

Application 3: Closest Pairs and Clustering

Algorithmic Thinking (Part2)

Shamsuddin Rehmani

July 15, 2019

Application 3 Description

In Project 3, you implemented two methods for clustering sets of data. In this Application, we will analyze the performance of these two methods on various subsets of our county-level cancer risk data set. In particular, we will compare these two clustering methods in three areas:

- Efficiency - Which method computes clusterings more efficiently?
- Automation - Which method requires less human supervision to generate reasonable clusterings?
- Quality - Which method generates clusterings with less error?

Answer to Question 1

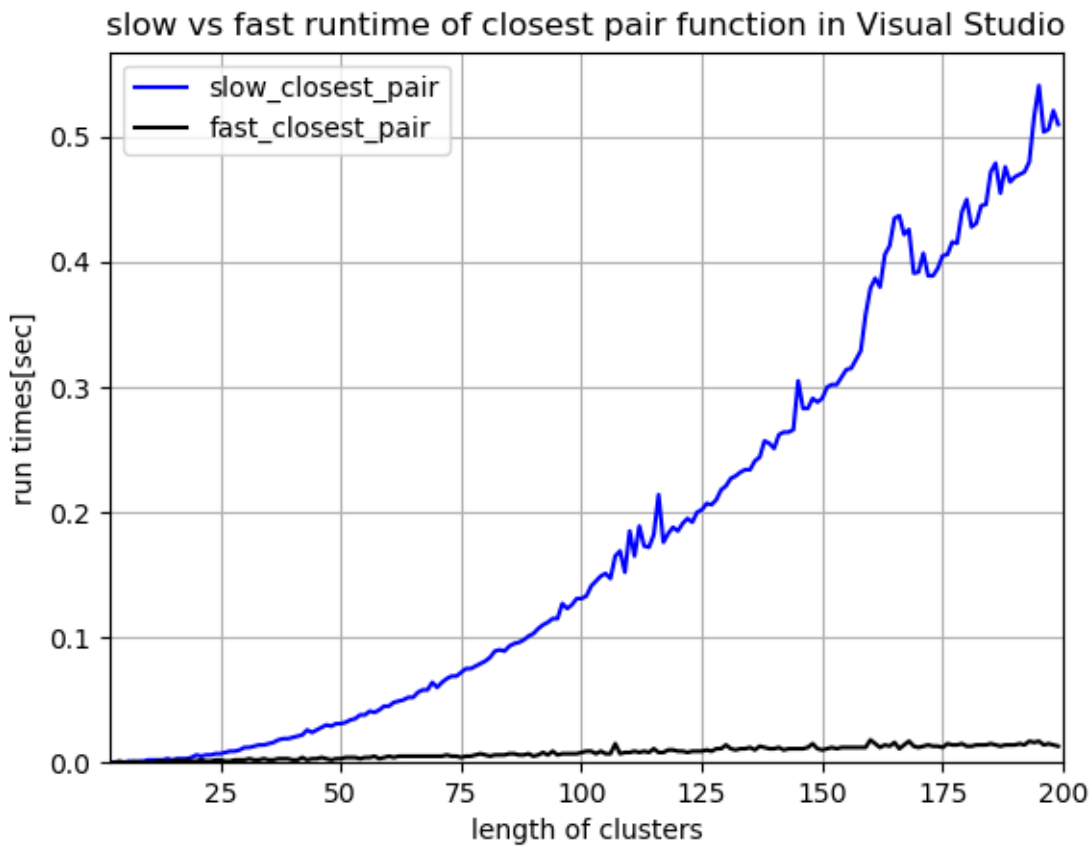


Figure 1: Fast vs slow closest pair function runtime comparison

Answer to Question 2

