

- Q) why Deep learning is too much famous?
- 1) Applicable in each and every domain
 - 2) able to solve vast amount of problems.
 - 3) DL model can automatically learn relevant feature from the data, reducing the need for manual feature engineering.
 - DL models can scale effectively with large amounts of data and computational resources.

- * feature extraction in deep learning
- feature extraction refers to the process of automatically learning and extracting relevant features or patterns from raw data
 - Unlike traditional machine learning approaches that often rely on manual feature engineering, deep learning model can automatically learn intricate representations directly from the input data.
 - feature Extraction is critical step in deep learning that allows the model to understand and interpret complex patterns within the data.

⊗ deep learning Model →

↳ Image classification → - ResNet

- google Net

↳ object detection: → - YOLO

- SSD

- fast and faster RCNN

3) object segmentation → i) Mask RCNN

- U-Net

- V-Net

4) Image translation → - pix to Pix

5) Speech generation → - Webnet

6) Text classification → - BERT

* Transfer learning →

- In transfer learning a pre-trained model, often trained on a large dataset for a related task, serves as a starting point for training new model.
- Instead of initializing the new model from scratch, the weight and parameters of the pre-trained model are utilized as foundation.

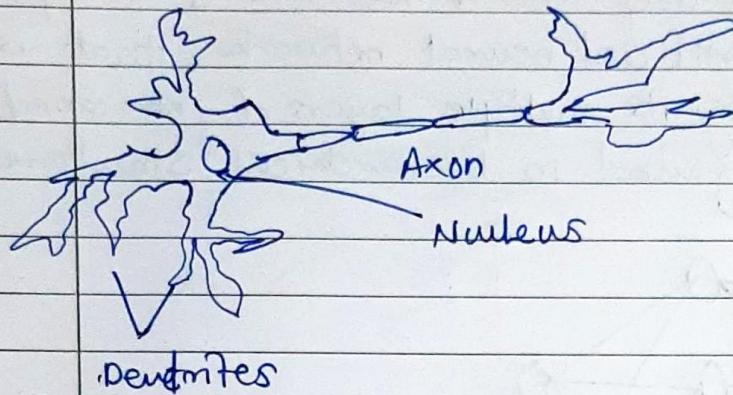
* ways to apply transfer learning in deep learning:

① feature extraction:

- This involve using pre-trained model as a fixed feature extractor. the earlier layers of the pre-trained model are used to extract relevant features from the input data, and this features are then fed into a new, custom classifier or set of layers tailored to specific task

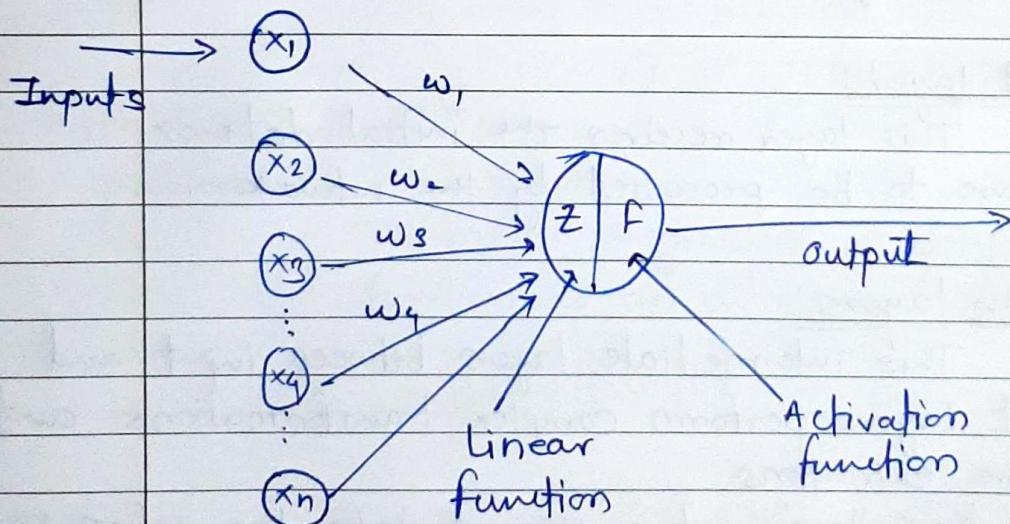
② fine-tuning: The pre-trained model's are further adjusted or fine-tuned during the training process on the new dataset. the goal is to allow to adapt its learned representations to the nuance of the new data while still retaining the knowledge gained from the original task.

