

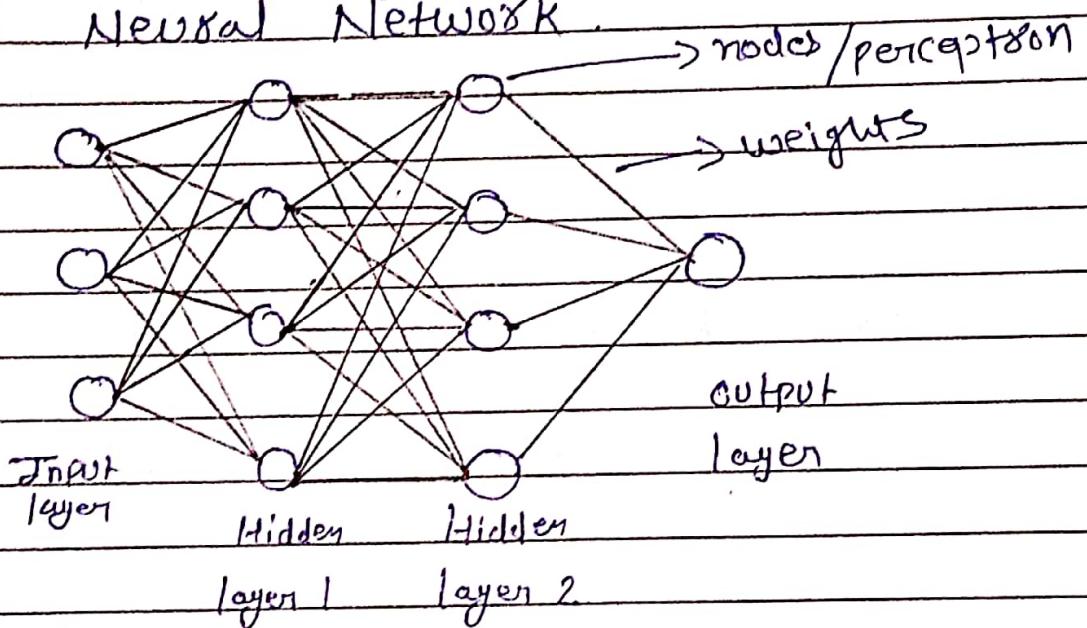


"Deep Learning"

## # What is deep learning ?

Deep Learning is a subfield of Artificial Intelligence and Machine Learning that is inspired by the structure of human brain.

Deep learning algorithms attempt to draw similar conclusions as humans would by continually analyzing data with a given with a given logical structure called Neural Network.



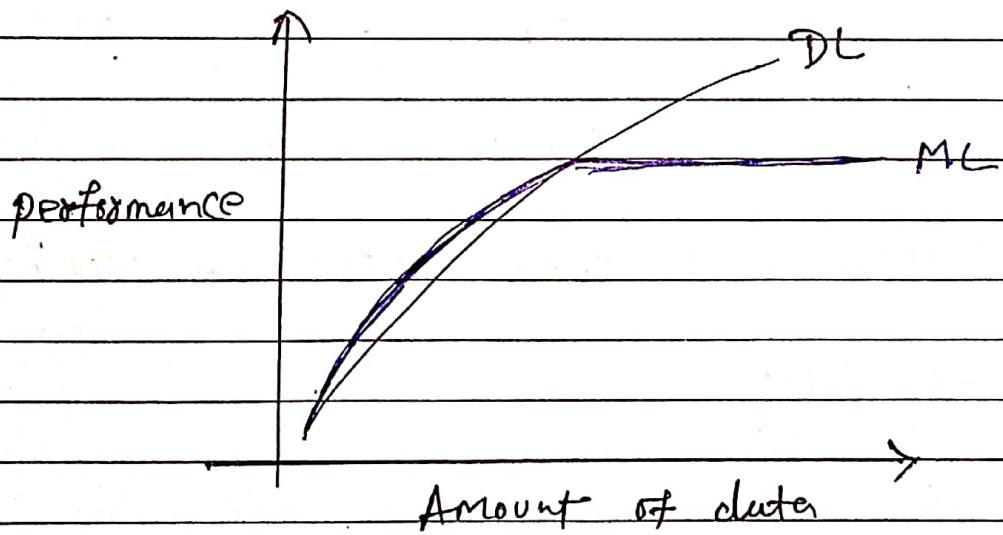
There are some types of NN :-

- \* ANN
- \* CNN
- \* RNN
- \* GAN
- \* etc.

Deep learning is a part of broader family of ML methods depend on ANN with representation / feature learning.

Deep learning algorithms uses multiple layers to progressively extract higher-level features from the raw input. for eg , in image processing , lower layer may identify edges , while higher layers may identify the concept relevant to a human such as digits or letters or faces.

Deep learning Vs ML :-



## ML

## DL

- |   |  |   |
|---|--|---|
| ① | ML requires less data. After a limit of data ML performance are fixed.                           | DL is data hungry. <del>After</del> its performance improves linearly to data.                                      |
| ② | Does not require powerful hardware's i.e. it can work on CPU, so it is budget friendly.          | Requires powerful hardware like GPU's So it is expensive.   |
| ③ | Training time is low.  | Training time is high (Cause big data)  |
| ④ | In ML we provide some features so that model generate result by using extract the given features | but in DL uses representation / feature learning which allows system to extract feature's according to requirement. |
| ⑤ | ML is a Superset of DL   | DL is a Subset of ML  |
| ⑥ | The data represented in ML is quite different compared to DL as it uses structured data          | The data represented in DL is quite diff as it uses ANN.  |

Why It's So Popular & Why Not

\* Datasets

\* Hardware

\* Frameworks

\* Libraries

\* People learning

Types of Neural Networks -

There are various types of NN and here we discuss some main types.

① Neural networks are computational model that mimic the way biological neurons work in the human brain process info.

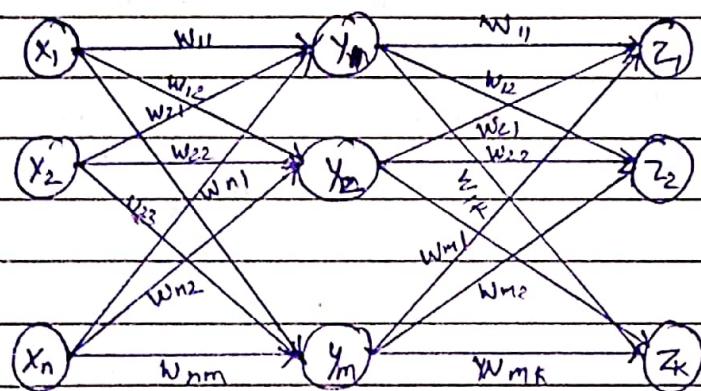
The first & last & central that function is to take input data into meaningful output through a series of mathematical operations.

## ① Feedforward Neural Networks :-

Feedforward neural networks are a form of ANN where without forming any cycles bet<sup>n</sup> layers or nodes means i/p can pass data through those nodes within the hidden level to the o/p nodes.

### Applications :-

in visual and voice recognition, NLP, financial forecasting and recommending system.

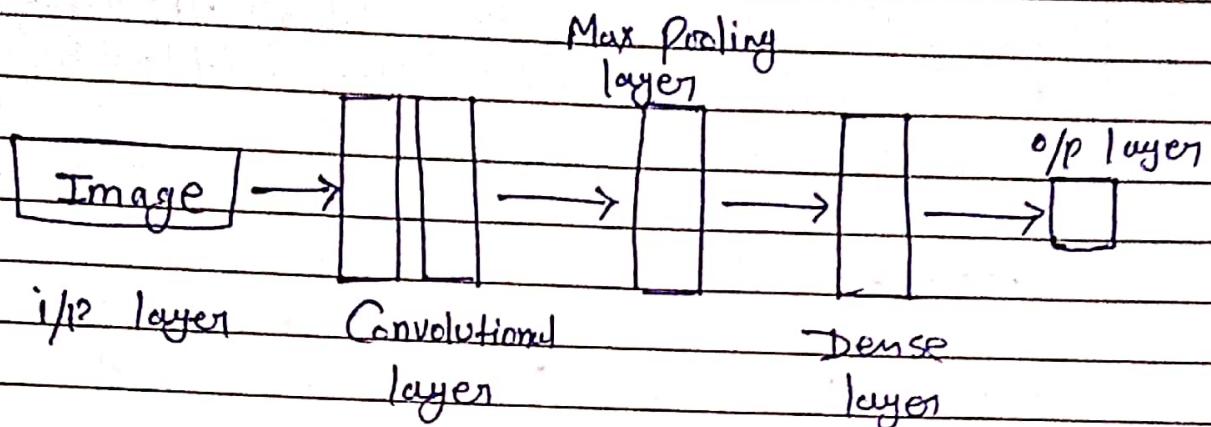


## ② Convolutional neural network :-

CNN structure is focused on processing the grid type data like images & videos by using convolutional layers filtering driving the patterns & spatial hierarchies.

## Application :-

Used for classification of images, object detection, autonomous driving and visualization in augmented reality.

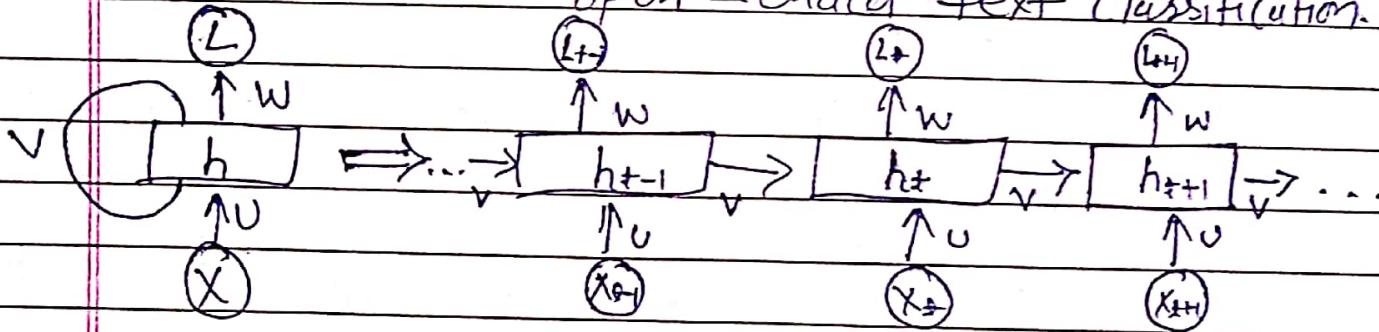


## (B) Recurrent Neural Network (RNN) :-

RNN handles sequential data in which the current o/p is a result of previous i/p by looping over themselves to hold internal state.

