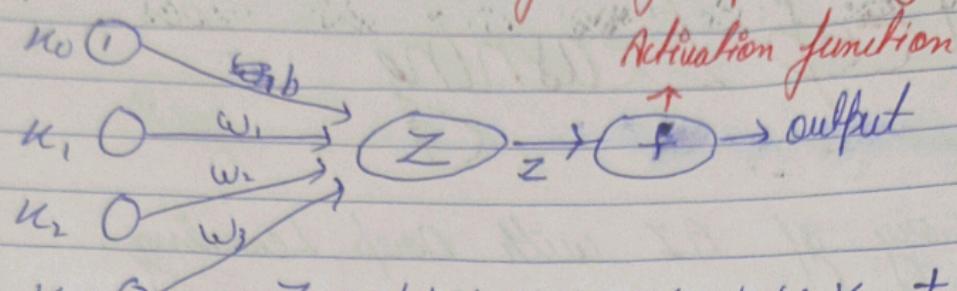


Deep Learning

History of AI with Deep Learning.

- (1) 1940's → Artificial neuron
 - (2) 1950's → Turing Test
 - (3) 1956 → birth of AI
 - (4) 1957 → Perceptron
 - (5) 1960 → ADALINE
 - (6) 1969 → XOR Problem
 - (7) 1980 → Neocognitron
 - (8) 1986 → Back propagation
 - (9) 1989 → VAT
 - (10) 1995 → SVM's
 - (11) 1998 → CNN
 - (12) 2006 → RBM Initialization
 - (13) 2012 → Alex Net
 - (14) 2014 → GAN
 - (15) 2017 → Transformer
 - (16) 2020 → Gpt-3
 - (17) 2022 → Chat-Gpt
- I First Golden Age
- ↑ First Dark Age
- I Second golden age
- ↑ Second dark age
- ↑ Third golden Age

Perception \rightarrow First idea in the field of deep learning.

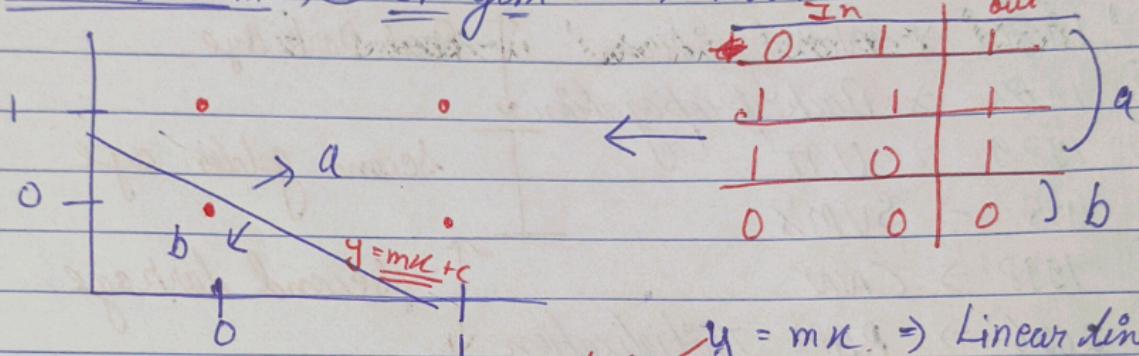


$$z = w_0 K_0 + w_1 K_1 + w_2 K_2 + w_3 K_3 + b$$

* Perception map your input with a linear functions giving us some linear output we can use for train.

\hookrightarrow this is the weight that it gives us

KOR Problem \rightarrow ① OR fun \rightarrow Truth table

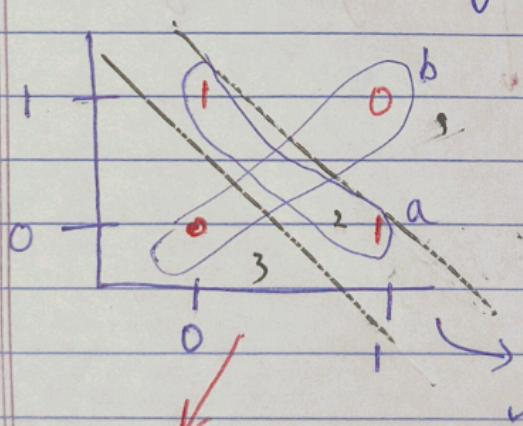


$$y = mx \Rightarrow \text{Linear line}$$

So How can i draw a linear line to classify data points

② KOR fun \rightarrow Truth Table

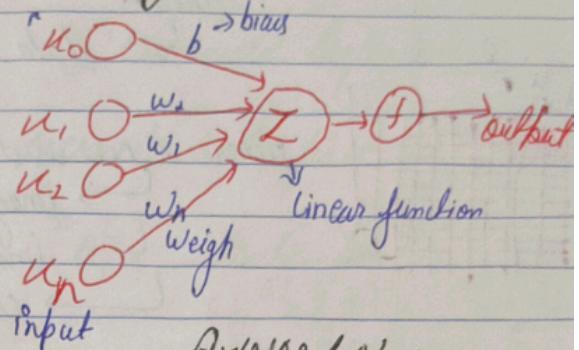
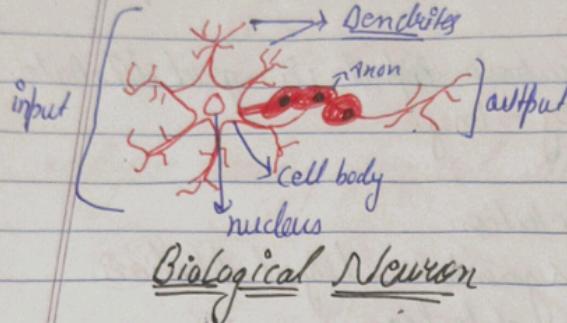
in	out
1 0	1
0 1	1
1 1	0
0 0	0



Here linear function is failing,
linear fun can divide in two classes
only.

