

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from download_data import download_data
4
5
6 MAX_ITERATIONS = 50;
7 MIN_SAMPLE_CHANGE = 0
8
9
10 colors=['red', 'blue', 'black', 'brown', 'c', 'm', 'y', 'k', 'w', 'orange']
11
12
13 def run_kmeans_clustering(k_value=2, showPlots=False, data=None):
14     """ Calculate K means, given K """
15     # pick "samples" random x,y coordinates from 0 to "samples"
16
17     # number of columns of data
18     samples_size = len(data)
19     data_dimensions = len(data[0])
20     xy_samples = data
21
22
23     # make sure not to show plots if more than 2 dimensions
24     if data_dimensions != 2:
25         showPlots = False
26
27     if showPlots:
28         plt.scatter(xy_samples[:, 0], xy_samples[:, 1], c='g')
29         plt.title("Data Points")
30         plt.show()
31
32     def indexes_to_list_array(indexes):
33         """ convert from indexes to an array of the x-y coordinates at each index """
34         items = []
35         for index in indexes:
36             items.append(xy_samples[index, :])
37         return np.array(items)
38     # initialize dictionary
39     grouped_coordinates_indexes = initialize_dict(k_value)
40
41     current_iteration = 0
42     # main iterations
43     while(current_iteration < MAX_ITERATIONS):
44
45         if current_iteration == 0:
46             # pick 2 random sets of x,y coordinates to be centroids to start off
47             centroid_points = np.random.randint(samples_size, size=(k_value, data_dimensions))
48             if showPlots:
49                 for i in range(k_value):
50                     plt.scatter(centroid_points[i, 0], centroid_points[i, 1], c=colors[i], marker='x')
51                     plt.scatter(xy_samples[:, 0], xy_samples[:, 1], c='g')
52                     plt.title("Initial centroids")
53                     plt.show()
54             else:
55                 # find centroids based on the current memberships
56                 centroid_points = _get_centroids(data_dimensions, xy_samples, grouped_coordinates_indexes)
57                 if showPlots:
58                     _display_plot("Calculate centroids", centroid_points, xy_samples, grouped_coordinates_indexes)
59
60
61             # with new centroids, calculate distances and assign points to new groups
62             new_grouped_coordinates_indexes = assign_points_to_groups(k_value, xy_samples, centroid_points)
63
64
65             if showPlots:
66                 _display_plot("Assign points to two centroids", centroid_points, xy_samples, new_grouped_coordinates_indexes)
67
68             max_changed = check_minimum_changes_met(k_value, current_iteration, grouped_coordinates_indexes, new_grouped_coordinates_indexes)
69             if max_changed >= 0:
70                 # already shown, but show if "showPlots" if false
71                 if data_dimensions == 2 and not showPlots:
72                     _display_plot("Assign points to two centroids", centroid_points, xy_samples, new_grouped_coordinates_indexes)
73
74                 print("Found due to {} minimum changes".format(max_changed))
75                 break;
76
77             # reassign new index group to existing
78             grouped_coordinates_indexes = new_grouped_coordinates_indexes.copy()
79
80             current_iteration += 1
81
82     print("Found after {} iterations".format(current_iteration + 1))
83     result_arrays = []
84     for i in range(k_value):
85         # return List of List of <centroid, grouped points>
86         result_arrays.append([centroid_points[i], indexes_to_list_array(grouped_coordinates_indexes[i])])
87
88     return result_arrays
89
90
91 def initialize_dict(k_value):
92     dict={}
93     for i in range(k_value):
94         dict[i] = []
95     return dict
96
97 def assign_points_to_groups(k_value, xy_samples, centroid_points):
98     """
99     assign each sample to the centroid it is closest to

```

```

100     returns:
101         dictionary of centroid (index) mapped to grouping of samples (indexes) that are closest
102     """
103     grouped_sample_indexes = initialize_dict(k_value)
104
105     for index in range(len(xy_samples)):
106         current_xy_samples = xy_samples[index, :]
107         distances_to_centroid = []
108         # calculate distance of x/y coordinate from center
109         for centroid_index in range(len(centroid_points)):
110             distances_to_centroid.append(calc_euclidean_dist_vector(centroid_points[centroid_index, :], current_xy_samples))
111
112         minimum_index = distances_to_centroid.index(min(distances_to_centroid))
113         grouped_sample_indexes[minimum_index].append(index)
114     return grouped_sample_indexes
115
116
117 def calc_euclidean_dist_vector(vector1, vector2):
118     #####placeholder # start #####
119     result =
120     #####placeholder # end #####
121     return result
122
123
124 def check_minimum_changes_met(k_value, current_iteration, old_grouped_samples, new_grouped_samples):
125     if current_iteration > 0:
126         # check to see if points within groups changed or not
127         unchanged_coordinates = []
128         for i in range(k_value):
129             original_group_length = len(old_grouped_samples[i])
130             unchanged_group_coordinates = set(new_grouped_samples[i]).intersection(old_grouped_samples[i])
131             unchanged_coordinates.append(abs(original_group_length - len(unchanged_group_coordinates)))
132             # find array that has the most number of samples that have changed
133
134         max_changed_index = unchanged_coordinates.index(max(unchanged_coordinates))
135         max_changed = unchanged_coordinates[max_changed_index]
136         if max_changed <= MIN_SAMPLE_CHANGE:
137             return max_changed
138     return -1
139
140
141 def _get_centroids(dimensions, xy_coordinates, groups):
142     xy_centroids = []
143     for i in range(len(groups)):
144         xy_centroids.append(xy_coordinates[groups[i], :])
145     centroid_points = get_centroids(dimensions, xy_centroids)
146     return centroid_points
147
148
149 def get_centroids(dimensions, xy_groups):
150     """ Takes tuple of coordinates
151     returns:
152         2,2 array of centroid points
153     """
154
155     def get_point_mean(array_values):
156         """ take care of issue of empty list
157         """
158         if len(array_values) == 0:
159             return 0
160         return int(array_values.mean())
161
162     length = len(xy_groups)
163     centroid_points = np.zeros((length, dimensions))
164     # for each group, get the average point
165     #####placeholder # start #####
166     for i in range(length):
167         centroid_points[i, :] =
168     #####placeholder # end #####
169     return centroid_points
170
171
172
173 def _display_plot(title, centroid_points, xy_coordinate_samples, grouped_coordinate_sample_indexes):
174     for i in range(len(centroid_points)):
175         plt.scatter(centroid_points[i, 0], centroid_points[i, 1], c=colors[i], marker='x')
176         plt.scatter(xy_coordinate_samples[grouped_coordinate_sample_indexes[i], 0], xy_coordinate_samples[grouped_coordinate_sample_indexes[i], 1], c=colors[i])
177     plt.title(title)
178     plt.show()
179
180
181 def get_sum_of_squares(dimensions, center, samples):
182     """ Get the sum of squared error, given centroid and all of its grouped points """
183     #####placeholder # start #####
184
185     sse =
186     #####placeholder # end #####
187     return sse
188
189
190 if __name__ == "__main__":
191     # Load data
192     data = download_data("cities_life_ratings.csv").values
193
194     dimensions = len(data[0])
195
196     # evaluating Ks
197
198     k_values = [3, 6, 8, 10] # if go higher than 10, need to add to "colors" List
199     k_errors = []

```

```

200 for k in k_values:
201     result_arrays = run_kmeans_clustering(k, showPlots=False, data=data)
202     # step 6: calculate the sum of squared errors (SSE)
203     sse_total = 0
204     for i in range(k):
205         center = result_arrays[i][0]
206         samples = result_arrays[i][1]
207         sse_total += get_sum_of_squares(dimensions, center, samples)
208     k_errors += [sse_total]
209
210 plt.plot(k_values, k_errors)
211 plt.title("K-Means")
212 plt.xlabel("K values")
213 plt.ylabel("errors")
214 plt.show()
215
216 for i in range(len(k_values)):
217     print("K = {}".format(k_values[i]))
218     print("SSE = {}".format(k_errors[i]))

```