

## Clustering

clustering is the process of grouping a set of data objects into multiple groups or clusters so that objects within in a cluster have highly similarity, but are very dissimilar to objects in other clusters. Dissimilarities are similarities are accessed based on the attribute values describing the objects and often involve distance measures. Clustering as a data mining tool has its roots in many application areas such as biology, security, business intelligence, and web search.

You will learn several basic clustering techniques organizing into the following categories:

- ① Partitioning Methods
- ② Hierarchical Methods
- ③ Density based Methods
- ④ Grid based Methods
- ⑤ Model based clustering Methods
- ⑥ Clustering of high dimensional data

Cluster Analysis:- cluster analysis or simply clustering is the process of partitioning a set of data objects (or observations) into subsets. Each subset is a cluster such that objects in a cluster are similar to one another yet dissimilar to object in other clusters. The set of clusters resulting from a cluster analysis can be referred to as clustering. In this context, different clustering methods may generate different clustering on the same data set. The partitioning is not performed by humans but by clustering algorithms.

Hence, clustering is useful in that can lead to the discovery of previously unknown groups within the data.

Requirements for cluster Analysis : → The following are the typical requirements of clustering in data mining

- \* Scalability
- \* Ability to deal with different types of attribute
- \* Requirements for domain knowledge to determine I/P parameters
- \* Ability to deal with noisy data
- \* Incremental clustering and Insensitivity to I/P order
- \* Capability of clustering high dimensional data
- \* Constraint based clustering
- \* Interpretability and usability

The following are orthogonal aspects with which clustering methods can be compared:

- \* The partitioning criteria
- \* Separation of clusters
- \* Similarity Measure
- \* Clustering space:-

Overview of Basic clustering Methods :-

Method	General characteristics
② Partitioning Methods	<ul style="list-style-type: none"><li>- find mutually exclusive clusters of spherical shape</li><li>- Distance based</li><li>- May use mean or median to represent</li><li>- Effective for small to medium data set</li></ul>

## Hierarchical Method

- Clustering is a hierarchical decomposition (i.e multiple levels)
- Cannot correct erroneous merges or splits
- May incorporate other technique like microclustering or consider object "linkages"

## Density based Method

- Can find arbitrarily shaped clusters
- Clusters are dense regions of objects in space that are separated by low density regions.
- Cluster density, each point must have a minimum number of points within its neighbourhood.
- May filter out outliers

## Grid based Methods

- use a multi-resolution grid based structure
- fast processing time (typically independent of the number of data objects, yet dependent on grid size)

Partitioning Methods: - Given a set of  $n$  objects, a Partitioning Method constructs  $K$  partitions of the data, where each partition represents a cluster and  $K \leq n$ , that is, it divides the data into  $K$  groups such that each group must contain at least one object.

The clusters are formed to optimize one objective partitioning criterion such as a dissimilarity function based on distance, so that the objects within a cluster are "similar" to one another and "dissimilar" to object in other clusters in terms of data set attributes.

The most well known and commonly used partitioning methods are (i) K-Means (ii) K-Medoids

K-Means: A centroid based technique : - Suppose a dataset  $D$ , contains  $n$  objects in  $D$  into  $K$ -clusters,  $C_1, C_2 \dots C_K$  that is  $c_i \in D$  and  $c_i \cap c_j = \emptyset$  for  $(i \leq k, j \leq k)$ . An objective function is used to assess the partitioning quality so that objects within a cluster are similar to one another but dissimilar to objects in other clusters. This is, the objective function aims for high intracluster similarity and low intercluster similarity.

A centroid based partitioning Method uses the centroid of a cluster,  $c_i$  to represent that cluster. Conceptually, the centroid of a cluster is a central point. The centroid can be defined in various ways such as by the Mean or the Medoid of the objects (or points) assigned to the cluster.

How does the K-Means algorithm work : - The K-Means algorithm defines the centroid of a cluster as the mean value of the points within the cluster. It proceeds as follows. first, it randomly selects  $K$  of the objects in  $D$ , each of which initially represents a cluster mean or center. for each of the remaining objects, an object is assigned to the cluster to which it is the most similar, based on the Euclidean distance b/w the object and the cluster mean. The K-Means algorithm then iteratively improves the within-cluster variation. for each cluster, it computes the new mean using the objects assigned to the cluster in the previous iteration. All the objects are then reassigned using the updated means as the new cluster centers. The iterations continue until the assignment is stable, that is, the clusters formed in the current round are the same as those formed in the previous round.

