

# DATA MINING

COURSE: B.Sc.(H)-VI Semester

TEACHER: MS. SONAL LINDA

## Solved Examples and Exercises

### Chapter 8. Cluster Analysis

#### Solved Examples

1. For the given data, compute two clusters using K-means algorithm for clustering where initial cluster centers are (1.0, 1.0) and (5.0, 7.0). Execute for two iterations.

Record Number	A	B
R1	1.0	1.0
R2	1.5	2.0
R3	3.0	4.0
R4	5.0	7.0
R5	3.5	5.0
R6	4.5	5.0
R7	3.5	4.5

#### Solution:

##### **Algorithm 8.1** Basic K-means algorithm.

- 1: Select  $K$  points as initial centroids.
- 2: **repeat**
- 3:   Form  $K$  clusters by assigning each point to its closest centroid.
- 4:   Recompute the centroid of each cluster.
- 5: **until** Centroids do not change.

**Initialization:** Number of clusters ( $K$ ) = 2, centroid for cluster1 ( $C_1$ )= (1.0, 1.0) and centroid for cluster2 ( $C_2$ ) = (5.0, 7.0). We use Euclidean distance to find closest point to centroids.

#### **Iteration1:**

Record Number	Close to $C_1(1.0, 1.0)$	Close to $C_2(5.0, 7.0)$	Assign to cluster
R1(1.0,1.0)	dist(R1, $C_1$ )=0.0	dist(R1, $C_2$ )=7.21	Cluster1
R2(1.5,2.0)	dist(R2, $C_1$ )=1.12	dist(R2, $C_2$ )=6.12	Cluster1
R3(3.0,4.0)	dist(R3, $C_1$ )=3.61	dist(R3, $C_2$ )=3.61	Cluster1
R4(5.0,7.0)	dist(R4, $C_1$ )=7.21	dist(R4, $C_2$ )=0.0	Cluster2
R5(3.5,5.0)	dist(R5, $C_1$ )=4.12	dist(R5, $C_2$ )=2.5	Cluster2
R6(4.5,5.0)	dist(R6, $C_1$ )= 5.31	dist(R6, $C_2$ )=2.06	Cluster2
R7(3.5,4.5)	dist(R7, $C_1$ )=4.30	dist(R7, $C_2$ )=2.92	Cluster2

Thus, we obtain two clusters containing:

Cluster1 {R1, R2, R3} and Cluster2 {R4, R5, R6, R7}.

Their new centroids are:

$$\begin{aligned} C1 &= (1.0+1.5+3.0)/3, (1.0+2.0+4.0)/3 & C2 &= (5.0+3.5+4.5+3.5)/4, (7+5+5+4.5)/4 \\ &= 5.5/3, 7.0/3 & &= 16.5/4, 21.5/4 \\ &= 1.83, 2.33 & &= 4.12, 5.37 \end{aligned}$$

**Iteration2:**

Record Number	Close to C1(1.83, 2.33)	Close to C2(4.12, 5.37)	Assign to cluster
R1(1.0,1.0)	dist(R1, C1)=1.57	dist(R1, C2)=5.37	Cluster1
R2(1.5,2.0)	dist(R2, C1)=0.47	dist(R2, C2)=4.27	Cluster1
R3(3.0,4.0)	dist(R3, C1)=2.04	dist(R3, C2)=1.77	Cluster2
R4(5.0,7.0)	dist(R4, C1)=5.64	dist(R4, C2)=1.85	Cluster2
R5(3.5,5.0)	dist(R5, C1)=3.15	dist(R5, C2)=0.72	Cluster2
R6(4.5,5.0)	dist(R6, C1)=3.78	dist(R6, C2)=0.53	Cluster2
R7(3.5,4.5)	dist(R7,C1)=2.74	dist(R7, C2)=1.07	Cluster2

Therefore, new clusters are:

Cluster1 {R1, R2} and Cluster2 {R3, R4, R5, R6, R7}.

