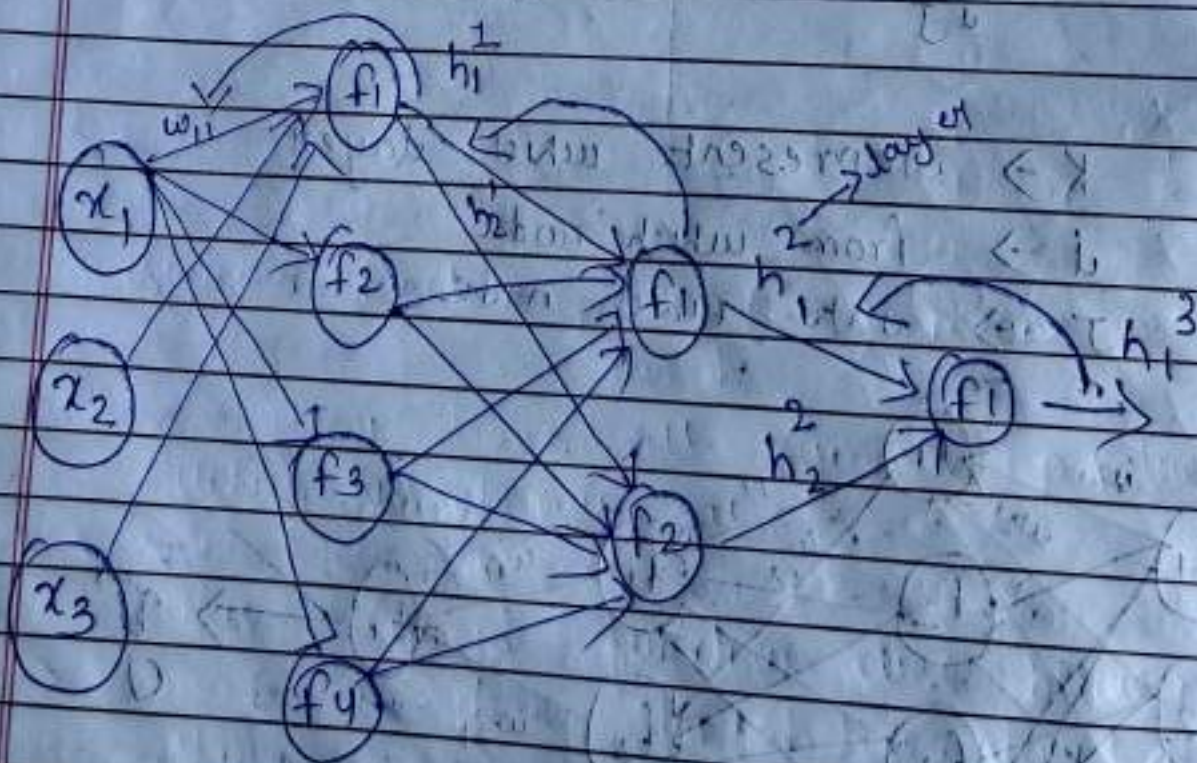
So as the name suggest we have to go backwards you ~~also~~ have to go back to back and you to traverse in the path that you moved forward in. that is basically back propagation.

Look Again in the optimization formula

$$w_{\text{new}} = w_{\text{old}} - \frac{dL}{dw_{\text{old}}}$$

$$w_{\text{new}} = w_{\text{old}} - \frac{dL}{dh_1^3} \times \frac{dh_1^3}{dw_{\text{old}}}$$

$$w_{\text{new}} = w_{\text{old}} \rightarrow w_{11}^3$$





$$\frac{dL}{dw_{11}^2} = \frac{dh_1^3}{dh_1^3} \times \frac{dh_1^2}{dh_1^2} \times \frac{dh_1^2}{dh_1^2} \times \frac{dh_1^1}{dw_{11}^1}$$

$$\frac{dL}{dh_1^3} \times \frac{dh_1^3}{dh_2^2} \times \frac{dh_2^2}{dh_2^1} \times \frac{dh_2^1}{dw_{11}^1}$$

So, in short we have forward feed where you go forward in the layers one by one and you have backward propagation where you go back find the differentiation one by one and update the weights and bias. So this whole (forward feed + Backward propagation) together for 1 epoch.

And many of this epoch are done one after the other. So that we reach the final ~~inter~~ model where you have minimum loss.

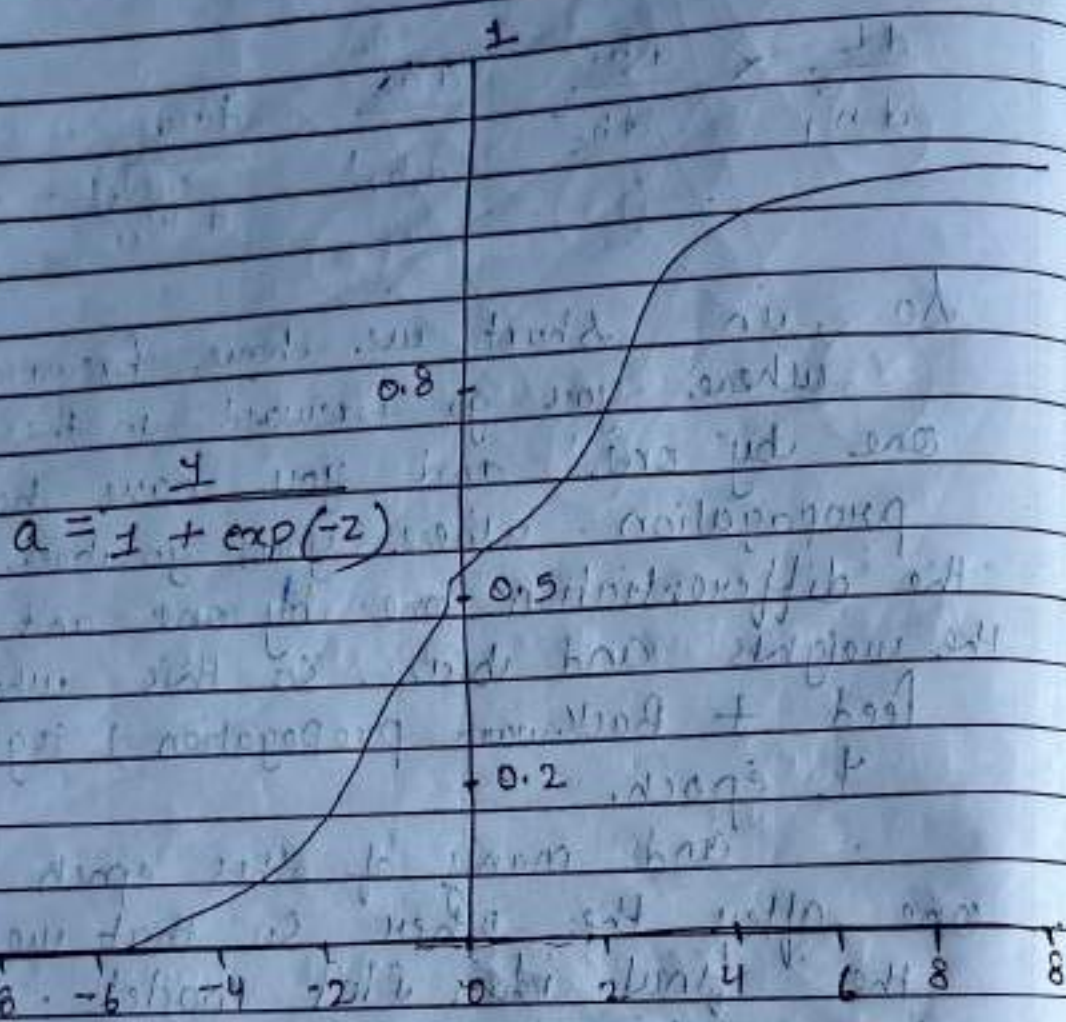
### \* Activation function in depth

Activations function  $\rightarrow$  you had this neuron here, you had all inputs going in. You first had a summation where the input multiplied with different weights and after that this output here which was known as  $Z$  was passed through a function that is an activation function to give the final output.