10.2

figure: The K-Means Partitioning algorithm:-

Algorithm: K-Means. The K-Means algorithm for Partitioning, where each cluster's center is represented by the mean value of the objects in the cluster.

Input:

- \*  $K$ : The number of clusters,
- \*  $D$ : a data set containing  $n$  objects

Output: A set of  $K$  clusters

Method:

- ① Arbitrarily choose  $K$  objects from  $D$  as the initial cluster centers;
- ② repeat
- ③ (Re) assign each object to the cluster to which the object is the most similar, based on the mean value of the objects in the cluster;
- ④ update the cluster means, that is, calculate the mean value of the objects for each cluster;
- ⑤ until no change;

Example:— Consider a set of objects located in 2-D space, as depicted in figure 10.3 (a). Let  $K=3$ , i.e., the user wanted like the objects to be partitioned into three clusters.

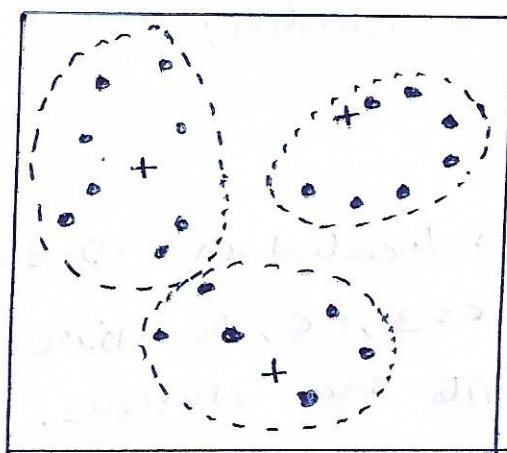
According to the Algorithm in fig 10.2, we arbitrarily choose three objects as the three initial cluster centers, where cluster centers are marked by a +. Each object is assigned to a cluster based on the cluster center to which it is nearest. Such a distribution forms silhouettes ex.

dotted curves, as shown in figure 10.3 (a).

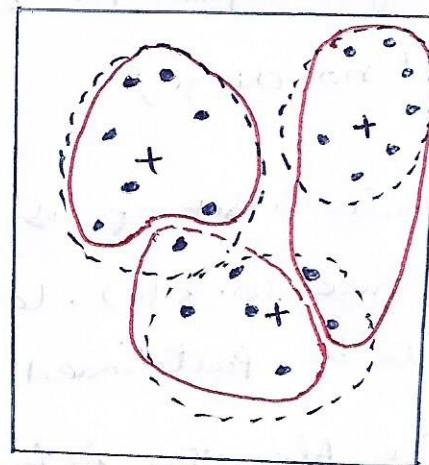
Next the cluster centers are updated. i.e., The Mean Value of each cluster is recalculated based on the current objects in the cluster. Using the new cluster centers, the objects are re-distributed to the clusters based on which cluster center is nearest. Such a redistribution forms new Silhouettes encircled by dashed curves, as shown in figure 10.2 (b).

This process iterates, leading to figure 10.3 (c). The process of iteratively reassigning objects to clusters to improve the partitioning is referred to as iterative relocation. Eventually, no ~~more~~ reassignment of the objects in any cluster occurs and so the process terminates.

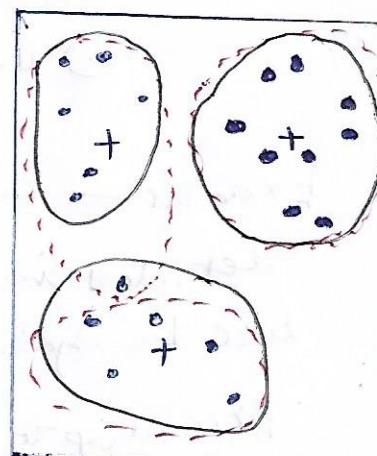
The resulting clusters are returned by the clustering process. The time complexity of the k-means also is  $O(nkt)$ , where  $n$  is the total number of objects,  $k$  is the number of clusters, and  $t$  is the number of iterations. Normally  $k \leq n$  and  $t \leq n$ .



