

Cluster analysis or clustering

Clustering is the task of grouping a set of objects in such a way that objects in the same group (called a cluster) are more similar (in some sense) to each other than to those in other groups (clusters). Clustering is not one specific algorithm, but the task of grouping data. This can be achieved by various algorithms that differ significantly in their notion of what constitutes a cluster and how to efficiently find them. Out of 100 published clustering algorithms there are common clustering algorithm,

Hard Clustering:

Connectivity-based clustering (Hierarchical Clustering) Centroid-based Clustering (k-means clustering) Distribution-based Clustering Density-based Clustering

Soft Clustering:

Fuzzy Clustering

Fuzzy Clustering

In **Fuzzy clustering** each element has a set of membership coefficients corresponding to the degree of being in a given cluster. Points close to the center of a cluster, may be in the cluster to a higher degree than points in the edge of a cluster. The degree, to which an element belongs to a given cluster, is a numerical value varying from 0 to 1.

Fuzzy clustering algorithms seeks to **minimize cluster memberships** and distances, but we will focus on Fuzzy C-Means Clustering algorithm.

Fuzzy c-means developed in 1973 and improved in 1981. It's very similar to k-means algorithm in a structure way:

1. Choose number of clusters.
2. Assign coefficients randomly to each data point for being in the clusters.
3. Repeat until algorithm converged (Objective Function C minimizes cluster memberships and distances):
4. Compute the centroid for each cluster
5. Compute each data points' coefficients of being in the clusters.

The main difference with k-means cluster is that **objective function for fuzzy c-means algorithm** allows different cluster membership with probability values, where k-means cluster has strict objective function allows only one cluster membership.

▼ Importing Libraries

```
1 import pandas as pd # reading all required header files
2 import numpy as np
3 import random
4 import operator
5 import math
6 import matplotlib.pyplot as plt
7 from scipy.stats import multivariate_normal
8 from sklearn import datasets
9 from sklearn.decomposition import PCA
10 from mpl_toolkits.mplot3d import Axes3D
```

```
1 # import some data to play with
2 iris = datasets.load_iris()
```

```
1 iris
```

