

# Deep Learning

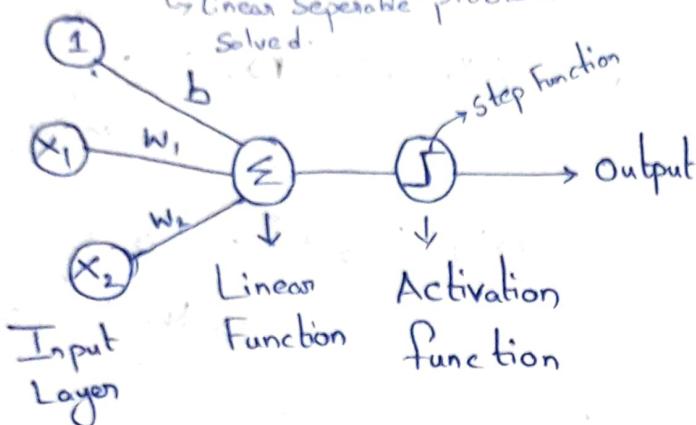
Supervised Algo

Feedforward Neural Network

Perceptron → Binary Classification

↳ Linear Separable Problems can be

Solved.



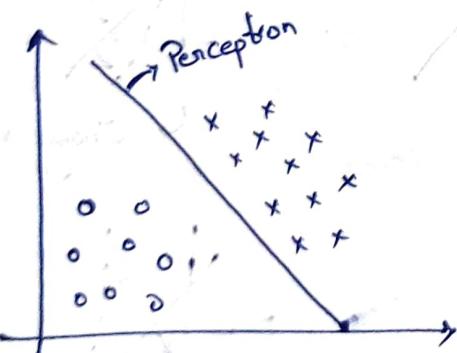
$$Z = w_1 x_1 + w_2 x_2 + b$$

$$Z = \sum_{i=1}^n w_i x_i + b \quad \text{Generalized Fn}$$

$x_1$  &  $x_2$  → Input / Features

$b$  → Bias / Intercept

$w$  → Weight



\* Works best / only on  
Linear Equations, doesn't work  
on non-linear Equations

- Perceptron is nothing but a line.  
It creates regions to classify

$$\hat{y} = f(\mathbf{z}) = \begin{cases} 1 & z \geq 0 \\ 0 & z < 0 \end{cases}$$

↳ Step Function

Activation function

↳ It tries to normalize

(or) Standardise the  $Z$  value. So that the output will be either 0/1 or -1/-1.

↳ It tries to bring the values under a range.

Half a page of notes

\* What does weight Convey?

↳ It tells us which feature is more important.

↳ On which input the output is more dependent.

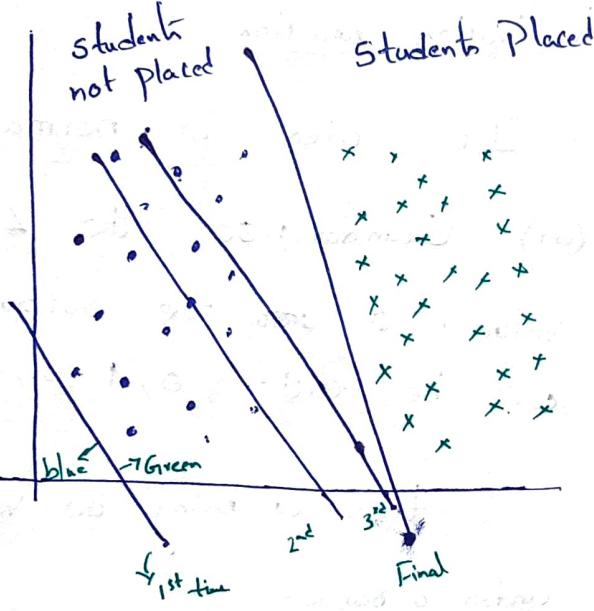
Line (2D)

Plane (3D)

Hyperplane ( $> 3D$ )

# Perception Trick

Epochs  $\rightarrow$  The no. of times a loop runs



What is happening here is that

we run a loop 1000 times (or) again and again and in each loop we select a student in random and we check if the student is in their respective regions with respect to the line.

If the student is not then we move the line to the right place. we keep on doing this till the loop exits (or) till there are no more mismatched points/students.

Hence we are updating  $W$  &  $b$  in every iteration ~~for~~ till we reach right region/place.

