

# 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: Yunong Ye, UID: 004757414

## 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](https://docs.conda.io/projects/conda/en/latest/user-guide/tasks/manage-environments.html) (<https://docs.conda.io/projects/conda/en/latest/user-guide/tasks/manage-environments.html>).

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 [25]: 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
```

The autoreload extension is already loaded. To reload it, use:  
%reload\_ext autoreload

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.

ID	Conference Name	Ground Truth Label	Algorithm output Label
1	IJCAI	3	2
2	AAAI	3	2
3	ICDE	1	3
4	VLDB	1	3
5	SIGMOD	1	3
6	SIGIR	4	4
7	ICML	3	2
8	NIPS	3	2
9	CIKM	4	3
10	KDD	2	1
11	WWW	4	4
12	PAKDD	2	1
13	PODS	1	3
14	ICDM	2	1
15	ECML	3	2
16	PKDD	2	1
17	EDBT	1	2
18	SDM	2	1
19	ECIR	4	4
20	WSDM	4	4

### Questions (please include intermediate steps)

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

```

In [21]: ground_truth = [3,3,1,1,1,4,3,3,4,2,4,2,1,2,3,2,1,2,4,4]
output = [2,2,3,3,3,4,2,2,3,1,4,1,3,1,2,1,2,1,4,4]

def info(truth, ouptut):
    result = {'count':0, 'tp':0, 'fp':0, 'tn':0, 'fn':0}
    for i in range(len(truth)):
        for j in range(len(output)):
            # avoid repetitive pair
            if i < j:
                if truth[i] == truth[j]:
                    if output[i] == output[j]:
                        result['tp'] += 1
                    else:
                        result['fn'] += 1
                else:
                    if output[i] == output[j]:
                        result['fp'] += 1
                    else:
                        result['tn'] += 1
            result['count'] += 1
    return result

def table(truth, output):
    result = dict()
    truth_table = dict()
    output_table = dict()
    for i in range(len(truth)):
        if truth[i] not in truth_table:
            truth_table[truth[i]] = list()
        truth_table[truth[i]].append(i)
    for j in range(len(output)):
        if output[j] not in output_table:
            output_table[output[j]] = list()
        output_table[output[j]].append(j)
    result['truth'] = truth_table
    result['output'] = output_table
    return result

def nmi(table, total):
    # I
    i = 0.
    for truth_label in table['truth'].keys():
        for output_label in table['output'].keys():
            truth_values = table['truth'][truth_label]
            output_values = table['output'][output_label]
            common = float(len(list(set(truth_values).intersection(output_values))))
            cur = common / total * np.log(total*common/len(truth_values)/len(output_values))
            i += cur

    # calculate entropy
    h_truth = 0.
    for truth_label in table['truth'].keys():
        truth_length = float(len(table['truth'][truth_label]))
        h_truth -= (truth_length/total*np.log(truth_length/total))

```

```

h_output = 0.
for output_label in table['output'].keys():
    output_length = float(len(table['output'][output_label]))
    h_output -= (output_length/total*np.log(output_length/total))

nmi = i/np.sqrt(h_truth*h_output)
return nmi

table = table(ground_truth, output)
info = info(ground_truth,output)
print(info)
print(table)
print(nmi(table, len(output)))

