

```
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)):
```

```

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                 noise.discard(point1)
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 +

```

```
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)

```

```

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    # -----
108    def Estep(self):
109        # Update self.W
110        for i in range(len(self.dataSet)):
111            denominator = 0.0
112            for j in range(self.K):
113                gaussian = multivariate_normal(self.mean[j],
self.coVariance[j])
                # Compute numerator for self.W[i][j] below

```

```

115         # =====#
116         # STRART YOUR CODE HERE #
117         # =====#
118         numerator = self.w[j] * gaussian.pdf([self.dataSet[i].x,
self.dataSet[i].y])
119         # =====#
120         # END YOUR CODE HERE #
121         # =====#
122         self.W[i][j] = numerator
123         denominator += numerator
124
125         # normalize W[i][j] into probabilities
126         # =====#
127         # STRART YOUR CODE HERE #
128         # =====#
129         for j in range(self.K):
130             self.W[i][j] /= denominator
131         # =====#
132         # END YOUR CODE HERE #
133         # =====#
134         # -----#
135     def Mstep(self):
136         for j in range(self.K):
137             denominator = 0.0
138             numerator_x = 0.0
139             numerator_y = 0.0
140             cov_xy = 0.0
141             updatedMean_x = 0.0
142             updatedMean_y = 0.0
143
144             # update self.w[j] and self.mean
145             for i in range(len(self.dataSet)):
146                 denominator += self.W[i][j]
147                 updatedMean_x += self.W[i][j] * self.dataSet[i].x
148                 updatedMean_y += self.W[i][j] * self.dataSet[i].y
149
150             self.w[j] = denominator / len(self.dataSet)
151
152             #update self.mean
153             # =====#
154             # STRART YOUR CODE HERE #
155             # =====#
156             self.mean[j][0] = updatedMean_x / denominator
157             self.mean[j][1] = updatedMean_y / denominator
158             # =====#
159             # END YOUR CODE HERE #
160             # =====#
161
162             # update covariance matrix
163             for i in range(len(self.dataSet)):
164                 numerator_x += self.W[i][j] * pow((self.dataSet[i].x -
self.mean[j][0]), 2)
165                 numerator_y += self.W[i][j] * pow((self.dataSet[i].y -
self.mean[j][1]), 2)
166                 # Compute conv_xy +=?
167                 # =====#
168                 # STRART YOUR CODE HERE #
169                 # =====#
170                 cov_xy += self.W[i][j] * (self.dataSet[i].x - self.mean[j]

```

```

171         # =====#
172         #   END YOUR CODE HERE   #
173         # =====#
174
175         self.stdDev[j][0] = numerator_x / denominator
176         self.stdDev[j][1] = numerator_y / denominator
177
178
179         self.coVariance[j][0][0] = self.stdDev[j][0]
180         self.coVariance[j][1][1] = self.stdDev[j][1]
181         self.coVariance[j][0][1] = self.coVariance[j][1][0] = cov_xy /
denominator
182         # -----
183         def Likelihood(self):
184             likelihood = 0.0
185             for i in range(len(self.dataSet)):
186                 numerator = 0.0
187                 for j in range(self.K):
188                     gaussian = multivariate_normal(self.mean[j],
self.coVariance[j])
189                     numerator += self.w[j] * gaussian.pdf([self.dataSet[i].x,
self.dataSet[i].y])
190                 likelihood += math.log(numerator)
191             return likelihood
192         # -----
193         def printArray(self, mat):
194             for i in range(len(mat)):
195                 for j in range(len(mat[i])):
196                     print(str(mat[i][j]) + " "),
197                 print("")
198         # -----
199         def print3D(self, mat):
200             for i in range(len(mat)):
201                 print("For Cluster : " + str((i + 1)))
202                 for j in range(len(mat[i])):
203                     for k in range(len(mat[i][j])):
204                         print(str(mat[i][j][k]) + " "),
205                     print("")
206                 print("")
207
208         # =====
209         if __name__ == "__main__":
210             g = GMM()
211             dataname = "dataset1.txt"
212             g.main(dataname)

```

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: Ali Mirabzadeh, **UID:** 305179067

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 [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.

| 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.

