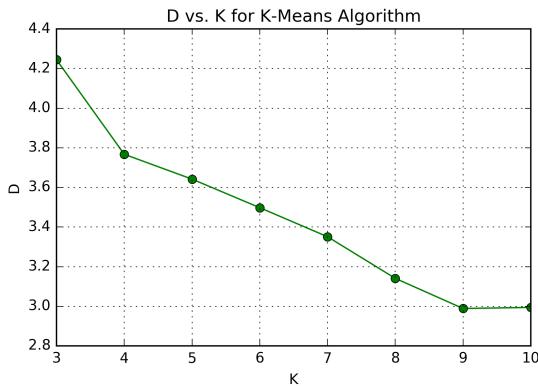


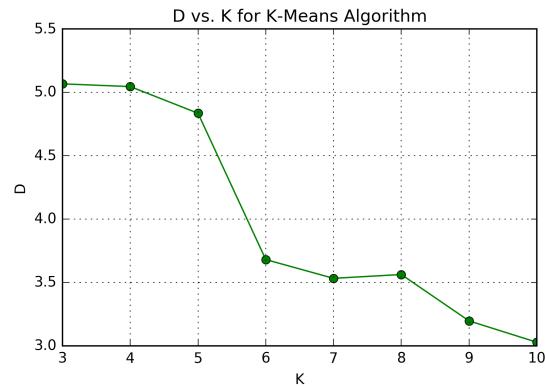
1 Clustering Result Comparison

I. Lloyd's (K-means) Algorithm

- In this part, with the Lloyd's algorithm for k-means clustering, we choose the best distortion D as the the objection function. The change of D versus cluster number K is shown in Fig. 1, where the result for clustering.txt is shown in Fig. 1a and the result for bigClusteringData.txt is show in Fig. 1b.



(a) clustering.txt



(b) bigClusteringData.txt

Figure 1: Change of Distortion versus Cluster Number K for K-Means Algorithm

- The scatter plot of the clustering result for clustering.txt is shown in Fig. 2 and the scatter plot of the clustering result for bigClusteringData.txt is shown in Fig. 3. The cluster centroids are clearly marked and different clusters are denoted by different colors.