# (1) First Activation function

## (1) Sigmoid function



Sigmoid function,  $\exp(-z)$  means  $e^{-z}$

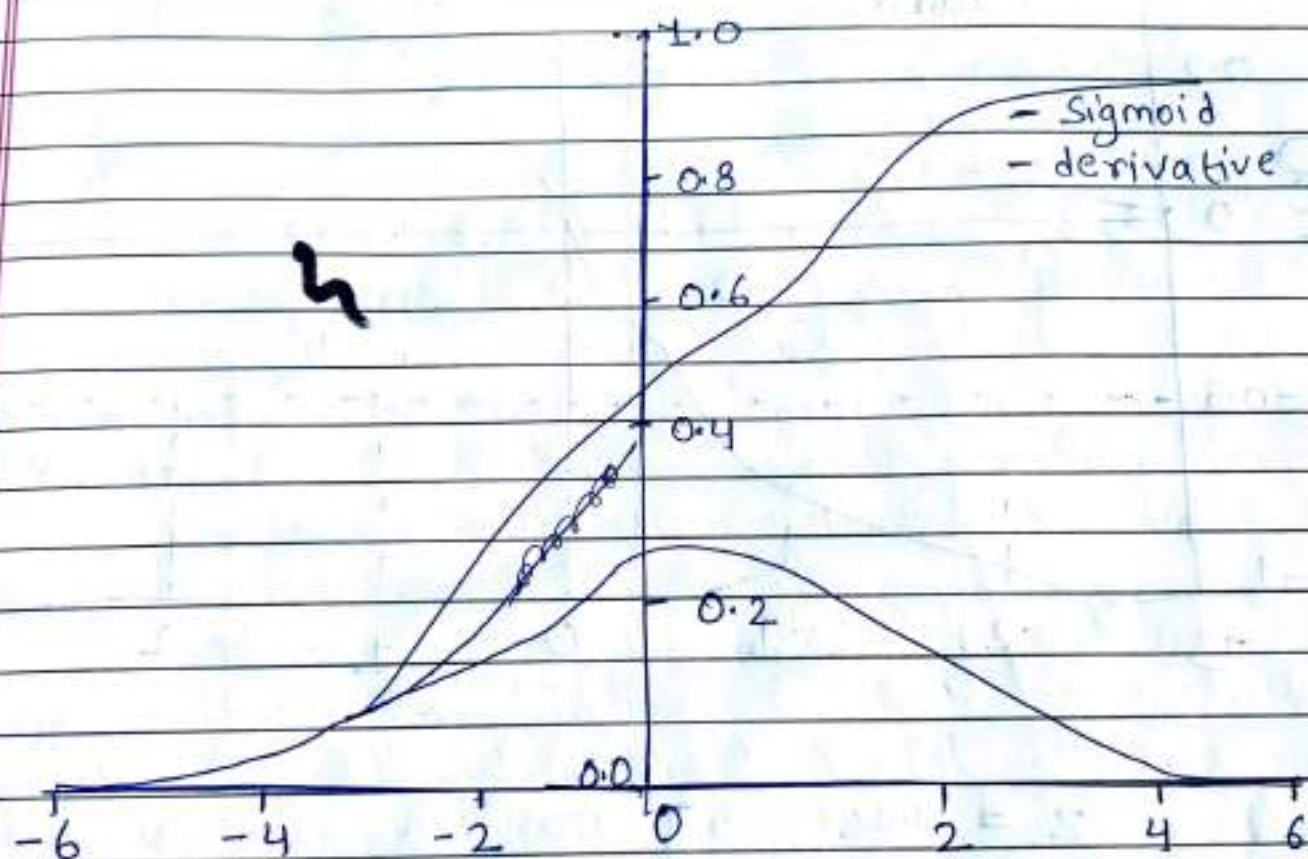
$z$  = you have all summed up value which is known as  $z$ .

This Sigmoid function lies between 0 to 1. It is a very good thing because this can be used in classifiers then to give a summed value. Some values



and this values we can take as prediction of or the probabilities of how this happening or not happening

### \* Derivative of Sigmoid Function



### \* why we need Derivative

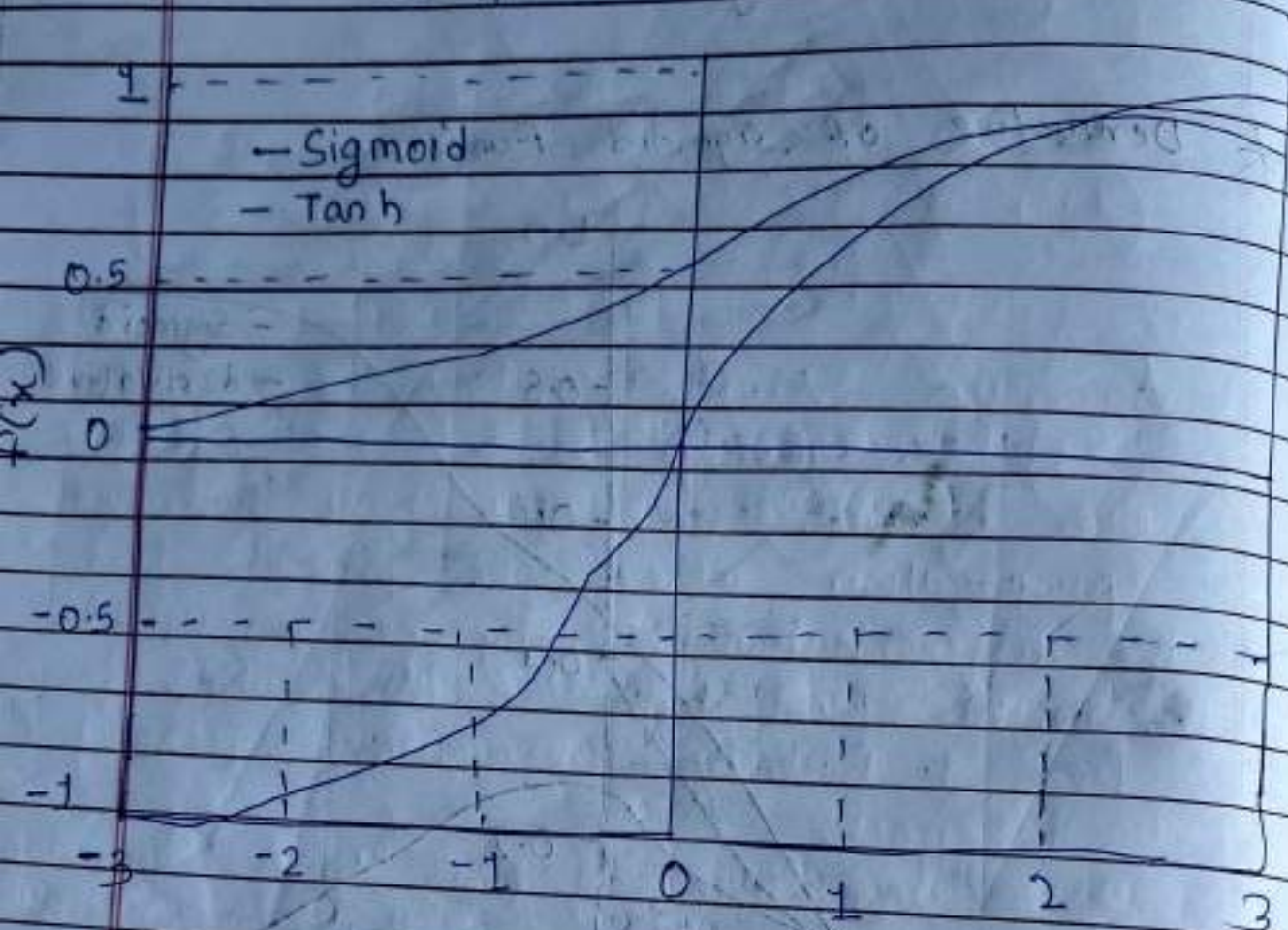
If we had optimiser and in optimiser and in optimiser ~~are~~ ~~are~~ we use to calculate the loss with respect to the weights