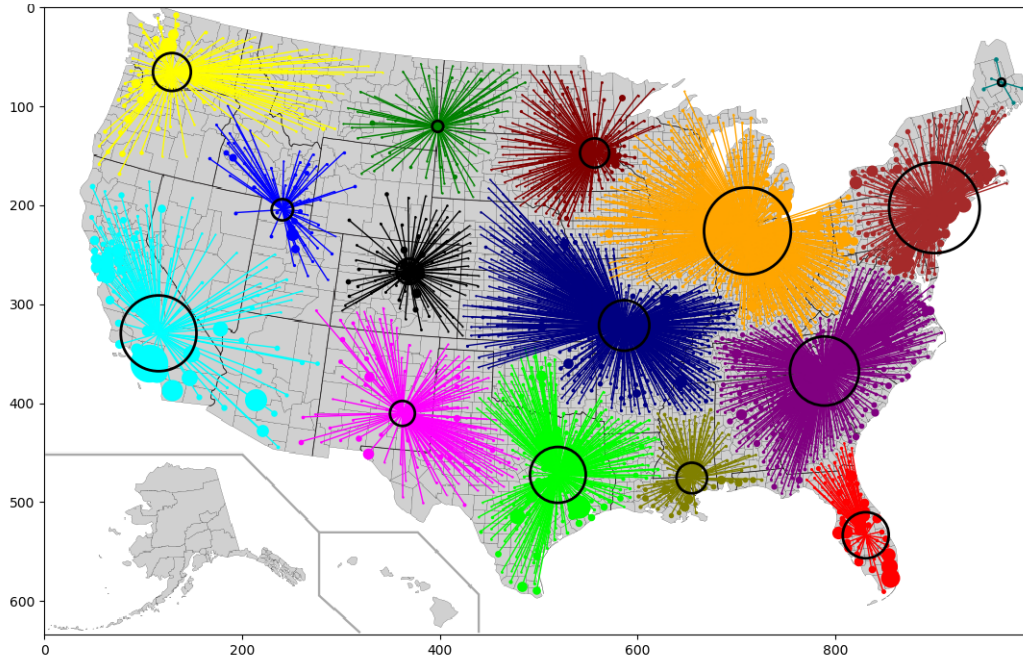


Figure 2: 15 clusters generated by hierarchical clustering for 3108 counties

Answer to Question 3

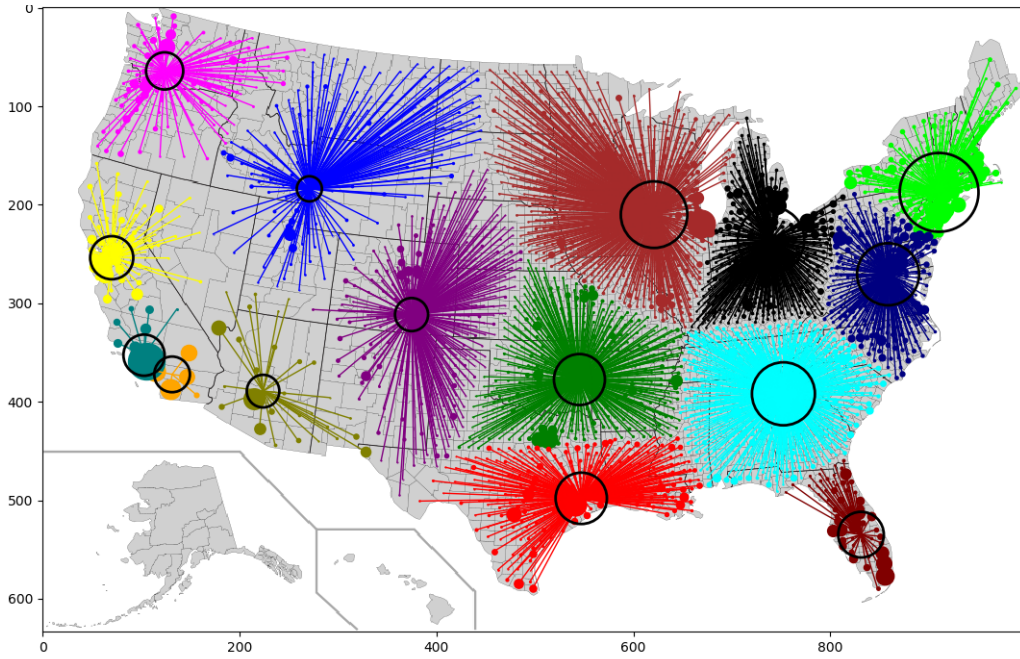


Figure 3: 15 clusters with 5 iterations generated by k-mean clustering for 3108 counties

Answer to Question 4

let k = number of clusters

n = size of the input cluster list

q = number of iterations (for k-mean clustering)

Then for hierarchical clustering, the time complexity is $O((n - k)(n \log(n) + n \log(n)^2))$ which is $O(n^2(\log n)^2)$ if k is small compared to n as stated in the above question

On the other hand, the time complexity of k-mean clustering is $O(qnk)$. Since k is small compared to n and q is also a small fixed number, time complexity is $O(n)$ which is much more efficient than heirarchical clustering

