

Deep Learning

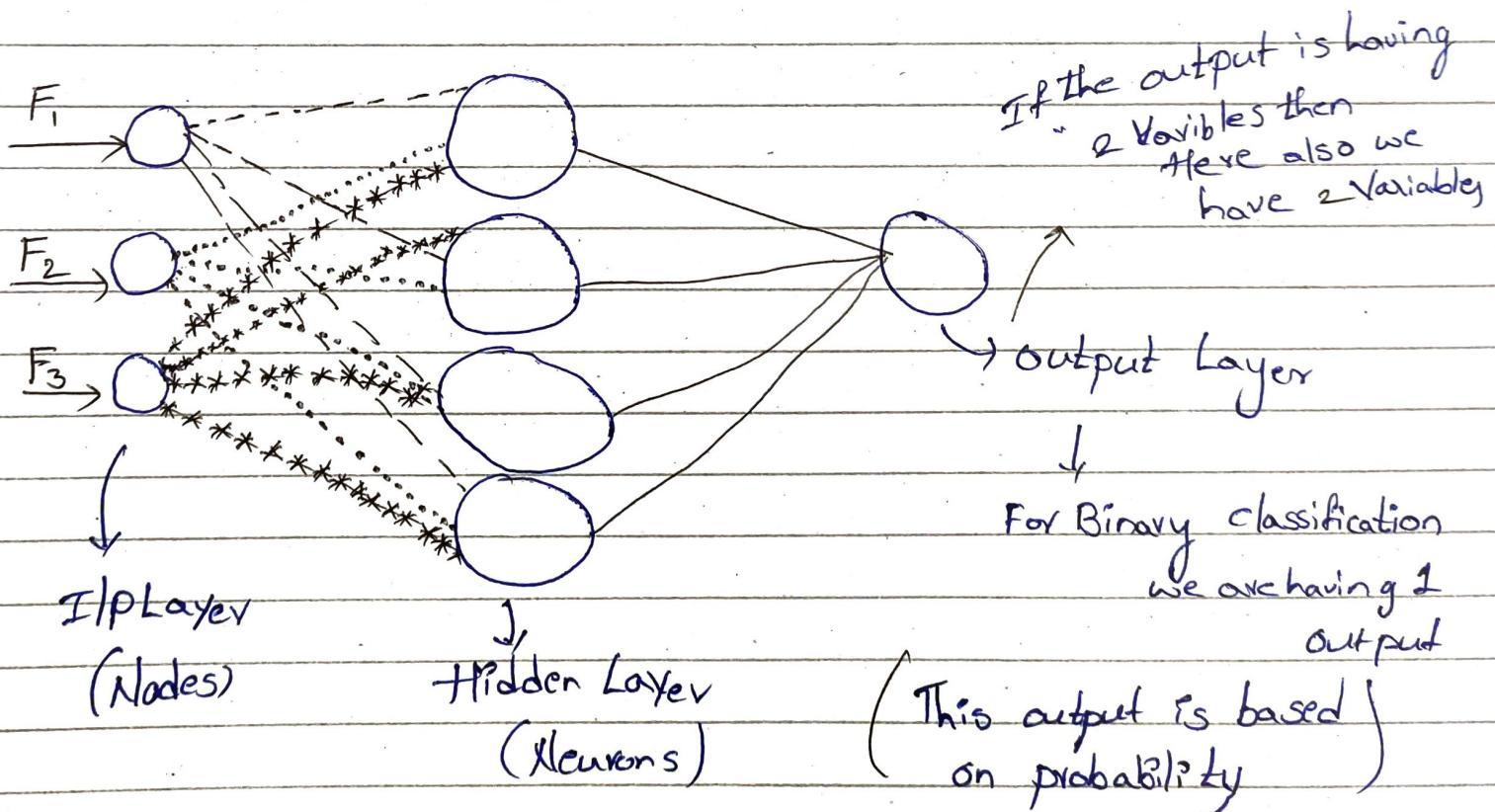
→ Scientists wants to invent a
→ A machine which learns and do things similar to human.

