

## Intro to ML

Page No.	Date
70/144	10/10/2024

> Machine Learning - A Computer Science Technology where a machine itself has a complex range OR knowledge that allows it to take certain data input and use complex statistical analysis to create output value. "Machine Learning" Focuses on building a system that can access data and use it to learn for themselves.

> Why ML is important? -

ML and data mining , a component of ml are crucial tools to process data and gain insights from massive amount of data.

→ Increase In Data Generation -

→ Improved Decision Making -

→ Uncover patterns & trends in data -

→ Solve Complex problems -

> Steps Involved.

- 1 Define the objective of problem statement.
- 2 Data Gathering
- 3 Data preparation
- 4 Exploratory Data Analysis
- 5 Building Machine Learning Model
- 6 Model Evaluation & optimization.

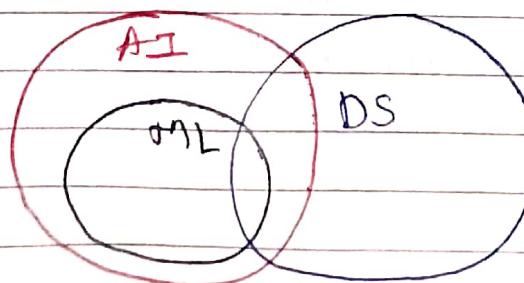
## Comparison of ML with Traditional Programming.

ML

TP

- 1) Not a manual process Algorithm Manual process, programmer automatically predicts Creates program without anyone has to manually formulate
- 2) Data → Computation → Program  
Result → Data → Computation → Result.  
Result program
- 3) As a Subset of AI ml is We write down the exact motivated by human learning steps required to solve problem behaviour.
- 4) ml algorithm takes an input Traditional algo takes input and output and gives some and some logic in the form logic which works on new of code and gives output.  
Input

### ML Data Science and AI



- ML is subset of AI wherein Computer System learn from the environment.
- Data Science is processing, analysis and extraction of relevant assumptions from data.
- ML uses statistical models. AI uses logic and decision tree. Data Science deals with structured data.

## Machine Learning

## AI

## Data Science

Focuses on learning from data and get experience to improve overtime. Giving machines cognitive & intellectual Abilities similar to human Focuses on extracting information needs from data for decision making.

Uses Statistical model

uses logic & decision trees

Deals with Structured unstructured data.

A form of Analytics Development of Computer program which learns from data and finds pattern. Process of fusing advanced analytics to simulate human intelligent extract relevant info.

objective is to maximize accuracy

objective is to maximize the chance of success

objective is to extract actionable insights from data

ML can be done through Collection of intelligence Supervised, Unsupervised or reinforcement learning Concepts, including Predicates, logic, perception, data wrangling

recommendation system Such as Spotify

chatbots and voice assistant.

Fraud detection and Healthcare

### > Types of Learning

#### Machine Learning

##### Supervised

- classification
- Regression

##### Unsupervised

- clustering
- analysis

##### Reinforcement

## > Supervised Learning -

- > It's a machine learning paradigm in which algorithm learns from labeled training data to make predictions or decisions. The goal is to learn mapping from input data to the correct output labels so that the algorithm can make accurate predictions on new unseen data.
- > The algorithm uses the iterative optimization of an Objective Function to predict the output that will be associated with new inputs.
- > Supervised learning algorithm uses classification and regression techniques to develop predictive models. Regression technique predicts continuous responses ex. Changes in temperature.
- > Supervised learning are trained using labelled dataset where the model learns about each type of data.
- > Give examples - Adv & disadv

## > UnSupervised Learning -

- > It refers to learning from unlabeled data. It is based on more similarities and differences than on anything else. In this learning similar items are clustered together in a particular class.
- > Many Unsupervised algorithms create similar hierarchical arrangements based on similarity based mapping.
- > The task of hierarchical clustering is to arrange a set of objects into hierarchy such that similar objects are grouped together.
- > A learner is fed with a set of scattered points, and it generates two clusters with representative centroid after learning. Clusters show that points with similar properties and closeness grouped together.

## Types -

1) Clustering

2) Association - Kmeans, KNN, Neural Network  
Apriori Algo

### Supervised

> Trained on labeled data

### Unsupervised

trained on unlabeled data

> Model takes feedback to check if it's predicting correctly or not

> Predicts the output

Finds hidden pattern in data

> Data is provided to the model along with output

> Goal is to train model that can predict when new data is inputted

Goal is to find pattern and useful insight from unknown dataset

> Needs supervision to train

Does not need any supervision

> Categorized into classification & regression

Categorize into clustering & association problem.

> Model produces accurate output

Less accurate compared to supervised learning

> Linear regression, Logistic, SVM

clustering, KNN, Apriori Algorithm.

Distance models - Distance metrics aka Similarity or dissimilarity metrics are mathematical functions used to quantify the similarity or dissimilarity between two points in a feature space. These metrics provide a way to measure the distance between data points which is fundamental concept in various ML algo.

→ Euclidean Distance - It's one of the well-known distance metrics and is used in various applications. Calculates the straight line distance between 2 points in Cartesian co-ordinate sys.

$$\text{Formula: } d(P, q) = \sqrt{\sum_{i=1}^n (P_i - q_i)^2}$$

Euclidean distance is used in KNN and hierarchical clustering.

→ Manhattan Distance - Measures the distance between two points by summing the absolute differences between their coordinates.

$$\text{Formula: } d(P, q) = \sum_{i=1}^n |P_i - q_i|$$

Manhattan used when movement between points can only occur along grid lines.

→ Cosine Similarity - Cosine Similarity measures the cosine of the angle between two vectors, indicating their similarity regardless of vector lengths. It is often used for comparing text documents, where each dimension corresponds to TF-IDF value.

$$\text{Formula: } \text{Cosine Similarity}(P, q) = \frac{P \cdot q}{|P| |q|}$$

→ Hamming Distance - It is used to measure the difference between 2 strings of equal length often binary strings. It calculates the number of positions at which the corresponding elements are different. It is used in error detection and correction codes in digital communication.

If 10001001 and 10110001 are 2 codewords, then the corresponding bits differ in these two codewords is 3 bits. The number of bit positions in which two codewords differ is called hamming distance.

$$\text{Equation } 2^H \geq M + H + 1 \quad M = \text{No. of bits in msg} \\ H = \text{Hamming bits}$$

**A** Reinforcement Learning - RL is a subset of ml that deals with how agents should take actions in an environment to maximize a cumulative reward. RL involves learning through interaction with an environment without explicit supervision.

Key Components of RL -

- 1) Agent - Agent interacts with its environment. The agent takes actions, observes the environment, and learns to make better decisions over time.
- 2) Environment - External System with which the agent interacts. Representing the current situation and it can change based on agent's actions.
- 3) States - Representation of the current configuration or situation of the environment at a given time. State gives essential info for the agent to make decisions.
- 4) Actions - Set of all possible moves or decisions that an agent can make in the environment.