(a) Initial clustering



(b) Iterate



(c) Final clustering

Fig:- Clustering of a set of objects using the k-means method for update cluster centers and reassign objects accordingly.

## K-Means clustering Algorithm :-

one of the most frequently used unsupervised algorithms is K-Means. K-Means cluster is exploratory data analysis technique. This is non-hierarchical method of grouping objects together.

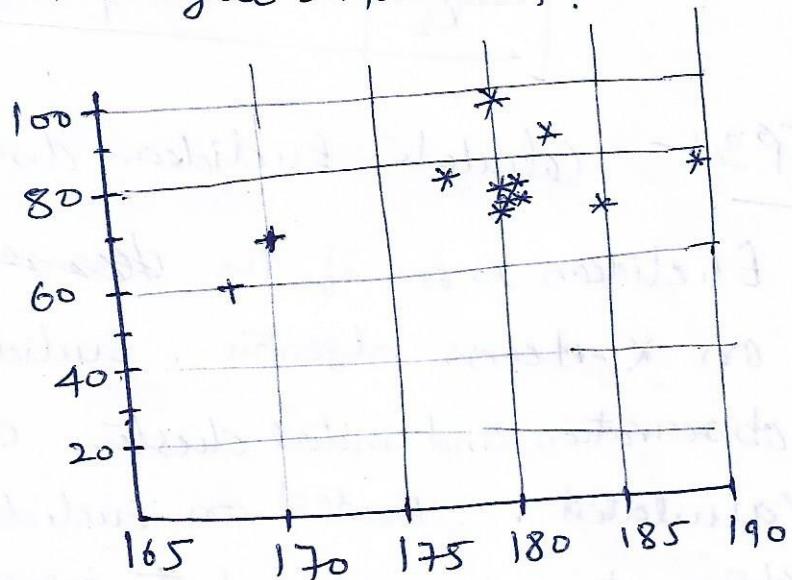
Cluster analysis or clustering is the task of grouping a set of objects in such a way that objects in the same group (called a cluster) are more similar (in some sense or another) to each other than to those in other groups (clusters).

For Ex:-

We have height and weight information, using these two variables, we need to group the objects based on height and weight information.

Data Sample

Height	Weight
185	72
170	56
168	60
179	68
182	72
188	77
180	71
180	70
183	84
180	88
180	67
177	76



If you look at the above diagram, you are able to see two visible clustered segments and we want these to be joined using K-Means algorithm.

Step 1:- <sup>Input</sup> Data set, Clustering Variables and maximum no. of clusters ( $K$  in Means clustering) for this dataset, only two variables - height and weight - are considered for clustering.  $K=2$

Step 2:- initialize cluster centroid; in this example, value of  $K$  is considered as 2. Cluster Centroids are initialized with first 2 observations

Cluster	Initial Centroid	
	Height	Weight
$K_1$	185	72
$K_2$	170	56

Step 3:- Calculate Euclidean distance

Euclidean is one of the ~~distance~~ distance measures used in K-Means algorithm. Euclidean distance b/w an observation and initial cluster centroid 1 and 2 is calculated. Based on Euclidean distance each observation is assigned to one of the clusters - based on minimum distance.

$$\text{Euclidean distance} = \sqrt{(x_H - H_1)^2 + (x_W - W_1)^2}$$

where! —  $x_H$ : Observation value of variable Height

$H_1$ : Centroid value of cluster 1 for variable Height

$x_W$ : Observation Value of variable Weight

$W_1$ : Centroid value of cluster 1 for variable weight

$x_W$ : Observation value of variable Weight

First Two Observations

Height	Weight
185	72
170	56

Step 3: calculate

Now Initial Cluster Centroid are:-

Cluster	Updated Centroid	
	Height	Weight
$k_1$	185	72
$k_2$	170	56

Euclidean distance calculation from each of the clusters is calculated

Euclidean distance from Cluster 1	Euclidean distance from Cluster 2	Assignment
$(185-185)^2 + (72-72)^2$ = 0	$(185-170)^2 + (72-56)^2$ = 21.93	1
$(170-185)^2 + (56-72)^2$ = 21.93	$(170-170)^2 + (56-56)^2$ = 0	2

We have considered two observations for assignment only because we knew the assignment, and there is no change in centroids as these two observations were only considered as initial centroids.