Their new centroids are:

$$\begin{aligned} C1 &= (1.0+1.5)/2, (1.0+2.0)/2 & C2 &= (3.0+5.0+3.5+4.5+3.5)/5, (4+7+5+5+4.5)/5 \\ &= 2.50/2, 3.0/2 & &= 19.5/5, 25.5/5 \\ &= 1.25, 1.5 & &= 3.9, 5.1 \end{aligned}$$

2. Use the **distance matrix** in Table1 to perform single link and complete link hierarchical clustering. Show your results by drawing a dendrogram. The dendrogram should clearly show the order in which the points are merged.

Table 1 Distance matrix

	P1	P2	P3	P4	P5
P1	0.00	0.10	0.41	0.55	0.35
P2	0.10	0.00	0.64	0.47	0.98
P3	0.41	0.64	0.00	0.44	0.85
P4	0.55	0.47	0.44	0.00	0.76
P5	0.35	0.98	0.85	0.76	0.00

### Solution:

---

**Algorithm 8.3** Basic agglomerative hierarchical clustering algorithm.

---

- 1: Compute the proximity matrix, if necessary.
  - 2: **repeat**
  - 3:   Merge the closest two clusters.
  - 4:   Update the proximity matrix to reflect the proximity between the new cluster and the original clusters.
  - 5: **until** Only one cluster remains.
- 

For the single link or MIN version of hierarchical clustering, the proximity of two clusters is defined as the minimum of the distance (maximum of the similarity) between any two points in the two different clusters.

Steps:

- Using graph terminology, start with all points as singleton clusters.
- Add links between points one at a time (shortest links first).
- These single links combine the points into clusters.

	P1	P2	P3	P4	P5
P1	0.00	0.10	0.41	0.55	0.35
P2	0.10	0.00	0.64	0.47	0.98
P3	0.41	0.64	0.00	0.44	0.85
P4	0.55	0.47	0.44	0.00	0.76
P5	0.35	0.98	0.85	0.76	0.00

Combine P1 and P2:

$$\begin{aligned} \text{dist}(\{P1, P2\}, \{P3\}) &= \min(\text{dist}(P1, P3), \text{dist}(P2, P3)) \\ &= \min(0.41, 0.64) \\ &= 0.41 \end{aligned}$$

$$\begin{aligned} \text{dist}(\{P1, P2\}, \{P4\}) &= \min(\text{dist}(P1, P4), \text{dist}(P2, P4)) \\ &= \min(0.55, 0.47) \\ &= 0.47 \end{aligned}$$

$$\begin{aligned} \text{dist}(\{P1, P2\}, \{P5\}) &= \min(\text{dist}(P1, P5), \text{dist}(P2, P5)) \\ &= \min(0.35, 0.98) \\ &= 0.35 \end{aligned}$$

	P12	P3	P4	P5
P12	0.00	0.41	0.55	0.35
P3	0.41	0.00	0.44	0.85
P4	0.55	0.44	0.00	0.76
P5	0.35	0.85	0.76	0.00

Combine P12 and P5:

$$\begin{aligned} \text{dist}(\{P12, P5\}, \{P3\}) &= \min(\text{dist}(P12, P3), \text{dist}(P5, P3)) \\ &= \min(0.41, 0.85) \\ &= 0.41 \end{aligned}$$

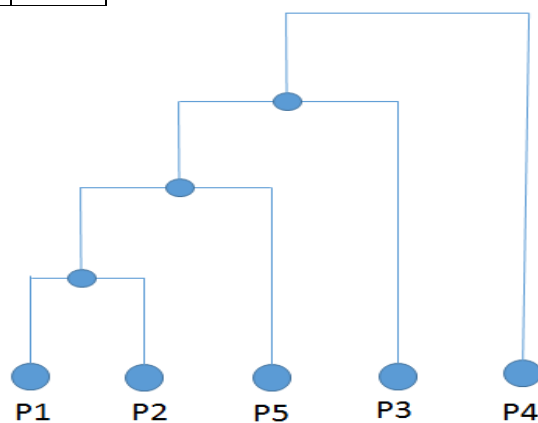
$$\begin{aligned} \text{dist}(\{P12, P5\}, \{P4\}) &= \min(\text{dist}(P12, P4), \text{dist}(P5, P4)) \\ &= \min(0.55, 0.76) \\ &= 0.55 \end{aligned}$$

