

```

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)

```