

Application 13

Unsupervised Machine Learning

Clustering using K Mean Algorithm

Consider below Machine Learning application which implements K mean Algorithm for randomly generated dataset.

```
1 import numpy as np
2 import pandas as pd
3 from copy import deepcopy
4 from matplotlib import pyplot as plt
5
6 def MarvellousKMean():
7     # Set three centers, the model should predict similar results
8     center_1 = np.array([1,1])
9     print(center_1)
10
11     center_2 = np.array([5,5])
12     print(center_2)
13
14     center_3 = np.array([8,1])
15     print(center_3)
16
17     # Generate random data and center it to the three centers
18     data_1 = np.random.randn(7, 2) + center_1
19     print("Elements of first cluster with size"+str(len(data_1)))
20     print(data_1)
21
22     data_2 = np.random.randn(7,2) + center_2
23     print("Elements of second cluster with size"+str(len(data_2)))
24     print(data_2)
25
26     data_3 = np.random.randn(7,2) + center_3
27     print("Elements of third cluster with size"+str(len(data_3)))
28     print(data_3)
29
30     data = np.concatenate((data_1, data_2, data_3), axis = 0)
31     print("Size of complete data set"+str(len(data)))
32
33     plt.scatter(data[:,0], data[:,1], s=7)
34     plt.title('Marvellous Infosystems : Input Dataset')
35     plt.show()
36
37     # Number of clusters
38     k = 3
39
40     # Number of training data
41     n = data.shape[0]
42     print("Total number of elements are",n)
43
44     # Number of features in the data
45     c = data.shape[1]
46     print("Total number of features are",c)
47
48     # Generate random centers, here we use sigma and mean to ensure it represent the whole data
49     mean = np.mean(data, axis = 0)
50     print("Value of mean",mean)
51
52     # Calculate standard deviation
53     std = np.std(data, axis = 0)
54     print("Value of std",std)
55
56     centers = np.random.randn(k,c)*std + mean
```

```

57 print("Random points are",centers)
58
59 # Plot the data and the centers generated as random
60 plt.scatter(data[:,0], data[:,1],c='r', s=7)
61 plt.scatter(centers[:,0], centers[:,1], marker='*', c='g', s=150)
62 plt.title('Marvellous Infosystems : Input Dataset with random centroid *')
63 plt.show()
64
65 centers_old = np.zeros(centers.shape) # to store old centers
66 centers_new = deepcopy(centers) # Store new centers
67
68 print("Values of old centroids")
69 print(centers_old)
70
71 print("Values of new centroids")
72 print(centers_new)
73
74 data.shape
75 clusters = np.zeros(n)
76 distances = np.zeros((n,k))
77
78 print("Initial distances are")
79 print(distances)
80
81 error = np.linalg.norm(centers_new - centers_old)
82
83 # When, after an update, the estimate of that center stays the same, exit loop
84 while error != 0:
85
86     # Measure the distance to every center
87     print("Measure the distance to every center")
88     for i in range(k):
89         print("Iteration number ",i)
90         distances[:,i] = np.linalg.norm(data - centers[i], axis=1)
91
92     # Assign all training data to closest center
93     clusters = np.argmin(distances, axis = 1)
94
95     centers_old = deepcopy(centers_new)
96
97     # Calculate mean for every cluster and update the center
98     for i in range(k):
99         centers_new[i] = np.mean(data[clusters == i], axis=0)
100     error = np.linalg.norm(centers_new - centers_old)
101 # end of while
102 centers_new
103
104 # Plot the data and the centers generated as random
105 plt.scatter(data[:,0], data[:,1], s=7)
106 plt.scatter(centers_new[:,0], centers_new[:,1], marker='*', c='g', s=150)
107 plt.title('Marvellous Infosystems : Final data with Centroid')
108 plt.show()
109
110 def main():
111     print("----- Marvellous Infosystems by Piyush Khairnar-----")
112
113     print("Unsuervised Machine Learning")
114
115     print("Clustering using K Mean Algorithm")
116
117     MarvellousKMean()
118
119 if __name__ == "__main__":
120     main()
121