## Application :-

language translation, open-ended text classification.

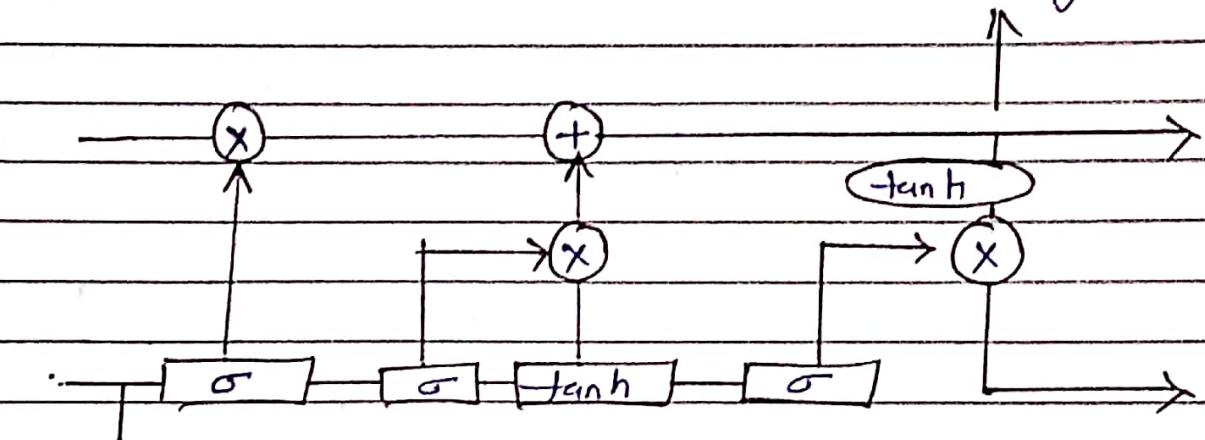


#### ④ LSTM :-

Long - Short - term memory

Networks are a variant of RNNs  
They exhibit memory cells to solve  
the disappearing gradient issue and  
keep large ranges of info in their  
memory.

Applications :- language translation and  
time series forecasting.

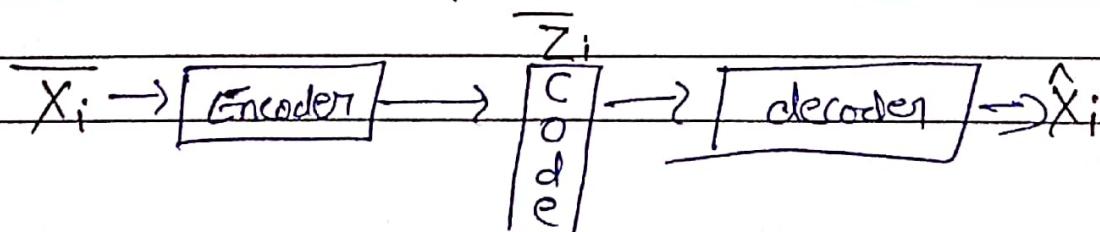


#### ⑤ Autoencoders (AE) :-

Autoencoders are feedforward networks (ANN) that are trained to acquire the most helpful presentations of the info through the process of re-coding the i/p data.

Types :-

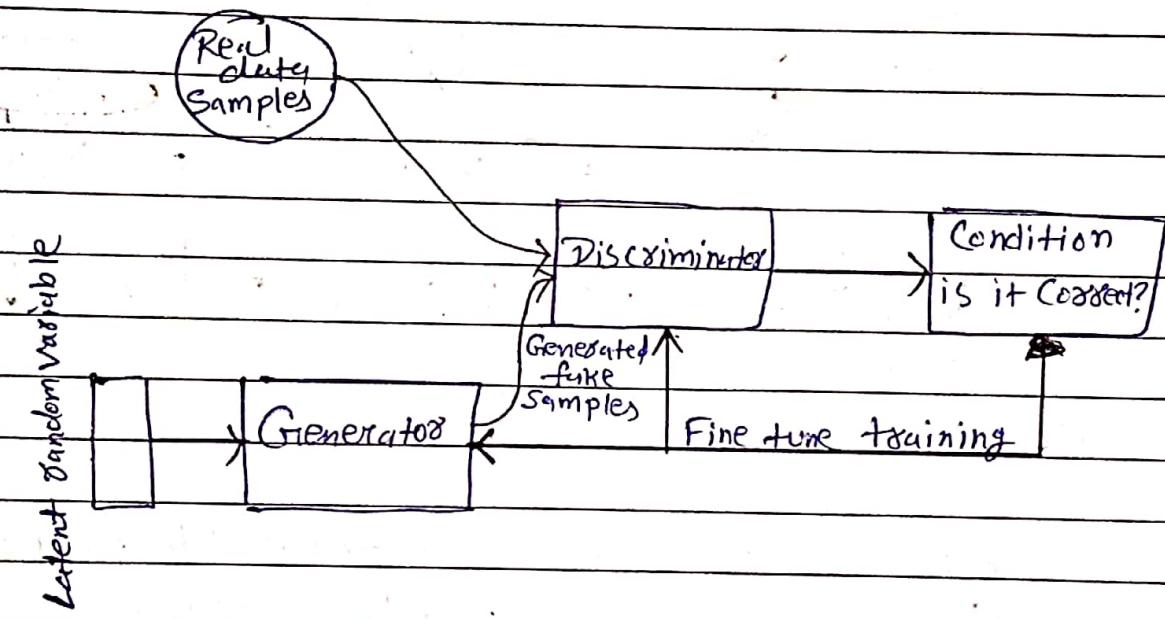
Variants include Undercomplete,  
Overcomplete, and  
Variational autoencoders.



## ⑥ GANs :-

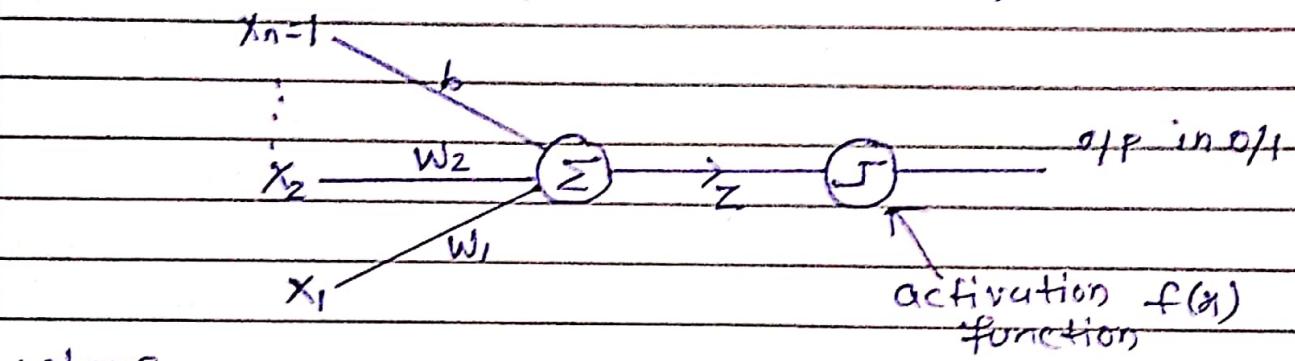
Generative Adversarial Networks has made up of two neural networks, the generator and discriminator, which compete against each other.

The generator creates a fake generated data and the discriminator learns to differentiate the real from and fake data.



## # Perceptron :-

- \* Perceptron is a fundamental unit (<sup>building</sup> block) of NN.
- \* Multilayer Perceptron is also known as ANN.
- \* Basically Perceptron is a algo that is work on supervised learning.
- \* Design / mathematical model / function :-



where,

$x_1$  &  $x_2$  are inputs

$w_1$  &  $w_2$  are weights

$b$  is bias

$$\Sigma = w_1x_1 + w_2x_2 + b$$

$$f(x) = \begin{cases} 1 & \Sigma \geq 0 \\ 0 & \Sigma < 0 \end{cases}$$

- \* Weights Showing how strongly Connected input. (feature importance)

## Geometric Intuition :-

if

$$z = w_1x_1 + w_2x_2 + b$$

$$w_1 \in A, w_2 \in B, b \in C$$
$$x_1 \in X, x_2 \in Y$$

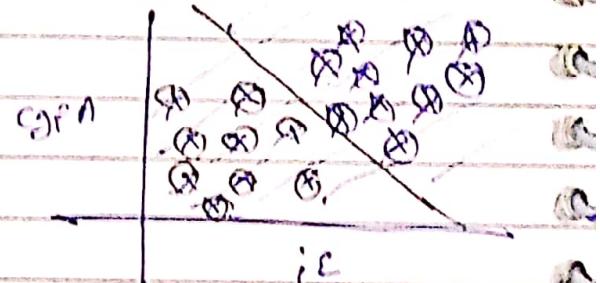
$$(z = Ax + By + C) \text{ & this is}$$

$$Ax + By + C = 0 \text{ --- the eqn of line}$$

$$f(z) = \begin{cases} 1 & z > 0 \\ 0 & z \leq 0 \end{cases}$$

i.e.  $Ax + By + C > 0$  shows it's generate a region.

So perceptron is nothing but line & its work to create a region & divide class.



because of these perceptron is called binary classifier.

in 2D perception work like a line  
in 3D — plane  
in 4D — hyperplane

but always Perception divide data into 2 regions so it is called binary classifier.

Perception can be use only on linear or sort of linear data.

## # Loss function :-

it is a way to define how over model is work (good or bad).

formula of Loss function :

$$L(w_1, w_2, b) = \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) + \alpha R(w)$$

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

$$\rightarrow f(x_i) = w_1 x_1 + w_2 x_2 + b$$

So

$$L(w_1, w_2, b) = \frac{1}{n} \sum_{i=1}^n \max(0, -y_i (w_1 x_1 + w_2 x_2 + b))$$

where

n is no of rows in data

Explanation of loss function :-

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

With the help of eg

	$x_1$	$x_2$	$y$
1	$x_{11}$	$x_{12}$	$y_1$
2	$x_{21}$	$x_{22}$	$y_2$
⋮	⋮	⋮	⋮
$n$	⋮	⋮	⋮

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

So if we have only 2 data than

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

## # McCulloch - Pitts Neural Model :-

The McCulloch and Pitts (MCP) model holds a place of historic significance.