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

| cluster | Truth | Name | cluster | output label | Name |
|---------|-------|--------|---------|--------------|--------|
| w_1 | 3 | FJCAI | C_1 | 2 | FJCAI |
| | 3 | AAAI | | 2 | AAAI |
| | 3 | ICML | | 2 | ICML |
| | 3 | NIPS | | 2 | NIPS |
| | 3 | ECML | | 2 | EDBT |
| w_2 | 1 | ICDE | C_2 | 2 | ECML |
| | 1 | SIGMOD | | 3 | ICDE |
| | 1 | PODS | | 3 | VLDB |
| | 1 | EDBT | | 3 | SIGMOD |
| | 1 | VLDB | | 3 | CIKM |
| w_3 | 2 | KDD | C_3 | 3 | PODS |
| | 2 | PAKDD | | 4 | STAIR |
| | 2 | ICDM | | 4 | WWW |
| | 2 | PKDD | | 4 | ECIR |
| | 2 | SDM | | 4 | WSDM |
| w_4 | 4 | STAIR | C_4 | 1 | KDD |
| | 4 | CIKM | | 1 | PAKDD |
| | 4 | WWW | | 1 | ICDM |
| | 4 | ECIR | | 1 | PKDD |
| | 4 | WSDM | | 1 | SDM |

① $C_k \in \{C_1, C_2, C_3, C_4\}$, $w_j \in \{w_1, w_2, w_3, w_4\}$

$C_1 \rightarrow w_3$: majority 5, $C_2 \rightarrow w_1$: majority 5

$C_3 \rightarrow w_4$: majority 4, $C_4 \rightarrow w_2$: majority 5

$$\text{Purity} = \frac{5+5+5+1}{20} = \frac{16}{20} = \underline{\underline{0.8}}$$

answer 2, 3, 4

- Random Index (RI) = $(TP+TN)/(TP+FP+FN+TN)$
- F-measure: $2Precision*Recall/(Precision+Recall)$
 - Precision = $TP/(TP+FP)$
 - Recall = $TP/(TP+FN)$

| class | | | |
|-------|-----------------|--------------------|------------------------|
| 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 |

| | Same cluster | Different clusters |
|-------------------|--------------|--------------------|
| Same class | TP | FN |
| Different classes | FP | TN |

we group by ID's
and follow table
above to get the
below results

grouping by
ID's

TP
34

TN
143

FP
7

FN
5

②

$$Precision = \frac{34}{34+7} = \frac{34}{41} = \underline{\underline{0.83}}$$

③

$$Recall = \frac{34}{34+5} = \frac{34}{39} = \underline{\underline{0.87}}$$

④

$$F1-score = \frac{34+143}{34+143+7+5} = \frac{177}{189} = \underline{\underline{0.94}}$$

answer 5

⑤

| | C1 | C2 | C3 | C4 | sum |
|-----|----|----|----|----|-----|
| 1 | 0 | 0 | 0 | 5 | 5 |
| 2 | 6 | 0 | 0 | 0 | 6 |
| 3 | 0 | 5 | 0 | 0 | 5 |
| 4 | 0 | 0 | 4 | 0 | 4 |
| sum | 6 | 5 | 4 | 5 | 20 |

$$I(C, D) = \frac{5}{20} \log_2 \left(\frac{100}{25} \right) + \frac{6}{20} \log_2 \left(\frac{120}{36} \right) + \frac{5}{20} \log_2 \left(\frac{100}{25} \right) + \frac{4}{20} \log_2 \left(\frac{80}{16} \right) = 0.679$$

$$H(C) = -\frac{6}{20} \log_2 \frac{6}{20} - \frac{5}{20} \log_2 \frac{5}{20} - \frac{4}{20} \log_2 \frac{4}{20} - \frac{5}{20} \log_2 \frac{5}{20}$$