Answer to Question 5

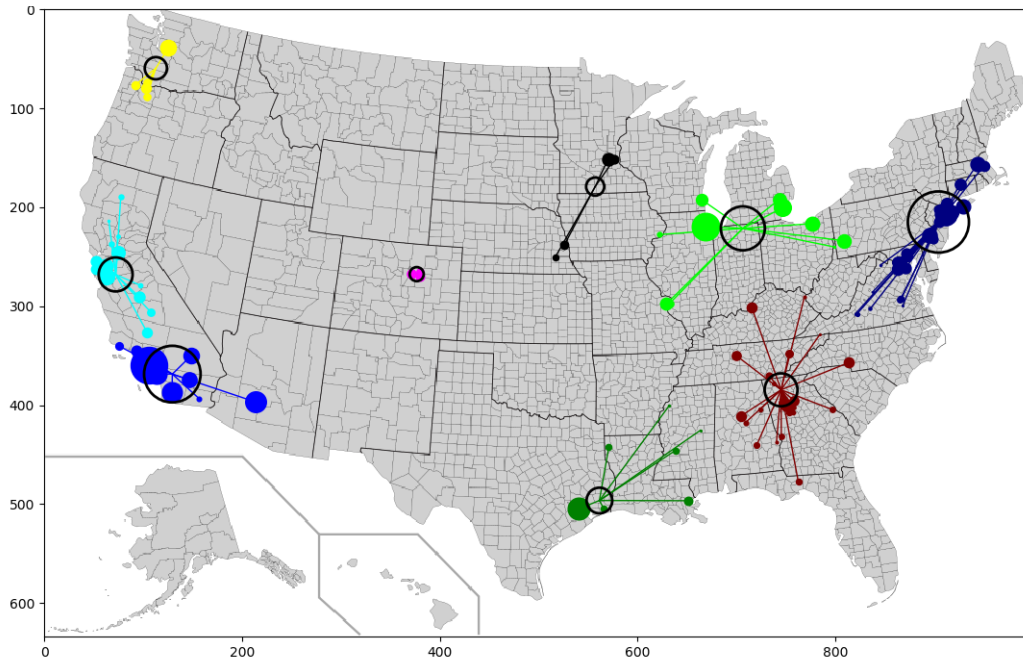


Figure 4: 9 clusters generated by hierarchical clustering for 111 counties

Answer to Question 6

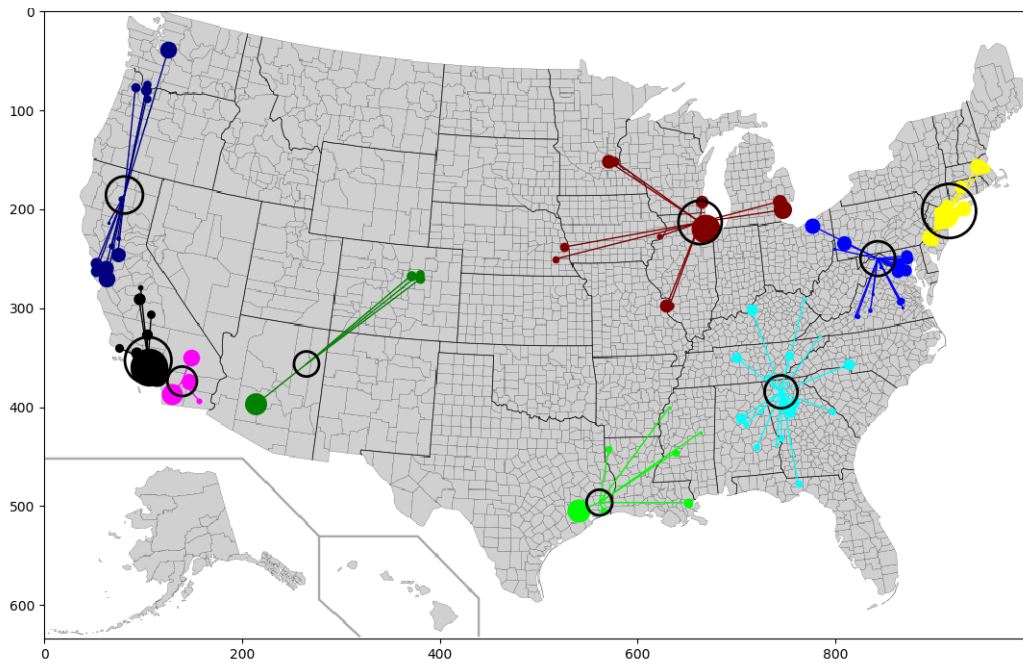


Figure 5: 9 clusters with 5 iterations generated by k-mean clustering for 111 counties

Answer to Question 7

Distortion of hierarchical clustering for 111 counties: 1.75×10^{11}

Distortion of k-mean clustering for 111 counties: 2.71×10^{11}

Answer to Question 8

For the k-mean clustering(Figure 5) we can see that the 3 cluster centers are located in the california region (one in northern while two in southern region) while of the 3 centers for the hierarchical clustering one is located in Washington state and two in the california region. The high distortion for the k-mean is due to the fact that the counties in the cluster with center in the northern California is distributed in Washington and southern California i.e the counties are much further than the center. The difference between the distortion values is because the intial clustering method for k-mean involves clustering around counties with highest population. Due to this all 3 counties(3 largest circles are black and pink) in the southern California were included in the intial clustering while none from Washington, Oregon or Northern California was selected. Thus this resulted in relatively higher distortion.

Answer to Question 9

Based on Question 8, we say that hierarchical clustering requires less human supervision as it requires only choosing the number of ouput clusters. While on the other hand, for k-mean we need a good choice of the intial cluster centers.

Answer to Question 10

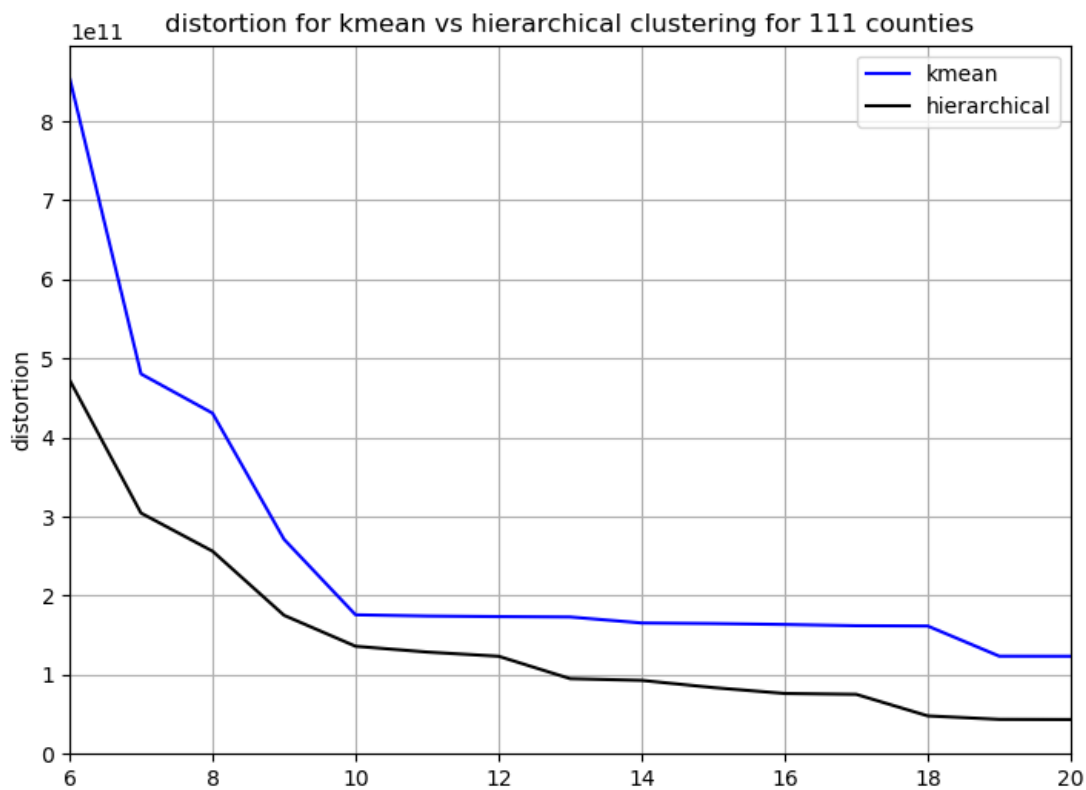


Figure 6: 9 clusters with 5 iterations generated by k-mean clustering for 111 counties

