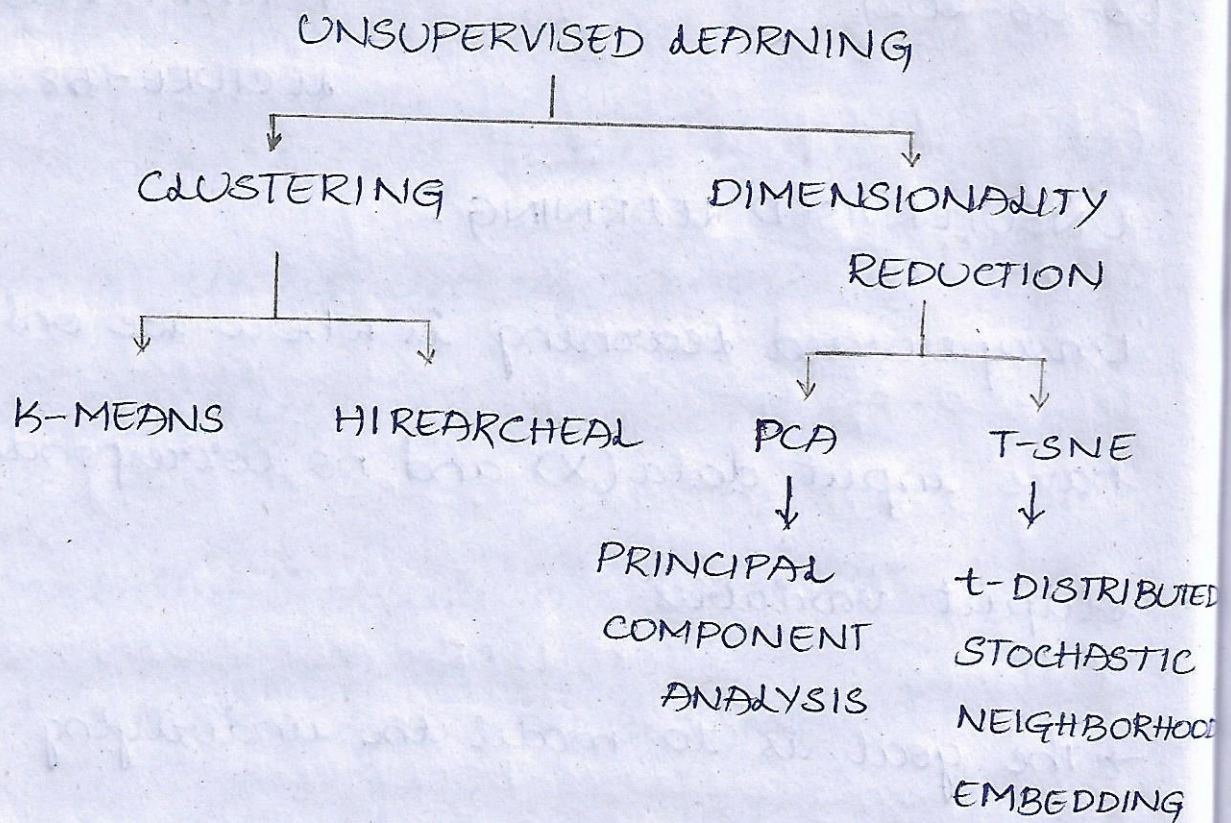## UNSUPERVISED LEARNING :

Unsupervised learning is where we only have input data (X) and no corresponding output variables.

→ The goal is to model the underlying structure (or) distribution in the data in order to learn more about the data.

## TECHNIQUES OF UNSUPERVISED LEARNING :

Common Algorithms used in unsupervised learning include

→ Clustering

→ Anomaly Detection

→ Neural Networks & approaches for learning latent variable models.

UNSUPERVISED LEARNING

- CLUSTERING
  - K-MEANS
  - HIREARCHEAL
- DIMENSIONALITY REDUCTION
  - PCA → PRINCIPAL COMPONENT ANALYSIS
  - T-SNE → t-DISTRIBUTED STOCHASTIC NEIGHBORHOOD EMBEDDING

## GUASSIAN MIXTURE MODEL (GMM):

The traditional GMM for pattern recognition is an unsupervised learning method.

→ GMM ~~model~~ are a probabilistic model for representing normally distributed subpopulat- -ions within a over population.

→ Estimating the parameters of the individual normal distribution components

is a canonical problem in modelling data with GMM's.

CLUSTERING : NOT A PREDICTION

clustering is a machine learning technique that involves the grouping of data points.

→ In theory, Data points that are in the same group should have similar properties and/or features, while data points in different groups should have highly dissimilar properties and/or features.

WHY CLUSTERING IS USED?

→ Useful for explorating / exploring data.

# K-MEANS CLUSTERING:

An unsupervised learning algorithm, which groups unlabeled dataset into different clusters.

→ Its an iterative algorithm that divides the unlabeled dataset into 'k' different clusters in such a way that each dataset belongs only one group that has similar Properties.

The Hyperparameter in k-means clustering is 'k'.

\* K-means & Hirearcheal clustering are HARD CLUSTERING & GMM is a SOFT CLUSTERING.

INTER CLUSTER DISTANCE: AS MAX. AS POSSIBLE

It shows the distance between data point with cluster center.