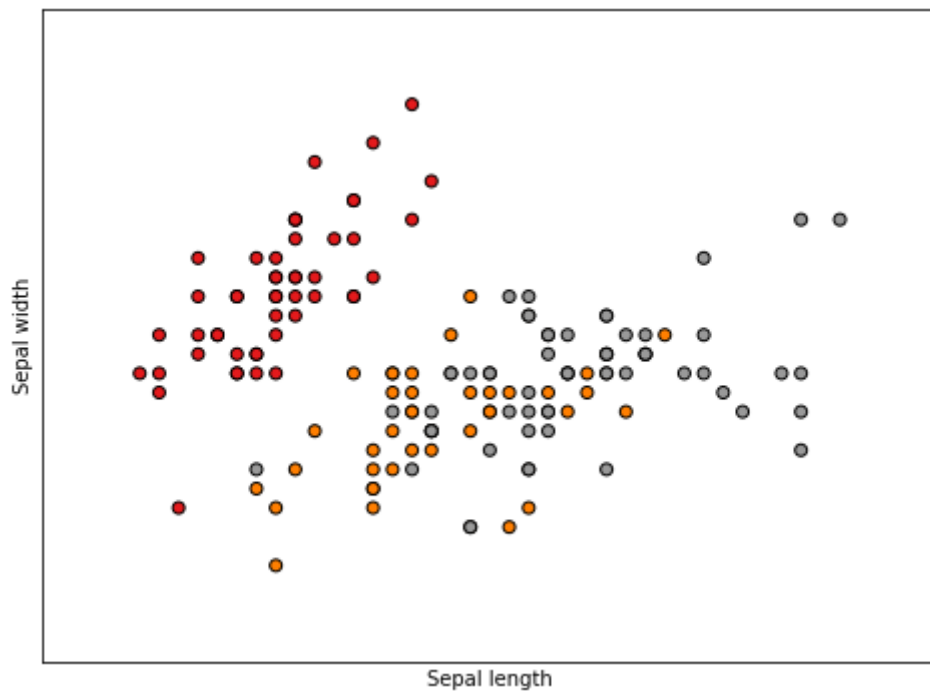


```
'feature_names': ['sepal length (cm)',  
                  'sepal width (cm)',  
                  'petal length (cm)',  
                  'petal width (cm)'],  
'filename': '/usr/local/lib/python3.6/dist-packages/sklearn/datasets/data/iris.csv',  
'target': array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 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, 0, 0, 0, 0, 0, 0,  
                 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,  
                 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,  
                 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 2, 2, 2, 2, 2, 2, 2, 2,  
                 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,  
                 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2]),  
'target_names': array(['setosa', 'versicolor', 'virginica'], dtype='<U10')}
```



```
1 X = iris.data[:, :2] # we only take the first two features.
2 y = iris.target
3
4 x_min, x_max = X[:, 0].min() - .5, X[:, 0].max() + .5
5 y_min, y_max = X[:, 1].min() - .5, X[:, 1].max() + .5
6
7 plt.figure(2, figsize=(8, 6))
8 plt.clf()
9
10 # Plot the training points
11 plt.scatter(X[:, 0], X[:, 1], c=y, cmap=plt.cm.Set1,
12             edgecolor='k')
13 plt.xlabel('Sepal length')
14 plt.ylabel('Sepal width')
15
16 plt.xlim(x_min, x_max)
17 plt.ylim(y_min, y_max)
18 plt.xticks(())
19 plt.yticks(())
20
21 # To get a better understanding of interaction of the dimensions
22 # plot the first three PCA dimensions
23 fig = plt.figure(1, figsize=(8, 6))
24 ax = Axes3D(fig, elev=-150, azim=110)
25 X_reduced = PCA(n_components=3).fit_transform(iris.data)
26 ax.scatter(X_reduced[:, 0], X_reduced[:, 1], X_reduced[:, 2], c=y,
27            cmap=plt.cm.Set1, edgecolor='k', s=40)
28 ax.set_title("First three PCA directions")
29 ax.set_xlabel("1st eigenvector")
30 ax.w_xaxis.set_ticklabels([])
31 ax.set_ylabel("2nd eigenvector")
32 ax.w_yaxis.set_ticklabels([])
33 ax.set_zlabel("3rd eigenvector")
34 ax.w_zaxis.set_ticklabels([])
35
36 plt.show()
```





▼ Installing package fuzzy-c-means

```
!pip install fuzzy-c-means
```

```
1 !pip install fuzzy-c-means
```

```
Collecting fuzzy-c-means
```

```
  Downloading https://files.pythonhosted.org/packages/99/fa/55219e166bb52dd7ed8e35d74
  Requirement already satisfied: scipy>=1.1.0 in /usr/local/lib/python3.6/dist-packages
  Requirement already satisfied: numpy>=1.15.4 in /usr/local/lib/python3.6/dist-packages
  Installing collected packages: fuzzy-c-means
  Successfully installed fuzzy-c-means-0.0.6
```

▼ Importing Libraries

```
1 #https://pypi.org/project/fuzzy-c-means/
2 from fcmeans import FCM
3 from sklearn.datasets import make_blobs
```

```

3 from sklearn.datasets import make_blobs
4 from matplotlib import pyplot as plt
5 from seaborn import scatterplot as scatter

```

```

➤ /usr/local/lib/python3.6/dist-packages/statsmodels/tools/_testing.py:19: FutureWarning
import pandas.util.testing as tm

```

▼ Create artificial dataset

```

1 #create artifitial dataset
2 n_samples = 50000
3 n_bins = 3 # use 3 bins for calibration_curve as we have 3 clusters here
4 centers = [(-5, -5), (0, 0), (5, 5)]

```

Applying Fuzzy C-means

```

1 X,_ = make_blobs(n_samples=n_samples, n_features=2, cluster_std=1.0,
2                 centers=centers, shuffle=False, random_state=42)
3
4 # fit the fuzzy-c-means
5 fcm = FCM(n_clusters=3)
6 fcm.fit(X)

```

```

➤ <fcmeans.fcm.FCM at 0x7f9f79fe79e8>

```

```

1 # outputs
2 fcm_centers = fcm.centers
3 fcm_labels = fcm.u.argmax(axis=1)