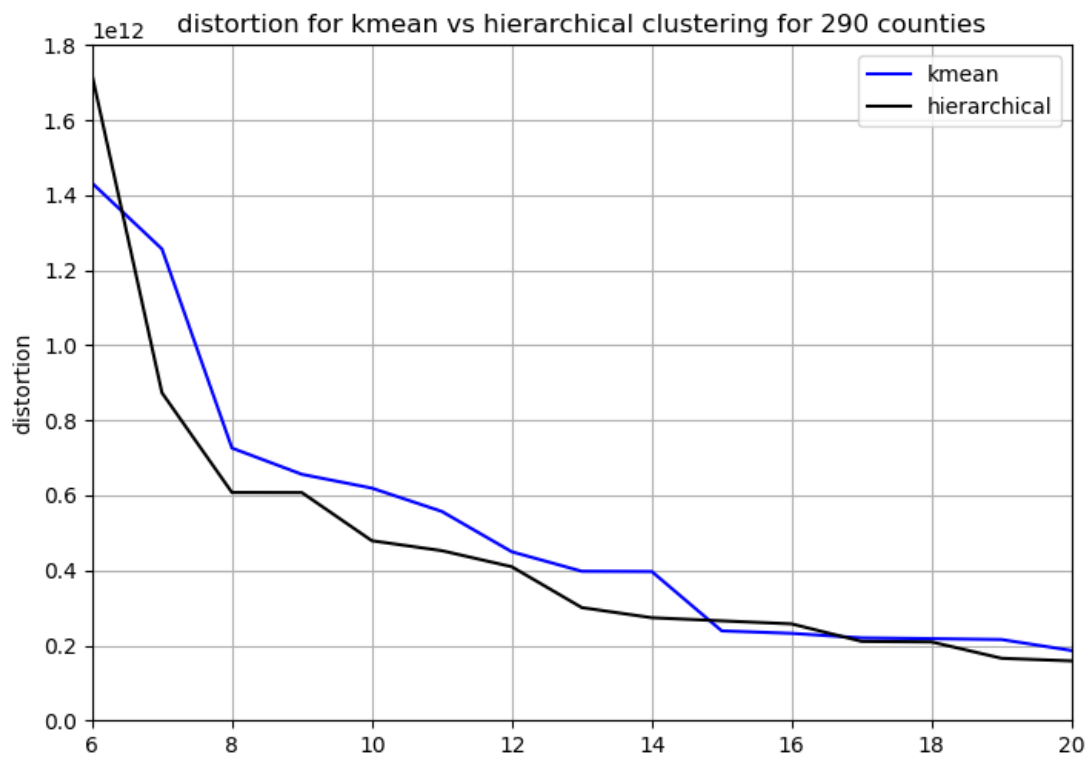


Figure 7: 9 clusters with 5 iterations generated by k-mean clustering for 111 counties

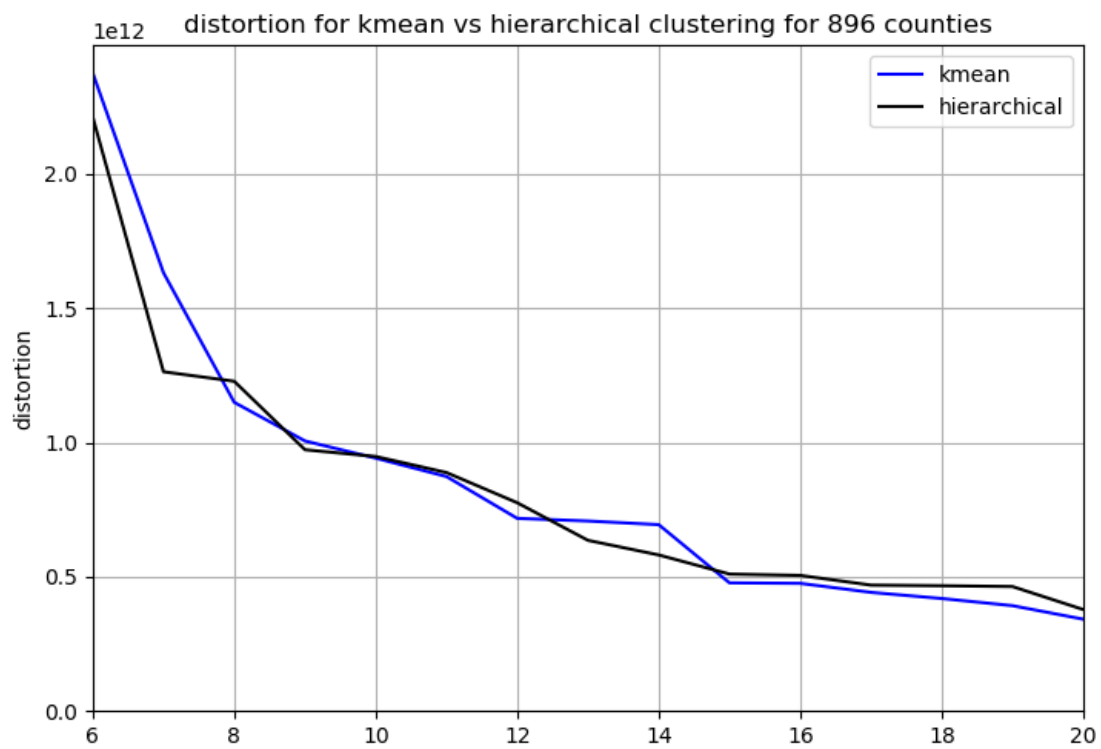


Figure 8: 9 clusters with 5 iterations generated by k-mean clustering for 111 counties