```

Output of above application

```

(base) MacBook-Pro-de-MARVELLOUS:Clustering marvellous$ python3 Clustering.py
---- Marvellous Infosystems by Piyush Khairnar ----
Unsuervised Machine Learning
Clustering using K Mean Algorithm
[1 1]
[5 5]
[8 1]
Elements of first cluster with size7
[[ 1.78591077  0.52204584]
 [ 1.06552455  1.58010131]
 [ 2.14878607 -0.60838652]
 [ 0.40137758  1.64781455]
 [ 0.39167985 -0.55167997]
 [ 0.05705409 -1.23997492]
 [ 0.72317674  0.3989142 ]]
Elements of second cluster with size7
[[ 5.13100873  5.68867927]
 [ 4.63122904  5.4776318 ]
 [ 5.08787673  4.79849352]
Elements of third cluster with size7
[[ 7.91639025 -0.08093321]
 [ 7.6395752  2.36452259]
 [ 8.55221324  0.68666628]
 [10.88013951  2.43256217]
 [ 7.86321404 -0.29989975]
 [ 7.10880253  1.65059962]
 [ 8.49230784  1.55607116]]
Size of complete data set21
Total number of elements are 21
Total number of features are 2
Value of mean [4.95998816 2.40135806]
Value of std [3.17645752 2.5908486 ]
Random points are [[ 9.91497286  3.02323184]
 [ 5.5055685  5.63498908]
 [-0.77374946 -0.56501315]]
Values of old centroids
[[0. 0.]
 [0. 0.]
 [0. 0.]]

```

```
[ 0.  0.  0.]  
[ 0.  0.  0.]  
[ 0.  0.  0.]  
[ 0.  0.  0.]  
[ 0.  0.  0.]  
[ 0.  0.  0.]  
[ 0.  0.  0.]  
[ 0.  0.  0.]  
[ 0.  0.  0.]  
[ 0.  0.  0.]
```

