

# Neural Network

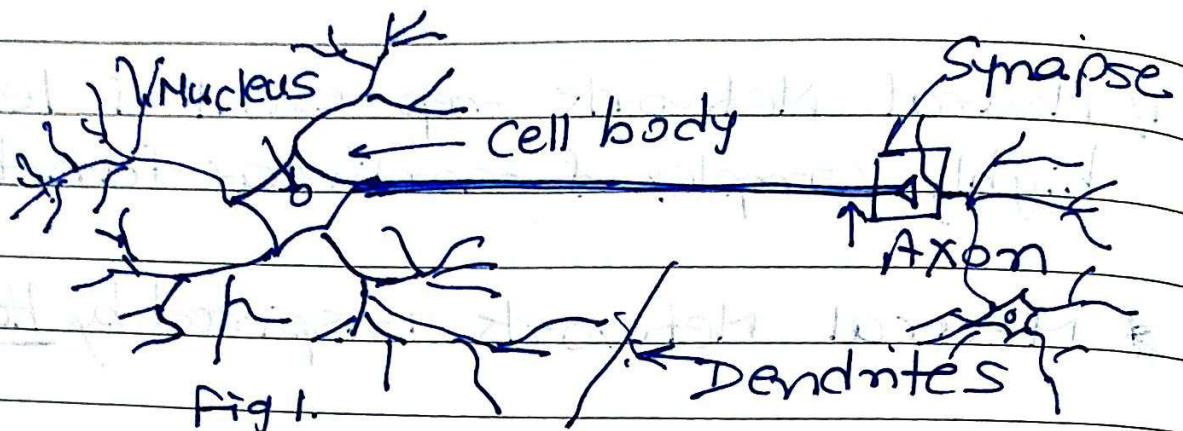
- 1) Neural Network represent & learn highly complex & nonlinear function.

\* Neural Network Inspired by Human Brain

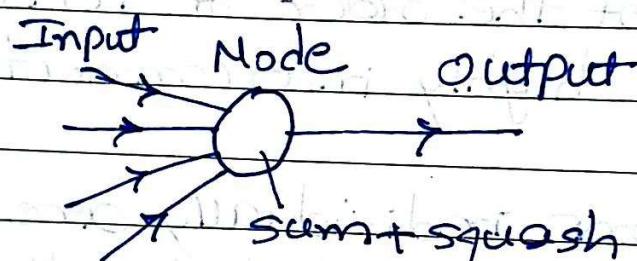
## Human Brain

- 1) Human brain contains a number of neurons, of the order of tens of billions, which are highly interconnected.
- 2) These neurons individually are simple computing units, but together they can perform very complex tasks.
- 3) Some characteristics of neural human brain incorporated to form Neural Network  
Following are characteristics.
  - 1) massive Parallelism
  - 2) Neurons are highly inter connected to each other & solve complex tasks together
  - 3) Model distributive associative memory through weights in synaptic connection.

## Human - Neuron



## Neural Unit

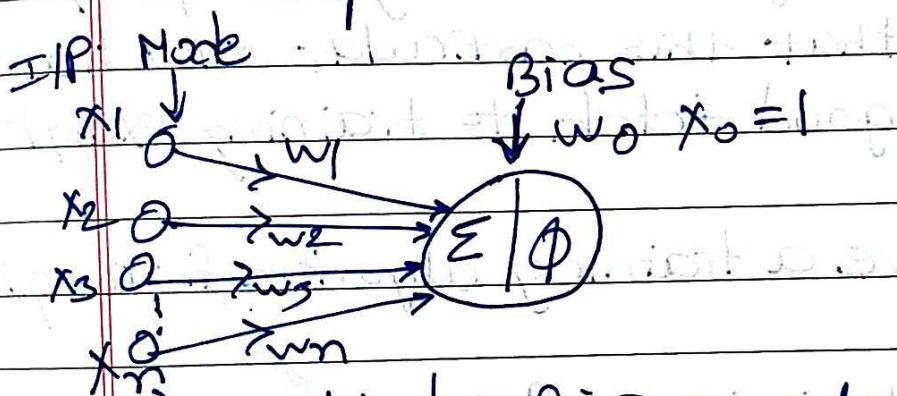


- 1) As shown in fig 1., there is Nucleus & dendrites through which IP is expected
- 2) Through Axon is the tail through which the output is given.
- 3) The fig 1 structure is of the neuron is simulated by a neural network unit as shown in fig 2. which has inputs & then some simple computation are at node. & OP unit
- 4) for computation use weighted sum of the inputs & then applies a function as squashing function (for ex Sigmoidfun<sup>n</sup>)

- 5) Through the dendrites the IIP is accepted & through the axon the output is transmitted through the electric impulses through the synaptic layer to other neurons
- 6) So neural NW inspired by neuron human brain.

## Single Neural Network

### Perceptron NW.



1) weighted IIP is computed along with bias & transfer through transfer function

1) for example thresholding transfer function

$$\phi(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

2) If weighted IIP & bias summation is greater than threshold OIP is high(1) otherwise low(0)

2) Some nonlinear function

1) Sigmoid function 2) Tanh function

$$y_{in} = \sum_{i=1}^N x_i w_i + b$$

$$\therefore y_{out} = f(y_{in})$$

Consider training example

$x_1 \quad y_1$  Training the network means  
 $x_2 \quad y_2$  learning these weights  $w_0, w_1, w_2, w_3, \dots, w_n$   
 $x_3 \quad y_3$  for given training example, so  
 $x_{in} \quad y_{in}$  that this particular net has a  
good fit to training example.

So we have a training algorithm for Perceptron

#### \* Perceptron training Rule.

1) Initialise weights  $w_0, w_1, w_2, \dots, w_n$  with initial values of the weights. At each iteration we will update the weights.

