

```

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)

```