

CS145 Homework 4

****Important Note:**** HW4 is due on **11:59 PM PT, Nov 20 (Friday, Week 7)**. Please submit through GradeScope.

Print Out Your Name and UID

****Name:** Kevin Li, **UID:** 405200619******

Before You Start

You need to first create HW4 conda environment by the given `cs145hw4.yml` file, which provides the name and necessary packages for this tasks. If you have `conda` properly installed, you may create, activate or deactivate by the following commands:

```
conda env create -f cs145hw4.yml
conda activate hw4
conda deactivate
```

OR

```
conda env create --name NAMEOFOURCHOICE -f cs145hw4.yml
conda activate NAMEOFOURCHOICE
conda deactivate
```

To view the list of your environments, use the following command:

```
conda env list
```

More useful information about managing environments can be found [here](#).

You may also quickly review the usage of basic Python and Numpy package, if needed in coding for matrix operations.

In this notebook, you must not delete any code cells in this notebook. If you change any code outside the blocks (such as some important hyperparameters) that you are allowed to edit (between START/END YOUR CODE HERE), you need to highlight these changes. You may add some additional cells to help explain your results and observations.

```
In [1]: import numpy as np
import pandas as pd
import sys
import random
import math
import matplotlib.pyplot as plt
from scipy.stats import multivariate_normal
```

```
%load_ext autoreload
%autoreload 2
```


If you can successfully run the code above, there will be no problem for environment setting.

1. Clustering Evaluation

This workbook will walk you through an example for calculating different clustering metrics.

Note: This is a "question-answer" style problem. You do not need to code anything and you are required to calculate by hand (with a scientific calculator).

Questions

Suppose we want to cluster the following 20 conferences into four areas, with ground truth label and algorithm output label shown in third and fourth column. Please evaluate the quality of the clustering algorithm according to four different metrics respectively. 

Questions (please include intermediate steps)

1. Calculate purity.
2. Calculate precision.
3. Calculate recall.
4. Calculate F1-score.
5. Calculate normalized mutual information.

Your answer here:

Note: you can use several code cells to help you compute the results and answer the questions. Again you don't need to do any coding.

Please type your answer here!

```
In [29]: Y = [3, 3, 1, 1, 1, 4, 3, 3, 4, 2, 4, 2, 1, 2, 3, 2, 1, 2, 4, 4]
Yhat = [2, 2, 3, 3, 3, 4, 2, 2, 3, 1, 4, 1, 3, 1, 2, 1, 2, 1, 4, 4]

for cluster in range(1, 5):
    for y, yhat in zip(Y, Yhat):
        if yhat == cluster:
            print(f'Cluster {cluster}. True label {y}')
    print()
```

```
Cluster 1. True label 2
Cluster 1. True label 2
Cluster 1. True label 2
Cluster 1. True label 2
Cluster 1. True label 2
```

```
Cluster 2. True label 3
Cluster 2. True label 3
Cluster 2. True label 3
Cluster 2. True label 3
Cluster 2. True label 3
Cluster 2. True label 1
```

```
Cluster 3. True label 1
Cluster 3. True label 1
```

Cluster 3. True label 1
 Cluster 3. True label 4
 Cluster 3. True label 1

Cluster 4. True label 4
 Cluster 4. True label 4
 Cluster 4. True label 4
 Cluster 4. True label 4

1

Purity is $\frac{1}{N} \sum_k \max_j |c_k \cap \omega_j|$

$N = 20$ since we have 20 points

For output cluster 1, we have true labels 2, 2, 2, 2 and 2 so all five are correct

For output cluster 2, we have true labels 3, 3, 3, 3, 3, 1 so five out of six are correct

For output cluster 3, we have true labels 1, 1, 1, 1, 4 so four out of five are correct

For output cluster 4, we have all true labels 4 so all four are correct

So, $\sum_k \max_j |c_k \cap \omega_j| = 5 + 5 + 4 + 4 = 18$

Purity is $\frac{18}{20} = 0.9$

```
In [31]: n = len(Y)
         tp, fp, fn, tn = 0, 0, 0, 0
         for i in range(n):
             for j in range(i + 1, n):
                 if Y[i] == Y[j]:
                     if Yhat[i] == Yhat[j]:
                         tp += 1
                     else:
                         fn += 1
                 else:
                     if Yhat[i] == Yhat[j]:
                         fp += 1
                     else:
                         tn += 1

         print(tp, fp, tn, fn)
```

32 9 141 8

2

$Precision = TP / (TP + FP)$

$TP = 32, FP = 9$

$Precision = 32/41 \approx 0.78$

3

$$Recall = TP / (TP + FN)$$

$$TP = 32 \quad FP = 8$$

$$Recall = 32/40 = 0.80$$

4

$$F_1 = 2 \cdot Precision \cdot Recall / (Precision + Recall)$$

$$F_1 = 2 \cdot .78 \cdot .8 / (.78 + .8)$$

$$F_1 \approx 0.79$$

5

NMI Table:

	C1	C2	C3	C4	Sum
L1	0	1	4	0	5
L2	5	0	0	0	5
L3	0	5	0	0	5
L4	0	0	1	4	5
Sum	5	6	5	4	20

We then apply the NMI formula

$$NMI(C, \Omega) = \frac{I(C, \Omega)}{\sqrt{H(C)H(\Omega)}}$$

To get $I(C, \Omega)$, we calculate the following:

$$\sum_k \sum_j \frac{|c_k \cap \omega_j|}{N} \ln \frac{N |c_k \cap \omega_j|}{|c_k| \cdot |\omega_j|}$$

To get $H(\Omega)$, we calculate the following:

$$\sum_j \frac{|\omega_j|}{N} \ln \frac{|\omega_j|}{N}$$

To get $H(C)$, we calculate the following:

$$\sum_k \frac{|c_k|}{N} \ln \frac{|c_k|}{N}$$

Once we calculate the values for I , H_Ω and H_C , we get the following:

$$I(C, \Omega) \approx 1.13$$

$$H(\Omega) = \frac{5}{20} \ln \frac{5}{20} + \frac{5}{20} \ln \frac{5}{20} + \frac{5}{20} \ln \frac{5}{20} + \frac{5}{20} \ln \frac{5}{20} \approx -1.39$$

$$H(\Omega) = \frac{5}{20} \ln \frac{5}{20} + \frac{6}{20} \ln \frac{6}{20} + \frac{5}{20} \ln \frac{5}{20} + \frac{4}{20} \ln \frac{4}{20} \approx -1.38$$

Once calculated, we combine them and get

$$NMI(C, \Omega) \approx 0.82$$

2. K-means

In this section, we are going to apply K-means algorithm against two datasets (dataset1.txt, dataset2.txt) with different distributions, respectively.

For each dataset, it contains 3 columns, with the format: x1 \t x2 \t cluster_label. You need to use the first two columns for clustering, and the last column for evaluation.

```
In [3]: from hw4code.KMeans import KMeans
        k = KMeans()
        # As a sanity check, we print out a sample of each dataset
        dataname1 = "data/dataset1.txt"
        dataname2 = "data/dataset2.txt"
        k.check_dataloader(dataname1)
        k.check_dataloader(dataname2)
```

```
For dataset1: number of datapoints is 150
      x      y  ground_truth_cluster
0 -0.163880 -0.219869                1
1 -0.886274 -0.356186                1
2 -0.978910 -0.893314                1
3 -0.658867 -0.371122                1
4 -0.072518  0.399157                1
```

```
For dataset2: number of datapoints is 200
      x      y  ground_truth_cluster
0  1.068587  0.136921                1
1  0.705440  0.393068                1
2  0.840811 -0.054906                1
3 -0.923447  0.598501                1
4  0.784353  0.724743                1
```

2.1 Coding K-means

Complete the `reassignClusters` and `getCentroid` function in `KMeans.py`.

Print out each output cluster's size and centroid (x,y) for dataset1 and dataset2 respectively.

```
In [16]: k = KMeans()
        #=====#
        # STRART YOUR CODE HERE #
        #=====#
        k.main(dataname1)
        k.main(dataname2)
        #=====#
        #  END YOUR CODE HERE  #
        #=====#
```

```

For dataset1
Iteration :3
Cluster 0 size :50
Centroid [x=2.5737264423871213, y=-0.027462568841232993]
Cluster 1 size :50
Centroid [x=-0.4633368646347212, y=-0.46611409698195794]
Cluster 2 size :50
Centroid [x=0.9888766205736857, y=2.010478965197201]

```

```

For dataset2
Iteration :4
Cluster 0 size :102
Centroid [x=1.2708406269481844, y=-0.08583389704900128]
Cluster 1 size :98
Centroid [x=-0.2018593506236788, y=0.5726963240559536]

```

2.2 Purity and NMI Evaluation

Complete the `compute_purity` function in `KMeans.py`.

In order to compute NMI, you need to firstly compute NMI matrix and then do the calculation. That is to complete the `getNMIMatrix` and `calcNMI` functions in `KMeans.py`.

Print out the purity and NMI for each dataset respectively.

```

In [33]: k = KMeans()
#=====#
# STRART YOUR CODE HERE #
#=====#
k.main(dataname1, isevaluate=True)
k.main(dataname2, isevaluate=True)
#=====#
# END YOUR CODE HERE #
#=====#

```

```

For dataset1
Iteration :3
Purity is 1.000000
NMI is 1.000000
Cluster 0 size :50
Centroid [x=2.5737264423871213, y=-0.027462568841232993]
Cluster 1 size :50
Centroid [x=-0.4633368646347212, y=-0.46611409698195794]
Cluster 2 size :50
Centroid [x=0.9888766205736857, y=2.010478965197201]

```

```

For dataset2
Iteration :4
Purity is 0.760000
NMI is 0.205096
Cluster 0 size :102
Centroid [x=1.2708406269481844, y=-0.08583389704900128]
Cluster 1 size :98
Centroid [x=-0.2018593506236788, y=0.5726963240559536]

```

2.3 Visualization

The clustering results for KMeans are saved as `KMeans_dataset1.csv` and

`KMeans_dataset2.csv` respectively under your root folder. Plot the clustering results for the two datasets, with different colors representing different clusters.

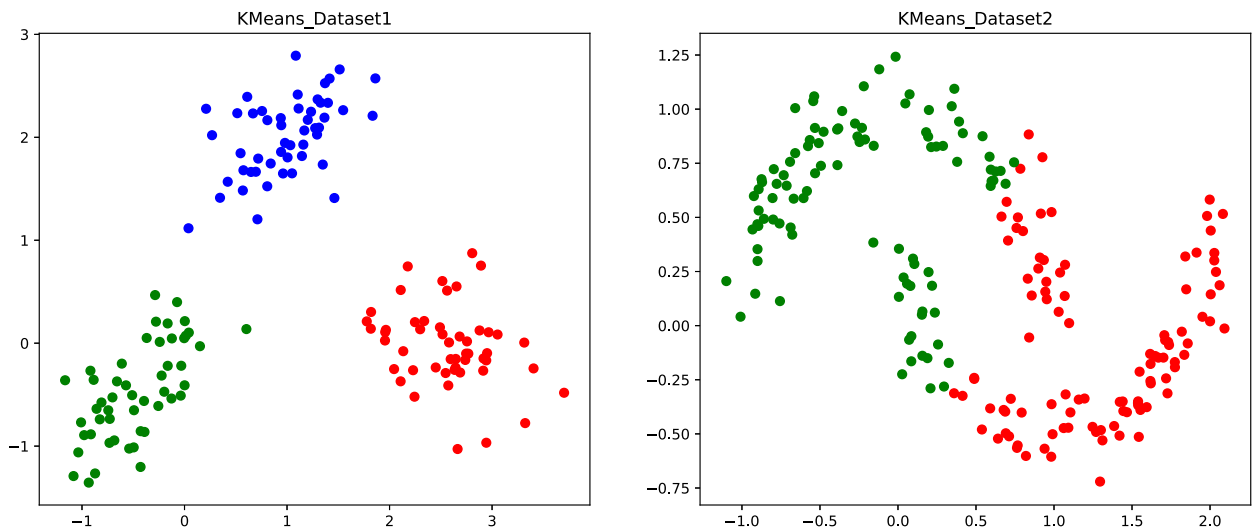
```
In [45]: CSV_FILE_PATH1 = 'Kmeans_dataset1.csv'
CSV_FILE_PATH2 = 'Kmeans_dataset2.csv'

df1 = pd.read_csv(CSV_FILE_PATH1,header=None,names=['x','y','pred'])
df2 = pd.read_csv(CSV_FILE_PATH2,header=None,names=['x','y','pred'])
fig, [ax0,ax1] = plt.subplots(1, 2, figsize=(15, 6))
ax0.title.set_text("KMeans_Dataset1")
ax1.title.set_text("KMeans_Dataset2")

#=====#
#  STRART YOUR CODE HERE  #
#=====#

colors = np.array(['r', 'g', 'b', 'orange'])
ax0.scatter(df1.x, df1.y, color=colors[df1.pred])
ax1.scatter(df2.x, df2.y, color=colors[df2.pred])

#=====#
#   END YOUR CODE HERE   #
#=====#
plt.show()
```



Question

Give the pros and cons of K-means algorithm. (At least one for pro and two for cons to get full marks)

Your answer here

[Please type your answer here!](#)

Pros

K-means is quick to train, taking $O(tkn)$ time where n is the number of points, k is the number of clusters, and t is the number of iterations. Generally, k and t are normally far less than n as well.

Additionally, K-means is simple to implement and understand

Cons

K-means can only be applied to continuous spaces (k-modes must be used for categorical data)

Very poor at discovering non-convex shapes (as we just showed!)

3 DBSCAN

In this section, we are going to use DBSCAN for clustering the same two datasets.

3.1 Coding DBSCAN

Complete the `dbscan` function in `DBSCAN.py`. Print out the purity, NMI and cluster size for each dataset respectively.

```
In [41]: from hw4code.DBSCAN import DBSCAN
d = DBSCAN()
#=====#
# STRART YOUR CODE HERE #
#=====#
d.main(dataname1)
d.main(dataname2)
#=====#
# END YOUR CODE HERE #
#=====#
```

```
For dataset1
Esp :0.3560832705047313
Number of clusters formed :4
Noise points :11
Purity is 0.940000
NMI is 0.959065
Cluster 0 size :49
Cluster 1 size :41
Cluster 2 size :47
Cluster 3 size :4
```

```
For dataset2
Esp :0.18652096476712493
Number of clusters formed :3
Noise points :3
Purity is 0.985000
NMI is 0.817349
Cluster 0 size :99
Cluster 1 size :51
Cluster 2 size :47
```

3.2 Visualization

The clustering results for DBSCAN are saved as `DBSCAN_dataset1.csv` and

`DBSCAN_dataset2.csv` respectively under your root folder. Plot the clustering results for the two datasets, with different colors representing different clusters.

```
In [44]: CSV_FILE_PATH1 = 'DBSCAN_dataset1.csv'
CSV_FILE_PATH2 = 'DBSCAN_dataset2.csv'

df1 = pd.read_csv(CSV_FILE_PATH1,header=None,names=['x','y','pred'])
df2 = pd.read_csv(CSV_FILE_PATH2,header=None,names=['x','y','pred'])
fig, [ax0,ax1] = plt.subplots(1, 2, figsize=(15, 6))
ax0.title.set_text("DBSCAN_Dataset1")
ax1.title.set_text("DBSCAN_Dataset2")
```



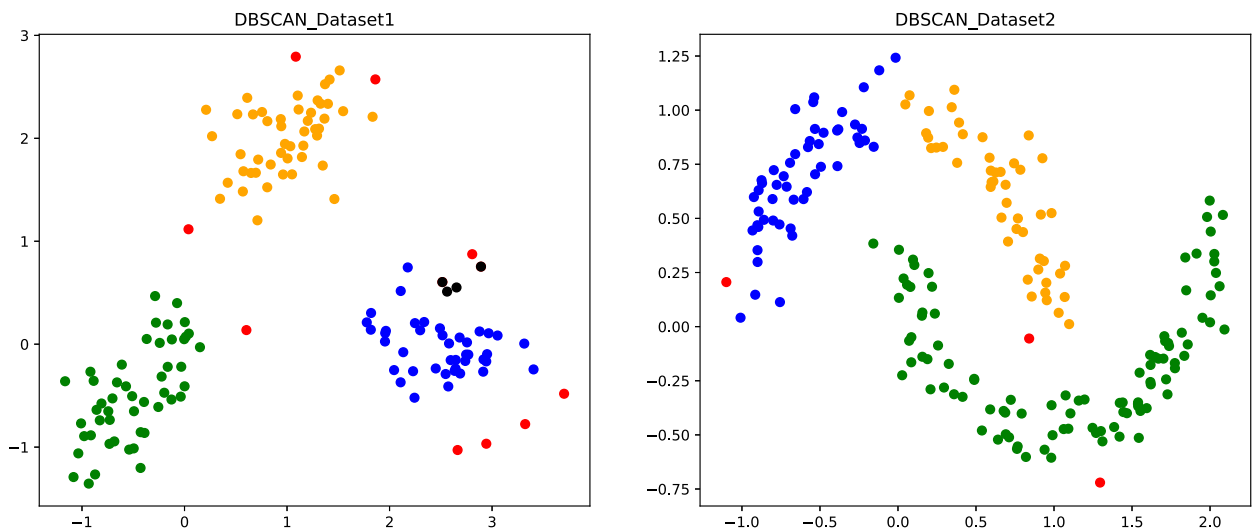
```

#=====#
#  START YOUR CODE HERE  #
#=====#

colors = np.array(['r', 'g', 'b', 'orange', 'black'])
ax0.scatter(df1.x, df1.y, color=colors[df1.pred])
ax1.scatter(df2.x, df2.y, color=colors[df2.pred])

#=====#
#  END YOUR CODE HERE  #
#=====#
plt.show()

```



Question

Give the pros and cons of DBSCAN algorithm. (At least two for pro and one for cons to get full marks)

Your answer here

Please type your answer here!

Pros

DBSCAN can perform well discovering non-convex shapes (as demonstrated above)

DBSCAN is also very resilient to noise

Cons

Since ϵ must be chosen, it may be very hard to find a reasonable value for it if one doesn't know the scale of the features

Additionally, ϵ and $minPts$ are static values; therefore, clusters with different densities cannot easily and correctly both be identified by DBSCAN.

One last issue that exists for any model that uses a distance metric such as Euclidean is the curse of dimensionality which essentially makes the distance between data points very precise (they almost all end up around the same values) making selection of ϵ virtually impossible

4 GMM

In this section, we are going to use GMM for clustering the same two datasets.

4.1 Coding GMM

Complete the `Estep` and `Mstep` function in `GMM.py`. Print out the purity, NMI, final mean, covariance and cluster size for each dataset respectively.

```
In [64]: from hw4code.GMM import GMM
g = GMM()
#=====#
# STRART YOUR CODE HERE #
#=====#
g.main(dataname1)
g.main(dataname2)
#=====#
# END YOUR CODE HERE #
#=====#
```

For dataset1
Number of Iterations = 22

After Calculations
Final mean =
-0.46247285694404044
-0.4638749980764899

0.9898929396029765
2.011802723814242

2.57342634413319
-0.027108746076609493

Final covariance =
For Cluster : 1
0.14918910487220216
0.1173463005433889

0.1173463005433889
0.215548612531075

For Cluster : 2
0.16028233507625483
0.07486967581052754

0.07486967581052754
0.13939774162738802

For Cluster : 3
0.18039223672749394
-0.04672614559811056

-0.04672614559811056

```
0.15206459963738586
```

```
Purity is 1.000000
NMI is 1.000000
Cluster 0 size :50
Cluster 1 size :50
Cluster 2 size :50
```

```
For dataset2
Number of Iterations = 95
```

```
After Calculations
Final mean =
0.7464905663922625
0.45649665848541027
```

```
0.28287851889390897
-0.05970560727188742
```

```
Final covariance =
For Cluster : 1
0.7692790765358334
-0.28782809642382123
```

```
-0.28782809642382123
0.1901249384356512
```

```
For Cluster : 2
0.6828574757628687
-0.30058915994390495
```

```
-0.30058915994390495
0.17583559485120057
```

```
Purity is 0.690000
NMI is 0.107406
Cluster 0 size :106
Cluster 1 size :94
```

4.2 Visualization

The clustering results for GMM are saved as `GMM_dataset1.csv` and `GMM_dataset2.csv` respectively under your root folder. Plot the clustering results for the two datasets, with different colors representing different clusters.

```
In [65]: CSV_FILE_PATH1 = 'GMM_dataset1.csv'
          CSV_FILE_PATH2 = 'GMM_dataset2.csv'

          df1 = pd.read_csv(CSV_FILE_PATH1,header=None,names=['x','y','pred'])
          df2 = pd.read_csv(CSV_FILE_PATH2,header=None,names=['x','y','pred'])
          fig, [ax0,ax1] = plt.subplots(1, 2, figsize=(15, 6))
          ax0.title.set_text("GMM_Dataset1")
          ax1.title.set_text("GMM_Dataset2")

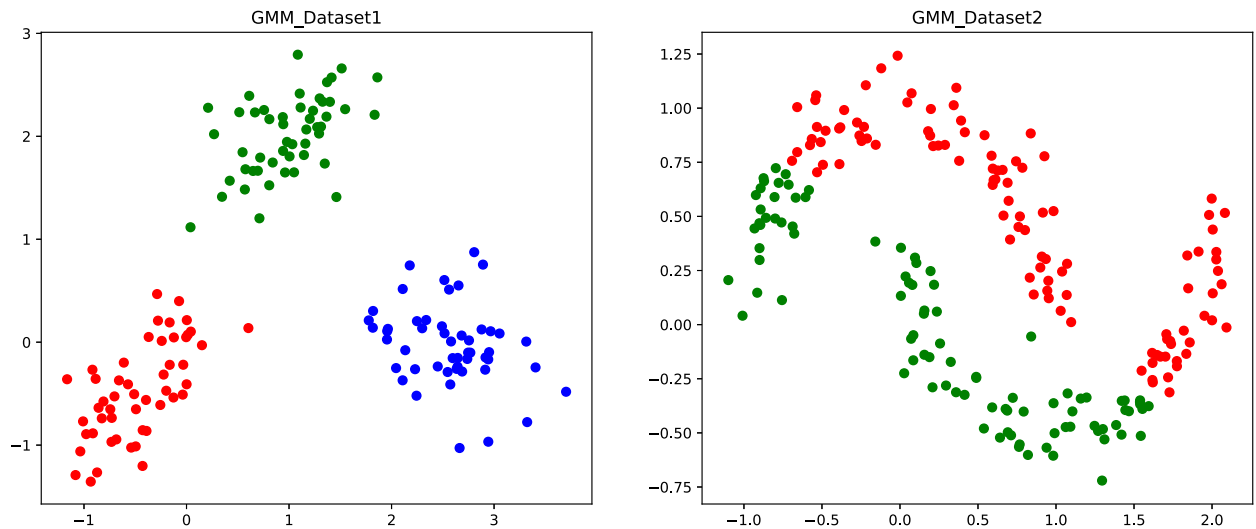
          #=====#
          # START YOUR CODE HERE #
          #=====#
```

```

colors = np.array(['r', 'g', 'b', 'orange', 'black'])
ax0.scatter(df1.x, df1.y, color=colors[df1.pred])
ax1.scatter(df2.x, df2.y, color=colors[df2.pred])

#=====#
#  END YOUR CODE HERE  #
#=====#
plt.show()

```



Questions

1. Give the pros and cons of GMM algorithm. (At least two for pro and two for cons to get full marks)
2. Compare the visualization results from three algorithms, analyze for each dataset why these algorithms would produce such result.

Your answer here:

[Please type your answer here!](#)

Pros of GMM

One pro of mixture models is that they can generalize different densities and different cluster sizes simultaneously

Another pro is that is is a simple model with only a few (and fairly explainable) parameters

Cons of GMM

GMM can only create round circular/ovalish shapes.

It is also computationally expensive to train if there are many distributions

Reasoning over dataset1

All three models do fairly well. This makes sense, because all three clusters are consistently dense (good for DBSCAN), have similarly sized clusters (good for K-means), and have roughly circular/oval shaped clusters (good for K-means and GMM)

Reasoning over dataset2

We see GMM and K-means doing pretty poorly while DBSCAN performs excellently. This makes sense, because the shapes are non-convex, oblong, non-circular shapes. Thus K-means and GMM will struggle significantly to cluster them while DBSCAN can do a good job since the data has a nice, consistent density in the clusters.

5 Bonus Question

Prove that KMeans algorithm would guarantee convergence. (**Hint: prove for each step the loss would decrease.**)

In the first step, we fix the centers and find the assignment of w_{ij} that minimizes J . If there is no better assignment, we leave the assignments as they are.

We note that the only way J changes is if the cluster label of a point changes. If it does, then we know J will decrease because we only change if we can pick an assignment for which the J_i for that point i decreases (and then we pick the label that maximally decreases J_i).

In the second step, we fix the assignments of the points w_{ij} , but then find centers for the clusters that minimize J . We note that the solution to this is closed; we can simply set $\frac{\partial J}{\partial c_j} = 0$. So, we will prove now that when doing this, we end up with the solution we run (where we recenter the clusters based on their center of mass). Thus, we're minimizing J and therefore either decreasing it or leaving it the same.

$$J = \sum_{j=1}^k \sum_i w_{ij} \|x_i - c_j\|^2$$

$$\frac{\partial J}{\partial c_j} = -2 \sum_i w_{ij} (x_i - c_j)$$

Setting the derivative equal to zero, we get

$$0 = -2 \sum_i w_{ij} (x_i - c_j)$$

$$c_j = \frac{\sum_i w_{ij} x_i}{\sum_i w_{ij}}$$

This is precisely what we do; therefore in step two we minimize J again and monotonically decrease it.

End of Homework 4 :)

After you've finished the homework, please print out the entire `ipynb` notebook and four `py` files into one PDF file. Make sure you include the output of code cells and answers for questions. Prepare submit it to GradeScope. Also this time remember assign the pages to the questions on GradeScope

```

1 from hw4code.DataPoints import DataPoints
2 import random
3 import sys
4 import math
5 import pandas as pd
6 import numpy as np
7
8 # =====
9 def sqrt(n):
10     return math.sqrt(n)
11
12 # =====
13 def getEuclideanDist(x1, y1, x2, y2):
14     dist = sqrt(pow((x2 - x1), 2) + pow((y2 - y1), 2))
15     return dist
16 # =====
17 def compute_purity(clusters, total_points):
18     # Calculate purity
19
20     # Create list to store the maximum union number for each output cluster.
21     maxLabelCluster = []
22     num_clusters = len(clusters)
23     # =====#
24     # STRART YOUR CODE HERE #
25     # =====#
26     for cluster in clusters:
27         vals, counts = np.unique([dp.label for dp in cluster], return_counts=True)
28         maxLabelCluster.append(max(counts))
29     # =====#
30     # END YOUR CODE HERE #
31     # =====#
32     purity = 0.0
33     for j in range(num_clusters):
34         purity += maxLabelCluster[j]
35     purity /= total_points
36     print("Purity is %.6f" % purity)
37
38 # =====
39 def compute_NMI(clusters, noOfLabels):
40     # Get the NMI matrix first
41     nmiMatrix = getNMIMatrix(clusters, noOfLabels)
42     # Get the NMI matrix first
43     nmi = calcNMI(nmiMatrix)
44     print("NMI is %.6f" % nmi)
45
46
47 # =====
48 def getNMIMatrix(clusters, noOfLabels):
49     # Matrix shape of [num_true_clusters + 1, num_output_clusters + 1] (example under
50     week6's slide page 9)
51     nmiMatrix = [[0 for x in range(len(clusters) + 1)] for y in range(noOfLabels +
52     1)]
53     clusterNo = 0
54     for cluster in clusters:
55         # Create dictionary {true_class_No: Number of shared elements}
56         labelCounts = {}
57         # =====#
58         # STRART YOUR CODE HERE #
59         # =====#
60         classes, counts = np.unique([dp.label for dp in cluster], return_counts=True)

```

```

59     labelCounts = dict(zip(classes, counts))
60     # =====#
61     #   END YOUR CODE HERE   #
62     # =====#
63     labelTotal = 0
64     labelCounts_sorted = sorted(labelCounts.items(), key=lambda item: item[1],
reverse=True)
65     for label, val in labelCounts_sorted:
66         nmiMatrix[label - 1][clusterNo] = labelCounts[label]
67         labelTotal += labelCounts.get(label)
68     # Populate last row (row of summation)
69     nmiMatrix[noOfLabels][clusterNo] = labelTotal
70     clusterNo += 1
71     labelCounts.clear()
72
73     # Populate last col (col of summation)
74     lastRowCol = 0
75     for i in range(noOfLabels):
76         totalRow = 0
77         for j in range(len(clusters)):
78             totalRow += nmiMatrix[i][j]
79         lastRowCol += totalRow
80         nmiMatrix[i][len(clusters)] = totalRow
81
82     # Total number of datapoints
83     nmiMatrix[noOfLabels][len(clusters)] = lastRowCol
84
85     return nmiMatrix
86
87     # =====
88 def calcNMI(nmiMatrix):
89     # Num of true clusters + 1
90     row = len(nmiMatrix)
91     # Num of output clusters + 1
92     col = len(nmiMatrix[0])
93     # Total number of datapoints
94     N = nmiMatrix[row - 1][col - 1]
95     I = 0.0
96     HOmega = 0.0
97     HC = 0.0
98
99     for i in range(row - 1):
100         for j in range(col - 1):
101             # Compute the log part of each pair of clusters within I's formula.
102             logPart_I = 1.0
103             # =====#
104             #   STRART YOUR CODE HERE   #
105             # =====#
106             logPart_I = N * nmiMatrix[i][j] / (nmiMatrix[-1][j] * nmiMatrix[i][-1])
107             # =====#
108             #   END YOUR CODE HERE   #
109             # =====#
110
111             if logPart_I == 0.0:
112                 continue
113             I += (nmiMatrix[i][j] / float(N)) * math.log(float(logPart_I))
114     # Compute HOmega
115     # =====#
116     #   STRART YOUR CODE HERE   #
117     # =====#