2) For  $x_i$  select  $w_i$  for that o/p  $y_i$ , check if o/p  $y$  is wrong need to update weight. So that o/p is closer to given o/p.

$$\therefore \text{New } w_i = \text{Old } w_i + \Delta w_i$$

$$\Delta w_i = \eta (y - \hat{y}) x_i$$

$\eta$  is learning rate, it controls how much you change based on error. the weight updation.

By taking the example one by one. It can be shown that this perceptron learning converges, if  $D$  is linearly separable.

That means, there is a linear decision surface that separates the positive & negative examples. If such a linear separation surface exists then by applying this perceptron training rule, the learning algorithm will converge to hypothesis which will separate the positive & negative example.

→ But if data is not linearly separable then this perceptron learning algorithm will not converge, it will not work.

→ Here use Gradient Descent Algorithm for weight updation.

$$E = \frac{1}{2} \sum (y_d - \hat{y}_d)^2$$

$$\phi(z) = z$$

Stochastic Gradient

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$

gradient descent give direction for local minima

So to handle nonlinear data we need to go for multilayer network.

Multilayer network connected with each other

In order to represent complex functions

we want non-linear unit, which is

differentiable, so we go for a transfer function which is differentiable & non-linear.

Often used transfer function is

Sigmoid Function (Logistic function)

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

Differential of sigmoid function is

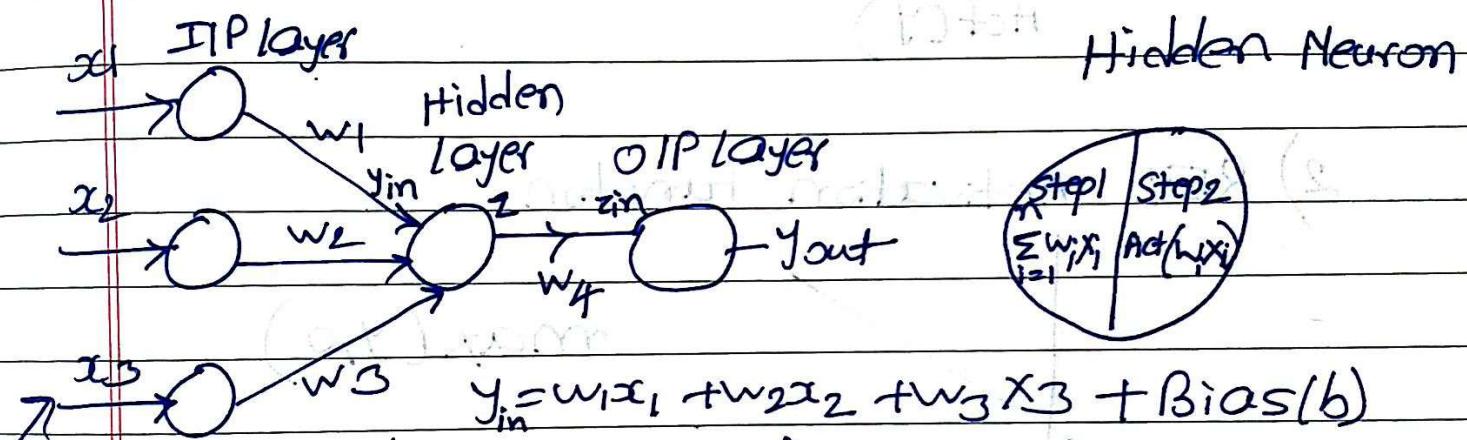
$$\phi'(z) = \phi(z)(1-\phi(z))$$

Error is  $E = \frac{1}{2} \cdot \sum (y_i - \hat{y}_i)^2$

$$\frac{\partial E}{\partial w_i} = \frac{1}{2} \sum \frac{\partial E}{\partial y_i} \cdot \frac{\partial y_i}{\partial w_i}$$

# Basic Architecture of Neural Network

How Neural Network works



Forward Propagation  $Z = \text{Act}(y)$

Weight & Activation fun<sup>n</sup> playing important role.

For ex. - Activation function Sigmoid =  $\frac{1}{1+e^{-y}}$

It transform OIP in bet<sup>n</sup> 0 to 1

→ when OIP of activation fun<sup>n</sup> is less than

→ 0.5 means OIP is zero means neuron is not activated

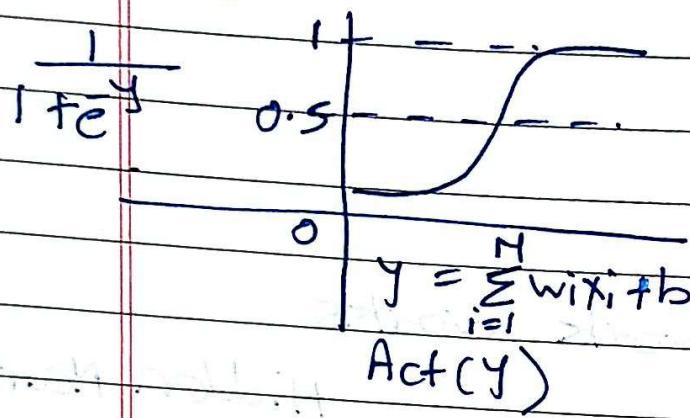
If it is greater than 0.5 the neuron gets activated.