$$WX_1 + WX_2 + WX_3 + b$$

$\downarrow$  Resultant Eq

$$AX_1 + BX_2 + CX_3 + b = 0 \rightarrow \text{line Eq}$$

$$\text{if } AX_1 + BX_2 + CX_3 + b > 0$$

$\hookrightarrow$  Positive Region

$$AX_1 + BX_2 + CX_3 + b < 0$$

$\hookrightarrow$  Negative Region

$\hookrightarrow$  Reference  $\rightarrow$  Desmos.com/calculator

~~$$AX_1 + BX_2 + CX_3 + D = 0$$~~

$\downarrow$  Initial line

~~$$AX_1 + BX_2 + CX_3 + D' = 0$$~~

$\hookrightarrow$  If change in  $C$

~~$$AX_1 + BX_2 + CX_3 + D'' = 0$$~~

$\downarrow$  After 1st

~~$$AX_1 + BX_2 + CX_3 + D''' = 0$$~~

$\hookrightarrow$  Change in  $A$

~~$$AX_1 + BX_2 + CX_3 + D'''' = 0$$~~

$\downarrow$  After 2nd

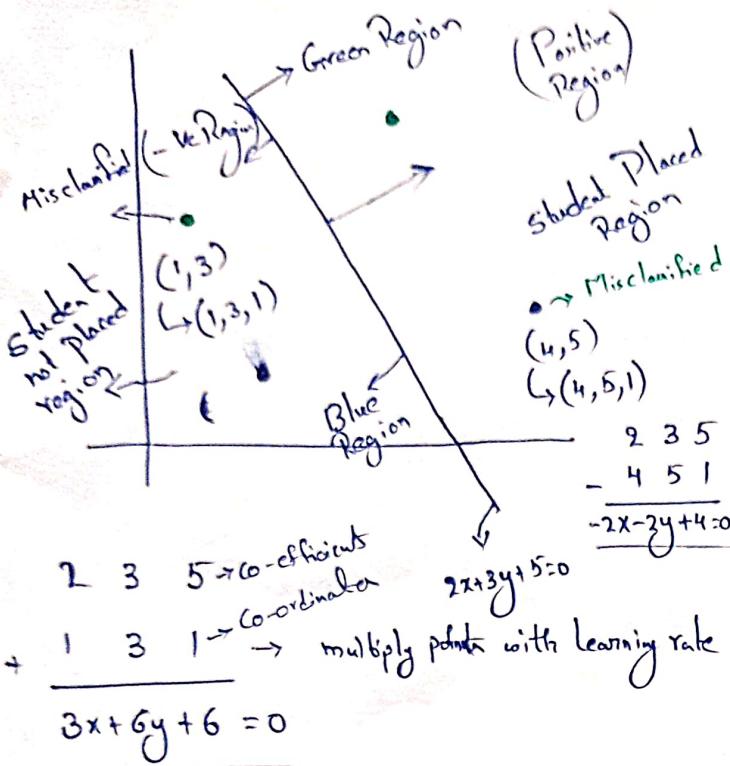
~~$$AX_1 + BX_2 + CX_3 + D''''' = 0$$~~

$\hookrightarrow$  Change in  $B$

$$AX_1 + BX_2 + CX_3 + D'''''' = 0 \rightarrow \text{Line}$$

After substituting the points line, if the result is zero, that means the point lie on the line.

## Transformations



## Formula

$$\text{coeff}_{\text{new}} = \text{coeff}_{\text{old}} - \eta \text{ coordinates}$$

## Algorithm (Pseudo Code)

$$\text{epoch} = 1000, \eta = 0.01$$

for i in range(epoch):

randomly select a student

if  $x_i \in N$  &  $\sum_{i=0}^n w_i x_i \geq 0$ :  $\rightarrow \text{eq}_1$

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

if  $x_i \in P$  &  $\sum_{i=0}^n w_i x_i < 0$ :  $\rightarrow \text{eq}_2$

$$w_{\text{new}} = w_{\text{old}} + \eta x_i$$

$$y \quad \hat{y} \quad y_i - \hat{y}_i$$

$$1 \quad 1 \quad 0$$

$$0 \quad 0 \quad 0$$

$$1 \quad 0 \quad 1$$

$$0 \quad 1 \quad -1$$

+ Let Simplify the Algo

Instead of having two eqns let have one eqn

$$w_{\text{new}} = w_{\text{old}} + \eta (y - \hat{y}) x_i$$

Case(i)  $y = 1, \hat{y} = 1 \rightarrow \text{Matched/Proper classification}$

$$w_{\text{new}} = w_{\text{old}} + \eta (1 - 1) x_i$$

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

Case(ii)  $y = 0, \hat{y} = 0 \rightarrow \text{Proper classification}$

$$w_{\text{new}} = w_{\text{old}} + \eta (0 - 0) x_i$$

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

Case(iii)  $y = 1, \hat{y} = 0 \rightarrow \text{Misclassified}$

$$w_{\text{new}} = w_{\text{old}} + \eta (1 - 0) x_i$$

$$w_{\text{new}} = \text{Some new value}$$

Case(iv)  $y = 0, \hat{y} = 1 \rightarrow \text{Misclassified}$

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

## Loss Function

→ It is a mathematical function that measures how well a model's predictions match the actual outcomes.

## In Simple Terms:-

→ It tells us how wrong your model.

Small Loss → Better model is performing

# Perceptron    Loss    Function

↳ sklearn → SGD

$$L(w_1, w_2, b) = \frac{1}{n} \sum_{i=1}^n \max(0, -y_i f(x_i))$$

$\hookrightarrow w_1 x_{i1} + w_2 x_{i2} + b$

$n \rightarrow$  no. of rows

$y_i \rightarrow$  Output/target

$x_1$	$x_2$	$y$
$x_{11}$	$x_{12}$	$y_1$
$x_{21}$	$x_{22}$	$y_2$

$$f(x_i) = w_1 x_{i1} + w_2 x_{i2} + b$$

$\therefore x_{ij}$   
row col

$$\max(0, -y_i f(x_i))$$

$$= \max(0, x)$$

$x = -y_i f(x_i)$

$x \geq 0 \quad x = x \text{ value}$

$x < 0 \quad x = 0$

$$L = \frac{1}{2} \left[ \max(0, -y_1 f(x_1)) + \max(0, -y_2 f(x_2)) \right]$$

↳ Since we have only two points

$y$	$y$	$\max(0, -y_i f(x_i))$
1	1	$\max(0, -ve) = 0$
-1	-1	$\max(0, -ve) = 0$
1	-1	$\max(0, +ve) = y_i f(x_i)$
-1	1	$\max(0, +ve) = y_i f(x_i)$

We need to find the value of

$$f(x_i)$$

$$f(x_i) = w_1 x_{i1} + w_2 x_{i2} + b$$

We need to find  $w_1, w_2, b$

## Gradient Descent

$$w_1 = w_1 + \eta \frac{\delta L}{\delta w_1}$$

$$w_2 = w_2 + \eta \frac{\delta L}{\delta w_2}$$

$$b = b + \eta \frac{\delta L}{\delta b}$$

$$\frac{\delta L}{\delta w_1} = \frac{\delta L}{\delta f(x_i)} \times \frac{\delta f(x_i)}{\delta w_1}$$

$$\frac{\delta L}{\delta f(x_i)} = \begin{cases} 0 & \text{if } y_i f(x_i) \geq 0 \\ -y_i & \text{if } y_i f(x_i) < 0 \end{cases} \quad \frac{\delta f(x_i)}{\delta w_1} = x_{i1}$$

