

# Machine Learning

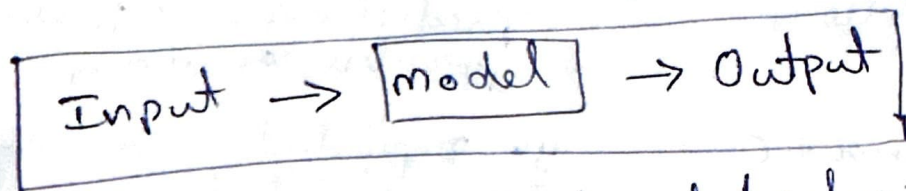


Algorithms

Linear Regression  
Logistic Regression  
K-Means Clustering  
Decision Tree  
Random Forest  
KNN Algorithm  
Support Vector Machine  
Naive Bayes Classifier  
Hierarchical Clustering

Learnings

Supervised → Labelled data  
Unsupervised → Unlabelled data  
Recognize hidden patterns  
Reinforcement → Reward system



Application: - Healthcare, Sentiment Analysis, Fraud Detection, E-commerce

Definition: - Machine learning is the science of making computers learn and act like humans by feeding data and information without being explicitly programmed.

Steps Involved: - Define Objective → Collect Data → Prepare Data  
Select Algorithm → Train Model → Test Model  
Predict → Deploy

Regression → Predict Quantity

Classification → To predict category Anomaly Detection → Detect Anomaly

Clustering → Discover structure in unexplored data

## \* Linear Regression

Independent Variable  $\rightarrow$  Is not affected by others (Rain)

Dependent Variable  $\rightarrow$  Is affected by others. (crops).

Application  $\rightarrow$  Economic growths, Product Pricing, Housing Sales, Score Preden.

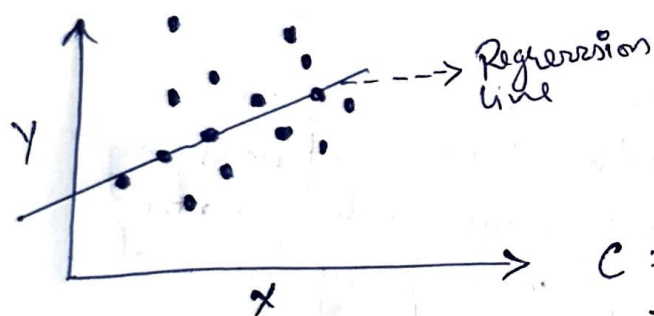
Defn:- Linear Regression is a statistical model used to predict the relationship between independent and dependent variables.

which variables in particular are significant predictors of the outcome variables. How significant is the Regression line to make predictions with highest possible accuracy.

Eqn:-  $y = mx + c$

$y$ :- Dependent  
 $x$ :- Independent

$m$ :- Slope  $= \frac{y_2 - y_1}{x_2 - x_1}$



$$m = \frac{(n * \sum (x * y)) - (\sum (x) * \sum (y))}{(n * \sum (x^2)) - (\sum (x))^2}$$

$$c = \frac{((\sum (y) * \sum (x^2)) - (\sum (x) * \sum (x * y)))}{((n * \sum (x^2)) - (\sum (x))^2)}$$

For Errors To Minimize

the Distance of line & Data Points / Finding the Best fit line

$\hookrightarrow$  use  $\rightarrow$  Sum of Sq errors, Sum of Absolute errors, Root Mean Sq errors, etc

Multiple Linear Regression

$$y = m_1 * x_1 + m_2 * x_2 + m_3 * x_3 + \dots + m_n * x_n + c$$



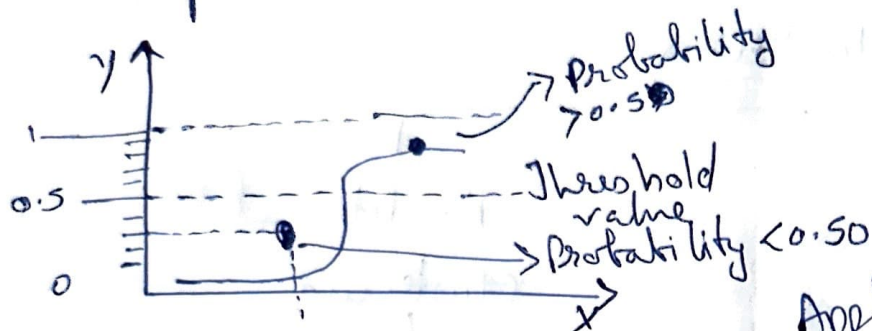
## \* Logistic Regression

Defn:- It is a classification algorithm, used to predict binary outcomes for a given set of independent variables. The dependent variable's output is discrete.