For this Scientists created Neural Networks

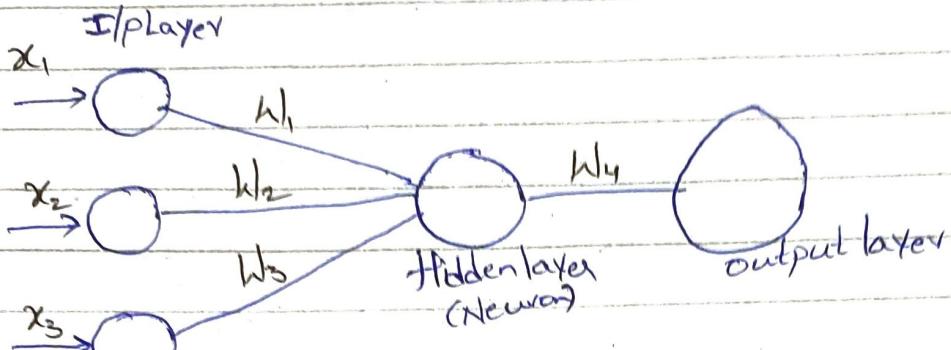
1st Simple type of Neural Network is called perceptron

→ This perceptron is having some issues then Jeffry Hinton discovered/created Back propagation.

Simple Neural Network Architecture →



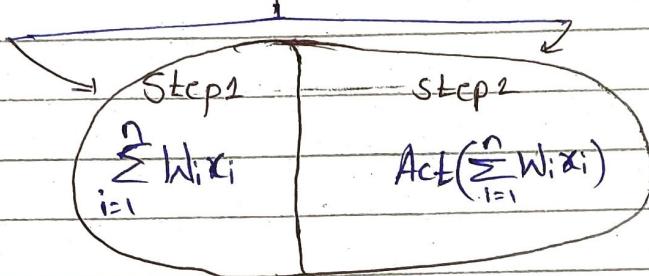
How the Neural Network Works.



($w_1, w_2, w_3, w_4 \rightarrow$ are called weights)

In the hidden layer (Neuron)

↳ there are 2 parts



Ex:-

→ Suppose if i keep normal water on my hand then

it won't affect.

→ Suppose if i keep very hot water on hand then it

will immediately neurons sends information.

→ Here in the both examples neurons are transmitting information but this weights are making impact.

→ When i kept hot water then the weights might be high & immediately sends information.

In the Neuron

Step-1:

→ In the Neuron 1st step

$$y_i = w_1x_1 + w_2x_2 + w_3x_3$$

adding bias

$$Y = w_1x_1 + w_2x_2 + w_3x_3 + \text{bias}$$

$$Y = \sum_{i=1}^n w_i x_i + \text{bias}$$

{ Each feature
is multiplied by
weights

Step-2:

apply activation function.

$$Z = \text{Act}(Y)$$

At the output: - (z)

$$z = Z \times w_4$$

$$= \text{Act}(Z)$$

→ So here again we applied Activation Function.

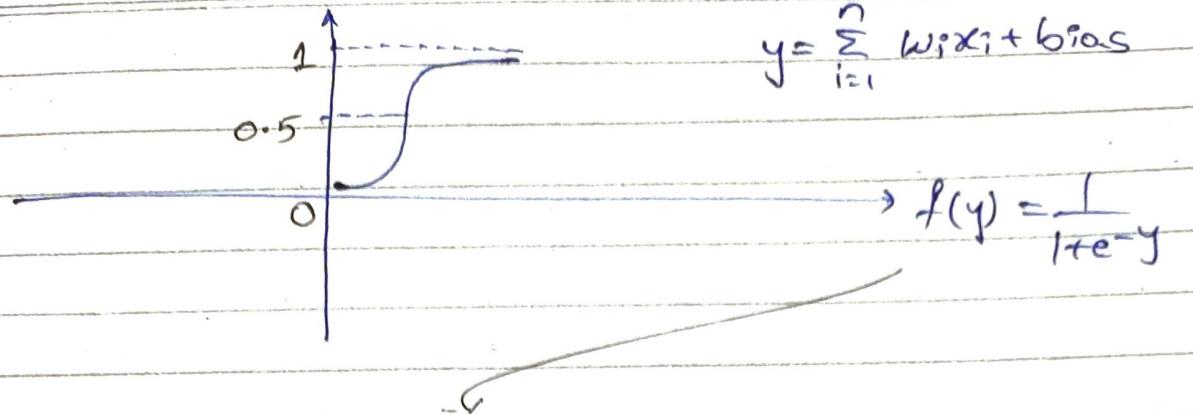
Activation Functions

1. Sigmoid Function: $\frac{1}{1+e^{-Y}}$ → transform values -1 to 1

In simple terms, when you kept normal water on hand then the neuron won't send information to brain.

But when you kept hot water neuron will be activated and passes information to brain quickly.

→ In this case neuron will be deactivated.



If the Value $< 0.5 \rightarrow 0 \rightarrow$ neuron not activated

If the Value $> 0.5 \rightarrow 1 \rightarrow$ neuron activated.