If we want to represent a complex problem like NOR we
want some complex eqⁿ

Neural Networks And its types



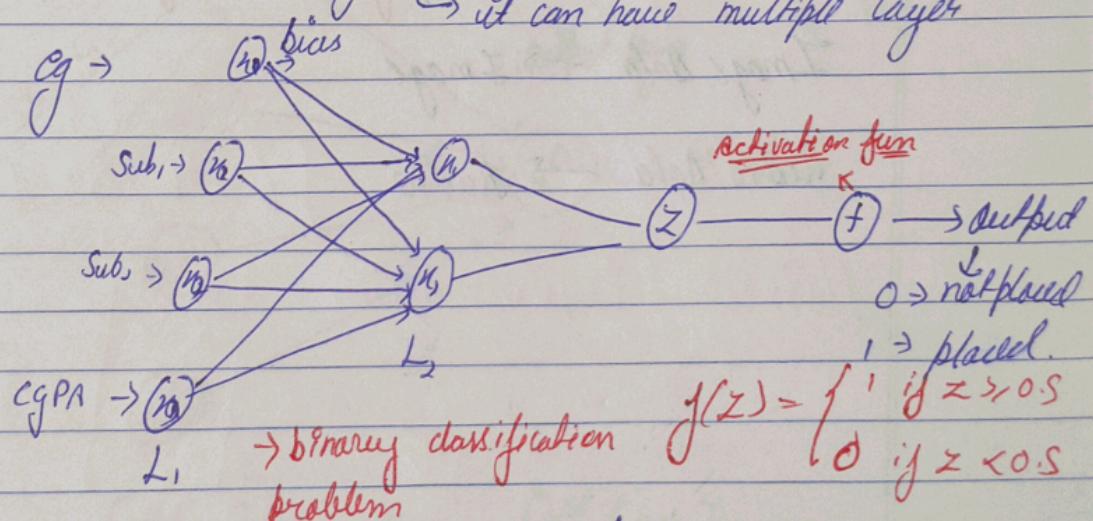
Types of neural network

(i) Feed Forward Neural Network \Rightarrow FNN or FFNN

Perception = Most Basic Feed Forward Neural Network

\hookrightarrow Why basic?

\hookrightarrow it can have multiple layers

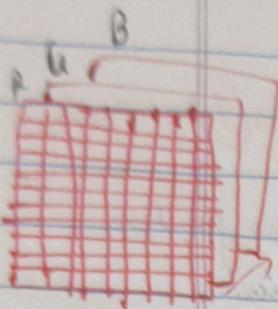


$$z = f(x) + b \quad \left\{ \begin{array}{l} \mathbf{x} = [x_0, u_1, u_2, u_3] \\ b = x_0 \cdot b \end{array} \right.$$

(3)

Convolutional Neural Network (CNN)

→ special Neural Network to handle complex data



where CNN used:

→ image detection

→ image segmentation / object detection

→ cancer detection

(3)

Recurrent Neural Network (RNN)

→ good for sequential data

→ eg → audio, text, speech, Time series

Transformers are also a type of RNN → multilayered RNN

(4)

GAN → Generative Adversarial Network

→ used to generate synthetic / artificial Data

Image Data → Image

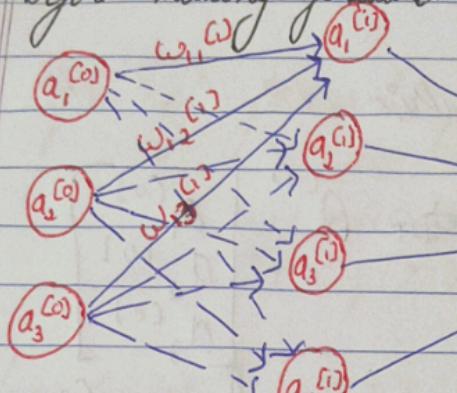
Audio Data → Audio

• Learning Mechanism of Neural Network

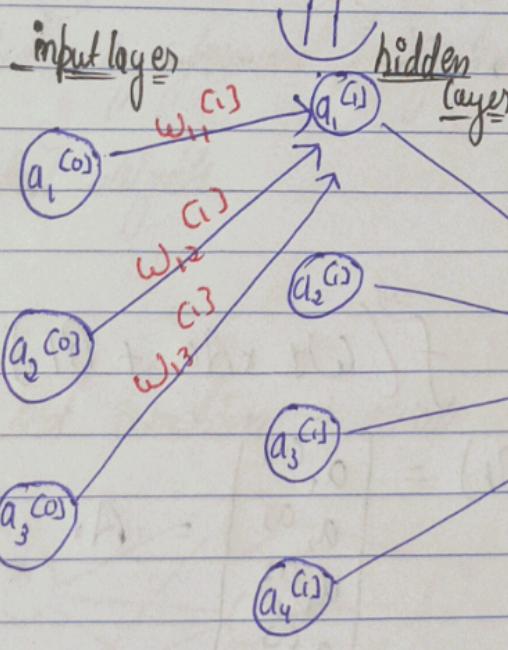
① Forward Propagation

② Backward Propagation

- Before moving forward let's understand this first :-



$$a_i^{(l)} = \text{activation fn} (w_{i,1}^{(l)} \times a_1^{(l-1)} + w_{i,2}^{(l)} \times a_2^{(l-1)} + w_{i,3}^{(l)} \times a_3^{(l-1)} + b_r)$$



$\left[\begin{matrix} 0^{\text{th}} \text{ layer} \\ 1^{\text{st}} \text{ layer} \\ 2^{\text{nd}} \text{ layer} \end{matrix} \right]$

w_{ij} → next layer to which weighted edge is connected

neuron no. of next layer

neuron no. of previous layer.