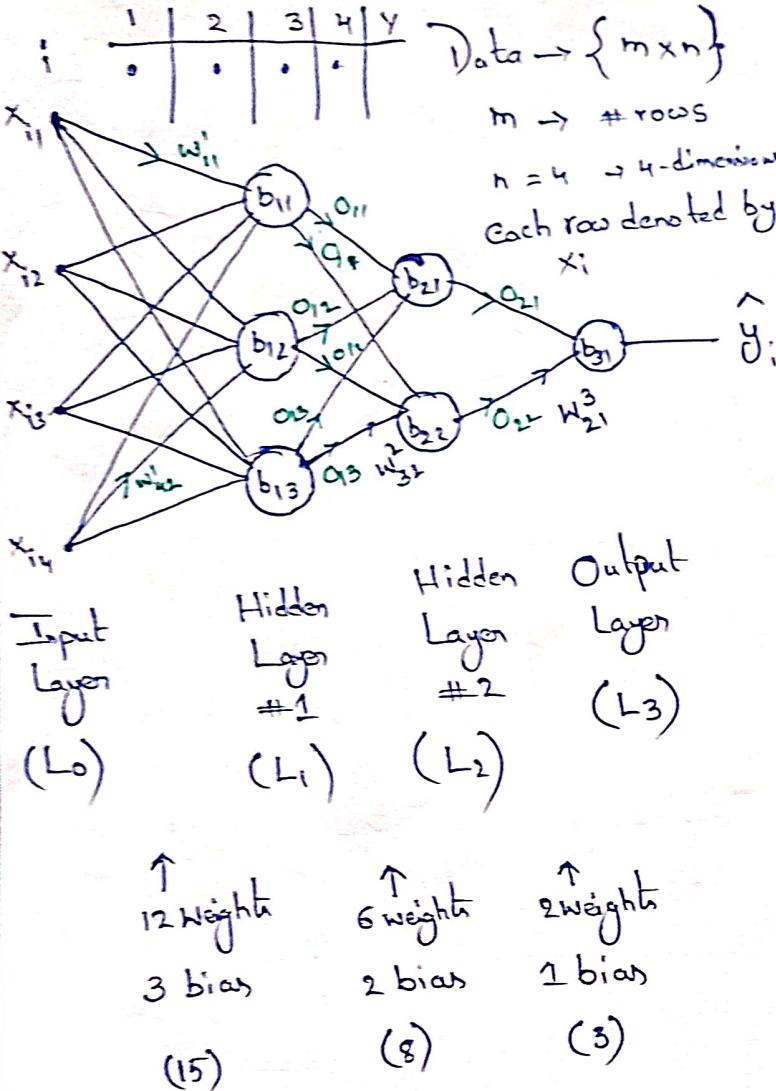
$$\frac{\delta L}{\delta w_1} = \begin{cases} 0 & \text{if } y_i f(x_i) \geq 0 \\ -y_i x_{i1} & \text{if } y_i f(x_i) < 0 \end{cases}$$

$$\frac{\delta L}{\delta w_2} = \begin{cases} 0 & \text{if } y_i f(x_i) \geq 0 \\ -y_i x_{i2} & \text{if } y_i f(x_i) < 0 \end{cases}$$

$$\frac{\delta L}{\delta b} = \begin{cases} 0 & \text{if } y_i f(x_i) \geq 0 \\ -y_i & \text{if } y_i f(x_i) < 0 \end{cases}$$

<u>Loss function</u>	<u>Activation function</u>	<u>Output</u>
Hinge Loss	Step	Act as Perception $\rightarrow$ binary Classification (0,1)
Log loss (binary Cross Entropy)	Sigmoid	Logistic Regression $\rightarrow$ binary classification (0,1)
Categorical Cross Entropy	Softmax	Softmax Regression $\rightarrow$ Probability
MSE	Linear	Linear Regression $\rightarrow$ Number

# Multi Layer Perceptron

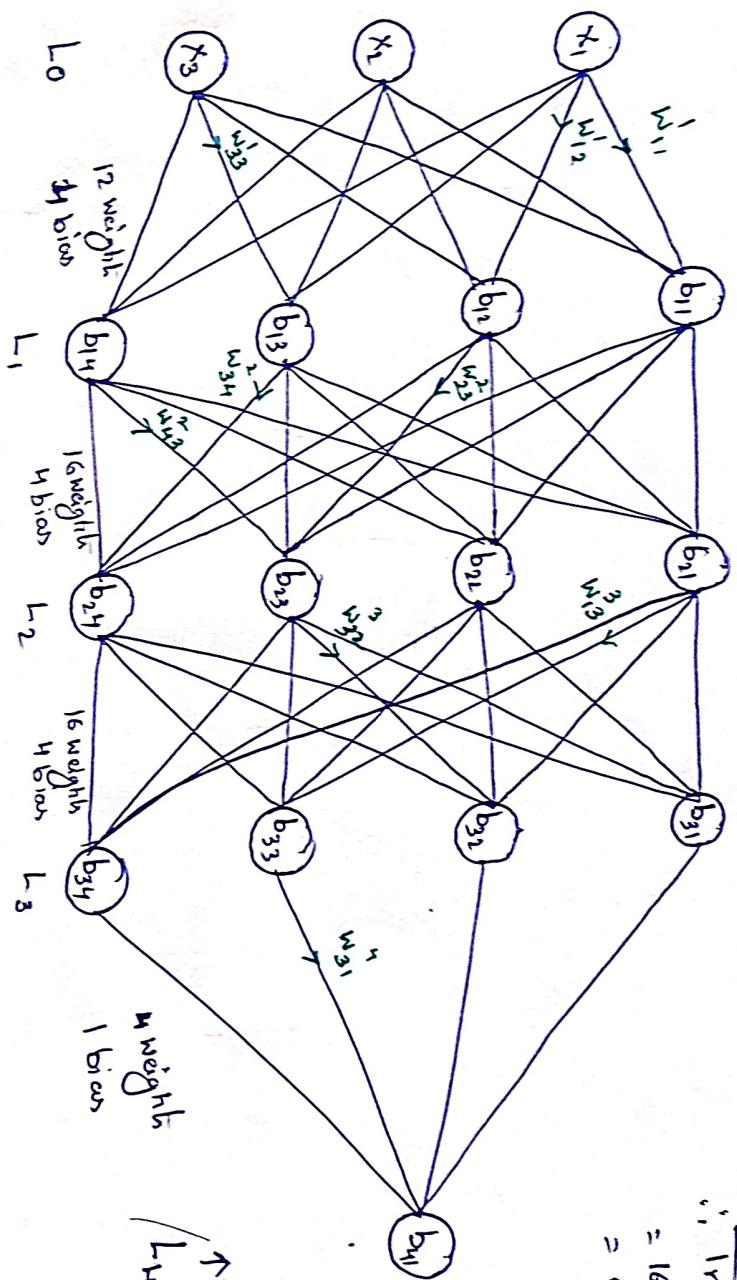


When you train the neural network it tries to find these 26 parameters.