$\left\{ \begin{array}{l} \text{We will Consider only when } > 0.5 \\ 0.5 \rightarrow \text{is the threshold} \end{array} \right\}$

2. Relu Activation:-

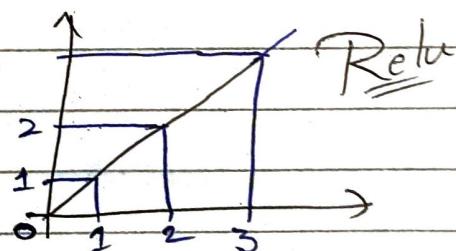
$$\text{Max}(y, 0)$$

$$\rightarrow y = w_1 x_1 + w_2 x_2 + w_3 x_3 + \text{bias} = \sum_{i=1}^n w_i x_i + \text{bias}$$

\rightarrow If $y = +ve$ then $\text{Max}(+ve, 0) = +ve \equiv y$

\rightarrow If $y = -ve$ then $\text{Max}(-ve, 0) = 0 =$

\rightarrow This will always give me maximum value.



Suppose you are solving classification problem then in the hidden layer we will use Relu and in the output

layer we use Sigmoid function (Because we applied Activation Functions in the hidden & output

Basically why we prefer Sigmoid at the output layer is because Sigmoid will classify the output nicely because it has values 0 to 1. (According to the probability we can choose).

(and you know Sigmoid will treat the outliers problem also
 (Converted to same range 0 to 1)

Neural Network Training

→ I am having dataset which contains 3 columns & 1 output

Play x_1	study x_2	Sleep x_3	output y	output → 1 (or) 0
2h	4h	8h	1	✓ Pass

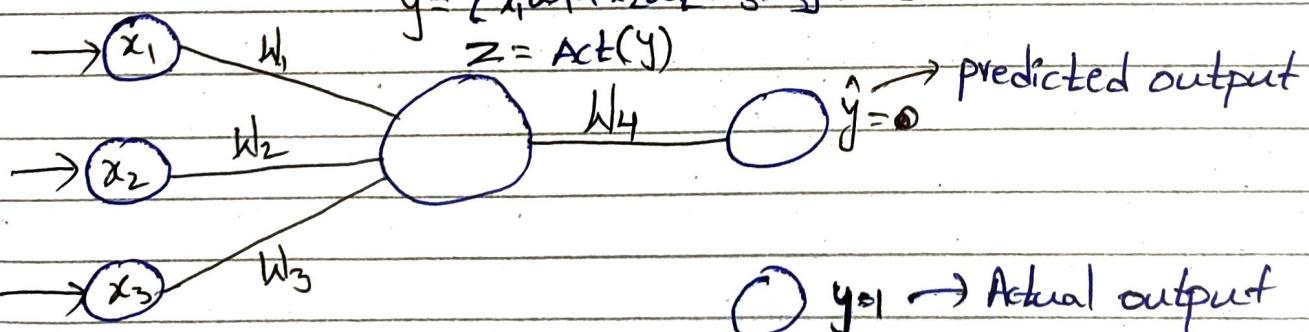
Input → play, study, sleep

→ So my neural network has 3 nodes (input) & 1 output layer.

No. of hours he is playing (or) studying (or) sleeping.

$$y = [x_1 w_1 + x_2 w_2 + x_3 w_3] + \text{bias}$$

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



So, in this example, i have predicted my output $\hat{y}=0$ but my output is $y=0$

↳ I need to compare my predicted output with the Actual output

↳ For that i need Loss Function

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

Q) Why Loss Function is having square

A) suppose if the Loss Function = $y - \hat{y}$

$$\text{Case(i)}: y=0, \hat{y}=1$$

$$\Rightarrow \text{Loss} = (0-1)$$

$$= -1$$

$$\text{Case(ii)}: y=1, \hat{y}=0$$

$$\Rightarrow \text{Loss} = 1-0$$

$$= 1$$

→ Here we are getting negative, that's why are eliminating negative by using Square

Q) What if you have multiple outputs

Ans: Then i will sum it up

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

↗ It won't be called as Loss Function

↗ Called as Cost Function.

