

```

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_dataloader(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 found
276 def reassignClusters(self, dataSet, c, clusters):
277     # reassign points based on cluster and continue till stable clusters
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
```