$a_i^{(l)}$ → representation of neuron
layer no.

no. neuron no.

$$a_1^{(1)} = \text{activation function} \left(\begin{matrix} a_1^{(0)} \\ a_2^{(0)} \\ a_3^{(0)} \\ a_4^{(0)} \end{matrix} \right) \times \begin{pmatrix} w_{11}^{(1)} & w_{12}^{(1)} & w_{13}^{(1)} \\ w_{21}^{(1)} & w_{22}^{(1)} & w_{23}^{(1)} \\ w_{31}^{(1)} & w_{32}^{(1)} & w_{33}^{(1)} \\ w_{41}^{(1)} & w_{42}^{(1)} & w_{43}^{(1)} \end{pmatrix} + b_1^{(1)}$$

This can also be written as

$$a_1^{(1)} = w \times A_0 + B_1$$

w can be written in form of matrix

$$W_1 = \begin{matrix} \text{eq} \\ n-y \\ \downarrow \end{matrix} \left[\begin{matrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \\ w_{41} & w_{42} & w_{43} \end{matrix} \right] \begin{matrix} \text{eq} \\ n-x \\ \leftarrow \quad \rightarrow \quad (4 \times 3) \end{matrix} A^{(0)} \begin{matrix} \text{eq} \\ n-y \\ \downarrow \end{matrix} \left[\begin{matrix} a_1^{(0)} \\ a_2^{(0)} \\ a_3^{(0)} \\ a_4^{(0)} \end{matrix} \right] \begin{matrix} \text{eq} \\ 3 \times 1 \end{matrix}$$

$n-x \rightarrow$ no. of col

$n-y \rightarrow$ no. of rows

Activation fun $\rightarrow f(w_1 \times A_0 + B_1)$

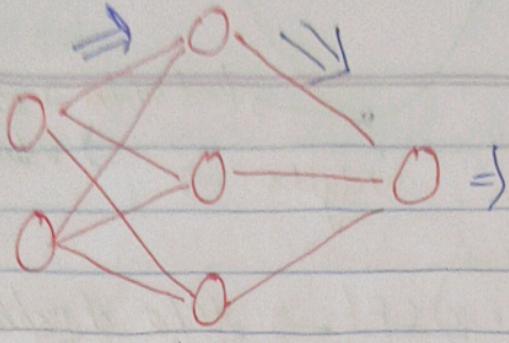
~~$f(w_1 \times A_0 + B_1) = \begin{bmatrix} a_1^{(1)} \\ a_2^{(1)} \\ a_3^{(1)} \\ a_4^{(1)} \end{bmatrix} = A_1$~~

$$A_1 = f(w_1 \times A_0 + B_1)$$

$$A_2 = f(w_2 \times A_1 + B_2)$$

Here $w_0 \rightarrow (w_0)$ is matrix containing all weights associated with first layer. i^{th} layer

$A_f \rightarrow$ matrix containing neurons of j^{th} layer



forward propagation

Date _____
Page No. _____

Here we are going forward by giving random weight at start & updating it from an activation fun.

- then we calculate loss/cost

→ This is also called forward propagation

Eqⁿ of forward propagation

$$Z_1 = w_1 * A_0 + b_1$$

~~$A_1 = f(Z_1)$~~

$$Z_2 = w_2 * A_1 + b_2$$

$$A_2 = f(Z_2)$$

Random weights

$$Z_3 = w_3 * A_2 + b_3$$

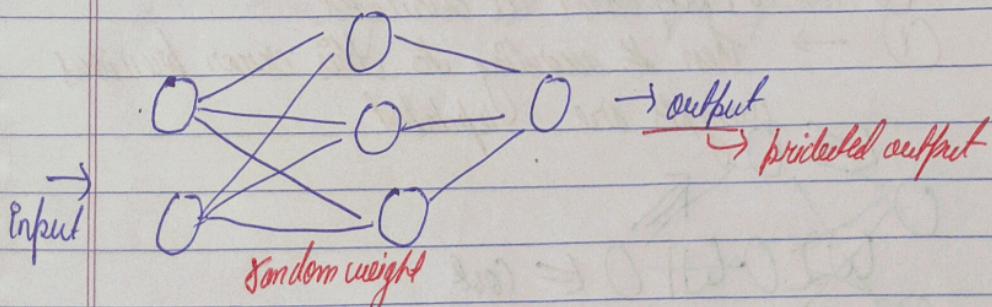
\Downarrow
A₃ = f(Z₃) → output prediction

Trained using forward propagation

\Downarrow

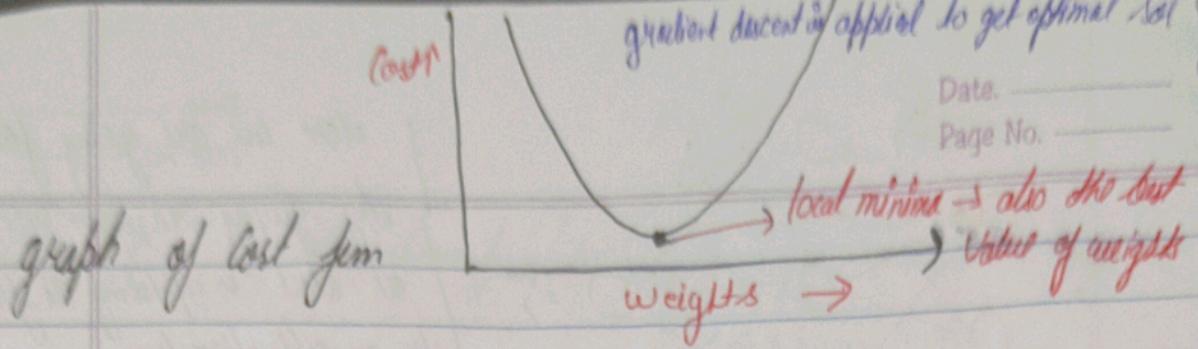
Trained weights.