→ So, we always want to have Minimize Loss Function

→ Because, when my outputs are correct then my Loss Function is zero

→ If my outputs are not correct then my Loss Function is 1

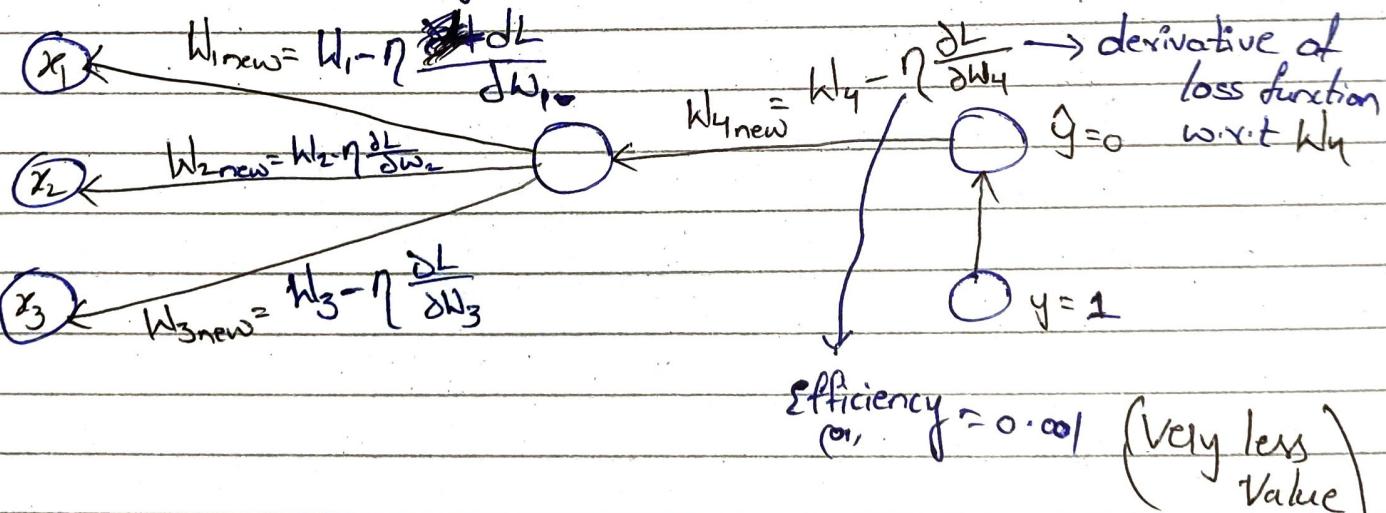
$$\Rightarrow (y - \hat{y})^2 = (1-0)^2 = 1$$

→ That's why we always want to minimize loss function

↗ Such that actual output and predicted o/p should be same.

↳ In order to reduce loss function we need optimizer.

↳ This optimizer will help us in adjusting weights in Back propagation.