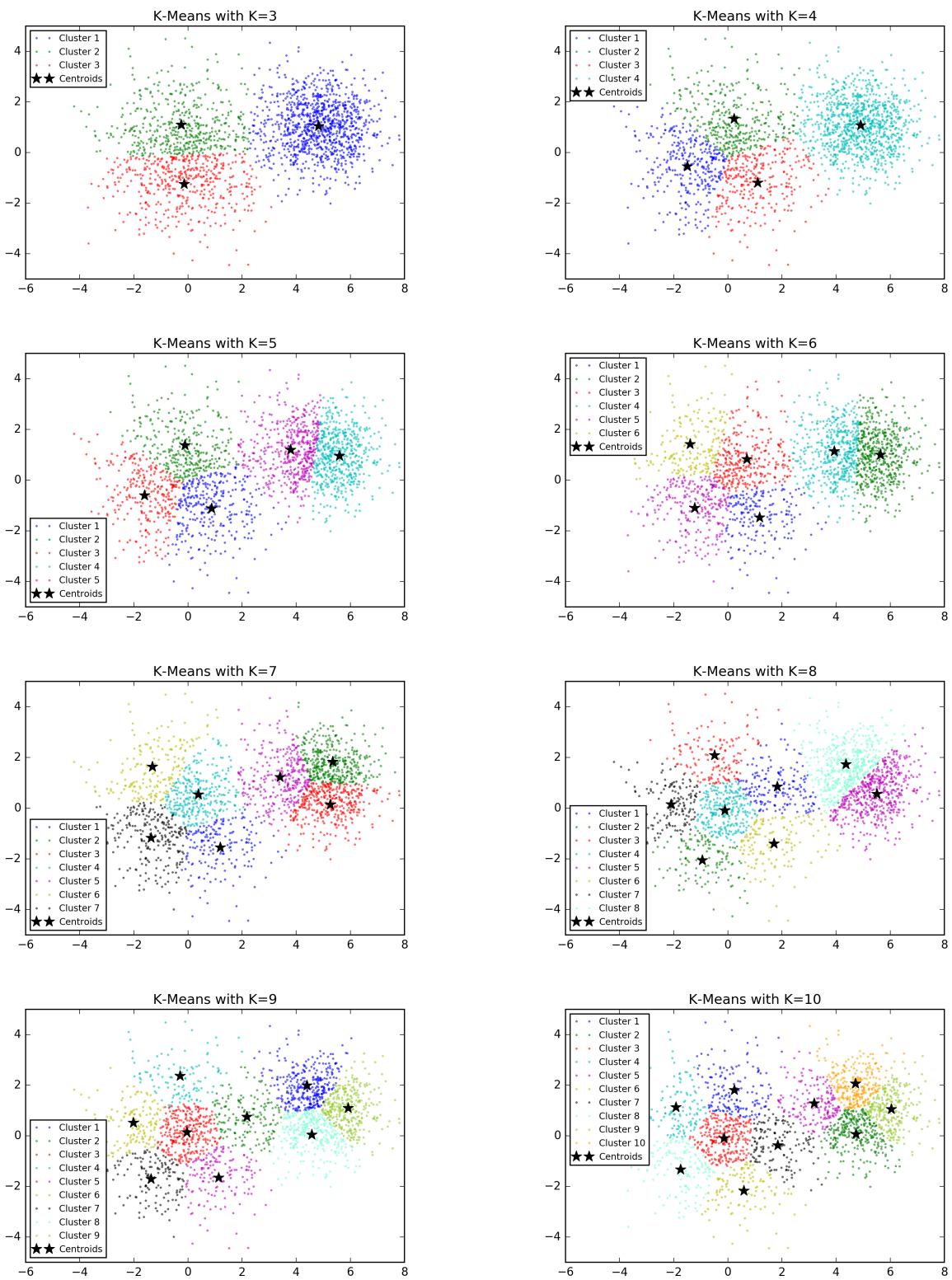


Figure 2: Clustering Result for `clustering.txt` with K-Means Algorithm

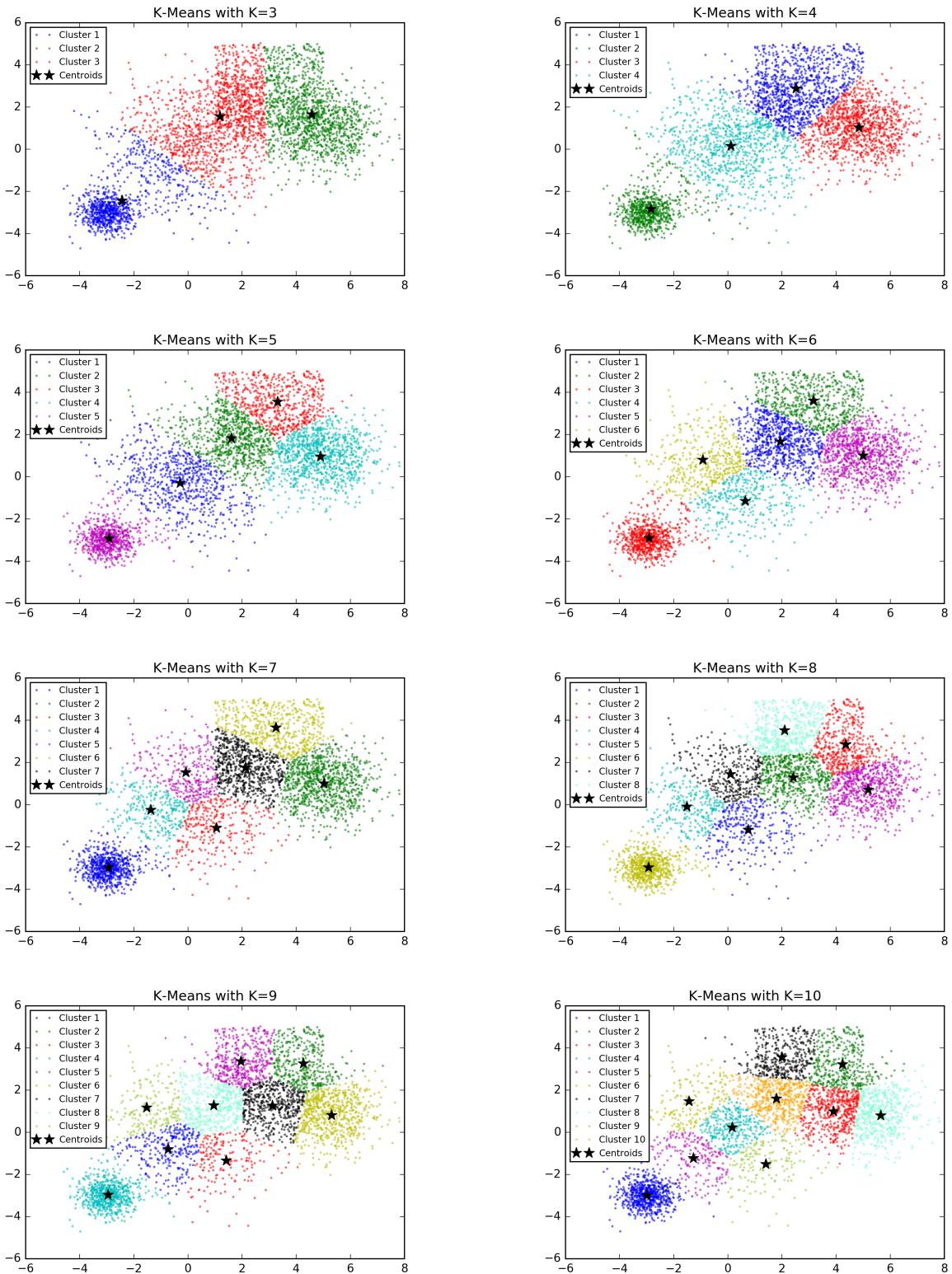


Figure 3: Clustering Result for `bigClusteringData.txt` with K-Means Algorithm

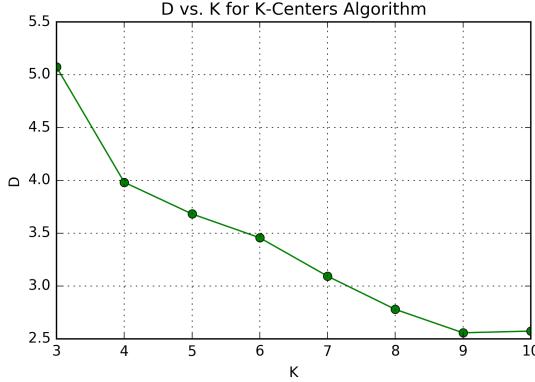
3. The Python code used for K-Means Algorithm is shown in Listing 1.

```
1 import numpy as np
2 import time
3
4 def kMeans(X, K, tol=0.00001, random_state=None, verbose=True):
5     """ function to implement the Lloyd's algorithm for k-means problem """
6     np.random.seed(random_state)
7     t0 = time.time()
8
9     N, d = X.shape # number of observations and dimensions
10    index = np.random.choice(range(N), size=K, replace=False)
11    Y = X[index, :]
12    C = np.zeros(N)
13    D = 100
14    count = 0
15    diff = 100 # difference between D1 and D0
16
17    while diff >= tol:
18        D0 = D
19        for i in range(N):
20            # assign centers to ith data
21            C[i] = np.argmin(np.sum((Y - X[i, :]) ** 2, axis=1))