$$= 1.98$$

$$H(D) = 4 \left(-\frac{1}{4} \log_2 \frac{1}{4} \right) = 2$$

$$\sqrt{H(C)H(D)} = 1.98$$

$$NMI = \frac{0.679}{1.98} = \underline{\underline{0.342}}$$

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 [2]: 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 [3]: k = KMeans()
#=====#
#  STRART YOUR CODE HERE  #
#=====#
k.main(dataname1)
k.kmeans()
k.main(dataname2)
k.kmeans()
#=====#
#   END YOUR CODE HERE   #
#=====#
```

```
For dataset1
Iteration :3
Cluster 0 size :50
Centroid [x=2.5737264423871222, y=-0.027462568841232965]
Cluster 1 size :50
Centroid [x=-0.46333686463472107, y=-0.46611409698195816]
Cluster 2 size :50
Centroid [x=0.988876620573686, y=2.0104789651972013]
Iteration :3
Cluster 0 size :50
Centroid [x=2.5737264423871222, y=-0.027462568841232965]
Cluster 1 size :50
Centroid [x=-0.46333686463472107, y=-0.46611409698195816]
Cluster 2 size :50
Centroid [x=0.988876620573686, y=2.0104789651972013]

For dataset2
Iteration :4
Cluster 0 size :102
Centroid [x=1.2708406269481842, y=-0.08583389704900131]
Cluster 1 size :98
Centroid [x=-0.20185935062367868, y=0.5726963240559536]
Iteration :4
Cluster 0 size :102
Centroid [x=1.2708406269481842, y=-0.08583389704900131]
Cluster 1 size :98
Centroid [x=-0.20185935062367868, 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 [3]: k = KMeans()
#=====#
#  STRART YOUR CODE HERE  #
#=====#
k.main(dataname1)
k.kmeans(True)
k.main(dataname2)
k.kmeans(True)
#=====#
#  END YOUR CODE HERE  #
#=====#
```

```
For dataset1
Iteration :3
Cluster 0 size :50
Centroid [x=2.5737264423871222, y=-0.027462568841232965]
Cluster 1 size :50
Centroid [x=-0.46333686463472107, y=-0.46611409698195816]
Cluster 2 size :50
Centroid [x=0.988876620573686, y=2.0104789651972013]
Iteration :3
Purity is 1.000000
NMI is 1.000000
Cluster 0 size :50
Centroid [x=2.5737264423871222, y=-0.027462568841232965]
Cluster 1 size :50
Centroid [x=-0.46333686463472107, y=-0.46611409698195816]
Cluster 2 size :50
Centroid [x=0.988876620573686, y=2.0104789651972013]

For dataset2
Iteration :4
Cluster 0 size :102
Centroid [x=1.2708406269481842, y=-0.08583389704900131]
Cluster 1 size :98
Centroid [x=-0.20185935062367868, y=0.5726963240559536]
Iteration :4
Purity is 0.760000
NMI is 0.145025
Cluster 0 size :102
Centroid [x=1.2708406269481842, y=-0.08583389704900131]
Cluster 1 size :98
Centroid [x=-0.20185935062367868, 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.