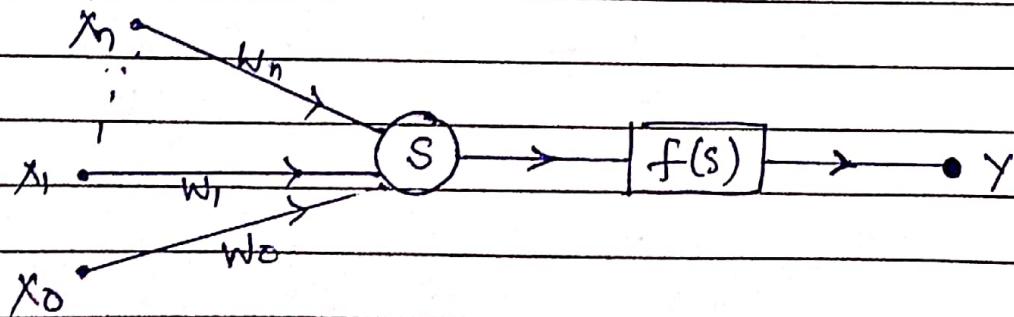
Proposed in 1943 by Warren McCulloch and Walter Pitts, this model laid the foundation framework for understanding how neural network can mimic the brain's processing abilities.

Definition :-

The MCP model is a mathematical representation of a biological neuron.

It Simplifies the complex operations of a neuron into a basic Computational Unit.

Each ~~input~~ neuron receives multiple binary inputs, processes them through weighted summation, and produces a binary o/p based on a threshold function.



## Mathematical foundation :-

The MCP model operates on a straight forward yet powerful set of principles:

### 1. Inputs and Weights :

Each neuron receives binary inputs  $x_1, x_2, \dots$  and assigns weights  $w_1, w_2, \dots$  to these i/p

### 2. Summation :

The neuron computes the weighted sum of its i/p

$$S = \sum_{i=1}^n w_i x_i$$

### 3. Activation Function :

The sum  $S$  is compared to a threshold  $\theta$ . using a step function :

$$y = \begin{cases} 1 & \text{if } S > 0 \\ 0 & \text{if } S \leq 0 \end{cases}$$

## Eg Calculation :-

let's, Consider a neuron with three binary i/p ( $x_1, x_2, x_3$ ) and weights  $w_1 = 0.5, w_2 = 0.3, \text{ and } w_3 = 0.2$  & a threshold  $\theta = 0.7$   
for inputs  $x_1 = 1, x_2 = 1, x_3 = 0$

$$S = w_1x_1 + w_2x_2 + w_3x_3$$

$$S = 0.5 \times 1 + 0.3 \times 1 + 0.2 \times 0$$

$$S = 0.5 + 0.3 = 0.8$$

Since  $0.8 \geq 0.7$  the o/p  $y = 1$

\* The MCP model is one of the first attempts to create an artificial neuron. Think of it as a very simple brain cell in a computer.

It helped lay the groundwork for today's AI, but it has its limits.

Limits of the MCP model :-

① Basic o/p :-

This model only gives two possible results : yes or no  
0 or 1

## ② Simple problems only :

It can only solve problems that are straightforward and can be separated by a single line. More complex issues, like the XOR problem are beyond its reach.

## ③ Fixed weights :

In the RCD model the weight of each IP is fixed.

This means it doesn't learn or get better over time by adjusting these weights.

## ④ No learning ability :

This model doesn't have a way to learn from experience. It can't adjust itself to improve its performance.

## Truth table Example :

Let's implement an AND Logic gate  
Using this model

Conditions :

- Inputs :  $x_1$  and  $x_2$  (binary)
- Weights :  $w_1 = 1$ ,  $w_2 = 1$
- Threshold :  $\theta = 2$

Now apply the model logic :

$x_1$	$x_2$	weighted sum ( $x_1w_1 + x_2w_2$ )	$\text{Sum} \geq \theta$	O/P
0	0	0	No	0
0	1	1	No	0
1	0	1	No	0
1	1	2	yes	1

Ques : Explain the structure and functioning of a biological neuron. How is it modeled in Artificial networks?

Ans : A biological neuron is the basic unit of the human brain responsible for receiving and transmitting info using electrical and chemical signals.

Biological Neural Network (BNN) is a structure that consists of Synapse, dendrites, cell body and axon.

In this neural network the processing is carried out by neurons. Dendrites receive signals from other neurons, Soma sums all the incoming signals and axon transmits the signals to other cells.

Dendrites → Cell body (Soma) → Axon → Synapse  
→ next neuron.

Dendrites : Branch-like structures that receive signals from other neurons

Soma : Also called the cell body processes i/p signals.

Axon : A long fiber that sends signals to other neurons.

Synapse : Gap bet" two neurons where signal transmission occurs

How is it modeled in ANNs ?

Artificial neurons in neural networks are mathematical models that mimic biological neurons.

Diagram from previous notes.

Ques :

What are Computational units in DL?

Ans :

in deep learning, Computational units are the basic building blocks of a neural network.

These are mathematical functions or models that take inputs, process them through a function and produce an o/p.

Think of a computational unit as an artificial version of a biological neuron — it computes an o/p based on the weighted sum of i/p's and an activation function.

Structure of a Computational Unit :

Each Computational Unit typically performs the following Steps:

① Takes inputs :  $x_1, x_2, x_3, \dots, x_n$

② Applies weights and Bias :

$$z = (w_1x_1) + (w_2x_2) + (w_3x_3) + \dots + (w_nx_n) + b$$

- weights define the importance of each i/p
- Bias allows shifting the activation function left or right.

③ Activation function is applied.

This adds non-linearity so the model can learn complex patterns.

$$y = f(z)$$

Mathematical Summary  $\Rightarrow$

$$y = f\left(\sum_{i=1}^n w_i x_i + b\right)$$

Eg of Computational units

1. Perceptron :

The Simplest Computational Unit, proposed by Rosenblatt

$$y = \text{Sign}(w \cdot x + b)$$

it o/p +1 or -1 depending on the weighted sum.

2. Sigmoid Unit:

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

Used in older networks of binary classification problems.

### 3. ReLU Unit :

$$y = \max(0, z)$$

Fast, efficient, and widely used in CNNs, RNNs etc.

### 4. LSTM Unit :

Advanced Computational Unit for Sequential data.

it contains Gater, Memory Cell Helps remember long-term dependencies.

### 5 Convolutional Unit :

Used in CNNs for images

- it performs Convolutional operations
- Helps extract features like edges, textures.

Ques What is a linear perceptron? Derive its learning ~~rule~~ rule.

The linear perceptron is the simplest form of a neural network - a single-layer binary classifier.

It is designed to separate i/p data into two classes using a linear decision boundary (a line, plane or hyperplane depending on the no of i/p features).

The perception learn this boundry by adjusting its internal weights during training so that it can make correct classification.

It is suitable only for linearly separable datasets.

perception diagram

Derivation of the learning rule :

To adjust the weights in such a way that the perception reduces the classification error & makes better

Predictions over time.

Let's

inputs  $x = [x_1, x_2, \dots, x_n]$

weight vector  $w = [w_1, w_2, \dots, w_n]$

Bias :  $b$

Target o/p :  $y \in \{-1, 1\}$

Predicted o/p :  $\hat{y} = \text{Sign}(w \cdot x + b)$

Learning rate :  $\eta$

Step 1 : Calculate error

$$\text{Error} = y - \hat{y}$$

Step 2 : Update Rule for weights

If the prediction  $\hat{y}$  is wrong, update the weights  $w$ :

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

Step 3 : Update Rule for Bias

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

## Convergence Theorem :-

If the training data is linearly separable the perceptron learning algo is guaranteed to converge to a sol' in a finite no of steps.

If the data is not linearly Separable (like XOR problem) it will never Converge.

Ques Explain the Perceptron learning algo with steps and example ?

The Perceptron learning algo is used to train a single-layer perceptron to classify input data into two classes (+1 or -1)

It works by updating weights based on misclassifications.

The aim is to find the best-fit hyperplane that separates the data.

Algorithm STEPS :-

① Initialize the

Perceptron is a linear Supervised ML algo  
it is used for binary classification

Perceptron learning algo is also understood  
as an artificial neuron or neural  
network unit that helps to detect  
certain i/p data computations in  
business intelligence.

The perceptron learning algorithm  
is treated as the most straight  
forward ANN. Hence, it is a single  
layer neural network with four  
main parameters, i.e., input values,  
weights and bias, net sum and  
an activation function.

There are 4 significant steps in  
a PLA :-

- ① First, multiply all i/p values with  
corresponding weight values and  
then add them to determine  
the weighted sum. Mathematically.

$$\sum w_i \cdot x_i = w_1x_1 + w_2x_2 + \dots + w_n x_n$$

Then add another essential term  
called bias 'b' to the weighted sum  
to improve the model performance.

$$\sum w_i \cdot x_i + b$$

- ② An activation function is applied to this weighted sum, producing a binary or a continuous valued o/p

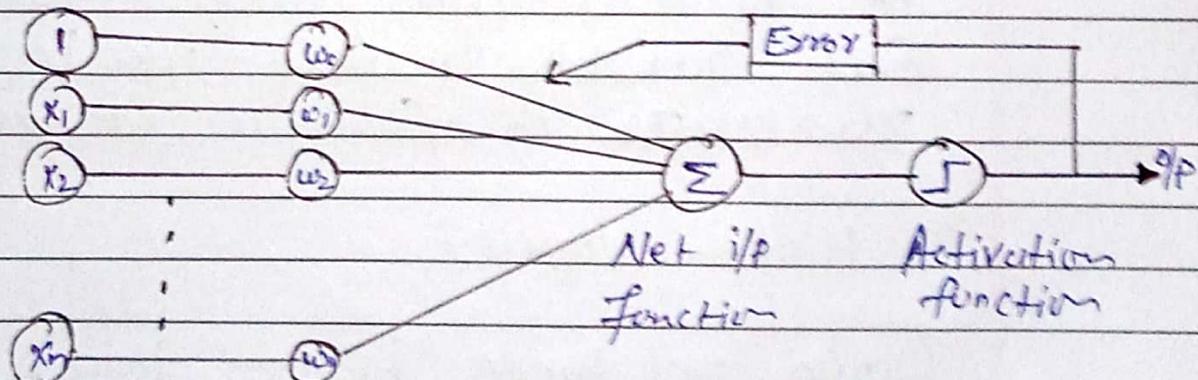
$$y = f(\sum w_i \cdot x_i + b)$$

- ③ Now the diff. betw this o/p and the actual target value is computed to get the error term, E. generally in terms of mean squared error.

$$E = (y - y_{\text{actual}})^2$$

- ④ we optimize this error (loss function) using an optimization algo.