Measure the distance to every center

Iteration number 0

Iteration number 1

Iteration number 2

Measure the distance to every center

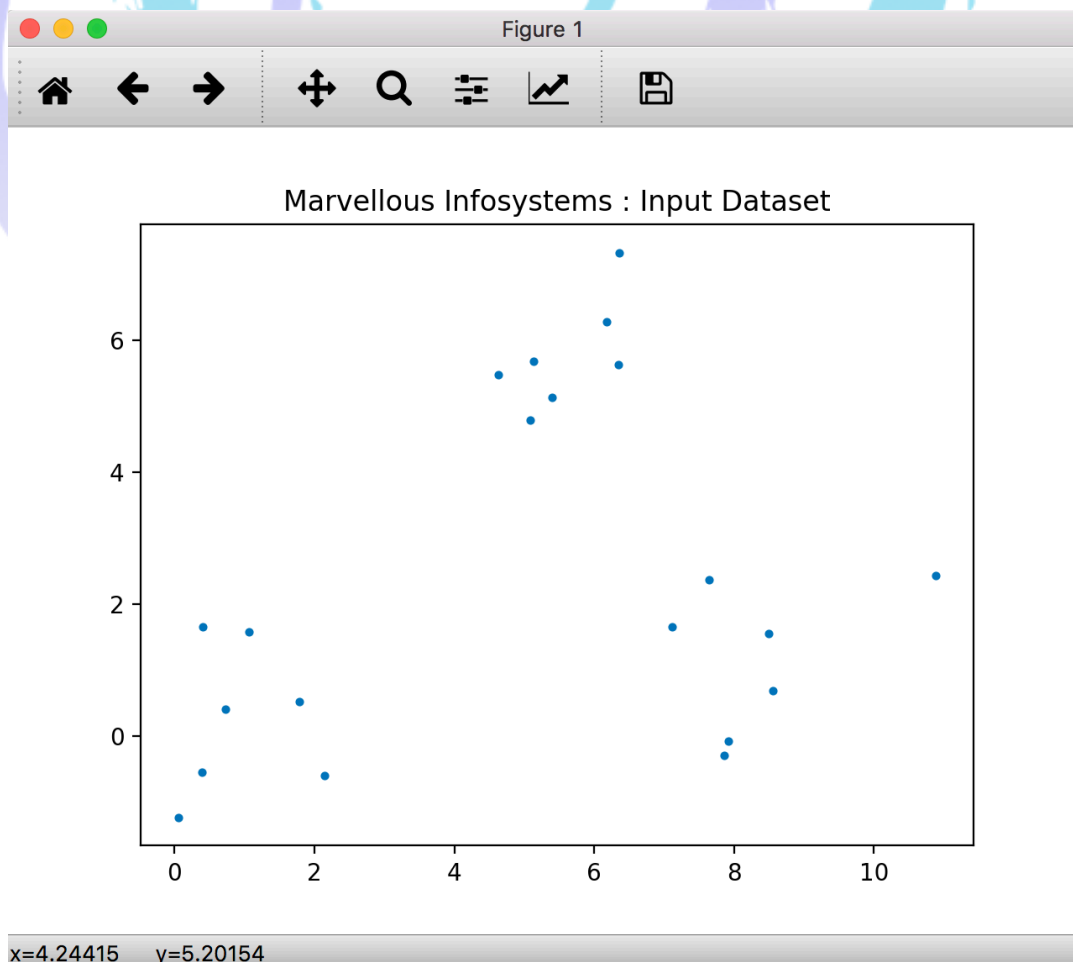
Iteration number 0

Iteration number 1

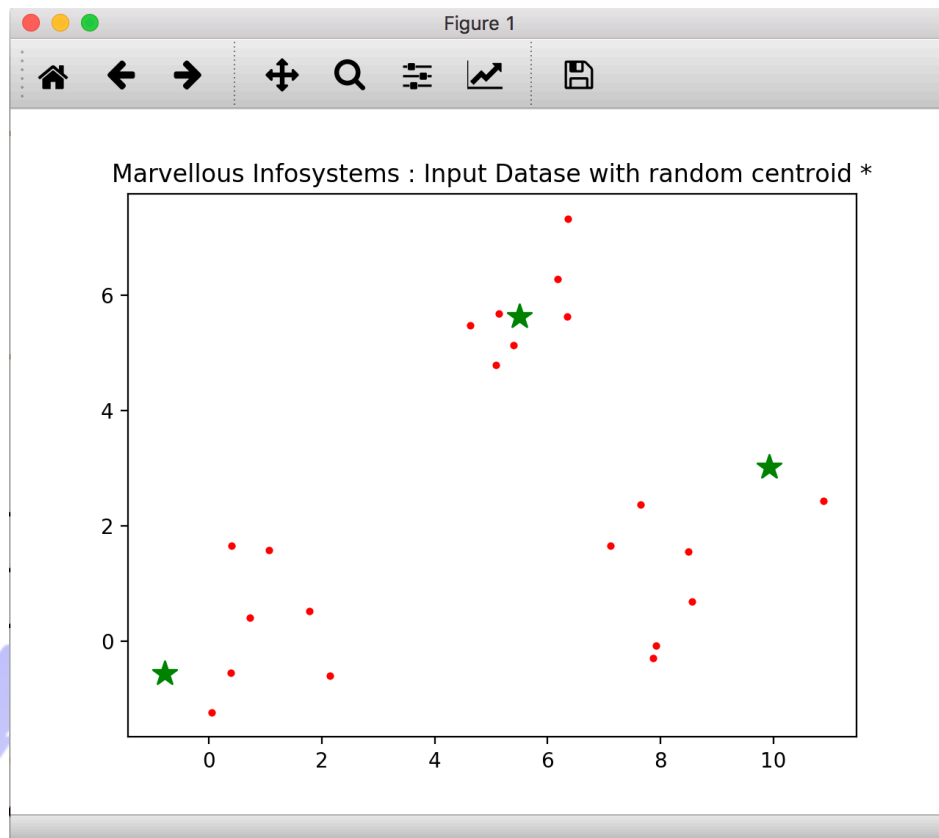
Iteration number 2

(base) MacBook-Pro-de-MARVELLOUS: Clustering marvellous\$ █

Plotting of input data set



Plotting of input data set with random centroids



Plotting of data elements after clustering