	P125	P3	P4
P125	0.00	0.41	0.55
P3	0.41	0.00	0.44
P4	0.55	0.44	0.00

Combine P125 and P3:

$$\begin{aligned} \text{dist}(\{P125, P3\}, \{P4\}) &= \min(\text{dist}(P125, P4), \text{dist}(P3, P4)) \\ &= \min(0.55, 0.44) \\ &= 0.44 \end{aligned}$$

	P1235	P4
P1235	0.00	0.44
P4	0.44	0.00



Single Link Dendrogram

For the complete link or MAX version of hierarchical clustering, the proximity of two clusters is defined as the maximum of the distance (minimum of the similarity) between any two points in the two different clusters.

Steps:

- Using graph terminology, start with all points as singleton clusters.
- Add links between points one at a time (shortest links first).

- Group points until all the points are completely linked, i.e., clique.

-	P1	P2	P3	P4	P5
P1	0.00	0.10	0.41	0.55	0.35
P2	0.10	0.00	0.64	0.47	0.98
P3	0.41	0.64	0.00	0.44	0.85
P4	0.55	0.47	0.44	0.00	0.76
P5	0.35	0.98	0.85	0.76	0.00

Combine P1 and P2:

$$\begin{aligned} \text{dist}(\{P1, P2\}, \{P3\}) &= \max(\text{dist}(P1, P3), \text{dist}(P2, P3)) \\ &= \max(0.41, 0.64) \\ &= 0.64 \end{aligned}$$

$$\begin{aligned} \text{dist}(\{P1, P2\}, \{P4\}) &= \min(\text{dist}(P1, P4), \text{dist}(P2, P4)) \\ &= \min(0.55, 0.47) \\ &= 0.47 \end{aligned}$$

$$\begin{aligned} \text{dist}(\{P1, P2\}, \{P5\}) &= \min(\text{dist}(P1, P5), \text{dist}(P2, P5)) \\ &= \min(0.35, 0.98) \\ &= 0.35 \end{aligned}$$

	P12	P3	P4	P5
P12	0.00	0.64	0.98	0.98
P3	0.64	0.00	0.44	0.85
P4	0.98	0.44	0.00	0.76
P5	0.98	0.85	0.76	0.00

Combine P3 and P4:

$$\begin{aligned} \text{dist}(\{P3, P4\}, \{P12\}) &= \max(\text{dist}(P3, P12), \text{dist}(P4, P12)) \\ &= \max(0.64, 0.98) \\ &= 0.98 \end{aligned}$$

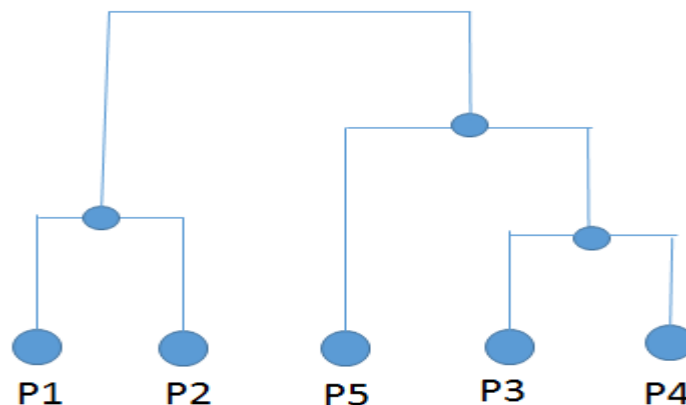
$$\begin{aligned} \text{dist}(\{P3, P4\}, \{P5\}) &= \min(\text{dist}(P3, P5), \text{dist}(P4, P5)) \\ &= \min(0.85, 0.76) \\ &= 0.76 \end{aligned}$$

