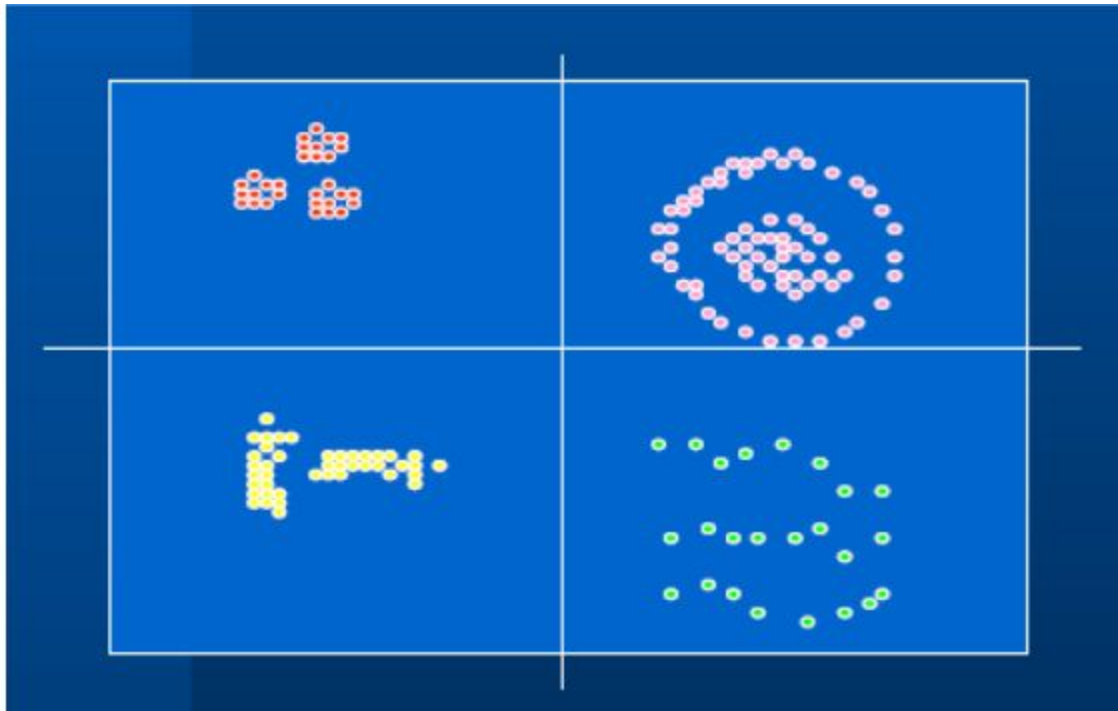


Unsupervised Learning

There are various methods through which we can process the unlabeled datasets. But the main and most common method used to analyze is Clustering. In this we divide data in two different chunks.

Clustering:

A loose definition of clustering could be “the process of organizing objects into groups whose members are similar in some way”. A cluster is therefore a collection of objects which are “similar to each other and are “dissimilar” to the objects belonging to other clusters. Cluster analysis is also used to form descriptive statistics to ascertain whether or not the data consists of a set distinct subgroups, each group representing objects with substantially different properties. The latter goal requires an assessment of the degree of difference between the objects assigned to the respective clusters.



Central to clustering is to decide what constitutes a good clustering. This can only come from subject matter considerations and there is no absolute “best” criterion which would be independent of the final aim of the clustering. For example, we could be interested in finding representatives for homogeneous groups (data reduction), in finding “natural clusters” and

describe their unknown properties (“natural” data types), in finding useful and suitable groupings (“useful” data classes) or in finding unusual data objects (outlier detection).

Clustering algorithms:

Clustering algorithms may be classified as listed below:

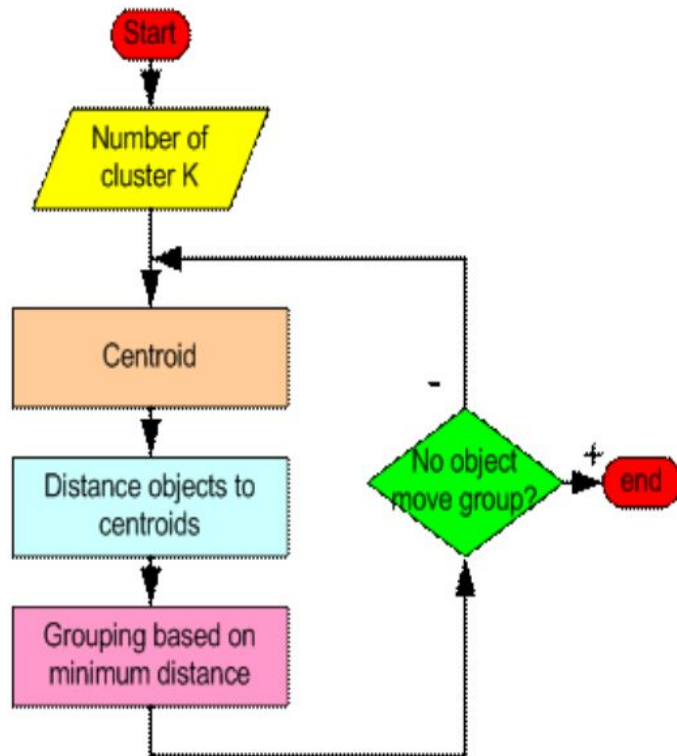
- Exclusive Clustering: In exclusive clustering data are grouped in an exclusive way, so that a certain datum belongs to only one definite cluster. K-means clustering is one example of the exclusive clustering algorithms.
- Hierarchical Clustering: Hierarchical clustering algorithm has two versions: agglomerative clustering and divisive clustering. Agglomerative clustering is based on the union between the two nearest clusters. The beginning condition is realized by setting every datum as a cluster. After a few iterations it reaches the final clusters wanted. Basically, this is a bottom-up version. Divisive clustering starts from one cluster containing all data items. At each step, clusters are successively split into smaller clusters according to some dissimilarity. Basically this is a top-down version.

K-means Clustering:

K-means clustering is an algorithm to classify or to group your objects based on attributes/features into K number of group. K is positive integer number. The grouping is done by minimizing the sum of squares of distances between data and the corresponding cluster centroid. Thus the purpose of K-mean clustering is to classify the data.

Objects	Attribute 1 (X):weight index	Attribute 2 (Y): pH
Medicine A	1	1
Medicine B	2	1
Medicine C	4	3
Medicine D	5	4

We also know before hand that these objects belong to two groups of medicine (cluster 1 and cluster 2). The problem now is to determine which medicines belong to cluster 1 and which medicines belong to the other cluster.



Example:

- The basic step of k-means clustering is simple. In the beginning we determine number of cluster K and we assume the centroid or center of these clusters. We can take any random objects as the initial centroids or the first K objects in sequence can also serve as the initial centroids.

- Then the K means algorithm will do the three steps below until convergence.

Step 1. Begin with a decision on the value of k = number of clusters.

Step 2. Put any initial partition that

classifies the data into k clusters. You may assign the training samples randomly, or systematically as the following:

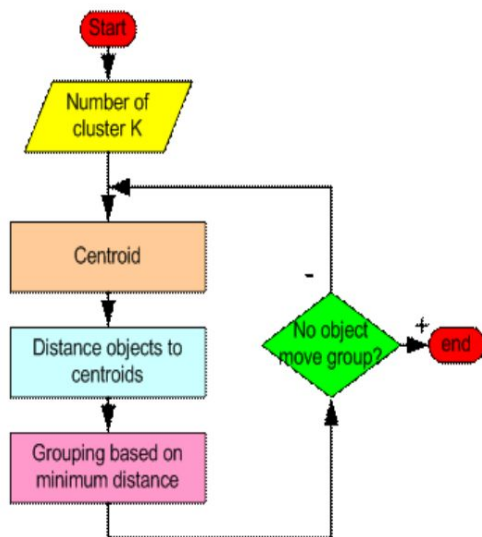
1) Take the first k training sample as single-element clusters.

2) Assign each of the remaining $(N-k)$ training sample to the cluster with the nearest centroid. After each assignment, recompute the centroid of the gaining cluster.

Step 3 . Take each sample in sequence and compute its distance from the centroid of each of the clusters. If a sample is not currently in the cluster with the closest centroid, switch this sample to that cluster and update the centroid of the cluster gaining the new sample and the cluster losing the sample.

Step 4 . Repeat step 3 until convergence is achieved, that is until a pass through the training sample causes no new assignments.

- If the number of data is less than the number of cluster then we assign each data as the centroid of the cluster. Each centroid will have a cluster number.



- If the number of data is bigger than the number of cluster, for each data, we calculate the distance to all centroid and get the minimum distance. This data is said belong to the cluster that has minimum distance from this data.

