

Experiment NO : 07

Instructor: Mr.Katkar Atish R.

Name: Suyog Rawas, Roll No: CS3257

AIM

Implement and demonstrate the K-means clustering algorithm.

INTRODUCTION TO K-MEANS CLUSTERING ALGORITHM

The K-means clustering algorithm is a popular unsupervised machine learning technique used for grouping data into clusters based on their similarity. It is an iterative algorithm that aims to partition a dataset into K distinct clusters, where K is a predefined number.

The algorithm works by assigning data points to the nearest centroid, which represents the center of each cluster. Initially, K centroids are randomly selected from the dataset. Then, the algorithm iteratively updates the centroids and reassigns data points to the nearest centroid until convergence is achieved.

The steps involved in the K-means clustering algorithm :

1. Initialization: Randomly select K data points as initial centroids.
2. Assignment: Assign each data point to the nearest centroid based on a distance metric, commonly the Euclidean distance.
3. Update: Recalculate the centroids by taking the mean of all data points assigned to each centroid.
4. Iteration: Repeat the assignment and update steps until convergence or a maximum number of iterations is reached.
5. Convergence: Stop the algorithm when the centroids no longer change significantly or when the maximum number of iterations is reached.

The algorithm aims to minimize the within-cluster sum of squares, also known as the inertia or distortion. It seeks to create compact and well-separated clusters by minimizing the total distance between data points and their respective centroids.

One challenge in using K-means clustering is determining the optimal value of K, the number of clusters. Various methods, such as the elbow method or silhouette analysis, can be employed to find an appropriate K value.

K-means clustering is widely used in various applications, including customer segmentation, image segmentation, anomaly detection, and recommender systems. Its simplicity and efficiency make it a popular choice for exploratory data analysis and clustering tasks.

DATASET DESCRIPTION :

In this code, the data list represents your dataset, where each element is a data point with its corresponding features. The K variable determines the the number of clusters you want to create.

The max-iterations sets the maximum number of iterations for convergence. The code initializes centroids randomly, assigns data points to the closest centroids, updates the centroids based on the mean of the assigned points, and repeats this process until convergence or reaching the maximum number of

IMPLEMENTATION CODE FOR K-MEANS CLUSTERING ALGORITHM

```

1
2 import random
3 import math
4 # Euclidean distance between two points
5 def euclidean_distance ( point1 , point2 ) :
6     squared_distance = 0
7     for i in range ( len ( point1 ) ) :
8         squared_distance += ( point1 [ i ] - point2 [ i ] ) ** 2
9     return math . sqrt ( squared_distance )
10
11 # Initialize centroids randomly
12 def initialize_centroids ( data , k ) :
13     centroids = random . sample ( data , k )
14     return centroids
15
16 # Assign each data point to the closest centroid
17 def assign_clusters ( data , centroids ) :
18     clusters = [[] for _ in centroids ]
19     for point in data :
20         distances = [ euclidean_distance ( point , centroid ) for centroid in
21             centroids ]
22         closest_centroid_index = distances . index ( min ( distances ) )
23         clusters [ closest_centroid_index ] . append ( point )
24     return clusters
25
26 # Update centroids based on the mean of the assigned points
27 def update_centroids ( clusters ) :
28     centroids = []
29     for cluster in clusters :
30         centroid = [ sum ( dim ) / len ( cluster ) for dim in zip ( * cluster ) ]
31         centroids . append ( centroid )
32     return centroids
33
34 # K- means clustering algorithm
35 def k_means ( data , k , max_iterations ) :
36     centroids = initialize_centroids ( data , k )
37     for _ in range ( max_iterations ) :
38         clusters = assign_clusters ( data , centroids )
39         new_centroids = update_centroids ( clusters )
40
41         # Check for convergence
42         if new_centroids == centroids :
43             break
44
45         centroids = new_centroids
46
47     return clusters , centroids
48
49 # Example usage
50 data = [
51     [2 , 10] ,
52     [2 , 5] ,
53     [8 , 4] ,
54     [5 , 8] ,
55     [7 , 5] ,
56     [6 , 4] ,
57     [1 , 2] ,
58     [4 , 9]
59 ]
60 k = 2
61 max_iterations = 100

```

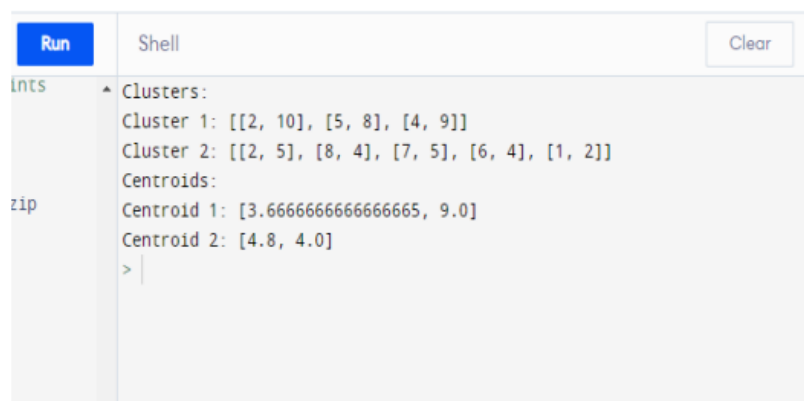
```

62
63 clusters , centroids = k_means ( data , k , max_iterations )
64
65 # Print the results
66 print ( " Clusters :")
67 for i , cluster in enumerate ( clusters ) :
68     print ( f" Cluster {i +1}: { cluster }")
69 print ( " Centroids :")
70 for i , centroid in enumerate ( centroids ) :
71     print ( f" Centroid {i +1}: { centroid }")

```

OUTPUT

1 Output Of K-means Clustering Algorithm



```

Run Shell Clear
Clusters:
Cluster 1: [[2, 10], [5, 8], [4, 9]]
Cluster 2: [[2, 5], [8, 4], [7, 5], [6, 4], [1, 2]]
Centroids:
Centroid 1: [3.6666666666666665, 9.0]
Centroid 2: [4.8, 4.0]
>

```

Figure 1: K-means Clustering Output

2 Work done by the Algorithm

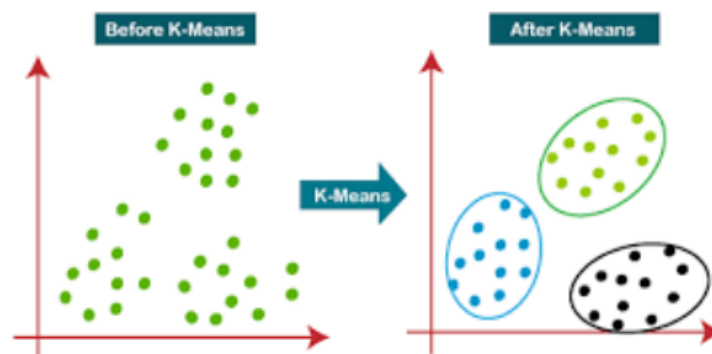


Figure 2: Work done by the Algorithm

CONCLUSION

In conclusion, the K-means clustering algorithm is a popular unsupervised machine learning technique used for grouping similar data points into clusters. It is a simple and efficient algorithm that can handle large datasets and is widely used in various applications such as customer segmentation, image compression, and anomaly detection.