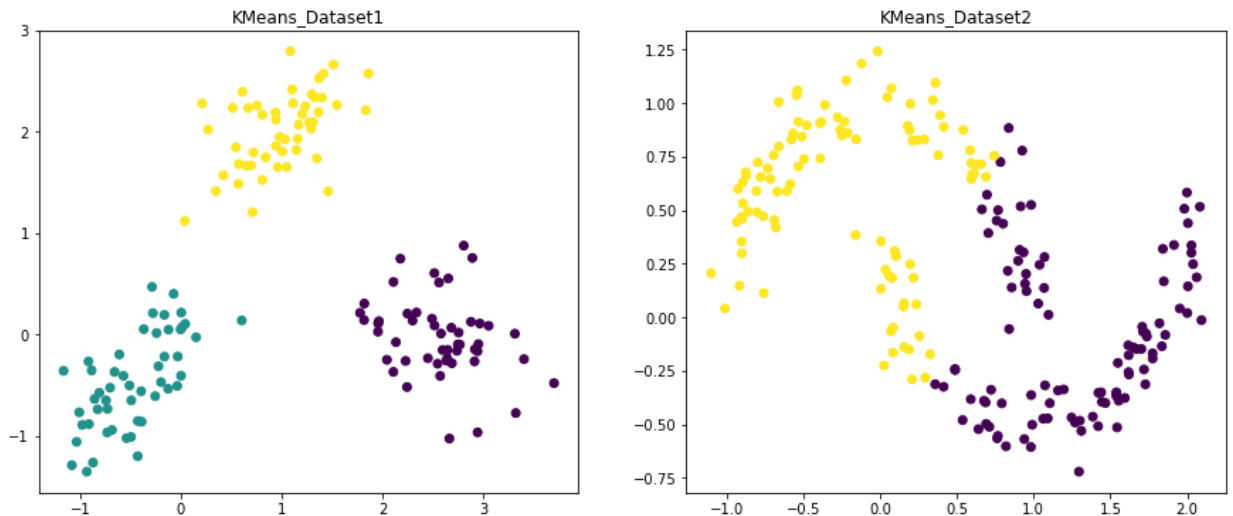

```

In [22]: 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 #
#####
ax0.scatter(df1.iloc[:, 0], df1.iloc[:, 1], c=df1.iloc[:, 2])
ax1.scatter(df2.iloc[:, 0], df2.iloc[:, 1], c=df2.iloc[:, 2])
#####
# 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. It's efficient as it has a linear run time

2. It's easy to interpret

Cons: Not suitable to discover clusters with non-convex shapes. 2 It's sensitive to noisy data

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 [10]: 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 :9
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
```

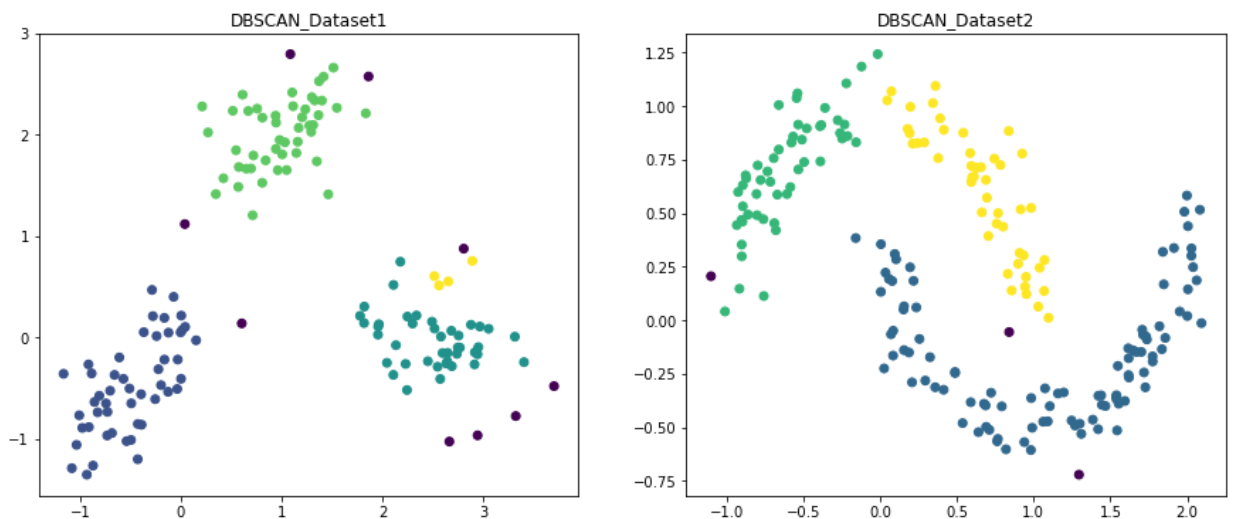
3.2 Visualization

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 [11]: 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 #
#####
ax0.scatter(df1.iloc[:, 0], df1.iloc[:, 1], c=df1.iloc[:, 2])
ax1.scatter(df2.iloc[:, 0], df2.iloc[:, 1], c=df2.iloc[:, 2])
#####
# 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: It can find clust with aribtarely shapes. 2 It's robust to outliers