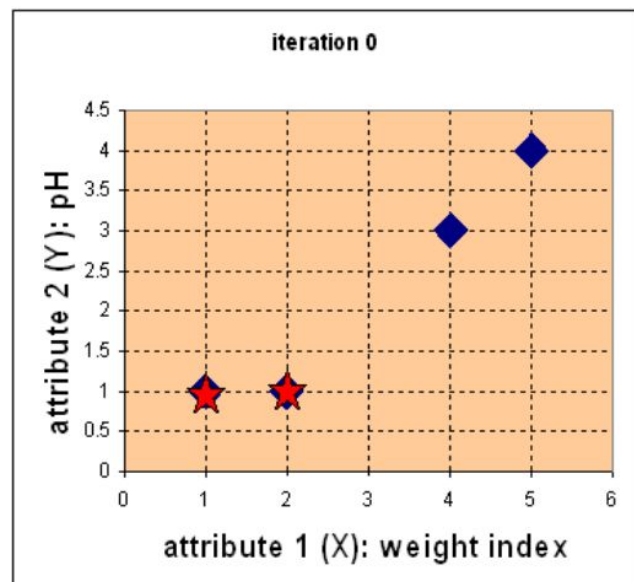
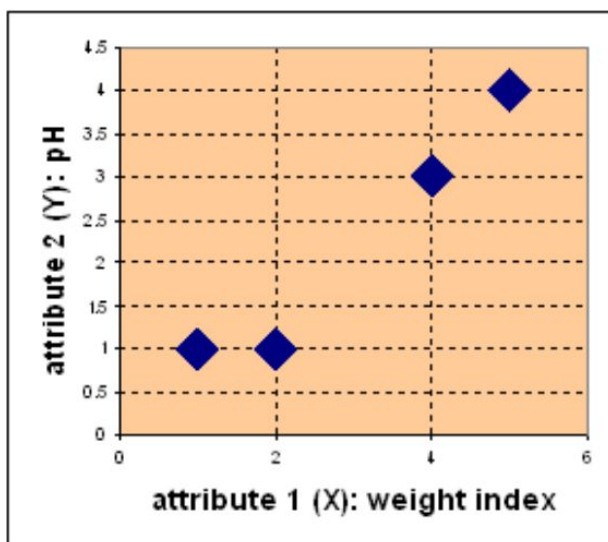
Iterate until stable (= no object move group):

- 1) Determine the centroid coordinate.
- 2) Determine the distance of each object to the centroids.

3) Group the object based on minimum distance.

Objects	Attribute 1 (X):weight index	Attribute 2 (Y): pH
Medicine A	1	1
Medicine B	2	1
Medicine C	4	3
Medicine D	5	4

Suppose we have several objects (4 types of medicines) and each object have two attributes or features as shown in table below. Our goal is to group these objects into K=2 group of medicine



based on the two features (pH and weight index).

Each medicine represents one point with two attributes (X, Y) that we can represent it as coordinate in an attribute space as shown in the figure.

- Initial value of centroids: Suppose we use medicine A and medicine B as the first centroids. Let and denote the coordinate of the centroids,

$$\mathbf{c}_1 = (1, 1)$$

$$\mathbf{c}_2 = (2, 1)$$

$$\mathbf{D}^0 = \begin{bmatrix} 0 & 1 & 3.61 & 5 \\ 1 & 0 & 2.83 & 4.24 \end{bmatrix} \quad \begin{matrix} \mathbf{c}_1 = (1, 1) & \text{group} - 1 \\ \mathbf{c}_2 = (2, 1) & \text{group} - 2 \end{matrix}$$

	A	B	C	D	
X	1	2	4	5	
Y	1	1	3	4	

- Objects-Centroids distance : Lets calculate the distance between cluster centroid to each object. Let us use Euclidean distance, then we have distance matrix at iteration 0 is

Each column in the distance matrix symbolizes the object. The first row of the distance matrix corresponds to the distance of each object to the first centroid and the second row is the distance each object to the second centroid.

Each column in the distance matrix symbolizes the object. The first row of the distance matrix corresponds to the distance of each object to the first centroid and the second row is the distance of each object to the second centroid.