{'count': 190, 'tp': 32, 'fp': 9, 'tn': 141, 'fn': 8}
{'truth': {3: [0, 1, 6, 7, 14], 1: [2, 3, 4, 12, 16], 4: [5, 8, 10, 18, 19], 2: [9, 11, 13, 15, 17]}, 'output': {2: [0, 1, 6, 7, 14, 16], 3: [2, 3, 4, 8, 12], 4: [5, 10, 18, 19], 1: [9, 11, 13, 15, 17]}}
0.8152212305376372

```

### **Purity**

output cluster 1: 10,12,14,16,18 ----> 5 matched with ground truth 2

output cluster 2: 1,2,7,8,15,17 ----> 5 matched with ground truth 3

output cluster 3: 3,4,5,9,13 ----> 4 matched with ground truth 1

output cluster 4: 6,11,19,20 ----> 4 matched with ground truth 4

$$\text{Purity} = \frac{1}{N} \sum_k \max |c_k \cap \omega_j| = 18/20 = 0.9$$

### **Precision**

$$\text{Precision} = \text{TP}/(\text{TP}+\text{FP}) = 32/(32+9) = 32/41 = 0.78049$$

### **Recall**

$$\text{Recall} = \text{TP}/(\text{TP}+\text{FN}) = 32/(32+8) = 32/40 = 0.8$$

### **F1-score**

$$\text{F1-score} = 2\text{precisionrecall} / (\text{precision}+\text{recall}) = 0.79012$$

**NMI** We used script above to calculate NMI = 0.81522

**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!

answer 1: Purity is 0.9

answer 2: Precision is 0.7809

answer 3: Recall is 0.8

answer 4: F1-score is 0.79012

answer 5: NMI is 0.81522

## 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 [32]: 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_data_loader(dataname1)
k.check_data_loader(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 [33]: k = KMeans()
#=====#
# STRART YOUR CODE HERE #
#=====#
k.main(dataname1)
k.main(dataname2)
#=====#
#  END YOUR CODE HERE  #
#=====#
```

```
For dataset1
Iteration :4
Cluster 0 size :50
Centroid [x=2.5737264423871213, y=-0.027462568841232982]
Cluster 1 size :50
Centroid [x=-0.4633368646347211, y=-0.466114096981958]
Cluster 2 size :50
Centroid [x=0.9888766205736857, y=2.0104789651972013]
```

```
For dataset2
Iteration :3
Cluster 0 size :102
Centroid [x=1.2708406269481842, y=-0.08583389704900131]
Cluster 1 size :98
Centroid [x=-0.2018593506236787, 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 [34]: k = KMeans()
#=====#
# STRART YOUR CODE HERE #
#=====#
k.main(dataname1, True)
k.main(dataname2, True)
#=====#
# END YOUR CODE HERE #
#=====#
```

```
For dataset1
Iteration :4
Purity is 0.040000
NMI is 1.000000
Cluster 0 size :50
Centroid [x=2.5737264423871213, y=-0.027462568841232982]
Cluster 1 size :50
Centroid [x=-0.4633368646347211, y=-0.466114096981958]
Cluster 2 size :50
Centroid [x=0.9888766205736857, y=2.0104789651972013]
```

```
For dataset2
Iteration :3
Purity is 0.015000
NMI is 0.205096
Cluster 0 size :102
Centroid [x=1.2708406269481842, y=-0.08583389704900131]
Cluster 1 size :98
Centroid [x=-0.2018593506236787, 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.



```

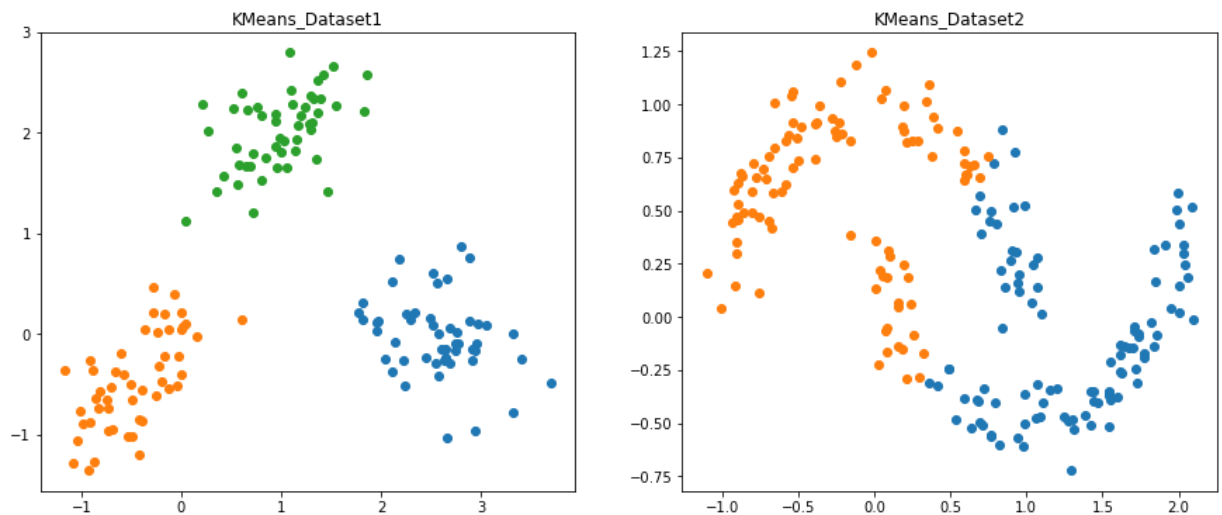
In [36]: CSV_FILE_PATH1 = 'Kmeans_dataset1.csv'
CSV_FILE_PATH2 = 'Kmeans_dataset2.csv'

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")

#=====#
# START YOUR CODE HERE #
#=====#
groups1 = df1.groupby('pred')
for name, group in groups1:
    ax0.plot(group["x"], group["y"], marker="o", linestyle="", label=name)

groups2 = df2.groupby('pred')
for name, group in groups2:
    ax1.plot(group["x"], group["y"], marker="o", linestyle="", label=name)
#=====#
# 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:

1. Efficient: time complexity for K-means is  $O(tkn)$  which is close to  $O(n)$  as  $t, k \ll n$  usually.

2. Simple: it is easy to implement the algorithm.

Cons:

1. Sensitive to noise and outlier data.
2. Not suitable for cluster with non-convex shapes
3. Needs to specify the number of k before training

## 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 [37]: 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 0.060000
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 0.020000
NMI is 0.817349
Cluster 0 size :99
Cluster 1 size :51
```

### 3.2 Visualization

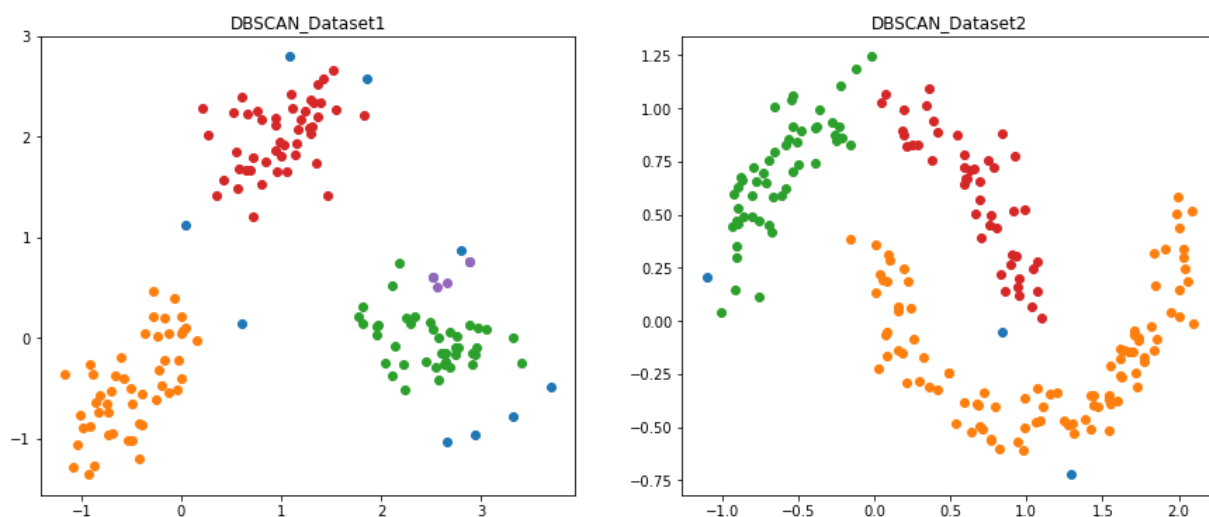
The clustering results for DBSCAN are saved as `DBSCAN_dataset1.csv` and `DBSCAN_dataset2.csv` respectively under your root folder. Plot the clustering results for the two datasets, with different colors representing different clusters.

```
In [38]: 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")

#=====#
# STRART YOUR CODE HERE #
#=====#
groups1 = df1.groupby('pred')
for name, group in groups1:
    ax0.plot(group["x"], group["y"], marker="o", linestyle="", label=name)

groups2 = df2.groupby('pred')
for name, group in groups2:
    ax1.plot(group["x"], group["y"], marker="o", linestyle="",label=name)
#=====#
# 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:

1. Does not need to specify number of  $k$  before.
2. Robust to outlier and noise.

Cons:

1. Has requirement of high density of dataset, also has problem with dataset with varying densities.

## 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 [40]: 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  
0.1173463305432000

