

(Unsupervised Learning)

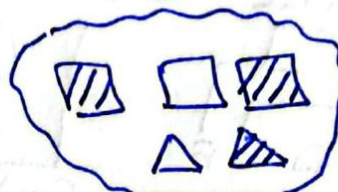
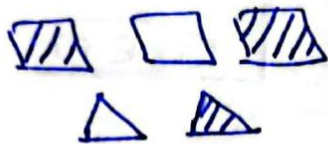
→ For data that is

- unclassified
- Unlabelled
- More complex
- Moderately accurate but reliable results.

→ Used for

- Finding patterns (clustering)
- Anomaly detection
- ~~Pattern~~

Example



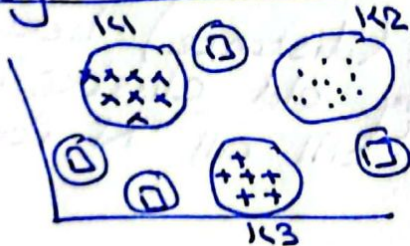
Based on shape
(similarity)

or



Based on color

Anomaly Detection



- Fault detection
- Intrusion "
- System fault "

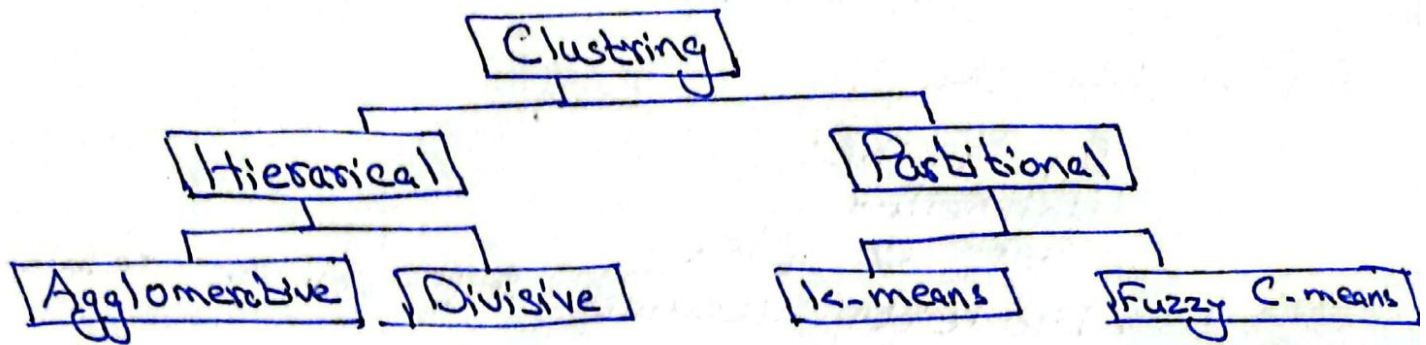
→ Need for clustering

in unlabelled, structured & unstructured data

- To determine intrinsic grouping.
- To organize data into clustering showing the internal structure of data.
- To partition the data points.
- To understand and exhibit value from large sets of structured & unstructured data.

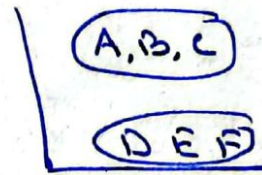
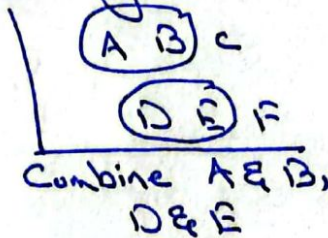
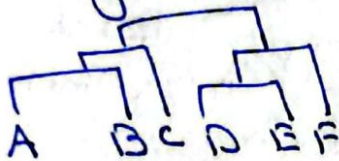
→ Types of Clustering:

(2)



→ Hierarchical Clustering

• Occupies hierarchy, a structure more informative than the unstructured set of clusters returned by flat clustering.