Cons: Its outputs relies on its parameters and it's difficult to find the optimal parameters

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

```
For Cluster : 1
0.7692790765358335
-0.28782809642382123

-0.28782809642382123
0.1901249384356512
```

```
For Cluster : 2
0.6828574757628689
-0.30058915994390517

-0.30058915994390517
0.17583559485120062
```

```
Purity is 0.690000
NMI is 0.075948
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.

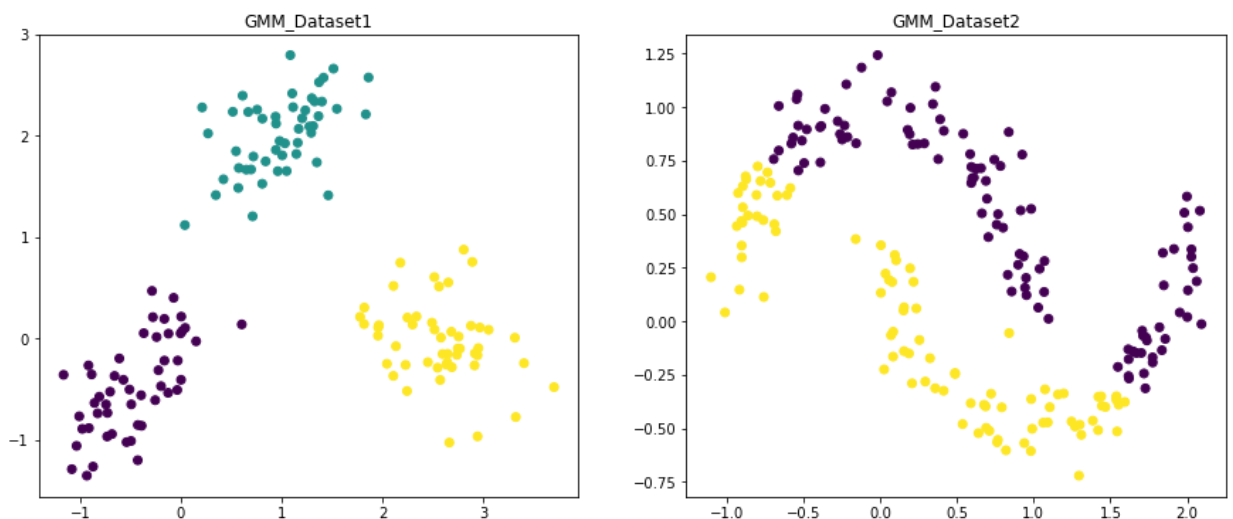
```

In [13]: 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")

#####
# STRART YOUR CODE HERE #
#####
ax0.scatter(df1.iloc[:, 0], df1.iloc[:, 1], c=df1.iloc[:, 2])
ax1.scatter(df2.iloc[:, 0], df2.iloc[:, 1], c=df2.iloc[:, 2])
#####
# 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. GMM models are more general than partitioning: different densities and sizes of clusters. 2. Clusters can be characterized by a small number of parameters

Cons of GMM: 1. Converge to local optimal. 2. Hard to estimate the number of clusters

Reasoning over dataset1: Kmeans and GMM resulted in 1,0 purity whereas DBSCAN in 0.94. And we can perfectly see in the visualization as there are three perfect clusters. However, for DBSCAN we can see there are some points identified as Noise and there is a fourth cluster as well. Maybe by tuning DBSCAN parameters we can get a better purity and cluster creation for DBSCAN