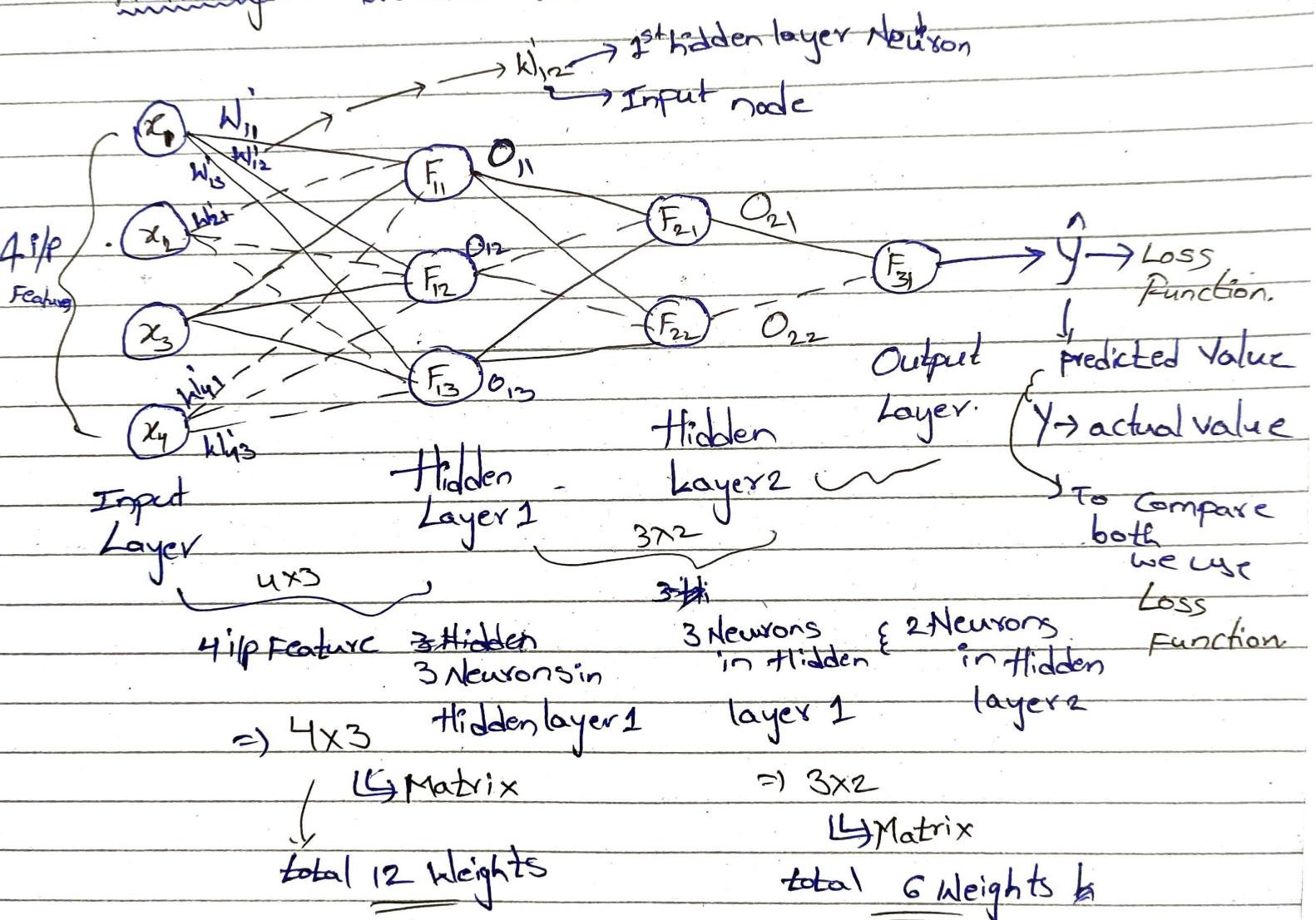


{ Here $w_1, w_2, w_3, w_4 \rightarrow$ old,

↳ Using optimizer we are updating w_1, w_2, \dots, w_4 to new weights.

After the weights are adjusted we need to do again forward propagation. This process will go until some number of epochs or the loss value completely reduced.

Multi layer Neural Network (MLN)



Note: $x_1, O_{11}, O_{12}, O_{13}, O_{21}, \dots, O_{22} \rightarrow$ all these are outputs.

- In the First hidden layer we wrote $F_{11}, F_{12}, F_{13} \rightarrow$ all these will apply activation Functions. (lets say $w_{11}x_1 + \dots + w_{13}x_4 = y$)
- on this we will apply activation Function (y_i) \Rightarrow we get O_{11} .

Loss Function:-

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

→ We want to reduce the loss function.



In order to reduce we use optimizer to update the weights.

↳ (Back propagation)

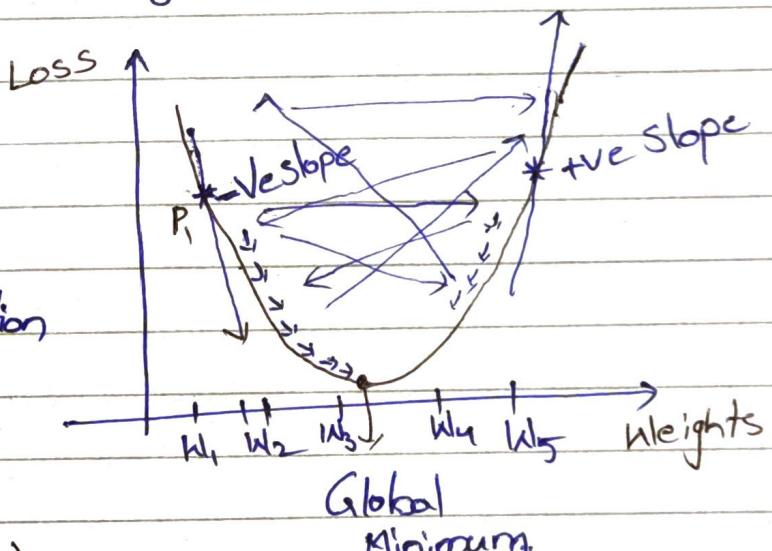
Gradient Descent

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w_{\text{old}}} \quad \begin{array}{l} \text{loss function} \\ \text{derivative of Loss Function} \\ \text{w.r.t to } w_{\text{old}} \end{array}$$

↳ learning Rate

→ Derivative is nothing but Slope.

→ We want to reach Global minimum then only our Loss Function will be less.



