

## Program – 8

### AIM: To implement k-mean clustering in Python

#### Introduction and Theory

---

**k-means clustering** is a method of vector quantization, originally from signal processing, that is popular for cluster analysis in data mining. k-means clustering aims to partition  $n$  observations into  $k$  clusters in which each observation belongs to the cluster with the nearest mean, serving as a prototype of the cluster. This results in a partitioning of the data space into Voronoi cells.

The problem is computationally difficult (NP-hard); however, efficient heuristic algorithms converge quickly to a local optimum. These are usually similar to the expectation-maximization algorithm for mixtures of Gaussian distributions via an iterative refinement approach employed by both k-means and Gaussian mixture modeling. They both use cluster centers to model the data; however, k-means clustering tends to find clusters of comparable spatial extent, while the expectation-maximization mechanism allows clusters to have different shapes.

The algorithm has a loose relationship to the k-nearest neighbor classifier, a popular machine learning technique for classification that is often confused with k-means due to the name. Applying the 1-nearest neighbor classifier to the cluster centers obtained by k-means classifies new data into the existing clusters. This is known as nearest centroid classifier or Rocchio algorithm.

To process the learning data, the K-means algorithm in data mining starts with a first group of randomly selected centroids, which are used as the beginning points for every cluster, and then performs iterative (repetitive) calculations to optimize the positions of the centroids

It halts creating and optimizing clusters when either:

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

```
1 Initialize k means with random values
2
3 For a given number of iterations:
4     Iterate through items:
5         Find the mean closest to the item
6         Assign item to mean
7         Update mean
```

The k-mean pseudo-code

#### Code

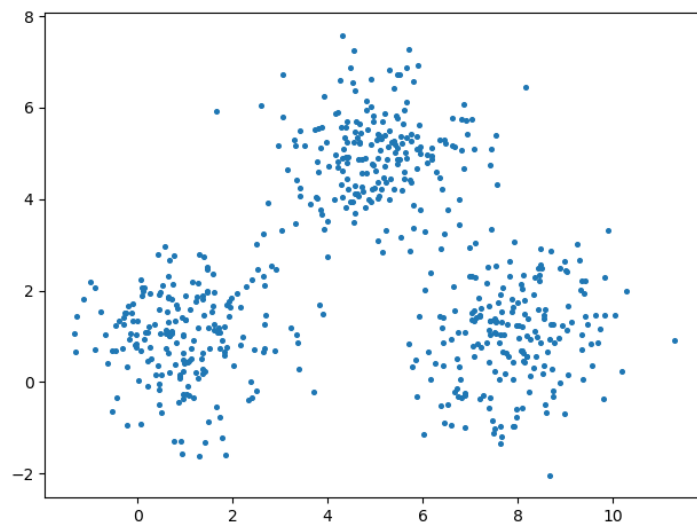
```
1 # Import necessary libraries
2 from copy import deepcopy
3 import numpy as np # linear algebra
4 import pandas as pd # data processing, CSV file I/O (e.g.
5 pd.read_csv)
6 from matplotlib import pyplot as plt
7
8 # Set three centers, the model should predict similar results
9 center_1 = np.array([1,1])
```

