

Assignment 4

Roll No.-41449

Title : Assignment on K-Means Clustering

Problem Statement : We have given a collection of 8 points. $P1=[0.1,0.6]$, $P2=[0.15,0.71]$, $P3=[0.08,0.9]$ $P4=[0.16, 0.85]$, $P5=[0.2,0.3]$, $P6=[0.25,0.5]$, $P7=[0.24,0.1]$, $P8=[0.3,0.2]$. Perform the k-mean clustering with initial centroids as $m1=P1 = \text{Cluster\#1}=C1$ and $m2=P8=\text{cluster\#2}=C2$.

Answer the following

- 1] Which cluster does P6 belong to?
- 2] What is the population of cluster around $m2$?
- 3] What is updated value of $m1$ and $m2$?

Objective : To understand how k-means clustering algorithm works on the given dataset

Outcome: To implement k-means clustering algorithm

Concept related theory:

K-means clustering is one of the simplest and popular unsupervised machine learning algorithms. Typically, unsupervised algorithms make inferences from datasets using only input vectors without referring to known, or labelled, outcomes. You'll define a target number k , which refers to the number of centroids you need in the dataset. A centroid is the imaginary or real location representing the center of the cluster. Every data point is allocated to each of the clusters through reducing the in-cluster sum of squares. In other words, the K-means algorithm identifies k number of centroids, and then allocates every data point to the nearest cluster, while keeping the centroids as small as possible. The '*means*' in the K-means refers to averaging of the data; that is, finding the centroid.

How the K-means algorithm works

K-means algorithm is an iterative algorithm that tries to partition the dataset into K pre-defined distinct non-overlapping subgroups (clusters) where each data point belongs to **only one group**. It tries to make the inter-cluster data points as similar as possible while also keeping the clusters as different (far) as possible. It assigns data points to a cluster such that the sum of the squared distance between the data points and the cluster's centroid (arithmetic mean of all the data points that belong to that cluster) is at the minimum. The less variation we have within clusters, the more homogeneous (similar) the data points are within the same cluster. Finally, this algorithm aims at minimizing an objective function known as squared error function given by:

$$J(V) = \sum_{i=1}^c \sum_{j=1}^{c_i} (\|x_i - v_j\|)^2$$

where,

' $\|x_i - v_j\|$ ' is the Euclidean distance between x_i and v_j .

' c_i ' is the number of data points in i th cluster.

' c ' is the number of cluster centers.

Algorithmic steps for k-means clustering

Let $X = \{x_1, x_2, x_3, \dots, x_n\}$ be the set of data points and $V = \{v_1, v_2, \dots, v_c\}$ be the set of centers.

- 1) Randomly select 'c' cluster centers.
- 2) Calculate the distance between each data point and cluster centers.
- 3) Assign the data point to the cluster center whose distance from the cluster center is minimum of all the cluster centers..
- 4) Recalculate the new cluster center using:

$$v_i = (1/c_i) \sum_{j=1}^{c_i} x_j$$

where, 'ci' represents the number of data points in *ith* cluster.

- 5) Recalculate the distance between each data point and new obtained cluster centers.
- 6) If no data point was reassigned then stop, otherwise repeat from step 3).

It halts creating and optimizing clusters when either:

- 1.The centroids have stabilized — there is no change in their values because the clustering has been successful.
- 2.The defined number of iterations has been achieved.

Conclusion:

Successfully implemented k-means clustering on given dataset.

Output:



```
In [37]: 1 count=0
2 for i in range(len(labels)):
3     if (labels[i]==1):
4         count=count+1

In [38]: 1 print('No of population around cluster 2:',count-1)

No of population around cluster 2: 3

In [39]: 1 new_centroids = model.cluster_centers_

In [40]: 1 print('Previous value of m1 and m2 is:')
2 print('M1==',centroids[0])
3 print('M1==',centroids[1])

Previous value of m1 and m2 is:
M1== [0.1 0.3]
M1== [0.6 0.2]

In [41]: 1 print('updated value of m1 and m2 is:')
2 print('M1==',new_centroids[0])
3 print('M1==',new_centroids[1])

updated value of m1 and m2 is:
M1== [0.2475 0.275 ]
M1== [0.1225 0.765 ]

In [ ]: 1
```