$$\frac{dL}{dw} \rightarrow \frac{dL}{dh} \times \frac{dh}{dw}$$



→ Next function

Tan h



z = -ve

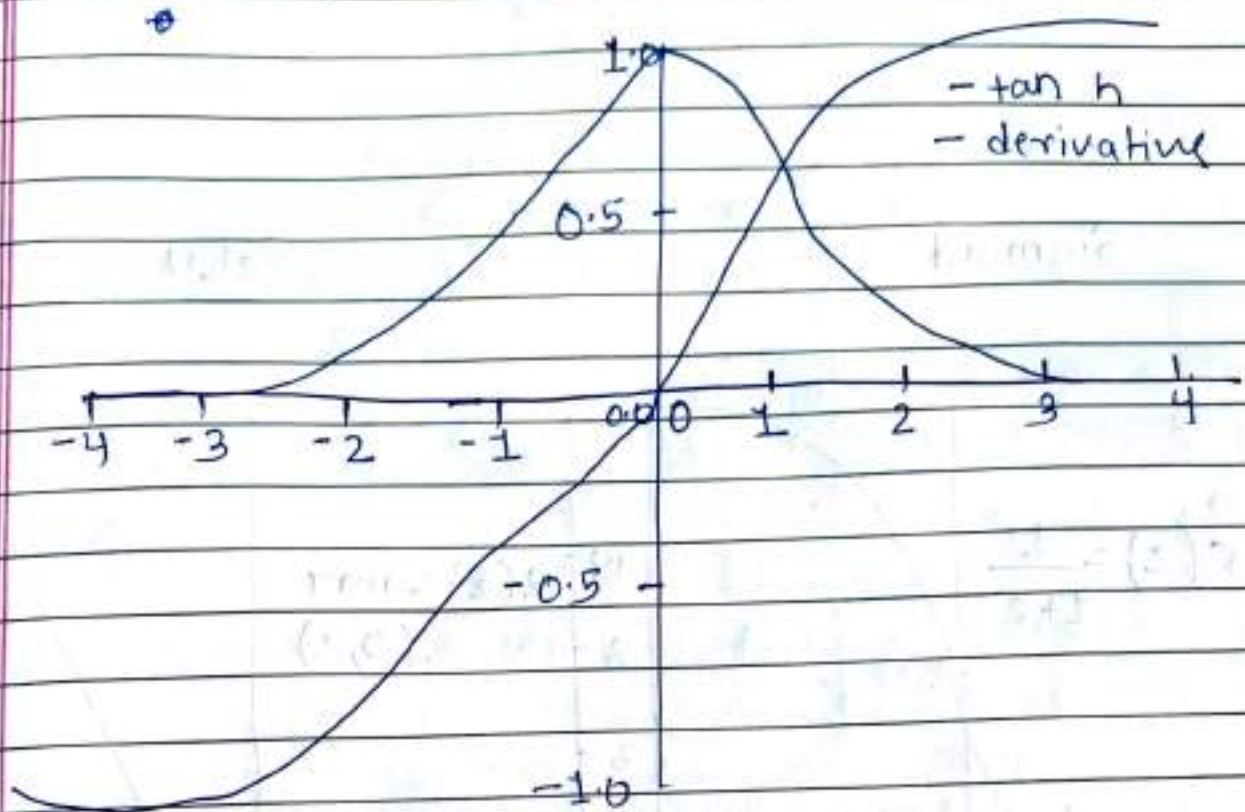
Tanh

equation of Tanh  $\rightarrow \frac{e^z - e^{-z}}{e^z + e^{-z}}$

→ Tanh function lies between  $[-1, 1]$



## - Tanh Derivative



$$\text{derivative} = 1 - a^2$$

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

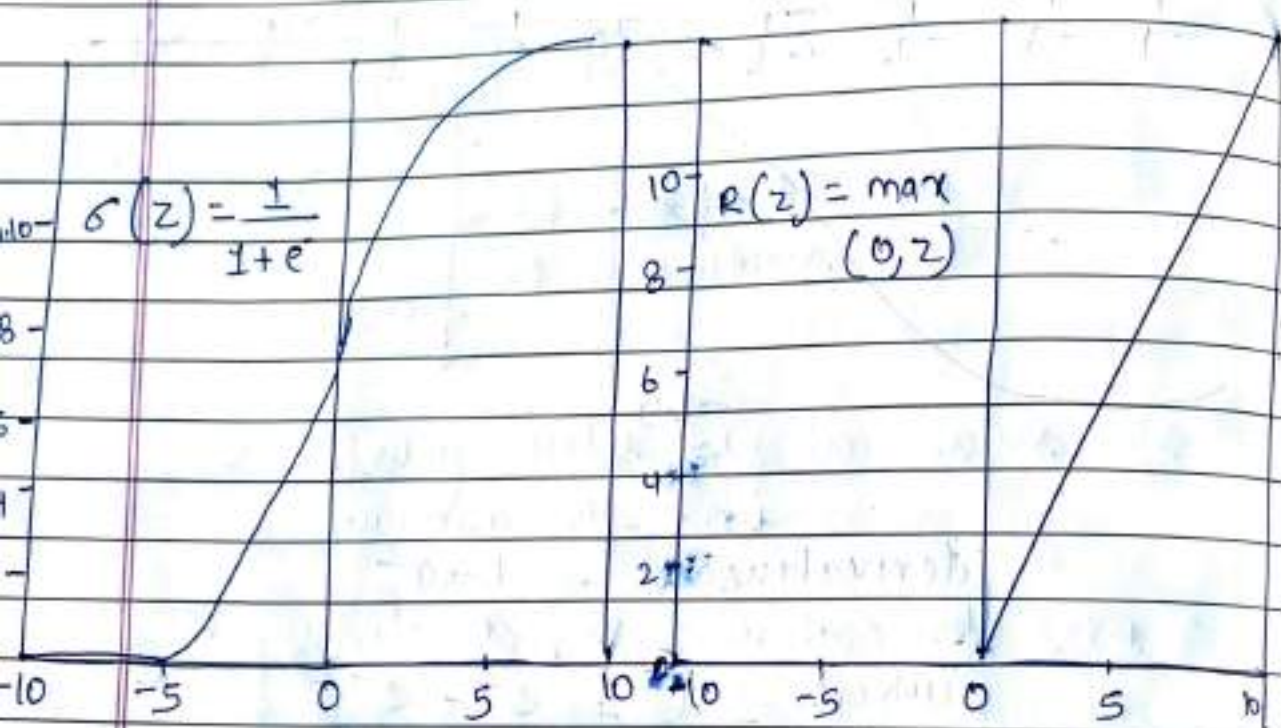
~~Redo~~



# \* Rectified Linear units (ReLU)

Sigmoid

ReLU



- ~~der~~ for any positive number derivative will be one

- for any negative numbers derivative will be zero

\* problems with Relu

→ Dying Relu

→ Dying neurons



~~dl~~  
~~dw~~

→ Suppose for any case you have  $z$  is less than zero, so in that case equally forms a zero. It comes on equal the value of zero.

now, whatever you multiply with zero, then whole things become zero.

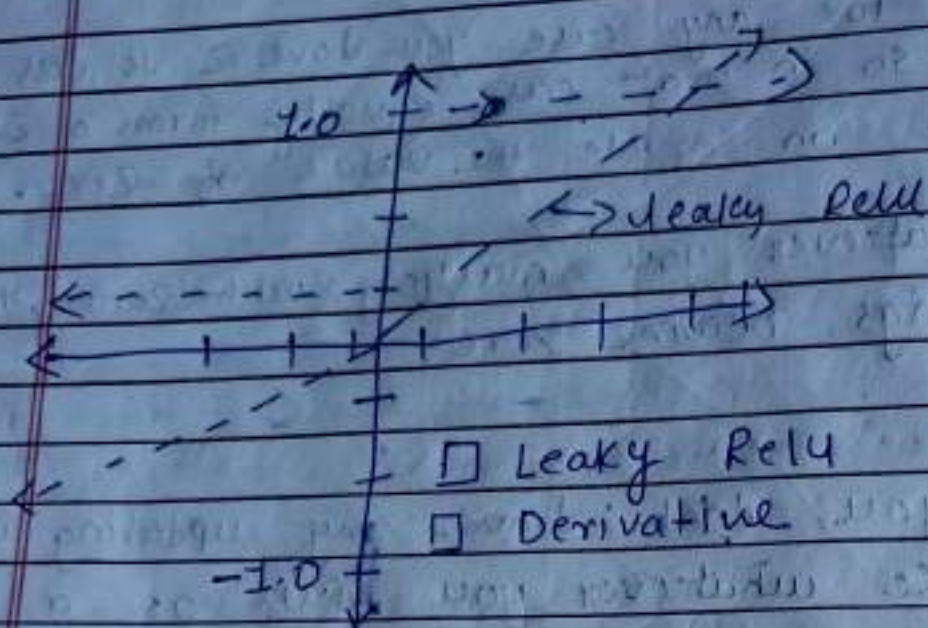
→ So you don't have any updating weight here. So whatever you gives as a random weight remains same weight even after a lot of calculations. So this problem is known as dying ReLU or Dying Neuron problems.

