

## (II) Hierarchical Methods

(149)

Hierarchical clustering method works by grouping data objects into a hierarchy or "tree" of clusters.

This method is useful for data summarization & visualization that represent the data in compressed way  
for eg:- organize the employees into major groups such as executives, managers & staff. Staff can be further divided into subgroups of senior officers, officers & trainees.

Hierarchical clustering methods :-

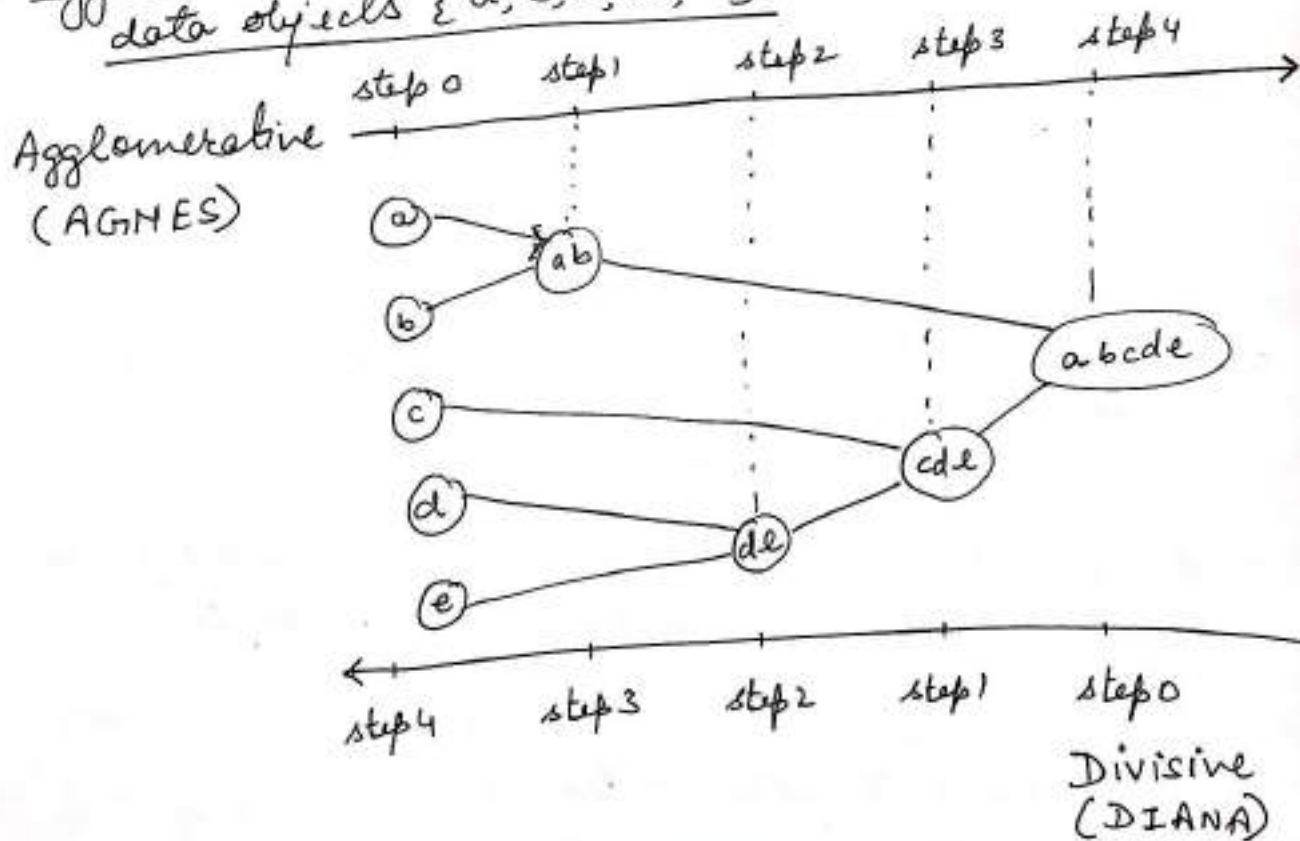
⇒ Agglomerative clustering method:-

- uses a bottom-up strategy.
- start with individual objects as clusters, which are iteratively merged to form larger clusters, until all objects are in a single cluster or certain termination conditions are satisfied.
- single cluster becomes the hierarchy root.
- for the merging step, it finds two clusters that are closest to each other (acc. to some similarity measure) & combines the two to form one cluster.
- Because the two clusters are merged per iteration, where each cluster contain at least one object.
- The method requires at most  $n$  iterations.
- Application:- AGNES (Agglomerative NESTing)