A Python code used to answer the Application Questions

```
1 """
2 Contains the answers to all the questions for
3 Application #3 – Comparision of Clustering Algorithm
4 """
5
6 import random
7 import copy
8 import matplotlib.pyplot as plt
9 import time
10 import alg_cluster
11 import alg_clusters.matplotlib as clust_plt
12 import alg_project3_solution as pj3_sol
13 import alg_project3_viz as pj3_viz
14
15 #####
16 ##### Code for Q1 Solution #####
17 # Q1: Write a function gen_random_clusters(num_clusters) that creates
18 # a list of clusters where each cluster in this list corresponds to one
19 # randomly generated point in the square with corners (+/-1,+/-1). Use
20 # this function and your favorite Python timing code to compute the
21 # running times of the functions slow_closest_pair and fast_closest_pair
22 # for lists of clusters of size 2 to 200. Once you have computed the running
23 # times for both functions, plot the result as two curves combined in a single
24 # plot. (Use a line plot for each curve.) The horizontal axis for your plot
25 # should be the the number of initial clusters while the vertical axis should
26 # be the running time of the function in seconds. Please include a legend in
27 # your plot that distinguishes the two curves.
28
29 def gen_random_clusters(num_clusters):
30     """[summary]
31
32     Arguments:
33         num_clusters {[type]} — [description]
34
35     Returns:
36         [type] — [description]
37     """
38
39     # initialize the cluster list to be returned
40     cluster_lst = num_clusters * [0]
41
42     for idx in range(num_clusters):
43         # generate a random point between -1 and 1
44         horz_center = random.uniform(-1.0, 1.0)
45         vert_center = random.uniform(-1.0, 1.0)
46         # create the cluster and add it to the list
47         cluster_lst[idx] = alg_cluster.Cluster(set(), horz_center, vert_center, 1, 0)
48
49     return cluster_lst
50
51 # initialize the range of cluster size to be used
52 clust_lens = range(2, 200)
53 # initialize the fast and slow pair function times
54 time_slow_closest = []
55 time_fast_closest = []
56
57 # clsts = gen_random_clusters(4)
58 # print len(clsts)
59 # print pj3_sol.slow_closest_pair(clsts)
60
61 # calculate the time to run the slow and fast functions for the closest pair
62 for clust_len in clust_lens:
63
64     # generate the cluster list of size clust_len
65     cluster_list = gen_random_clusters(clust_len)
66
67     # calculate the closest pair in the cluster using slow algorithm and
68     # store the time it takes to run the function for given cluster size
69     start = time.time()
70     pj3_sol.slow_closest_pair(cluster_list)
71     end = time.time()
72     time_slow_closest.append((end - start))
73     # calculate the closest pair in the cluster using fast algorithm and
74     # store the time it takes to run the function for given cluster size
75     start = time.time()
76     pj3_sol.fast_closest_pair(cluster_list)
77     end = time.time()
78     time_fast_closest.append((end - start))
79
80 # plot the graphs of resilience vs number of nodes removed for each of the 3 graphs
81 #plt.figure(1)
82 plt.plot(clust_lens, time_slow_closest, '-b', label = 'slow_closest_pair')
83 plt.plot(clust_lens, time_fast_closest, '-k', label = 'fast_closest_pair')
84 plt.title('slow vs fast runtime of closest pair function in Visual Studio')
85 plt.xlabel('length of clusters')
86 plt.ylabel('run times[sec]')
87 plt.legend(loc = 'upper left')
88 plt.xlim(2, 200)
89 plt.ylim(0, None)
90 plt.grid()
91 plt.show()
92 # uncommet to save the plot
93 #plt.savefig("Q1-closest-pair-comparision.png")
94
```



```

95
96 #####
97 # ##### Code for Q2 Solution #####
98 # Use alg_project3_viz to create an image of the 15 clusters generated by applying
99 # hierarchical clustering to the 3108 county cancer risk data set. You may submit
100 # an image with the 3108 counties colored by clusters or an enhanced visualization
101 # with the original counties colored by cluster and linked to the center of their
102 # corresponding clusters by lines. You can generate such an enhanced plot using our
103 # alg_clusters_matplotlib code by modifying the last parameter of plot_clusters to be
104 # True. Note that plotting only the resulting cluster centers is not acceptable
105
106 # load the data
107 data_table = pj3_viz.load_data_table(pj3_viz.DATA_3108_URL)
108
109 # generate cluster from the data
110 singleton_list = []
111 for line in data_table:
112     singleton_list.append(alg_cluster.Cluster(set([line[0]]), line[1], line[2], line[3], line[4]))
113
114 # create a deep copy since hierarchical_clustering modifies the references to the list
115 singleton_list_cpy = copy.deepcopy(singleton_list)
116
117 # form clusters based on hierarchical clustering algorithm
118 cluster_list = pj3_sol.hierarchical_clustering(singleton_list_cpy, 15)
119 print "Displaying", len(cluster_list), "hierarchical clusters"
120
121 # generate the image for the 15 clusters formed using hierarchical clustering algorithm
122 clust_plt.plot_clusters(data_table, cluster_list, True)
123
124 #####
125 # ##### Code for Q3 Solution #####
126 # Use alg_project3_viz to create an image of the 15 clusters generated by applying 5
127 # iterations of k-means clustering to the 3108 county cancer risk data set. You may
128 # submit an image with the 3108 counties colored by clusters or an enhanced visualization
129 # with the original counties colored by cluster and linked to the center of their corresponding
130 # clusters by lines. As in Project 3, the initial clusters should correspond to the 15 counties
131 # with the largest populations.
132
133 # form clusters based on k-mean clustering algorithm
134 cluster_list = pj3_sol.kmeans_clustering(singleton_list, 15, 5)
135 print "Displaying", len(cluster_list), "k-mean clusters"
136
137 clust_plt.plot_clusters(data_table, cluster_list, True)
138
139 #####
140 # ##### Q4 Solution #####
141 # Which clustering method is faster when the number of output clusters is either a small
142 # fixed number or a small fraction of the number of input clusters? Provide a short
143 # explanation in terms of the asymptotic running times of both methods. You should assume
144 # that hierarchical_clustering uses fast_closest_pair and that k-means clustering always
145 # uses a small fixed number of iterations.
146 #
147 # Ans: let k = number of clusters
148 #       n = size of the input cluster_list
149 #       q = number of iterations (for k-mean clustering)
150
151 # Then for hierarchical clustering, the time complexity is  $O((n - k) * (n * \log n + n * (\log n)^2))$ 
152 # which is  $\sim O(n^2 * (\log n)^2)$  if k is small compared to n as stated in the above question
153
154 # On the other hand, the time complexity of k-mean clustering is  $\sim O(q * n * k)$ . Since k is small
155 # compared to n and q is also a small fixed number, time complexity is  $O(n)$  which is much more
156 # efficient than heirarchiacl clustering
157
158 #####
159 # ##### Q5 Solution #####
160 # Use alg_project3_viz to create an image of the 9 clusters generated by applying hierarchical
161 # clustering to the 111 county cancer risk data set. You may submit an image with the 111
162 # counties colored by clusters or an enhanced visualization with the original counties colored
163 # by cluster and linked to the center of their corresponding clusters by lines.
164
165 # load the data
166 data_table = pj3_viz.load_data_table(pj3_viz.DATA_111_URL)
167
168 # generate cluster from the data
169 singleton_list = []
170 for line in data_table:
171     singleton_list.append(alg_cluster.Cluster(set([line[0]]), line[1], line[2], line[3], line[4]))
172
173 # create a deep copy since hierarchical_clustering modifies the references to the list
174 singleton_list_cpy = copy.deepcopy(singleton_list)
175 # form clusters based on hierarchical clustering algorithm
176 cluster_list_hierarc = pj3_sol.hierarchical_clustering(singleton_list_cpy, 9)
177
178 # generate the image for the 9 clusters formed using hierarchical clustering algorithm
179 print "Displaying", len(cluster_list_hierarc), "hierarchical clusters"
180 clust_plt.plot_clusters(data_table, cluster_list_hierarc, True)
181
182 #####
183 # ##### Q6 Solution #####
184 # Use alg_project3_viz to create an image of the 9 clusters generated by applying 5
185 # iterations of k-means clustering to the 111 county cancer risk data set. You may
186 # submit an image with the 111 counties colored by clusters or an enhanced visualization
187 # with the original counties colored by cluster and linked to the center of their
188 # corresponding clusters by lines. As in Project 3, the initial clusters should correspond
189 # to the 9 counties with the largest populations.
190