η → very less (0.001)

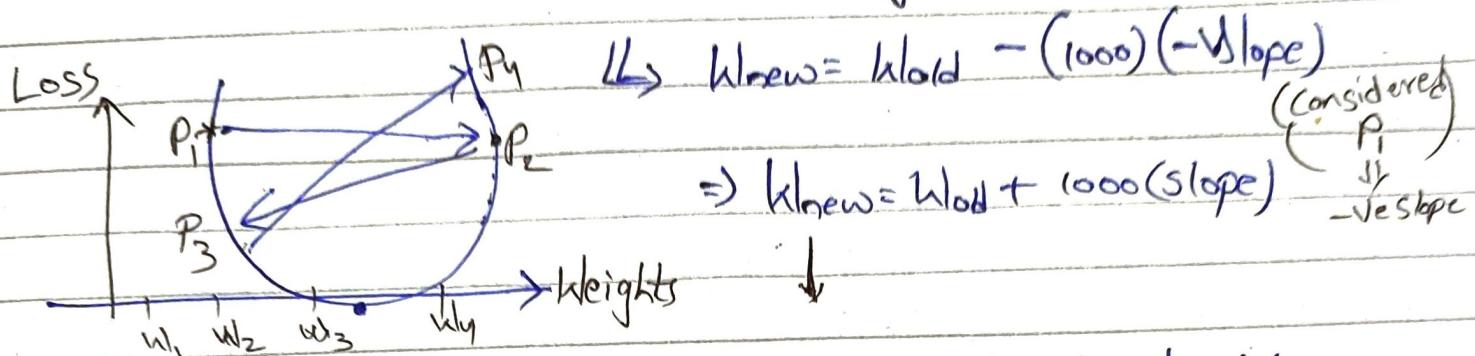
⇒ Suppose My loss value is at P_1 , then i calculated slope using derivative & it is -ve (Pointing down)

→ My equation becomes $w_{\text{new}} = w_{\text{old}} - \eta (-\text{ve}) = w_{\text{old}} + \eta (\text{slope})$
=) So, My weight is increasing some-ve value

→ So, here my weight is increasing at a speed of η

$\eta =$ Very high.

⇒ If I keep $\eta = 10000 \Rightarrow$ then my weight increases rapidly



In this case we never reach global minimum

↳ It's like weights will be updated from

$$P_1 \rightarrow P_2 \rightarrow P_3 \rightarrow P_4.$$

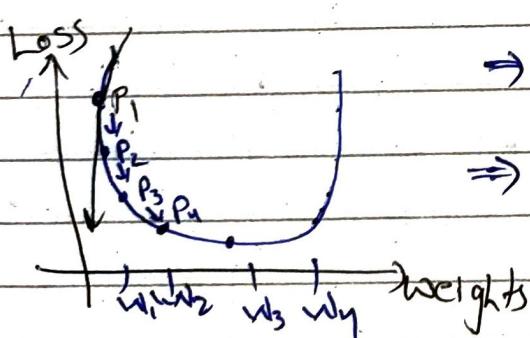
↳ So it will never converge.

→ $\eta =$ Very low.