* Artificial Neuron \leftrightarrow (Perception)



Human Neuron.

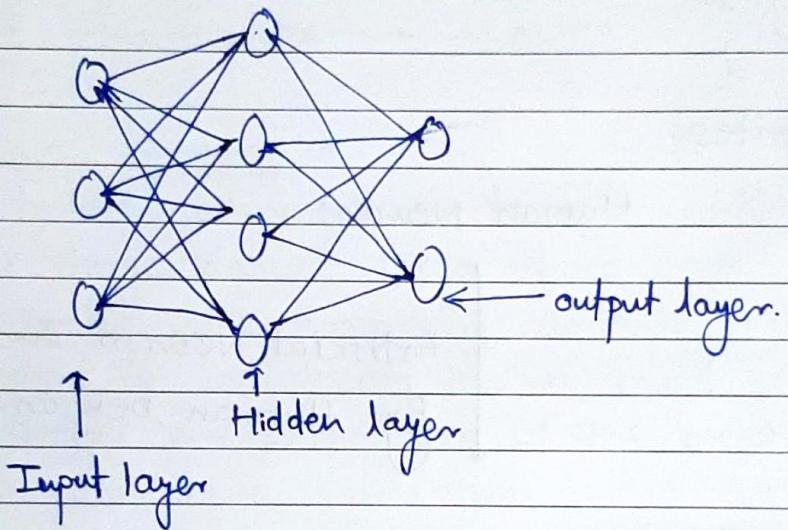
↓
Artificial Neuron Inspired
By Human neuron.



Artificial Neuron [Perception]

* DNN → Deep Neural Network

- Basic unit of deep Neural Network is perceptron.
- is type of artificial neural network, that is characterized by its multiple layers of interconnection of neuron organized in hierarchical structure.



* Input layer :-

This layer receives the initial data or feature to be processed by the network.

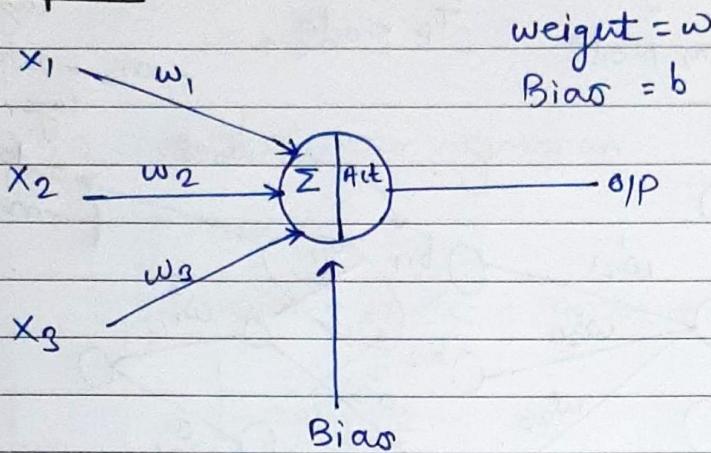
* Hidden layers :-

This intermediate layers between input and output layers perform complex transformations and feature extractions.

- The term 'Deep' in deep Neural Networks refers to the presence of multiple hidden layers.

* Output layers :- produce final output of the network.