Step 4:- Move on to next observation and calculate Euclidean Distance.

height	weight
168	60

Euclidean distance from cluster 1	Euclidean distance from cluster 2	Assignment
$(168-185)^2 + (60-72)^2 \Rightarrow 20.808$	$(170-168)^2 + (56-60)^2 \Rightarrow 4.472$	2

Since distance is Minimum from cluster 2, so the observation is assigned to cluster 2. Now revise Cluster Centroid - Mean value height and weight as cluster centroids. Addition is only to cluster 2, so centroid of cluster 2 will be updated

updated cluster centroids

cluster	updated centroid	
	height	weight
k=1	185	72
k=2	$(170+168)/2 \Rightarrow 169$	$(56+60)/2 \Rightarrow 58$

Step 5:- calculate Euclidean distance for the next observation, assign next observation based on minimum Euclidean distance and update the cluster centroids  
Next observation

height	weight
179	68

# Euclidean Distance Calculation and assignment

Euclidean distance from cluster 1	Euclidean distance from cluster 2	Assignment
70.211	14.142	1

update cluster centroid

Cluster	update centroid	
	height	weight
K=1	182	70.667
K=2	169	59

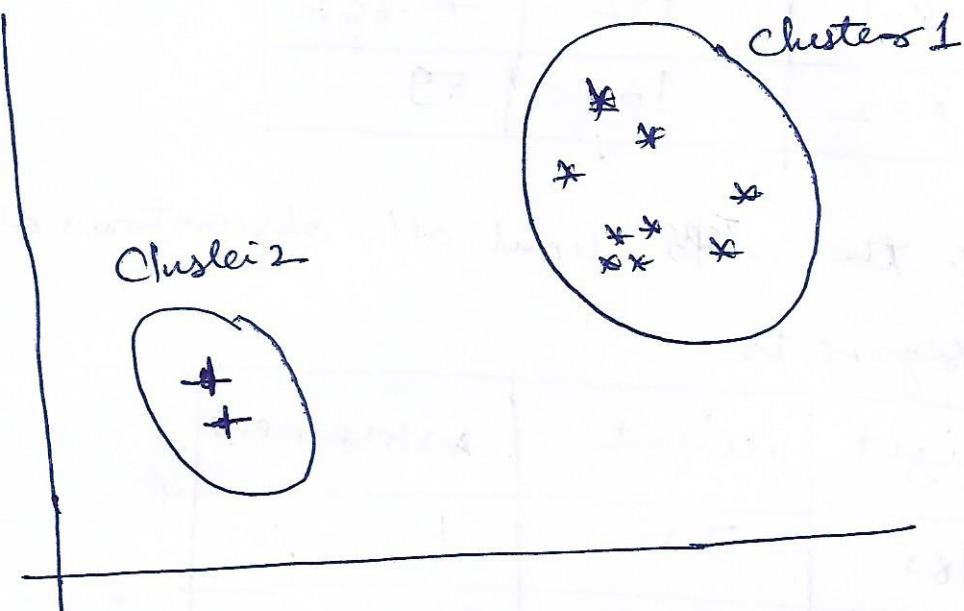
Continue the steps until all observations are assigned  
final assignments

Height	Weight	Assignment
185	72	1
170	56	2
168	60	2
139	68	1
182	72	1
188	77	1
180	71	1
180	70	1
183	84	1
180	88	1
180	67	1
170	76	1

## final cluster centroids

cluster	updated centroid	
	height	weight
K=1	182.8	72
K=2	169	58

that is what was expected based on two dimensional plot.



Important considerations in K-means:-

- ① Scale of Measurements influences Euclidean distance, so Variable Standardisation becomes necessary
- ② Depending on expectations - you may require outlier treatment
- ③ K-means Clustering may be biased on initial centroids - called cluster seeds.
- ④ Maximum clusters is typically inputs and may also impacts the clusters settings created.