```

```

191 # form clusters based on k-mean clustering algorithm
192 cluster_list.kmean = pj3.sol.kmeans_clustering(singleton_list, 9, 5)
193
194 # generate the image for the 9 clusters formed using k-means clustering algorithm
195 print "Displaying", len(cluster_list_hierarc), "k-mean clusters"
196 clust_plt.plot_clusters(data_table, cluster_list_kmean, True)
197
198 #####
199 ##### Q7 Solution #####
200 # Write a function compute_distortion(cluster_list) that takes a list of clusters and
201 # uses cluster_error to compute its distortion. Now, use compute_distortion to compute
202 # the distortions of the two clusterings in questions 5 and 6
203
204 def compute_distortion(cluster_list, data_table):
205
206     distortion = sum([cluster.cluster_error(data_table) for cluster in cluster_list])
207
208     return distortion
209
210 print compute_distortion(cluster_list_hierarc, data_table)
211 print compute_distortion(cluster_list_kmean, data_table)
212
213 #####
214 ##### Q8 Solution #####
215 # Examine the clusterings generated in Questions 5 and 6. In particular, focus your
216 # attention on the number and shape of the clusters located on the west coast of the USA.
217 # Describe the difference between the shapes of the clusters produced by these two methods
218 # on the west coast of the USA. What caused one method to produce a clustering with a much
219 # higher distortion? To help you answer this question, you should consider how k-means
220 # clustering generates its initial clustering in this case.
221 #
222 # Ans: For the k-mean clustering (Figure 5) we can see that the 3 cluster centers are located
223 # in the California region (one in northern while two in southern region) while of the 3
224 # centers for the hierarchical clustering one is located in Washington state and two in the
225 # California region. The high distortion for the k-mean is due to the fact that the counties
226 # in the cluster with center in the northern California is distributed in Washington and
227 # southern California i.e the counties are much further than the center. The difference between
228 # the distortion values is because the initial clustering method for k-mean involves clustering
229 # around counties with highest population. Due to this all 3 counties (3 largest circles are black
230 # and pink) in the southern California were included in the initial clustering while none from
231 # Washington, Oregon or Northern California was selected. Thus this resulted in relatively
232 # higher distortion.
233
234 #####
235 ##### Q9 Solution #####
236 # Based on your answer to Question 8, which method (hierarchical clustering or k-means
237 # clustering) requires less human supervision to produce clusterings with relatively
238 # low distortion?
239 #
240 # Ans: based on Q8, we say that hierarchical clustering requires less human supervision
241 # as it requires only choosing the number of output clusters. While on the other hand
242 # for k-mean we need a good choice of the initial cluster centers.
243
244 #####
245 ##### Q10 Solution #####
246
247 hierarc_distortion = 15 * [0]
248 kmean_distortion = 15 * [0]
249 num_clusters_range = range(20, 5, -1)
250
251 ## do the calculations for 111 data set
252 # load the data
253 data_table = pj3.viz.load_data_table(pj3.viz.DATA.111.URL)
254
255 # generate cluster from the data
256 singleton_list = []
257 for line in data_table:
258     singleton_list.append(alg_cluster.Cluster(set([line[0]]), line[1], line[2], line[3], line[4]))
259
260 singleton_list.cpy = copy.deepcopy(singleton_list)
261
262 for num_clusters in num_clusters_range:
263     cluster_list_hierarc = pj3.sol.hierarchical_clustering(singleton_list.cpy, num_clusters)
264     cluster_list_kmean = pj3.sol.kmeans_clustering(singleton_list, num_clusters, 5)
265     hierarc_distortion[num_clusters - 6] = compute_distortion(cluster_list_hierarc, data_table)
266     kmean_distortion[num_clusters - 6] = compute_distortion(cluster_list_kmean, data_table)
267
268 # reverse the num_clusters_range
269 num_clusters_range.reverse()
270
271 # plot the graphs for distortion for two clustering method for 111 counties
272 plt.figure(2)
273 plt.plot(num_clusters_range, kmean_distortion, '-b', label = 'kmean')
274 plt.plot(num_clusters_range, hierarc_distortion, '-k', label = 'hierarchical')
275 plt.title('distortion for kmean vs hierarchical clustering for 111 counties')
276 plt.xlabel('number of clusters')
277 plt.ylabel('distortion')
278 plt.legend(loc = 'upper right')
279 plt.xlim(6, 20)
280 plt.ylim(0, None)
281 plt.grid()
282 plt.show()
283
284 ## do the calculations for 290 data set
285 # load the data
286 data_table = pj3.viz.load_data_table(pj3.viz.DATA.290.URL)