* perception \Rightarrow



$$\text{weight} = \omega$$

$$\text{Bias} = b$$

$$x_1 w_1 + x_2 w_2 + x_3 w_3 + b \Rightarrow \omega^T x + b$$

$$\text{ACT}(\omega^T x + b) \Rightarrow \text{O/P}$$

Similar to $y = mx + c$

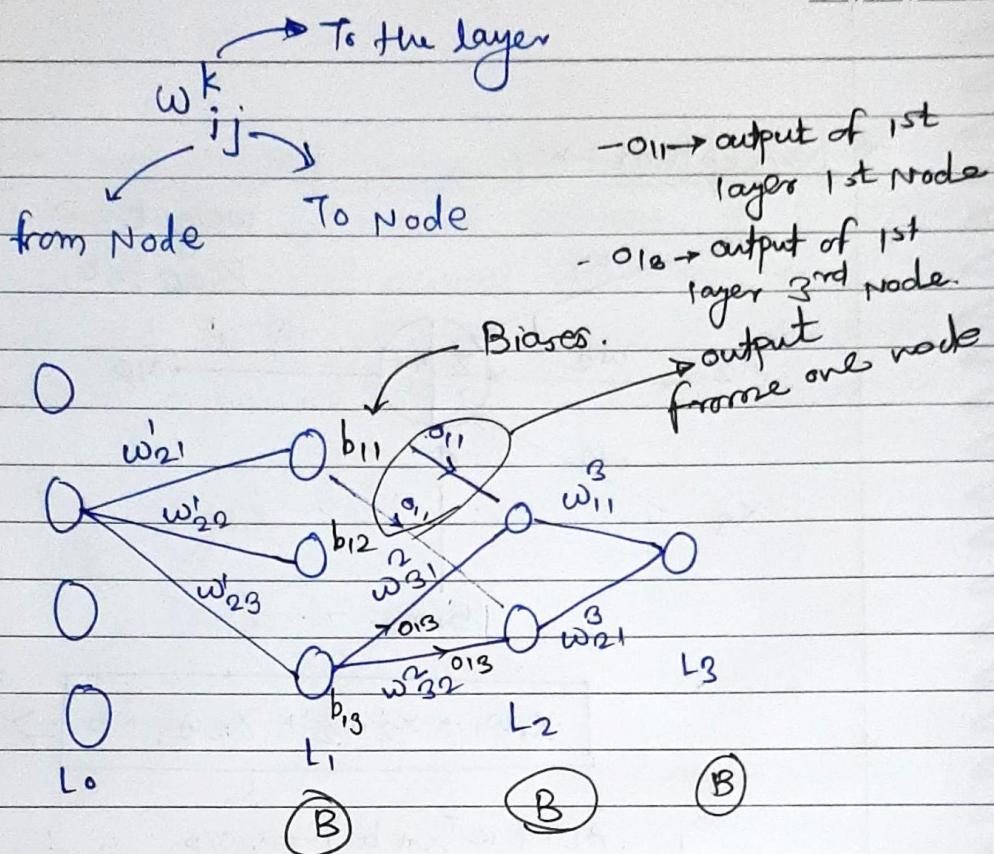
$$y = m_1 x_1 + m_2 x_2 + m_3 x_3 + c$$

for classification
problem

In logistic regression we apply sigmoid function as an activation function.