Cost function → error representation



$$\text{Cost fm} \downarrow = \text{Actual output} - \text{Predicted output}$$

$\left\{ \begin{array}{l} \text{CF} \uparrow \rightarrow \text{output} \rightarrow \text{wrong} \\ \text{CF} \downarrow \rightarrow \text{output} \rightarrow \text{correct} \end{array} \right.$
 We minimize CF to get best output.



$$\omega' = \omega - \alpha \cdot \frac{\partial C_F}{\partial \omega} \quad \begin{matrix} \rightarrow \text{slope of graph} \\ \text{learning rate} \end{matrix}$$

So our Random weights gets updated to optimal weights using gradient descent.

formulae for gradient descent

Repeat {

$$\omega' = \omega - \alpha \cdot \frac{\partial C_F}{\partial \omega}$$

$$\beta' = \beta - \alpha \cdot \frac{\partial C_F}{\partial \beta}$$

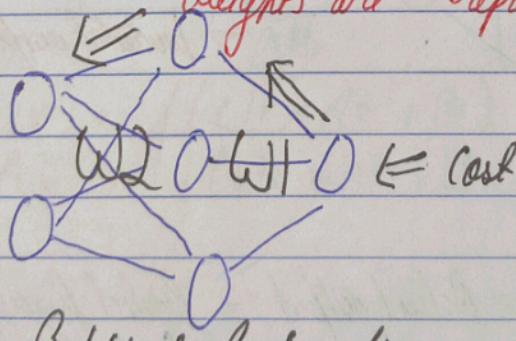
}

How How weights are updated?

↳ They are updated in backward dir.

① → Cost/error is calculated

② → Then according to this error previous weights are updated.



Backward Propagation

How A neural network is trained

- ① → initialise random weights
- ② → Forward Propagation

$$Z = \omega_1 \times A_0 + b_1$$

$$A_2 = f(Z)$$

- ③ → find the cost/error

$$\Delta = \text{pred output} - \text{actual output}$$

$$\delta = \frac{\alpha \cdot \Delta_{\text{cost}}}{\text{no. of samples}}$$

- ④ Backward Propagation

$$\omega' = \omega - \alpha \cdot \frac{\partial C}{\partial \omega}$$

$$b' = b - \alpha \cdot \frac{\partial C}{\partial b}$$

- ⑤ → Repeat ③ & ④ until we get optimal sol.

1 epoch = Step 1 + Step 2 + Step 3 + Step 4
for training we do multiple epochs

* Activation functions \rightarrow

$$Z_1 = \omega_1 \times A_0 + B_1$$

$$A_1 = f(\text{Activation fun}(z_1))$$

if let's say this function is linear and not learning anything but the same old input

$$Z_1 = \omega_1 \times A_0 + B_1$$

$$A_1 = f(z_1)$$

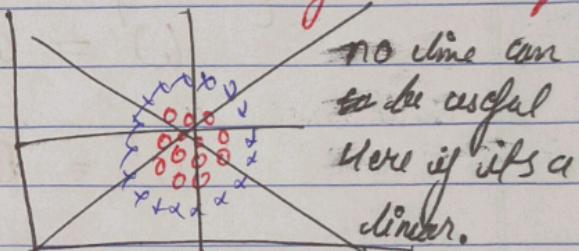
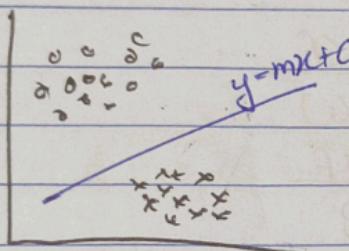
$$Z_2 = \omega_2 \times A_1 + B_2$$

$$A_2 = \omega_2 \times (\omega_1 \times A_0 + B_1) + B_2$$

$$A_2 = \underbrace{\omega_2 \times \omega_1 \times A_0}_{\text{1st}} + \underbrace{\omega_2 \times B_1 + B_2}_{\text{2nd}}$$

w_1

b' This is again a linear eqn



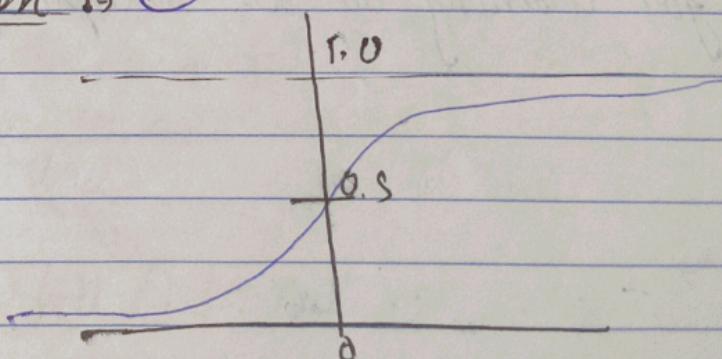
That's why we need a non-linear activation fun.