This is forward  $Z_{in} = Z \times w_4$

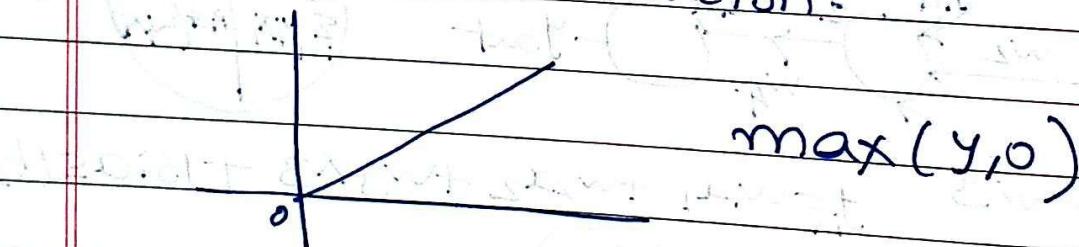
Activation function  $y_{out} = f(Z_{in})$

## \* Activation Functions

### 1) Sigmoid Activation Function.



### 2) REL Activation Function.



When  $y$  is -ve O/P is zero

When  $y$  is +ve O/P is w.r.t that

# Geoffrey Hinton

CLASSMATE

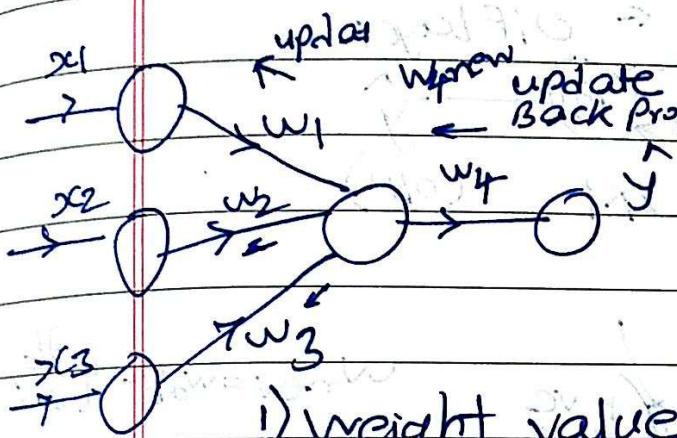
Date \_\_\_\_\_

Page \_\_\_\_\_

## \* Train Neural Network with Backpropagation

$$\hat{y} = [x_1 w_1 + x_2 w_2 + x_3 w_3] + b_1$$

$$z = \text{Act}(y)$$



$x_1, x_2, x_3$  Pass/Fail

Play Study Sleep O/P

update  $w_1$   $w_2$   $w_3$   $w_4$  Back Propagate Optimization is used to

minimize Loss

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

1) weight value should be adjusted to reduce loss value.

- 2) To reduce the loss value we need to backpropagate & update the weight value.
- 3) Update  $w_4$  by backpropagate.

$$w_4 = w_4 - \eta \frac{\partial L}{\partial w_4}$$

Find derivative of loss wrt.  $w_4$

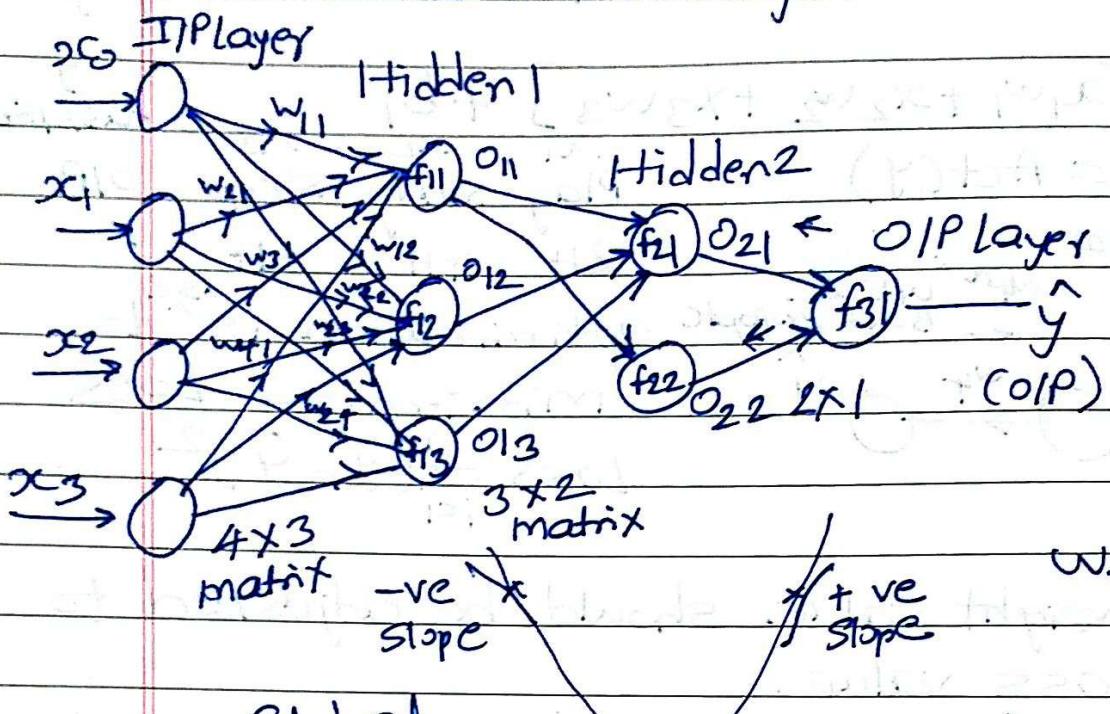
$\eta$  = learning rate - Small value it will help us to achieve global minima.

- 4) Same way update  $w_1, w_2, w_3$  also

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

- 5) Once we get updated value forward propagation happen to find predicate  $\hat{y}$
- 6) This will repeat the process epochs/iterations until loss gets reduced.

## \* How to train multilayer Neural Networks



- 1) IIP matrix of size  $4 \times 3$
- 2)  $f_{11}, f_{12}, f_{13}, f_{21}, f_{22}, f_{31}$  Neuron function which perform 2 steps
  - 1) IIP multiply with weight & Summation
  - 2) Apply this Summation to activation function
- 3) After all layer execution we get predicted OIP  $\hat{y}$ , by which we find out Loss function
- 4) Main aim is to reduce loss function to achieve actual output.
- 5) To reduce loss function we have to use optimizer
- 6) Here we use Gradient Descent  
Stochastic Gradient Descent

why derivative to adjust the  
slope.