$b_{ij} \rightarrow \text{bias}$ ,  $i \rightarrow \text{layer}$ ,  $j \rightarrow \text{node}$

Output  $\rightarrow O_{ij}$

$w_{ij}^k \rightarrow \text{Weight}$   
 $k \rightarrow \text{which layer (current layer)}$   
 $i \rightarrow \text{Originating node}$   
 $j \rightarrow \text{Destination node}$



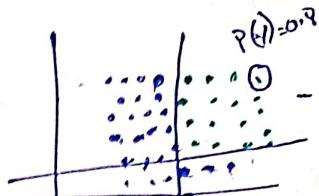
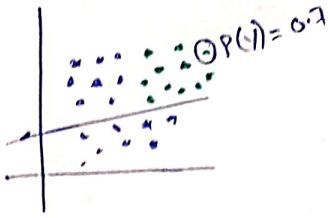
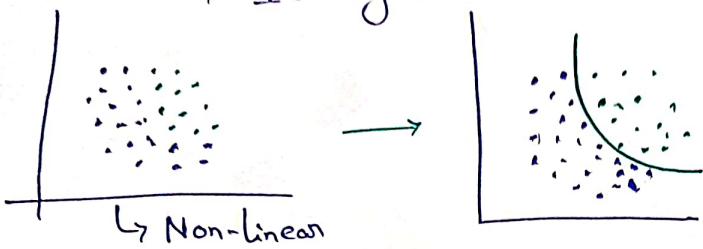
# Multi Layer Perception Intuition

↳ Tensorflow Playground

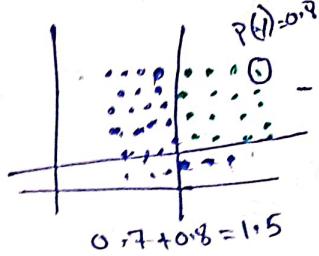
## Perception with Sigmoid

### Problem with Perceptron

↳ It only works on linear Data



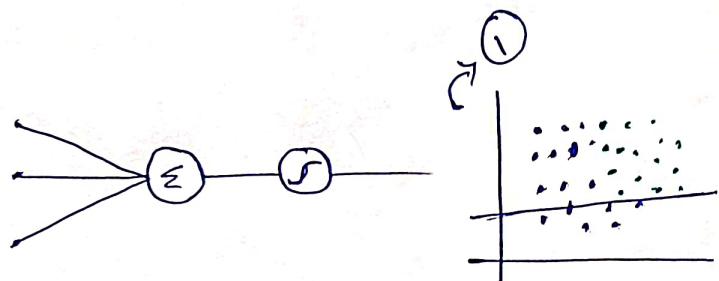
=



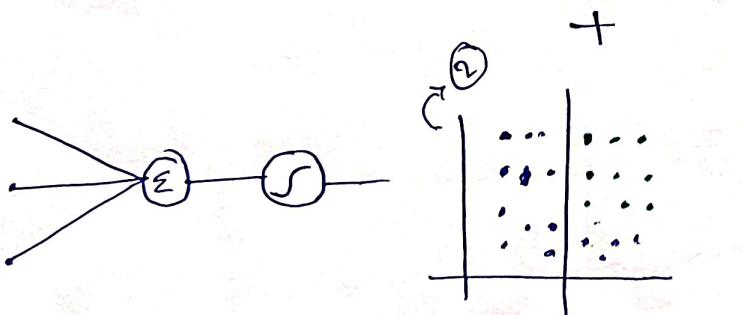
$$\sigma(z) = \frac{1}{1+e^{-z}}$$

$$= 0.8$$

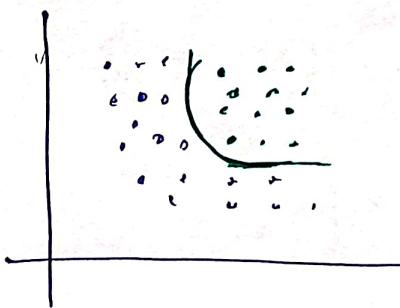
## Perception with Sigmoid



We are summing doing Linear Combination  
of two perceptions to find out  
new probability.

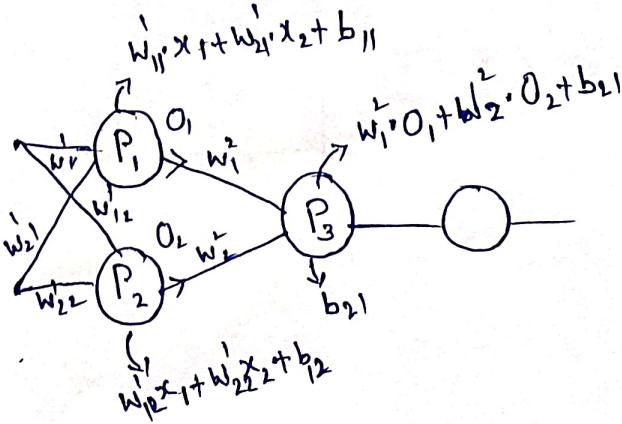
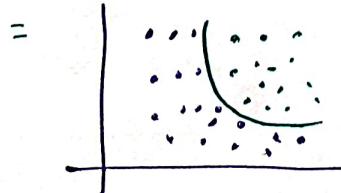