① Sigmoid Function \rightarrow σ

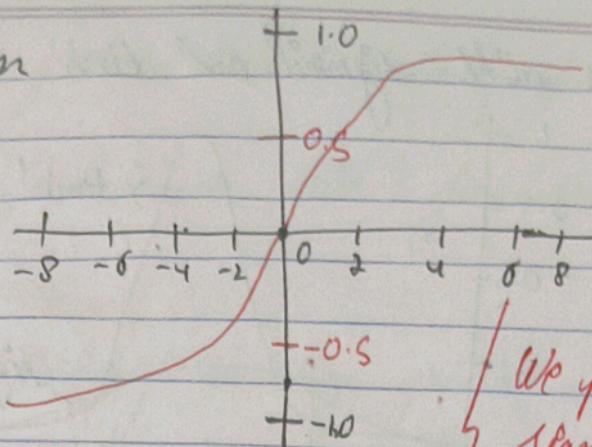
$$\sigma(y) = \frac{1}{1 + e^{-x}}$$

Used \rightarrow

• Binary Classification



② Tanh function

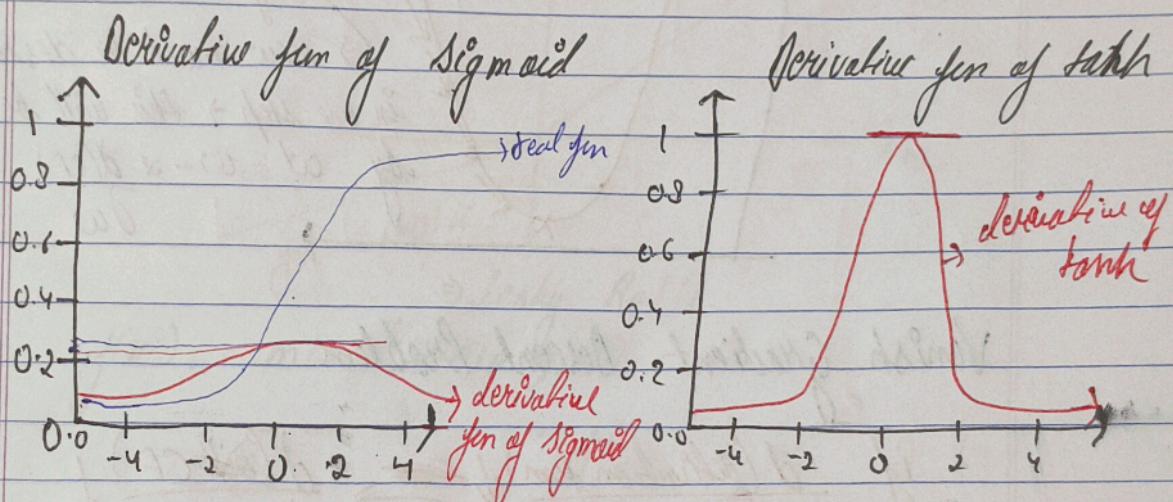


$$\tanh \text{ fun } (y) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

We prefer tanh over sigmoid in the hidden layer because?
→ avoids back

So we also need derivative of Activation fun as well
So the cost fun depends on Activation fun.

propagation → we do derivative of (cost fun) wrt (c)



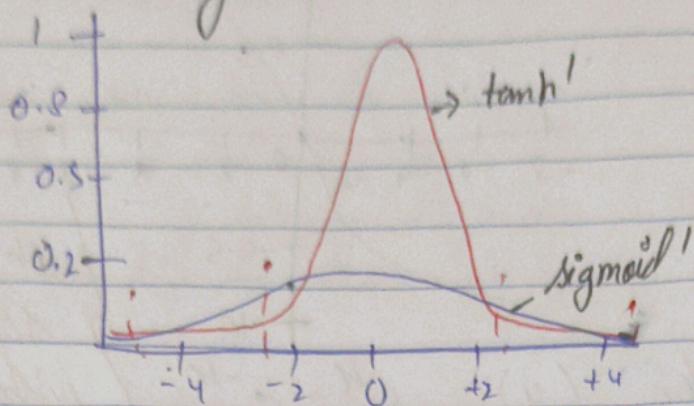
peak value of derivative of sigmoid ≈ 0.2
peak value of derivative of tanh ≈ 1

$$\omega' = \omega - \alpha \times \frac{\partial C_F}{\partial \omega} \rightarrow \text{Here derivative value gets multiplied by the learning rate}$$

Why tanh over sigmoid

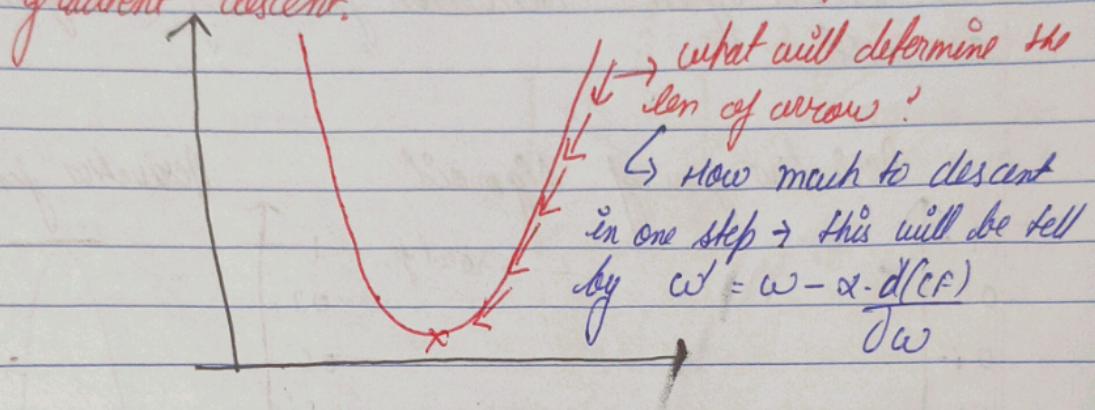
- Centered around zero.
- So how (tanh fun) is more useful for fast learning

Problems with sigmoid and tanh



Problem with these function is as the value of K increases \rightarrow the derivative of both fun get almost zero

\hookrightarrow this arises a problem of vanishing gradient descent.



Vanish Gradient Descent Problem

if $\frac{d(\text{activation fun})}{dw} \rightarrow$ ~~$\frac{d(\text{act CF})}{dw}$~~ \downarrow
(very low) $\quad \quad \quad$ $\frac{d(\text{act CF})}{dw} \downarrow$
(also very low)

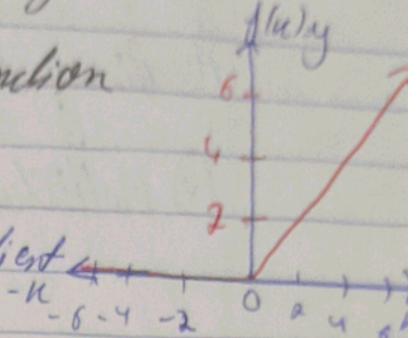
\hookrightarrow So we will move very slow to the minima
 \hookrightarrow This is you both sigmoid and tanh function.