CLASSMATE

Date \_\_\_\_\_

Page \_\_\_\_\_

to update the weight

\* Gradient Descent Optimizer  
use to update (optimize) weight value  
to reduce loss function.

This will done by backpropagation.

\* Using backpropagation we need to update  
weight.

\* If loss function is not reduced means  
there is problem with learning rate &  
derivative value & it will not converges  
properly.

So need to adjust properly learning rate

minimize loss

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

require optimizer

$$y = [x_1 w_1 + x_2 w_2 + x_3 w_3] + b$$

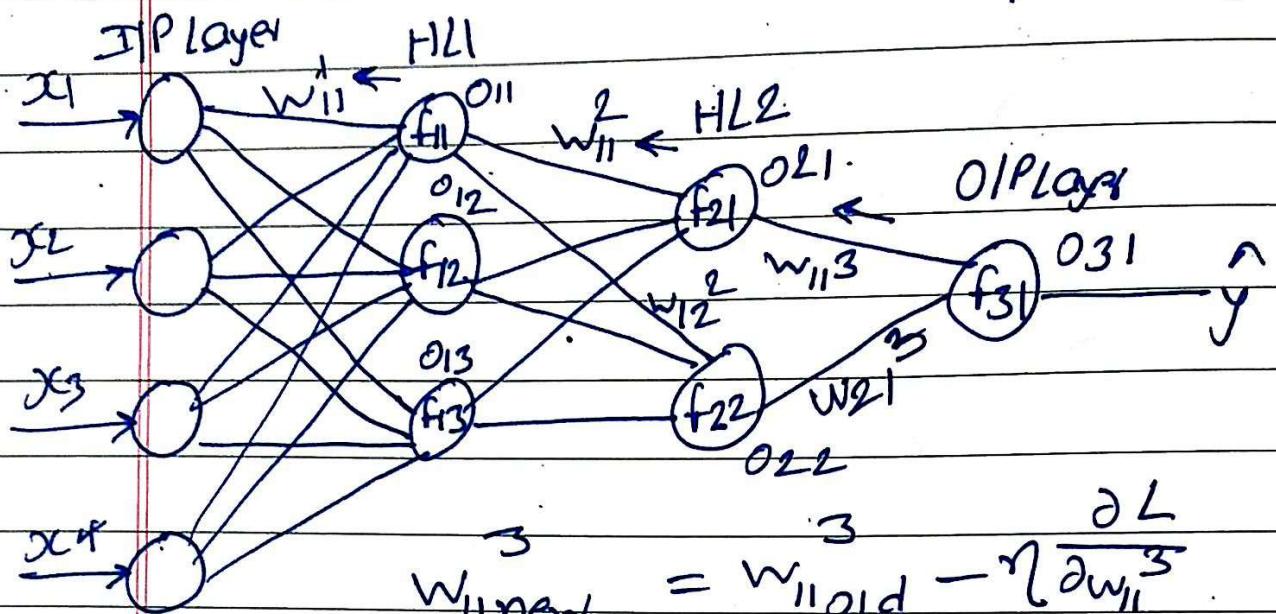
$$z = \text{Act}(y)$$

$$w_3 \text{new} = w_3 - \eta \frac{\partial L}{\partial w_3}$$

$$w_3 \text{new} = w_3 \text{old} - \eta \frac{\partial L}{\partial w_3}$$

$\eta$  = Learning Rate

## \* Chain Rule In Backpropagation



$$w_{11\text{new}}^3 = w_{11\text{old}}^3 - \eta \frac{\partial L}{\partial w_{11}^3}$$

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

Reduce Loss we need Optimizer  
(for ex: Gradient Descent, Stochastic Gradient)

through backpropagation we update each & every weight of hidden layer

Now through backpropagation we need to update weight from back.

1) first find loss function slope of weight

$$2) w_{11\text{new}}^3 = w_{11\text{old}}^3 - \eta \frac{\partial L}{\partial w_{11}^3} \quad \begin{matrix} \downarrow \\ n : \text{Learning Rate.} \end{matrix}$$

$$\frac{\partial L}{\partial w_{11}^3} = \frac{\partial L}{\partial o_{31}} \cdot \frac{\partial o_{31}}{\partial w_{11}^3} \quad \text{--- Chain Rule}$$

o) for  $w_{21\text{new}}^3 = w_{21\text{old}} - \eta \frac{\partial L}{\partial w_{21}^3}$

$$\frac{\partial L}{\partial w_{21}^3} = \frac{\partial L}{\partial w_{21}^3} \times \frac{\partial w_{31}}{\partial w_{21}^3}$$

1)  $w_{11\text{new}}^2 = w_{11\text{old}}^2 - \eta \frac{\partial L}{\partial w_{11}^2}$

$$\frac{\partial L}{\partial w_{11}^2} = \left[ \frac{\partial L}{\partial w_{31}} \times \frac{\partial w_{31}}{\partial w_{21}} \times \frac{\partial w_{21}}{\partial w_{11}^2} \right]$$

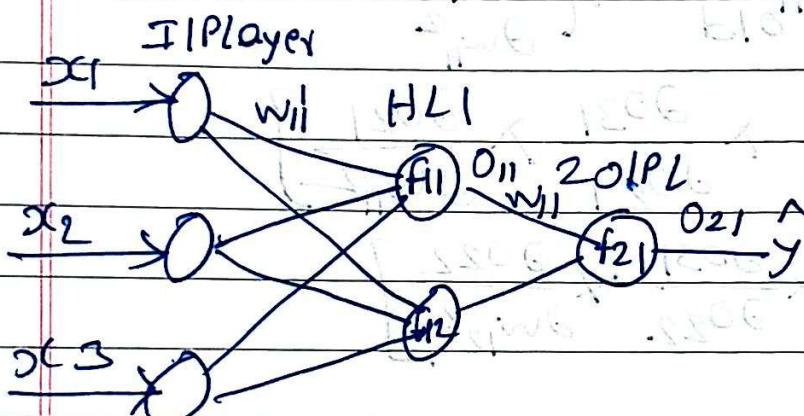
$$+ \left[ \frac{\partial L}{\partial w_{31}} \times \frac{\partial w_{31}}{\partial w_{22}} \times \frac{\partial w_{22}}{\partial w_{11}^2} \right]$$

This derivative means weights are updated unless & centill we reach global minima.