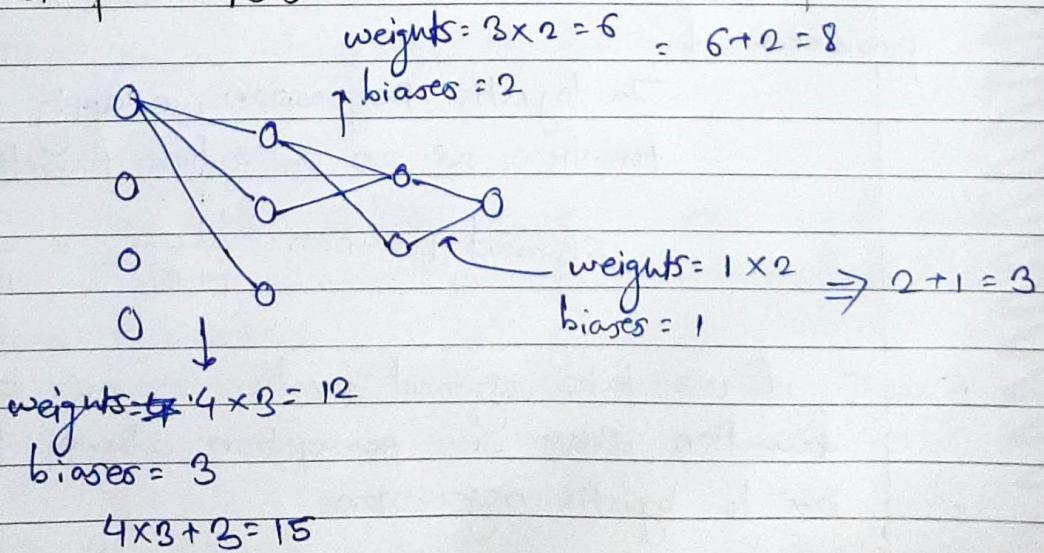
$$\therefore \text{sigmoid function} = \frac{1}{1+e^{-x}} = \frac{1}{1+e^{-(y)}}$$

- if we take sigmoid function as our activation function then our perception above become similar to logistic regression.



- with respect to every perceptron we have bias. it have only to hidden layers and output layers
- Biases are constant.

* Trinable parameters:



∴ Total trainable parameter =

$$15 + 8 + 3 = 26$$

- Now, in case of linear regression.

$$\hat{y} = mx + c$$

$$\hat{y} = m_1x_1 + m_2x_2 + m_3x_3 + m_4x_4 + c$$

↓ optimizer.

The parameter going to change or update for greater accuracy will be m_1, m_2, m_3, m_4, c getting minimum loss.

∴ these 5 parameters are trainable parameters.

* Forward propagation →

Data:-

CGPA	iq	10 th	12 th	placed
-	-	-	-	-

Hyperparameter
- Hidden layer.
- No. of Neuron.

Input layer = No. of feature

CGPA

$$x_1 \rightarrow 0$$

→ forward propagation

0

IQ

$$x_2 \rightarrow 0$$

0

10th

$$x_3 \rightarrow 0$$

0

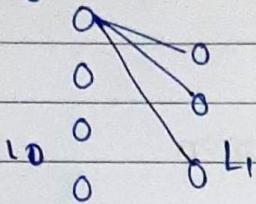
12th

$$x_4 \rightarrow 0$$

0

← Backward propagation.

- * forward propagation calculation.
- for Layer 1 →



$$\begin{bmatrix} w'_{11} & w'_{12} & w'_{13} \\ w'_{21} & w'_{22} & w'_{23} \\ w'_{31} & w'_{32} & w'_{33} \\ w'_{41} & w'_{42} & w'_{43} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} + \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \end{bmatrix}$$

↓ Take the transpose

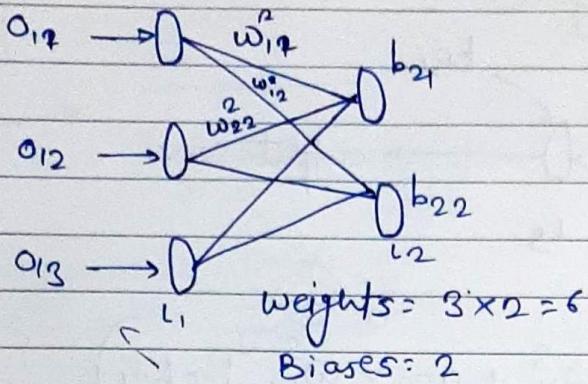
$$\begin{bmatrix} w'_{11} & w'_{21} & w'_{31} & w'_{41} \\ w'_{21} & w'_{22} & w'_{32} & w'_{42} \\ w'_{31} & w'_{32} & w'_{33} & w'_{43} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} + \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \end{bmatrix}$$

$$\begin{bmatrix} w'_{11}x_1 + w'_{21}x_2 + w'_{31}x_3 + w'_{41}x_4 \\ w'_{21}x_1 + w'_{22}x_2 + w'_{32}x_3 + w'_{42}x_4 \\ w'_{31}x_1 + w'_{32}x_2 + w'_{33}x_3 + w'_{43}x_4 \end{bmatrix} + \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \end{bmatrix} = \begin{bmatrix} o_{11} \\ o_{12} \\ o_{13} \end{bmatrix}$$

↓ This output of 1st layer goes to
The 2nd layer.
- forwarded by the 3 perceptron -

- That's how our forward propagation got calculated.

- for the 2nd layer



$$\begin{bmatrix} w_{11}^2 & w_{12}^2 \\ w_{21}^2 & w_{22}^2 \\ w_{31}^2 & w_{32}^2 \end{bmatrix} \times \begin{bmatrix} O_{11} \\ O_{12} \\ O_{13} \end{bmatrix} + \begin{bmatrix} b_{21} \\ b_{22} \end{bmatrix} =$$

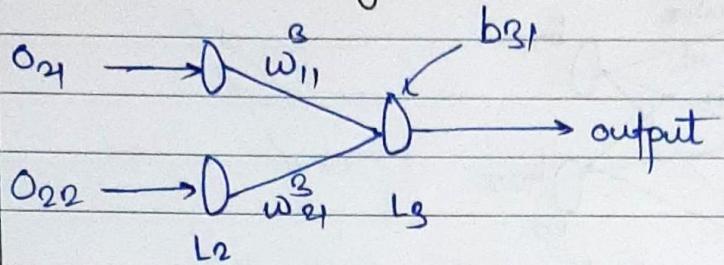
↓
Transpose

$$\begin{bmatrix} w_{11}^2 & w_{12}^2 & w_{31}^2 \\ w_{12}^2 & w_{22}^2 & w_{32}^2 \end{bmatrix} \times \begin{bmatrix} O_{11} \\ O_{12} \\ O_{13} \end{bmatrix} + \begin{bmatrix} b_{21} \\ b_{22} \end{bmatrix}$$

$$\left(\begin{bmatrix} w_{11}^2 O_{11} + w_{21}^2 O_{12} + w_{31}^2 O_{13} + b_{21} \\ w_{12}^2 O_{11} + w_{22}^2 O_{12} + w_{32}^2 O_{13} + b_{22} \end{bmatrix} \right) = \begin{bmatrix} O_2 \\ O_{22} \end{bmatrix}$$

after applying activation
functions then you'll get
the output for that
particular layer.

* for the final layer



$$\left([O_{21} \times w_{11}^3 + O_{22} \times w_{21}^3] + b_{31} \right)$$

↓ activation function

final output $\rightarrow O_{31}$

- final output may be prediction, classification can be anything w.r.t. to the problem statement.

* loss function \Rightarrow

- difference between loss and cost :

Height	weight	BMI	BMI - Predicted
170	65	22	2
185	68	23	2.5
175	71	24	0.5

$\rightarrow \text{loss} = 22 - 21 = 1$

$\Rightarrow \text{actual} - \text{predicted}$

- loss will take only one sample data at a time.

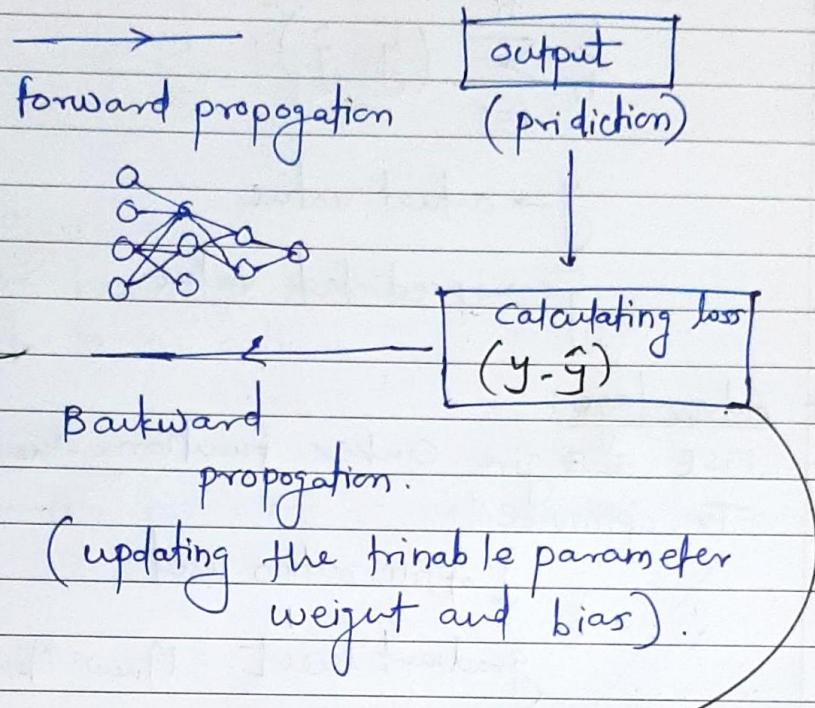
$\text{cost} = \frac{1}{n} \sum_{i=1}^n (\text{actual} - \text{pred})$

- cost will take whole data into consideration.

forward propagation
again forward propagation
with updated parameter.

The number of cycle divided
epoch [Not loop]

- loss function:



- the loss is like a feedback
- if the loss is more then we have to update the trainable parameter to reduce the loss.

- loss function

Regression

- ① MSE
- ② MAE

classification

- ① Binary Cross Entropy or log loss.
- ② Categorical Cross Entropy
- ③ Sparse Categorical Cross Entropy.

* MSE → Mean Square Error.

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

y → actual value

\hat{y} → predicted value.

+ Advantage:

- MSE ~~is~~ give convex function. therefore easy to optimize.

↓ optimization used

$$\text{gradient descent} = \mathbf{m}_{\text{new}} = \mathbf{m}_{\text{old}} - \alpha \frac{\partial L}{\partial \mathbf{m}}$$

loss function
learning rate

- easy to interpret

+ disadvantage:

- not robust to the outlier
- for real value we have to go again take not of the same.

* MAE: Mean Absolute Error.

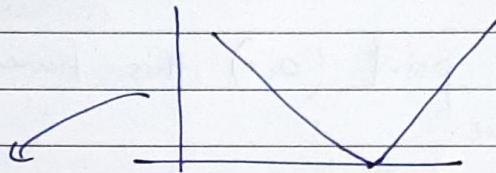
$$\frac{1}{n} \sum_{j=1}^n |y - \hat{y}|$$

* Advantage:

- give the loss in particular range.
- robust to the outliers.

* Disadvantage:

- The function for MAE is not convex function



- Non-differentiable

- we can't use gradient descent to optimize and find the global minima.

* Backorder Propagation \Rightarrow

- The goal is to minimize the loss by adjusting the weights of neural network.
- Backorder propagation is a process of updating the weights by propagating the error backward through the network.

* Calculation of minima and maxima.

$$z = (x^2 + y^2)$$

Condition for Calculating minima

$$\frac{\partial z}{\partial x} = 0 \quad \text{and} \quad \frac{\partial z}{\partial y} = 0$$

$$\frac{\partial (x^2 + y^2)}{\partial x} = 0 \quad \text{and} \quad \frac{\partial (x^2 + y^2)}{\partial y} = 0$$

$$2x = 0$$

$$\boxed{x = 0}$$

$$2y = 0$$

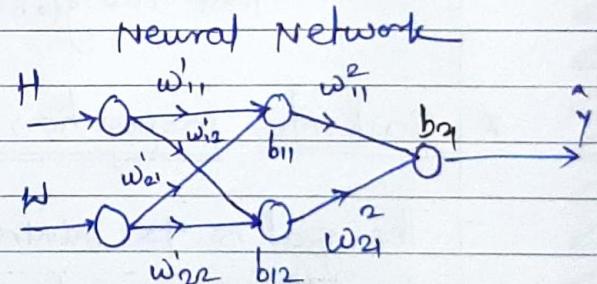
$$\boxed{y = 0}$$

∴ at point $(0,0)$ this function 'z' have minimum value.