Solution for gradient descent vanishing fun is ReLU

(3) ReLU \rightarrow (Rectified Linear Unit) function

$$f(u) = \max(0, u)$$

solve the problem of vanishing gradient



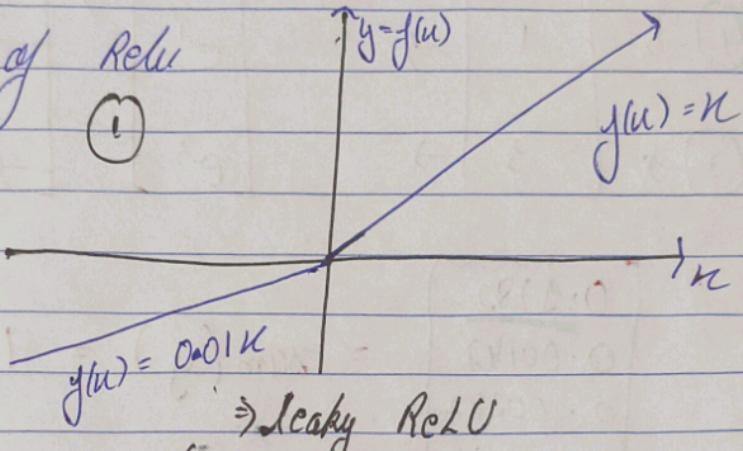
$$\frac{d}{du} f(u) = \begin{cases} 1 & u > 0 \\ 0 & u \leq 0 \end{cases}$$

This function is piece wise linear

↳ gives advantage of both linear - non-linear functions.

Variants of ReLU

(1)

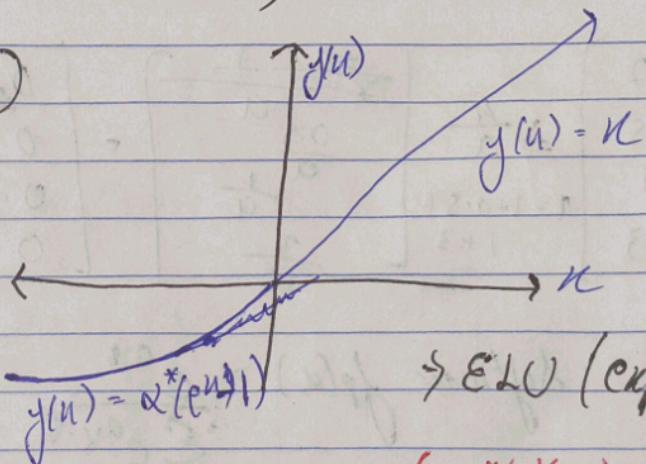


$$y(u) = 0.01u$$

↳ Leaky ReLU

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

(2)



$$y(u) = \alpha * (e^u - 1)$$

↳ ELU (exponential linear unit)

$$f(u) = \max(\alpha * (e^u - 1), u)$$

which one to use when

① Binary Classification \rightarrow Sigmoid

② Multi-class classification \rightarrow ? \hookrightarrow Softmax

④ Softmax function

$$a = e^7 + e^{0.5}, e^1 + e^3$$

$$\begin{array}{l} \stackrel{\geq 7}{\Rightarrow} \stackrel{7}{\rightarrow} \left[\begin{array}{c} e^7 \\ e^{0.5} \\ e^1 \\ e^3 \end{array} \right] \stackrel{\frac{e^7}{e^7+e^{0.5}+e^1+e^3}}{\rightarrow} \\ \stackrel{\geq 0.5}{\Rightarrow} \stackrel{0.5}{\rightarrow} \left[\begin{array}{c} e^7 \\ e^{0.5} \\ e^1 \\ e^3 \end{array} \right] \stackrel{\frac{e^{0.5}}{e^7+e^{0.5}+e^1+e^3}}{\rightarrow} \\ \stackrel{\geq 1}{\Rightarrow} \stackrel{1}{\rightarrow} \left[\begin{array}{c} e^7 \\ e^{0.5} \\ e^1 \\ e^3 \end{array} \right] \stackrel{\frac{e^1}{e^7+e^{0.5}+e^1+e^3}}{\rightarrow} \\ \stackrel{\geq 3}{\Rightarrow} \stackrel{3}{\rightarrow} \left[\begin{array}{c} e^7 \\ e^{0.5} \\ e^1 \\ e^3 \end{array} \right] \stackrel{\frac{e^3}{e^7+e^{0.5}+e^1+e^3}}{\rightarrow} \end{array}$$

$$y = \begin{bmatrix} 0.9782 \\ 0.00142 \\ 0.00147 \\ 0.00171 \end{bmatrix} = \text{sum}(y) = 1$$

$$\begin{bmatrix} 7 \\ 0.5 \\ 1 \\ 3 \end{bmatrix} \stackrel{a=7+0.5+1+3}{\Rightarrow} \frac{7}{a} \begin{bmatrix} \cancel{7} \\ \cancel{0.5} \\ \cancel{1} \\ \cancel{3} \end{bmatrix} \stackrel{\frac{7}{a}}{=} \begin{bmatrix} 0.6086 \\ 0.0435 \\ 0.0861 \\ 0.2608 \end{bmatrix}$$