Q → what was the problem with sigmoid and tanh?

Ans → well it was lots of calculation and sigmoid have values lies between 0 and 1 while tanh had it -1 to 1.



## \* Next Function Leaky ReLU



→ Leaky ReLU try to solve the function of ReLU

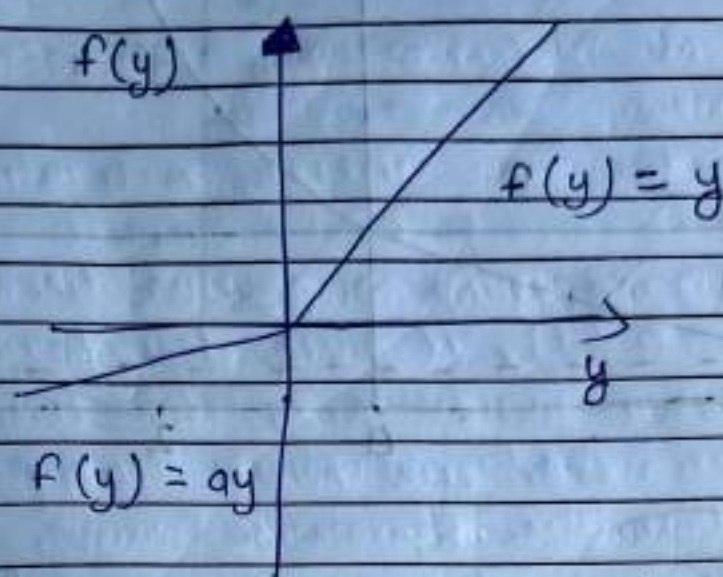
→ what is the function of Leaky ReLU

example:  $\max(0.01z, z)$

→ it trying to eliminate the problem of dying ReLU and dying Neuron using ~~Leaky~~  $\max(0.01z, z)$  function



\* Next parametrized Relu.



Here the function is

$$\max(\alpha z, z) \quad \alpha - \text{Learning rate}$$



alpha can any number  
that has been trained

by the model  
model has to decide

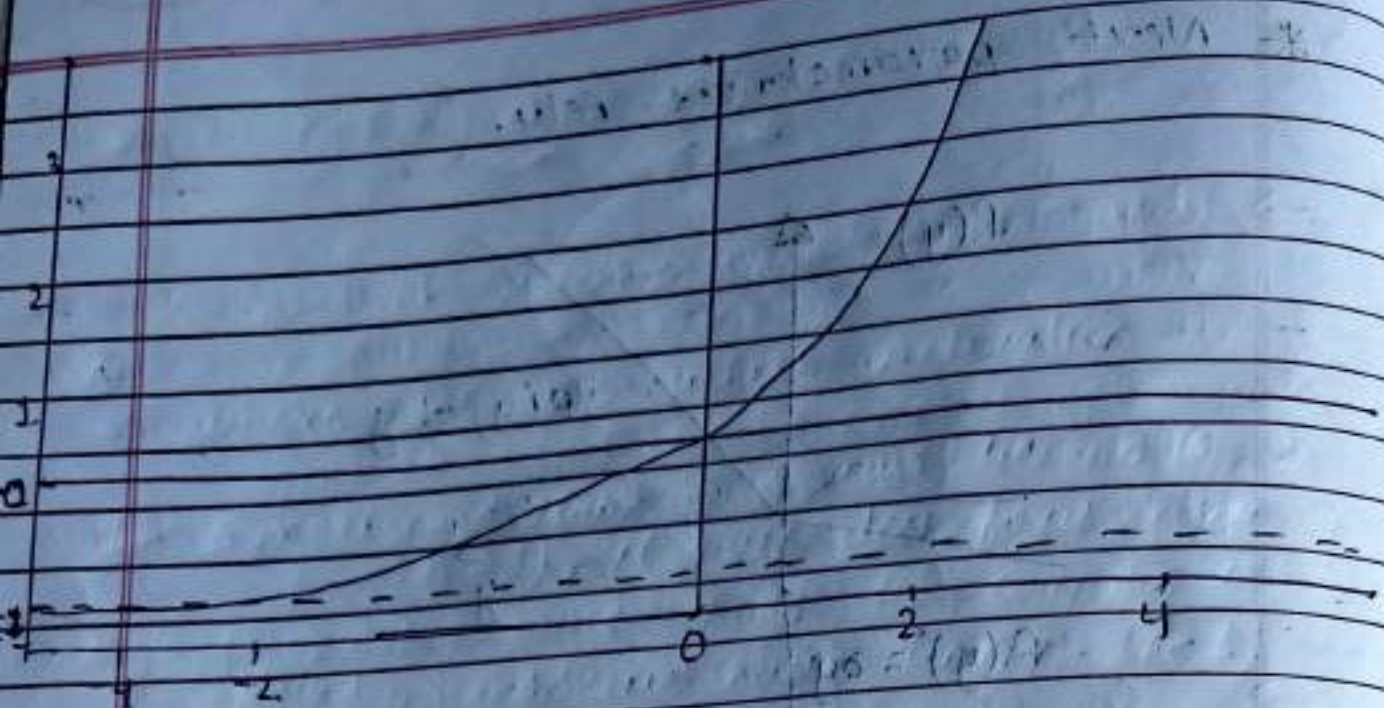
so we have alpha if  $z < 0$  and 1  
if  $z \geq 0$ .

\* ELU  $\rightarrow$  exponential Relu

ELU activation function ( $\alpha = 1$ )