↳ then my weight moves very slowly

$$\eta = 0.001$$

$$\Rightarrow w_{\text{new}} = w_{\text{old}} - (0.001)(-\text{slope})$$



$$\Rightarrow w_{\text{new}} = w_{\text{old}} + 0.001\text{slope}$$

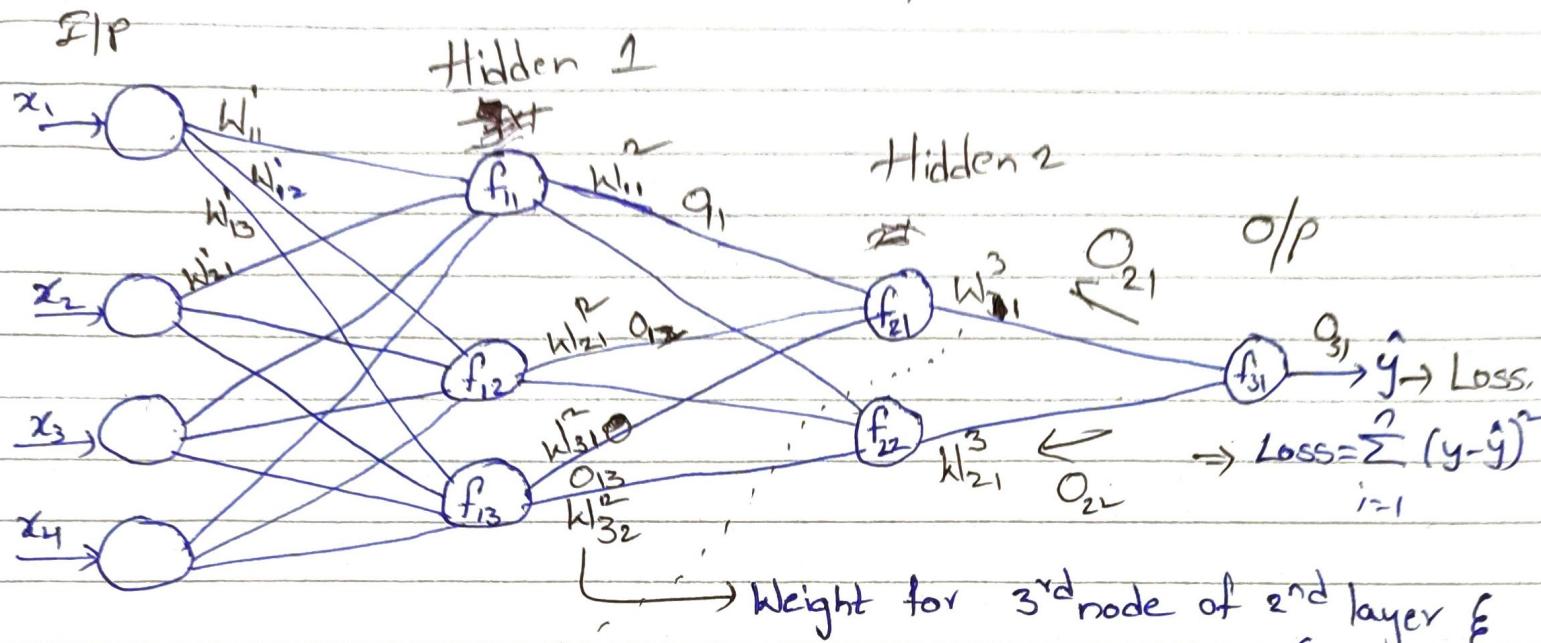
⇒ w_{new} increased at very low pace/speed.

→ So here weights will be updated slowly

and like $(P_1 \rightarrow P_2 \rightarrow P_3 \rightarrow P_4)$

→ at it will reach Global Minimum (Converges)

Chain Rule in Back propagation:



Back propagation j

$$w_{11}^3 \text{ new} = w_{11}^3 \text{ old} - \left\{ \begin{array}{l} \frac{\partial L}{\partial w_{11}^3} \rightarrow \text{slope} \\ \text{learning Rate} \end{array} \right\} \quad (\text{optimizer})$$

Derivative of O/P layer $O_{31} \rightarrow (w_{11}^3, w_{21}^3)$

$$\frac{\partial L}{\partial w_{11}^3} = \frac{\partial L}{\partial O_{31}} \times \frac{\partial O_{31}}{\partial w_{11}^3}$$

\downarrow
O/P at output layer

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

\downarrow
 $O_{31} \rightarrow O_{21} (w_{11}^3)$
impacts

\downarrow
 $O_{31} \rightarrow O_{22} (w_{21}^3)$
impact

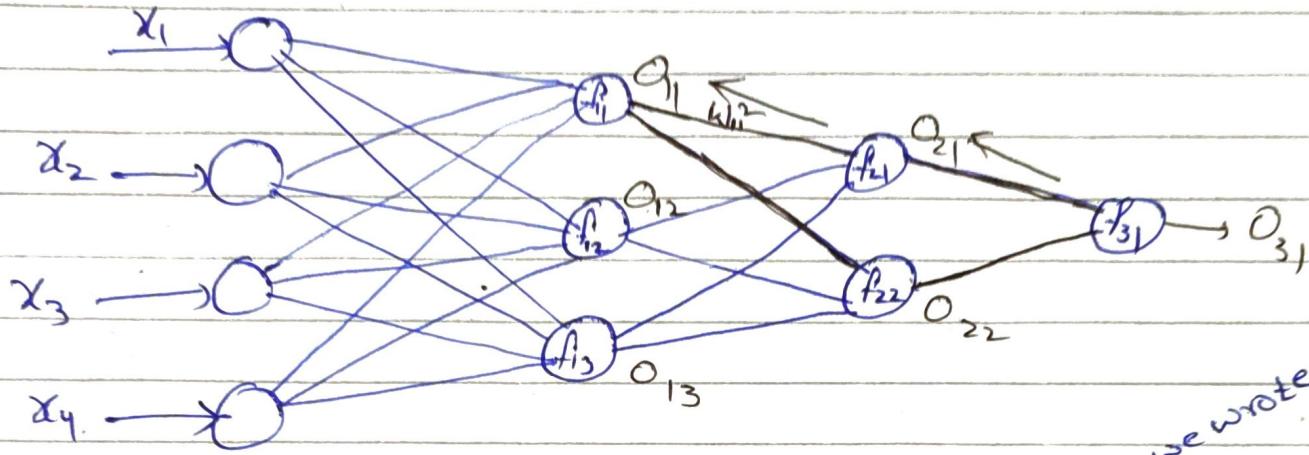
$$w_{11}^2 \text{ new} = w_{11}^2 \text{ old} - \eta \frac{\partial L}{\partial w_{11}^2}$$