22
23        D = 0
24        # re-compute the new centers
25        for j in range(K):
26            Y[j, :] = np.mean(X[C == j, :], axis=0)
27
28        # compute the loss
29        loss = np.zeros((N, K))
30        for i in range(K):
31            loss[:, i] = np.sqrt(np.sum((X - Y[i, :]) ** 2, axis=1))
32        D = np.max(np.min(loss, axis=1))
33        diff = abs(D - D0)
34        count += 1
35
36    if verbose is True:
37        t = np.round(time.time() - t0, 4)
38        print('K-Means finished in ' + str(t) + 's, ' + str(count) + ' iterations')
39
40    return Y, C, D
```

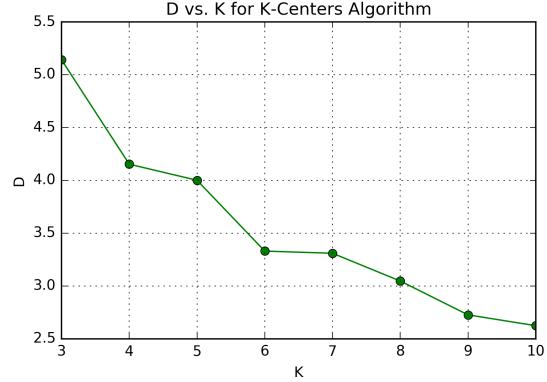
Listing 1: K-Means Algorithm Python Code

II. Greedy K-centers Algorithm

1. In this part, with the Lloyd's algorithm for k-means clustering, we choose the best distortion D as the the objection function. The change of D versus cluster number K is shown in Fig. 4, where the result for clustering.txt is shown in Fig. 4a and the result for bigClusteringData.txt is show in Fig. 4b.



