

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from download_data import download_data
4 from sklearn.cluster import KMeans
5 import random
6
7 MAX_ITERATIONS = 50;
8 MIN_SAMPLE_CHANGE = 0
9
10
11 colors=['red', 'blue', 'black', 'brown', 'c', 'm', 'y', 'k', 'w', 'orange']
12
13
14 def run_kmeans_clustering(k_value=2, showPlots=False, data=None):
15     """ Calculate K means, given K """
16     # pick "samples" random x,y coordinates from 0 to "samples"
17
18     # number of columns of data
19     samples_size = len(data) #329
20     data_dimensions = len(data[0]) #9
21     xy_samples = data #329*9 all data points
22
23     def indexes_to_list_array(indexes):
24         """ convert from indexes to an array of the x-y coordinates at each index """
25         items = []
26         for index in indexes:
27             items.append(xy_samples[index, :])
28         return np.array(items)
29     #returned a np.array type of indexed points' coordinates in all samples
30
31     # initialize dictionary
32     grouped_coordinates_indexes = initialize_dict(k_value) #a dictionary with k types
33
34     current_iteration = 0
35
36     # main iterations
37     while(current_iteration < MAX_ITERATIONS):
38
39         if current_iteration == 0:
40             # pick 2 random sets of x,y coordinates to be centroids to start off
41             #10*9 np.array this is random generated centroid points, performance is really bad!!!
42             #centroid_points = np.random.randint(samples_size, size=(k_value, data_dimensions))
43
44             #below is randomly picked centroid points, better performance
45             centroid_points = []
46             random_samples = random.sample(range(samples_size), k_value)
47             for index in random_samples:
48                 centroid_points.append(xy_samples[index, :])
49             centroid_points = np.array(centroid_points)
50         else:
51             #find centroids based on the current memberships
52             #          9          329*9          dict
53             centroid_points = _get_centroids(data_dimensions, xy_samples, grouped_coordinates_indexes)
54
55
56             # with new centroids, calculate distances and assign points to new groups
57             new_grouped_coordinates_indexes = assign_points_to_groups(k_value, xy_samples, centroid_points)
58
59
60             max_changed = check_minimum_changes_met(k_value, current_iteration, grouped_coordinates_indexes, new_grouped_coordinates_indexes)
61
62             if max_changed >= 0:
63                 # already shown, but show if "showPlots" if false
64                 print("Found due to {} minimum changes".format(max_changed))
65                 break;
66
67             # reassign new index group to existing
68             grouped_coordinates_indexes = new_grouped_coordinates_indexes.copy()
69
70             current_iteration += 1
71
72     print("Found after {} iterations".format(current_iteration + 1))
73
74     result_arrays = []
75
76     for i in range(k_value):
77         # return list of list of <centroid, grouped points>
78         result_arrays.append([centroid_points[i], indexes_to_list_array(grouped_coordinates_indexes[i])])
79
80     return result_arrays
81
82
83 def initialize_dict(k_value):
84     dict={}
85     for i in range(k_value):

```

```

86     dict[i] = []
87     return dict
88
89 def assign_points_to_groups(k_value, xy_samples, centroid_points):
90     """
91     assign each sample to the centroid it is closest to
92     returns:
93     dictionary of centroid (index) mapped to grouping of samples (indexes) that are closest
94     """
95     grouped_sample_indexes = initialize_dict(k_value)
96
97     for index in range(len(xy_samples)):
98         current_xy_samples = xy_samples[index, :]
99         current_xy_samples = np.array(current_xy_samples)
100         distances_to_centroid = []
101         # calculate distance of x/y coordinate from center
102         for centroid_index in range(len(centroid_points)):
103             distances_to_centroid.append(calc_euclidean_dist_vector(centroid_points[centroid_index, :], current_xy_samples))
104
105         minimum_index = distances_to_centroid.index(min(distances_to_centroid))
106         grouped_sample_indexes[minimum_index].append(index)
107     return grouped_sample_indexes
108
109
110 def calc_euclidean_dist_vector(vector1, vector2):
111     #####placeholder # start #####
112     result = np.linalg.norm(vector1-vector2)
113     #####placeholder # end #####
114     return result
115
116 # 10 x
117 def check_minimum_changes_met(k_value, current_iteration, old_grouped_samples, new_grouped_samples):
118     if current_iteration > 0:
119         # check to see if points within groups changed or not
120         unchanged_coordinates = []
121         for i in range(k_value):
122             original_group_length = len(old_grouped_samples[i])
123             unchanged_group_coordinates = set(new_grouped_samples[i].intersection(old_grouped_samples[i]))
124             unchanged_coordinates.append(abs(original_group_length - len(unchanged_group_coordinates)))
125         # find array that has the most number of samples that have changed
126
127         max_changed_index = unchanged_coordinates.index(max(unchanged_coordinates))
128         max_changed = unchanged_coordinates[max_changed_index]
129         if max_changed <= MIN_SAMPLE_CHANGE:
130             return max_changed
131     return -1
132
133
134 def _get_centroids(dimensions, xy_coordinates, groups):
135     xy_centroids = []
136     for i in range(len(groups)):
137         xy_centroids.append(xy_coordinates[groups[i], :])
138     centroid_points = get_centroids(dimensions, xy_centroids)
139     return centroid_points
140
141
142 def get_centroids(dimensions, xy_groups):
143     """ Takes tuple of coordinates
144     returns:
145     2,2 array of centroid points
146     """
147
148     def get_point_mean(array_values, dimensions):
149         """ take care of issue of empty list
150         """
151         if len(array_values) == 0:
152             return np.zeros(dimensions)
153         return array_values.mean(axis = 0)
154
155     length = len(xy_groups) #10
156     centroid_points = np.zeros((length, dimensions)) #10*9
157     # for each group, get the average point
158     #####placeholder # start #####
159     for i in range(length):
160         centroid_points[i, :] = get_point_mean(xy_groups[i], dimensions)
161     #####placeholder # end #####
162     return centroid_points
163
164
165 def get_sum_of_squares(dimensions, center, samples):
166     """ Get the sum of squared error, given centroid and all of its grouped points """
167     #####placeholder # start #####
168     length = len(samples)
169     sse = 0
170     for i in range(length):
171         sse = sse + np.linalg.norm(center - samples[i])**2

```

```

172 #####placeholder # end #####
173 return sse
174
175
176 if __name__ == "__main__":
177     # Load data
178     data = download_data("cities_life_ratings.csv").values
179
180     dimensions = len(data[0])
181
182     # evaluating Ks
183
184     k_values = [3, 6, 8, 10] # if go higher than 10, need to add to "colors" list
185     k_errors = []
186     for k in k_values:
187         result_arrays = run_kmeans_clustering(k, showPlots=False, data=data)
188         # step 6: calculate the sum of squared errors (SSE)
189         sse_total = 0
190         for i in range(k):
191             center = result_arrays[i][0]
192             samples = result_arrays[i][1]
193             sse_total += get_sum_of_squares(dimensions, center, samples)
194         k_errors += [sse_total]
195
196     plt.plot(k_values, k_errors)
197     plt.title("K-Means")
198     plt.xlabel("K values")
199     plt.ylabel("errors")
200     plt.show()
201
202     for i in range(len(k_values)):
203         print("K = {}".format(k_values[i]))
204         print("SSE = {}".format(k_errors[i]))

```