## Program – 8

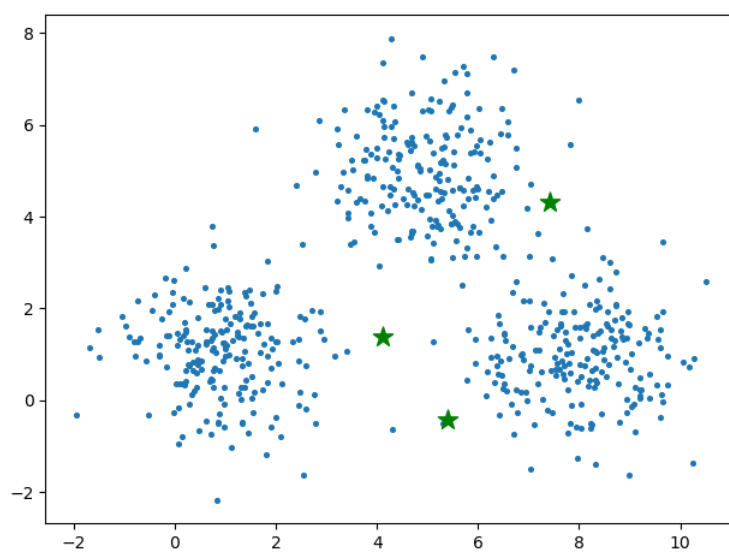
```
10 center_2 = np.array([5,5])
11 center_3 = np.array([8,1])
12
13 # Generate random data and center it to the three centers
14 data_1 = np.random.randn(200, 2) + center_1
15 data_2 = np.random.randn(200,2) + center_2
16 data_3 = np.random.randn(200,2) + center_3
17
18 data = np.concatenate((data_1, data_2, data_3), axis = 0)
19
20 plt.scatter(data[:,0], data[:,1], s=7)
21 plt.show()
22
23 # Number of clusters
24 k = 3
25 # Number of training data
26 n = data.shape[0]
27 # Number of features in the data
28 c = data.shape[1]
29
30 # Generate random centers, here we use sigma and mean to ensure it
31 represent the whole data
32 mean = np.mean(data, axis = 0)
33 std = np.std(data, axis = 0)
34 centers = np.random.randn(k,c)*std + mean
35
36 # Plot the data and the centers generated as random
37 plt.scatter(data[:,0], data[:,1], s=7)
38 plt.scatter(centers[:,0], centers[:,1], marker='*', c='g', s=150)
39 plt.show()
40 centers_old = np.zeros(centers.shape) # to store old centers
41 centers_new = deepcopy(centers) # Store new centers
42 print(data.shape)
43 clusters = np.zeros(n)
44 distances = np.zeros((n, k))
45
46 error = np.linalg.norm(centers_new - centers_old)
47
48 # When, after an update, the estimate of that center stays the same,
49 exit loop
50 while error != 0:
51     # Measure the distance to every center
52     for i in range(k):
53         distances[:, i] = np.linalg.norm(data - centers[i], axis=1)
54     # Assign all training data to closest center
55     clusters = np.argmin(distances, axis=1)
56     centers_old = deepcopy(centers_new)
57     # Calculate mean for every cluster and update the center
58     for i in range(k):
59         centers_new[i] = np.mean(data[clusters == i], axis=0)
60     error = np.linalg.norm(centers_new - centers_old)
61 print(centers_new)
62 # Plot the data and the centers generated as random
63 plt.scatter(data[:,0], data[:,1], s=7)
64 plt.scatter(centers_new[:,0], centers_new[:,1], marker='*', c='g',
65 s=150)
66 plt.show()
```

## Program – 8

### Results & Outputs

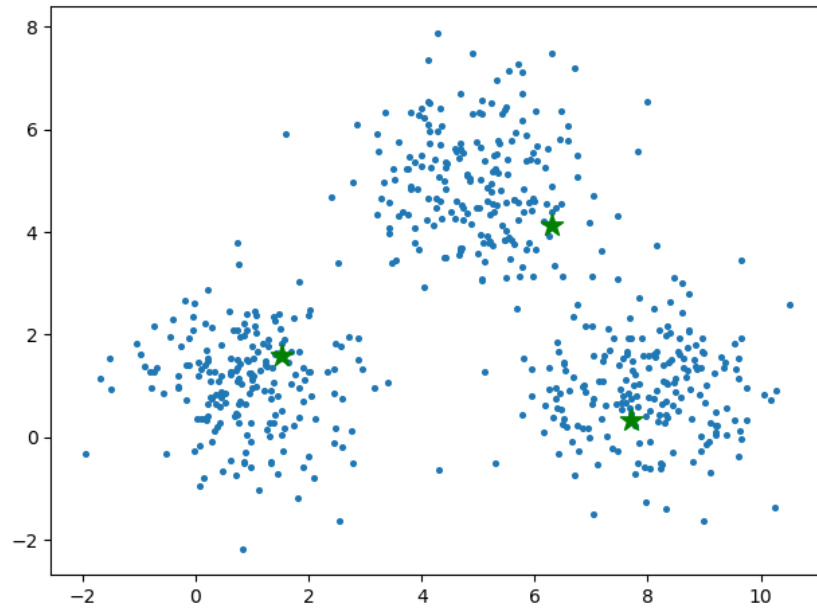


Sample data



Random centroid initialization

## Program – 8



Final cluster assignment

### Findings and Learnings:

1. K-means is one of the simplest unsupervised learning algorithms that solve the well-known clustering problem.
2. We have successfully implemented k-means clustering in Python