INTRA CLUSTER DISTANCE: AS MIN. AS POSSIBLE

It shows the distance between the data point of one cluster with the other data point in other cluster.

STEPS INVOLVED IN K-MEANS CLUSTERING:

→ Initialise 'k':

  -Randomly pick 'k' points from data & assign them as centroids $c_1, c_2, ....., c_k$
° "We generally take 'K' value as '3'."
→ Assignment Step:

  -Iterate i.e., for each data point (x) in data

      * Select nearest centroid $c_J$ where $J \in [1,2,3]$
      * Add 'x' in $S_J$ → list of $S_J$ [    ]

→ Recompute the centroids

$$\Rightarrow c_i = \frac{1}{n} \sum x_j \quad \text{where } x_j \in s_i$$

→ Repeat step 2 & 3, until it reaches convergence

NOTE:

So, In step-3:

- The minimum distance in the data is considered

$$c_1 \rightarrow s_1 [ \quad ] , \quad c_2 \rightarrow s_2 [ \quad ] , \quad c_3 \rightarrow s_3 [ \quad ]$$

From $c_1, c_2, c_3$ pick the middle values.

* → Centroids doesn't move if they are at the center.

WHAT DOES K-SCORE MEAN?

The K-means objective is to reduce the sum of squares of the distances of points from their respective cluster centroids.

→ It has other names like J-squared error function, J-score (or) within-cluster sum of squares.

→ This value tells how internally coherent the clusters are.

$$\Rightarrow \min \sum_{i=1}^{K} \sum \| x - c_i \|^2 \text{ where } x \in S_i$$

for all(V) data points (x)     Nearest centroid

## ADVANTAGES:

↳ If variables are huge, then K-means most of the times computationally faster than hierarchical clustering, if we keep 'K' small.

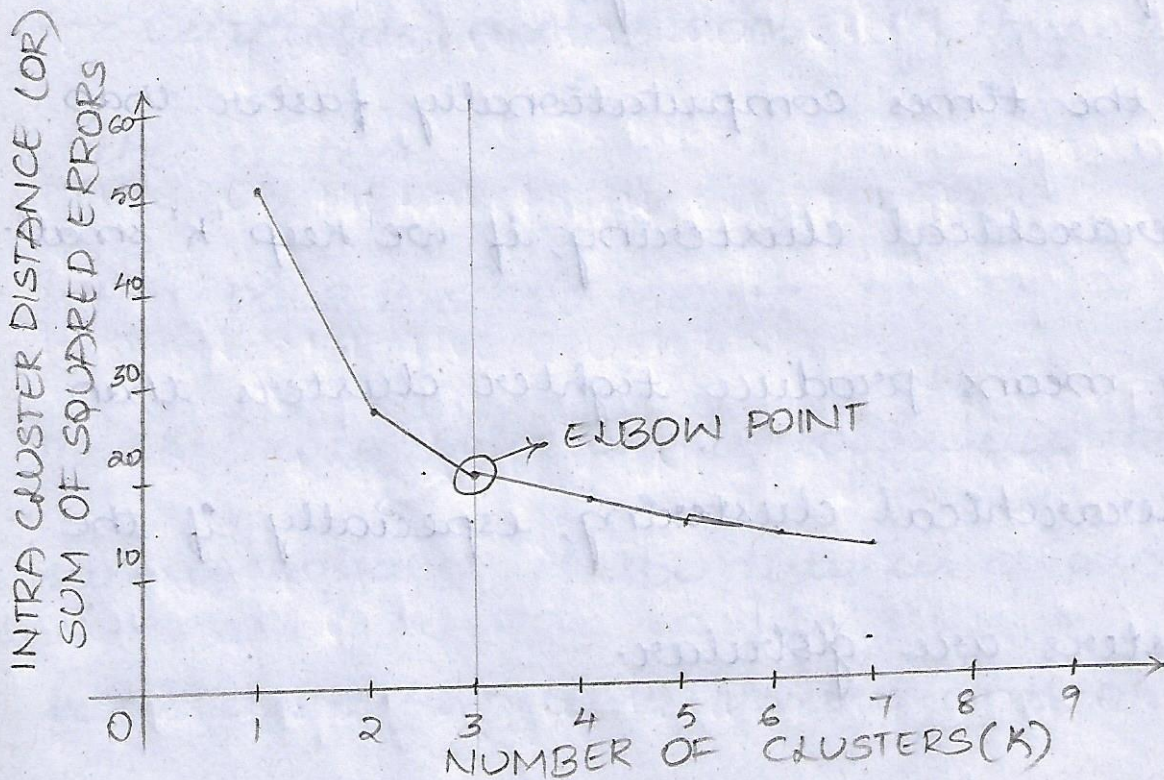↳ K-means produce tighter clusters than hierarchical clustering, especially if the clusters are globular.

DISADVANTAGES :

→ Difficult to predict K-value.

HOW TO FIND THE BEST VALUE OF 'K' in K-MEANS?

1. Compute clustering algorithm for different
   ^K-means

   values of K.

2. For each K, calculate the total within-

   cluster sum of square (WSS).

3. Plot the curve of "WSS" according to the

   number of clusters K.

ELBOW POINT :

the point which defines the optimal number of clusters is known as "Elbow point".

→ It can be used as a visual measure to find out the best pick for the value of K.

NOTE :

As the K increases, the intra cluster distance decreases.

(value)