Reasoning over dataset2: DBSCAN gets the best purity, 0.98, then Kmeans with 0.88 and lastly GMM with 0.69 as the lowest one. I think DBSCAN performed the best because clusters are dense and separable. Kmeans is still perform relatively good and the reason it's not as good. Lastly we can see why GMM is not performing well due to its constraint on performing on non-convex shapes

5 Bonus Question

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

[Please type your answer here!](#)

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
7 # =====
8 def sqrt(n):
9     return math.sqrt(n)
10
11 # =====
12 def getEuclideanDist(x1, y1, x2, y2):
13     dist = sqrt(pow((x2 - x1), 2) + pow((y2 - y1), 2))
14     return dist
15 # =====
16 def compute_purity(clusters, total_points):
17     # Calculate purity
18
19     # Create list to store the maximum union number for each output cluster.
20     maxLabelCluster = []
21     num_clusters = len(clusters)
22     # =====#
23     # STRART YOUR CODE HERE #
24     # =====#
25     for i in range(num_clusters):
26         labelCounts = {}
27         for point in clusters[i]:
28             if not point.label in labelCounts:
29                 labelCounts[point.label] = 0
30                 labelCounts[point.label] += 1
31         max_union = -sys.maxsize - 1
32         for label in labelCounts:
33             if max_union < labelCounts[label]:
34                 max_union = labelCounts[label]
35         maxLabelCluster.append(max_union)
36     # =====#
37     # END YOUR CODE HERE #
38     # =====#
39     purity = 0.0
40     for j in range(num_clusters):
41         purity += maxLabelCluster[j]
42     purity /= total_points
43     print("Purity is %.6f" % purity)
44
45 # =====
46 def compute_NMI(clusters, noOfLabels):
47     # Get the NMI matrix first
48     nmiMatrix = getNMIMatrix(clusters, noOfLabels)
49     # Get the NMI matrix first
50     nmi = calcNMI(nmiMatrix)
51     print("NMI is %.6f" % nmi)
52
53
54 # =====
55 def getNMIMatrix(clusters, noOfLabels):
56     # Matrix shape of [num_true_clusters + 1, num_output_clusters + 1]
57     (example under week6's slide page 9)
58     nmiMatrix = [[0 for x in range(len(clusters) + 1)] for y in
59 range(noOfLabels + 1)]

```

```

59     for cluster in clusters:
60         # Create dictionary {true_class_No: Number of shared elements}
61         labelCounts = {}
62         # =====#
63         # STRART YOUR CODE HERE #
64         # =====#
65         for point in cluster:
66             if not point.label in labelCounts:
67                 labelCounts[point.label] = 0
68                 labelCounts[point.label] += 1
69         # =====#
70         # END YOUR CODE HERE #
71         # =====#
72         labelTotal = 0
73         labelCounts_sorted = sorted(labelCounts.items(), key=lambda item:
item[1], reverse=True)
74         for label, val in labelCounts_sorted:
75             nmiMatrix[label - 1][clusterNo] = labelCounts[label]
76             labelTotal += labelCounts.get(label)
77         # Populate last row (row of summation)
78         nmiMatrix[noOfLabels][clusterNo] = labelTotal
79         clusterNo += 1
80         labelCounts.clear()
81
82     # Populate last col (col of summation)
83     lastRowCol = 0
84     for i in range(noOfLabels):
85         totalRow = 0
86         for j in range(len(clusters)):
87             totalRow += nmiMatrix[i][j]
88         lastRowCol += totalRow
89         nmiMatrix[i][len(clusters)] = totalRow
90
91     # Total number of datapoints
92     nmiMatrix[noOfLabels][len(clusters)] = lastRowCol
93
94     return nmiMatrix
95
96     # =====#
97 def calcNMI(nmiMatrix):
98     # Num of true clusters + 1
99     row = len(nmiMatrix)
100    # Num of output clusters + 1
101    col = len(nmiMatrix[0])
102    # Total number of datapoints
103    N = nmiMatrix[row - 1][col - 1]
104    I = 0.0
105    H0mega = 0.0
106    HC = 0.0
107
108    for i in range(row - 1):
109        for j in range(col - 1):
110            # Compute the log part of each pair of clusters within I's
formula.
111            logPart_I = 1.0
112            # =====#
113            # STRART YOUR CODE HERE #
114            # =====#
115            logPart_I = (float(N) * nmiMatrix[i][j]) / (float(nmiMatrix[i]

```