(a) clustering.txt



(b) bigClusteringData.txt

Figure 4: Change of Distortion versus Cluster Number K for K-Center Algorithm

2. The scatter plot of the clustering result for clustering.txt is shown in Fig. 5 and the scatter plot of the clustering result for bigClusteringData.txt is shown in Fig. 6. The cluster centroids are clearly marked and different clusters are denoted by different colors.

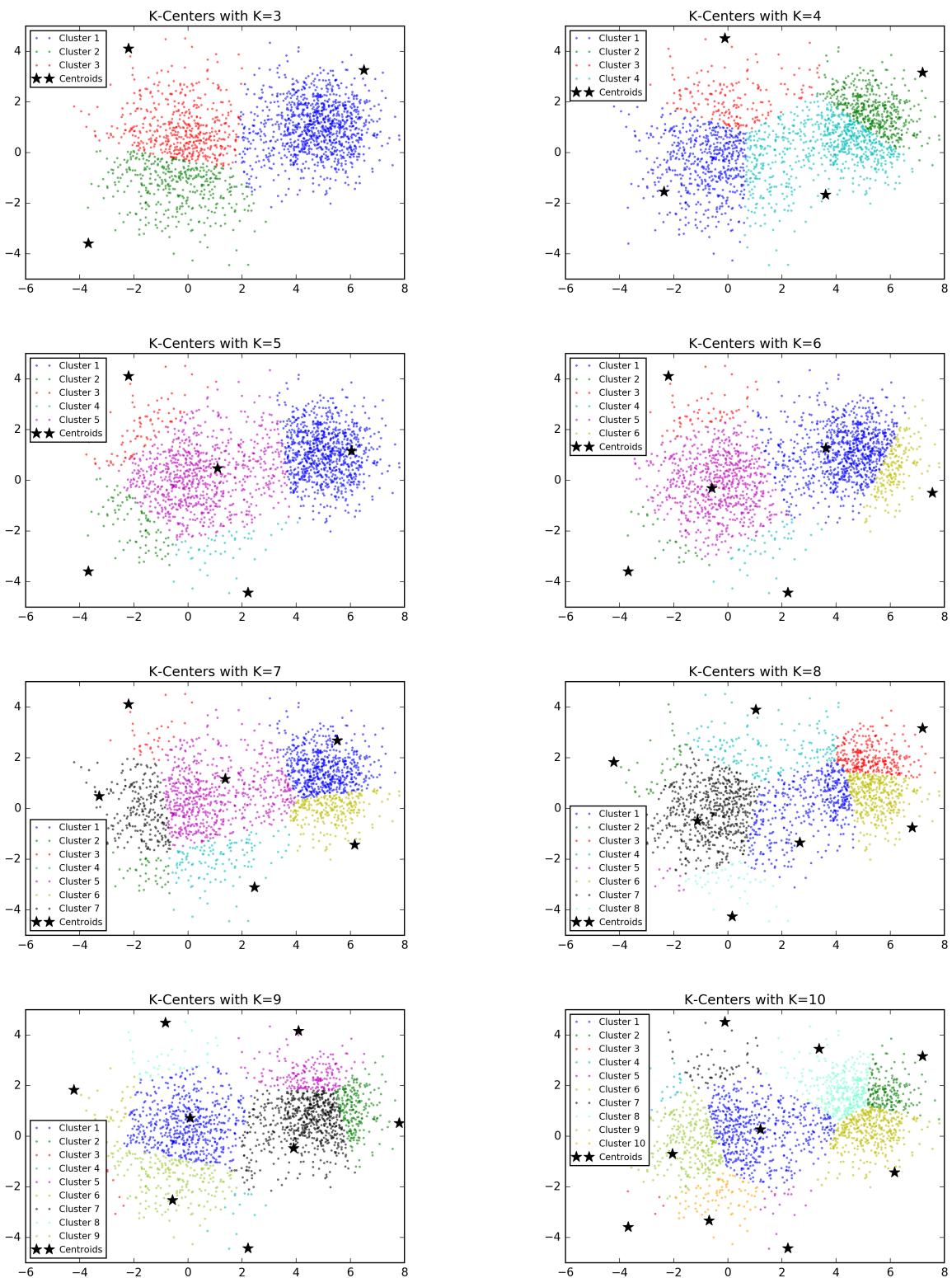


Figure 5: Clustering Result for clustering.txt with K-Center Algorithm

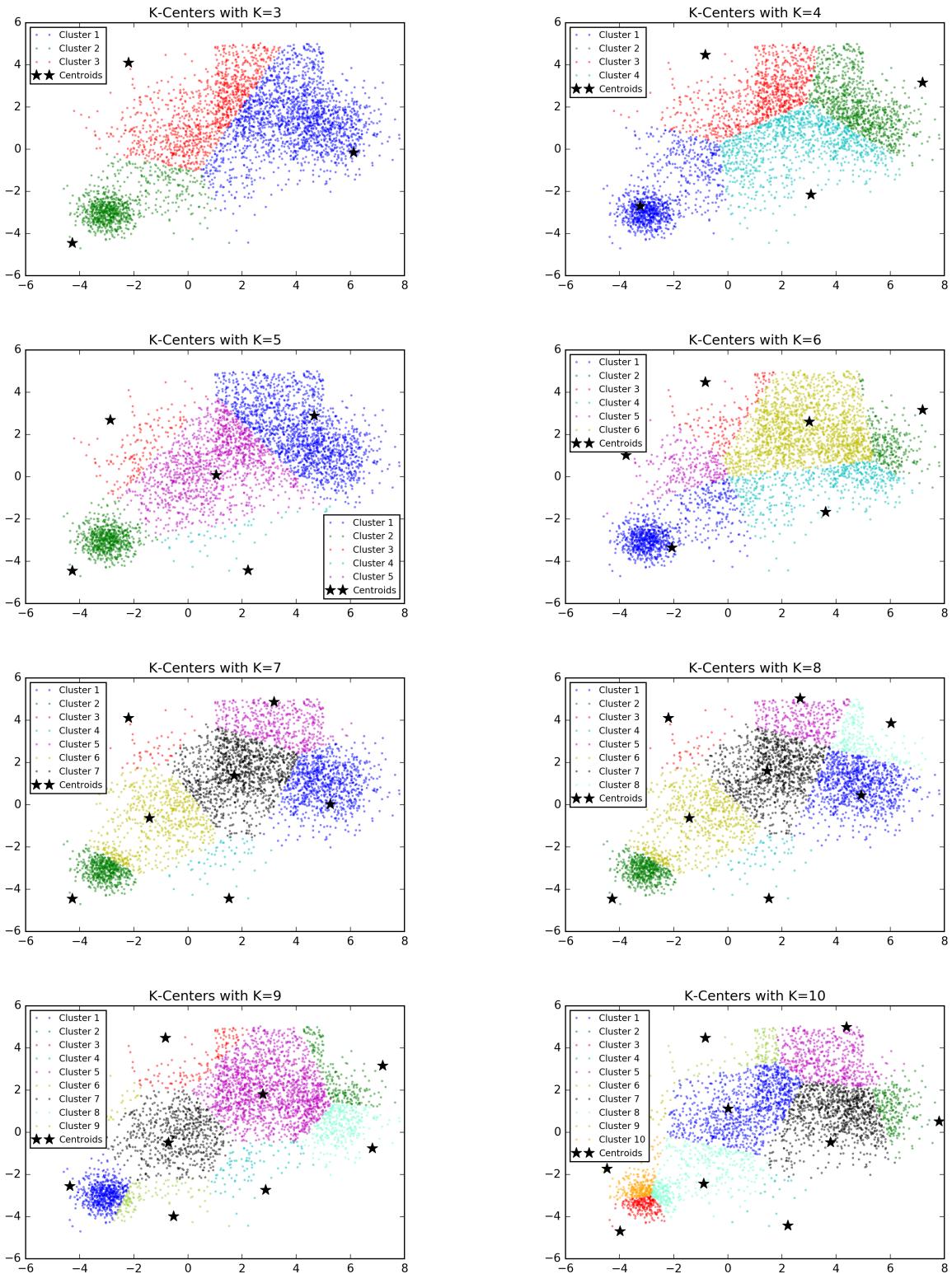


Figure 6: Clustering Result for bigClusteringData.txt with K-Center Algorithm

3. The Python code used for K-Centers Algorithm is shown in Listing 2.

```
1 import numpy as np
2 import time
3
4 def kCenters(X, K, random_state=None, verbose=True):
5     """ function to implement the greedy k-centers algorithm """