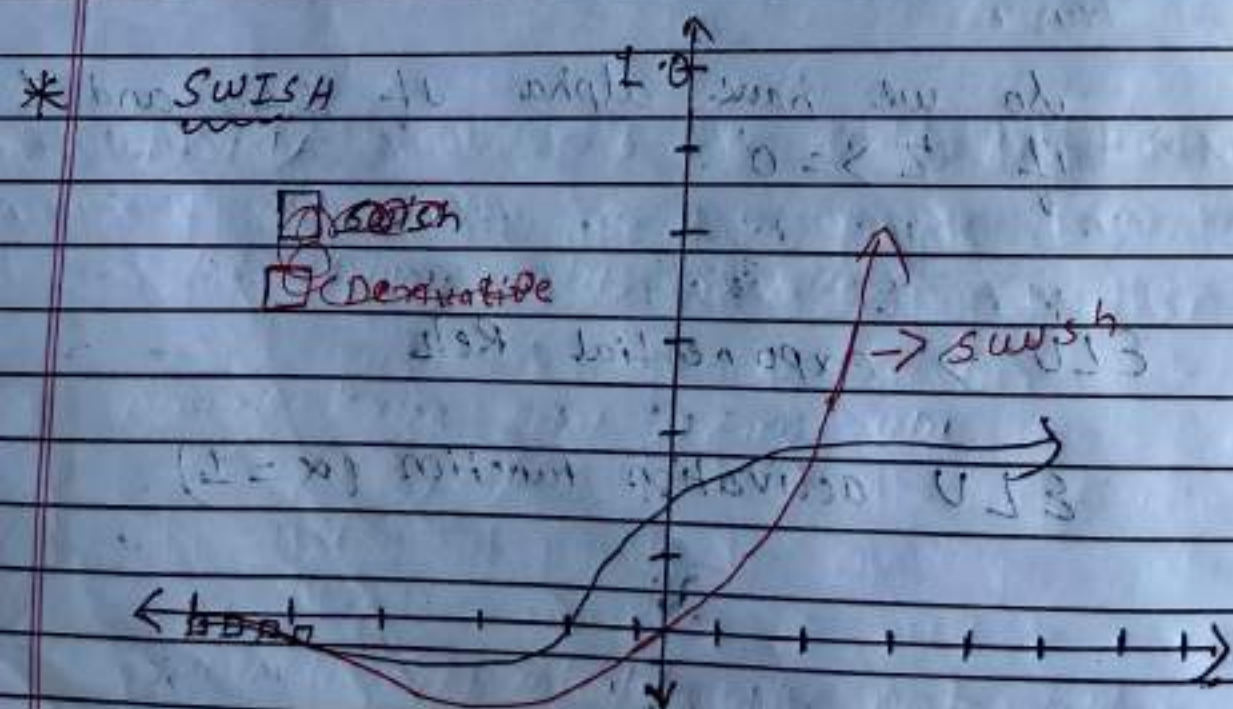


$$f(x) = (e^x - 1), x \geq 0$$

$$f(x) = (e^x - 1), x < 0$$

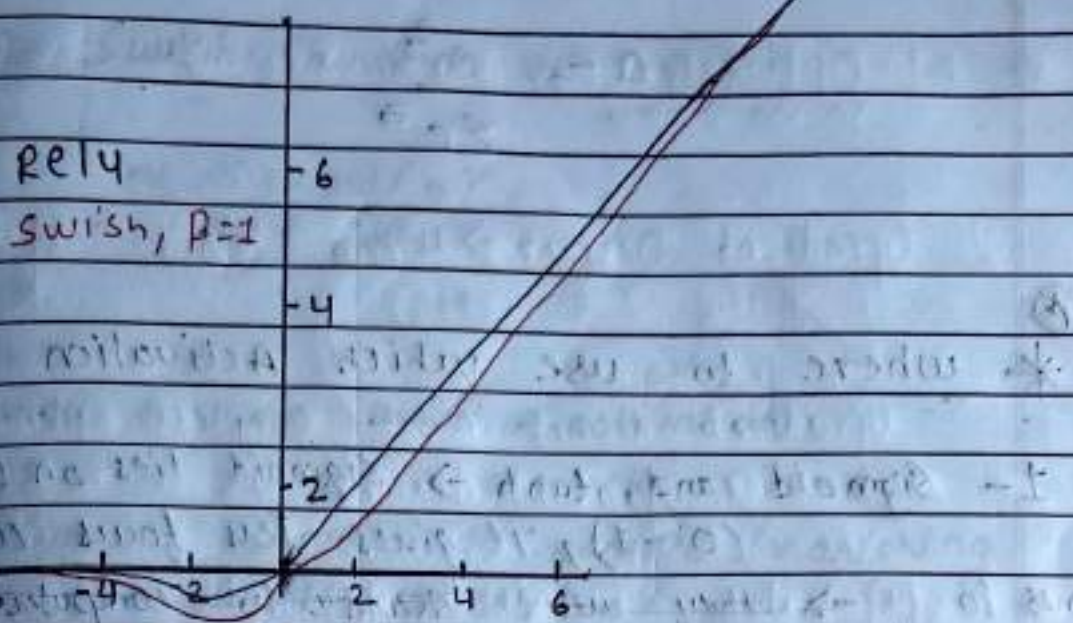
$$f(x) = e^x \rightarrow x < 0$$

$$f(x) = 1 \rightarrow x \geq 0$$

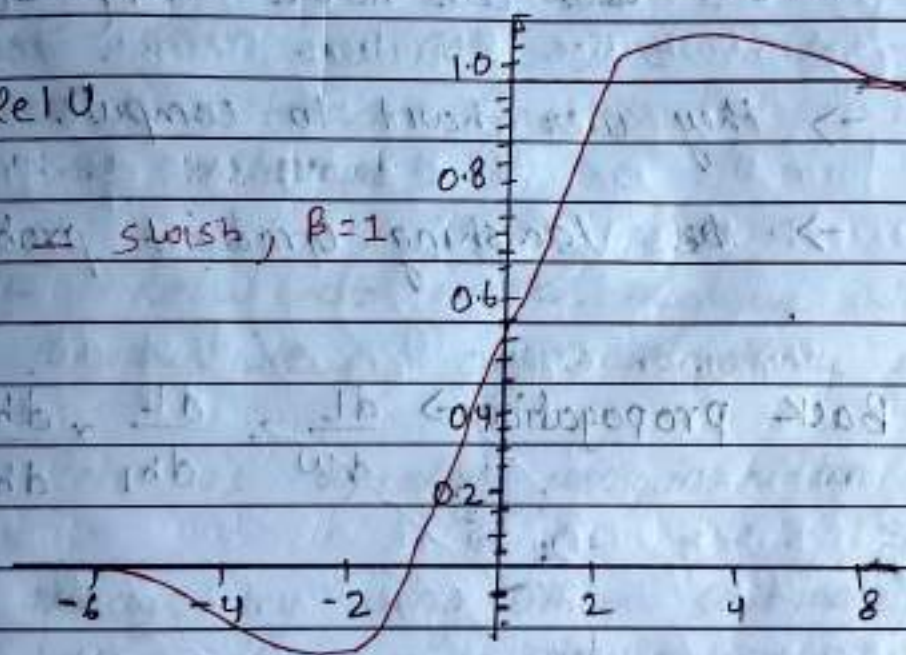




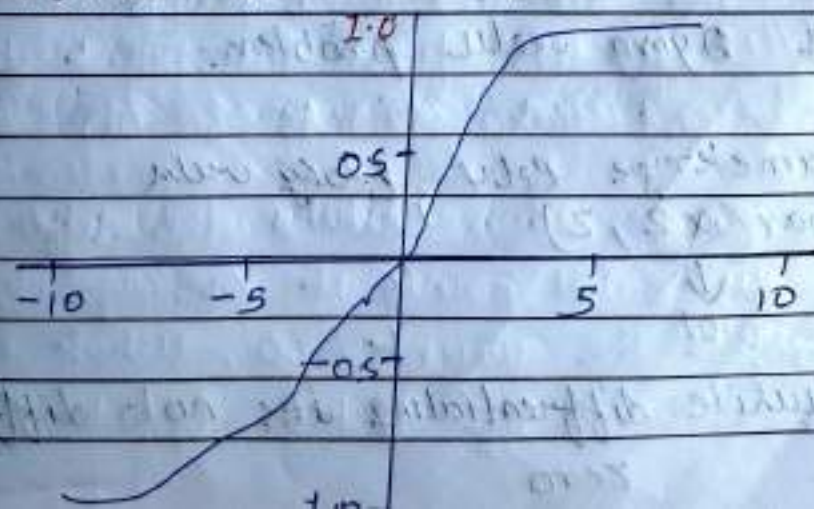
ReLU

Swish,  $\beta=1$ 

\* ReLU

Swish,  $\beta=1$ 

\* Softmax Activation function.





$$a = \frac{e^z}{\sum e^z}$$

⑧

\* where to use which Activation

1- Sigmoid and tanh  $\rightarrow$  Sigmoid lies on ~~(0,1)~~ (0-1), it gives less focus when it comes to  $\rightarrow$  they were hard to compute negative values. ~~we~~ we have tanh (-1 to 1)

$\rightarrow$  They were hard to compute.

$\rightarrow$  ~~the~~ Vanishing Gradient problem.

\* Back propagation  $\rightarrow \frac{dL}{dw} = \frac{dL}{dh_1} \times \frac{dh_1}{dh_2} \times \frac{dh_2}{dw}$

\* 2 ReLU

the Dying ReLU problem.

3- parametrize ReLU, leaky ReLU  
 $\rightarrow \max(\alpha z, z)$

$\downarrow$   
0.01

$\rightarrow$  while differentiating it's not differentiable at zero



## 4 → Sigmoid Function

$$x. \text{Sigmoid}(n)$$

\* where to use which activation

- Sigmoid and tanh avoided (vanishing gradient)

→ ~~ReLU~~ Relu used it's a very versatile activation function and it's used in most problems.