```

116         # =====#
117         #   END YOUR CODE HERE   #
118         # =====#
119
120         if logPart_I == 0.0:
121             continue
122         I += (nmiMatrix[i][j] / float(N)) * math.log(float(logPart_I))
123     # Compute H0mega
124     # =====#
125     # STRART YOUR CODE HERE #
126     # =====#
127     H0mega += nmiMatrix[row - 1][j]/float(N) * math.log(nmiMatrix[row
- 1][j] / float(N))
128     # =====#
129     #   END YOUR CODE HERE   #
130     # =====#
131
132     #Compute HC
133     # =====#
134     # STRART YOUR CODE HERE #
135     # =====#
136     HC += nmiMatrix[i][col - 1]/float(N) * math.log(nmiMatrix[i][col -
1]/float(N))
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     # -----

```

```

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_data_loader(self, dataname):
205
206         df = pd.read_table(dataname, sep = "\t", header=None, names=
207         ['x', 'y', 'ground_truth_cluster'])
208         print("\nFor " + dataname[5:-4] + ": number of datapoints is %d" %
209         df.shape[0])
210         print(df.head(5))
211
212     # -----
213     def kmeans(self, isevaluate=False):
214         clusters = []
215         k = 0
216         while k < self.K:
217             cluster = set()
218             clusters.append(cluster)
219             k += 1
220
221         # Initially randomly assign points to clusters
222         i = 0
223         for point in self.dataSet:
224             clusters[i % k].add(point)
225             i += 1
226
227         # calculate centroid for clusters
228         centroids = []
229         for j in range(self.K):
230             centroids.append(self.getCentroid(clusters[j]))

```

```

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) +
271 "\n")
272     f.close()
273
274     # -----
275     def reassignClusters(self, dataSet, c, clusters):
276         # reassign points based on cluster and continue till stable clusters
277         found
278         dist = [0.0 for x in range(self.K)]
279         for point in dataSet:
280             for i in range(self.K):
281                 dist[i] = getEuclideanDist(point.x, point.y, c[i].x, c[i].y)
282
283             minIndex = self.getMin(dist)
284             # assign point to the closest cluster
285             # =====#
286             # STRART YOUR CODE HERE #
287             # =====#
288             for cluster in clusters:
289                 if point in cluster:
290                     cluster.remove(point)

```

```
290         # =====#
291         #   END YOUR CODE HERE   #
292         # =====#
293     # -----#
294     def getMin(self, dist):
295         min = sys.maxsize
296         minIndex = -1
297         for i in range(len(dist)):
298             if dist[i] < min:
299                 min = dist[i]
300                 minIndex = i
301         return minIndex
302
303     # -----#
304     def getCentroid(self, cluster):
305         # mean of x and mean of y
306         cx = 0
307         cy = 0
308         # =====#
309         # STRART YOUR CODE HERE #
310         # =====#
311         for data_point in cluster:
312             # print(data_point)
313             cx += data_point.x
314             cy += data_point.y
315         cx /= len(cluster)
316         cy /= len(cluster)
317         # =====#
318         #   END YOUR CODE HERE   #
319         # =====#
320         return Centroid(cx, cy)
321     # -----#
322     @staticmethod
323     def readDataSet(filePath):
324         dataSet = []
325         with open(filePath) as f:
326             lines = f.readlines()
327         lines = [x.strip() for x in lines]
328         for line in lines:
329             points = line.split('\t')
330             x = float(points[0])
331             y = float(points[1])
332             label = int(points[2])
333             point = DataPoints(x, y, label)
334             dataSet.append(point)
335         return dataSet
336
```