

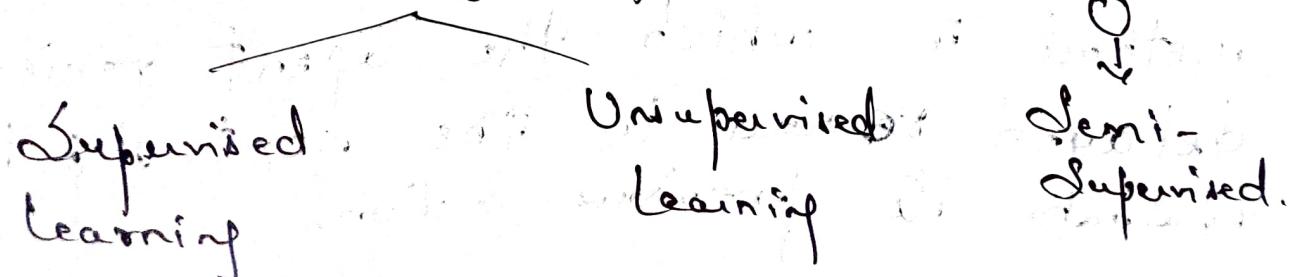
Machine Learning

Machine learning serves as the foundational technical infrastructure for data mining. Its primary function is to extract information from raw data stored in databases, transforming this information into a comprehensible form that can be applied across various situations.

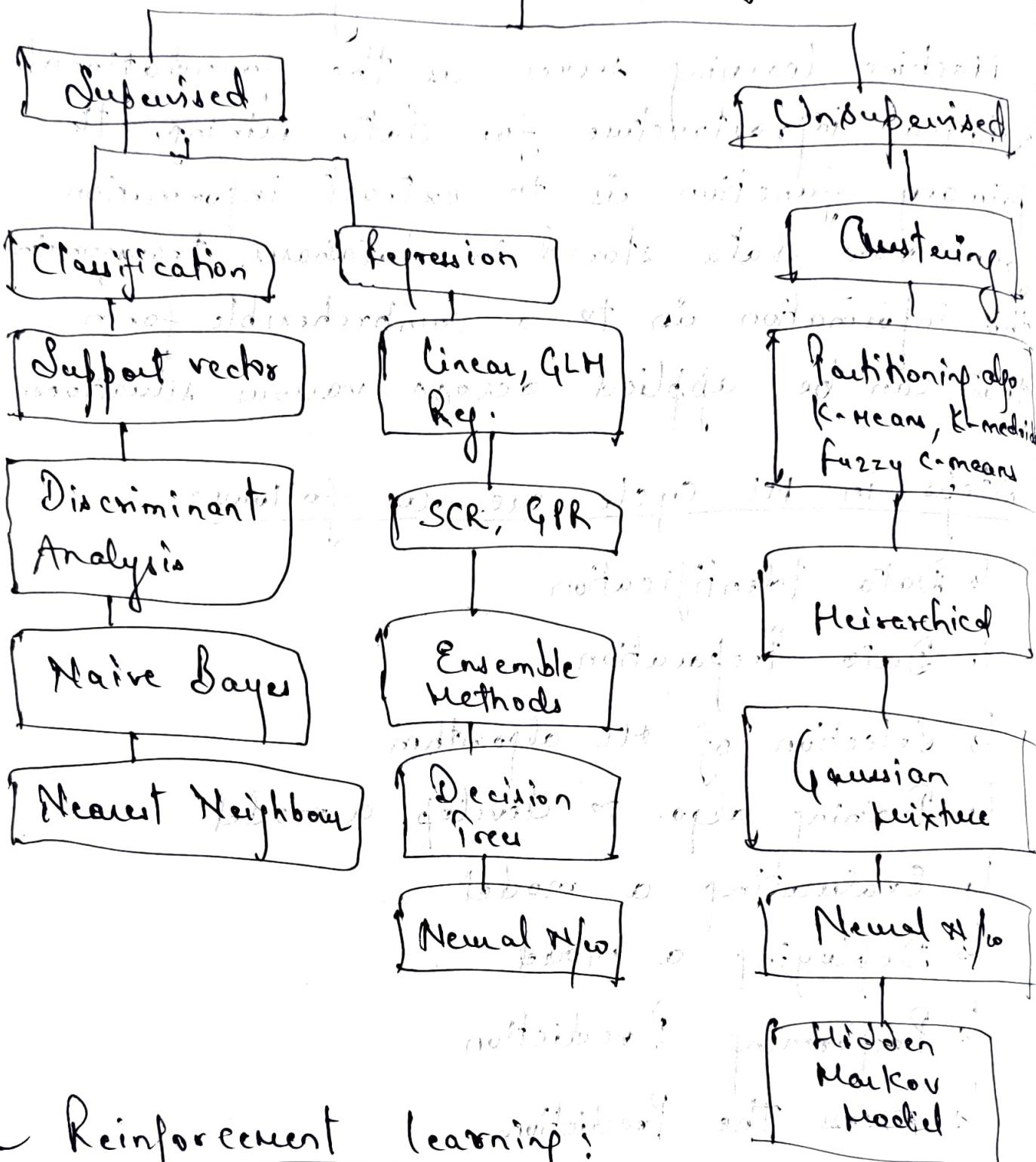
Steps in ML cycle are as follows:

- ↳ Data Identification
- ↳ Data Preparation
- ↳ Selection of ML algorithm
- ↳ Training algo. to develop a model
- ↳ Evaluating a model
- ↳ Deploying a model
- ↳ Performing Prediction
- ↳ Assess the Predictions

Classes of learning algo.



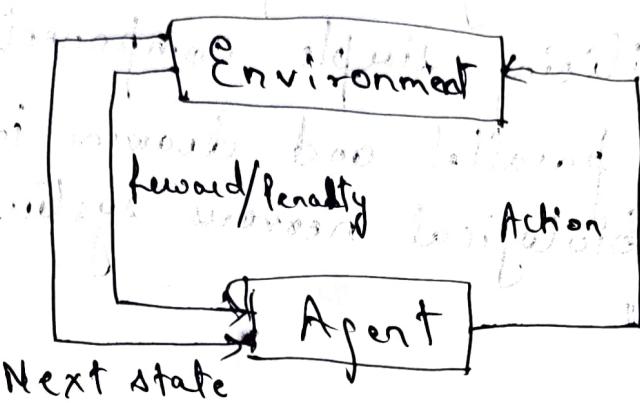
Machine Learning Algorithms



Reinforcement learning:

refer to a kind of machine learning method in which the agent receives delayed reward in the next step to evaluate its previous action.

RL setup is composed of two components: an agent and an environment.



It may be broken down into two categories: Model-free and Model-based.

Deep learning: serves as a method for machines to acquire knowledge by utilising deep neural networks. It is used to solve practical problems in fields like computer vision (image), natural language processing (text), and automated speech (audio).

Neural Networks: are comprised of neurons. Each neuron receives signals as input, multiplies them by weights, adds them together, and then applies a non-linear function. These neurons are organised in layers and closely stacked to each other.

- A deep neural network is characterised by the integration of multiple non linear processing layers, utilising simple components that operate in parallel and drawing inspiration from the biological nervous system of living organisms.
- Convolutional Neural Net's (CNN): The idea that underlies them is that rather than linking each neuron with all of the ones that come after it, we just connect it with a select few of those that come after it.
- Recurrent Neural Net's: find application in time series forecasting due to their suitability for temporal data.
- Recursive Neural Net's: set up in a tree-like manner. As a result, they can simulate the hierarchical structures of training datasets.
- Auto-Encoders: type of unsupervised technique that is used to reduce dimensionality & compress data.

Restricted Boltzmann Machines (RBMs): are stochastic neural nets that can learn a probability distribution over their inputs & so have generative capabilities. They take the input & create a representation of it in the forward phase of the training. They rebuild the original output from the representation in the backward pass.

Generative Adversarial N/w (GAN): construct two models: make up fake data (generator) & tell difference b/w actual and fake data (discriminator) and turned them against one another.

Graph Neural N/w: Deep learning doesn't operate well with unstructured data in general. GNNs are designed for modelling graph data, allowing them to identify and represent the relationships b/w nodes in a n/w as embeddings typically in the form of integers. This enables their apps in various ML models for tasks such as clustering, classification & more.