### Result

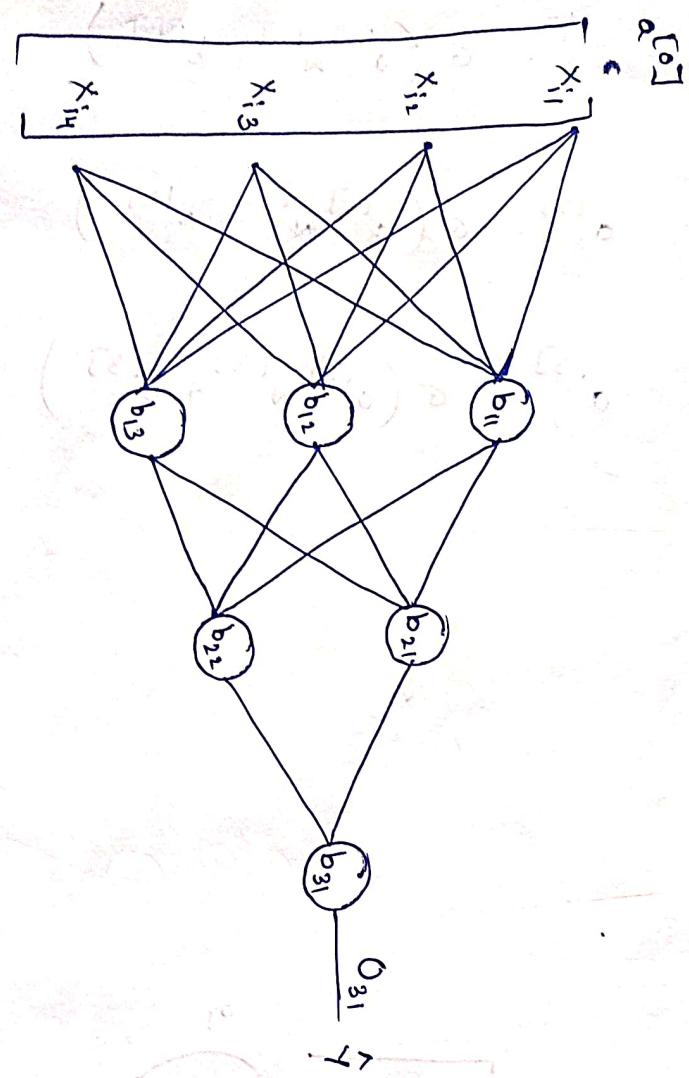


$$(1) + (2) = \begin{array}{|c|c|} \hline & + \\ \hline \end{array}$$

Smoothing



## Forward Propagation



## Layer 1

$$\begin{array}{c}
 \xrightarrow{\quad} \\
 \left[ \begin{array}{ccc} w_{11}^1 & w_{12}^1 & w_{13}^1 \\ w_{21}^1 & w_{22}^1 & w_{23}^1 \\ w_{31}^1 & w_{32}^1 & w_{33}^1 \\ w_{41}^1 & w_{42}^1 & w_{43}^1 \end{array} \right] \xleftarrow{\quad} \\
 \left[ \begin{array}{c} x_{i1} \\ x_{i2} \\ x_{i3} \\ x_{i4} \end{array} \right] + \left[ \begin{array}{c} b_{11} \\ b_{12} \\ b_{13} \end{array} \right]
 \end{array}$$

Dimensions:  $4 \times 3$  (Weights) and  $4 \times 1$  (Inputs)

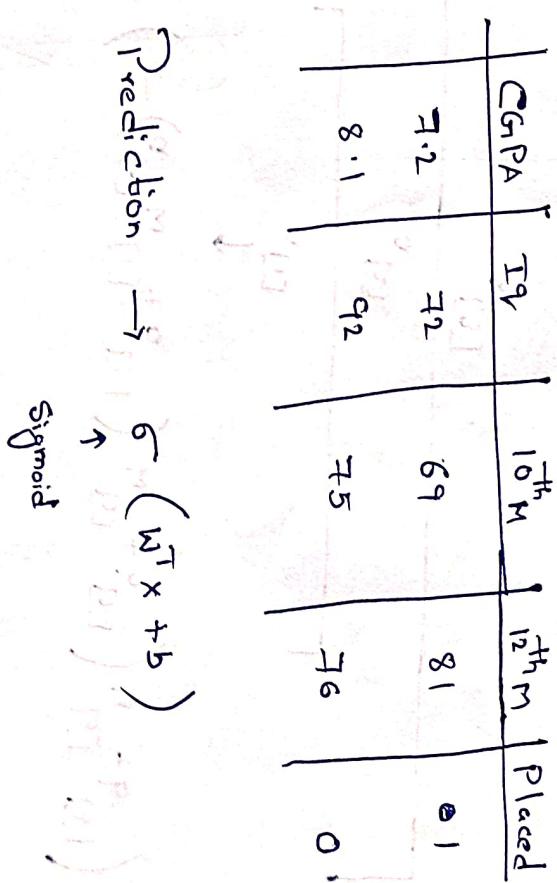
Row wise → Weights Originating from Single node  
Eg:  $x_{i1}$

Column wise → Weights Converging to a Single node  
Eg:  $b_{11}$

$$\Rightarrow \left[ \begin{array}{l} w_{11}^1 x_{i1} + w_{21}^1 x_{i2} + w_{31}^1 x_{i3} + w_{41}^1 x_{i4} + b_{11} \\ w_{12}^1 x_{i1} + w_{22}^1 x_{i2} + w_{32}^1 x_{i3} + w_{42}^1 x_{i4} + b_{12} \\ w_{13}^1 x_{i1} + w_{23}^1 x_{i2} + w_{33}^1 x_{i3} + w_{43}^1 x_{i4} + b_{13} \end{array} \right]$$

$$\Rightarrow \left[ \begin{array}{c} 6 \\ 6 \\ 6 \end{array} \right]$$

$$\Rightarrow \left[ \begin{array}{c} 0_{11} \\ 0_{12} \\ 0_{13} \end{array} \right]$$



Layer 2

$$\begin{aligned}
 & \left[ \begin{array}{cc} w_{11}^2 & w_{12}^2 \\ w_{21}^2 & w_{22}^2 \\ w_{31}^2 & w_{32}^2 \end{array} \right]^T \left[ \begin{array}{c} o_{11} \\ o_{12} \\ o_{13} \end{array} \right] + \left[ \begin{array}{c} b_{21} \\ b_{22} \end{array} \right], \\
 & \quad 3 \times 2 \xrightarrow{T} 2 \times 3 \quad 3 \times 1 \\
 & = \sigma \left( \begin{bmatrix} w_{11}^2 o_{11} + w_{21}^2 o_{12} + w_{31}^2 o_{13} + b_{21} \\ w_{12}^2 o_{11} + w_{22}^2 o_{12} + w_{32}^2 o_{13} + b_{22} \end{bmatrix} \right) \\
 & = \begin{bmatrix} o_{21} \\ o_{22} \end{bmatrix}
 \end{aligned}$$