→ if Relu not work then we shift to leaky Relu, or parameterize ReLU

→ Relu is usually used in hidden layers

→ always try using either softmax/sigmoid when it comes to classification problem

\* Optimisers in depths 1



## Optimisers

Gradient Descent  $\rightarrow w_{\text{new}} = w_{\text{old}} - \eta \frac{dL}{dw_{\text{old}}}$

Gradient descent most simplest optimiser and even gradient descent is a 3 types

- (1) Batch Gradient Descent
- (2) mini Batch Gradient Descent
- (3) Stochastic Gradient Descent

formula for all 3 types of gradient descent is same,

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{dL}{dw_{\text{old}}}$$

$\Rightarrow$  Batch GD  $\rightarrow dL$

it calculates the differentiation  $\frac{dL}{dw}$  for all the data points. it calculates for the first data point then second data points so on.

once it calculated the differentiation of all data points it average out and updates the weight. So, what is the problem of Batch Gradient descent



→ it becomes kind of slow

## 2- Mini Batch Gradient Descent

As the name suggest it will not use the all data points. it will randomly select data points it will calculate the  $\frac{dL}{dw}$  of these 100 data points and it will update the weight.

## 3- Stochastic Gradient Descent

Here, it ~~use~~ randomly select one data points and update the weight  $\frac{dL}{dw}$  and weights are change base on  $\frac{dL}{dw}$  the differential of that one data points it take so much time

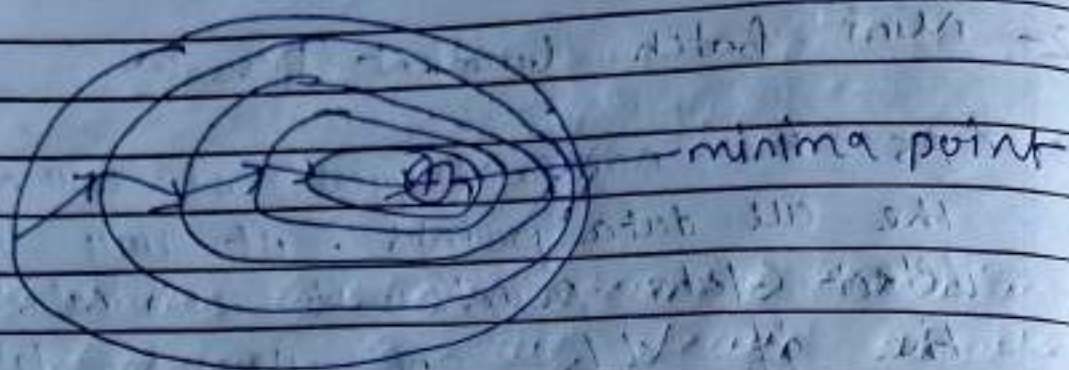
## 4- Batch Gradient Descent



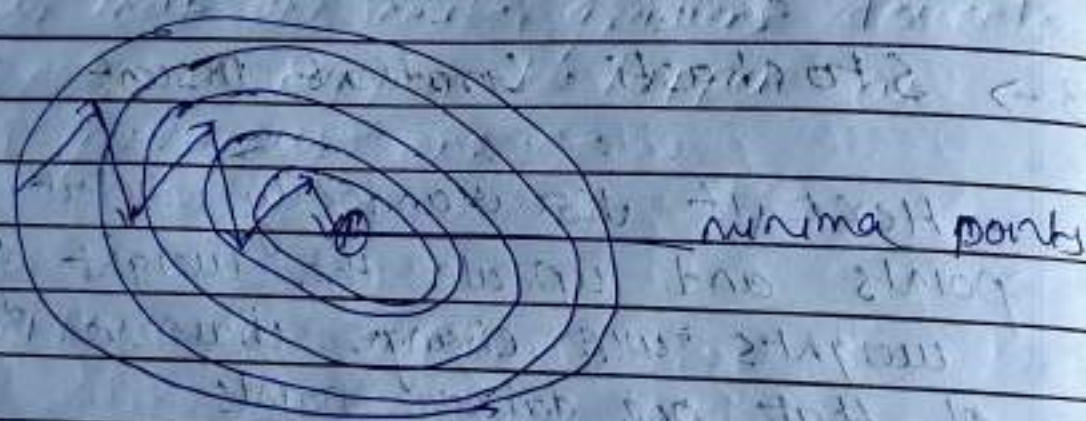
minima point



## \* Mini Batch Gradient Descent



## \* Stochastic Gradient Descent



So what can we do, so that we don't have to use so many data points and it requires so much of times but at the same time we can not be this happy.

So we use Stochastic Gradient Descent or Mini batch Gradient Descent with momentum.





S.G.D

Stochastic Gradient Descent



Stochastic Gradient Descent with momentum

$$\text{So } w_{\text{new}} = w_{\text{old}} - \eta \frac{dL}{dw_{\text{old}}}$$

$$w_t = w_{t-1} - \eta \frac{dL}{dw_{t-1}}$$

$$w_t = w_{t-1} - \eta \frac{dL}{dw_{t-1}}$$

$$v_t = v_{t-1} + \eta g_t$$

$$= v_t = v_{t-1} + \eta g_t$$

$$= w_t = w_{t-1} - \eta v_t$$



$$w_t = w_{t-1} - \eta g_t - \gamma \eta g_{t-1} = \gamma^2 \eta g_{t-2} - \gamma \eta g_{t-1}$$

\* Stochastic Gradient Descent with Momentum  $\rightarrow$   
 is giving so, we learn about the term momentum, so in short basically means giving weightage to the previous calculated gradient descent so that the next gradient descent is not as happened to happen as previous one.

\* Optimizers in depth 2

1) Adagrad  $\rightarrow$

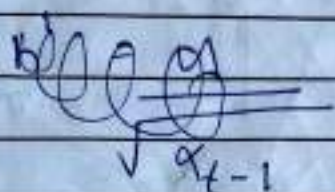
$$w_t = w_{t-1} - \eta g_t$$

idea  $\rightarrow$  So main idea behind Adagrad was, we change instead of changing of gradient descent at all time what if some how come with a method that can change the learning rate so, we need to change the learning rate with every epoch or iteration



- \* To change learning rate with every epoch so that if we get closer to the minima in decreases it self, so we don't have to change the learning rate.

$$w_t = w_{t-1} - \eta g_t$$



problem  $\rightarrow$  Since  $x_{t-1}$  is sum of all the gradient squares

- ② Ada delta  $\rightarrow$  ~~as~~ optimising the weight so ~~as~~ that we can keep it in control ~~as~~ our learning rate does not decrease very fast. So what is the formula behind ada delta

$$w_t = w_{t-1} - \eta g_t$$

$$\eta = \frac{\eta_0}{\sqrt{g_{t-1}} + \epsilon}$$

ada  $\rightarrow$  exponential decay

$$g_{t-1} = \gamma g_{t-2} + (1-\gamma) g_{t-1}^2$$



Suppose  $\rightarrow \gamma = 0.01$  or  $0.99$

$$eda_{t-1} = 0.01 g_{t-1}^2 + 0.99 eda_{t-2}$$

$$= 0.01 g_{t-1}^2 + 0.99 [0.01 g_{t-2}^2 + 0.99 eda_{t-3}]$$

$$= 0.01 g_{t-1}^2 + 0.99 \times 0.01 g_{t-2}^2 + 0.99^2 eda_{t-3}$$

$\rightarrow$  So in one sense this is kind of like momentum term

EDA  $\rightarrow$  Exponential decay average

\* Adam <sup>prop</sup> Adaptive Moment ~~ex~~ estimation

So, what Adam came across was along with using gradient square we will also use gradient that means instead of using EDA doing EDA of just gradient square we will also do EDA for gradient.



$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$$

$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

$$w_t = w_{t-1} - \frac{\alpha \hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon}$$

→ Adam has been proven to be a very good method it is used in most of the cases. Now let us look into which optimizers to use well.

\* Some general rules

\* ~~Sparse~~ Sparse data → like data from bag of words. In that case go for adam function. (means any ada can do working fine) most likely we adam.