$$\text{Softmax } f_i(a) = \frac{e^{a_i}}{\sum_k e^{a_k}}$$

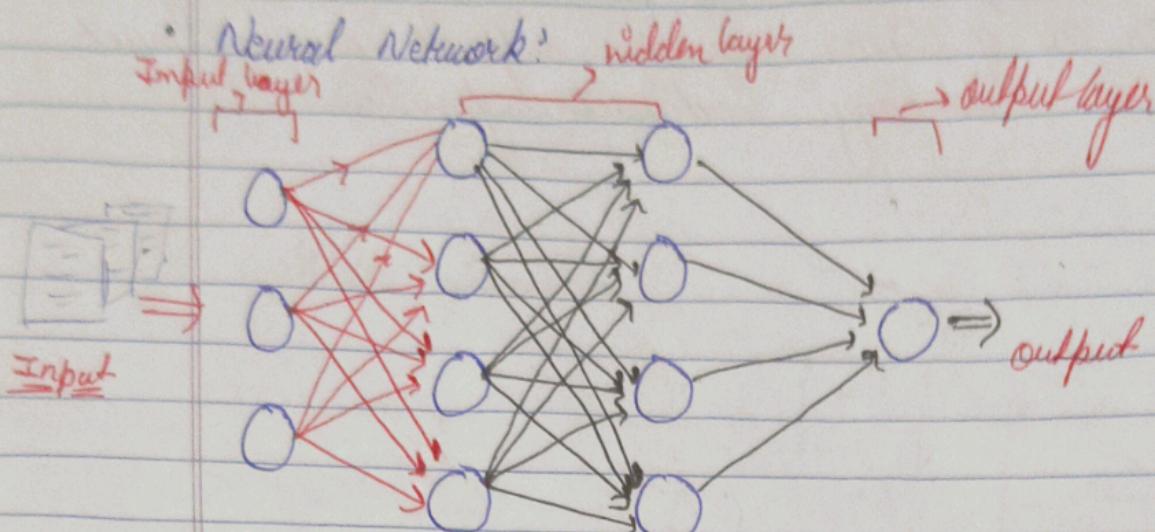
③ House price pred.

→ Your output can be (0 to 1000000) any number

→ So output neuron can use any activation function

In case of linear regression output neuron must
not use ~~any~~ any activation function

Multi Layer Neural Network



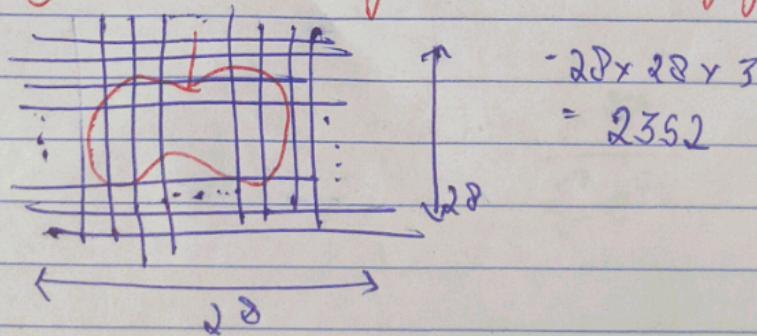
Fully Connected multilayer NN looks like

What is neuron? \rightarrow (it is ~~a function~~ a function that takes input
 Input \rightarrow $f(w)$ \rightarrow output and gives output for that
 \rightarrow Neurons are also called as activation functions.

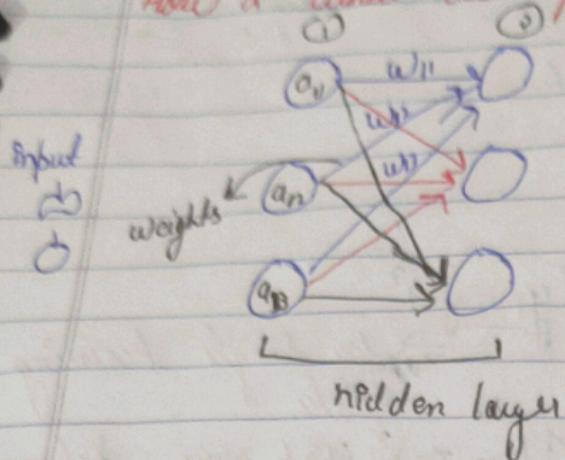
Note \rightarrow Most of the learning happens in hidden layers

How do know the no. of neurons to be used in input layer.

Basic thinking \rightarrow No. of Neuron = No. of features



How & what is happening to the hidden layer.



1st layer has to be responsible for identifying the shape

2nd layer is responsible for identifying the color.

How many layers? → More Neuron = More Layer = More Calculation
= More links = More Complex problem

Connection b/w neurons are called as weight

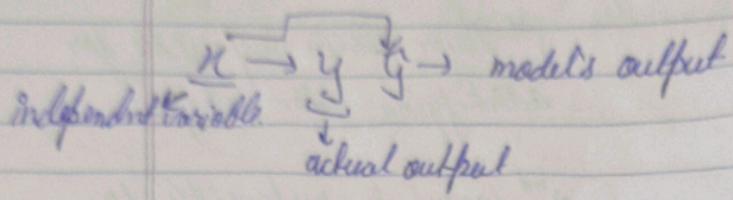
$$a_{21} = \{ w_{11} \times a_{11} + w_{12} \times a_{12} + w_{13} \times a_{13} + b \}$$

\curvearrowright bias

Here we have two values that model learns {weight, bias}

$$a_{21} = \text{fun}(w_{11} \times a_{11} + w_{12} \times a_{12} + w_{13} \times a_{13} + b)$$

Loss Function & Cost Function

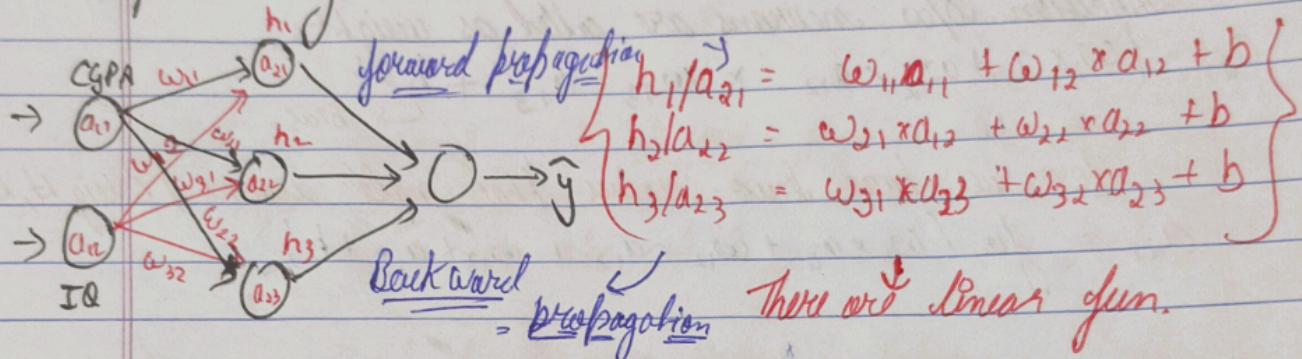


In linear reg. we have \rightarrow error = $(y_i - \hat{y})^2$ for single training data.

