

Implementing K-Means Clustering

Mayisha Farzana

dept.Computer Science and Engineering
Ahsanullah University of Science and Technology
Dhaka, Bangladesh
160204028@aust.edu

Abstract—This experiment aims to learn about the K-Means clustering algorithm, which is one of the most commonly used clustering algorithms. This is a method of learning that is not supervised. This suggests that the class label is undefined. However, in order to determine the case, we need to know the ground truth in order to test the algorithm's efficiency.

Index Terms—Unsupervised learning; clustering; K-means clustering; python code

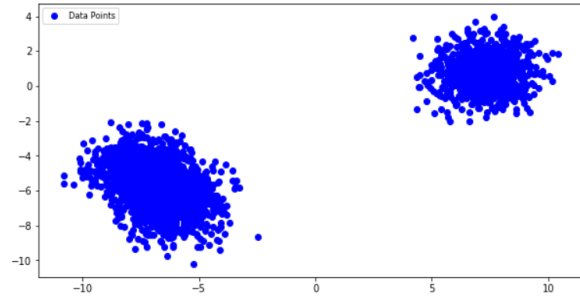
I. INTRODUCTION

K-means is one of the most basic unsupervised learning algorithms for solving the well-known problem of clustering. The procedure employs a basic and straightforward method for classifying a given data set using a predetermined number of clusters (assume k clusters). In the beginning, we determine the 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. Since different locations produce different results, these centroids should be positioned with caution. As a result, the best option is to put them as far apart as possible. The next step is to connect each point in a data set with the centroid that is closest to it. If there are no pending points, the first stage is over, and an early grouping is completed. At this stage, we must re-calculate k new centroids to serve as barycenters for the clusters formed in the previous phase. Once we have created these k new centroids, we will need to rebind the same data set points to the nearest new centroid. A loop has been generated. As a result of this loop, we may notice that the k centroids change their location step by step until no more changes are done. In other words, centroids do not move anymore.

II. EXPERIMENTAL DESIGN / METHODOLOGY

A. Choosing the initial centroids: Since in this experiment we choose to cluster the whole dataset in three cluster we have to choose three centroids. At first we choose it randomly but the centroid should be as far as possible so we fixed the range.

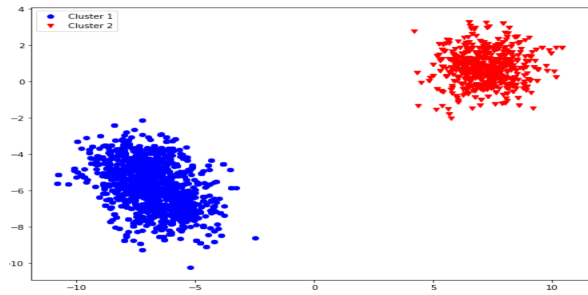
B. Calculating the distance matrix: Here in this experiment we calculate the Euclidian distance between the sample points and centroid. After that this distance matrix will be used to assign groups to each column of the distance matrix. In which row it will be determined by the minimum distance stored in the distance matrix. After assigning the group now we have to count that how many sample are in each group/cluster.



Then we will calculate the new centroids for the next iteration. For distance calculation, we will use euclidean distance.

$$\text{Euclidean distance} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

C. Calculating new centroid: The new centroid will be all the data samples average which corresponding row has a 1 in the group matrix. Thus we will get new set of centroid. After that the whole procedure will continue until each group has the same number of samples. In our experiment, the centroid value will not change much after some iteration. It will change sometimes cause the initial centroid are chosen randomly. So it is possible to get slight different output but the difference won't be very high.



III. RESULT ANALYSIS

From the test points, we will get different classes for each value. We will see that we found three sets of red, blue coded centroids and tuned them to perfection until we found the same set of matrices to represent the clusters.

IV. CONCLUSION

The K-means clustering algorithm is observed in this experiment. This is a commonly used and popular clustering strategy. If we want to calculate the cluster purity output in this algorithm, we need to know the ground truth.

V. ALGORITHM IMPLEMENTATION / CODE

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import random
5 random.seed(1)
6 train_dataset = np.loadtxt('data_k_mean.txt')
7 print(train_dataset)
8 length=len(train_dataset)
9 print(length)
10 f, ax = plt.subplots()
11 f.set_figheight(5)
12 f.set_figwidth(10)
13 ax.scatter(train_dataset[:,0],train_dataset[:,1],
14            marker='o',color='b',label='Data Points')
15 legend = ax.legend(loc='upper left', shadow=False,
16                   fontsize='small',labelspacing=0.5)
17 legend.get_frame().set_facecolor('None')
18 plt.show()
19 random.shuffle(train_dataset)
20 print(train_dataset)
21 def dis(a,b):
22     return np.sqrt(((a[0]-b[0])**2)+((a[1]-b[1])**2))
23
24 def calculate_mean(cluster):
25     x_sum=0
26     y_sum=0
27     for p in cluster:
28         x_sum+=train_dataset[p][0]
29         y_sum+=train_dataset[p][1]
30
31     n=max(1,len(cluster))
32     return (x_sum)/n, (y_sum)/n
33 #####
34 def get_new_cluster(means,k):
35     new=[]
36     for i in range(k):
37         x=[]
38         new.append(x)
39
40     for tp in range(len(train_dataset)):
41         mn=10000000000#distance big value
42         nc=-1#new cluster
43         for i in range(k):
44             d=dis(means[i],train_dataset[tp])
45             if(d<mn):
46                 mn=d#new distance
47                 nc=i
48         new[nc].append(tp)
49     return new
50
51 def check(c_old,c_new,k):
52     for i in range(k):
53         if(len(c_old[i])!=len(c_new[i])):
54             return False
55     for i in range(k):
56         c_old[i].sort()
57         c_new[i].sort()
58         for j in range(len(c_old[i])):
59             if(c_old[i][j]!=c_new[i][j]):
60                 return False
61     return True
62
63 def k_means(k):
64     ar=[]
65     for i in range(k):
```

```
61         x=[]
62         ar.append(x)
63
64     for i in range(k):
65         ar[i].append(i)
66     prev=ar.copy()
67     while(True):
68         means=[]
69         for c in range(k):
70             means.append(calculate_mean(ar[c]))
71         curr=get_new_cluster(means,k)
72         if(check(ar,curr,k)):
73             break
74         ar=curr.copy()
75     return curr
76 K=int(input('Enter a value for K:'))
77 final_cluster=k_means(K)
78
79 for i in range(K):
80     print(final_cluster[i])
81 x_point=[]
82 y_point=[]
83 for i in range(K):#
84     x=[]
85     y=[]
86     x_point.append(x)
87     y_point.append(y)
88 for i in range(K):
89     for x in final_cluster[i]:
90         x_point[i].append(train_dataset[x][0])
91         y_point[i].append(train_dataset[x][1])
92 fig,ax=plt.subplots()
93 colors=['blue','red','green','yellow','orange','c','m',
94         'y','k']
95 marker=['o','v','l','^','s','p','+','D','X','P','2',
96         'H','3','d','4']
97 for i in range(K):
98     name="Cluster "+str(i+1)
99     ax.scatter(x_point[i],y_point[i],marker=marker[i],
100              color=colors[i],label=name)
101 fig.set_figheight(8)
102 fig.set_figwidth(10)
103 legend.get_frame().set_facecolor('None')
104 ax.legend()
105 plt.show()
```