$$\frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial O_{31}} \times \frac{\partial O_{31}}{\partial O_{21}} \times \frac{\partial O_{21}}{\partial w_{11}^2}$$

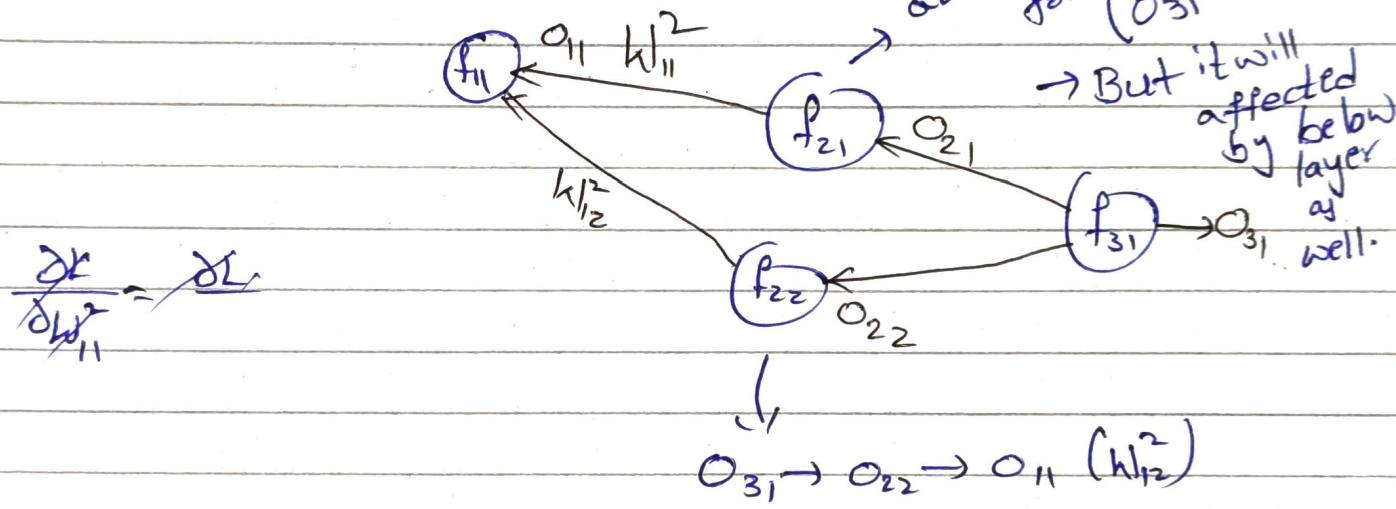
\downarrow
1 $O_{31} \rightarrow O_{21} \rightarrow O_{11} (w_{11}^2)$
impacts impacts

This is chain rule

at Eq① (Hidden 2 to hidden 1)



at Eq① we wrote only
for this $O_{21} \rightarrow O_{11}$
(w_{11}^2)



Here in both ways O_{11} is affecting. So, we need to add both

$$\frac{\partial L}{\partial w_{11}^2} = \left[\frac{\partial L}{\partial O_{31}} \times \frac{\partial O_{31}}{\partial O_{22}} \times \frac{\partial O_{22}}{\partial w_{11}^2} \right] + \left[\frac{\partial L}{\partial O_{31}} \times \frac{\partial O_{31}}{\partial O_{22}} \times \frac{\partial O_{22}}{\partial w_{12}^2} \right]$$

\downarrow
upper

\downarrow
lower.

→ until we reach global minimum, this weights will be updated in Back propagation (or) it will do for some epochs.

→ Epochs! → 1 Forward & 1 Back propagation \Rightarrow i/p $\xrightarrow{\text{O/p}} (1 \text{ epoch})$
↳ is called 1 epoch.