$$a^{[1]} = \sigma(a^{[0]} w^{[1]} + b^{[1]})$$

$$a^{[2]} = \sigma(a^{[1]} w^{[2]} + b^{[2]})$$

$$a^{[3]} = \sigma(a^{[2]} w^{[3]} + b^{[3]})$$

Layer 3

$$\begin{aligned}
 & \left[ \begin{array}{c} w_{11}^3 \\ w_{21}^3 \end{array} \right]^T \left[ \begin{array}{c} o_{21} \\ o_{22} \end{array} \right] + \left[ \begin{array}{c} b_{31} \end{array} \right], \\
 & \quad 2 \times 1 \xrightarrow{T} 2 \times 2 \quad 2 \times 1 \\
 & = \sigma \left( \begin{bmatrix} w_{11}^3 o_{21} + w_{21}^3 o_{22} + b_{31} \end{bmatrix} \right)
 \end{aligned}$$

$$y_i = o_{31}$$

$$\begin{array}{c}
 \sigma \left( \sigma \left( \sigma \left( a^{[0]} w^{[1]} + b^{[1]} \right) w^{[2]} + b^{[2]} \right) w^{[3]} + b^{[3]} \right) \\
 \uparrow \\
 a^{[1]} \\
 a^{[2]} \\
 a^{[3]}
 \end{array}$$

## Loss Function

- It is a method of evaluating how well your algorithm is modelling your dataset.
- Measures performance of your algorithm.

High  $\rightarrow$  Poor Performance

Low  $\rightarrow$  Great Performance

### MSE

$$L = \sum_i (y_i - \hat{y}_i)^2$$

$y_i = mx_i + c$

$$L(m, c) = \min \sum_i (y_i - mx_i - c)^2$$

### Why Loss fn is important?

You can't improve what you can't measure.

## Loss Functions in Deep Learning

### Regression

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

### Classification

- (i) Binary Cross Entropy
- (ii) Categorical Cross Entropy
- (iii) Hinge loss

## Loss Function Vs Cost Function

### Error Function

Loss Function  $\rightarrow$  Single Training Example

$$\Rightarrow (y_i - \hat{y}_i)^2$$

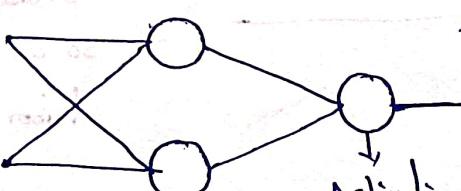
Cost Function  $\rightarrow$  Complete Training Data

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

### ① MSE / Squared Loss / L2 Loss in DL

$$\Rightarrow (y_i - \hat{y}_i)^2$$

↳ Not robust for outliers



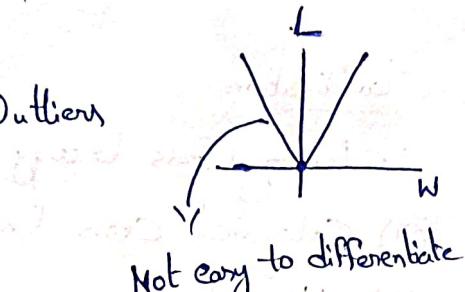
↳ should be linear

② Mean Absolute Error  $\rightarrow$  L1 Loss

$$L = |y_i - \hat{y}_i|$$

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

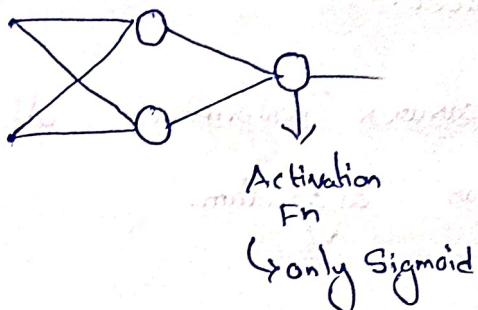
$\hookrightarrow$  Robust for Outliers



④ Binary Cross Entropy / Log loss

$$L = -y \log \hat{y} - (1-y) \log (1-\hat{y})$$

$$CF = \frac{1}{n} \left[ \sum_{i=1}^n y \log \hat{y} + (1-y) \log (1-\hat{y}) \right]$$

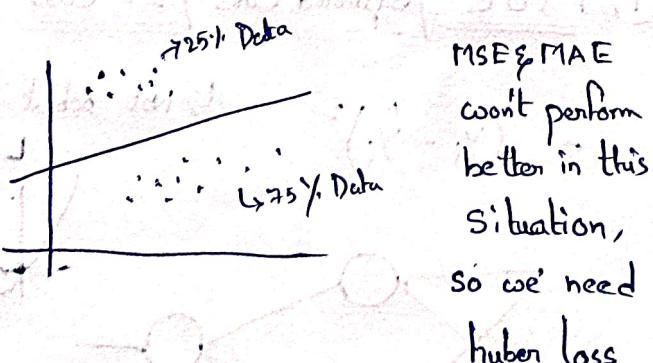
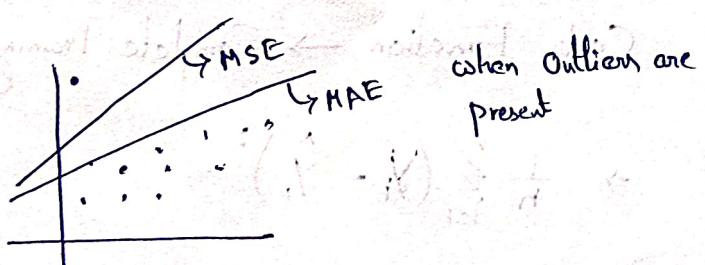


⑤ Huber Loss

$$L = \begin{cases} \frac{1}{2} (y - \hat{y})^2 & \text{for } |y - \hat{y}| \leq \delta \\ \delta |y - \hat{y}| - \frac{1}{2} \delta^2 & \text{otherwise} \end{cases}$$

$\hookrightarrow$  (MAE)

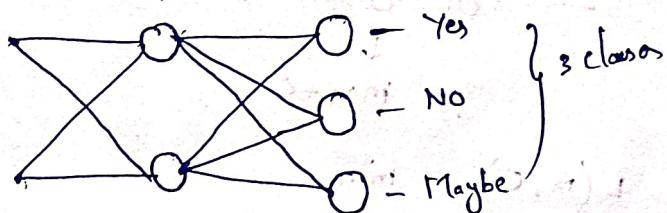
no outliers (MSE)  
Hyper parameter



⑤ Categorical Cross Entropy

$$L = - \sum_{j=1}^K y_i \log (\hat{y}_i)$$

$$CF = -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^K y_{ij} \log (\hat{y}_{ij})$$



$\downarrow$   
Activation  
 $F_n$   
 $\downarrow$   
softmax

- \* One Hot Encoding is used in this.