```



```

118     HOmega += nmiMatrix[i][-1] / N * np.log(nmiMatrix[i][-1] / N)
119     # =====#
120     #   END YOUR CODE HERE   #
121     # =====#
122
123     #Compute HC
124     # =====#
125     # STRART YOUR CODE HERE  #
126     # =====#
127     c_row = np.array(nmiMatrix[-1][: -1])
128     HC = np.sum((c_row / N) * np.log(c_row/N))
129     # =====#
130     #   END YOUR CODE HERE   #
131     # =====#
132
133     return I / math.sqrt(HC * HOmega)
134
135
136
137
138
139 # =====#
140 class Centroid:
141     # -----#
142     def __init__(self, x, y):
143         self.x = x
144         self.y = y
145     # -----#
146     def __eq__(self, other):
147         if not type(other) is type(self):
148             return False
149         if other is self:
150             return True
151         if other is None:
152             return False
153         if self.x != other.x:
154             return False
155         if self.y != other.y:
156             return False
157         return True
158     # -----#
159     def __ne__(self, other):
160         result = self.__eq__(other)
161         if result is NotImplemented:
162             return result
163         return not result
164     # -----#
165     def toString(self):
166         return "Centroid [x=" + str(self.x) + ", y=" + str(self.y) + "]"
167     # -----#
168     def __str__(self):
169         return self.toString()
170     # -----#
171     def __repr__(self):
172         return self.toString()
173
174
175
176
177

```

```

178
179
180 # =====
181 class KMeans:
182     # -----
183     def __init__(self):
184         self.K = 0
185     # -----
186     def main(self, dataname, isevaluate=False):
187         seed = 71
188         self.dataname = dataname[5:-4]
189         print("\nFor " + self.dataname)
190         self.dataSet = self.readDataSet(dataname)
191         self.K = DataPoints.getNoOFLabels(self.dataSet)
192         random.Random(seed).shuffle(self.dataSet)
193         self.kmeans(isevaluate)
194
195     # -----
196     def check_dataloader(self, dataname):
197
198         df = pd.read_table(dataname, sep = "\t", header=None, names=
199 ['x', 'y', 'ground_truth_cluster'])
200         print("\nFor " + dataname[5:-4] + ": number of datapoints is %d" %
201 df.shape[0])
202         print(df.head(5))
203
204     # -----
205     def kmeans(self, isevaluate=False):
206         clusters = []
207         k = 0
208         while k < self.K:
209             cluster = set()
210             clusters.append(cluster)
211             k += 1
212
213         # Initially randomly assign points to clusters
214         i = 0
215         for point in self.dataSet:
216             clusters[i % k].add(point)
217             i += 1
218
219         # calculate centroid for clusters
220         centroids = []
221         for j in range(self.K):
222             centroids.append(self.getCentroid(clusters[j]))
223
224         self.reassignClusters(self.dataSet, centroids, clusters)
225
226         # continue till converge
227         iteration = 0
228         while True:
229             iteration += 1
230             # calculate centroid for clusters
231             centroidsNew = []
232             for j in range(self.K):
233                 centroidsNew.append(self.getCentroid(clusters[j]))
234
235             isConverge = False
236             for j in range(self.K):

```

```

236         if centroidsNew[j] != centroids[j]:
237             isConverge = False
238         else:
239             isConverge = True
240     if isConverge:
241         break
242
243     for j in range(self.K):
244         clusters[j] = set()
245
246     self.reassignClusters(self.dataSet, centroidsNew, clusters)
247     for j in range(self.K):
248         centroids[j] = centroidsNew[j]
249     print("Iteration :" + str(iteration))
250
251     if isevaluate:
252         # Calculate purity and NMI
253         compute_purity(clusters, len(self.dataSet))
254         compute_NMI(clusters, self.K)
255
256     # write clusters to file for plotting
257     f = open("Kmeans_" + self.dataname + ".csv", "w")
258     for w in range(self.K):
259         print("Cluster " + str(w) + " size :" + str(len(clusters[w])))
260         print(centroids[w].toString())
261         for point in clusters[w]:
262             f.write(str(point.x) + "," + str(point.y) + "," + str(w) + "\n")
263     f.close()
264
265     # -----
266     def reassignClusters(self, dataSet, c, clusters):
267         # reassign points based on cluster and continue till stable clusters found
268         dist = [0.0 for x in range(self.K)]
269         for point in dataSet:
270             for i in range(self.K):
271                 dist[i] = getEuclideanDist(point.x, point.y, c[i].x, c[i].y)
272
273             minIndex = self.getMin(dist)
274             # assign point to the closest cluster
275             # =====#
276             # STRART YOUR CODE HERE #
277             # =====#
278
279             for cluster in clusters:
280                 try:
281                     cluster.remove(point)
282                 except KeyError:
283                     pass
284
285             clusters[minIndex].add(point)
286
287             # =====#
288             # END YOUR CODE HERE #
289             # =====#
290     # -----
291     def getMin(self, dist):
292         min = sys.maxsize
293         minIndex = -1
294         for i in range(len(dist)):
295             if dist[i] < min:

```

```
296         min = dist[i]
297         minIndex = i
298     return minIndex
299
300     # -----
301 def getCentroid(self, cluster):
302     # mean of x and mean of y
303     cx = 0
304     cy = 0
305     # =====#
306     # STRART YOUR CODE HERE #
307     # =====#
308     for dp in cluster:
309         cx += dp.x
310         cy += dp.y
311
312     cx /= len(cluster)
313     cy /= len(cluster)
314     # =====#
315     # END YOUR CODE HERE #
316     # =====#
317     return Centroid(cx, cy)
318     # -----
319 @staticmethod
320 def readDataSet(filePath):
321     dataSet = []
322     with open(filePath) as f:
323         lines = f.readlines()
324     lines = [x.strip() for x in lines]
325     for line in lines:
326         points = line.split('\t')
327         x = float(points[0])
328         y = float(points[1])
329         label = int(points[2])
330         point = DataPoints(x, y, label)
331         dataSet.append(point)
332     return dataSet
333
```

```
1 from hw4code.KMeans import KMeans, compute_purity, compute_NMI, getEuclideanDist
2 from hw4code.DataPoints import DataPoints
3 import random
4
5
6 class DBSCAN:
7     # -----
8     def __init__(self):
9         self.e = 0.0
10        self.minPts = 3
11        self.noOfLabels = 0
12    # -----
13    def main(self, dataname):
14        seed = 71
15
16        self.dataname = dataname[5:-4]
17        print("\nFor " + self.dataname)
18        self.dataSet = KMeans.readDataSet(dataname)
19        random.Random(seed).shuffle(self.dataSet)
20        self.noOfLabels = DataPoints.getNoOfLabels(self.dataSet)
21        self.e = self.getEpsilon(self.dataSet)
22        print("Esp :" + str(self.e))
23        self.dbscan(self.dataSet)
24
25
26    # -----
27    def getEpsilon(self, dataSet):
28        distances = []
29        sumOfDist = 0.0
30        for i in range(len(dataSet)):
31            point = dataSet[i]
32            for j in range(len(dataSet)):
33                if i == j:
34                    continue
35                pt = dataSet[j]
36                dist = getEuclideanDist(point.x, point.y, pt.x, pt.y)
37                distances.append(dist)
38
39            distances.sort()
40            sumOfDist += distances[7]
41            distances = []
42        return sumOfDist/len(dataSet)
43    # -----
44    def dbscan(self, dataSet):
45        clusters = []
46        visited = set()
47        noise = set()
48
49        # Iterate over data points
50        for i in range(len(dataSet)):
51            point = dataSet[i]
52            if point in visited:
53                continue
54            visited.add(point)
55            N = []
56            minPtsNeighbours = 0
57
58            # check which point satisfies minPts condition
59            for j in range(len(dataSet)):
60                if i==j:
```

```

61         continue
62     pt = dataSet[j]
63     dist = getEuclideanDist(point.x, point.y, pt.x, pt.y)
64     if dist <= self.e:
65         minPtsNeighbours += 1
66         N.append(pt)
67
68     if minPtsNeighbours >= self.minPts:
69         cluster = set()
70         cluster.add(point)
71         point.isAssignedToCluster = True
72
73         j = 0
74         while j < len(N):
75             point1 = N[j]
76             minPtsNeighbours1 = 0
77             N1 = []
78             if not point1 in visited:
79                 visited.add(point1)
80                 for l in range(len(dataSet)):
81                     pt = dataSet[l]
82                     dist = getEuclideanDist(point1.x, point1.y, pt.x, pt.y)
83                     if dist <= self.e:
84                         minPtsNeighbours1 += 1
85                         N1.append(pt)
86                 if minPtsNeighbours1 >= self.minPts:
87                     self.removeDuplicates(N, N1)
88
89             # Add point1 is not yet member of any other cluster then add it
to cluster
90             # Hint: use self.isAssignedToCluster function to check if a point
is assigned to any clusters
91             # =====#
92             # STRART YOUR CODE HERE #
93             # =====#
94             if not point1.isAssignedToCluster:
95                 cluster.add(point1)
96                 point1.isAssignedToCluster = True
97             # =====#
98             # END YOUR CODE HERE #
99             # =====#
100             j += 1
101
102             # add cluster to the list of clusters
103             clusters.append(cluster)
104
105         else:
106             noise.add(point)
107
108
109     # List clusters
110     print("Number of clusters formed :" + str(len(clusters)))
111     print("Noise points :" + str(len(noise)))
112
113     # Calculate purity
114     compute_purity(clusters, len(self.dataSet))
115     compute_NMI(clusters, self.noOfLabels)
116     DataPoints.writeToFile(noise, clusters, "DBSCAN_" + self.dataname + ".csv")
117     # -----
118     def removeDuplicates(self, n, n1):

```

```
119     for point in n1:
120         isDup = False
121         for point1 in n:
122             if point1 == point:
123                 isDup = True
124                 break
125         if not isDup:
126             n.append(point)
127
128
```

```

1 from hw4code.DataPoints import DataPoints
2 from hw4code.KMeans import KMeans, compute_purity, compute_NMI
3 import math
4 from scipy.stats import multivariate_normal
5
6 # =====
7 class GMM:
8     # -----
9     def __init__(self):
10         self.dataSet = []
11         self.K = 0
12         self.mean = [[0.0 for x in range(2)] for y in range(3)]
13         self.stdDev = [[0.0 for x in range(2)] for y in range(3)]
14         self.coVariance = [[[0.0 for x in range(2)] for y in range(2)] for z in
range(3)]
15         self.W = None
16         self.w = None
17     # -----
18     def main(self, dataname):
19
20         self.dataname = dataname[5:-4]
21         print("\nFor " + self.dataname)
22         self.dataSet = KMeans.readDataSet(dataname)
23         self.K = DataPoints.getNoOfLabels(self.dataSet)
24         # weight for pair of data and cluster
25         self.W = [[0.0 for y in range(self.K)] for x in range(len(self.dataSet))]
26         # weight for pair of data and cluster
27         self.w = [0.0 for x in range(self.K)]
28         self.GMM()
29
30     # -----
31     def GMM(self):
32         clusters = []
33         # [num_clusters,2]
34         self.mean = [[0.0 for y in range(2)] for x in range(self.K)]
35         # [num_clusters,2]
36         self.stdDev = [[0.0 for y in range(2)] for x in range(self.K)]
37         # [num_clusters,2]
38         self.coVariance = [[[0.0 for z in range(2)] for y in range(2)] for x in
range(self.K)]
39         k = 0
40         while k < self.K:
41             cluster = set()
42             clusters.append(cluster)
43             k += 1
44
45         # Initially randomly assign points to clusters
46         i = 0
47         for point in self.dataSet:
48             clusters[i % self.K].add(point)
49             i += 1
50
51         # Initially assign equal prior weight for each cluster
52         for m in range(self.K):
53             self.w[m] = 1.0 / self.K
54
55         # Get Initial mean, std, covariance matrix
56         DataPoints.getMean(clusters, self.mean)
57         DataPoints.getStdDeviation(clusters, self.mean, self.stdDev)
58         DataPoints.getCovariance(clusters, self.mean, self.stdDev, self.coVariance)