Steps:-

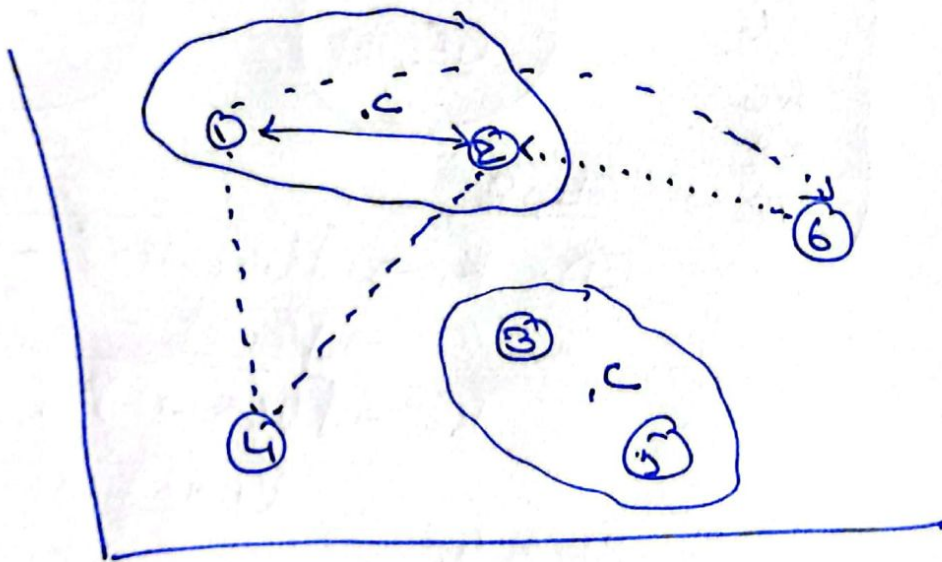
- 1) Assign each item to its own cluster, (eg. if there are N items, you will have N clusters)
- 2) Find the closest (most similar pair) of clusters & combine them.
- 3) Compute similarities (distance) between the new clusters & every old cluster. Then combine
- 4) Repeat step 2 & 3 till all ' N ' items are in single cluster.

→ Partitional Clustering

• Division of data into non-overlapping clusters where a data object is only in one set (cluster).

→ Distance Measure (3)

- 1) Complete Linkage Clustering
↳ maximum possible distance between points.
- 2) Single Linkage Clustering
↳ Minimum possible distance between points.
- 3) Mean Linkage Clustering
↳ find all possible pair-wise distance between pair-wise two clusters & then calculate the average distance.
- 4) Centroid Linkage clustering
↳ find centroids of each cluster & calculate the distance between them.



→ K-Means Clustering (4)

Step 1: Choose cluster's ($k=2$ eg. k_1, k_2) Centroids.

Step 2: Calculate Euclidean Distance of each point(item)

$$ED = \sqrt{(x_p - x_c)^2 + (y_p - y_c)^2}$$

Step 3: Put the point(item) with smallest (nearest) ED in respective cluster.

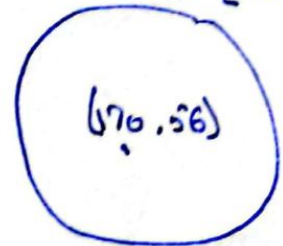
Step 4: Recalculate the respective cluster's centroid with new addition.

Step 5: Repeat step 2~4.