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

High error = bad model trained
Low error = good model trained.

in deep learning



There are linear fun.

Motive of training \rightarrow minimise the cost function.

\hookrightarrow "we cannot improve what you can't measure"

Types of loss fun. for regression problems

① MSE $\rightarrow \sum_{i=1}^n (y_i - \hat{y})^2$ • Not robust to outliers

② easy to interpret = $\frac{\partial (y - \hat{y})^2}{\partial w} = \left\{ 2 \frac{\partial y}{\partial w} - 2 \frac{\partial \hat{y}}{\partial w} \right\}$

③ Differentiable used in gradient descent

④ Only one global minimum

$$\textcircled{2} \text{ MAE} = \sum_{i=1}^n |y_i - \hat{y}_i|$$

- ↳ more robust to outliers
- ↳ Not differentiable

$$\textcircled{3} \text{ Huber Loss} \rightarrow \text{Combo of MAE \& MSE } e = y - \hat{y}$$

- ↳ differentiable
- ↳ robust to outliers

$$L = \begin{cases} \frac{1}{2}e^2 & \text{if } |e| \leq \delta \\ \delta(|e| - \frac{1}{2}\delta) & \text{if } |e| > \delta \end{cases}$$

$\delta = \text{threshold}$

Loss/Cost function types for Classification problem.

① for binary classification

① Binary Cross entropy

also called as
log loss in logistic
reg.

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

→ it can be easily differentiate

② for multiple classes. (we use softmax as activation
fun in output neuron)

① Categorical Cross entropy

$$\text{Loss} = \sum_{j=1}^K y_j \log(\hat{y}_j)$$

③ Hinge Loss → used for SVM

$$\hookrightarrow [0, 1]$$

$$L = \max(0, 1 - y \cdot f(u)) \quad w^T u + b = \text{model out.}$$

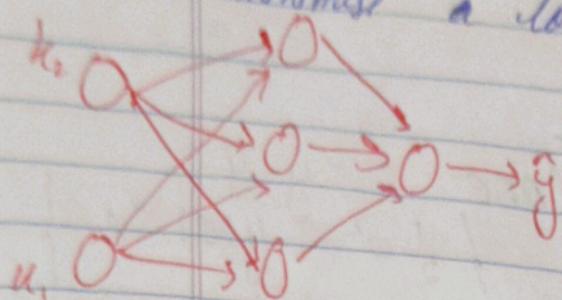
if ($y \cdot f(u) \geq 1$) \rightarrow Loss = 0 \rightarrow confident correct

if ($y \cdot f(u) > 0$ & $y \cdot f(u) < 1$) ($\text{Loss} > 0$) \rightarrow correct output but need improvement

if ($y \cdot f(u) < 0$) \rightarrow wrong pred.

Optimizers

→ optimizers are algorithm that adjusts model's parameters during training to minimize a loss function and improve model's performance.

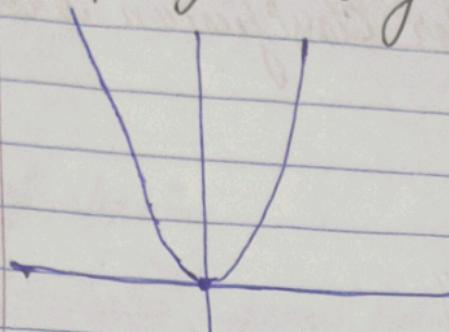


$$\text{Loss} = (y - \hat{y})^2$$

$$h_1, \hat{y} = w_1 u_1 + w_2 u_2 + b_1$$

$$h_2 = w_3 u_3 + w_4 u_4 + b_2$$

Optimizers → ① Gradient Descent



$$w_{\text{new}} = w_{\text{old}} - \delta \frac{\partial \text{Loss}}{\partial w_{\text{old}}}$$

$$b_{\text{new}} = b_{\text{old}} - \delta \frac{\partial \text{Loss}}{\partial b_{\text{old}}}$$

train → 1 epoch → taking entire data. and training for one (w,b)

→ 1 epoch = 100 iteration

iteration → smaller version of epoch. it happens for one single datapoint

→ n iteration → 1 epoch

n epochs → 1 batch

Gradient Descent) will take entire dataset & loss as input and gives better {weights and bias}

- Concerns → ① Computationally expensive.
- ② Higher training time

(2) SGD → Stochastic Gradient Descent

in case of GD

1000 record

100 → Batch

→ 10 iteration = 1 epoch

We propagate back the loss for each iteration of data.

But in SGD → if we have ~~1000~~ SGD ~~and~~ we send
only one clip per iteration & we
calculate loss for that and update weight vector.

→ for each iteration we update weight vector

but it not solve the problem of GD

(3) mini-batch Gradient Descent

1000 records of 100 Batch size

$\frac{1000}{100} \Rightarrow 10$ iteration = 1 epoch

after every batch the value of weights updated ~~not~~
not like GD always it update for entire
data set and not like SGD always it updates
after single data points

Solving problem of GD & SGD

But these GD are also have a problem that
they are not very smooth to reach minima
for that \downarrow we get exponential average

$$V_t = \beta V_{t-1} + (1-\beta) \nabla f$$

momentum optimizer

$$0 < \beta < 1 \rightarrow \beta = 0.9$$

smoothing of descent

(4) Adaptive Gradient Boosting } Adagrad

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

$$\eta = 0.01 \text{ (fixed before training)}$$

Here $\eta \rightarrow$ adaptive / when far from minimum
 η is large
when close to minimum
 η is small

(5) ADAM \rightarrow most effective in industry.

↳ maximum used in real-time in industry

ADAM

↳ momentum + RMS propagation
Better generalizability.

Note: in every optimizer ~~is~~ it is in the updates of parameters.

↳ aim is same to reach the minima as soon as possible.