Generally some form of GDA (Gradient Descent Algo) is used to find the optimal values of the hyperparameters like learning rate, weight, bias etc.



## # Feed Forward Network :-

Feed forward Neural Network (FFNN) is a type of ANN in which info flows in a single direction from the i/p layer through hidden layers to the o/p layer. - without loops or feedback.

It is mainly used for pattern recognition tasks like image & speech classification.

Structure of FFNN :-

FFNN have a structured layered design where data flow sequentially through each layer.

① Input layer :

The i/p layer consists of neurons that receive the i/p data. Each neuron in the i/p layer represents a feature of the i/p data.

② Hidden layer :

One or more hidden layers are placed bet<sup>n</sup> the i/p & o/p layers.

These layers are responsible for

learning the complex patterns in the data.

Each neuron in a hidden layer applies a weighted sum of i/p followed by a non-linear activation function.

⑧ O/P layer :

The O/P layer provides the final O/P of the network.

Diagram

Activation function  $\circ \rightarrow$

Common activation functions :

$\rightarrow$  Sigmoid :

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

$\rightarrow$  Tanh :

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

$\rightarrow$  ReLU :

$$\text{ReLU}(x) = \max(0, x)$$

## # Forward Propagation :

in Forward propagation i/p data moves through each layer of neural network where each neuron applies weighted sum, adds bias, passes the result through an activation function & making predictions.

~~This~~ This process is crucial before backpropagation updates the weights.

it determines the o/p of neural network with a given set of input and current state of model parameters (weights and bias).

Mathematical explanation of Forward Pn:  
let's target o/p is 0.5 & the learning rate is 1.

### ① Initial Calculation

The weighted sum at each node is calculated using:

$$a_j = \sum (w_{i,j} * x_i)$$

Where,

$a_j$  = is the weighted sum of all the i/p weights at each node.

- $w_{ij}$  represents the weights bet<sup>n</sup> the  $i^{th}$  i/p and  $j^{th}$  neuron.
- $x_i$  represents the value of the  $i^{th}$  i/p

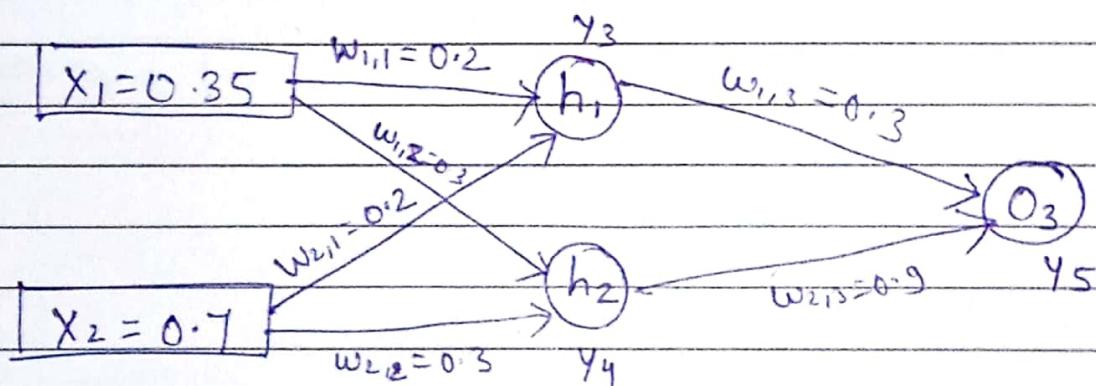
→ After applying the activation function to  $a_i$ , we get the o/p of the neuron

$$o_j = \text{activation function}(a_j)$$

## ② Sigmoid function :

The Sigmoid function returns a value bet<sup>n</sup> 0 and 1, introducing non-linearity into the model.

$$y_j = \frac{1}{1 + e^{-a_j}}$$



## ③ Computing Outputs :

At  $h_1$  node

$$a_1 = (w_{1,1} \cdot x_1) + (w_{2,1} \cdot x_2)$$

$$a_1 = (0.2 \times 0.35) + (0.2 \times 0.7)$$

$$a_1 = 0.21$$

Once we calculated the  $a_1$  value, we can now proceed to find  $y_3$

$$y_j = F(a_{1j}) = \frac{1}{1 + e^{-a_1}}$$

$$y_3 = \cancel{F(0.21)}$$

$$y_3 = \frac{1}{1 + e^{-0.21}} = 0.56$$

Similarly find the values of  $y_4$  at  $h_2$  and  $y_5$  at  $0_3$

$$a_2 = (w_{1,2} \times 2_1 + w_{2,2} \times 2_2) = (0.3 \times 0.85 + 0.3 \times 0.7)$$

$$a_2 = 0.315$$

So

$$y_4 = F(0.315) = \frac{1}{1 + e^{-0.315}} = 0.59$$

&

$$\boxed{y_5 = 0.67}$$

#### ④ Error Calculation

Let our actual O/P is 0.5 but we obtained 0.67.

$$\begin{aligned} \text{Error} &= y_{\text{target}} - y_5 \\ &= 0.5 - 0.67 \end{aligned}$$

$$\text{Error} = -0.17$$

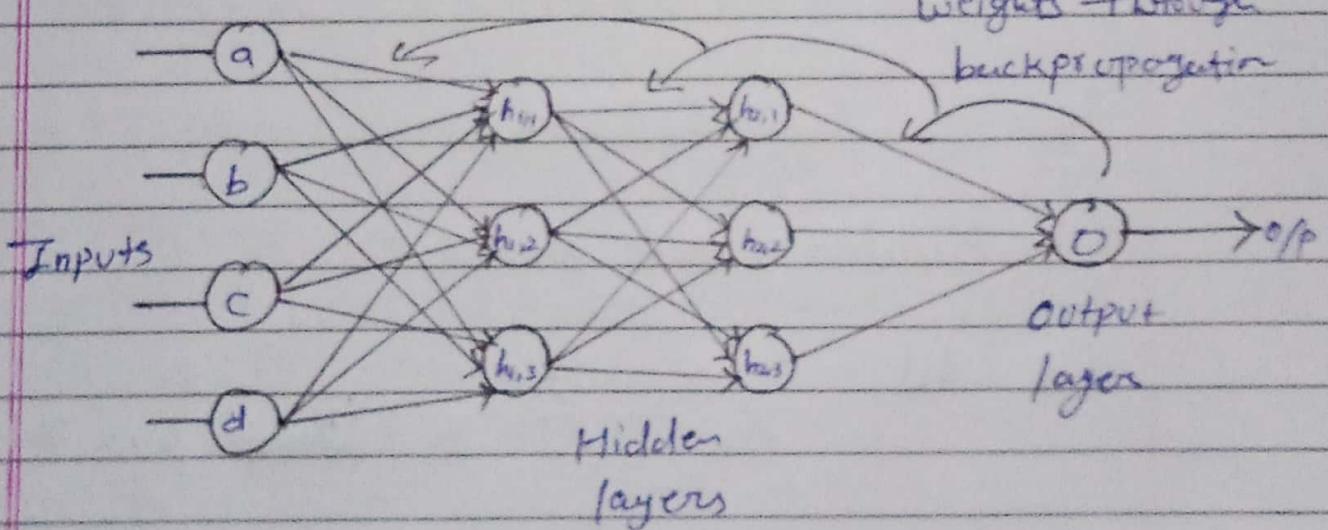
## # Backpropagation :

Backpropagation is a technique used in DL to train ANNs particularly Feed forward networks.

It works iteratively to adjust weights and bias to minimize the cost function

Backpropagation often uses optimization algo's like Gradient Descent or Stochastic gradient descent.

The algo computes the gradient using the chain rule from Calculus allowing it to effectively navigate complex layers in the neural network to minimize the cost function.



## Backpropagation :-

### ① Calculating Gradients $\Rightarrow$

The change in each weight is calculated as:

$$\Delta w_{ij} = \eta \times \delta_j \times o_j$$

Where,

- $\delta_j$  is the error term for each i/p
- $\eta$  is the learning rate.

### ② O/P unit Error $\Rightarrow$

For  $O_3$ :

$$\begin{aligned}\delta_5 &= y_5(1-y_5)(y_{\text{target}} - y_5) \\ &= 0.67(1-0.67)(0.5-0.67) \\ &= -0.0376\end{aligned}$$

### ③ Hidden Unit Error $\Rightarrow$

For  $h_1$ :

$$\begin{aligned}\delta_3 &= y_3(1-y_3)(\cancel{y_{\text{target}}} + w_{1,3} \times \delta_5) \\ &= 0.56(1-0.56)(0.3 + (-0.0376)) \\ &= -0.0027\end{aligned}$$

For  $h_2$ :

$$\begin{aligned}\delta_4 &= y_4(1-y_4)(w_{2,3} * s_5) \\ &= 0.59(1-0.59)(0.9 \times (-0.0376)) \\ &= -0.0819\end{aligned}$$

#### ④ weight updates

For the weights from hidden to o/p layer

$$\Delta w_{2,3} = 1 * (-0.0376) \times 0.59 = -0.022184$$

New weight:

$$w_{2,3}(\text{new}) = -0.022184 + 0.9 = 0.877816$$

Similarly other weights are updated:

- $w_{1,1}(\text{new}) = 0.200945$
- $w_{1,2}(\text{new}) = 0.273225$
- $w_{1,3}(\text{new}) = 0.086615$
- $w_{2,1}(\text{new}) = 0.269445$
- $w_{2,2}(\text{new}) = 0.18534$

After updating the weights the forward pass is repeated:

$$y_3 = 0.57$$

$$y_4 = 0.56$$

$$y_5 = 0.61$$

So now

$$\text{Error} = y_{\text{target}} - y_5$$

$$= 0.5 - 0.61$$

$$= -0.11$$

## Unit 02

Date \_\_\_\_\_  
Page \_\_\_\_\_

### # Multilayer Perceptron :-

Multilayer perceptron (MLP) is an ANN widely used for solving classification & regression tasks.

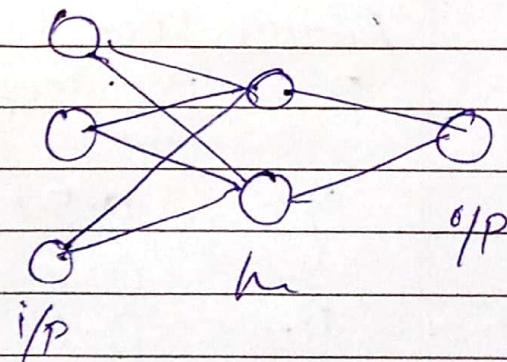
MLP consists of fully connected dense layers that transform i/p data from one dimension to another.