Step 2 $k_1 \{1, 4, 5\}$



$k_2 \{2, 3\}$



Step 3

$$ED = k_1 \rightarrow \sqrt{(168 - 185)^2 + (60 - 72)^2}$$

$$\rightarrow \boxed{20.8}$$

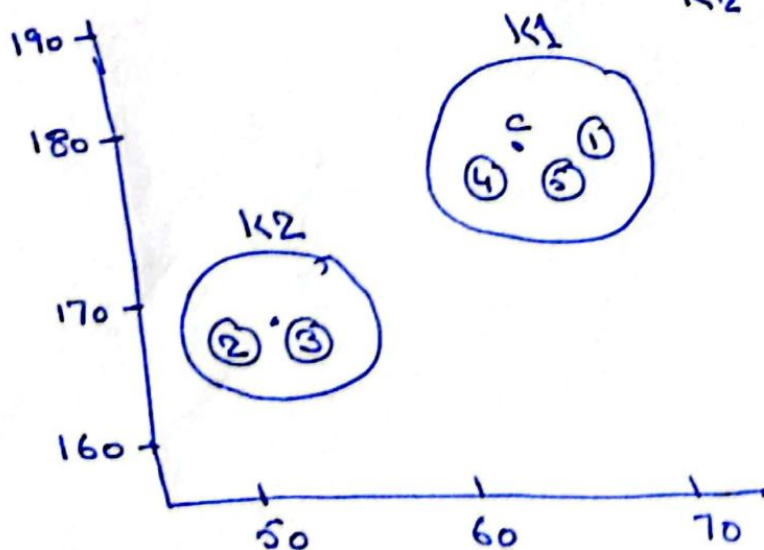
$$k_2 \rightarrow \sqrt{(168 - 170)^2 + (60 - 56)^2}$$

$$\boxed{4.48}$$

Step 4:

$$k_2 \text{ Centroid} = \left(\frac{170 + 168}{2}, \frac{60 + 56}{2} \right)$$

$$= (169, 58)$$



→ Optimum Cluster Number: (5)

- 1) Within-Sum-of-Squares (WSS): Total distance of data points from their respective centroids.
- 2) Total-Sum-of-Squares (TSS): Total distance of data points from global mean for a given data point. This is constant.
- 3) Between-Sum-of-Squares (BSS): Total weighted distance of various cluster centroids to the global mean of data.
- 4) R^2 : R-square is the total variance (BSS/TSS)
- 5) Sum-of-Square-Error (SSE): is Euclidean distance of each point to its closest centroid.

→ How many clusters:

- It is a fundamental issue in k-means clustering.
- If SSE plotted against k , you will see the error decreases as k increases becuz their size decreases & hence distortion is also small.
- The goal of the 'Elbow Method' is to choose k , where SSE decreases abruptly.

From $k_1 \dots k_2$
(WSS for each value of k)



→ Three most popular

- 1) Silhouette Coefficient
- 2)
- 3)

Silhouette Coefficient (SC)

(6)

→ We have to compute
Step 1) SC of each point

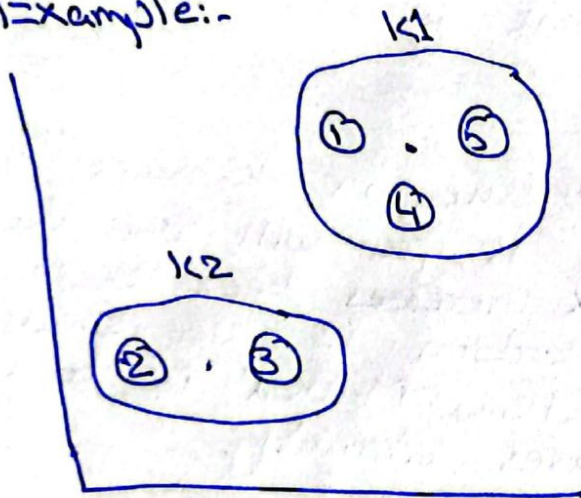
$$1 - \frac{a}{b}$$

(avg dist of a point to all other points in a cluster) / (Minimum avg dist of a point to all points in another cluster)

Step 2) SC of each cluster

Step 3) of all clusters.

→ Example:-



Step 1

$$a = \frac{\{(1 \rightarrow 5) + (4 \rightarrow 5)\}}{2 \text{ (No. of points)}}$$

$$b = \frac{\{(1 \rightarrow 2) + (1 \rightarrow 3)\}}{2}$$

$$\text{SC of } ① = 1 - \frac{a}{b}$$

Step 2: SC of each cluster

Let's suppose SC of ② & ③ is x & y.
respectively

$$\text{SC of } K2 = \frac{x+y}{2}$$

Step 3: Overall SC

$$\text{SC of } K1 \text{ \& } K2 = \frac{(\text{SC of } K1) + (\text{SC of } K2)}{2}$$