Reinforcement learning widely used in various applications robotics, game playing, autonomous vehicles, recommendation systems, NLP and more.

A) Decision Tree - It's a Supervised ml algo that is used for both Classification and regression tasks. It is a versatile and interpretable model that is particularly useful when you need to make decision based on a set of conditions or features.

key characteristics & concepts -

1) Tree Structure - Tree like Structure composed of nodes. The tree begins with a root node, which represents the entire dataset and branches out into child nodes.

2) Nodes - Root Node - Representing entire dataset. Internal Nodes - Repre. Feature conditions and lead to child nodes.

Leaf Nodes - Terminal nodes that make predictions.

3) Edges - Connection between nodes represent the outcomes of feature conditions.

4) Feature Selection - At each internal node, the algo selects a feature and a threshold value that best splits the data into more homogeneous subsets.

5) Decision Rules - Each path from root node to leaf node represents decision rule. For classification tasks, these rules lead to class labels. They lead to predicted values.

6) Impurity Measures -

→ Gini impurity - measures probability of misclassifying a randomly chosen element.

→ Entropy -

→ Mean Squared Error (MSE) - measures the variance of the target variable.

M	T	W	T	F	S	S
Page No.:						
Date:						YOUVA

## Parametric method

## Non Parametric method.

- 1) USE a fixed number of parameters to build the model.
- 2) Parametric Analysis is for testing for testing medians group means.
- 3) applicable only for variables
- 4) Always considers Strong assumptions about data.
- 5) Require lesser data than Non parametric
- 6) Handles intervals data or relation data
- 7) Follow normal distribution
- 8) Output generated can be easily affected by outliers.
- 9) Parametric methods have more statistical power than Non parametric.
- 10) Ex: Logistic Regression, Naive Bayes
- 1) USES flexible number of parameters to build model.
- 2) Generally considers fewer assumptions and data.
- 3) Requires much more data than parametric method
- 4) handles original data.
- 5) There's no assumed distribution
- 6) Output cannot be easily affected by outliers.
- 7) have less statistical power than parametric methods.
- 8) ex. KNN, Decision Tree, etc.

## — Feature Extraction —

- \* **Feature** - Features are individual independent variables that act like input in the system. Feature is an attribute of a data set and used in ml process.
- > Feature extrac<sup>n</sup> is the process of Selecting, transforming or creating relevant features from raw data. It aims to reduce the dimensionality of data.
- > Feature Engineering - Process of creating new features or modifying existing one to improve performance.
- > Feature Selection is a process that chooses a subset of features from the original feature space is optimally reduced according to a certain criterion.

### \* Subset Selection -

It is a technique in feature selec<sup>n</sup> where you choose a subset of most relevant features for use in a ml model. The goal is to reduce dimensionality of dataset by retaining the most relevant features

> There are 2 types of subset selec<sup>n</sup> -

> Forward selec<sup>n</sup> - It starts with an empty set of features and iteratively adds one feature at a time, selecting the one that improves model performance

> Backward selec<sup>n</sup> - It starts with all variables and remove them one by one, at each step removing the one that decreases the error the most, until any further removal increases the error significantly.

> Sequential Forward Search - SFS is the simplest greedy search algo. It starts from empty set, sequentially add the feature  $x^+$ . SFS performs best when optimal subset is small.

> The main disadvantage of SFS is that it is unable to remove features that become obsolete after adding new feature.

4) Sequential Backward Selection - Starts with all variables and remove them one by one, at each step removing the one that decreases the error the most, until any further removal increases the error significantly.

> SBS works best when the optimal feature subset is large, since SBS spends most of its time visiting large subsets.

## # Preprocessing of Data -

Data preprocessing is a data mining technique that involves transforming raw data into understandable format.

> Data which capture from various source is not pure. It contains some noise known as dirty data.

> Incomplete, noisy and inconsistent data are commonplace properties of large real world database.

> Steps during Pre-processing -

1) Data Cleaning - It is cleaned through processes such as filling in missing values, smoothing the noisy data or resolving inconsistencies in the data.

2) Data Integration - Data with different representation are put together and conflicts within the data are resolved.

3) Data Transformation - normalization, scaling.

4) Data Reduction - To present reduced representation of dataset.

> Normalisation & Scaling -

> Normalization is a data preparation technique frequently used in ml. The process of transforming the columns in a dataset to the same scale is referred to as normalization.

- > Normalization makes features more consistent with each other, which allows the model outputs more accurately.
- > Normalization refers to rescaling real valued numerical attributes into a 0 to 1 range. KNN, SVM uses normalization.
- > Most widely used normalization

### 1) Min-max Scaling -

$$x_{\text{normalize}} = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

### 2) Standardization (Z-score) -

$$x_{\text{std}} = \frac{x - \mu}{\sigma}$$

$$\sigma = \sqrt{\text{Var}(x)}$$

$$\sigma^2 = \frac{(x_i - \bar{x})^2}{N}$$

## # Dimensionality Reduction -

Two Components - Feature Selection & Extraction.

- Missing value, low variance, Decision Tree, Random Forest

## # Principle Component Analysis -

> It's Dimensionality reduction technique commonly used in machine learning. Primary goal is to reduce the number of features in a dataset while preserving as much of original info as possible.

> PCA is UnSupervised learning algo technique used to examine interrelations among a set of variables.

PCA reduce the dimensionality of a dataset by finding new variables.

> Variance Capture - PCA identifies a set of orthogonal axes called Principal Components, that captures the maximum variance in the data.

> Orthogonal Transformation - PCA performs an orthogonal linear transformation of original features ensuring

that new components are uncorrelated with each other.  
PCA often uses data visualization

> STEPS Involved -

Step 1 - Data Pre-Processing

Start with dataset that contain m data points and m features. Standardize or normalize the data to ensure that features with different scales do not dominate analysis.

Step 2 - Calculate Covariance Matrix.

measures the strength of joint variability between 2 or more variables, indicating how much they change in relation to each other.

$$\text{Cov}(\mathbf{x}_1, \mathbf{x}_2) = \frac{1}{m-1} \sum_{i=1}^m (\mathbf{x}_{1i} - \bar{\mathbf{x}}_1)(\mathbf{x}_{2i} - \bar{\mathbf{x}}_2)$$

Step 3 - Calculate eigen vectors and Eigen values.

> performs eigen value decomposition on covariance matrix  $\Sigma$  to calculate the eigen vectors & values.

> eigenvectors represent the directions in which the data varies the most, and eigenvalue indicates how much variance is explained along with principle component.

> Given square matrix A find eigen value  $\lambda$  and eigen vector  $X$  such that  $Ax = \lambda x$

Rearrange eqn  $Ax - \lambda x = 0$ .

>  $x(A - \lambda I) = 0$  I is identity matrix.