Eqn of sigmoid func

$$p(x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

$\rightarrow \log\left(\frac{p(x)}{1-p(x)}\right)$

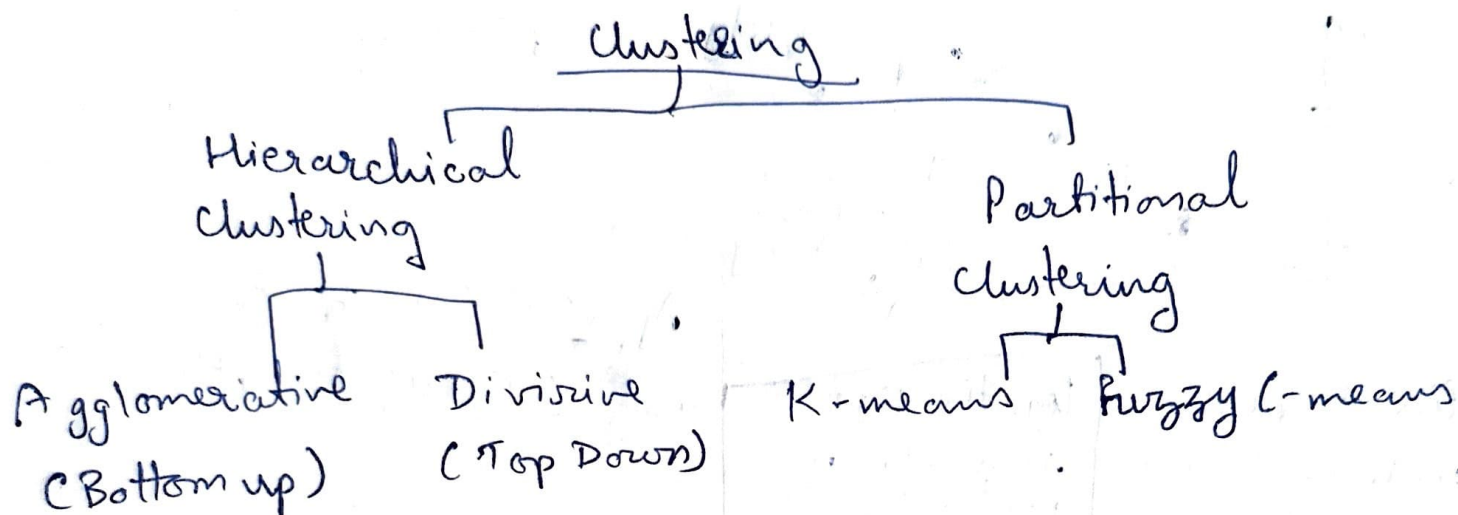


Applications:- Weather Prediction  
Image Categorization, Healthcare

Linear Regn	Logistic Regn
<ul style="list-style-type: none"><li>• Solves Regression problems</li><li>• Response variable continuous in nature</li><li>• Estimates the dependent variable in the independent variable</li><li>• Straight line</li></ul>	<ul style="list-style-type: none"><li>• Solves Classification problems</li><li>• Response variable categorical in nature</li><li>• Helps calculating the probability of an event happening</li><li>• S-curve.</li></ul>

## \* K-Means Clustering

Defn: K-means performs division of objects into clusters which are "similar" between them and are "dissimilar" to the objects belonging to another cluster.



