

Clustering
→ K-mean[^] Algo. aims at minimizing an objective function
→ known as squared error function

$$J(V) = \sum_{i=1}^C \sum_{j=1}^{C_i} (||x_i - v_j||)^2$$

where $||x_i - v_j||$ = Euclidean distance b/w x_i & v_j

C_i = # data points in i th cluster

C = # cluster centers.

Algorithm :-

Let $X = \{x_1, x_2, \dots, x_n\}$ be data points & $V = \{v_1, v_2, \dots, v_c\}$ be set of centers

- ① Randomly select 'c' cluster centers. (as far away as possible) in cunning way
- ② Calculate the distance b/w each data points & cluster center
- ③ Assign the data point to the cluster center whose distance from the cluster center is minimum of all the cluster centers.
- ④ Recalculate the new cluster center using:

$$v_i = \frac{1}{C_i} \sum_{j=1}^{C_i} x_j$$

[Take mean of all data points assigned to centroid cluster]

- ⑤ Recalculate the distance b/w each data point & new obtained cluster centers.
- ⑥ If no data point was reassigned then stop, otherwise repeat from step-③.
- ⑦ The results of K-mean clustering algorithm are:
 - ① The centroids of K-cluster, which can be used to label new data
 - ② Labels for training data (each data point is assigned to a single cluster).

→ Clustering is task of dividing the datapoints into # groups/# clusters that data points in the same group are more similar to other data points in the same group than those in other group.

→ Application of clustering :-

- Recommendation engine
- Market segmentation
- Social N/w Analysis
- Search result grouping
- Medical Imaging
- Image Segmentation
- Anomaly detection

Adv of K-mean :-

- ① Fast, robust, easier to understand
- ② T.C. $O(tKnd)$
 $t = \# \text{ iterations}$, $K = \# \text{ clusters}$
 $n = \# \text{ objects}$, $d = \# \text{ dimensions each obj}$
 $K, t, d < n$
- ③ Gives best result when data set are distinct (or) well separated to each other.

Dis of K-mean :-

Hierarchical clustering :- 2 types

i) Agglomerative HCA (Bottom-up Approach)

ii) Divisive HCA (top-down approach) → Reverse of Bottom-up approach

Algo :-

Given set of N items to be clustered, $N \times N$ distance (similarity) matrix.

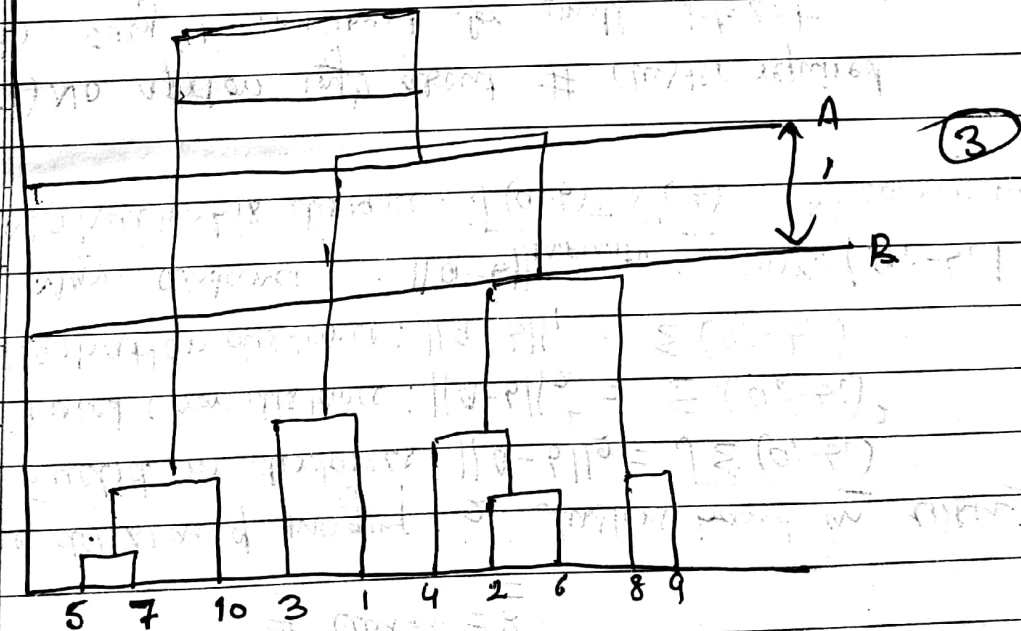
(i) Start by assigning each item to its own cluster, so if we have N items, means we have N clusters, each containing just one item. Let distance b/w the clusters is equal distance (similarity) b/w items they contain.

(ii) Find the closest (most similar) clusters and merge them into one cluster. So now we have one less cluster.

(iii) Compute distance (similarities) b/w new cluster & each of old clusters.

(iv) Repeat step-(i) & (ii) until all items are clustered into single cluster of size N .

→ Dendrogram → Tree diagram, especially in showing taxonomic relationship.



→ The height of dendrogram at which 2 clusters are merged represent distance b/w 2 clusters.

→ The best choice of # clusters is # vertical lines in the dendrogram cut by a horizontal line that can traverse the max distance vertically without intersecting a cluster.

cluster = 3

→ The decision of merging 2 clusters based on either

- i) Euclidean distance: $\|a-b\|_2 = \sqrt{\sum (a_i - b_i)^2}$
- ii) Squared error distance: $\|a-b\|_2^2 = \sum (a_i - b_i)^2$
- iii) Manhattan distance: $\|a-b\|_1 = \sum (a_i - b_i)$
- iv) Max. distance: $\|a-b\|_{\infty} = \max_i (a_i - b_i)$
- v) Mahalanobis distance: $\sqrt{(a-b)^T S^{-1} (a-b)}$ S : covariance matrix

Adv:-

- i) No apriori info about # cluster required
- ii) Easy to implement for small dataset

Dis:-

- i) Algo. can never undo what was done previously
 - ii) T.C = $O(n^2 \log n)$, $n = \# \text{ data point}$.
 - iii) Sometimes difficult to identify correct # clusters
 - iv) HCA can't handle big data but k-mean can
- $O(n^2)$ $O(n)$