This eqn will only hold if  $(A - \lambda I)$  is singular matrix  
Calculate determinant of  $(A - \lambda I)$  set it to equal to zero.

$$\det(A - \lambda I) = 0$$

> Solve  $\lambda$  to obtain eigenvalue of matrix A

For each eigenvalue  $\lambda_i$ , you can find corresponding  $x_i$  by solving eqn  $Ax_i = \lambda_i x_i$ .

### Step 4 - Project onto Principal Component -

- > Project the original data onto Selected K principle Components to obtain new reduced dimension dataset.
- > Project is performed by multiplying the standardized data by the K- Selected eigen vectors.

### \* Local Binary Pattern -

- > It is a Simple yet very efficient texture operator which labels the pixels of an image by thresholding the neighbourhood of each pixel and consider the result as binary number.
- > LBP operator is expressed by a sequence of binary values obtained by Comparing central pixel with its neighbours in Circular manner.
- > Compare the intensity value of central pixel with the intensity value of neighbor. for each neighbor if intensity value is greater than or equal to that of central pixel assign a binary '1' otherwise '0'

5	4	3	Threshold	1	1	1	11101001
4	3	1	→	1	X	0	Decimal 233
2	0	3		0	0	1	

- > The Decimal values obtained from LBP are used as texture descriptors for each Pixel and they can be further processed for various CV tasks.

- > Texture classifiers  
Face recogn  
Object Detect  
Image retrieval

# Feature Selection Techniques -

## Feature Selection

Filter method

wrapper

Embedded

### Filter method -

Set all Features → Selecting the Best Subset → Learning Algo → Performance.

→ Filter method are generally used as preprocessing Step.

The Selection of Features is independent of any ml algo. Features are Selected on the basis of their Scores in various Statistical tests

Some Common techniques.

→ Information Gain , Chi Square Test , Anova

→ Correlation based Feature Selection is a simpler filter algo that ranks features Subsets according to a Correlation based on heuristic evaluation Function.

### Wrapper method -

Set of Features.

In Wrapper method,

Selection of Features is done

by considering it as

Search problem, in

which different combination

are made, evaluated and

compared with other combination

It trains algorithm by using subset of features iteratively.

methods of wrapper technique are.

- forward Selection
- Backward Elimination

## ❖ Embedded method -

It is combined with the advantages of both Filter and wrapper methods by considering the interaction of features along with less computational cost. These are fast processing methods similar to the filter but more accurate than the filter method.

Set of Features



→ Generate Subset



Algo + Performance ←

- In embedded method feature selection algo is blended as part of learning algo, thus having its own built in feature selection method
- Some techniques are -
  - Regularization - This method adds penalty to different parameters of ml model to avoid overfitting of model.
  - Tree based methods - These methods such as Random forest, Gradient Boosting provides usefulness feature importance as a way to select feature as well.

## ML Supervised Learning.

### \* Underfitting -

- > Underfitting occurs when model is too simple to capture underlying patterns in training data.
- > Underfitting destroys the accuracy of machine learning model, model or Algorithm does not fit the data well enough.
- > It usually happens when model is trained on less data.
- > Techniques to Reduce underfitting.
  - 1) Increase model Complexity
  - 2) Increasing no. of features perform Feature enginnering
  - 3) Remove noise from the data.
  - 4) Increase no. of epochs.

### \* Overfitting -

- > A Statistical model is said to be overfitted, when the model is trained on too much data, it starts learning from the noise & inaccurate data in dataset.
- > Then the model does not categorize the data correctly because of too many details & noise. The cause of overfitting are non-parametric & non-linear methods because these models have freedom in building the model on dataset.
- > Techniques -
  - 1) Cross validation
  - 2) Feature Selection
  - 3) Regularizer - L1 or L2
  - 4) Ensemble methods
  - 5) Simplify model.

## Bias variance Tradeoff -

> Bias - Refers to the error introduced by approximating real world problem too simplistically. It is the inability of model because of that there is some difference or error occurring between model's predicted value and actual value.

> ways to reduce -

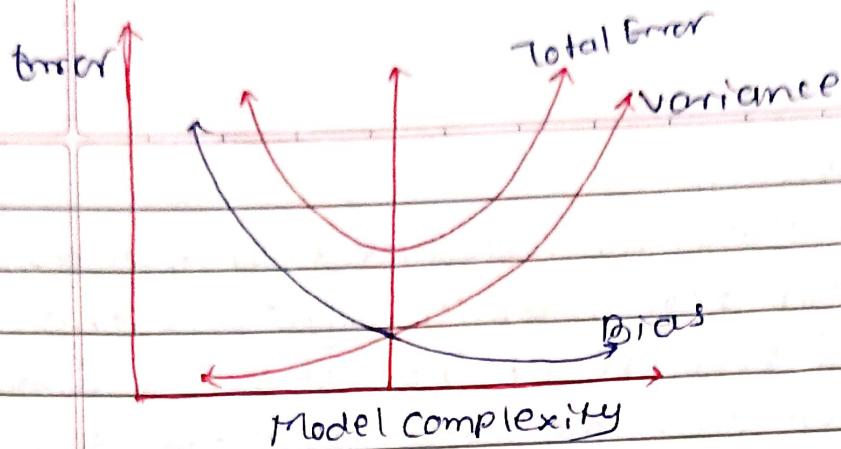
- 1) Use complex model
- 2) Increase number of features
- 3) Reduce Regularization
- 4) Increase the size of training data.

> Variance - Variance is measure of spread in data from its mean position. Variance is the amount by which the performance of a predictive model changes when it is trained on different subsets of training data.

> ways to reduce

- 1) Cross validation
- 2) Feature selection
- 3) Regularization -  $L_1$  &  $L_2$
- 4) Ensemble methods.
- 5) Simplifying model.

> If the algorithm is too simple then it may be on high bias and low variance conditions. If algorithm fit too complex bias and others it may be on high variance & low bias. There's something between both of these conditions known as Trade-off or Bias variance trade-off. The goal is to find the right level of complexity in a model minimize both bias & variance.



**Linear Regression** - Linear regression is a type of supervised learning algo which computes the linear relationship b/w dependent & more independent variables. The goal of algorithm is to find best linear equation that can predict the value of dependent value based on independent values.

$$\text{eqn } y = \beta_0 + \beta_1 x + \epsilon$$

$y$  - dependent variable.  $\beta_0$  -  $y$ , intercept.

$x$  - independent  $\beta_1$  - Slope

$\epsilon$  - error term.

$$y = ax + b$$

$$a = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2} \quad b = \frac{1}{n} (\sum y - a * \sum x)$$

**Regularization** - It is a technique used in machine learning to prevent overfitting and improve the generalization of model. It introduces a penalty term to the function that the model is trying to minimize during training. The magnitude of the coefficients is reduced by using different types of regularization techniques like:

- 1) Lasso Regression /  $L_1$  regularization
- 2) Ridge Regression /  $L_2$  regularization

> **Ridge Regression** - It's a linear regression technique that adds penalty term of sum of squared values of model's coefficients. Goal is to prevent overfitting.

$$\text{eqn} \quad \text{Ridge } R = \text{Loss} + \lambda \sum_{i=1}^n (w_i)^2$$

- 1) Loss Term represents the standard linear regression, loss often MSE
- 2)  $\lambda$  is regularization Strength.
- 3) Ridge regression shrinks the coefficient towards zero but does not force them to exact zero
- 4) Larger value of  $\lambda$  results more shrinkage
- 5) It discourages model from relying too much on any single variable when multiple variables are highly correlated.