It is called "multi layer" because it contains an input layer, one or more hidden layers and an o/p layers.

The purpose of an MLP is to model complex relationships bet" i/p & o/p making it a powerful tool for various ML tasks.

Key Components (Architecture) of MLP :-

- i/p layer
- hidden layers
- o/p layer



## Working of MLP :-

Step 1 : Forward Propagation

→ Weighted Sum

→ Activation function

Step 2 : loss function Calculate

Step 3 : Backpropagation

→ Gradient Calculation

→ Error Propagation

→ Gradient Descent

Step 4 : Optimization

## Applications of MLP :-

→ Tabular Data Classification & Regression

→ Handwritten Digit and Character Recognition

→ Function Approximation & Control

→ Time Series forecasting

→ Feature Extraction & pre-training

→ Anomaly Detection

Advantages :-

- Versatility
- Non-linearity
- Parallel Computation

Disadvantages :-

- Computationally Expensive
- Prone to Overfitting
- Sensitivity to Data Scaling

## # Gradient Descent Algo :-

Gradient Descent is an optimization algo used in ML to minimize the cost function by iteratively adjusting parameters in the direction of the negative gradient, aiming to find the optimal set of parameters.

The cost function represents the discrepancy bet' the predicted o/p of the model and the actual o/p.

Gradient descent aims to find the parameters that minimize the discrepancy and improve the model performance.

The algo operates by calculating the gradient of the cost function, which indicates the direction and magnitude of the ~~cost~~ function steepest ascent.

By iteratively updating the model's parameters in the negative gradient direction, gradient descent gradually converges towards the optimal set of parameters that yields the lowest cost.

Gradient descent can be applied to various ML algo including linear regression, logistic regression, neural networks & SVM.

Cost function :

It is a function that measures the performance of a model for any given data.

Cost function quantifies the error b/w predicted values and expected values and presents in the form of a real no.

Hypothesis :  $h_{\theta}(x) = \theta_0 + \theta_1 x$

Parameters :  $\theta_0, \theta_1$

Cost function :  $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

Goal : minimize ( $J(\theta_0, \theta_1)$ )

Given a cost function  $J(\theta)$ , where  $\theta$  represents the model parameters, the update rule for gradient descent is:

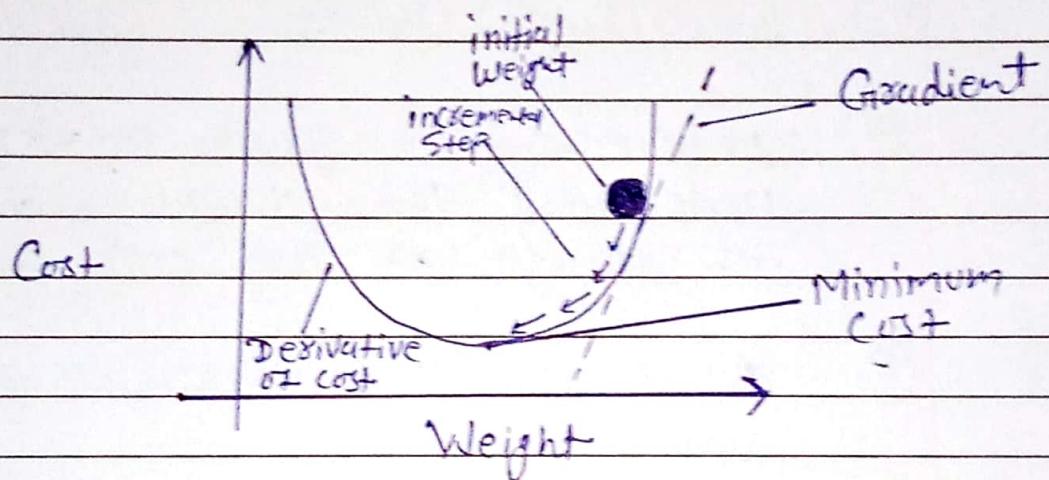
$$\theta := \theta - \alpha \nabla J(\theta)$$



Where.

$\theta$  are the parameters  
 $\alpha$  is learning rate  
 $\nabla J(\theta)$  is the gradient of the Cost function with respect to  $\theta$ .

Gradient descent was originally proposed by CAUCHY in 1847. it is also known as steepest descent.



The goal of the GDA is to minimize the Cost function. To achieve this goal, it performs two steps iteratively:

- ① Compute the gradient
- ② Make a move in the opposite direction of the gradient.

## GDAs steps :

- ① Start with initial guesses for the parameters  $\theta$ .
- ② Use the current parameter to make predictions.
- ③ Evaluate the cost function  $J(\theta)$  to measure the error between prediction & actual values.
- ④ Determine the gradient  $\nabla J(\theta)$  to understand the direction and rate of change of the cost function.
- ⑤ Adjust the parameters using the update rule:  
$$\theta := \theta - \alpha \nabla J(\theta)$$
- ⑥ Repeat step 2-5 until convergence.

## Variants of GD :

→ Batch Gradient Descent

→ Stochastic Gradient Descent (SGD)

→ Mini-Batch

## # Batch, Stochastic & Mini-batch GD :-

### ① BGD :-

Batch Gr. D. Computes the gradient using the entire dataset before updating model parameters.

While ~~this~~ this approach provides accurate gradient estimates, it can be computationally expensive for large datasets.

- Pros  $\Rightarrow$  Accurate updates, Stable convergence
- Cons  $\Rightarrow$  High memory consumption, slow for large datasets.

### ② Stochastic Gr. D. :-

SGD updates model parameters after computing the gradient from a single data point. This approach is computationally cheaper but introduces noise, leading to an unstable convergence path.

- Pros  $\Rightarrow$  Fast updates, works well for online learning
- Cons  $\Rightarrow$  High variance

## (3) MBGD :-

MBGD combines the benefits of both batch and stochastic gradient descent. Instead of using the entire dataset or a single sample, MBGD processes small batches of data at a time. This balances computational efficiency & model stability.

**Pros :-** Faster training than batch gradient descent, more stable than SGD.

**Cons :-** Choosing the right batch size can be challenging.

Aspect	TGID	SGD	MBGD
Data processed per update	Entire training dataset	Single training ex	Small subset of training data
Update frequency	Once per epoch	After each training ex	Affter each mini batch
Convergence Speed	Can be slow due to processing the entire dataset	Typically faster but may fluctuate	Faster & more stable than SGD

Convergence Stable Convergence May oscillate More stable than  
Stability to the min. around the min SGD, less than  
due to noise Batch.

Memory requirement High Low Moderate

Computational less efficient Highly eff. balance eff. & performance  
Efficiency for large dataset

Typical Small to medium scale or ML applications.  
use cases Size dataset online learning  
where precision Scenarios  
is crucial

## # Empirical risk minimization :-

In SL, the goal is to find ~~a function~~  
a function  $f$  that maps inputs  $x$   
to o/p  $y$  such that the predictions  
 $f(x)$  are as close as possible to the  
actual o/p's.

The 'risk' or 'loss' quantifies the  
discrepancy betw predicted & actual  
o/p typically define as :

$$R(f) = E_{(x,y)} \sim P [l(f(x), y)]$$

Here,  $l$  is a loss function &  $P$  is the true but unknown data distribution.

Since  $P$  is unknown, we approximate the risk using the empirical risk over a training dataset  $\{(x_i, y_i)\}_{i=1}^n$ :

$$\hat{R}_n(f) = \frac{1}{n} \sum_{i=1}^n l(f(x_i), y_i)$$

ERM involves selecting the function  $f$  from a hypothesis space  $F$  that minimizes this empirical risk:

$$\hat{f}_n = \arg \min_{f \in F} \hat{R}_n(f)$$

This approach assumes that minimizing the loss on training data will lead to good performance on unseen data.

Applications :-

- SVM
- Linear Regression
- logistic regn

Advantages :-

- Theoretically Sound
- Wide Applicability
- Flexibility

Challenges :-

- Overfitting
- No Guarantee of Generalization
- Sensitive to Noise

## III Regularization :-

Regularization in DL is a technique used to prevent overfitting by adding a penalty term to the loss function, discouraging the model from becoming overly complex.

This helps the model generalize better to unseen data.

Overfitting happens when a model learns the training data too well, including the noise and outliers which causes it to perform poorly on new data.

Types of Regularization :-

## ① Lasso Regression :-

A regression model which uses the L1 Regularization technique is called LASSO (Least Absolute Shrinkage and Selection Operator) regression.

Lass regression adds the "absolute value of magnitude" of the Coefficient as a penalty term to the loss function (L).

$$\text{Cost} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{i=1}^m |w_i|$$

where ,

$m$  = no of features

$n$  = no of Eg

$y_i$  = Actual Target Value

$\hat{y}_i$  = Predicted Target Value

## ② Ridge Regression :-

A regression model that uses the L2 regularization technique is called Ridge regre .

R.R. adds the "squared magnitude" of the Coefficient as a penalty term to the loss function (L).

$$\text{Cost} = \frac{1}{n} \sum_{i=1}^m (y_i - \hat{y}_i)^2 + \lambda \sum_{i=1}^m w_i^2$$

Where,

$m = \text{---}$ ,  $n = \text{---}$ ,  $y_i = \text{---}$ ,  $\hat{y}_i = \text{---}$   
 $w_i = \text{Coefficients of the features}$

$\lambda = \text{Regularization parameter}$

### ③ Elastic Net Regression :

Elastic Net Regression is a Combination of both L1 as well as L2 regularization

That implies that we add the absolute norm of the weights as well as the Squared measure of the weights.

With the help of an extra hyperparameter that Controls the ratio of the L1 & L2 regular

$$\text{Cost} = \frac{1}{n} \sum_{i=1}^m (y_i - \hat{y}_i)^2 + \lambda ((1-\alpha) \sum_{i=1}^m |w_i| + \alpha \sum_{i=1}^m w_i^2)$$

where,

$\alpha = \text{Mixing Parameter (where } 0 \leq \alpha \leq 1\text{)}$

$\rightarrow \alpha = 1$  Corresponds to L1 Regular

$\rightarrow \alpha = 0$   $\longrightarrow$  L2

$\rightarrow$  Values bet<sup>n</sup> 0 & 1 provide a balance of both L1 & L2 regular

## # Bias :

Bias measures the error introduced by approximating a real world problem by a much simpler model.