## 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.

```

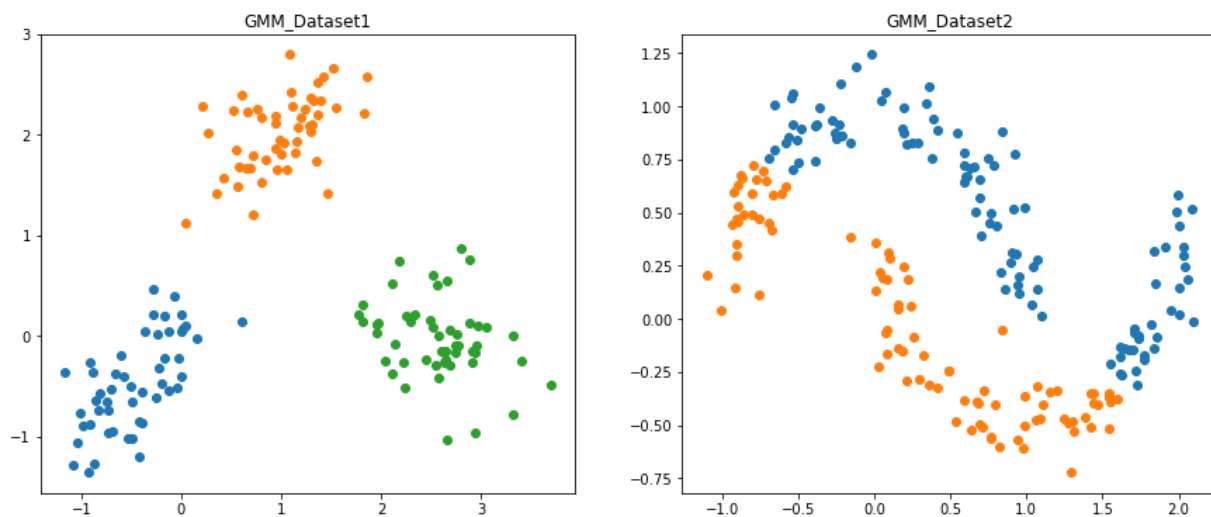
In [41]: CSV_FILE_PATH1 = 'GMM_dataset1.csv'
        CSV_FILE_PATH2 = 'GMM_dataset2.csv'

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 #
#=====#
groups1 = df1.groupby('pred')
for name, group in groups1:
    ax0.plot(group["x"], group["y"], marker="o", linestyle="", label=name)

groups2 = df2.groupby('pred')
for name, group in groups2:
    ax1.plot(group["x"], group["y"], marker="o", linestyle="",label=name)
#=====#
# 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:

1. Models are more general, so it can deal with different densities and sizes of cluster.
2. Clusters can be characterized by a small number of parameters.
3. The result may satisfy the statistical assumptions of the generative models.

Cons of GMM:

1. It will converge to local optimal instead of global.
2. Can only deal with spherical clusters.
3. Hard to estimate the number of clusters.

Reasoning over dataset1:

1. K-means: each of the three clusters are relatively far from each other while close inside the cluster, so the local optimum k-means reach is the global one this case.
2. DBSCAN: same reason that the three clusters are relatively far from each other. DBSCAN marks some outlier of each cluster and some points that are relatively between the cluster as noise as they are not density-reachable for given eps and minpts.
3. GMM: same reason that the clusters are far from each other making the final result the only possible local optimal for GMM.

Reasoning over dataset2:

1. K-means: the points are relatively uniformly distributed which makes k-means tend to create 2 cluster by just breaking the points in half by a linear separator.
2. DBSCAN: There is break in the upper half of the points which make DBSCAN cluster the upper swirl into two clusters while the bottom one is more densely-connected so clustered together. And there is certain space between the upper and bottom swirl so DBSCAN captures that and makes some points in between as noise.
3. GMM: similar reason to K-means. K-means reach the local optimal by simply breaking the points in two halves by a relatively linear separator.

## 5 Bonus Question

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

$J = \sum_{j=1}^k \sum_i w_{ij} ||x_i - c_j||^2$  as  $w_{ij} = 1$  if  $x_i$  belongs to cluster  $j$  is the loss function.

Let  $J_i$  be the loss in  $i$ -th iteration,  $c_{ix}$  be the center assigned to  $x$  in  $i$ -th iteration.

Then  $J_{i+1} - J_i = \sum ||x - c_{(i+1)x}||^2 - ||x - c_{ix}||^2$

As  $c_{(i+1)x}$  is created by  $\operatorname{argmin}_c(\operatorname{argmin}_w(c_{ix}))$ , then  $\sum ||x - c_{(i+1)x}||^2 - ||x - c_{ix}||^2 \leq 0$  and only =0 when  $c_{(i+1)x} = c_{ix}$  for all x which is the stopping condition.

Because J always  $\geq 0$ . Then the algorithm can converge as loss function is monotonically decreasing until it reaches the stopping/optimal/minimal condition.

## 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

In [ ]:



```

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

```

```

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

```

```

117         if logPart_I == 0.0:
118             continue
119         I += (nmiMatrix[i][j] / float(N)) * math.log(float(logPart_I))
120     # Compute H0mega
121     # =====#
122     # STRART YOUR CODE HERE #
123     # =====#
124     p_wj = nmiMatrix[i][-1] / float(N)
125     H0mega += (p_wj * np.log(p_wj))
126     # =====#
127     # END YOUR CODE HERE #
128     # =====#
129
130     #Compute HC
131     # =====#
132     # STRART YOUR CODE HERE #
133     # =====#
134     for c in range(col-1):
135         p_cj = nmiMatrix[-1][c] / float(N)
136         HC += (p_cj * np.log(p_cj))
137     # =====#
138     # END YOUR CODE HERE #
139     # =====#
140
141     return I / math.sqrt(HC * H0mega)
142
143
144
145
146
147 # =====#
148 class Centroid:
149     # -----#
150     def __init__(self, x, y):
151         self.x = x
152         self.y = y
153     # -----#
154     def __eq__(self, other):
155         if not type(other) is type(self):
156             return False
157         if other is self:
158             return True
159         if other is None:
160             return False
161         if self.x != other.x:
162             return False
163         if self.y != other.y:
164             return False
165         return True
166     # -----#
167     def __ne__(self, other):
168         result = self.__eq__(other)
169         if result is NotImplemented:
170             return result
171         return not result
172     # -----#
173     def toString(self):
174         return "Centroid [x=" + str(self.x) + ", y=" + str(self.y) + "]"
175     # -----#
176     def __str__(self):

```

```

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

```

```

235 while True:
236     iteration += 1
237     # calculate centroid for clusters
238     centroidsNew = []
239     for j in range(self.K):
240         centroidsNew.append(self.getCentroid(clusters[j]))
241
242     isConverge = False
243     for j in range(self.K):
244         if centroidsNew[j] != centroids[j]:
245             isConverge = False
246         else:
247             isConverge = True
248     if isConverge:
249         break
250
251     for j in range(self.K):
252         clusters[j] = set()
253
254     self.reassignClusters(self.dataSet, centroidsNew, clusters)
255     for j in range(self.K):
256         centroids[j] = centroidsNew[j]
257     print("Iteration :" + str(iteration))
258
259     if isevaluate:
260         # Calculate purity and NMI
261         compute_purity(clusters, len(self.dataSet))
262         compute_NMI(clusters, self.K)
263
264     # write clusters to file for plotting
265     f = open("Kmeans_" + self.dataname + ".csv", "w")
266     for w in range(self.K):
267         print("Cluster " + str(w) + " size :" + str(len(clusters[w])))
268         print(centroids[w].toString())
269         for point in clusters[w]:
270             f.write(str(point.x) + "," + str(point.y) + "," + str(w) + "\n")
271     f.close()
272
273     # -----
274     def reassignClusters(self, dataSet, c, clusters):
275         # reassign points based on cluster and continue till stable clusters found
276         dist = [0.0 for x in range(self.K)]
277         for point in dataSet:
278             for i in range(self.K):
279                 dist[i] = getEuclideanDist(point.x, point.y, c[i].x, c[i].y)
280
281             minIndex = self.getMin(dist)
282             # assign point to the closest cluster
283             # =====#
284             # STRART YOUR CODE HERE #
285             # =====#
286             clusters[minIndex].add(point)
287             # =====#
288             # END YOUR CODE HERE #
289             # =====#
290
291     # -----
292     def getMin(self, dist):
293         min = sys.maxsize
294         minIndex = -1
295         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     cx = np.average([c.x for c in cluster])
309     cy = np.average([c.y for c in cluster])
310     # =====#
311     # END YOUR CODE HERE #
312     # =====#
313     return Centroid(cx, cy)
314     # -----
315 @staticmethod
316 def readDataSet(filePath):
317     dataSet = []
318     with open(filePath) as f:
319         lines = f.readlines()
320     lines = [x.strip() for x in lines]
321     for line in lines:
322         points = line.split('\t')
323         x = float(points[0])
324         y = float(points[1])
325         label = int(points[2])
326         point = DataPoints(x, y, label)
327         dataSet.append(point)
328     return dataSet
329
```