```



```

59
60     length = 0
61     while True:
62         mle_old = self.Likelihood()
63         self.Estep()
64         self.Mstep()
65         length += 1
66         mle_new = self.Likelihood()
67
68         # convergence condition
69         if abs(mle_new - mle_old) / abs(mle_old) < 0.000001:
70             break
71
72     print("Number of Iterations = " + str(length))
73     print("\nAfter Calculations")
74     print("Final mean = ")
75     self.printArray(self.mean)
76     print("\nFinal covariance = ")
77     self.print3D(self.coVariance)
78
79     # Assign points to cluster depending on max prob.
80     for j in range(self.K):
81         clusters[j] = set()
82
83     i = 0
84     for point in self.dataSet:
85         index = -1
86         prob = 0.0
87         for j in range(self.K):
88             if self.W[i][j] > prob:
89                 index = j
90                 prob = self.W[i][j]
91         temp = clusters[index]
92         temp.add(point)
93         i += 1
94
95     # Calculate purity and NMI
96     compute_purity(clusters, len(self.dataSet))
97     compute_NMI(clusters, self.K)
98
99     # write clusters to file for plotting
100    f = open("GMM_" + self.dataname + ".csv", "w")
101    for w in range(self.K):
102        print("Cluster " + str(w) + " size :" + str(len(clusters[w])))
103        for point in clusters[w]:
104            f.write(str(point.x) + "," + str(point.y) + "," + str(w) + "\n")
105    f.close()
106    # -----
107    def Estep(self):
108        # Update self.W
109        for i in range(len(self.dataSet)):
110            denominator = 0.0
111            for j in range(self.K):
112                gaussian = multivariate_normal(self.mean[j], self.coVariance[j])
113                # Compute numerator for self.W[i][j] below
114                numerator = 0.0
115                # =====#
116                # STRART YOUR CODE HERE #
117                # =====#
118

```

```

119         pt = self.dataSet[i]
120         numerator = self.w[j] * gaussian.pdf((pt.x, pt.y))
121
122         # =====#
123         #   END YOUR CODE HERE   #
124         # =====#
125         self.W[i][j] = numerator
126         denominator += numerator
127
128         # normalize W[i][j] into probabilities
129         # =====#
130         # STRART YOUR CODE HERE #
131         # =====#
132
133         row_weight = sum(self.W[i])
134         for j in range(self.K):
135             self.W[i][j] /= row_weight
136
137         # =====#
138         #   END YOUR CODE HERE   #
139         # =====#
140     # -----#
141     def Mstep(self):
142         for j in range(self.K):
143             denominator = 0.0
144             numerator_x = 0.0
145             numerator_y = 0.0
146             cov_xy = 0.0
147             updatedMean_x = 0.0
148             updatedMean_y = 0.0
149
150             # update self.w[j] and self.mean
151             for i in range(len(self.dataSet)):
152                 denominator += self.W[i][j]
153                 updatedMean_x += self.W[i][j] * self.dataSet[i].x
154                 updatedMean_y += self.W[i][j] * self.dataSet[i].y
155
156             self.w[j] = denominator / len(self.dataSet)
157
158             #update self.mean
159             # =====#
160             # STRART YOUR CODE HERE #
161             # =====#
162
163             self.mean[j] = [0, 0]
164
165             den = 0
166             for i, dp in enumerate(self.dataSet):
167                 self.mean[j][0] += self.W[i][j] * dp.x
168                 self.mean[j][1] += self.W[i][j] * dp.y
169
170                 den += self.W[i][j]
171
172             self.mean[j][0] /= den
173             self.mean[j][1] /= den
174
175             denominator = 0
176
177             # =====#
178             #   END YOUR CODE HERE   #

```

```

179         # =====#
180
181         # update covariance matrix
182         for i in range(len(self.dataSet)):
183             numerator_x += self.W[i][j] * pow((self.dataSet[i].x - self.mean[j]
[0]), 2)
184             numerator_y += self.W[i][j] * pow((self.dataSet[i].y - self.mean[j]
[1]), 2)
185
186             # Compute conv_xy +=?
187             # =====#
188             # STRART YOUR CODE HERE #
189             # =====#
190
191             cov_xy += self.W[i][j] * (self.dataSet[i].x - self.mean[j][0]) *
(self.dataSet[i].y - self.mean[j][1])
192             denominator += self.W[i][j]
193
194             # =====#
195             # END YOUR CODE HERE #
196             # =====#
197
198             self.stdDev[j][0] = numerator_x / denominator
199             self.stdDev[j][1] = numerator_y / denominator
200
201             self.coVariance[j][0][0] = self.stdDev[j][0]
202             self.coVariance[j][1][1] = self.stdDev[j][1]
203             self.coVariance[j][0][1] = self.coVariance[j][1][0] = cov_xy /
denominator
204         # -----
205         def Likelihood(self):
206             likelihood = 0.0
207             for i in range(len(self.dataSet)):
208                 numerator = 0.0
209                 for j in range(self.K):
210                     gaussian = multivariate_normal(self.mean[j], self.coVariance[j])
211                     numerator += self.w[j] * gaussian.pdf([self.dataSet[i].x,
self.dataSet[i].y])
212                 likelihood += math.log(numerator)
213             return likelihood
214         # -----
215         def printArray(self, mat):
216             for i in range(len(mat)):
217                 for j in range(len(mat[i])):
218                     print(str(mat[i][j]) + " "),
219                 print("")
220         # -----
221         def print3D(self, mat):
222             for i in range(len(mat)):
223                 print("For Cluster : " + str((i + 1)))
224                 for j in range(len(mat[i])):
225                     for k in range(len(mat[i][j])):
226                         print(str(mat[i][j][k]) + " "),
227                     print("")
228                 print("")
229
230         # =====
231         if __name__ == "__main__":
232             g = GMM()
233             dataname = "dataset1.txt"

```

```
234 | g.main(dataname)
```