## # Variance :

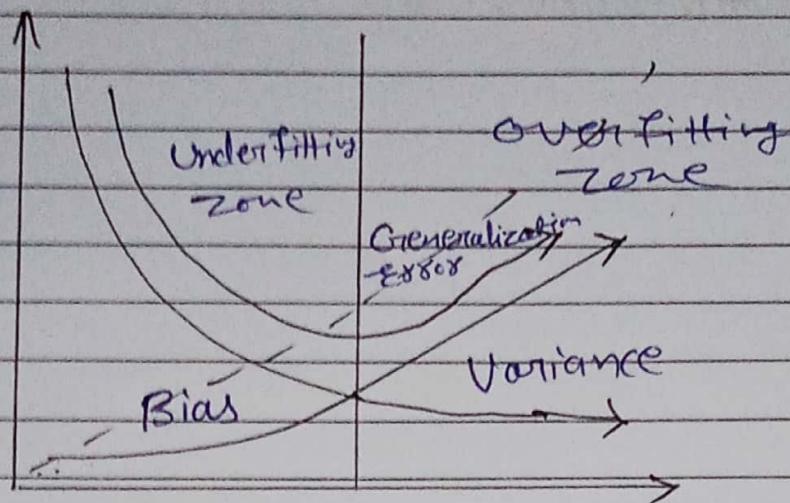
Variance measures how much function would change if we estimated it using a different training dataset drawn from the same distribution.

## # Bias Variance Tradeoff :

The Bias - Variance Tradeoff is a fundamental concept in ML.

It refers to the balance b/w bias & variance which affect predictive model performance.

The bias - Variance tradeoff demonstrates the inverse relationship b/w bias & variance. When one decreases the other tends to increase & vice versa.



### Benefits of Regularization :

- Prevents Overfitting
- improves interpretability
- Enhances Performance
- Stabilizes models
- Prevents Complexity
- Allows Fine - Tuning

## # Autoencoders :-

An autoencoder is a type of artificial neural network that learns to represent data in a compressed form and then reconstructs it as closely as possible to the original i/p.

Autoencoder consists of two Components:

- Encoder  $\Rightarrow$

This compresses the i/p into a compact representation & capture the most relevant features.

- Decoder  $\Rightarrow$

It reconstructs the i/p data from this compressed form to make it as similar as possible to the original i/p.

$\rightarrow$  Autoencoder aim to minimize reconstruction error which is the difference b/w the i/p & the reconstructed o/p.

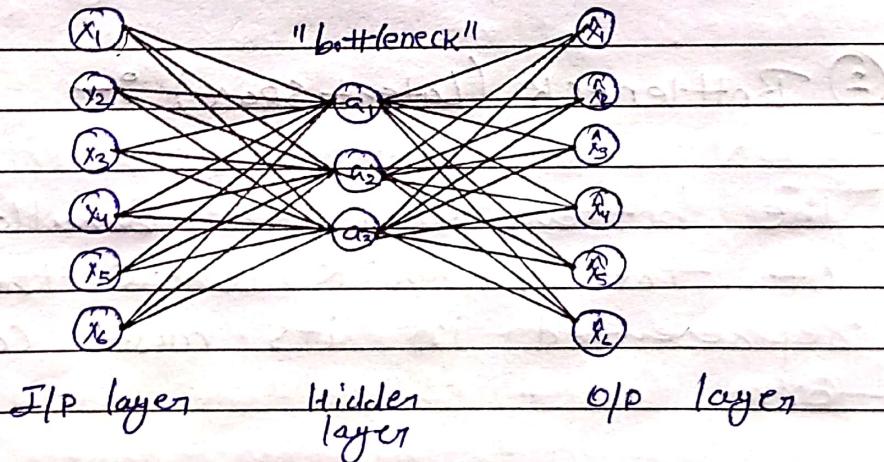
They use loss functions such as Mean Squared Error (MSE) or Binary Cross Entropy (BCE) and optimize through backpropagation and gradient descent.

They are used in applications like image processing, anomaly detection, noise removal and feature extraction.

Architecture of autoencoder  $\Rightarrow$

The architecture of an autoencoder consists of 3 main components:

The encoder, the bottleneck (latent space) and the Decoder.



1. Encoder  $\Rightarrow$

The encoder is the part of the network that takes the I/P data and compresses it into a smaller/lower dimensional representation.

• I/P layer  $\Rightarrow$

This is where the original data enters the network.

- Hidden layer  $\Rightarrow$

These layers apply transformations to the i/p data.

The encoder's goal is to extract essential features and reduce the data's dimensionality.

- O/P of Encoder (Latent Space)  $\Rightarrow$

The encoder o/p a compressed version of the data often called the latent representation or encoding.

### (2) Bottleneck (Latent Space) :-

The bottleneck is the smallest layer of the network where the data is represented in its most compressed form

### (3) Decoder :-

The decoder is responsible for taking the compressed representation from the latent space & reconstructing it back into original data form.

## Loss function in autoencoder Training :-

During training an autoencoder aims to minimize the reconstruction loss which measures the difference b/w the original i/p & the reconstructed o/p. The choice of loss function depends on the type of data:

### \* Mean Squared Error (MSE) :-

This is commonly used for continuous data. It measures the average squared differences b/w the i/p and the reconstructed data.

### \* Binary Cross-Entropy :-

Used for binary data (0 or 1 value). It calculates the diff. in probability b/w the o/p & the d/p.

## Types of autoencoders :-

→ Denoising Autoencoder

→ Sparse

→ Variational

→ Convolutional

## Unit :- 3

### Convolutional Networks

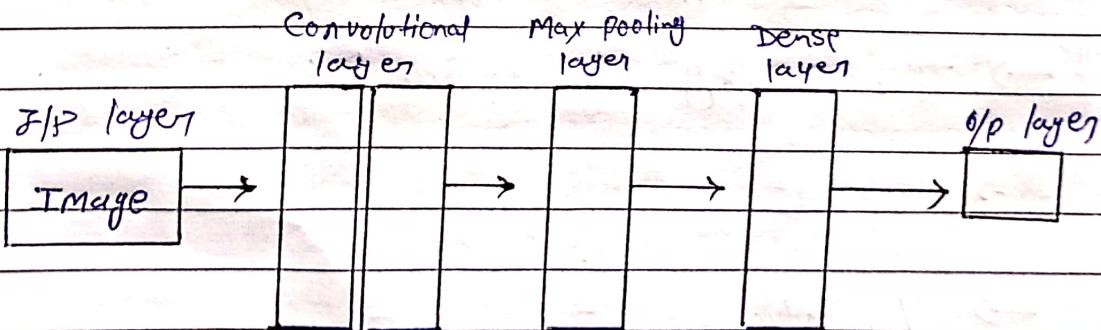
# CNN :-

Convolutional Neural Network (CNN) is an advanced version of ANN, primarily designed to extract features from grid-like matrix datasets.

This is particularly useful for visual datasets such as images or videos, where cluster patterns play a crucial role.

CNN are widely used in Computer Vision applications due to their effectiveness in processing visual data.

CNNs consist of multiple layers like the I/P layer, Convolutional layer, Pooling layer & fully Connected layers.



### Advantages :-

- Good at detecting patterns & features in images, Videos and audio Signals.
- Robust to translation, rotation and Scaling invariance.
- End-to-end training, no need for manual feature extraction.
- Can handle large amount of data & achieve high accuracy.

### Disadvantages :-

- Computationally expensive to train & require a lot of memory.
- Can be prone to overfitting if not enough data or proper regularization is used.
- Requires large amount of labeled data.
- Interpretability is limited, it's hard to understand what the network has learned.

## # Convolutional operation :-

What is Kernel / Filter =?

The Kernel is a rectangular small matrix, which slides over the image from left to right and top to bottom.

$$\begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix}$$

What is Stride ?

The no of pixels we slide over the i/p image by the Kernel is called a Stride.

What is Conv Opn =>

The Convolutional Operation is the process of implying a Combination of two functions that produce the third function as a result, employing filter across the entire i/p image allows the filter to discover that feature of the image, which is also called a feature map.

A jigsaw puzzle is a perfect example of a Convolutional operation.

In a jigsaw puzzle, each piece has a portion of an img that reveals something about the complete picture when assembled.

Like a jigsaw puzzle in Convolutional networks, multiple filters are taken to slice through the img and map them one by one and learn different portions of an input image.

Padding :

Adding extra borders around the i/p to control o/p size

## II Variants of basic Convolution function :

There are several variants of the basic Convolution function that are commonly used in DL, especially in CNNs.

Each variant is designed to improve performance, reduce computation or capture different features.

→ Dilated Convolution (Atrous Convolution)

→ Transposed Convolution

- Depthwise Separable Convolution
- Grouped Convolution
- $1 \times 1$  Convolution (Pointwise Convolution)
- Circular Convolution (Periodic padding)
- Dynamic / Deformable Convolution

## # Structured O/P in CNN :-

Convolutional Neural Networks are not just for classifying images (like saying "this is a cat").

They can also be trained to produce detailed, structured O/P - where the prediction isn't a single label but something more complex.

A common ex is image segmentation, it's goal to label every single pixel in the img with a class

That means the O/P of the model looks like the i/p image in size but each pixel contains a probability.

for each possible class.

How to keep o/p same size as i/p :

To make the o/p map the same size as the i/p, CNNs use same padding which means adding zeros around the i/p so that each convolution doesn't shrink it.

Smoothing the labels :

This is where the graphical models come in — they model the relationships b/w pixels.

A popular CNN architecture for Segmentation is the U-Net, especially in medical imaging. It's shaped like a "U" because it first shrinks the img to understand the big picture and then grows it back to get precise, pixel-level o/p.

## # Efficient Convolution Algorithm :

### ① Depthwise Separable Convolution :

Depthwise Separable Convolution is a technique used in CNN to reduce the Computational Cost and the no of parameters compared to standard Convolution operations. This technique is particularly useful in mobile-vision applications where Computational resources are limited.

Depthwise separable Convolutions break down the Convolution operation into two separate steps:

### ② Depthwise Convolution :

In depthwise Convolution, each i/p channel is convolved with a separate filter of size

Convolution is a very imp mathematical op. in ANN. CNNs can be used to learn features as well as classify data with the help of img frames.

There are many types of CNN's.

One class of CNN's are depth wise separable convolutional neural network. These type of CNN's are widely used because of the following two reasons :-

- They have lesser no of parameters to adjust as compared to the standard CNN's, which reduces overfitting.
- They are computationally cheaper because of fewer computations which makes them suitable for mobile vision app.

- - - - -

## # LeNet :-

LeNet is a pioneering Convolutional neural network architecture that laid the foundation for modern image recognition, introducing key concepts like Convolution, Pooling & hierarchical feature learning.

Overview of LeNet =>

LeNet, particularly LeNet-5, was developed by Yann LeCun and his collaborators at AT&T Labs in the late 1980s & early 1990s, during a time when traditional methods heavily relied on hand-crafted features for image processing.

LeNet was specifically designed for handwritten digit recognition, making it one of the first significant applications of CNN in practical scenarios such as banking & mail processing.

### Architecture of LeNet-5 :-

The Architecture of LeNet-5 consists of multiple layers that enable effective

feature extraction from images.

The primary layers include:

① Input layer  $\Rightarrow$

Accepts grayscale images  
typically sized at  $32 \times 32$  pixels.

② Convolutional layer ( $C_1$  &  $C_3$ )  $\Rightarrow$

The first Convolution layer ( $C_1$ ) uses 6 filters of size  $5 \times 5$ , producing  $28 \times 28$  feature maps.

The Second Conv layer ( $C_3$ ) applies K filters of the same size, outputting  $10 \times 10$  feature maps.

These layers automatically learn spatial hierarchies in the i/p data.

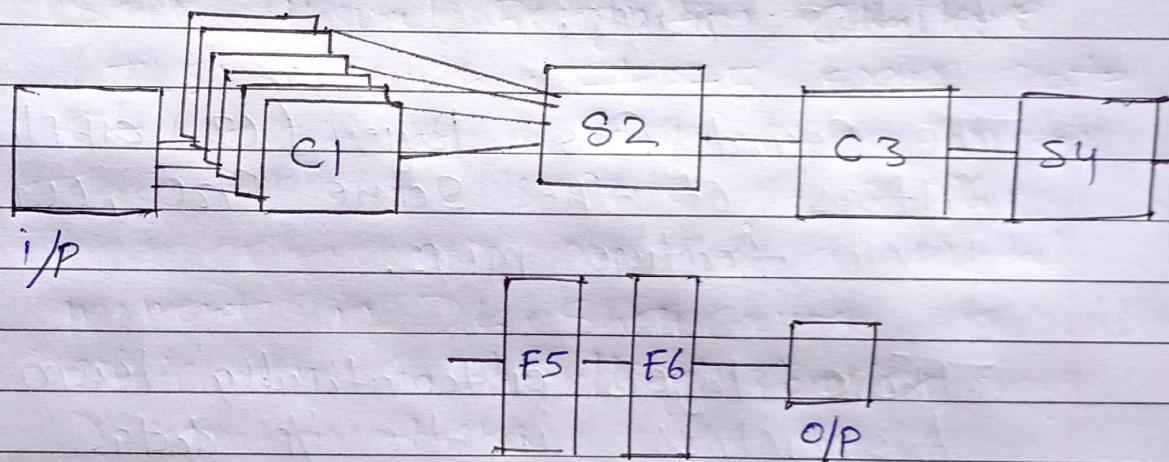
③ Pooling layers ( $S_2$  &  $S_4$ )  $\Rightarrow$

After each Convolutional layer, average pooling layers reduce the dimensionality of the feature maps, increasing translation invariance and decreasing computational load.  $S_2$  reduces the size to  $14 \times 14$  &  $S_4$  further reduces it to  $5 \times 5$ .

#### (4) Fully Connected layers (F5 & F6) :-

- The o/p from the pooling layers is flattened and passed through two Fully Connected layers, with 120 & 84 nodes, respectively.

The final layer o/p probabilities for each of the 10 classes using a Softmax activation function.



Significance in DL  $\Rightarrow$

LeNet's intro popularized several key innovations within CNNs, including:

$\rightarrow$  Convolution op  $\Rightarrow$

Allowing the network to automatically learn & extract relevant features from images.

→ Pooling Mechanism ⇒

Reducing the dimensionality of feature maps while retaining essential info, aiding in building translational invariance.

→ Hierarchical Feature Learning ⇒

Enabling layers to learn from raw pixel values to a ~~search~~ rich & abstract representation of visual data, which mimics human visual perception.

### Applications :-

→ Handwriting recognition in postal services & banking

→ Object & face recognition in img & videos.

→ Autonomous driving Systems for interpreting road signs & obj on the road.

## # AlexNet :-

AlexNet is a groundbreaking CNN architecture that revolutionized image classification by significantly outperforming previous models, highlighting the potential of DL.

Introduced in 2012 by Alex Krizhevsky, Ilya Sutskever & Geoffrey Hinton, AlexNet won the ImageNet Large Scale Visual Recognition Challenge (ILSVRC) with a top-5 error rate of 15.3% dramatically outperforming the runner-up's rate of 26.2%. This achievement marked a pivotal moment in DL, showcasing its efficiency in handling complex image recognition tasks.

### Key innovations :-

- ReLU Activation  $\Rightarrow$

Accelerated training by 6 times compared to tanh or Sigmoid function.

- Dropout  $\Rightarrow$

Prevented overfitting by randomly

dropping neurons during training.

- Data Augmentation  $\Rightarrow$

Enhanced model generalization by artificially increasing the no of images through techniques like img translation & reflections.

- Architecture  $\Rightarrow$

- ① i/p layer  $\Rightarrow$

Accepts RGB images sized  $227 \times 227$  pixels.

- ② Convolutional layers  $\Rightarrow$

The initial layer uses 96 filters of size  $11 \times 11$  followed by layers with increasing filter counts but smaller sizes, going down to  $3 \times 3$  in deeper layers. This design helps extract essential features from images effectively.

- ③ Pooling layers  $\Rightarrow$

Max - pooling is employed to downsample feature maps, reducing dimensionality while retaining Critical Spatial info.

## ④ Fully Connected Layers :-

The last 3 layers are fully connected culminating with a Softmax layer that o/p probabilities for 1000 different classes.

### Impact & Legacy :-

AlexNet was a big breakthrough in DL & inspired many other models like VGGNet, GoogleNet, and ResNet.

It made Computer vision easier by letting models learn imp features from data on their own, instead of relying on people to manually pick out those features.

### Applications :-

- Medical imaging
- Autonomous Vehicles
- Robotics

## Unit 8 4

## RNN

Recurrent Neural Network (RNNs) are specialized neural networks that excel in handling sequential data such as text, audio, or time-series.

Unlike traditional feedforward networks, RNNs have connections that loop back, allowing them to use o/p from previous time steps as i/p for future steps, thus integrating context from past data into current predictions.

Key Components & Architecture :-

- Memory  $\Rightarrow$

The core feature of RNNs is their internal memory, which allows them to remember previous i/p's.

- Recurrent Connections  $\Rightarrow$

RNNs are structured with loops in their architecture. At each time step, they process current i/p along with info retained from previous i/p stored in hidden states.

- Hidden States  $\Rightarrow$

This is a memory of

the network, continuously updated as the network processes each sequence element.

Types of RNNs :-

- One - to - One (A single i/p produce single o/p)
- One - to - Many
- Many - to - One
- Many - to - Many

Variants of RNNs :-

① LSTM  $\Rightarrow$

A special type of RNN that solves problems like forgetting info. too quickly. It uses gates to better remember things for a longer time.

② GRU  $\Rightarrow$  (Gated Recurrent Unit)

like LSTM but Simpler & Faster. It still helps the network remember imp info, but with fewer parts.

③ Bidirectional RNN  $\Rightarrow$

This type of

RNN looks at the data both forward and backward, so it understands the context better.

### Applications :-

- NLP
- Speech Recognition
- Time-Series Prediction.

### Challenges :-

- Vanishing & Exploding Gradient problems
- Computational Complexity

## # Bidirectional RNNs :-

A Bidirectional Recurrent Neural Network extends traditional RNN by allowing sequential data to be processed in two directions: forward & backward.

This dual processing capability enables the BRNN to capture info from both past & future contexts, providing a richer understanding of the i/p sequence.

How BRNNs Work :-

① Architecture  $\Rightarrow$

A BRNN Consists of two Separate RNN layers :

$\rightarrow$  forward layer :

processes i/p Sequence from the Start to the end.

$\rightarrow$  Backward layer :

processes the same i/p Sequence from the end to the Start.

② Combining O/P  $\Rightarrow$

The O/Ps from these two layers are typically combined to form the final O/P, which leverages the Contextual info from both directions.

③ Applications  $\Rightarrow$

$\rightarrow$  NLP

$\rightarrow$  Machine Translation

$\rightarrow$  Speech Recognition

Advantages :-

- Enhanced Context understanding.
- Better performance.
- Flexibility.

Limitations :-

- Increased Computational Cost
- Real - Time Processing Challenges
- Greater Complexity in Training

## # Deep Recurrent Network :-

Deep Recurrent Networks use memory and step-by-step processing to understand ordered data like time series, text or Audio. They stack multiple RNN layers, which helps the model learn more detailed patterns & understand info at different levels & time steps.

## Key Components of Deep RNN Architecture ↴

- Recurrent Neurons ↴

These are the core of RNNs and can remember past i/p's using a hidden state.

- Hidden State ↴

Updated at each step  
to keep track of previous info & maintain context

- Layer Stacking ↴

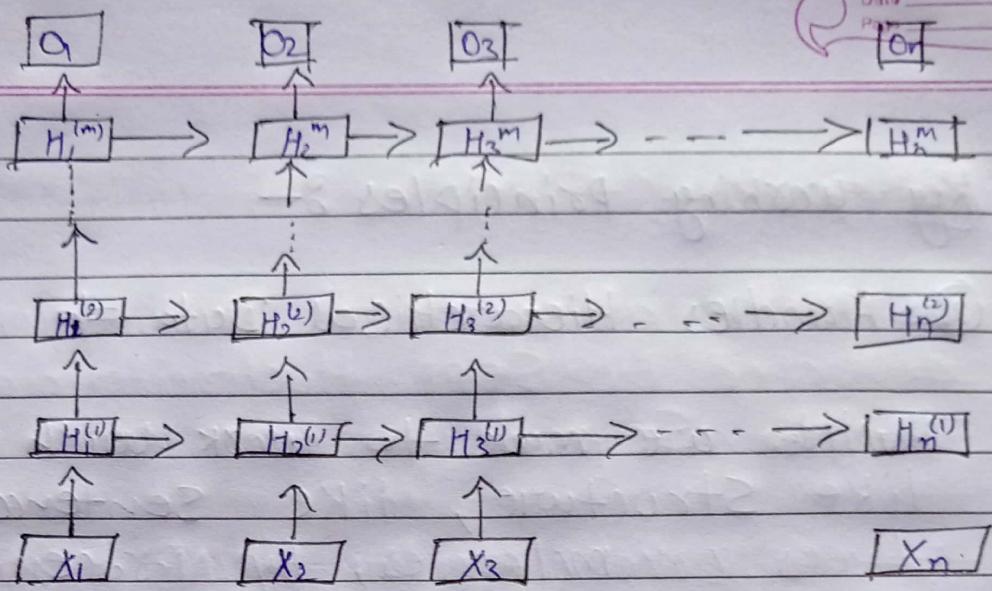
Multiple RNN layers can be stacked  
to learn more complex patterns

- Connection ↴

O/Ps from one step influence future steps, helping the model understand sequence flow.