Ensemble learning: that enhances predicting performance by amalgamating predictions from diverse models. It is further categorised as:

b) Bagging: Ensemble learning is the process of fitting "multiple" decision trees to various samples of same dataset and averaging the results. Popular algo. are based on this approach including:

b) Bagged Decision Tree (canonical bagging)
and Random Forest and Extra Trees.

c) Stacking: Seeks diversity among its members by employing various types of models fitted to the training data and utilising a model to aggregate predictions. Popular algo.:

b) Stacked Models (canonical stacking)

b) Blending

b) Super Ensemble

Boosting: It's designed to modify the training data, emphasising instances that were incorrectly predicted by previous models.

Popular EL algorithms include:
↳ AdaBoost (canonical Boosting)
↳ Gradient Boosting Machines
↳ Stochastic Gradient Boosting
(XGBoost and similar)

Supervised Learning: entails training a model using data where both inputs and outputs are already identified, enabling the model to make predictions about future outputs. E.g: Medical Imaging, speech recognition etc.

Unsupervised Learning: analysis data to uncover previously unknown patterns or structures. It is used to infer conclusions from sets of data that contain input but no tagged answers. E.g: Market Research etc.

Semi-Supervised Learning: algo. are trained on small sets of labelled data before being applied to unlabeled data, like in UL.

Naïve Bayes : represents a case of statistical classification, wherein it assesses the probability that a specific sample belongs to a particular group based on given sample. When faced to big Database, the Bayesian classification demonstrates both improved accuracy and increased speed.

Steps to perform naïve Bayes Algorithm:

- ↳ Handling data: Data is loaded from CSV file and spread into training & tested assets.
- ↳ Summarising the data: A particular prediction is made using a summarise of the data set to make a single prediction.
- ↳ Making a Prediction: Summarise the properties in the training dataset, to calculate the probabilities & make predictions and follow up, consider the next set.
- ↳ Making all the Predictions: Generate prediction given a test data set and a summarise data set.

6 Evaluate Accuracy: Accuracy of the prediction model for the test dataset, as a percentage correct out of them all, the predictions made

6 Putting all together: finally, we tie all the steps together and formed our own model of Naive Bayes Classifier.

K-Nearest Neighbours (K-NN): This method assigns items to the class to which they are closest in terms of proximity to their Neighbours. It involves the calculations of the distance b/w an item and a class to make such determinations. Classes are represented by centroid (central value) and individual points.

Decision Trees: Basic Steps are as follows,

6 Building the tree by using the training set dataset/database.

6 Applying the tree, to the new dataset/database.

Logistic Regression: A key component of the supervised learning approach. Its purpose is to forecast the categorical dependent variables.

- LR predicts the outcome of a dependent variable that has a "yes" or "no" answer.
- LR shares similarities with Linear Regression, yet their application differs.
- In LR, we fit a S-shaped Logistic function which predicts two maximum values, instead of a regression line (0 or 1).
- The curve from the logistic function shows how likely something is, like whether the cells are cancerous or not, whether a movie is overweight or not, etc.

Types of LR:

- Binomial: dependency variables can be either 0 or 1, pass or fail, etc.
- Multinomial: dependent variables can be one of three types that are not in order, such as cats, dogs or sheep.
- Ordinal: dependent variables can be ranked such as low, medium, high.

Support Vector Machine: a type of algo:

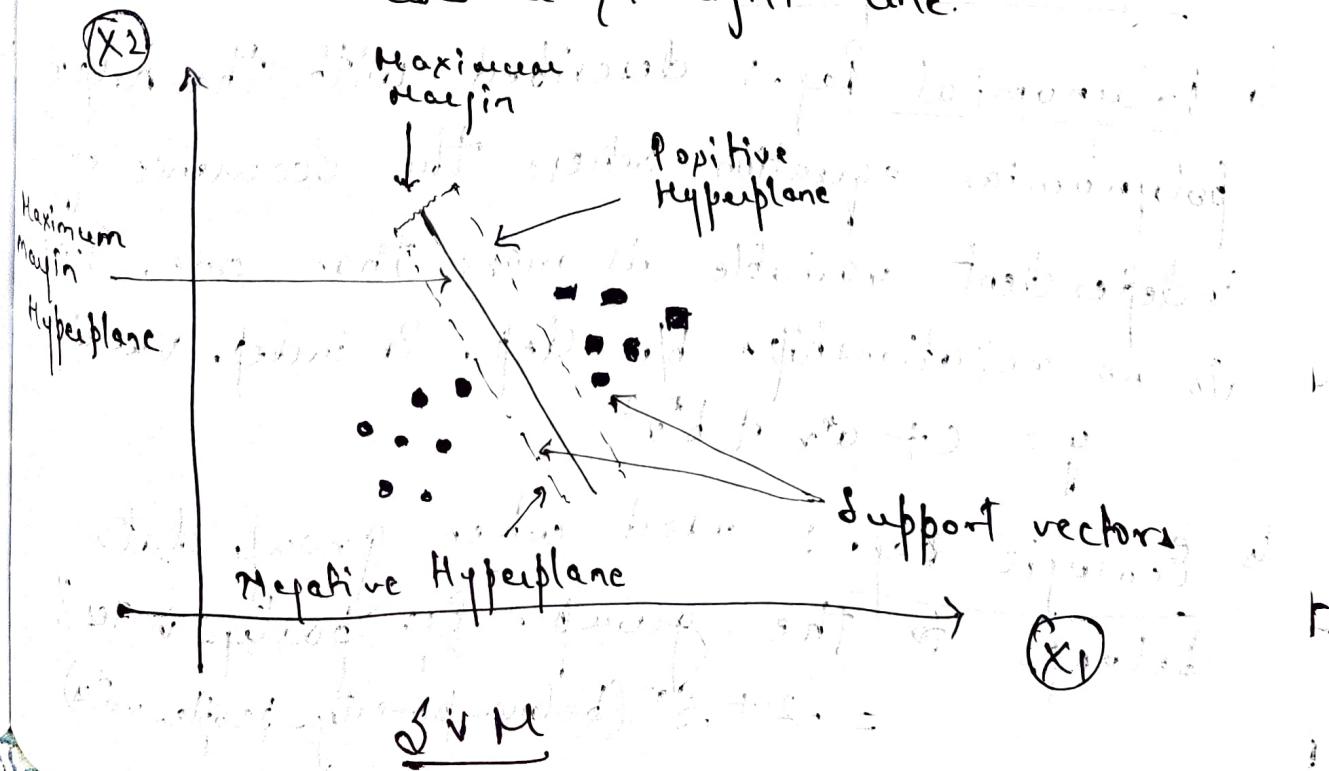
that transforms the primary training data into a new format that has a higher dimension by making use of a non-linear mapping.

Two main categories that can be applied to SVM: linear & non-linear classification.

Linear SVM: can use linear SVM if the data can be completely separated into linear categories.

Non-linear SVM: when the data is not linearly separable means we apply advanced techniques like kernel trick to categorise the data points that can't be divided into two classes by a straight line.

(X₂)



Various types of regression algorithms:

↳ Linear Regression: used when there is a single dependent variable, and there can be one or more independent variables. In linear reg. the relationships b/w the dependent and independent variables are linear, i.e. of type: $y_i = a + b^*x_i$; Eg.: Child height = $a + b^*$ (parent height)

↳ Multiple Linear Regression: when there is only one dependent variable and more than one independent variables, then it results in MLR. i.e., $y = a + b_1x_1 + c_2x_2 + d_3x_3$. Eg.: $a + b^*$ (daily meal) + c^* (daily exercise).

↳ Logistic Regression: Discussed earlier.

↳ Polynomial Reg.: described with the help of polynomial equation where the occurrence of independent variable is more than one. There is no relationship b/w dep. & indep. variables.

$$y = c + a^*x + b^*x^2$$

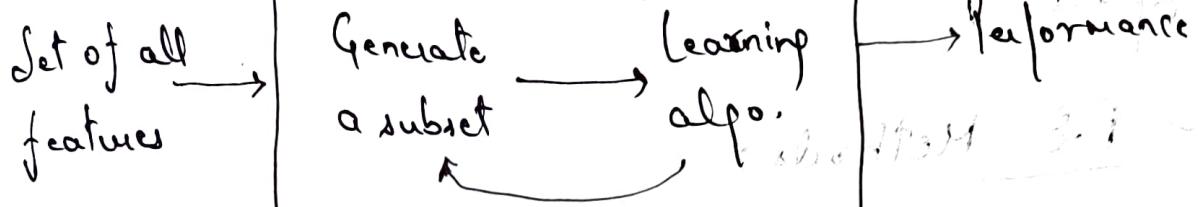
↳ Ecologic Reg.: used when group data belongs to the group. Eg.: party-votes% = $0.2 + 0.5^*$ (below-poverty-people-vote)