	P12	P34	P5
P12	0.00	0.98	0.98
P34	0.98	0.00	0.85
P5	0.98	0.85	0.00

Combine P34 and P5:

$$\begin{aligned} \text{dist}(\{P34, P5\}, \{P12\}) &= \max(\text{dist}(P34, P12), \text{dist}(P5, P12)) \\ &= \max(0.98, 0.98) \\ &= 0.98 \end{aligned}$$

	P12	P345
P12	0.00	0.98
P345	0.98	0.00



Complete Link Dendrogram

### Exercises

1. Find all well separated clusters in the set of points shown in Figure 1.



Figure 1: Points

2. Many partitional clustering algorithms that automatically determine the number of clusters claim that this is an advantage. List two situations in which this is not the case.
3. Identify the clusters in Figure 2 using the center-, contiguity-, and density-based definitions. Also indicate the number of clusters for each case and give a brief indication of your reasoning. Note that darkness or the number of dots indicates density. If it helps, assume center-based means k-means, contiguity-based single link, and density-based means DBSCAN.

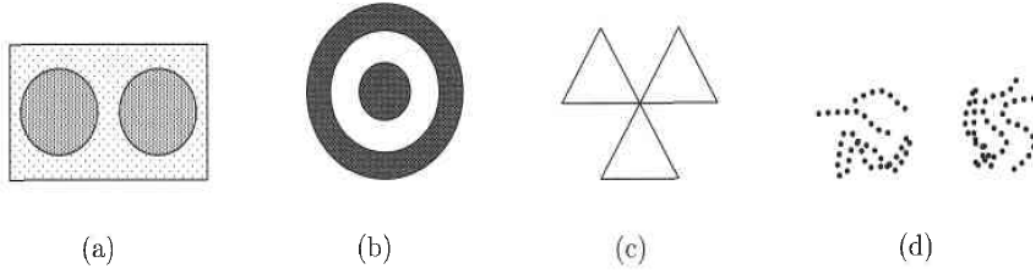


Figure 2: Clusters

4. Given the following points: 2, 4, 10, 12, 3, 20, 30, 11, 25. Given  $k = 3$ , and the initial means,  $\mu_1 = 2, \mu_2 = 4$  and  $\mu_3 = 6$ . Show the clusters obtained and new means after each iteration using the K-means algorithm.
5. Use the **distance matrix** in Table 2 to perform single link and complete link hierarchical clustering. Show your results by drawing a dendrogram. The dendrogram should clearly show the order in which the points are merged.

Table 2: Distance Matrix

	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00

6. Use the **similarity matrix** in Table 3 to perform single link and complete link hierarchical clustering. Show your results by drawing a dendrogram. The dendrogram should clearly show the order in which the points are merged.

Table 3: Similarity matrix

	p1	p2	p3	p4	p5
p1	1.00	0.10	0.41	0.55	0.35
p2	0.10	1.00	0.64	0.47	0.98
p3	0.41	0.64	1.00	0.44	0.85
p4	0.55	0.47	0.44	1.00	0.76
p5	0.35	0.98	0.85	0.76	1.00

7. Consider the mean of a cluster of objects from a binary transaction data set. What are the minimum and maximum values of the components of the mean? What is the interpretation of components of the cluster mean? Which component most accurately characterizes the objects in the cluster?
8. Differentiate between agglomerative and divisive methods of hierarchical clustering with the help of a diagram.
9. Explain the following terms with reference to the DBSCAN clustering algorithm:
  - (a) Core points
  - (b) Noise points
  - (c) Border points
10. Describe the following clustering algorithm in terms of:
  - (a) Shape of clusters
  - (b) Limitations:
    - i. K-means
    - ii. DBSCAN