- O/P layer ↴

Uses the final hidden state to make predictions, often with LSTM or GRU for better memory handling



Advantages :-

- Hierarchical learning
- Contextual Understanding
- Versatility

## # Recursive Neural Network :-

Recursive Neural Networks (RNNs) are models designed to work with tree structured data. Instead of processing data in a straight line, they combine smaller parts to build bigger parts. This helps them understand complex, nested relationships in the data.

## Key Working Principles :-

### ① Handles Hierarchical Data =>

RvNNs are made to work with tree like structure, like sentence parse trees or molecules, by processing smaller parts and building up to larger ones.

### ② Shared Parameters =>

They use the same set of rules at all levels, which helps the model work well across different parts of the structure.

### ③ Tree Traversal =>

RvNNs can move through the tree from the bottom up or top down, updating each part based on its connected parts.

### ④ Composition Function =>

This ~~function~~ Combines info. from smaller parts to form a bigger part helping the model understand complex relationships.

Applications :-

→ NLP

→ Image Analysis.

Difference bet<sup>n</sup> Recursive & Recurrent  
Neural Network =

① Type of Data =

Traditional RNN work with sequences  
like Time Series or text, using memory  
to track past steps.

Recursive N.N., on the other hand, handle  
tree-like or hierarchical data, such  
as sentence structures or molecules.

② How They Remember =

Traditional RNNs use loops to remember  
earlier info. RNN (Recursive) don't rely  
on time steps - instead, they build  
relationships by combining parts of  
a structure like putting pieces of  
a puzzle together.

## # LSTM :-

LSTM is an advanced architecture developed by Sepp Hochreiter and Jürgen Schmidhuber in 1997.

It improves regular RNNs by helping them remember info over long sequences, which is useful for tasks like language modeling, speech recognition and Time Series prediction. Unlike standard RNNs, LSTMs use memory cells and gates to decide what info to keep or forget.

## Architecture :-

The typical LSTM Unit includes a cell state and three types of gates:

(i) Input Gate  $\Rightarrow$

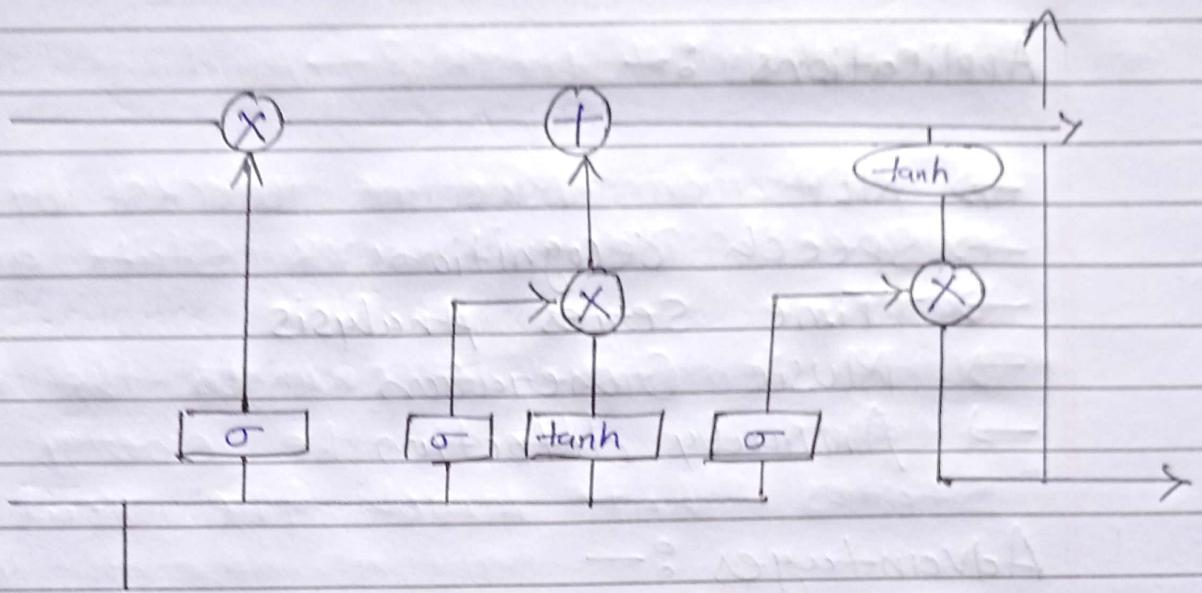
Controls what new info is added to the cell state.

(ii) Forget Gate  $\Rightarrow$

Decides which info to discard from the cell state.

(iii) Output Gate  $\Rightarrow$

Determines what info to o/p from the cell state.



Structure of an LSTM Network

### Working Mechanism :-

- Cell State  $\Rightarrow$  Carries imp info through the sequence.
- Gates  $\Rightarrow$  Use Special functions to decide what to keep or forget, with values bet<sup>n</sup> 0 & 1.

Together, these parts help LSTMs remember long term info, making them useful for many tasks.

Applications :-

- NLP
- Speech Recognition
- Time Series Analysis
- Music Generation
- Anomaly detection

Advantages :-

- Handle long term dependencies
- Flexible

Disadvantages :-

- Complexity
- Resource Intensive.

## Unit : 5

## # Boltzmann Machine :-

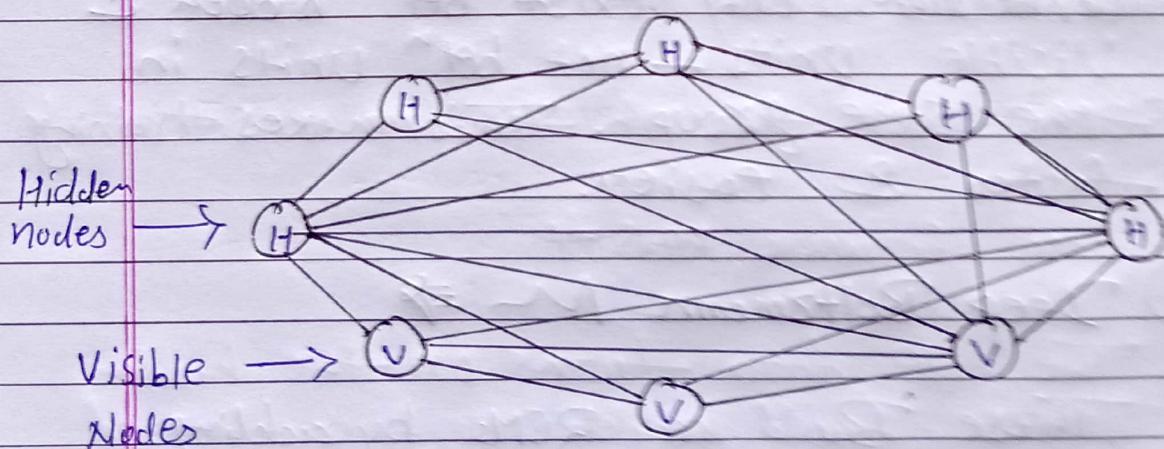
A Boltzmann Machine is a network of Symmetrically Connected Units that makes stochastic decisions about their states (on/off). Named after the physicist Ludwig Boltzmann, it operates based on principles from Statistical mechanics. The network includes two types of nodes:

(i) Visible Units  $\Rightarrow$

These units are used to i/p & o/p data.

(ii) Hidden units  $\Rightarrow$

These units learn patterns from i/p data, capturing the underlying structure in the data.



How does it work :-

A Boltzmann Machine works by giving each possible network state an energy value.

During learning, it changes the connection weights to lower the energy for states that match the real data.

It updates the units based on probabilities linked to these energy levels, often using a method called Gibbs Sampling to estimate these probabilities.

Types :-

① Restricted Boltzmann Machine (RBMs) :-

These are simpler versions where connections only exist between hidden & visible units - not between units in the same group. This makes training faster & easier.

② Deep Boltzmann M ↗

These build on RBMs by adding multiple hidden layers, helping the model learn more complex patterns.

### (3) Full Boltzmann $\curvearrowleft$ :-

- These allow all units to connect with each other, but they are much harder & slower to train.

#### Applications :-

- Feature Extraction
- Dimensionality Reduction
- Recommendation Systems
- Image Recognition.

#### Challenges :-

- Computationally intensive
- Training Complexity.

## # Restricted Boltzmann Machine :-

An RBM is a simple type of neural network with two layers:

- Visible layer  $\Rightarrow$  Takes in the i/p data.
- Hidden layer  $\Rightarrow$  Learns features or patterns from the data.

The word "restricted" means that units in the same layer don't talk to each other - only connections bet<sup>n</sup> visible and hidden units are allowed. This setup makes learning faster and easier. RBMs don't need labeled data, so they're great for tasks like reducing data size or finding useful patterns.

### How RBMs Work :-

RBM's learn through two main steps  $\Rightarrow$

#### ① Feed Forward pass $\Rightarrow$

The i/p data turns on the hidden units. Each hidden unit is turned on or off based on a probability calculated from the i/p.

## ii) Feed Backward Pass :-

The hidden layer tries to rebuild the original i/p. The RBM adjusts its connections to make this reconstruction better.

## Training RBMs :-

Training means improving the network by adjusting the weights bet<sup>n</sup> units. RBMs use Contrastive Divergence (CD) to do this efficiently. A more advanced version called persistent ~~cost~~ Contrastive Divergence (PCD) can improve learning by remembering previous states.

## Applications :-

- Collaborative filtering
- Image & Video processing
- NLP
- Bioinformatics
- Financial modeling.

## Advantages =>

- Can find useful features in data without labels.
- Can create new data that looks like the i/p
- Easier to train than more complex models like full BM.

## Disadvantages =>

- Training can be tricky & slow
- Needs careful tuning to work well
- Hard to understand what ~~model~~ exactly the model is learning.