1 Epoch consider both forward & backward propagation

## \* Vanishing Gradient Problem

\* In 1990 to 2000 research cannot create deep neural net because of vanishing gradient because at that time ReLU is not invented.



### Weight update

$$w_{11\text{new}} = w_{11\text{old}} - \eta \frac{\partial L}{\partial w_{11\text{old}}}$$

1) At Layer 1 in neuron we do a summation of iIP multiply by weight & add bias.

$$f_{11\text{in}} = (x_1 w_{11} + x_2 w_{12} + x_3 w_{13}) + \text{bias}$$

$$\therefore O_{11} = \text{Activation function } (f_{11})$$

Here we use Sigmoid function which gives 0IP in between of 1

2) After getting first transformation this output get transfer to second layer & weight get

assign here  $w_{11}^2$  & again ~~the~~ through optimizer we get output by reducing loss by updating weight.

Using chain rule & weight updation

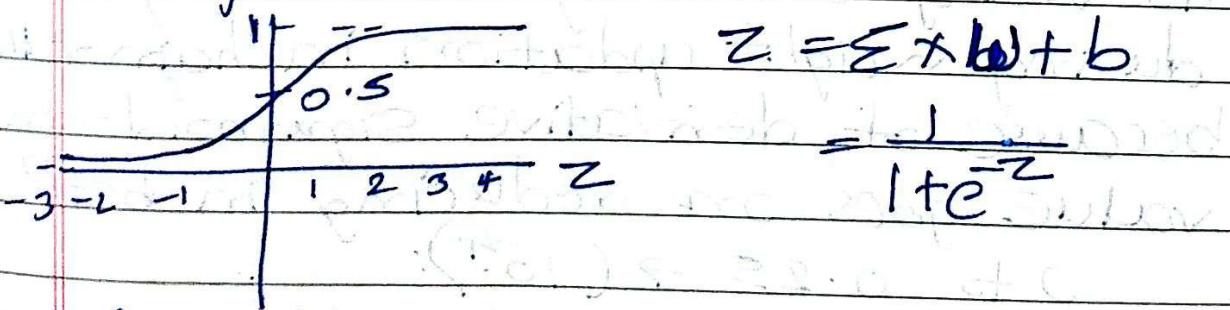
$$w_{11\text{new}} = w_{11\text{old}} - \eta \frac{\partial L}{\partial w_{11\text{old}}}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial o_{21}}{\partial o_{11}} \cdot \frac{\partial o_{11}}{\partial w_{11}} \leftarrow \text{chain rule}$$

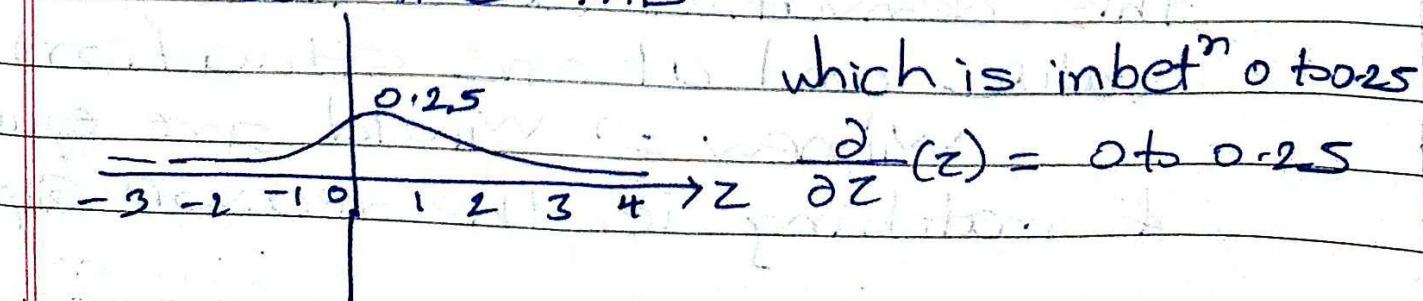
For neuron O/P we use Sigmoid activation function which transforms O/P in bet<sup>n</sup> 0 to 1

But derivative of Sigmoid specifically bet<sup>n</sup> 0 to 0.25

so Sigmoid is



But derivative of Sigmoid function will look like this



So derivative of Sigmoid is in range  
 $0 \leq g(z) \leq 0.25$

As in our eqn

$$\frac{\partial L}{\partial w_{ii}} = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial w_{ii}}$$

$\frac{\partial z}{\partial w_{ii}} = \frac{0.20}{0.02}$  value of gradient  
 Suppose using Sigmoid fun which come in bet' 0 to 0.25 after derivative.

it goes on reducing Suppose 0.04  
 for example  $\frac{\partial L}{\partial w_{ii}} = 0.04$

Original weight  $2.5 = w_{ii, old}$ .  $\eta = 1$

$$\therefore w_{ii, new} = w_{ii, old} - \eta \frac{\partial L}{\partial w_{ii, old}}$$

$$\begin{aligned} w_{ii, new} &= 2.5 - 1(0.04) \\ &= 2.46 \end{aligned}$$

As we go on adding hidden layer of during weight updation in chain rule because of derivative Sigmoid function value goes on reducing in bet' 0 to 0.25  $\rightarrow (10^{-4})$

\* So as the number of layers increase this derivative value is very smaller value and at one situation