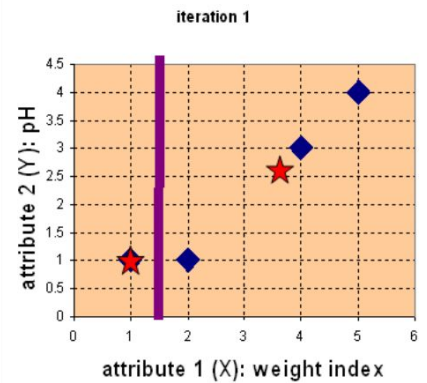
first centroid $\mathbf{c}_1 = (1, 1)$ is $\sqrt{(4-1)^2 + (3-1)^2} = 3.61$, and its distance to the second centroid $\mathbf{c}_2 = (2, 1)$ is $\sqrt{(4-2)^2 + (3-1)^2} = 2.83$

- For example, distance from medicine C = (4, 3) to the first centroid is , and its distance to the second centroid

- Objects clustering : Assign each object based on the minimum distance. Thus, medicine A is assigned to group 1, medicine B to group 2, medicine C to group 2 and medicine D to group 2. The element of Group matrix below is 1 if and only if the object is assigned to that group.

$$\mathbf{c}_2 = \left(\frac{2+4+5}{3}, \frac{1+3+4}{3} \right) = \left(\frac{11}{3}, \frac{8}{3} \right)$$

$$D^1 = \begin{bmatrix} 0 & 1 & 3.61 & 5 \\ 3.14 & 2.36 & 0.47 & 1.89 \\ A & B & C & D \\ \begin{bmatrix} 1 & 2 & 4 & 5 \\ 1 & 1 & 3 & 4 \end{bmatrix} & X & Y \end{bmatrix} \quad \begin{array}{l} \mathbf{c}_1 = (1,1) \text{ group-1} \\ \mathbf{c}_2 = (\frac{11}{3}, \frac{8}{3}) \text{ group-2} \end{array}$$



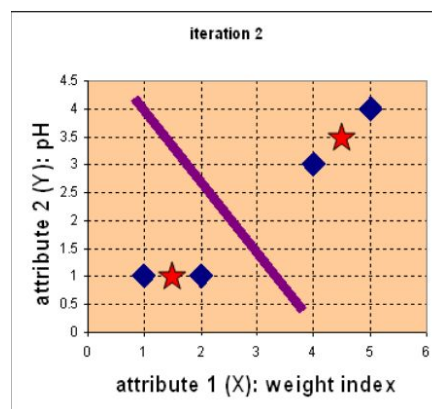
2.

Iteration-1, determine centroids: Knowing the members of each group, now compute the new centroid of each group based on these new memberships. Group 1 only has one member. Thus the centroid remains in . Group 2 now has three members, thus the centroid is the average coordinate among the three members

3. Iteration-1, Objects-Centroids distances : The next step is to compute the distance of all objects to the new centroids. Similar to step 2, we have distance matrix at iteration 1 is

4. Iteration- 2, determine centroids: Now we repeat step 4 to calculate the new centroids coordinate based on the clustering of previous iteration. Group1 and group 2 both has two members, thus the new centroids are and

$$G^1 = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ A & B & C & D \end{bmatrix} \quad \begin{array}{l} \text{group-1} \\ \text{group-2} \end{array} \quad \mathbf{c}_1 = (\frac{1+2}{2}, \frac{1+1}{2}) = (1\frac{1}{2}, 1) \text{ and } \mathbf{c}_2 = (\frac{4+5}{2}, \frac{3+4}{2}) = (4\frac{1}{2}, 3\frac{1}{2})$$



5. Iteration-2, Objects-Centroids distances: Repeat step 2 again, we have new

$$D^2 = \begin{bmatrix} 0.5 & 0.5 & 3.20 & 4.61 \\ 4.30 & 3.54 & 0.71 & 0.71 \end{bmatrix} \quad \begin{matrix} c_1 = (1\frac{1}{2}, 1) \text{ group-1} \\ c_2 = (4\frac{1}{2}, 3\frac{1}{2}) \text{ group-2} \end{matrix}$$

$$\begin{matrix} A & B & C & D \\ \begin{bmatrix} 1 & 2 & 4 & 5 \\ 1 & 1 & 3 & 4 \end{bmatrix} & X & Y \end{matrix}$$

distance matrix at iteration 2 as

6. Iteration-2, Objects clustering: Again, we assign each object based on the minimum distance.

$$G^2 = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \end{bmatrix} \quad \begin{matrix} \text{group-1} \\ \text{group-2} \end{matrix}$$

$$\begin{matrix} A & B & C & D \end{matrix}$$

$$G^2 = G^1$$

We obtain result that . Comparing the grouping of last iteration and this iteration reveals that the objects does not move group anymore. Thus, the computation of the k-mean clustering has reached its stability and no more iteration is needed. We get the final grouping as the results.

Final Grouping –As a Result:

Objects	Attribute 1 (X):weight index	Attribute 2 (Y): pH	Group (Result)
Medicine A	1	1	1
Medicine B	2	1	1
Medicine C	4	3	2
Medicine D	5	4	2