```

```

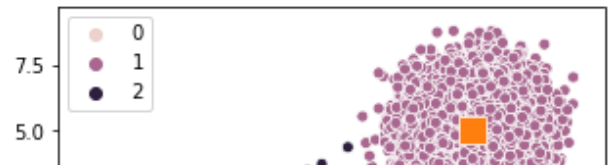
1 # plot result
2 %matplotlib inline
3 f, axes = plt.subplots(1, 2, figsize=(11,5))
4 scatter(X[:,0], X[:,1], ax=axes[0])
5 scatter(X[:,0], X[:,1], ax=axes[1], hue=fcm_labels)
6 scatter(fcm_centers[:,0], fcm_centers[:,1], ax=axes[1], marker="s", s=200)
7 plt.show()

```

```

➤

```



1



```
1 from google.colab import drive
2 drive.mount("/content/drive")
```



```
1 import pandas as pd
2 import numpy as np
3 import random
4 import operator
5 import math
6
7
8 df_full = pd.read_csv("/content/drive/My Drive/Colab Notebooks/clustering/SPECTF_M
9 columns = list(df_full.columns)
10 features = columns[:len(columns)-1]
11 class_labels = list(df_full[columns[-1]])
12 df = df_full[features]
13
14 # Number of Attributes
15 num_attr = len(df.columns) - 1
16
17 # Number of Clusters
18 k = 2
19
20 # Maximum number of iterations
21 MAX_ITER = 1000
22
23 # Number of data points
24 n = len(df)
25
26 # Fuzzy parameter
27 m = 2.00
28
29 def accuracy(cluster_labels, class_labels):
30     county = [0,0]
31     countn = [0,0]
32     tp = [0, 0]
33     tn = [0, 0]
34     fp = [0, 0]
35     fn = [0, 0]
36
37     for i in range(len(df)):
38         # Yes = 1, No = 0
39         if cluster_labels[i] == 1 and class_labels[i] == 'Yes':
40             tp[0] = tp[0] + 1
41         if cluster_labels[i] == 0 and class_labels[i] == 'No':
42             tp[1] = tp[1] + 1
```



```

42         tn[0] = tn[0] + 1
43     if cluster_labels[i] == 1 and class_labels[i] == 'No':
44         fp[0] = fp[0] + 1
45     if cluster_labels[i] == 0 and class_labels[i] == 'Yes':
46         fn[0] = fn[0] + 1
47
48     for i in range(len(df)):
49         # Yes = 0, No = 1
50         if cluster_labels[i] == 0 and class_labels[i] == 'Yes':
51             tp[1] = tp[1] + 1
52         if cluster_labels[i] == 1 and class_labels[i] == 'No':
53             tn[1] = tn[1] + 1
54         if cluster_labels[i] == 0 and class_labels[i] == 'No':
55             fp[1] = fp[1] + 1
56         if cluster_labels[i] == 1 and class_labels[i] == 'Yes':
57             fn[1] = fn[1] + 1
58
59     a0 = float((tp[0] + tn[0]))/(tp[0] + tn[0] + fn[0] + fp[0])
60     a1 = float((tp[1] + tn[1]))/(tp[1] + tn[1] + fn[1] + fp[1])
61     p0 = float(tp[0])/(tp[0] + fp[0])
62     try:
63         p1= float(tp[1])/(tp[1] + fp[1])
64     except ZeroDivisionError:
65         print ('0')
66         p1=0
67     r0 = float(tp[0])/(tp[0] + fn[0])
68     r1 = float(tp[1])/(tp[1] + fn[1])
69
70     accuracy = [a0*100,a1*100]
71     precision = [p0*100,p1*100]
72     recall = [r0*100,r1*100]
73
74     return accuracy, precision, recall
75
76
77 def initializeMembershipMatrix():
78     membership_mat = list()
79     for i in range(n):
80         random_num_list = [random.random() for i in range(k)]
81         summation = sum(random_num_list)
82         temp_list = [x/summation for x in random_num_list]
83         membership_mat.append(temp_list)
84     return membership_mat
85
86
87 def calculateClusterCenter(membership_mat):
88     cluster_centers = list(zip(*membership_mat))
89     #cluster_centers = list()
90     for j in range(k):
91         x = list(cluster_centers[j])
92         xraised = [e ** m for e in x]
93         denominator = sum(xraised)
94         temp_num = list()
95         for i in range(n):
96             data_point = list(df.iloc[i])
97             num = [xraised[i] * val for val in data_point]