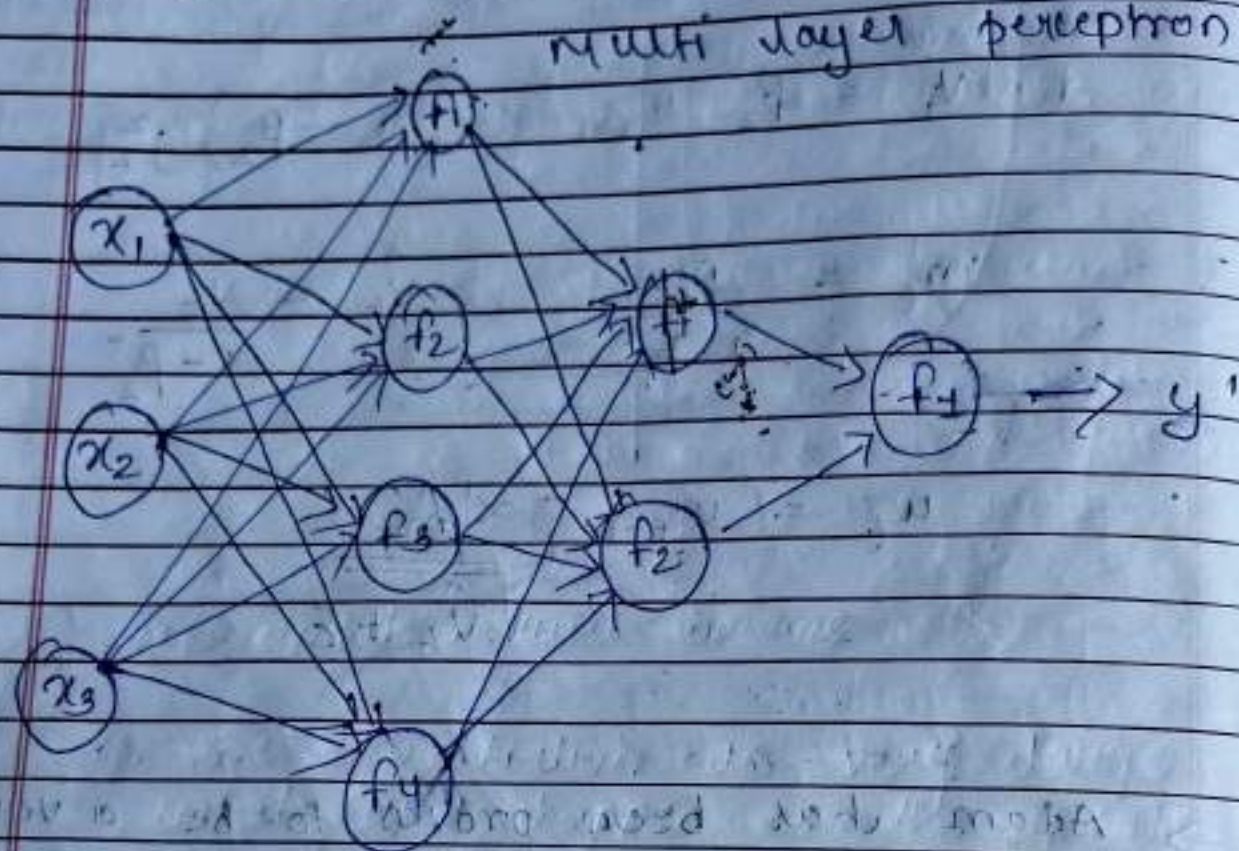
\* 1. Adam → we should use Adam most of the time. So, adam is very versatile. Go it convert its faster and we should use this its most of time in most of the cases.

→



\* weight initialiser  $\rightarrow$

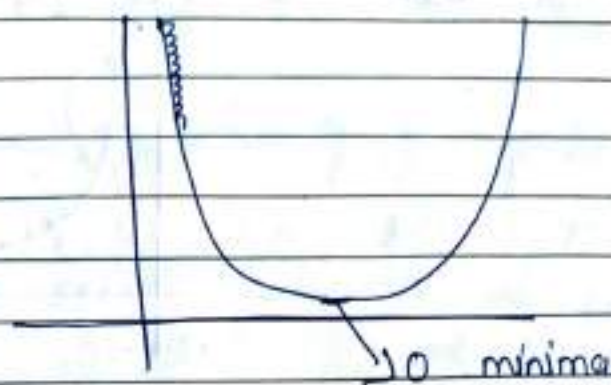
Multi layer perceptron



if you remember structure of a MLP and whenever we assigned weights what i mention was this weights randomly assign. now this statement partially true. so what that means is it's quite random assign yes, but it's not randomly assign from uniform distribution or something. so all the weights will lie in some sort of the distribution (normal distribution, uniform distribution or something).



→ we know that weight should ~~no~~ neither be too small and not too big so why is that.



So, we have the loss here which is against the weight. So what I mention is if we slope ~~no~~ very slowly our weights are not getting updated a lot and it will take very long time to converge or may ~~no~~ be not even reach here.

So, if differentiation is very small our weight will not get updated and this problem statement was known as vanishing problem.

this happens when weights are too small.

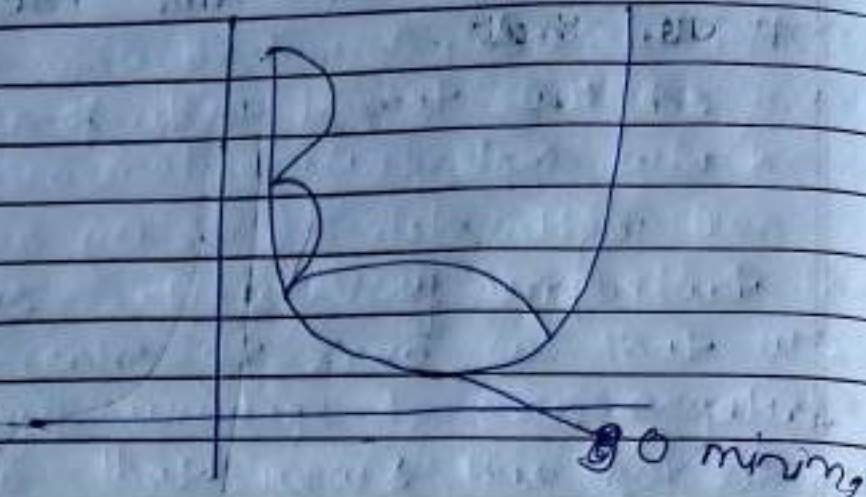
Q- what happens if the weights are too big?

Ans → we come across the problem of gradient explosion or exploding gradient. Just the opposite of



vanishing gradient. So

$\frac{dC}{dW}$



~~the~~ update will become

if this differentiation is very large, what will happen is the update will become chaphazardous. So

let say we are here. So ~~here~~ we go down in a very big step. In the next turn again we will

very big step and suppose we realise forward a bit so we want to go back but since

the steps are very big again go back this whole big time. So this problem is known as Exploding & Gradient problem, and this happens when weights are too big.



\* So we need to somehow clip these weights so we need to keep them in check so that do not cross the particular range. So for that we need a distribution.

\* weights should not be equal?

if all these weights are equal then updates will be equal. the new weights are also equal. hence there is no way of defining which input should be given more weightage compared to the other input. hence we should make sure the weights are (1) not too big, not too small.

(2) the weights are not to be equal.

So they should be randomise, but randomise through a distribution.



So let say this is a neuron it has three inputs and two outputs.

So we have a term called fan in. (Fan in) = number of inputs

fan out  $\rightarrow$  number of outputs



## \* Initialization

1  $\rightarrow$  Normal initialization  $\rightarrow$  So all the weights of this format will be initialized from a normal initialization which lies between  $(0, \sigma)$ , like this is very basic it is very important that we have centre at zero, so that we give enough wiggle room to even negative weights. So that means this is not at all related in any way to the output that we are getting.

2  $\rightarrow$  uniform distribution  $\rightarrow$  So the weights will be initialized from a uniform distribution which lies between  $-1$  to  $1$ .

$$\frac{1}{\sqrt{\text{fan in}}}$$

where it totally depends upon the input. So



uniform distribution

3  $\rightarrow$  Xavier/Glorot initialization  $\rightarrow$  So, these both of the names of Scientist and they



- its usually use for sigmoid function
- They again have two types of different initialization

①

1- Normal initialization  $\rightarrow$  where weight,

$$w_{ij}^k \sim N(0, \sigma) \quad \sigma = \frac{2}{fan_{in} + fan_{out}}$$

2  $\rightarrow$  uniform distribution  $\rightarrow$  the weights randomly initialize from uniform distribution which lies between

$$\left[ -\frac{\sqrt{6}}{\sqrt{fan_{in} + fan_{out}}}, \frac{\sqrt{6}}{\sqrt{fan_{in} + fan_{out}}} \right]$$

3  $\rightarrow$  He initialization

1- Normal initialization  $\rightarrow$  the weights randomly from a normal distribution which lies between

$$(0, \sigma) \quad \sigma = \sqrt{\frac{2}{fan_{in}}}$$

2  $\rightarrow$  uniform distribution  $\rightarrow$  the weights come from a uniform which lies between

$$\left[ -\frac{\sqrt{6}}{\sqrt{fan_{in}}}, \frac{\sqrt{6}}{\sqrt{fan_{in}}} \right]$$



Now, normal and uniform & totally depends on what type of data you are using. So the initialization totally depends on what type of data you are using. So the initialization is usually used with activation function like ReLU, Leaky ReLU,

### \* Losses in Neural Networks

Regression loss  $\rightarrow$

(1) MSE  $\rightarrow$

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

So you calculate the mean of this and you get mean squared error so this is the squared error. So this squared error, error is basically the difference between the actual value and the predicted value. So you have the square of the error and you have to calculate the mean for all the data point

(2) MAE  $\rightarrow$  it's same as

MSE but instead of doing the square, we will be doing

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



Q- why do we need this two function and how are they different

Ans- first of all MSE  $\rightarrow$  is very bad when it comes to finding the outliers. So it's not at all robust to outliers.

in a way like suppose you have a value that is way outside the dataset or way different from the actual distribution of the dataset. So in that case if you will subtract the value is going to be very large value and then again if you squared this is going to be a huge loss. So it actually skews the actual loss. So that is why mean squared error is not robust to outliers.

MSE  $\rightarrow$  it is Robust to outliers. when you have very small error so in that case the problem is the value become very small and the model will not actually learn anything because the loss is going to very small.



\*\*\* Huber loss  $\rightarrow$  is basically a combination of MSE and MAE taking a best of both.

$$\text{Huber loss} = \begin{cases} \frac{1}{2} (y_i - \hat{y}_i)^2 & \text{when} \end{cases}$$

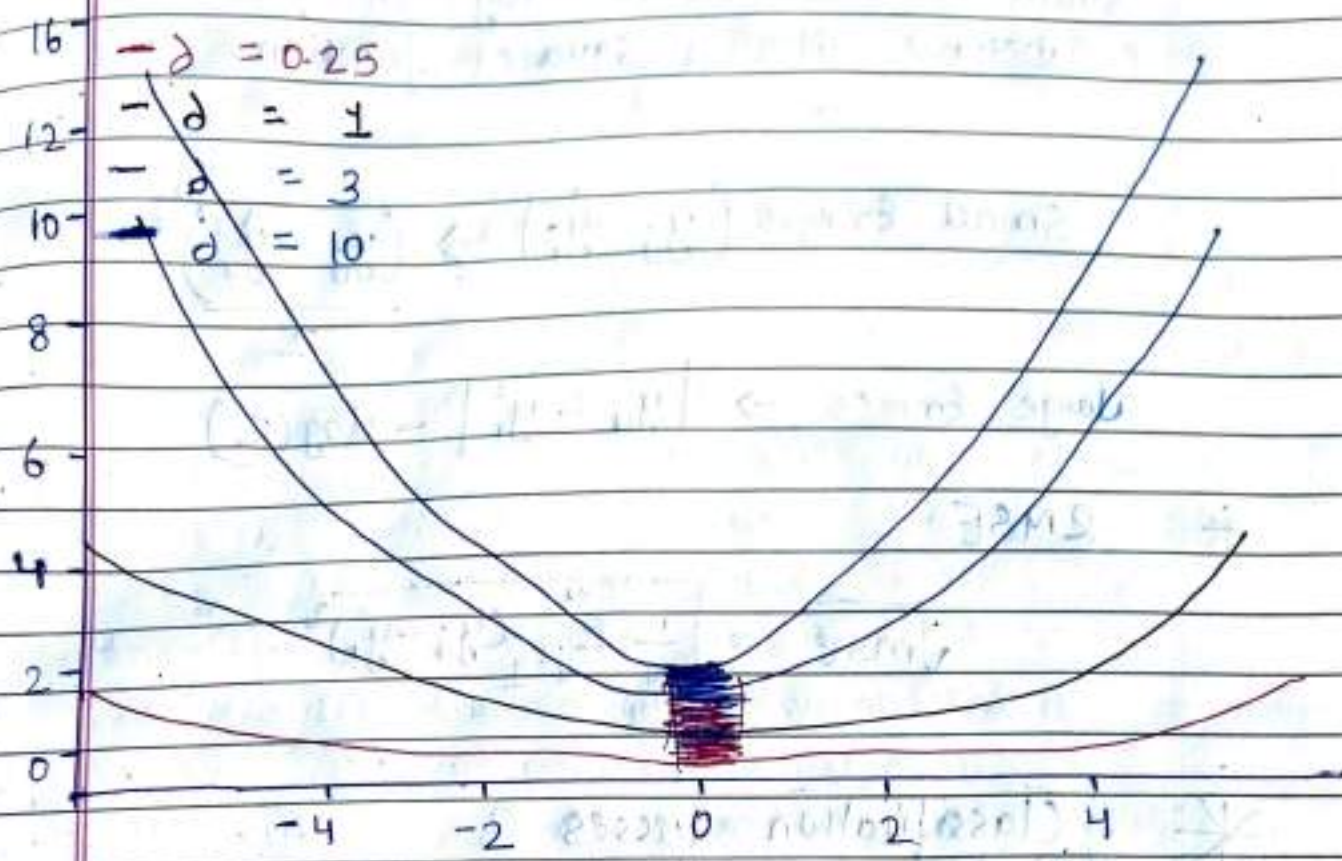
$$\text{mod} \left( |y_i - \hat{y}_i| \leq \delta \right)$$

$$\text{other wise } \delta |y_i - \hat{y}_i| - \frac{\delta^2}{2}$$

if your error is very small so at that case square the losses or the errors. so it will square the error. so your loss is not the very small and it will penalized the loss as well as the loss is very big it will not do the square part of it will just use the absolute value.



Huber loss for target = 0



predictions for target = 0

~~problem~~problem  $\rightarrow$  deciding the value of delta
 $\delta$   
 $\{$ 

\* log cosh  $\rightarrow$  it is a good loss function. ~~be~~  
 it is not exactly ~~you~~ you used a dot.

$$-\log(\cosh(y_i - \hat{y}_i))$$

for very small values  $y_i$ , for very



Small errors. ~~like~~ Log cash function behave like a squared function.

$$\text{Small Errors } (y_i - \hat{y}_i) \rightarrow \frac{(y_i - \hat{y}_i)^2}{2}$$

$$\text{Large Errors} \rightarrow |y_i - \hat{y}_i| - \log(2)$$

\* RMSE

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

\* Classification losses

(I) Binary classification problem

$$\text{(1) Binary cross entropy} \rightarrow -\frac{1}{N} \sum_{i=1}^N y_i \log \hat{y}_i +$$

labels  $\rightarrow (0, 1)$

$$(1 - y_i) \log (1 - \hat{y}_i)$$

$$(i) \text{ Entropy} = y_i \log y_i$$

$$(ii) \text{ Cross Entropy} = y_i \log \hat{y}_i$$

$$\text{(2) Hinge loss} \rightarrow \max(0, 1 - y_i \hat{y}_i)$$

labels are  $\rightarrow (-1, 1)$

Suppose your  $\hat{y}_i$  is very close to  $y_i$   
Let's say both are  $-1$ , so in that case



what will happen • you will have maximum off  
 $\max(0, 1 - (1 \times 1))$

$$\max(0, 1 - 1)$$

$$\max(0)$$

if your actual predicted value is very close to the actual value then your loss will be 0. in the same way.

## (II) Classification losses for multiclass problem

It is a same thing you don't have a lots of different as such so instead of binary cross entropy we have categorical cross entropy.

### (i) categorical cross entropy →

$$= \frac{1}{N} \sum_{j=1}^N \left[ \sum_{i=1}^C -y_{ij} \log y_{ij} \right]$$

(ii) KL Divergence → This is use when you have very different distribution and you need to find the difference between the distribution.



KL :

$$= \sum p(x) [\log(p(x)) - \log(q(x))]$$

- KL divergence is used in generative models.

We ~~discuss~~ discuss about the different classification functions of neural network and we discuss how similar this was to machine learning losses. So finally we completed one basic neural network structure, which is multi-layer perceptron. We went over all the concepts of optimisers, layers, activation function and everything. So we completed multi-layer perceptron.