- ↳ Lasso Regression: is a regularization technique employed in scenarios where data variables exhibit strong correlations.
- ↳ Lasso Regression: coefficients are penalized, & this involves the application of the least absolute shrinkage and selection operator.
- ↳ Logistic Reg.: predictor and response variable both are binary in nature & applicable to both classification and regression problem.
- ↳ Bayesian Reg.: also, relies on Bayesian statistics, utilizing random variables as parameters for estimates. In this approach, when data is unavailable, the algo. incorporates prior data as input.
- ↳ Quantile Reg.: when the focus is on determining the boundaries of quantiles. Eg.: Health Analysis (assessing over & underweight).
- ↳ Cox Reg.: used when output of a variable depends on set of independent variables.
Eg.: patient survival after surgery (Survived, Died) = (age, condition, BMI)

- Dimensionality Reduction: is a technique used to reduce the number of features in a dataset while retaining as much of the important information as possible.
- There are two main approaches to dimensionality reduction:
 - ↳ feature Selection
 - ↳ feature Extraction
- feature Selection: involves selecting a subset of the original features that are the most relevant to the problem at hand. The goal is to reduce the dimensionality of the dataset while retaining the most important features.
- feature Selection Methods:
 - filter Methods
 - Wrapper Methods
 - Embedded Methods
 - Hybrid Methods
- forward feature Selection: first step is to evaluate each individual feature and

choose the one that results in the most effective algo. model. This is referred to as forward feature selection.

Selecting the best subset



(iii) Sequential Feature Selection

- Procedure:
 - ↳ Train the model with each feature being treated as a separate entity, then evaluate its overall performance.
 - ↳ Select the variable that results in the highest level of performance.
 - ↳ Carry on with the process while gradually introducing each variable.
 - ↳ The variable that produced the greatest amount of improvement is the one that gets kept.
 - ↳ Perform the entire process once more until the performance of model doesn't show any meaningful signs.

→ feature Extraction: involves creating new features by combining or transforming the original features. The goal is to create a set of features that captures the essence of the original data in a lower dimensional space.

→ F.E Methods:

↳ Principal Component Analysis (PCA)

↳ Linear Discriminant Analysis (LDA)

↳ t-distributed stochastic neighbour embedding (t-SNE)

→ PCA: is a way to get important variables from a large set of variables in a dataset. PCA is more useful when you have data with three or more dimensions.

→ LDA: is the recommended linear classification method for situations involving two or more than two classes.

→ Problems associated to logistic regression:

↳ Instable, but with well-defined classes. instability when there are only few occurrences.

Association Rules: play a crucial role and find extensive applications in scenarios such as market basket analysis.

AR do not reveal the customer's preference for specific items. Instead, they identify relationships among items that are typically purchased together by customers.

Concept:

↳ **k-itemset**: a set of k items. Eg: 2-itemset can be {pencil, eraser} or {bread, butter}, etc, 3-itemset can be {bread, butter, milk}

↳ **Support**: $\text{Support}(x) = \frac{\text{No. of transaction containing } x}{\text{Total no. of considered transac.}}$

↳ **Confidence**:

$\text{Confidence}(y) = \frac{\text{No. of transac. containing } x \& y}{\text{No. of considered containing } x}$

↳ **frequent itemset**: An itemset that has a support value equal to or greater than the minimum support threshold is termed a frequent itemset.

Eg: Apriori Algorithm.

✓ FP Tree Growth: Despite a notable decrease in the number of candidate items, the Apriori Algo. may still face slowness as it involves other iterations of scanning of the entire transaction set. In FP growth, it employs divide and conquer strategy to generate frequent itemsets themselves.

Algorithm:

- ↳ Create a FP Tree, by compressing the transaction database. Along with preserving the information about the itemsets, the tree structure also retains the association among the itemsets.
- ↳ Divide transac. database into a set of conditional databases where each associated with one frequent item or 'pattern frequent' and examine each db separately.
- ↳ for each 'pattern frequent', examine its associated itemsets only.

Clustering: Cluster analysis, also known as clustering, is a technique used to partition a set of data objects into subsets based on individual observations.

Example of Clustering is outlier detection where Credit card fraud and criminal activities are monitored.