$w_{ii, new} = w_{ii, old}$  are equal & matching ie the reason Sigmoid

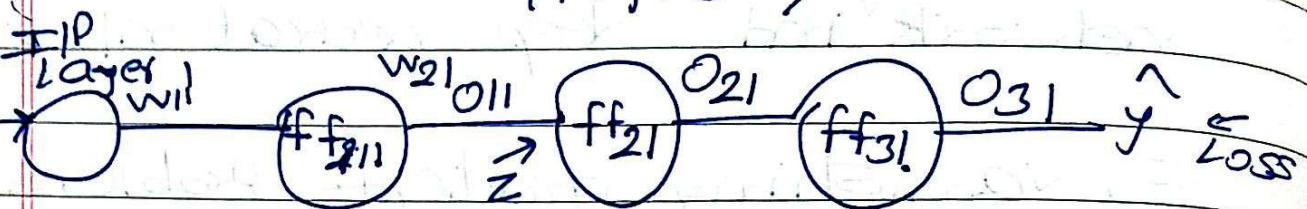
not work properly when we added more layer & convert neural network into deep neural nw.

So vanishing gradient problem is due to Sigmoid Activation function as value is in bet<sup>n</sup> 0 to 0.25

Due to this ReLU Sigmoid function comes.

for tanh transfer function also same problem comes.

## \* Exploding Gradient Problem (Derivative value become higher)



$$w_{1\text{new}} = w_{1\text{old}} - \eta \frac{\partial L}{\partial w_{11}}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial O_{11}} \cdot \frac{\partial O_{11}}{\partial w_{11}}$$

$$z = w_{21}O_{11} + b_2$$

$$\therefore O_{21} = \phi(z)$$

$\phi$  is sigmoid activation function

$$\therefore \frac{\partial O_{21}}{\partial O_{11}} = \frac{\partial \phi(z)}{\partial z} \times \frac{\partial z}{\partial O_{11}} \quad 0 \leq \phi(z) \leq 0.25$$

$$\text{for ex consider } = 0 \leq \phi(z) \leq 0.25 \times \frac{\partial (w_{21}O_{11} + b_2)}{\partial O_{11}}$$

$$= 0.5\phi(z) \leq 0.25 \times w_{21} \quad \text{assume } w_{21} = 500$$

so

If weights are initialized to very high value we get variance high is also if it never get convergence to global minima. if  $\frac{\partial L}{\partial w_{11}} = 200 \times 125 \times 100$  which is very important. due to which  $w_{11\text{new}}$  is -ve

as  $\frac{\partial L}{\partial w_{11}}$  is high & situation occurs of Exploding Gradient

## \* Drop Out & Regularization.

- 1) when we have many layer in artificial neural network & it is deep network then we have more no. of weights & bias parameter.
- 2) Due to with ANN <sup>tend to</sup> becomes overfit the dataset problem or a particular use case so we should have to set weight to overcome fitting problem.
- 3) In multilayer neural network we never face underfitting problem but overfitting may problem & facing high variance problem.

- A) There are two ways to avoid overfitting problem.
  - 1) Regularization by  $L_1$  &  $L_2$  Norm.
  - 2) Drop Out

- \* Dropout comes in 2014 by
  - 1) Nitish Srivastav Thesis.
  - 2) Geoffrey Hinton

for ex - Random Forest Tree

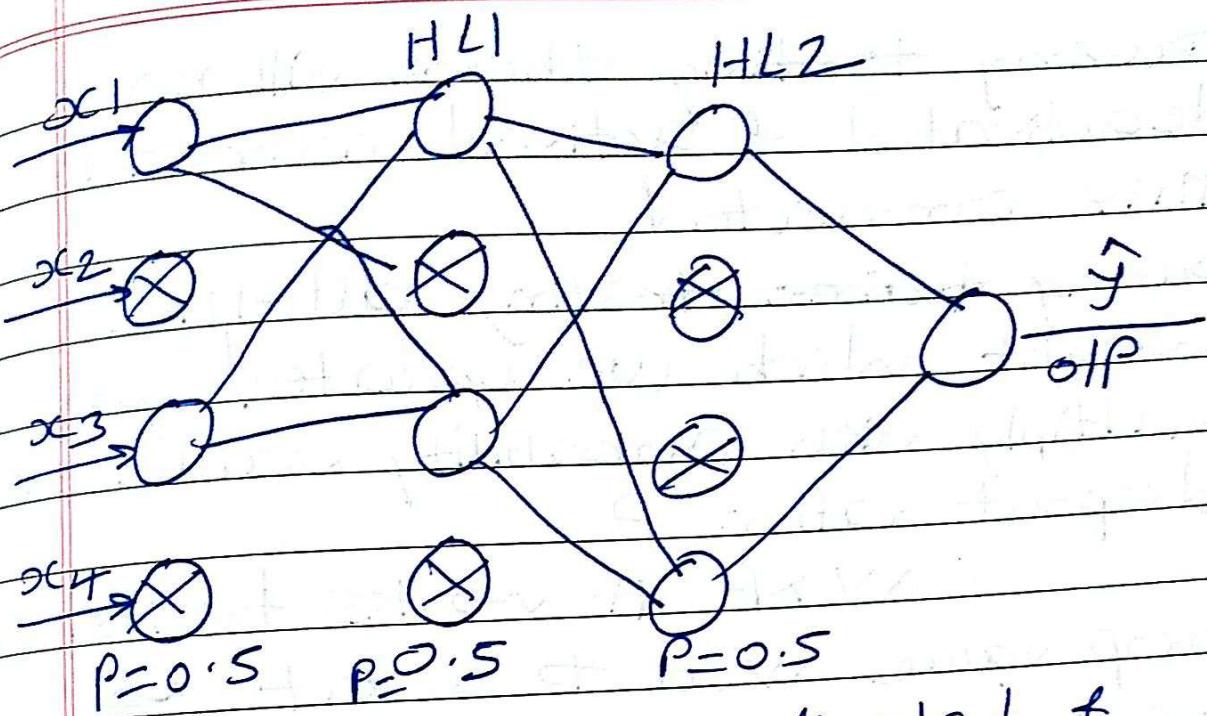
- i) In decision tree when we try to reach in depth overfitting occurs
- ii) So in random forest we use subset of features & form multiple tree → aside overfitting & improve accuracy

- \* Similarly in Dropout

To implement we select Dropout ratio  $\rightarrow P$  which is  $0 \leq P \leq 1$

- 1) We select subset of feature from i/p layer, similarly subset of feature from hidden layer
- 2) So we are not selecting all feature we are selecting subset of feature along with i/p layer to hidden layer

So as shown in fig. Suppose  
 $P$  i.e Dropout Ration  $P=0.5$  &  
 we select subset of feature.