## > LASSO Regression (Least Absolute Shrinkage & Selection Operator)

It is a form of linear regression that incorporates use of regularization to avoid overfitting. Term lasso reflects its ability to shrink coefficient towards zero and also achieves exact zero value.

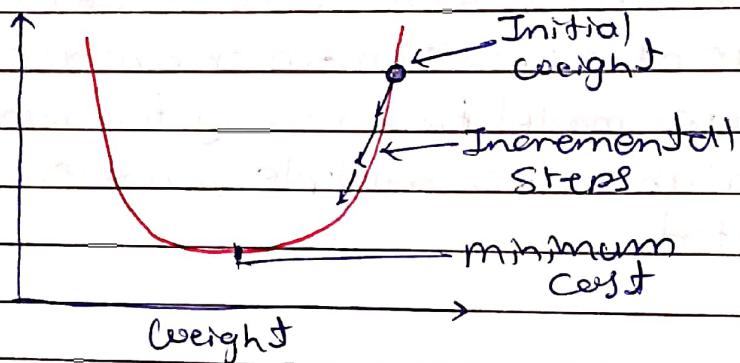
$$\text{eqn} \quad \text{Lasso } R = \text{Loss} + \lambda \sum_{i=1}^n |w_i|$$

- 1) Adds the sum of squared absolute values of model's coefficients as penalty term
- 2) Encourage Sparsity in model i.e. it tends to drive some coefficient to exactly zero
- 3) Useful when there's belief that many features are irrelevant or redundant.

## Gradient Descent Algo -

Gradient Descent is an optimization algorithm used to minimize a function iteratively by adjusting its parameters. It helps in finding local minimum of a function.

- > If we move towards negative gradient or away from the gradient of the Function at current point, it will give local minimum.
- > Move towards positive Gradient we will get local max



- > The main goal is to minimize cost function iteratively. To achieve this goal it performs 2 consecutive steps.
  - 1) Calculates 1<sup>st</sup> order derivative of the function to compute the gradient or slope of that function
  - 2) Move away from the direction of gradient (which means slope increased from current point by alpha time alpha is the learning rate)
  - 3) Gradient Descent requires function to follow certain conditions like
    - > Function should be differentiable &
    - > Function should be convex

### Steps involved -

- 1) Initialization
- 2) Compute Gradient
- 3) Update parameters
- 4) Repeat
- 5) Convergence

## > TYPES OF GD -

(BGD) 1) **Batch GD** - It is used to find the error for each point in training set and update the model after evaluating all training examples. This is known as training epoch.

(SGD) 2) **Stochastic GD** - SGD is a type of GD that runs one training example per iteration. It processes training epoch for each example within dataset and updates each training example's parameters one at a time.

3) **Mini Batch GD** - Combination of both BGD & SGD divides the training dataset into small batch sizes then perform the updates on those batches separately.

**Evaluation Metrics** - Evaluation metrics are measures used to assess the performance of a machine learning model. These metrics provide insights into how well the model is performing on a given task, such as classification, regression, clustering.

### > classification metrics

$$1) \text{Accuracy} = \frac{\text{No. of correct Prediction}}{\text{Total No. of Prediction}}$$

$$2) \text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

$$3) \text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$4) \text{F1 Score} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

> Confusion matrix - Performance measurement for machine learning classification problems consists of True positive, True negative, False positive, False negative

> AUC ROC - Area under the curve

ROC (Receiver operating characteristic) curve is the plot between False positive rate and True positive rate.

> Gini Coeff. - Used for classification problems. Can be derived straight away from AUC ROC number. It is ratio between the ROC curve & diagonal line & the area above triangle.

$$\text{Gini} = 2 * \text{AUC} - 1$$

> Mean Absolute Error (MAE) -

Diff. between Actual & predicted values

$$\text{MAE} = \frac{1}{N} \sum_j |y_i - \hat{y}_i|$$

Actual      Predicted.

Aim is to get minimum MAE

> Root Mean Squared Error (RMSE) -

$$\text{RMSE} = \sqrt{\text{MSE}}$$
$$= \sqrt{\frac{1}{n} \sum_i (y_i - \hat{y}_i)^2}$$

Actual      Predicted.

Lower RMSE expressed indicates better model performance  
AMSE is sensitive to outliers

> R Squared -  $R^2$  is independent of context.

R Squared is also known as Coeff. of Determination or Goodness of fit.

$$R^2 \text{ Squared} = 1 - \frac{SS_r}{SS_m} = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (\bar{y}_i - \hat{y}_i)^2}$$

$SS_r$  - Squared sum error of regression line.

$SS_m$  - Squared sum error of mean line.

R-Squared measures the strength of relationship between your model and the dependent variable on a convenient of ~~0-100% scale~~ 0-1 scale

> 1 when model perfectly fits the data.

> 0 when model does not predict any variability in the model and it does not learn any relationship between the dependent & independent variable.

## Supervised Learning: Classification

### KNN Algorithm (K - Nearest Neighbor)

The KNN algorithm is a simple, intuitive & versatile supervised machine learning algo for both classification & regression tasks. KNN is non-parametric and instance based learning algo means it doesn't make assumption about underlying data distribution & memorizes training dataset.

> The algorithm stores entire training dataset in memory.

> To make a prediction for a new data point, the algorithm finds the k-nearest neighbors in the training dataset based on Similarity metric.

Q Consider Following data to predict student pass or fail using Knn for value Phy = 6mark & chem = 8marks.

Phy	Chem	Result
4	3	Fa
6	7	P
7	8	P
5	5	F
8	8	P

— Distance =  $\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$   
 $(x_2, y_2) = (8, 6)$

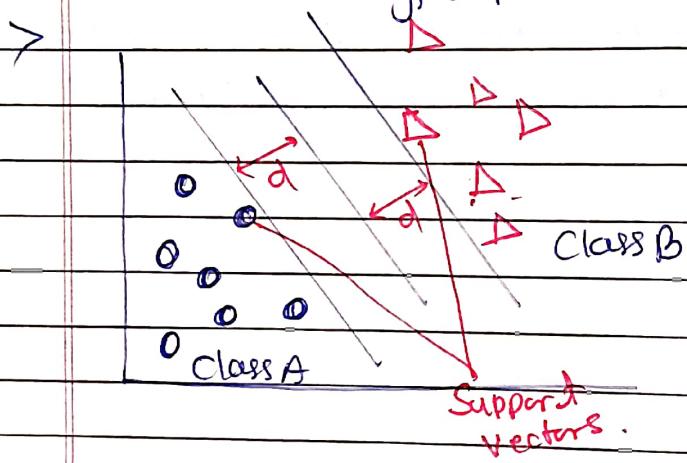
For (4, 3) dist =  $\sqrt{(8-4)^2 + (6-3)^2} = 5$

Phy	Chem	Res	dist.
4	3	F	5
6	7	P	1
7	8	P	2.23
5	5	F	3.16
8	8	P	2.

3 Smallest dist are 1, 2.23 & 2. The result of these distances is Pass so res of result of new data point is pass

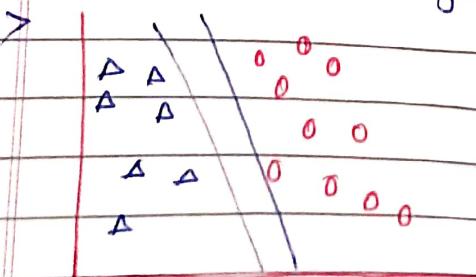
## \* SVM (Support vector machine) -

- > SVM are set of Supervised learning algorithm which learns from the dataset
- > SVM used for classification & regression tasks It performs both linear & non-linear classification
- > SVM finds a hyperplane to separate the inputs into Sequence separate groups There can be many hyperplanes that successfully separates input vectors.
- > Points closest to the hyperplane are known as support vectors. These are the points lie closest to the decision boundary influence the position & orientation of hyperplane.



> **Linear SVM** → Linear SVM aims to find a linear hyperplane in the input space that best separates data points of different classes.

> The decision boundary is straight line that maximizes the margin between classes.

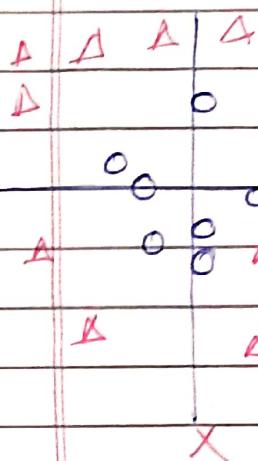


> Non Linear SVM - Non linear SVM used for non linearity. Sep separated data, which means if a dataset cannot be classified by using a straight line.

> Non linear uses kernel trick to implicitly map the input features into higher dimensional space where a linear decision boundary can be effective.

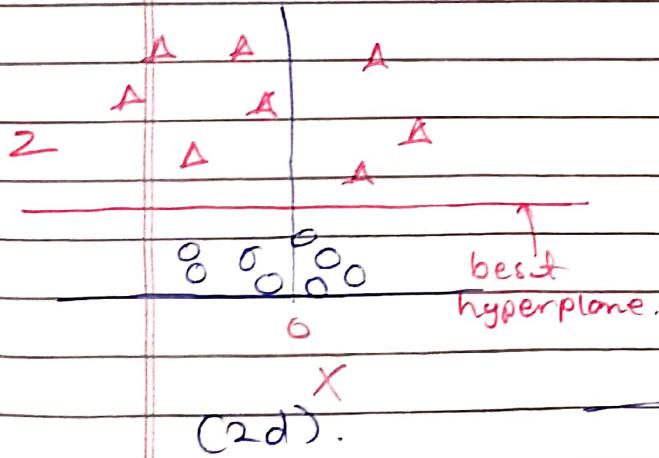
> It is suitable for dataset with complex & non-linear relationships between Features.

> There are 2 classes that can't

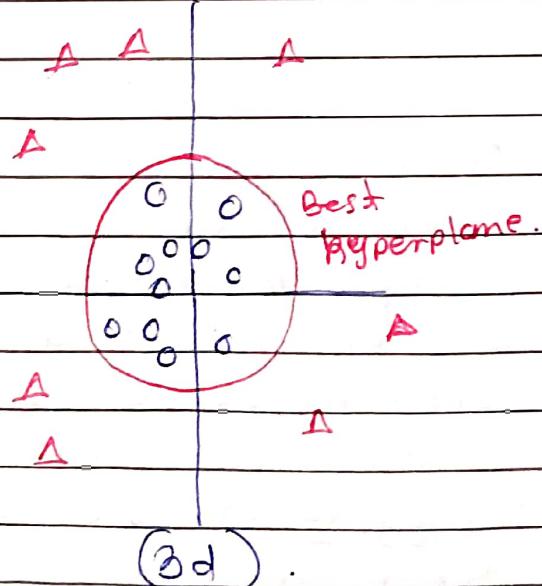


be separated by straight line. But a circular hyperplane can

separate them with the help of  $xc^2 + y^2$  where  $z = xc^2 + y^2$ . Now introducing 3rd dimension in graph.



(2d).



(3d).

> Separating hyperplane - used in binary classification which semaximally separates 2 classes of data points in feature space. Goal is to

Search optimally Separating hyperplane

$$w^T x + b = 0 \quad w \text{ is weight.}$$

$x$  is Feature vector.

$b$  is bias

$w^T x + b > 0$  corresponds to one class

$w^T x + b < 0$  corresponds to another class.

> Margin - Distance between separating hyperplanes and nearest datapoint from both classes. It represents safety buffer or width of corridor in which data points can reside without affecting classification result. Larger margin model is robust if it can tolerate noisy data. Smaller margin may lead to more accuracy but may overfit

$$\text{Margin} = \frac{1}{\|w\|} |w^T x_i + b|$$

\* Ensemble Learning - Ensemble learning is a machine learning technique that involves combining multiple models to improve overall performance and generalization. basic idea is to combine multiple weak models to create a strong & more accurate model.

> Bagging (Bootstrap Aggregating) -

> Bagging involves training multiple instances of the same learning algo on different subsets of training data.

> Each model in the ensemble is exposed to a different subset of the data introducing diversity among the models.

- > Random Forest is a popular ensemble method that uses bagging with decision trees.
- > Bagging can reduce overfitting, increase stability & improves accuracy & generalization.

### > Boosting -

- > Boosting aims to sequentially train a series of weak models, each focusing on correcting the errors made by predecessor.
- > Instances that are misclassified by one model are given higher weights in the training set of next model.
- > Predictions from all weak models are combined with different weights, and the final prediction is made by taking a weighted majority vote.
- > Ex - AdaBoost, Gradient Boosting, XGBoost.
- > It can achieve high accuracy. Effectively handles imbalanced datasets, adapts more complex relationships in data.

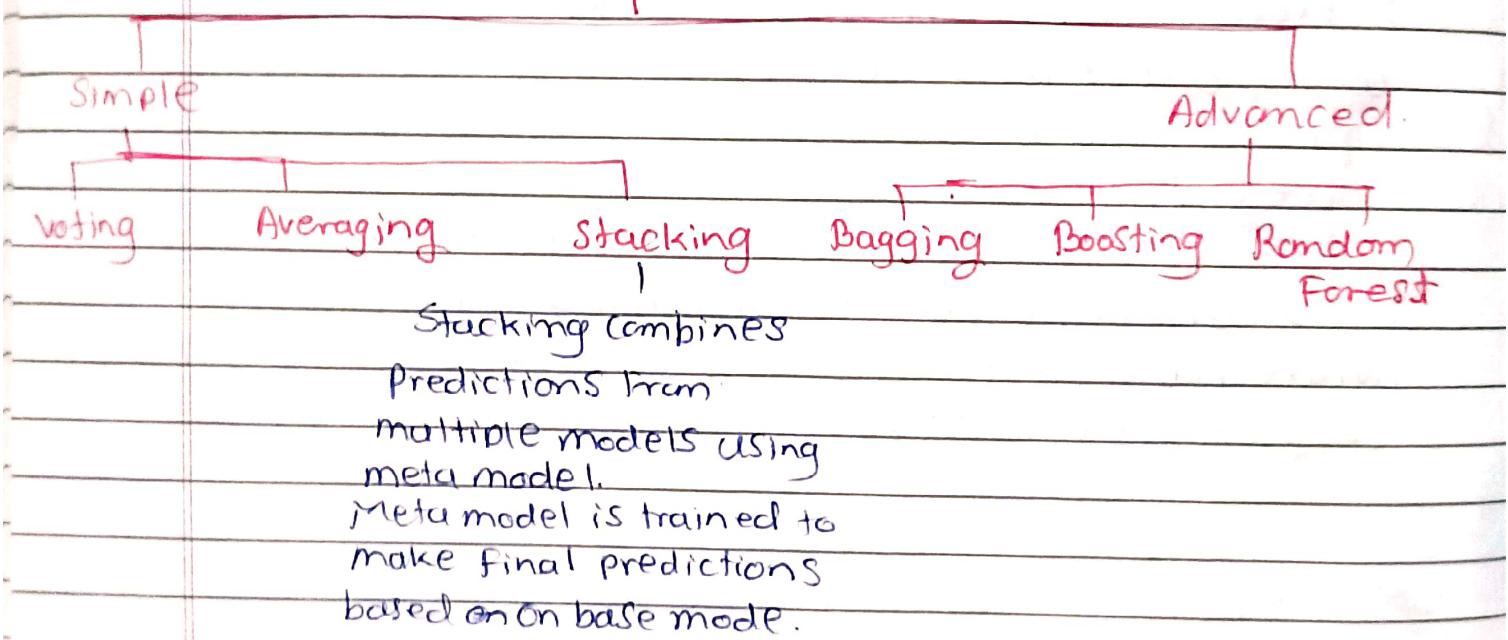
**# Random Forest -** Random Forest is an ensemble learning technique that belongs to family of bagging. It builds multiple decision trees during training and merges them together to obtain a more accurate and stable prediction.

Steps -

- 1) Random Forest starts creating multiple random samples with replacement from original training dataset. Each sample is used to train a separate decision tree.

- 2) At each node Random Forest randomly selects a subset of features. To include randomness diversity.
- 3) Decision tree are grown for each bootstrap sample using randomly selected features. Each tree expands until a certain condition is met.
- 4) For classification tasks the prediction from individual trees are combined through majority vote. For regression task prediction are averaged.

## \* Ensemble Learning methods -



\* Confusion Matrix - It is a table that is often used to evaluate the performance of a classification algo. It provides a summary of the classification results by showing the number of true positive, true negative, false positive, false negative predictions.

- > True Positive (TP) - No. of instances correctly predicted as +ve
  - > True-ve - No. of instances correctly predicted as -ve
  - > False +ve - No. of instances incorrectly predicted as +ve (Type I error)
  - > False-ve - No. of instances incorrectly predicted as -ve (Type II error).
  - > The importance of a confusion matrix lies in its ability to provide a detailed breakdown of a classifier's performance. Various performance metrics can be calculated including:
    - 1) Accuracy  $\Rightarrow \frac{TP + TN}{TP + FP + FN + TN}$
    - 2) Precision  $\Rightarrow \frac{TP}{TP + FP}$
    - 3) Recall  $\Rightarrow \frac{TP}{TP + FN}$
    - 4) F1 Score  $\Rightarrow \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$
- > These metrics help in assessing different aspects of a classifier's performance, such as ability to correctly classify positive instances, avoid false positives & overall accuracy.

- > F1 Score - It is a metric that combines precision and recall into a single value. It is particularly useful in situations where there is an uneven class distribution.
- > F1 Score ranges from 0 to 1, where 1 indicates perfect precision and recall, and 0 indicates the worst possible performance.
- > F1 Score is harmonic mean of precision & recall, giving more weight to lower values. This makes it's a suitable metric when there's an imbalance b/w the classes.

## UnSupervised Learning

\* **Clustering** - Clustering is a process of partitioning a set of data in a set of meaningful subclasses called as clusters.

\* **K-Means clustering** - It is a popular machine learning algo used for partitioning a dataset into  $k$  distinct, non overlapping Subgroups or clusters.

The goal of K-means is to group data points into clusters such that points within the same cluster are more similar to each other.

Steps

> Initialization - choose number of clusters  $k$  (centroids)

> Assignment - Assign each data point to the nearest cluster centroid. Usually done by using Euclidean distance.

> Update centroids - Recalculate the centroid of each cluster by taking mean of all data points assigned to that cluster.

4) Repeat step 2 & 3 until convergence, convergence.

Occurs when assignment of data points to clusters no longer changes significantly

> Algo takes the unlabeled dataset as input, divides the dataset into  $k$ -number of clusters and repeats the process until it doesn't find best clusters.

> Find value of  $K$

\* **Elbow method** - This method uses the concept of WCSS value (Within cluster sum of squares) which defines the total variations within cluster.

> It executes the Kmeans clustering on a given dataset for different  $K$  value from 1-10

> For each value of  $K$ , calculates the WCSS value.

> Plots a curve between calculated WCSS values and the number of clusters K.

> The sharp point of bend or a point of the plot looks like an arm, then the point is considered as the best value of K.

$$\text{Sum Squared Error} = \sum_{i=1}^n \sum_{j=1}^k \|x_i - c_j\|^2$$

n - no. of data pt.

k - no. of clusters

$x_i$  - is  $i^{th}$  data pt.

$c_j$  -  $j^{th}$  cluster.

**K-MEDOIDS** - It is a clustering algo similar to K-means but has key difference in how it defines cluster centres. While Kmeans uses the mean of the data points in a cluster as the cluster center,

K-medoids uses actual data point that is most centrally located in a cluster as the medoid.

The medoid is the data point with minimum average dissimilarity to all other data points in cluster.

→ Steps

- 1) Initialization - choose k initial data points in dataset.
- 2) Assignment - Assign each data point to the nearest medoid by using any common dist. metric (Euclidean dist. or Manhattan dist.).