- Applications in real-world scenarios:

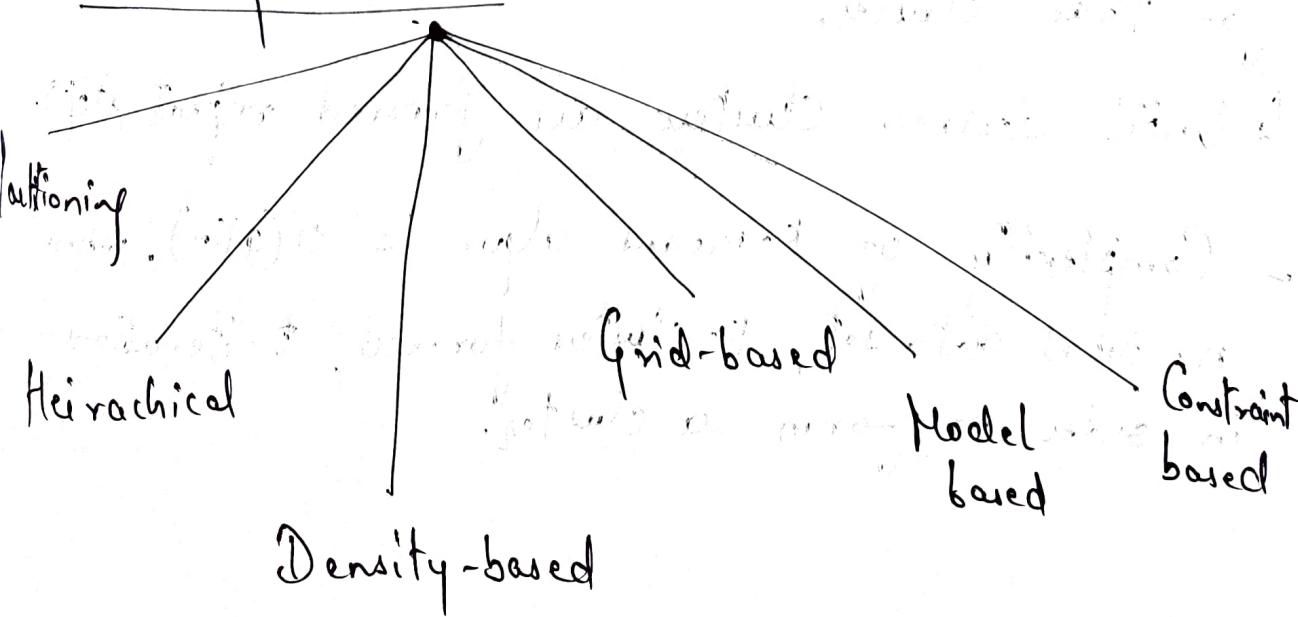
↳ Customer segmentation: telecom, ecommerce, sports etc.

↳ Document Clustering: cluster similar doc together

↳ Image Segmentation: club similar pixels in image

↳ Document clustering: document clustering

- Clustering Methods:



K-means Algorithm: primary objective is to determine the cluster locations, minimising the distance between the cluster and dataset. Also called (Lloyd's algorithm).

- Algorithm will be performed in following steps:
- ↳ Define no. of clusters (k), to be produced and identical data point centroids.
 - ↳ The distance from every data point to all the centroids is calculated and the point is assigned to the cluster with a minimum distance.
 - ↳ Following the above setup for all the data points, the average of the data points present in a cluster is calculated and can set new centroid for that cluster.
 - ↳ Until desired clusters are formed repeat step 2
 - Complexity of K-means algo. is $O(tKn)$, where n - total dataset, k - clusters formed, t - iterations required in order to form a cluster.

K-medoid or PAM (Partitioning Around Medoids):

In PAM, the medoid of the cluster has to be an input data point while this is not true for k-means clustering as the average of all the data points in a cluster may not belong to an input data point.

- Algo. is implemented into two steps:
 - ↳ Build: Initial Medoids are innermost objects.
 - ↳ Swap: A fm can be swapped by another fm until the fm can no longer be reduced.
- Algo:
 - ↳ Initially choose m random points as initial medoids from given dataset.
 - ↳ for every data point assign a closest medoid by distance metric.
 - ↳ Swapping cost is calculated for every chosen & unchosen object given as T_{ns} where s is selected and n is non-selected object.
 - ↳ If $T_{ns} < 0$, s is replaced by n
 - ↳ Until there is no change in medoids, repeat 2 and 3.

Characteristics:

