

Chapter-5 (Clustering) - I

* What is clustering?

- Many definitions for clusters have been proposed:
- Set of like elements. Elements from different clusters are not alike.
 - The distance between points in a cluster is less than the distance between a point in the cluster and any point outside it.

Given a database $D = \{t_1, t_2, \dots, t_n\}$ of tuples and an integer value K , the clustering problem is to define a mapping $f: D \rightarrow \{1, \dots, K\}$ where each t_i is assigned to one cluster K_j , $1 \leq j \leq K$. A cluster, K_j , contains precisely those tuples mapped to it. That is $K_j = \{t_i | f(t_i) = K_j, 1 \leq i \leq n\}$ and $t_i \in D\}$.

* What is application of clustering?

-
1. Clustering has been used in many application domains, including biology, medicine, anthropology, marketing and Economics.
 2. Clustering applications include plant and animal classification, disease classification, image processing, pattern recognition, and document retrieval.

1-(prioritously) 2-regular

* Write down the advantages of clustering. or write down the
→ basic features of clustering.

1. The (best) number of clusters is not known.
2. There may not be any a priori knowledge concerning the clusters.
3. Cluster results are dynamic.

* Write down the disadvantages of clustering.

1. Outlier handling is difficult.
2. Dynamic data in the database implies that cluster membership may change over time.
3. Interpreting the semantic meaning of each cluster may be difficult.
4. There is no correct answer to a clustering problem.
5. Clustering can be viewed as similar to unsupervised learning.

Diameter: The diameter is the square root of the average mean squared distance between all pairs of points in the cluster.

$$\text{diameter} = D_m = \sqrt{\frac{\sum_{i=1}^N \sum_{j=1}^N (t_{mi} - t_{mj})^2}{(N)(N-1)}}$$

* Discuss different kinds of method to measure the distance between clusters.



1. Single link: Smallest distance between an element in one cluster and an element in the other. we thus have

$$dis(K_i, K_j) = \min(dis(t_{il}, t_{jm})) \quad \forall t_{il} \in K_i \notin K_j \text{ and } \forall t_{jm} \in K_j \notin K_i$$

2. Complete link: largest distance between an element in one cluster and an element in the other. we thus have

$$dis(K_i, K_j) = \max(dis(t_{il}, t_{jm})) \quad \forall t_{il} \in K_i \notin K_j \text{ and } \forall t_{jm} \in K_j \notin K_i$$

3. Average: Average distance between an element in one cluster and an element in the other. we thus have

$$dis(K_i, K_j) = \text{mean}(dis(t_{il}, t_{jm})) \quad \forall t_{il} \in K_i \notin K_j \text{ and } \forall t_{jm} \in K_j \notin K_i$$

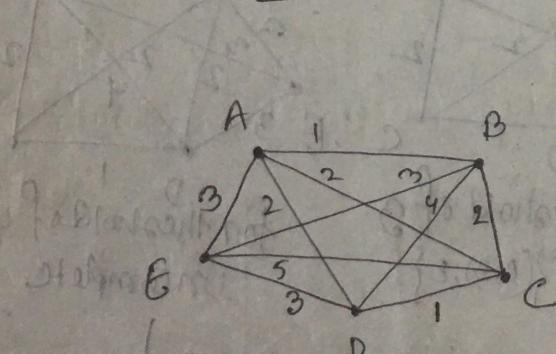
4. Centroid: If clusters have a representative centroid, then the centroid distance is defined as the distance between the centroids. we thus have $dis(K_i, K_j) = dis(c_i, c_j)$, where c_i is the centroid for K_i and similarly for c_j .

5. Medoid: Using a medoid to represent each cluster, the distance between the clusters can be defined by the distance between the medoids: $\text{dis}(k_i, k_j) = \text{dis}(m_i, m_j)$.

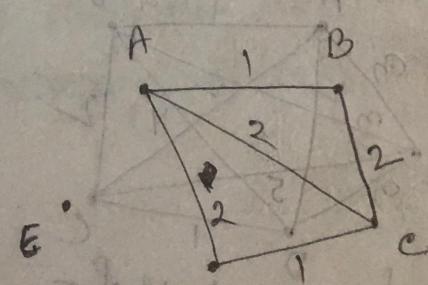
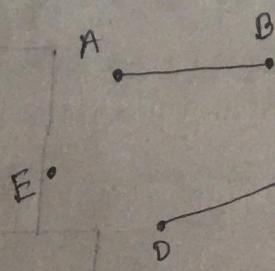
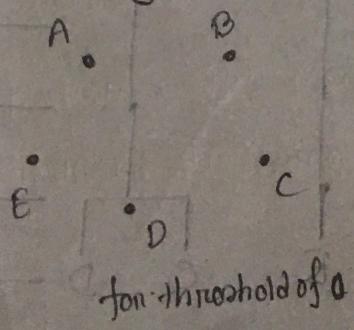
Example 5.3: for the below data draw the dendrogram for Single link, Complete link, Average link.

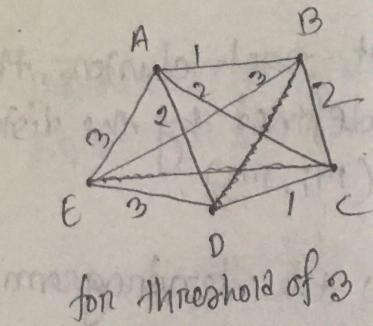
Item	A	B	C	D	E
A	0	1	2	2	3
B	1	0	2	4	3
C	2	2	0	1	5
D	2	4	1	0	3
E	3	3	5	3	0

Solution:



Single Link:-





for threshold of 3

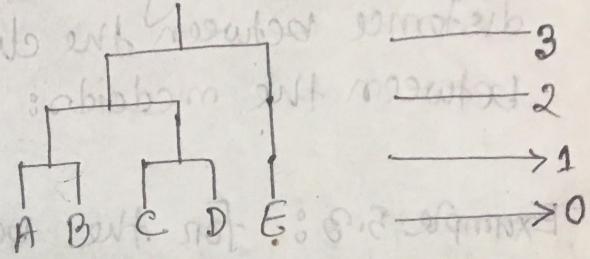
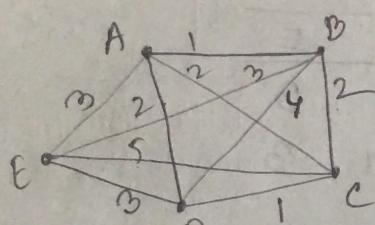
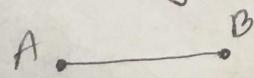


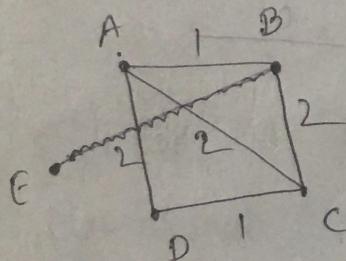
fig: dendrogram for single link.



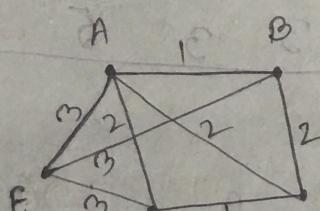
for threshold of 0
(complete)



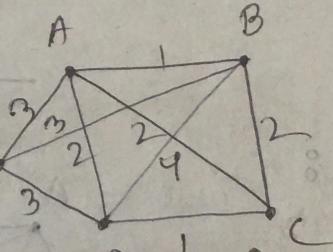
for threshold of 1
(complete)



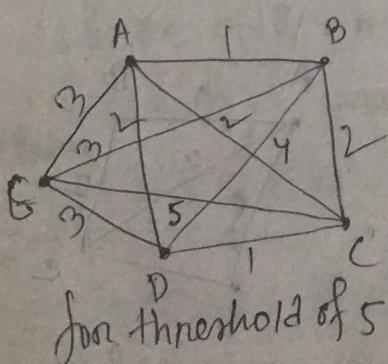
for threshold of 2
(uncomplete)



for threshold of 3
(complete)
{(A,B), E}



for threshold of 4
uncomplete



for threshold of 5

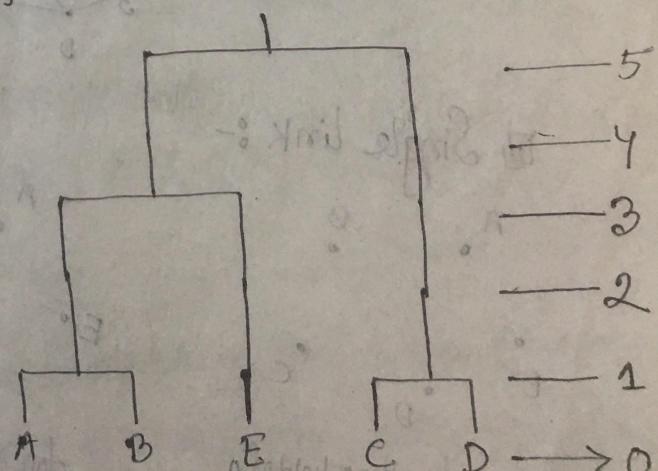
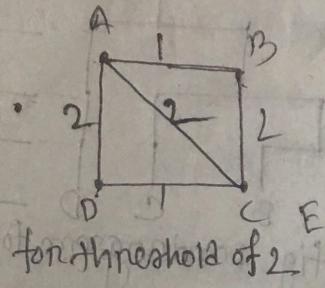
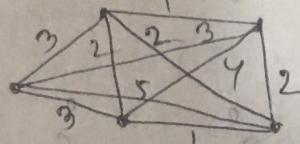


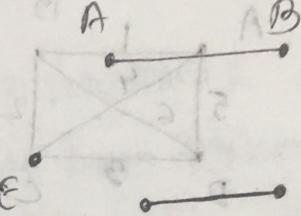
fig: Dendrogram for complete link

Average Link:-

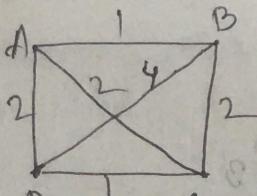


A
B
C
D
E

for threshold of 0



for threshold of 1
taken average = $0.5 \leq 1$



to complete the graph we need the path BD and the average $\frac{2+2+4+2}{4} = 2.5$

to complete the graph we need path BD and CE and the average is $\frac{3+3+5+3}{4} = 3.5$

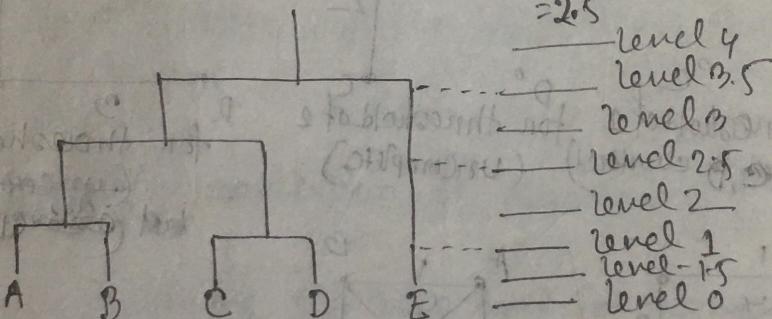


fig. Dendrogram for average link.

* show the dendrogram created by the single, complete and average link clustering algorithms using the following adjacency matrix:

Item	A	B	C	D
A	0	1	4	5
B		1	0	2
C	2	2	0	3
D	6	6	3	0

Example 5.4: if $m_1=2$, $m_2=4$, then the following data items are clustered into two (2) classes.

$$\{2, 4, 10, 12, 3, 20, 30, 11, 25\}$$

Solⁿ:

given that,

$$m_1=2$$

$$m_2=4$$

m_1	m_2	K_1	K_2
2	4	{2, 3}	{4, 10, 12, 20, 30, 11, 25}
2.5	16	{2, 3, 4}	{10, 12, 20, 30, 11, 25}
3	18	{2, 3, 4, 10}	{12, 20, 30, 11, 25}
4.75	19.6	{2, 3, 4, 10, 11, 12}	{20, 30, 25}
7	25	{2, 3, 4, 10, 11, 12}	{20, 30, 25}

$$\therefore K_1 = \{2, 3, 4, 10, 11, 12\} \text{ and } K_2 = \{20, 30, 25\}$$

বিঃ নি:- cluster two means $K=2$ (K_1, K_2). m_1 and m_2

প্ৰ. অনুঠঃ নড় কৈলো। এই মুভ অনুসৰি বেছৰ কৈ আমা প্ৰয়োগ। K_1 ক'ৰ কৈ। দেখো K_1 ও K_2 ক'ৰ মাঝে কৈলো এক ক'ৰ অধি' দূৰ্ব। অনুসৰি m_1 ও m_2 ক'ৰ দূৰ্ব, দেখো ক'ৰ অধি' দূৰ্ব। এখন ক'ৰ অধি' দূৰ্ব আৰু ক'ৰ অধি' দূৰ্ব ক'ৰ অধি' দূৰ্ব। Step-১: K_1 ও K_2 ক'ৰ element আৰু' আৰু' ১

* clustering the following data items for threshold value
 2 (using nearest neighbor Algorithm).

Item	A	B	C	D	E
A	0	1	2	2	3
B	1	0	2	4	3
C	2	2	0	1	5
D	2	4	1	0	3
E	3	3	5	3	0

Step 1:

$$\text{initially, } K_1 = \{A\}$$

now the distance between A to B is $1 \leq 2(threshold value)$

so B inserted into K_1 .

$$\text{now } K_1 = \{A, B\}$$

\therefore The distance between A to C is $= 2 \leq 2$

\therefore The distance between B to C is $= 2 \leq 2$

so C inserted into K_1

$$\therefore K_1 = \{A, B, C\}$$

\therefore The distance between A to D is $= 2 \leq 2$

\therefore The distance between B to D is $= 4 > 2$

\therefore The distance between C to D is $= 1 \leq 2$

so D inserted into K_1

$$\therefore K_1 = \{A, B, C, D\}$$

The distance between A to E is = 3

The distance between B to E is = 3

The distance between C to E is = 5

The distance between D to E is = 3

~~The distance between E~~ There is no value below of threshold value.

So E is inserted into new cluster (K_2) .

$$\therefore K_2 = \{E\}$$

$$\therefore K_1 = \{A, B, C, D\}, K_2 = \{E\}$$

* what is the usefulness or effectiveness of using large database in clustering algorithm?



A clustering algorithm should :

1. require no more than one scan of the database.
2. have the ability to provide status and "best" answer so far during the algorithm execution.
3. be suspendable, stoppable and resumable.
4. be able to update the results incrementally as data are added or removed from the database.