```

```

97         prod = [xraiseu[i] * val for val in data_point]
98         temp_num.append(prod)
99         numerator = map(sum, zip(*temp_num))
100         center = [z/denominator for z in numerator]
101         cluster_centers.append(center)
102     return cluster_centers
103
104
105 def updateMembershipValue(membership_mat, cluster_centers):
106     p = float(2/(m-1))
107     for i in range(n):
108         x = list(df.iloc[i])
109         distances = [np.linalg.norm(list(map(operator.sub, x, cluster_centers[j])))
110                     for j in range(k):
111                         den = sum([math.pow(float(distances[j]/distances[c]), p) for c in range(k)])
112                         membership_mat[i][j] = float(1/den)
113     return membership_mat
114
115
116 def getClusters(membership_mat):
117     cluster_labels = list()
118     for i in range(n):
119         max_val, idx = max((val, idx) for (idx, val) in enumerate(membership_mat[i]))
120         cluster_labels.append(idx)
121     return cluster_labels
122
123
124 def fuzzyCMeansClustering():
125     # Membership Matrix
126     membership_mat = initializeMembershipMatrix()
127     curr = 0
128     while curr <= MAX_ITER:
129         cluster_centers = calculateClusterCenter(membership_mat)
130         membership_mat = updateMembershipValue(membership_mat, cluster_centers)
131         cluster_labels = getClusters(membership_mat)
132         curr += 1
133     print(membership_mat)
134     return cluster_labels, cluster_centers
135
136
137 labels, centers = fuzzyCMeansClustering()
138 a,p,r = accuracy(labels, class_labels)
139
140 print("Accuracy = " + str(a))
141 print("Precision = " + str(p))
142 print("Recall = " + str(r))

```





```

1 import pandas as pd
2 import numpy as np
3 import random
4 import math
5
6 def eucledian_dist(a,b):
7     ans = float(np.sqrt(np.sum((a-b)**2)))
8     return ans
9
10 if __name__ == '__main__':
11     df = pd.read_csv("/content/drive/My Drive/Colab Notebooks/clustering/SPECTF_Ne
12     df = df.sample(frac=1)
13     #print df.head()
14     df = df.values
15     print (df.shape)
16     no_of_clusters = 2
17     m = 2
18     rows = random.sample(range(0,df.shape[0]),no_of_clusters)
19     centers = []
20     for i in range(no_of_clusters):
21         centers.append(df[rows[i],:df.shape[1]-1])
22     centers = np.array(centers,dtype=np.float)
23     #print( centers
24     membership_matrix = np.zeros((df.shape[0],no_of_clusters),dtype=np.float)
25
26     delta = 0.01
27     itr=0
28     flag=1
29
30     while itr<500 and flag==1:
31         print( "-----")
32         for i in range(membership_matrix.shape[0]):
33             for j in range(membership_matrix.shape[1]):
34                 numerator = eucledian_dist(df[i,:44],centers[j,:])
35                 if numerator==0:
36                     membership_matrix[i,j] = 1
37                     for k in range(j):
38                         membership_matrix[i,k] = 0
39                     break
40                 else:
41                     sumratios = 0
42                     for k in range(centers.shape[0]):
43                         dist = eucledian_dist(df[i,:44],centers[k,:])
44                         if dist==0:
45                             ratio = 0
46                         else:
47                             ratio = float(numerator)/dist
48                             ratio = math.pow(ratio,2.0/(m-1))
49                     sumratios += ratio
50                     membership_matrix[i,j] = 1/float(sumratios)
51

```

```
52     '''
53     for i in range(membership_matrix.shape[0]):
54         print( np.sum(membership_matrix[i])
55     '''
56     #print( centers
57     flag=0
58     for i in range(centers.shape[0]):
59         numerator=np.zeros(44)
60         denominator=0
61         for j in range(df.shape[0]):
62             #print( df[j,:44].shape
63             numerator = numerator + (df[j,:44] * math.pow(membership_matrix[j,
64             denominator += math.pow(membership_matrix[j,i],m)
65             if all(i<delta for i in centers[i]-(numerator/denominator)) is False:
66                 flag=1
67             centers[i] = numerator/denominator
68     itr+=1
69     print( membership_matrix)
70     print( itr)
71     print( "-----"
```