* Backward propagation \Rightarrow

Height weight BMI

160	60	21
170	65	22
180	70	20
190	71	18



Trainable parameter $\rightarrow 4 + 2 + 2 + 1 = 9$

\hat{y} depends upon this '9' parameter

understand the Backward propagation
and optimization of weight and bias
with example

* Step-1

Initialize weight and bias

Method \rightarrow Random and define

* Step-2

point (row)

* Step-3 :-

forward propagation

[Matrix Multiplication, dot product]

- output calculated \hat{y} .

* Step-4

loss function

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

↗ variable
 ↓ fix value

$$\text{let } \hat{y} = 24$$

$$\text{Loss} = (21 - 24)^2$$

$$\boxed{\text{Loss} = 9}$$

- our main aim to reduce the loss

* Step-5

gradient descent [updating the weight and biases]

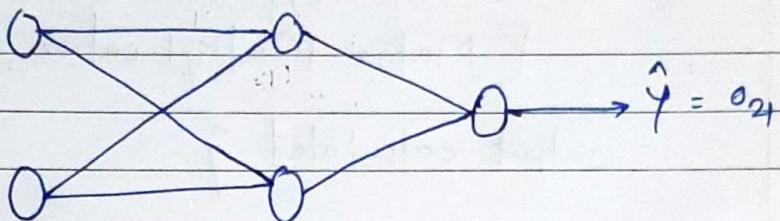
$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w}$$

[Python → Implement Backpropagation → code →
from scratch]

loss function.

$$b_{new} = b_{old} - \eta \frac{\partial L}{\partial b}$$

* chain rule of derivation \Rightarrow



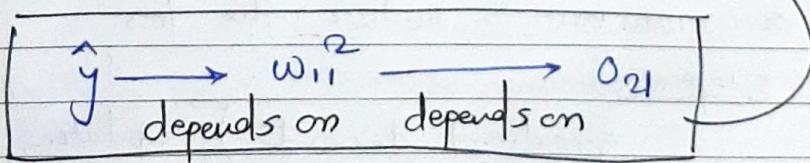
$$\hat{y} = o_{11} \times w_{11}^2 + o_{21} w_{21}^2 + b_{21}$$

$$\text{Loss} = y \cdot \hat{y}$$

for updating the w_{11}^2

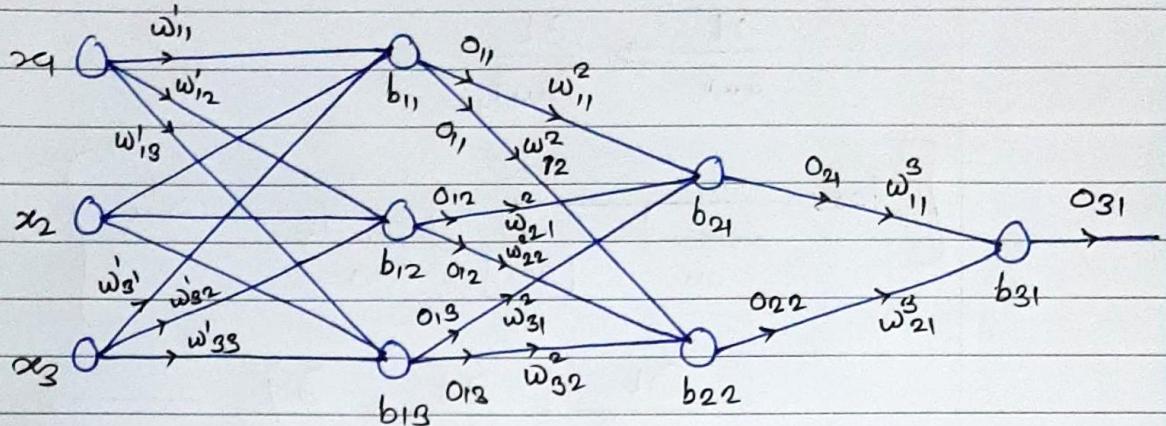
$$\frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial \hat{y}}$$

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



chain rule of
derivation.

updating trainable parameter with backpropagation \rightarrow



Trainable parameter \rightarrow

$$3 \times 3 + 3 = 12$$

$$(3 \times 2) + 2 = 8$$

$$(2 \times 1) + 1 = 3$$

Total
[23]

+ for updating \rightarrow gradient descent

Suppose we want to update ' w ' [weight]

Then -

$$\boxed{w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w_{\text{old}}}} \rightarrow \text{loss function}$$

Learning Rate

Now,

will focus on only $\boxed{\frac{\partial L}{\partial w}}$

and How chain rule of derivation work.

take $\omega = \omega_{11}^3$

$$\frac{\partial L}{\partial \omega} \Rightarrow \frac{\partial L}{\partial \omega_{11}^3}$$

$\left[\text{loss} \xrightarrow[\text{on } o_{31}]{} \hat{y} \xrightarrow[\text{on } \omega_{11}^3]{} \omega_{11}^3 \right]$

$$\therefore \frac{\partial L}{\partial \omega_{11}^3} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial \omega_{11}^3} \quad \begin{matrix} \text{by chain} \\ \text{rule of derivation} \end{matrix}$$

Now,

$\left[\omega_{11}^3 \xrightarrow[\text{on } o_{21}]{} \omega_{11}^2 \xrightarrow[\text{on } \omega_{11}^2]{} \omega_{11}^2 \right]$

$$\left\{ \begin{matrix} \text{can} \\ \text{include} \\ \omega_{11}^3, \end{matrix} \quad \therefore \frac{\partial L}{\partial \omega_{11}^2} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial o_{21}} \times \frac{\partial o_{21}}{\partial \omega_{11}^2} \right.$$

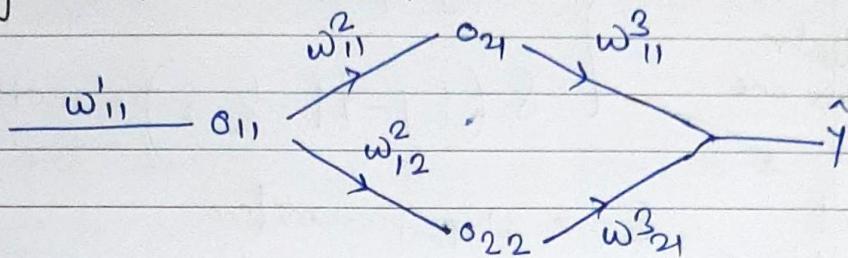
Now,

$\left[\text{loss} \rightarrow \omega_{11}^3 \rightarrow o_{21} \rightarrow \omega_{11}^2 \rightarrow o_{11} \rightarrow \omega_{11}' \right]$

again
can include
weight which was
not included.

$$\frac{\partial L}{\partial \omega_{11}'} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial o_{21}} \times \frac{\partial o_{21}}{\partial \omega_{11}^2} \times \frac{\partial \omega_{11}^2}{\partial o_{11}} \times \frac{\partial o_{11}}{\partial \omega_{11}'}$$

- * The error \hat{y} depends on two paths.



$\therefore \frac{\partial L}{\partial w_{11}}$ will take both path into consideration
 \rightarrow final output will be addition of two paths.

$$\therefore \frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial \hat{y}} \left[\frac{\partial \hat{y}}{\partial o_{21}} \times \frac{\partial o_{21}}{\partial o_{11}} \times \frac{\partial o_{11}}{\partial w_{11}} + \frac{\partial \hat{y}}{\partial o_{22}} \times \right.$$

$$\left. \frac{\partial o_{22}}{\partial o_{11}} \times \frac{\partial o_{11}}{\partial w_{11}} \right]$$

If you want to mention the weights