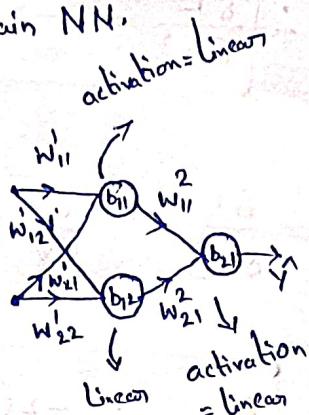
⑥ Sparse Categorical Cross Entropy

$\hookrightarrow$  Similar to Categorical Cross Entropy but One Hot Encoding is not used, instead Integer Encoding is used.

## Back Propagation

↳ An algorithm to train NN.

iq	CGPA	LPA
80	8	3
60	9	5
70	5	8



$$\hat{Y} = O_{21}$$

$$O_{21} = w_{11}^2 O_{11} + w_{21}^2 O_{12} + b_{21}$$

↳ From forward Propagation

Steps

0.) Initialize  $W, b$  Randomly  
 $W=1, b=0$

1.) You select a point (row)  
↳ student

2.) Predict ( $\hat{y}_{pq}$ )  $\rightarrow$  Forward Prop (DotProduct)

3.) Choose a Loss Function:

↳ it's regression problem  $\rightarrow$  MSE

$$L = (\hat{y}_i - y_i)^2$$

$$= (5 - 18)^2$$

↳ let  $\hat{y} = 18$

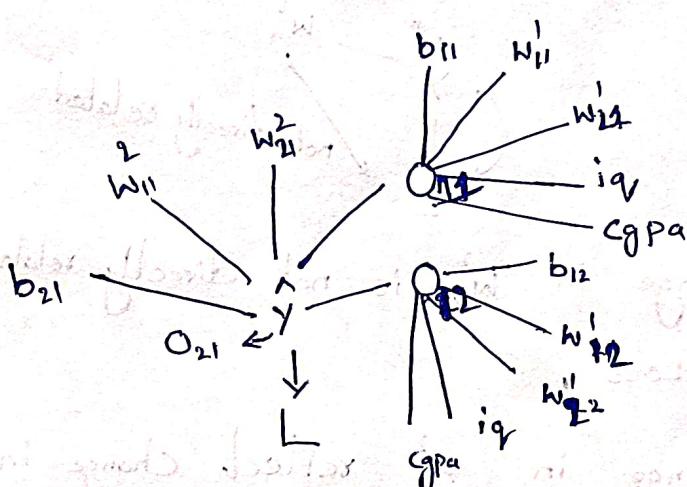
$$L = \text{High}$$

To reduce loss, we need to:

decrease / increase  $\hat{y}$

if  $\hat{y} = 18, Y = 5 \downarrow$

if  $\hat{y} = 1, Y = 5 \uparrow$



4.) Update Weights & biases

↳ Gradient Descent

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

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

$$w_{11}^2_{\text{new}} = w_{11}^2_{\text{old}} - \eta \frac{\delta L}{\delta w_{11}^2_{\text{old}}}$$

$$w_{21}^2_{\text{new}} = w_{21}^2_{\text{old}} - \eta \frac{\delta L}{\delta w_{21}^2_{\text{old}}}$$

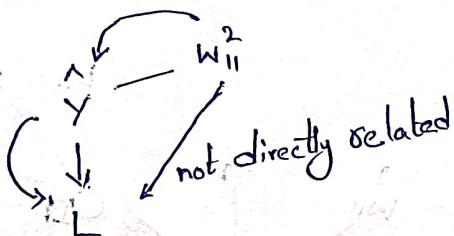
$$b_{21}^2_{\text{new}} = b_{21}^2_{\text{old}} - \eta \frac{\delta L}{\delta b_{21}^2_{\text{old}}}$$

$$\begin{array}{|c|c|} \hline ① - O_{21} & ② - O_{11} \\ \frac{\delta L}{\delta w_{11}^2}, \frac{\delta L}{\delta w_{21}^2}, \frac{\delta L}{\delta b_{21}} & \left| \frac{\delta L}{\delta w_{11}^1}, \frac{\delta L}{\delta w_{21}^1}, \frac{\delta L}{\delta b_{11}} \right. \\ \hline \end{array}$$

$$\therefore \frac{\delta L}{\delta w_{11}^2} = -2(y - \hat{y}) O_{11} \rightarrow 1$$

$$\frac{\delta L}{\delta w_{12}^1}, \frac{\delta L}{\delta w_{22}^1}, \frac{\delta L}{\delta b_{12}} \rightarrow O_{12}$$

We need to calculate these 9 derivatives



Change in  $w_{11}^2$  is not directly related to Loss

Change in  $w_{11}^2$  reflects change in

$\hat{y}$  & Change in  $\hat{y}$  results in

Change in  $L$  due to change in Loss w.r.t  $\hat{y}$

$$\therefore \frac{\delta L}{\delta w_{11}^2} = \boxed{\frac{\delta L}{\delta \hat{y}}} \times \boxed{\frac{\delta \hat{y}}{\delta w_{11}^2}} \rightarrow \text{change in } \hat{y} \text{ w.r.t change in } w_{11}^2$$

↳ Chain Rule of differentiation

$$\frac{\delta L}{\delta \hat{y}} = \frac{\delta}{\delta \hat{y}} (y - \hat{y})^2 = -2(y - \hat{y})$$