∴ 2 Neuron are activated &  
2 Neuron are deactivated

During forward & backward propagation  
based on P value. it will randomly  
deactivated some neuron

In backpropagation the neuron which  
are activated for that weight gets  
updated. Again in forward propagation  
randomly some neuron activated &  
deactivated.

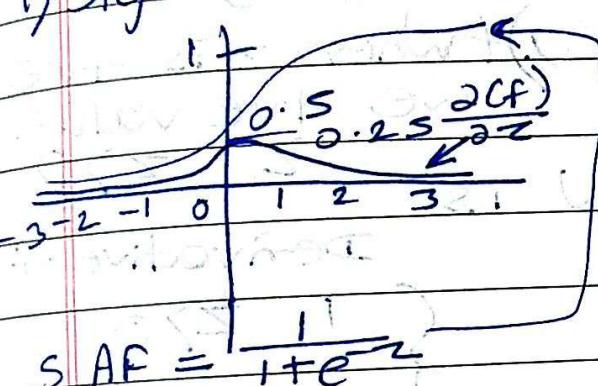
So to avoid overfitting we are  
deactivating & activating neuron or  
in feature by drop out ratio in  
training.

- \* During testing there will no deactivated & activated neuron even this connected
- \* During ~~train~~ testing all the weight which we updated gets multiply with probability value or dropout ration P
  - :  $w \times p$  it works for test data
- \* Drop value need to select as hyperparameter optimization

# \* Rectified Linear Unit & Leaky Rule

## \* Summary

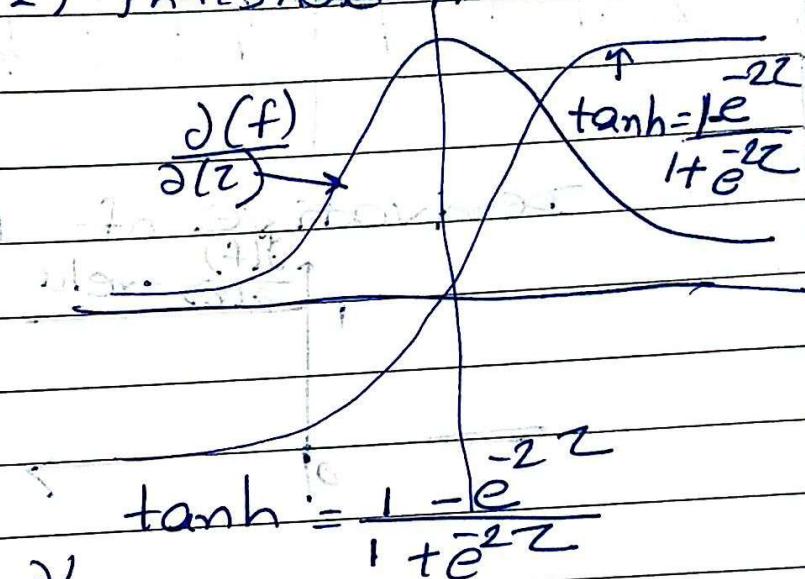
### 1) Sigmoid AF



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

$$0 \leq \frac{\partial f}{\partial z} \leq 0.25$$

### 2) Threshold Activation



$$\tanh = \frac{1-e^{-2z}}{1+e^{-2z}}$$

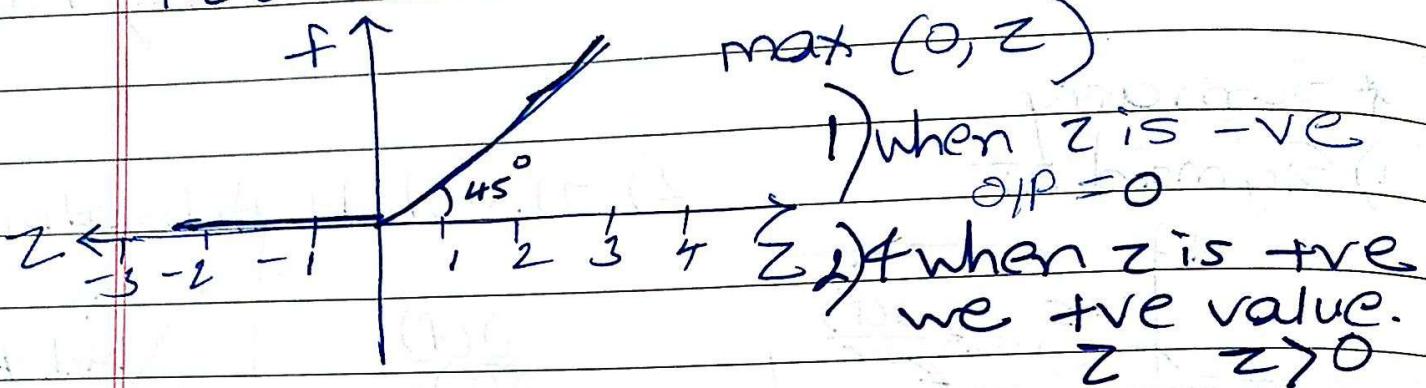
$$W_{new} = W_{old} + \eta \frac{\partial L}{\partial w}$$

Here derivative is less than one during backpropagation if  $W_{new} > W_{old}$

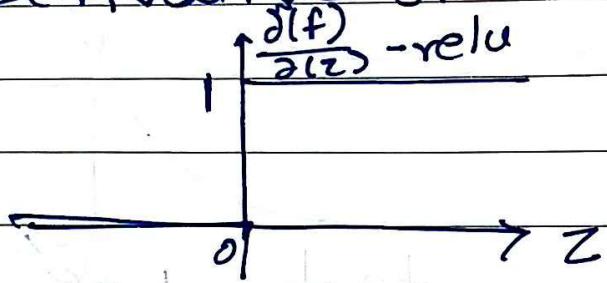
Because of this there is very minor change in  $W_{new}$  it will take much time to converges to global minimum.

So this vanishing gradient problem is handle by ReLU activation function

## ReLU (Rectified Linear Unit)



Derivative of ReLU is



Derivative of ReLU

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

So in sigmoid always derivative value is 0 to 0.25 & in tanh it is less than one. But in ReLU we have  $\frac{\partial f}{\partial z} = 1$  &  $0$ .

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

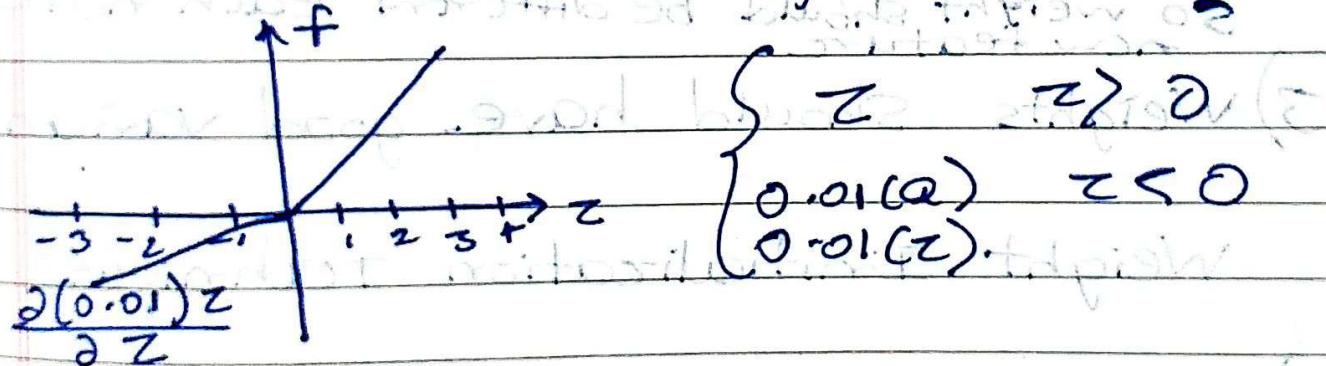
Because of this now  $\frac{\partial L}{\partial w}$  is 1 no vanishing gradient problem & our weights gets converges easily.

Suppose Derivitve Derivative is 0, then we get  $w_{\text{ad}} = w_{\text{new}}$  as  $\frac{\partial L}{\partial w} = 0$ . So it will create dead neuron or dead activation function.

means no weight updation process happens so to handle this problem Leaky Relu use.

### Leaky Relu

In Relu for -ve value of  $z$  o/p function is zero. But in place of zero during -ve value of  $z$  we are adding some value.



$\frac{\partial f}{\partial z} = 0.01$  So when  $z$  is -ve its derivative is 0.01, so it will not be zero & solving the problem of dead neuron

## Weight Initialization

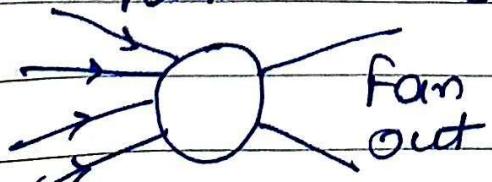
- 1) Weights Should be small (medium small)
- 2) weight should not be small due to which vanishing gradient occurs
- 3) Weights should not be same  
If weights are same the all neurons learn features in the same way.  
So weight should be different each neuron learns new feature.
- 4) Weights should have good variance

## Weight Initialization techniques

- 1) Uniform Distribution  
weights taking from uniform distribution in range  $a$  (lower) &  $b$  (higher)

$$W_{ij} \sim \text{Uniform} \left[ \frac{-1}{\sqrt{\text{fanin}}}, \frac{1}{\sqrt{\text{fanout}}} \right]$$

$\text{IP fanin}$



fanin: No of IPs

fanout: - No of OPs

## 2) Xavier / Gorat

i) Xavier Normal or Gorat Normal weight are taken from normal distribution with mean = 0,  $\sigma^2 = \frac{2}{\text{fan-in} + \text{fan-out}}$  variance

$$w_{ij} \sim N(0, \frac{2}{\text{fan-in} + \text{fan-out}})$$

$$\sigma = \sqrt{\frac{2}{\text{fan-in} + \text{fan-out}}}$$

## 2) Xavier Uniform

Taken from uniform distribution

$$w_{ij} \sim U\left[\frac{-\sqrt{6}}{\sqrt{\text{fan-in} + \text{fan-out}}}, \frac{\sqrt{6}}{\sqrt{\text{fan-in} + \text{fan-out}}}\right]$$

This will work well with Sigmoid function

## 3) He init : Use full for ReLU (RLU)

### i) He Uniform

weight are basically taken from Uniform distribution

$$w_{ij} \sim U\left[-\sqrt{\frac{6}{\text{fan-in}}}, \sqrt{\frac{6}{\text{fan-out}}}\right]$$

### ii) He Normal

$$w_{ij} \sim N(0, \frac{2}{\text{fan-in}})$$

$$\sigma = \sqrt{\frac{2}{\text{fan-in}}}$$