3) Swap (while the cost of configuration decreases).

→ For each medoid 'm' and each data pt 'o'

- a) Consider the swap of m and o and compute cost change.

- b) If the cost change is current best, remember this m and o combination.

- 2) Perform the best Swap mbest & Obest if it decreases the cost function

\* Hierarchical clustering - It is a unsupervised learning clustering algo builds a hierarchy of clusters where each node in the tree represents a cluster classified into 2 types

a) Agglomerative (bottom up) & divisive (top-down)

> Agglomerative - It is known as bottom up

) Initialization - approach or hierarchical agglomerative clustering. A Step It starts clustering by treating the individual data points as a single cluster then it is merged continuously based on similarity until it forms one big cluster containing all objects.

Steps -

- 1) Create each data point as Single cluster.
- 2) Take two closest data points or clusters and merge them to form one cluster.
- 3) Again take 2 closest data points clusters and merge them together to form one cluster.
- 4) Repeat Step 3 until one cluster left.
- 5) Once all clusters are combined into one big cluster develop the dendrogram & divide the clusters as per the problem.

> Divisive - This works just opposite of agglomerative clustering. It starts by considering all the data points into big single cluster and later on splitting them into smaller heterogeneous clusters continuously until all data points are in their own cluster. It follows a top down approach and is more efficient than Agglomerative clustering.

### Steps -

- 1) Split into clusters using any flat-clustering method say K-means.
- 2) choose the best cluster among the clusters to split further, choose the one that has the largest Sum of Squared Error.
- 3) Repeat steps 2 & 3 until a single cluster is formed.

### > Measure for the distance b/w 2 clusters -

- 1) Single Linkage - Shortest distance b/w the closest points of clusters.
- 2) Complete Linkage - Farthest distance between the 2 points of two diff. clusters.
- 3) Average Linkage - Distance between each pair of datasets is added up and then divided by the total number of datasets to calculate avg dist. b/w 2 clusters.
- 4) Centroid linkage - It is the linkage method in which the dist. b/w centroid of the clusters is calculated.

> Dendrogram - A dendrogram is a tree-like diagram that illustrates the hierarchical structure of clusters in hierarchical clustering. It's a visual representation of relationships and similarities b/w data pt. or

Key Features:

- 1) Vertical Lines - Represents data pt. or clusters.
- 2) Horizontal Lines - merging or splitting of clusters.
- 3) Height of Fusion - height of each fusion or division represents the dissimilarity b/w merged or split clusters.

1) Nodes - Points where vertical & horizontal lines meet are nodes.

2) Leaves - The individual data pt or small clusters at the bottom of dendrogram.

> Steps -

- 1) Start by treating each node data pt. as singleton cluster
- 2) Identify the two most similar cluster and merge them into new cluster.
- 3) Repeat the process, identifying & merging the next most similar clusters until all data pts. belong to a single cluster.

## A DBSCAN (Density Based Spatial Clustering of Applications with Noise)

> It groups densely grouped data points into a single cluster. It can identify cluster in large spatial datasets by looking at the local density of data points. It is robust to outliers.

> It also does not require the number of clusters to be told beforehand unlike k-Means.

> DBSCAN requires only two parameters

I) epsilon - radius of the circle to be created around each data point to check the density

II) minpoints - min. no. of data points required inside that circle for that data point to be classified as a core point.

> DBSCAN creates a circle of epsilon radius around every data point and classifies them into core point, Border point & Noise

- > If the no. of points is less than min points, then it is classified as Border point. And if there are no other data points around any data point within epsilon radius, then it is treated as Noise.
- > Direct density Reachable - A point A is directly density reachable from another point "B" if A is in  $\epsilon$ -neighborhood of "B" and "B" is core point.
- > Density Reachable - A point "A" is density reachable from B if there are a set of core points leading from "B" to "A".
- > Density Connected - A and B are density connected if there are a core point "C" such that both A and B are density reachable from "C"

### ~~\* Local Outlier Factor (LOF) -~~

- > LOF is an algorithm that identifies the outliers present in the dataset. Outlier detection method can be distribution based, clustering based and density based. LOF allows to define outliers by doing density based scoring.
- > LOF algo is an unsupervised anomaly detection method which computes the local density deviation of given dt data pt. with respect to its neighbors.
- > The LOF algo can be used for outlier detection & novelty detection. The diff. betw outlier & novelty detection lies in training dataset. Outlier detection includes outliers in training dataset. The algorithm fits the area with high density data & ignores the outliers & anomalies.

## \* OPTICS (Ordering Points To Identify clustering structure)

It is density based clustering algo that extends the concepts of DBSCAN to identify clusters of varying shapes and sizes in a dataset.

> The core concept in OPTICS is the reachability distance, which measures how easily one data point can be reached from another.

> The minimum requirement reachability distance required for a point to be considered as core point core points have at least specified no. of neighbors within specified radius.

> OPTICS constructs a reachability m plot. which is a sorted list of points and their corresponding reachability distances.

> Local maxima in reachability plot correspond to clusters in dataset. These maxima indicate points where the density changes suggesting presence of cluster.

## \* DENCLUE (Density based Clustering) -

It is a Density based clustering is a density based clustering algo focusing on modeling clusters as attractors in feature space. Aims to discover clusters with arbitrary shapes handle varying cluster densities & adapt noisy data.

> It has a solid mathematical foundation

2) It is definitely good for datasets with large noise

It allows compact mathematical descrip<sup>n</sup> of arbitrarily shaped clusters in high dimensional data sets.

- 3) It uses grid cells but only keeps info about grid cells that do actually contain data points and manages these cells in a tree based access structure.
- 4) The influence Func<sup>n</sup> is central to the gradient ascent process, where DENCLOF iteratively moves towards the maximum density points known as attractors or density peaks.
- 5) The clusters can be determined mathematically by identifying density attractors. Density attractors are local maxima of the overall density Func<sup>n</sup>.

## Intro to ANN

→ ~~Artificial~~ ANN

### Artificial Neural Network -

An ANN is a Computational Model inspired by the Structure and Functioning of the brain. It consists of Interconnected nodes also known as artificial neurons or perceptions, organized into layers.

ANNS are used for machine learning tasks, including pattern recognition, classification, regression.

> characteristics

1) ANNs consist of layers of nodes, including an input layer, one or more hidden layer and an output layer. The connections between nodes have associated weights that are adjusted during learning process.

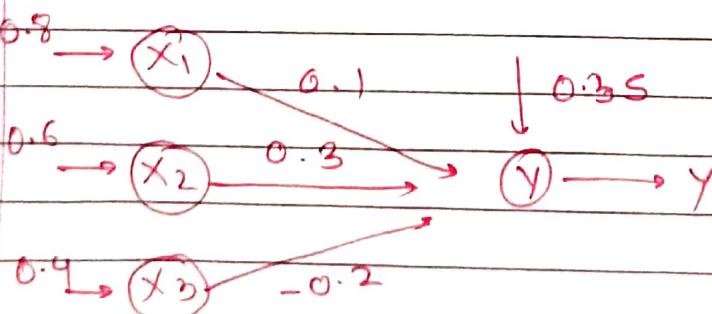
2) Each node in an ANN typically applies an activation function to its input, determining the node's output.