```
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         # w_ij = w_j*f(xi)
120         numerator = self.w[j]*gaussian.pdf([self.dataSet[i].x,
self.dataSet[i].y])
121         # =====#
122         #   END YOUR CODE HERE   #
123         # =====#
124         self.W[i][j] = numerator
125         denominator += numerator
126
127         # normalize W[i][j] into probabilities
128         # =====#
129         # STRART YOUR CODE HERE #
130         # =====#
131         self.W[i] = self.W[i] / denominator
132         # =====#
133         #   END YOUR CODE HERE   #
134         # =====#
135     # -----#
136     def Mstep(self):
137         for j in range(self.K):
138             denominator = 0.0
139             numerator_x = 0.0
140             numerator_y = 0.0
141             cov_xy = 0.0
142             updatedMean_x = 0.0
143             updatedMean_y = 0.0
144
145             # update self.w[j] and self.mean
146             for i in range(len(self.dataSet)):
147                 denominator += self.W[i][j]
148                 updatedMean_x += self.W[i][j] * self.dataSet[i].x
149                 updatedMean_y += self.W[i][j] * self.dataSet[i].y
150
151             self.w[j] = denominator / len(self.dataSet)
152
153             #update self.mean
154             # =====#
155             # STRART YOUR CODE HERE #
156             # =====#
157             self.mean[j][0] = updatedMean_x / denominator
158             self.mean[j][1] = updatedMean_y / denominator
159             # =====#
160             #   END YOUR CODE HERE   #
161             # =====#
162
163             # update covariance matrix
164             for i in range(len(self.dataSet)):
165                 numerator_x += self.W[i][j] * pow((self.dataSet[i].x - self.mean[j]
[0]), 2)
166                 numerator_y += self.W[i][j] * pow((self.dataSet[i].y - self.mean[j]
[1]), 2)
167
168                 # Compute conv_xy +=?
169                 # =====#
170                 # STRART YOUR CODE HERE #
171                 # =====#
172
173                 covar = (self.dataSet[i].x - self.mean[j][0]) * (self.dataSet[i].y -
self.mean[j][1])
174                 cov_xy += self.W[i][j] * covar
175                 # =====#

```

```

175         # END YOUR CODE HERE #
176         # =====#
177
178         self.stdDev[j][0] = numerator_x / denominator
179         self.stdDev[j][1] = numerator_y / denominator
180
181
182         self.coVariance[j][0][0] = self.stdDev[j][0]
183         self.coVariance[j][1][1] = self.stdDev[j][1]
184         self.coVariance[j][0][1] = self.coVariance[j][1][0] = cov_xy /
denominator
185         # -----
186         def Likelihood(self):
187             likelihood = 0.0
188             for i in range(len(self.dataSet)):
189                 numerator = 0.0
190                 for j in range(self.K):
191                     gaussian = multivariate_normal(self.mean[j], self.coVariance[j])
192                     numerator += self.w[j] * gaussian.pdf([self.dataSet[i].x,
self.dataSet[i].y])
193                 likelihood += math.log(numerator)
194             return likelihood
195         # -----
196         def printArray(self, mat):
197             for i in range(len(mat)):
198                 for j in range(len(mat[i])):
199                     print(str(mat[i][j]) + " "),
200                 print("")
201         # -----
202         def print3D(self, mat):
203             for i in range(len(mat)):
204                 print("For Cluster : " + str((i + 1)))
205                 for j in range(len(mat[i])):
206                     for k in range(len(mat[i][j])):
207                         print(str(mat[i][j][k]) + " "),
208                     print("")
209                 print("")
210
211         # =====
212         if __name__ == "__main__":
213             g = GMM()
214             dataname = "dataset1.txt"
215             g.main(dataname)

```