

## Types of clustering

### 1) Partitioning methods -

- (a) Given a database of  $n$  objects or data tuples a partitioning method constructs  $K$  partitions of data, where each partition represents a cluster and  $K \leq n$ .
- (b) Given  $K$ , the number of partitions to construct it creates an initial partitioning.
- (c) It then uses an iterative relocation technique that attempts to improve the partitioning by moving objects from one group to another.
- (d) The criteria of a good partitioning is that objects in the same cluster are close to each other, whereas objects of different clusters are far apart.
- (e) Ex - K-mean, K-medoid

### 2) Hierarchical methods

- (a) This method creates a hierarchical decomposition of the given set of data objects. It can be either agglomerative or divisive.
- (b) The agglomerative (bottom-up) approach starts with each object forming a separate group. It successively merges the objects that are close to one another, until all of the groups are merged into one.

- (c) The divisive (top-down) approach starts with all of the objects in the same cluster. On each successive iteration, a cluster is split into small clusters, until each object forms its own cluster or until a termination criterion meets.
- Termination criterion : BIRCH
- Ex - AGNES, DIANA,

- (3) Density based clustering
- (a) This method based on the notion of density.
- (b) For each data point within a given cluster, the neighbourhood of a given radius have to contain at least a minimum no. of points.
- (c) This method can be used to filter out noise and discover clusters of arbitrary shape.
- Ex - DBSCAN, OPTICS

- (4) Grid based methods
- (a) These methods quantize the object space into a finite no. of cells that form a grid structure.
- (b) All of the clustering operations are performed on the grid.
- (c) Advantage is fast processing time.
- Ex - STING

### (5) Model-based methods

- (a) This approach hypothesizes a model for each of the clusters and finds the best fit of the data to the given model.
  - (b) It also leads to a way of automatically determining the no. of clusters based on standard statistics.
  - (c) It takes noise or outliers into account, thereby contributing to the robustness of the approach.
- Ex - COBWEB, self-organizing feature maps

### (6) Methods for high dimensional data

- (a) High dimensional data can typically have many irrelevant dimensions.
  - (b) Subspace clustering methods - which search for cluster in subspaces of the data, rather than over the entire data ~~objects~~.
  - (c) Frequent pattern based clustering - It extracts distinct frequent patterns among subsets of dimensions that occur frequently.
- Ex - CLIQUE, PROCLUS

### (7) constraint based methods

- (a) These methods perform clustering by incorporating user-specified or application oriented constraint.

## K-means clustering

- Q. Suppose that the data mining is to cluster the following eight points onto 3-clusters.
- $A_1(2,10), A_2(2,5), A_3(8,4), B_1(5,8), B_2(7,5)$

$B_3(6,4), C_1(1,2), C_2(4,9)$  using Euclidean distance.

The distance function is Euclidean distance. Suppose initially we assign  $A_1, B_1$  and  $C_1$  as the center of each cluster, respectively. Use the K-means algorithm to show only

- the 3-cluster centers after the first round of execution

- The final three clusters.

Ans  $A_1(2,10)$

$$A_1 A_1 = \sqrt{(2-2)^2 + (10-10)^2} = 0$$

$$A_2 A_1 = \sqrt{(2-2)^2 + (5-10)^2} = \sqrt{25} = 5$$

$$A_3 A_1 = \sqrt{(8-2)^2 + (4-10)^2} = \sqrt{36+36} = \sqrt{72} = 8.485$$

$$B_2 A_1 = \sqrt{(7-2)^2 + (5-10)^2} = \sqrt{50} = 7.0710$$

$$B_3 A_1 = \sqrt{(6-2)^2 + (4-10)^2} = \sqrt{40} = 6.324$$

$$C_2 A_1 = \sqrt{(4-2)^2 + (9-10)^2} = \sqrt{5} = 2.236$$

Clusters  
 $\{A_1\}$

$(2,10)$

$\{B_1, A_3, (5,8), (8,4), (7,5), (6,4), (4,9)\}$

$(6,6)$

$B_1(5,8)$

$$A_2 B_1 = \sqrt{(2-5)^2 + (5-8)^2} = \sqrt{9+9} = 4.2426$$

$$A_3 B_1 = \sqrt{(8-5)^2 + (4-8)^2} = \sqrt{9+16} = \sqrt{25} = 5$$

$$B_2 B_1 = \sqrt{(7-5)^2 + (5-8)^2} = \sqrt{4+9} = 3.605$$

$$B_3 B_1 = \sqrt{(6-5)^2 + (4-8)^2} = \sqrt{1+16} = \sqrt{17} = 4.1231$$

$$C_2 B_1 = \sqrt{(4-5)^2 + (9-8)^2} = \sqrt{1+1} = \sqrt{2} = 1.414$$

$$C_1 B_1 = \sqrt{(1-5)^2 + (2-8)^2} = \sqrt{16+36} = \sqrt{52} = 7.2111$$

$$C_2 C_1 = \sqrt{(4-1)^2 + (9-2)^2} = \sqrt{9+49} = \sqrt{58} = 7.615$$

$$\{C_1, A_2\}$$

$$(1, 2)$$

$$(2, 5)$$

$$(5, 8) \rightarrow$$

$$(1.5, 3.5) \rightarrow$$

$$\text{New cluster head.}$$

$$(6, 6)$$

$$\{C_1, A_2\}$$

(C<sub>1</sub>)