3) Feed Forward phase, input data is propagated through the network to produce an output.

Backpropagation phase calculates the weight and errors and adjust it to minimize errors.

Q Obtain output of the neuron Y for the network in figure using activation function as

→ binary sigmoidal & → bipolar sigmoidal.



$$x_1 = 0.8 \quad x_2 = 0.6 \quad x_3 = 0.4$$

$$w_1 = 0.1 \quad w_2 = 0.3 \quad w_3 = -0.2$$

$$y_{in} = b + \sum_{i=1}^3 x_i w_i$$

$$= 0.35 + 0.08 + 0.18 - 0.08 = 0.53.$$

$$y_{in} = 0.53$$

1) Binary Sigmoidal activation Func<sup>n</sup>

$$y = F(y_{in}) = \frac{1}{1 + e^{-y_{in}}} = \frac{1}{1 + e^{-0.53}} = 0.625$$

2) Bipolar Sigmoidal activation Func<sup>n</sup>

$$y = F(y_{in}) = \frac{2}{1 + e^{-y_{in}}} - 1 \Rightarrow \frac{2}{1 + e^{-0.53}} - 1 = 0.259$$

## A Back Propagation Learning -

Backward propagation of errors is a supervised learning algorithm used to train an ANNs.

It is a key component of the learning process in neural networks is employed to minimize the error between the predicted output and actual output.

1) Propagation Steps -

- 1) Initialize weights and biases of Neural Net randomly
- 2) Input data is propagated forward through the network to compute the predicted output. Apply activation Func<sup>n</sup> act at each layer.
- 3) Computing the error with respect to the weights and biases using the chain rule of calculus.
- 4) Computing error or difference bet<sup>n</sup> predicted output & actual target using loss function.

- 4) Propagate the error backward through the network to update the weights and biases
- 5) Adjust weight and biases by using backpropagation or optimization algo like Gradient descent.
- 6) Repeat steps for multiple iterations or epochs.

### ~~Functional Link ANN~~ -

- > FLANN is a single layer ANN with less computational complexity which is used in different fields of application such as system identification, pattern recognition
- > FLANN architecture uses a single layer Feed Forward network.
- > Using functionally expanded features FLANN overcomes the non-linearity nature of problems it encounters in single layer network.
- > It FLANN eliminates the hidden layers. Due to its single layer construction, the FLANN structure offers less computational complexity and faster convergence than MLP.

## ~~#~~ Activation Function -

- 1) Activat<sup>n</sup> Func<sup>n</sup> also known as transfer Function is used to map input nodes in certain fashion.
- 2) It decides whether to activate the neuron or not and transferred to the next layer.
- 3) Activat<sup>n</sup> Func<sup>n</sup> help in normalizing the output b/w 0 to 1 or -1 to 1. It helps in the process of backpropaga<sup>n</sup>. During backpropaga<sup>n</sup> loss Func<sup>n</sup> gets updated and for activat<sup>n</sup> Func<sup>n</sup> helps these gradient descent curves to achieve their local minima.
- 4) Activat<sup>n</sup> Func<sup>n</sup> decides in any neural net if that given input or receiving info is relevant or not.

### 5) Sigmoid Activat<sup>n</sup> Func<sup>n</sup> -

$$\sigma(x) = \frac{1}{1+e^{-x}} \quad \text{Range } (0, 1)$$

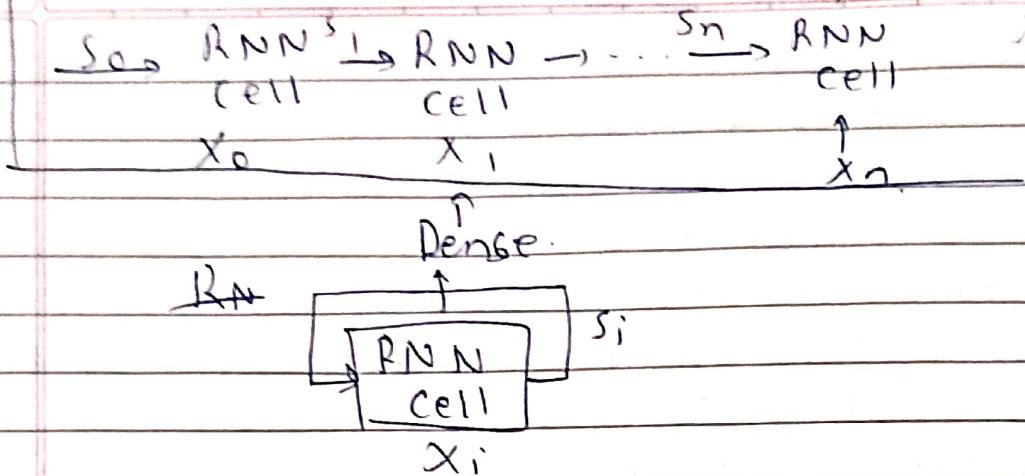
### 6) Tanh activat<sup>n</sup> Func<sup>n</sup> -

$$\tanh(x) = \frac{e^{2x} - 1}{e^{2x} + 1} \quad \text{Range } (-1, 1)$$

### 7) Rectified Linear Unit (ReLU)

$$f(x) = \max(0, x) \quad \text{Range } [0, \infty]$$

~~#~~ RNN - A RNN is a type of ANN designed for processing sequences of data. Unlike feed forward neural net, RNNs have connections that form directed cycles, allowing them to maintain a hidden state that captures info about previous input. Ability to capture temporal dependencies makes RNNs well suited for tasks involving sequences such as NLP, time series predict<sup>n</sup>. & speech recognit<sup>n</sup>. RNNs have hidden state which acts as a memory, capturing info from previous time steps.



Formula for calculating current state :

$$h_t = F(h_{t-1}, x_t)$$

Example LSTM.

**A CNN -** In CNN, the  $T$  is a class of deep neural network specifically designed for processing structured grid data such as images. CNNs are highly effective in image recognition.

Key components:

- > Convolutional layer - it is designed to perform convolution on input data. It is a mathematical operation that combines input data with a filter to extract local patterns or features.
- > Filters are learnable matrices. When convolved with input data it detects features like textures, edges.
- > Convolution operation slides across input to capture local patterns & create feature maps.
- > Activation Func.

## 2) Hidden layer - 1

- 1) Hidden layer refer to the layers between the convolution layer & output layer.
- 2) Hidden layers play a crucial role in learning hierarchical features from the input data.
- 3) Hidden layers also have learnable parameters including weights & biases which are updated during training to optimize network.
- 4) Activation Funcn. - To introduce non-linearity