Appln :- Academic Performance, Diagnostic Systems, Search Engines, Wireless Sensors Networks, Colour Compression

Euclidean Distance Measure :  $d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$

### Working:

start  
↓

Elbow Point (K)

↓

→ measure the Distance  
↓

① Grouping based on minimum distance  
↓

Reposition the centroids

Convergence

⊕  
IF clusters are stable

IF clusters are unstable

### Algorithm

Assuming we have inputs  $x_1, x_2, x_3, \dots$  and value of  $k$

Step 1: Pick  $k$  random points as cluster centers called centroids.

Step 2: Assign each  $x_i$  to nearest cluster by calculating its distance to each centroid

Step 3: Find new cluster center by taking the average of the assigned points.

Step 4: Repeat steps 2 & 3 until no or the clusters assignments change.



## \* Decision Tree

Problems that Decision Tree can solve :-

1) Classification.

A classification tree will determine a set of logical if-then cond<sup>n</sup>s to classify problems or continuous in nature.

2) Regression.

Regression tree is used when the target variable is numerical.

## Advantages

1) Simple to understand, interpret and visualize

2) Little effort required for data preparation

3) Nonlinear parameters don't affect its performance.

## Disadvantages

1) Overfitting occurs when the algorithm captures noise in the data.

2) The model can get unstable due to small variation in data

3) A highly complicated Decision tree tends to have a low bias which makes it difficult for the model to work with new data

## IMP Terms

Entropy - It is the measure of ~~decrease in~~ entropy of the randomness or unpredictability in the dataset

Information Gain - It is the measure of decrease in entropy after the dataset is split.

Leaf node - carries the classification or the Decision



## \* Random Forest

Appl<sup>n</sup> - Remote Sensing, Object Detection, Kinect (Gaming)

Why Random Forest?

- 1) No overfitting
- 2) High Accuracy
- 3) Estimates Missing Data

What is Random Forest?

Random Forest or Random Decision Forest is a method that operates by constructing multiple Decision Trees during training phase. The decision of majority of the trees is chosen by the random forest as the final decision.

## \* KNN

Def<sup>n</sup> :- K Nearest Neighbors, is one of the simplest supervised Machine Learning algorithm mostly used for classification.

How do we choose the factor  $k$ ?

KNN Algorithm is based on feature similarity. Choosing the right value of  $k$  is a process called parameter tuning, and it is important for better accuracy.

To choose value of  $k$ :-

- $\text{Sqrt}(n)$ , where  $n$  is the total number of data points.
- Odd value of  $k$  is selected to avoid confusion between two classes of data.



When do we use KNN?

- 1) Data is labelled
- 2) Data is noise free
- 3) Dataset is small.

## \* Support Vector Machine

Appl<sup>n</sup>: Face Detection, Text & Hyper Text Categorization, Classification of images, Bioinformatics.

Why SVM?

SVM is a supervised learning method that looks at data and sorts it into one of the two categories.

Advantages

- 1) High Dimensional Space
- 2) Sparse documents vectors.
- 3) Regularization parameters.

## \* Naive Bayes Classifier

\* Bayes Th<sup>m</sup> - 
$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

Where is Naive Bayes Used?

- 1) Face Recognition
- 2) Weather Prediction
- 3) Medical Diagnosis
- 4) News Classification

Advantages

- 1) Very simple & Easy to implement
- 2) Needs less training
- 3) Handles both continuous & discrete data
- 4) Highly scalable with numbers of predictors & data points
- 5) As it is fast, it can be used in real time predictions
- 6) Not sensitive to irrelevant features.





## \* Applications of Machine Learning

- 1) Virtual Assistants
- 2) Traffic Predictions
- 3) Social Media Personalization
- 4) Email Spam Filtering
- 5) Online Fraud Detection
- 6) Assistive Medical Tech
- 7) Automatic Translation

## \* What does ML Engineer Do?

- Creates & maintains ML solutions to solve business problems
- Optimizes these solutions for performance & scalability
- Solves business problems like reducing customer churn, running targeted marketing campaigns and improving product experience.
- Contributes to cutting edge research in AI & ML

## \* Math Required for ML (IMP Topics)

\* Probability & Statistics - Bayes Theorem, Probability Distribution, Sampling, Hypothesis testing

\* Linear Algebra - Matrices, Vectors

\* Calculus - Differential calculus, Integral calculus.