

```

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     def getMin(self, dist):
292         min = sys.maxsize
293         minIndex = -1
294         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
```