Q) Divisive hierarchical clustering method :-

- It employs a Top down strategy.
- starts by placing all objects in one cluster, which is the hierarchy root. It then divides the root cluster into several smaller subclusters & successively partitions those clusters into smaller ones.
- The partitioning process continues until each cluster at the lowest level is coherent enough - either containing only one object, or the objects within a cluster are sufficiently similar to each other.
- Application :- DIANA (Divisive ANAlysis)

⇒ Agglomerative & Divisive hierarchical clustering on data objects {a, b, c, d, e}.



### AGNES (Agglomerative method)

- Initially, AGNES, places each object into a cluster of its own.
- The clusters are then merged step-by-step acc. to some criterion.
- For eg:- clusters  $C_1$  &  $C_2$  may be merged if an object in  $C_1$  & an object in  $C_2$  form the minimum Euclidean distance between any two objects from different clusters.
- This is a single-linkage approach. In that each cluster is represented by all the objects in the cluster, & the similarity b/w the two clusters is measured by the similarity of the closest pair of data objects belonging to different clusters. Process continues until all the objects are merged to form one cluster.

### DIANA (Divisive method)

- Proceeds in contrasting way.
- All the objects form one initial cluster.
- The cluster is split according to some principle such as maximum Euclidean distance b/w the closest neighboring objects in the cluster.
- The cluster-splitting process repeats, until each new cluster contains only a single object.

## hierarchical methods (categories).

- we can improve the clustering quality of hierarchical methods is to integrate hierarchical clustering with other clustering techniques, resulting in multiple-phase (or multiphase) clustering.

- Two such methods are :-

- Ⓐ. BIRCH (Balanced Iterative Reducing & Clustering using hierarchies) :- first apply hierarchical clustering & then other methods on leaf & non-leaf.
- Ⓑ. Chameleon :- explores dynamic modeling in hierarchical structure.

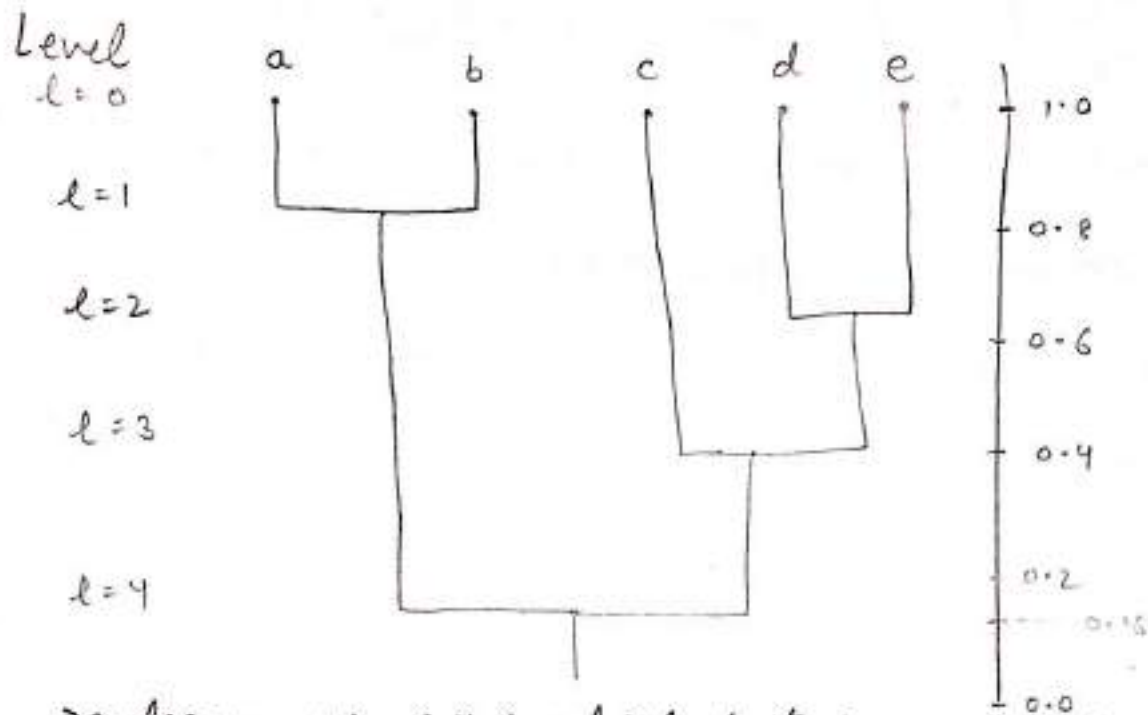
- Hierarchical methods can also be categorized as :-

- Ⓐ. Algorithmic methods (data objects are deterministic).
  - Agglomerative methods.
  - divisive methods.
  - multiphase methods.
- Ⓑ. Probabilistic methods (measures quality of clusters by fitting of models).
- Ⓒ. Bayesian methods. (computes distribution of possible clusterings).

## Example of tree structure - dendrogram

(151)

- It shows how objects are grouped together (in an agglomerative method) or partitioned (in a divisive method) step-by-step.



dendrogram rep. of Hierarchical clustering of data objects {a, b, c, d, e}

At  $l=0$  shows the five objects as singleton clusters at level 0.

At  $l=1$ , objects 'a' & 'b' are grouped together to form the first cluster & they stay together at all subsequent levels.

We use a vertical axis to show the similarity scale b/w clusters. For eg., when the similarity of two groups of objects  $\{a, b\}$  &  $\{c, d, e\}$  is roughly 0.16, they are merged together to form a single cluster.

## Disadvantages of Divisive methods

- Challenge with Divisive method is how to partition large cluster into several smaller ones.
- There are  $2^{n-1}$  possible ways to partition a set of 'n' objects into two exclusive subsets.
- when 'n' is large, it is computationally prohibitive to examine all possibilities.
- Consequently, a divisive method typically uses heuristics in partitioning, which can lead to inaccurate results.
- For the sake of efficiency, divisive methods do not backtrack on partitioning decisions that have been made.

Thus agglomerative methods are preferred over divisive methods.

## Distance Measures in Algorithmic Methods (152)

- In Algorithmic methods, a core need is to measure the distance between two clusters, where each cluster is a set of objects.
- Four widely used measures for distance between clusters are as follows where
  - $\Rightarrow |p - p'|$  is the distance b/w two objects or point
  - $\Rightarrow m_i$  is mean for cluster  $C_i$
  - $\Rightarrow n_i$  is the no. of objects in  $C_i$
- Also known as linkage measures.

① minimum distance :-

$$\text{dist}_{\min}(C_i, C_j) = \min_{p \in C_i, p' \in C_j} \{ |p - p'| \} \quad - (1)$$

② maximum distance :

$$\text{dist}_{\max}(C_i, C_j) = \max_{p \in C_i, p' \in C_j} \{ |p - p'| \} \quad - (2)$$

③ Mean distance :

$$\text{dist}_{\text{mean}}(C_i, C_j) = |m_i - m_j| \quad - (3)$$

④ Average distance :

$$\text{dist}_{\text{avg}}(C_i, C_j) = \frac{1}{n_i n_j} \sum_{\substack{p \in C_i \\ p' \in C_j}} |p - p'| \quad - (4)$$

Minimum Distance :- If  $\text{dist}_{\min}(C_i, C_j)$  measure is used, then algorithm is also called nearest-neighbor clustering algorithm.

- If the clustering process is terminated when the distance between nearest clusters exceeds a user-defined threshold, it is called single linkage algorithm.

- If we view datapoint as nodes of a graph, with edges forming a path b/w the nodes in a cluster, then the merging of two clusters  $C_i$  &  $C_j$  means adding an edge b/w the nearest pair of nodes in  $C_i$  &  $C_j$ .

- Because edges linking clusters always go between distinct clusters, the resulting graph will generally be a tree.

- An Agglomerative hierarchical clustering algorithm that uses the minimum distance measure is also called a minimal spanning tree Algorithm.



## Maximum distance

153

- when an Algorithm uses the maximum distance,  $\text{dist}_{\max}(C_i, C_j)$  to measure the distance between clusters, it is also called farthest-neighbor clustering Algorithm.

- If the clustering process is terminated when the maximum distance between nearest clusters exceeds a user-defined threshold, it is called a complete-linkage - algorithm.

- By viewing data points as a nodes of a graph, with edges linking nodes, we can think of each cluster as a complete subgraph, i.e., with edges connecting all the nodes in the clusters.

The distance b/w two clusters is determined by the most distant nodes in the two clusters.

# distance b/w two clusters can also be the distance b/w centroids.

## Mean & Average distance

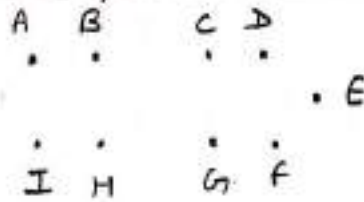
The use of mean or average distance is a compromise b/w the minimum & maximum distance & overcomes the outlier sensitivity problem.

Mean distance is simplest to compute (cannot compute mean for categorical data.)

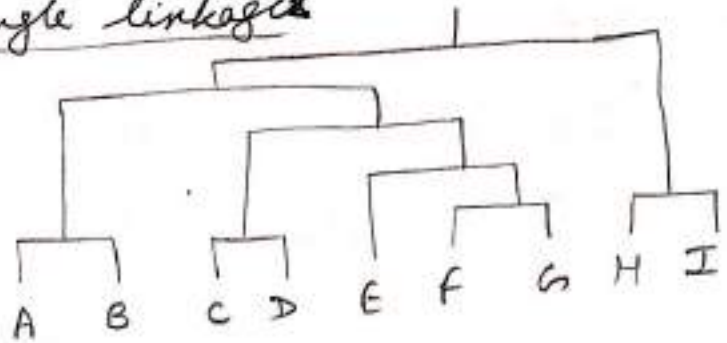
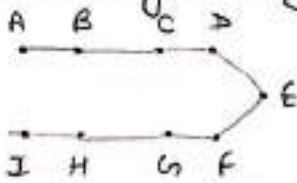
Average distance can handle categorical as well as numeric data.

# Single versus Complete linkages

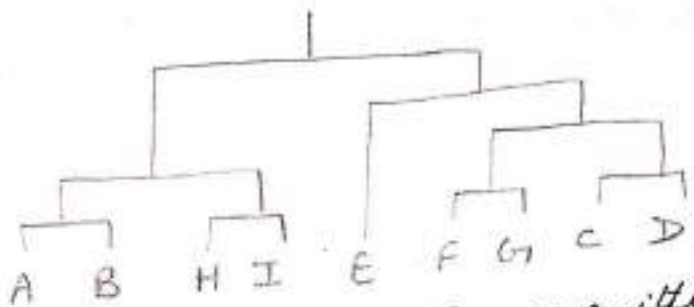
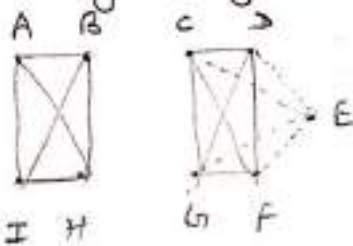
(a) Dataset



(b) clustering using single linkage



(c) clustering using complete linkage



[In Fig. (c),  
# Edges b/w  $\{A, B, I, H\}$  &  $\{C, D, G, F, E\}$  are omitted for ease of presentation.]

# using single linkages we can find hierarchical clusters defined by local proximity, whereas complete linkage tends to find clusters opting for global closeness.

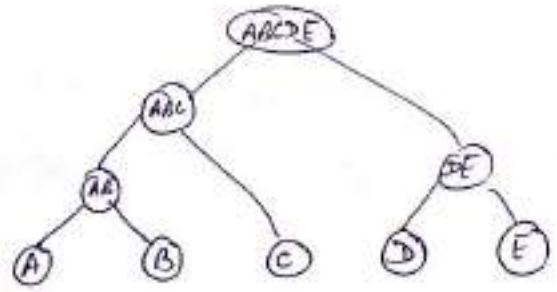
# ch-14. Hierarchical clustering

- create sequence of nested partitions which can be visualize as a tree or hierarchy of clusters, also called cluster dendrogram.
- from fine grained (lowest level) to coarse grained (highest level)
- meaningful clustering at intermediate levels.
- Two approaches
  - Agglomerative (Bottom up).
  - Divisive (Top down manner).

- Partition dataset  $D = \{x_1, x_2, \dots, x_n\}$  where  $x_i \in R^d$ , into  $k$  clusters  $C = \{C_1, C_2, \dots, C_k\}$  where  $C_i \in D$  such that  $C_i \cap C_j = \emptyset$  ( $i \neq j$ ) &  $\bigcup_{i=1}^k C_i = D$

- A clustering  $A = \{A_1, A_2, \dots, A_r\}$  is said to be nested in another clustering  $B = \{B_1, B_2, \dots, B_s\}$  iff  $r > s$ .  $\{A_i \in A, B_j \in B\}$  for each cluster  $A_i$  there exist a cluster  $B_j$  such that  $A_i \subseteq B_j$

- $D = \{A, B, C, D, E\}$
- $C = \{ABC, DE\}$
- $B = \{AB, C, DE\}$
- $A = \{A, B, C, D, E\}$
- or  $A = \{A_1, A_2, A_3, A_4, A_5\}$



Here A is nested in B because  $5 > 3$ .

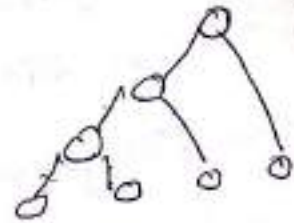
For each cluster A, there exist a cluster AB such that  $A \subseteq AB$ .

# Number of Hierarchical clustering

# means of different binary rooted trees or dendrograms with  $n$  leaves having distinct labels.

# A rooted tree with  $m$  leaves has

- $m-1$  internal nodes
- $m + m-1$  total no. of nodes.  $(2m-1)$
- $2m-2$  no. of edges.



- leaves = 4 ( $m$ )
- internal nodes = 3. ( $m-1$ )
- $4 + 3 = 7$  (total nodes)
- edges =  $2 \cdot 4 - 2 = 6$ .

# To count no. of possible dendrograms, extend leaves from  $m$  to  $m+1$ :

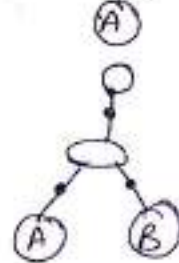
- we can add the extra leaf by splitting / branching any of  $2m-2$  edges.

Eg:-

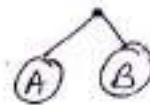
leaves  
 $m=1$



adding 1 more.  
← virtual root  
← place to split.



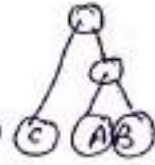
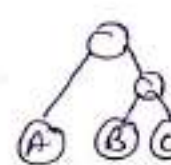
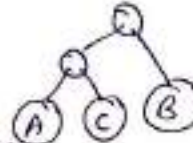
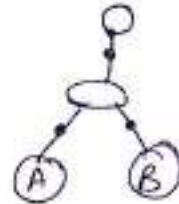
possible dendrogram.



No. of dendr.

1

$m=2$

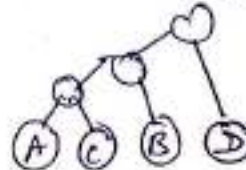
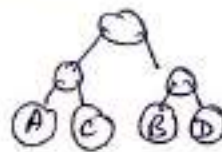
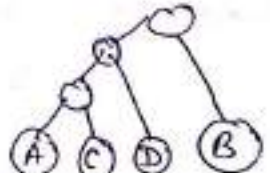
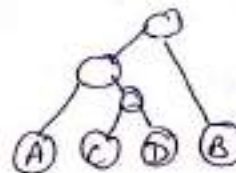
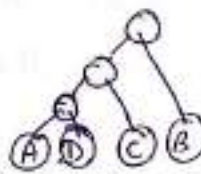
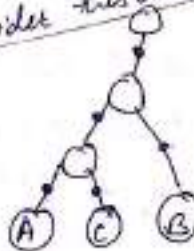


3

$m=3$



consider this one



5

# similarly for other two  $\therefore$  Total dendrograms possible are

$5 + 5 + 5 = 15$  by adding 1 leaf to existing 3 leaves

2. If we can also add new leaf as a child to the new root i.e.  $2^{m-2} + 1 = 2^{m-1}$  new dendrograms with  $m+1$  leaves.

# The total no. of different dendrograms with  $n$  leaves is given by:

$$\prod_{m=1}^{n-1} (2m-1) = 1 \times 3 \times 5 \times \dots \times (2n-3) = (2n-3)!!$$

### AGGLOMERATIVE hierarchical clustering

- Begin with each of 'n' points in a separate cluster.
- Repeatedly merge two closest clusters until all points are members of the same cluster.

- given set of clusters  $\mathcal{C} = \{C_1, \dots, C_m\}$

find closest pair of clusters  $C_i$  &  $C_j$ ,

merge them into a new cluster  $C_{ij} = C_i \cup C_j$ ,

update the set of clusters by removing  $C_i$  &  $C_j$  & adding  $C_{ij}$  as-

$$\mathcal{C} = (\mathcal{C} \setminus \{C_i, C_j\}) \cup \{C_{ij}\}$$

- Repeat it until  $\mathcal{C}$  has only one cluster.

- Agglomerative Clustering (D, k):

1.  $\mathcal{C} \leftarrow \{C_i = \{x_i\} \mid x_i \in D\}$  // each point in separate cluster.
2.  $\Delta \leftarrow \{d(x_i, x_j) : x_i, x_j \in D\}$  // compute distance matrix
3. repeat
4. find closest pair of clusters  $C_i, C_j \in \mathcal{C}$ .
5.  $C_{ij} \leftarrow C_i \cup C_j$  // merge the clusters.
6.  $\mathcal{C} \leftarrow (\mathcal{C} \setminus \{C_i, C_j\}) \cup \{C_{ij}\}$  // update the clustering
7. update distance matrix  $\Delta$  to reflect new clustering.
8. until  $|\mathcal{C}| = k$ .

## Distance between clusters

- Determine the closest pair of clusters.
- Compute distance b/w two points using Euc distance /  $L_2$ -norm.

$$s(x, y) = \|x - y\|_2 = \sqrt{\sum_{i=1}^d (x_i - y_i)^2}$$

- Distance measures

- ① Single link :- minimum distance b/w a point in  $C_i$  & a point in  $C_j$ .

$$s(C_i, C_j) = \min\{s(x, y) \mid x \in C_i, y \in C_j\}$$

- ② Complete link :- maximum distance b/w a point in  $C_i$  & a point in  $C_j$ .

$$s(C_i, C_j) = \max\{s(x, y) \mid x \in C_i, y \in C_j\}$$

- ③ Group Average :- average pairwise distance b/w points in  $C_i$  &  $C_j$ .

$$s(C_i, C_j) = \frac{\sum_{x \in C_i} \sum_{y \in C_j} s(x, y)}{n_i \cdot n_j}$$

eg. :-



$$s(C_1, C_2) = \frac{\sum_{x \in C_1} (s(x, y_1) + s(x, y_2))}{3 \cdot 2}$$

$$s(C_1, C_2) = s(x_1, y_1) + s(x_2, y_1) + s(x_3, y_1) + s(x_2, y_2) + s(x_2, y_2) + s(x_3, y_2) + s(x_1, y_2) + s(x_1, y_2)$$

where  $n_i = |C_i|$

no. of points in cluster  $C_i$ .

$$= 6 \cdot [3 \cdot 2]$$

Mean Distance :- Distance b/w the means or centroids of two clusters.

(3)

$$S(C_i, C_j) = S(\mu_i, \mu_j) \quad \text{where } \mu_i = \frac{1}{n_i} \sum_{x \in C_i} x.$$

(5) Minimum Variance :- (Ward's method) :-

Distance b/w two clusters is defined as increase in the SSE when two clusters are merged.

SSE for a given cluster  $C_i$  is given as -

$$\begin{aligned} SSE_i &= \sum_{x \in C_i} \|x - \mu_i\|^2 \\ &= \sum_{x \in C_i} (x^T \cdot x - 2x^T \cdot \mu_i + \mu_i^T \cdot \mu_i) \\ &= \sum_{x \in C_i} x^T \cdot x - 2n_i \mu_i \frac{\sum_{x \in C_i} x^T}{n_i} + \mu_i^T \sum_{x \in C_i} 1 \\ &= \sum_{x \in C_i} x^T \cdot x - 2n_i \mu_i \cdot \mu_i^T + \mu_i^T \cdot \mu_i \cdot n_i \\ &= \sum_{x \in C_i} x^T \cdot x - n_i \mu_i^T \cdot \mu_i \quad \text{[This is for single cluster 'i'] (1)} \end{aligned}$$

# The SSE for a clustering  $\mathcal{C} = \{C_1, C_2, \dots, C_m\}$  is given as

$$SSE = \sum_{i=1}^m SSE_i = \sum_{i=1}^m \sum_{x \in C_i} \|x - \mu_i\|^2$$

# Ward's measure defines the distance b/w two clusters  $C_i$  &  $C_j$  as the net change in SSE value b/w  $SSE_{ij}$  &  $(SSE_i \& SSE_j)$ .

$$\therefore S(C_i, C_j) = \Delta SSE_{ij} = SSE_{ij} - SSE_i - SSE_j \quad \text{--- (2)}$$

$$\text{Also, } C_{ij} = C_i \cup C_j \\ C_i \cap C_j = \emptyset,$$

$$\therefore |C_{ij}| = n_{ij} = |n_i| + |n_j|.$$

## updating Distance Matrix

(5)

- merge  $C_i$  &  $C_j \rightarrow C_{ij}$
- update distance of  $C_{ij}$  to all other clusters  $C_r$  ( $r \neq i, r \neq j$ ).
- Randall William eq. to recompute the distances for all of the cluster proximity measures is given as :-

$$S(C_{ij}, C_r) = \alpha_i S(C_i, C_r) + \alpha_j S(C_j, C_r) + \beta S(C_i, C_j) + \gamma |S(C_i, C_r) - S(C_j, C_r)|$$

Measure	$\alpha_i$	$\alpha_j$	$\beta$	$\gamma$
single link	$1/2$	$1/2$	0	$-1/2$
complete link	$1/2$	$1/2$	0	$1/2$
group Avg.	$\frac{n_i}{n_i + n_j}$	$\frac{n_j}{n_i + n_j}$	0	0
mean distance	$\frac{n_i}{n_i + n_j}$	$\frac{n_j}{n_i + n_j}$	$\frac{-n_i n_j}{(n_i + n_j)^2}$	0
Ward's measure	$\frac{n_i + n_r}{n_i + n_j + n_r}$	$\frac{n_j + n_r}{n_i + n_j + n_r}$	$\frac{-n_r}{n_i + n_j + n_r}$	0

## Complexity.

- $O(n^2)$  time to create pairwise distance matrix (unless specified as I/P)
- At merge step, distance from merged cluster to all other clusters is recomputed. if at step 't',  $O(n-t)$  distances are computed.
- To find closest pair in matrix, keep  $n^2$  distances in min heap. Thus  $O(1)$  time to find closest pair.
- Creation of heap takes  $O(n^2)$ .
- Deleting/updating distances from heap takes  $O(\log n)$  for each operation for a total time across all merge steps of  $O(n^2 \log n)$ .
- Computational complexity of hierarchical clustering is  $O(n^2 \log n)$ .