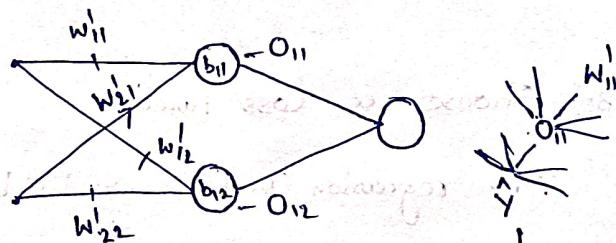
$$\frac{\delta \hat{y}}{\delta w_{11}^2} = \frac{\delta}{\delta w_{11}^2} [O_{11} w_{11}^2 + O_{12} w_{21}^2 + b_{21}] \\ = O_{11}$$

$$\frac{\delta L}{\delta w_{11}^2} = \frac{\delta L}{\delta \hat{y}} \times \frac{\delta \hat{y}}{\delta w_{11}^2}$$

$$\frac{\delta L}{\delta w_{11}^2} = -2(y - \hat{y}) O_{11} \rightarrow 2$$

$$\frac{\delta L}{\delta b_{21}} = \frac{\delta L}{\delta \hat{y}} \times \frac{\delta \hat{y}}{\delta b_{21}}$$

$$\frac{\delta L}{\delta b_{21}} = -2(y - \hat{y}) \rightarrow 3$$



$$\frac{\delta L}{\delta w_{11}^1} = \frac{\delta L}{\delta \hat{y}} \times \frac{\delta \hat{y}}{\delta O_{11}} \times \frac{\delta O_{11}}{\delta w_{11}^1}$$

$$\frac{\delta L}{\delta w_{21}^1} = \frac{\delta L}{\delta \hat{y}} \times \frac{\delta \hat{y}}{\delta O_{11}} \times \frac{\delta O_{11}}{\delta w_{21}^1} \rightarrow O_{11}$$

$$\frac{\delta L}{\delta b_{11}} = \frac{\delta L}{\delta \hat{y}} \times \frac{\delta \hat{y}}{\delta O_{11}} \times \frac{\delta O_{11}}{\delta b_{11}}$$

$$\frac{\delta L}{\delta w_{12}^1} = \frac{\delta L}{\delta \hat{y}} \times \frac{\delta \hat{y}}{\delta O_{12}} \times \frac{\delta O_{12}}{\delta w_{12}^1}$$

$$\frac{\delta L}{\delta w_{22}^1} = \frac{\delta L}{\delta \hat{y}} \times \frac{\delta \hat{y}}{\delta O_{12}} \times \frac{\delta O_{12}}{\delta w_{22}^1} \rightarrow O_{12}$$

$$\frac{\delta L}{\delta b_{12}} = \frac{\delta L}{\delta \hat{y}} \times \frac{\delta \hat{y}}{\delta O_{12}} \times \frac{\delta O_{12}}{\delta b_{12}}$$

$$\frac{\delta \hat{Y}}{\delta O_{11}} = \frac{\delta}{\delta O_{11}} [W_{11}^L O_{11} + W_{21}^L O_{12} + b_{21}] \\ = W_{11}^2$$

$$\frac{\delta L}{\delta W_{11}} = -2(Y - \hat{Y}) W_{11}^2 X_{i1} \rightarrow 4$$

$$\frac{\delta \hat{Y}}{\delta O_{12}} = \frac{\delta}{\delta O_{12}} [W_{11}^L O_{11} + W_{21}^L O_{12} + b_{21}] \\ = W_{21}^2$$

$$\frac{\delta L}{\delta W_{21}} = -2(Y - \hat{Y}) W_{21}^2 X_{i2} \rightarrow 5$$

$$\frac{\delta O_{11}}{\delta W_{11}^I} = \frac{\delta [i_q w_{11}^I + c_g p_a \cdot w_{21}^I + b_{11}]}{\delta W_{11}^I} \\ = i_q \rightarrow X_{i1}$$

$$\frac{\delta L}{\delta W_{11}^I} = -2(Y - \hat{Y}) W_{11}^2 X_{i1} \rightarrow 7$$

$$\frac{\delta L}{\delta W_{21}^I} = -2(Y - \hat{Y}) W_{21}^2 X_{i2} \rightarrow 8$$

$$\frac{\delta L}{\delta b_{12}} = -2(Y - \hat{Y}) W_{21}^2 \rightarrow 9$$

$$\frac{\delta O_{11}}{\delta W_{21}^I} = X_{i2}$$

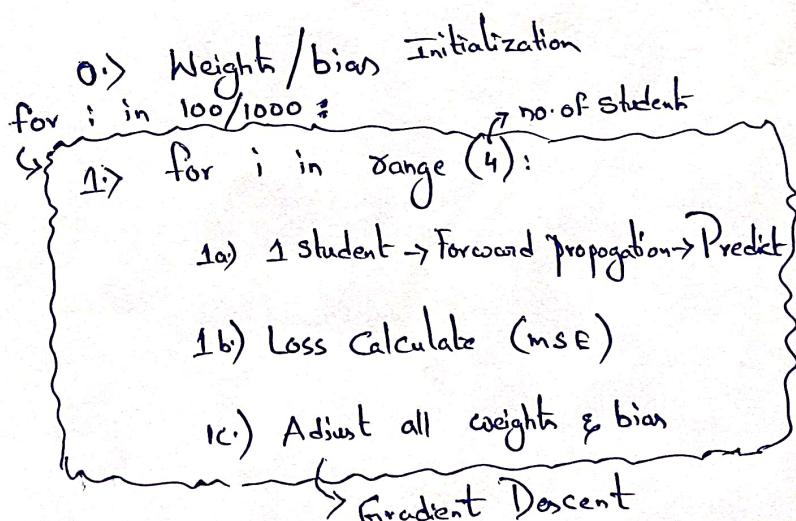
$$\frac{\delta O_{11}}{\delta b_{11}} = 1$$

$$\frac{\delta O_{12}}{\delta W_{12}^I} = \frac{\delta}{\delta W_{12}^I} [X_{i1} W_{12}^I + X_{i2} W_{22}^I + b_{12}] \\ = X_{i1}$$

$$\frac{\delta O_{12}}{\delta W_{22}^I} = \frac{\delta}{\delta W_{22}^I} [X_{i1} W_{12}^I + X_{i2} W_{22}^I + b_{12}] \\ = X_{i2}$$

$$\frac{\delta O_{12}}{\delta b_{12}} = \frac{\delta}{\delta b_{12}} [X_{i1} W_{12}^I + X_{i2} W_{22}^I + b_{12}] \\ = 1$$

Steps (once again)



Until loss becomes min, this keeps on repeating