```



```

287
288 # generate cluster from the data
289 singleton_list = []
290 for line in data_table:
291     singleton_list.append(alg_cluster.Cluster(set([line[0]]), line[1], line[2], line[3], line[4]))
292
293 num_clusters_range = range(20, 5, -1)
294 singleton_list_cpy = copy.deepcopy(singleton_list)
295
296 for num_clusters in num_clusters_range:
297     cluster_list_hierarc = pj3_sol.hierarchical_clustering(singleton_list_cpy, num_clusters)
298     cluster_list_kmean = pj3_sol.kmeans_clustering(singleton_list, num_clusters, 5)
299     hierarc_distortion[num_clusters - 6] = compute_distortion(cluster_list_hierarc, data_table)
300     kmean_distortion[num_clusters - 6] = compute_distortion(cluster_list_kmean, data_table)
301
302 # reverse the num-cluster-range
303 num_clusters_range.reverse()
304
305 # plot the graphs for distortion for two clustering method for 290 counties
306 #plt.figure(3)
307 plt.plot(num_clusters_range, kmean_distortion, '-b', label = 'kmean')
308 plt.plot(num_clusters_range, hierarc_distortion, '-k', label = 'hierarchical')
309 plt.title('distortion for kmean vs hierarchical clustering for 290 counties')
310 plt.xlabel('number of clusters')
311 plt.ylabel('distortion')
312 plt.legend(loc = 'upper right')
313 plt.xlim(6, 20)
314 plt.ylim(0, None)
315 plt.grid()
316 plt.show()
317
318 ## do the calculations for 896 data set
319 # load the data
320 data_table = pj3_viz.load_data_table(pj3_viz.DATA_896_URL)
321
322 # generate cluster from the data
323 singleton_list = []
324 for line in data_table:
325     singleton_list.append(alg_cluster.Cluster(set([line[0]]), line[1], line[2], line[3], line[4]))
326
327 num_clusters_range = range(20, 5, -1)
328 singleton_list_cpy = copy.deepcopy(singleton_list)
329
330 for num_clusters in num_clusters_range:
331     cluster_list_hierarc = pj3_sol.hierarchical_clustering(singleton_list_cpy, num_clusters)
332     cluster_list_kmean = pj3_sol.kmeans_clustering(singleton_list, num_clusters, 5)
333     hierarc_distortion[num_clusters - 6] = compute_distortion(cluster_list_hierarc, data_table)
334     kmean_distortion[num_clusters - 6] = compute_distortion(cluster_list_kmean, data_table)
335
336 # reverse the num-cluster-range
337 num_clusters_range.reverse()
338
339 # plot the graphs for distortion for two clustering method for 290 counties
340 #plt.figure(3)
341 plt.plot(num_clusters_range, kmean_distortion, '-b', label = 'kmean')
342 plt.plot(num_clusters_range, hierarc_distortion, '-k', label = 'hierarchical')
343 plt.title('distortion for kmean vs hierarchical clustering for 896 counties')
344 plt.xlabel('number of clusters')
345 plt.ylabel('distortion')
346 plt.legend(loc = 'upper right')
347 plt.xlim(6, 20)
348 plt.ylim(0, None)
349 plt.grid()
350 plt.show()
351
352 #####
353 ##### Q11 Solution #####
354 # For each data set (111, 290, and 896 counties), does one clustering method consistently
355 # produce lower distortion clusterings when the number of output clusters is in the range
356 # 6 to 20? Is so, indicate on which data set(s) one method is superior to the other.
357 #
358 # Ans: for 111 counties, hierarchical clustering consistently produces less clustering
359 # For other two data set, there is no one method that consistently produces lower distortion

```

B All functions for project 4 used in the application

```

1 """
2 All the algorithms for Project 3 for closest
3 pair and clustering
4 """
5
6 import alg_cluster
7
8 #####
9 # Code for closest pairs of clusters
10 def slow_closest_pair(cluster_list):
11     """
12     Uses the brute-force algorithm ( $O(n^2)$ ) to find the closest pair of clusters in a list
13
14     Arguments:
15         cluster_list {Cluster} — a list of Cluster objects
16
17     Returns:
18         tuple — returns a closest pair where the pair is represented by the tuple