$$\underline{A_1(2, 10)}$$

$$A_2 A_1 = 5$$

(C<sub>2</sub>)

$$\underline{B_1(5, 8)}$$

$$A_2 B_1 = 4.2428$$

(C<sub>3</sub>)

$$\underline{C_1(1, 2)}$$

$$A_2 C_1 = 3.1622$$

$$A_3 A_1 = 8.485$$

$$A_3 B_1 = 5$$

$$A_3 C_1 = 7.28$$

$$B_2 A_1 = 7.0710$$

$$B_2 B_1 = 3.605$$

$$B_2 C_1 = 6.708$$

$$B_3 A_1 = 6.324$$

$$B_3 B_1 = 4.1231$$

$$B_3 C_1 = 5.385$$

$$C_2 A_1 = 2.236$$

$$C_2 B_1 = 1.414$$

$$C_2 C_1 = 7.615$$

Clusters

$$\{A_1\}, \{B_1, A_3, B_2, B_3, C_2\} \quad \{C_1, A_2\}$$

update cluster head

$$(2, 10)$$

$$(5, 8)$$

$$(1, 2)$$

$$(8, 4)$$

$$(2, 5)$$

$$(7, 5)$$

$$(1.5, 3.5)$$

$$(6, 4)$$

$$(4, 9)$$

$$(6, 6)$$

Repeat the  
new cluster

same process from there  
already.

K-mean clustering for IRIS data.

IRIS data -  $(150 \times 4)$  (PL, PW, SL, SW)

contains - 3 classes of data.

1. Load the data.
2. choose randomly three cluster heads from 150 tuples.
3. calculate the Euclidean distance of 150 tuples from these cluster head and assign the tuple to that cluster from which the distance is minimum.
4. update the cluster heads.
5. Repeat the steps 2 and 3 with the new cluster heads.
6. Stopping criteria - until the no. of data objects does not change in each cluster.

## Comments on K-means

Strength Relatively efficient

$O(tkn)$ ,  $k, t \ll n$

$n$  = objects.

$k$  = no. of clusters.

$t$  = no. of iterations.

## Weakness

- 1) Terminating at a local solution.
- 2) Need to specify  $k$ , the no. of clusters in advance.
- 3) Unable to handle noisy and outliers.
- 4) Not suitable to discover clusters with non-convex shapes.

## K-Medoid (partitioning around Medoids) (PAM), CLARA, CLARANS

1. Initialize! select  $K$  of the  $n$  data points as the medoids
2. Associate each data point to the closest medoid.
3. While the cost of the configuration decreases
  - (a) For each medoid  $m$ , for each non-medoid data point,  $o$ 
    - > swap  $m$  and  $o$ , recompute the cost
    - > if the total cost of the configuration increased by the previous step, undo the step

Ex cluster the following data set of 10 objects into two clusters i.e  $K=2$

$x_1 2, 6$

$x_2 3, 4$

$x_3 3, 8$

$x_4 4, 7$

$x_5 6, 2$

$x_6 6, 4$

$x_7 7, 3$

$x_8 7, 4$

$x_9 8, 5$

$x_{10} 7, 6$

Step 1

Two observations  $c_1 = x_2 = (3, 4)$  and  $c_2 = x_8 = (7, 4)$  are randomly selected as medoids (cluster centers).

Step 2

Manhattan distances are calculated to each center to associate each data object to its nearest medoid.

$$d(c_{ij}) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{in} - x_{jn}|$$

PAM is more robust than K-mean.

Work well for small datasets.

$x_i$	Distance to	
	$c_1(3,4)$	$c_2(7,4)$
(2,6)	3✓	7
(3,4)	0✓	4
(3,8)	4✓	8
(4,7)	4✓	6
(6,2)	5	3✓
(6,4)	3	1✓
(7,3)	5	1✓
(7,4)	4	0✓
(8,5)	6	2✓
(7,6)	6	2✓
Cost	11	9

Since the points (2,6), (3,8) and (4,7) are closer to  $c_1$  hence they form one cluster while the remaining points form another cluster.

$$\text{Cluster 1} = \{(3,4), (2,6), (3,8), (4,7)\}$$

$$\text{Cluster 2} = \{(7,3), (6,2), (6,4), (7,4), (8,5), (7,6)\}$$

The total cost of this clustering is - the sum of the distance between a data point and its cluster center =

$$3 + 0 + 4 + 4 + 5 + 1 + 1 + 0 + 2 + 2 = 20.$$

Step 3

Select one of the non-medoids

$$\text{Let } x_7 = (7,3) = O'$$

so now medoids are  $c_1(3,4)$  and  $O'(7,3)$

calculate the total cost involved.

$x_i$	$c_1(3,4)$	$O'(7,3)$	Total Cost
$x_1$	3	8	3 + 8 = 11
$x_2$	4	9	4 + 9 = 13
$x_3$	4	7	4 + 7 = 11
$x_4$	4	2	4 + 2 = 6
$x_5$	5	2	5 + 2 = 7
$x_6$	3	8	3 + 8 = 11
$x_7$	4	1	4 + 1 = 5
$x_8$	4	3	4 + 3 = 7
$x_9$	6	3	6 + 3 = 9
$x_{10}$	6	9	6 + 9 = 15

so the cost of swapping medoid from  $c_2$  to  $O'$  is 22.

$s = 22 - 20 = 2 > 0$   
so moving  $O'$  would be a bad idea.  
so moving other non-medoidal points to get minimum distance.