```

```

19         (dist, idx1, idx2) with idx1 < idx2 where dist is the distance between
20         the closest pair cluster_list[idx1] and cluster_list[idx2].
21     """
22     # initialize the tuple that will store the closest pair of cluster distance and index
23     (dist, idx1, idx2) = (float("inf"), -1, -1)
24
25     for cluster_i_idx in range(len(cluster_list)):
26         for cluster_j_idx in range(len(cluster_list)):
27             if (cluster_i_idx != cluster_j_idx):
28                 curr_dist = cluster_list[cluster_i_idx].distance(cluster_list[cluster_j_idx])
29                 (dist, idx1, idx2) = (min(set([(dist, idx1, idx2), (curr_dist, cluster_i_idx, cluster_j_idx)]),
30                                     key = lambda tup: tup[0]))
31
32     if (idx2 > idx1):
33         return (dist, idx1, idx2)
34     else:
35         return (dist, idx2, idx1)
36
37 def fast_closest_pair(cluster_list):
38     """
39     Compute the distance between the closest pair of clusters in a list using fast algorithm
40     [O(n*(logn)^2)]
41
42     Arguments:
43         cluster_list {list} — list of clusters sorted based on the horizontal distance of
44         their centers in ascending order
45
46     Returns:
47         tuple — tuple of the form (dist, idx1, idx2) where the centers of the clusters
48         cluster_list[idx1] and cluster_list[idx2] have minimum distance dist.
49     """
50
51     # base case
52     if len(cluster_list) <= 3:
53         (dist, idx1, idx2) = slow_closest_pair(cluster_list)
54
55     # inductive case
56     else:
57         # divide the problem in half, solve it and then merge the results from both
58         idx_m = len(cluster_list) / 2
59         # solve for the left half of the cluster list
60         (d_left, idxi_l, idxj_l) = fast_closest_pair(cluster_list[:idx_m])
61         # solve for right half of the cluster list
62         (d_right, idxi_r, idxj_r) = fast_closest_pair(cluster_list[idx_m:])
63         # find the minimum of the left and right partition minimum distances
64         (dist, idx1, idx2) = (min(set([(d_left, idxi_l, idxj_l), (d_right, idxi_r + idx_m, idxj_r + idx_m)]),
65                                 key = lambda tup: tup[0]))
66         # find the horizontal position of the strip's vertical center line i.e the midpoint of the
67         # horizontal position of the last element of the left partition and first element of the right partition
68         horiz_center = 0.5 * (cluster_list[idx_m - 1].horiz_center() + cluster_list[idx_m].horiz_center())
69         # find the minimum of the minimum distance found earlier and the closest pair in the strip
70         (dist, idx1, idx2) = (min(set([(dist, idx1, idx2), closest_pair_strip(cluster_list, horiz_center, dist)]),
71                                 key = lambda tup: tup[0]))
72
73     return (dist, idx1, idx2)
74
75 def closest_pair_strip(cluster_list, horiz_center, half_width):
76     """
77     Helper function to compute the closest pair of clusters in a vertical strip
78
79     Arguments:
80         cluster_list {list} — a list of Cluster objects
81         horiz_center {integer} — the horizontal position of the strip's vertical center line
82         half_width {integer} — the half the width of the strip (i.e; the maximum horizontal distance
83         that a cluster can lie from the center line)
84
85     Returns:
86         tuple — returns a tuple of the form (dist, idx1, idx2) where the centers of the clusters
87         cluster_list[idx1] and cluster_list[idx2] lie in the strip and have minimum distance dist.
88     """
89
90     # store all cluster whose centres are within the vertical strip specified by horiz_center and half_width
91     strip_clust = [(cluster_list[clst_idx], clst_idx) for clst_idx in range(len(cluster_list))
92                   if abs(cluster_list[clst_idx].horiz_center() - horiz_center) < half_width]
93
94     # sort the cluster based on their vertical distances
95     strip_clust.sort(key = lambda cluster: cluster[0].vert_center())
96
97     # initialize the minimum distance and their indices
98     (dist, idx1, idx2) = (float("inf"), -1, -1)
99
100    # for each cluster inspect the next 3 ones and record the pair of cluster indices that
101    # corresponds to closest pair thus found
102    for clsti_idx in range(len(strip_clust) - 1):
103        for clstj_idx in range(clsti_idx + 1, min((clsti_idx + 4), len(strip_clust))):
104            curr_dist = strip_clust[clsti_idx][0].distance(strip_clust[clstj_idx][0])
105            (dist, idx1, idx2) = (min(set([(dist, idx1, idx2),
106                                         (curr_dist, strip_clust[clsti_idx][1], strip_clust[clstj_idx][1])]), key = lambda tup: tup[0]))
107
108    # return minimum distance and their indices (dist, idx1, idx2) where idx1 < idx2
109    if (idx2 > idx1):
110        return (dist, idx1, idx2)
111    else:
112        return (dist, idx2, idx1)
113

```

```

114 #####
115 # Code for hierarchical clustering
116 def hierarchical_clustering(cluster_list, num_clusters):
117     """
118     Compute a hierarchical clustering of a set of clusters
119     Note: the function mutates cluster_list to have length num_clusters
120
121     Arguments:
122         cluster_list {list} — a list of Cluster objects
123         num_clusters {integer} — integer number of clusters to be made from cluster_list
124
125     Returns:
126         list — a list of Cluster objects whose length is num_clusters
127     """
128
129     while len(cluster_list) > num_clusters:
130         # sort the cluster based on their horizontal distances
131         #print cluster_list
132         cluster_list.sort(key = lambda cluster: cluster.horiz_center())
133         #print cluster_list
134         (dummy_dist, idx1, idx2) = fast_closest_pair(cluster_list)
135         cluster_list[idx1].merge_clusters(cluster_list[idx2])
136         #print cluster_list
137         cluster_list.pop(idx2)
138
139     return cluster_list
140
141 #####
142 # Code for k-means clustering
143
144 def kmeans_clustering(cluster_list, num_clusters, num_iterations):
145     """
146     Compute the k-means clustering of a set of clusters
147     Note: the function does not mutate cluster_list to have length num_clusters
148
149     Arguments:
150         cluster_list {list} — a list of Cluster objects
151         num_clusters {integer} — integer number of clusters to be made from cluster_list
152         num_iterations {integer} — number of iterations
153
154     Returns:
155         list — a list of Cluster objects whose length is num_clusters
156     """
157     # compute an initial list of clusters with the property that each cluster consists of
158     # a single county chosen from the set of the num_cluster counties with the largest populations
159     cluster_list_cpy = list(cluster_list)
160     cluster_list_cpy.sort(key = lambda cluster: cluster.total_population())
161     cluster_len = len(cluster_list)
162     old_cluster = cluster_list_cpy[cluster_len - num_clusters : ]
163
164     for dummy_idx in range(num_iterations):
165         # initialize an empty cluster
166         new_cluster = [alg_cluster.Cluster(set(), 0.0, 0.0, 0, 0) for dummy_i in range(num_clusters)]
167         for idx_i in range(cluster_len):
168             min_dist = float("inf")
169             min_idx = -1
170             # find the cluster from old_cluster and its index that is closest to current cluster
171             # in the cluster_list
172             for idx_j in range(num_clusters):
173                 curr_dist = old_cluster[idx_j].distance(cluster_list[idx_i])
174                 if curr_dist < min_dist:
175                     min_dist = curr_dist
176                     min_idx = idx_j
177             # add the cluster to the new_cluster at index = min_idx
178             new_cluster[min_idx].merge_clusters(cluster_list[idx_i])
179
180         old_cluster = new_cluster
181
182     return new_cluster

```