6     np.random.seed(random_state)
7     t0 = time.time()
8
9     N, d = X.shape
10    # find the initial center
11    index = np.random.choice(range(N), size=1)
12    Q = np.zeros((K, d))
13    Q[0, :] = X[index, :]
14    idx = [index]
15
16    i = 1
17    while i < K:
18        distance = np.zeros((N, i))
19        for j in range(i):
20            distance[:, j] = np.sum((X - Q[j, :]) ** 2, axis=1)
21        min_distance = np.min(distance, axis=1)
22        new_index = np.argmax(min_distance)
23        idx.append(new_index)
24        Q[i, :] = X[new_index, :]
25        i += 1
26
27    loss = np.zeros((N, K))
28    for i in range(K):
29        loss[:, i] = np.sqrt(np.sum((X - Q[i, :]) ** 2, axis=1))
30    D = np.max(np.min(loss, axis=1))
31    C = np.argmin(loss, axis=1)
32
33    if verbose is True:
34        t = np.round(time.time() - t0, 4)
35        print('K-Centers is finished in ' + str(t) + 's')
36
37    return Q, C, D, idx
```

Listing 2: K-Centers Algorithm Python Code

III. Single-Swap Algorithm

1. In this part, with the Lloyd's algorithm for k-means clustering, we choose the best distortion D as the the objection function. The change of D versus cluster number K is shown in Fig. 7, where the result for clustering.txt is shown in Fig. 7a and the result for bigClusteringData.txt is show in Fig. 7b.

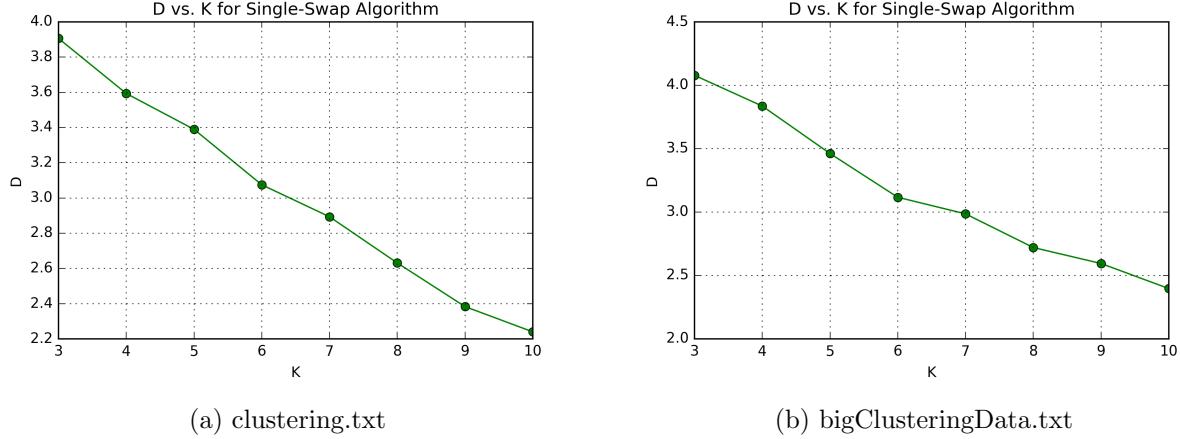


Figure 7: Change of Distortion versus Cluster Number K for Single-Swap Algorithm

2. The scatter plot of the clustering result for clustering.txt is shown in Fig. 8 and the scatter plot of the clustering result for bigClusteringData.txt is shown in Fig. 9. The cluster centroids are clearly marked and different clusters are denoted by different colors.

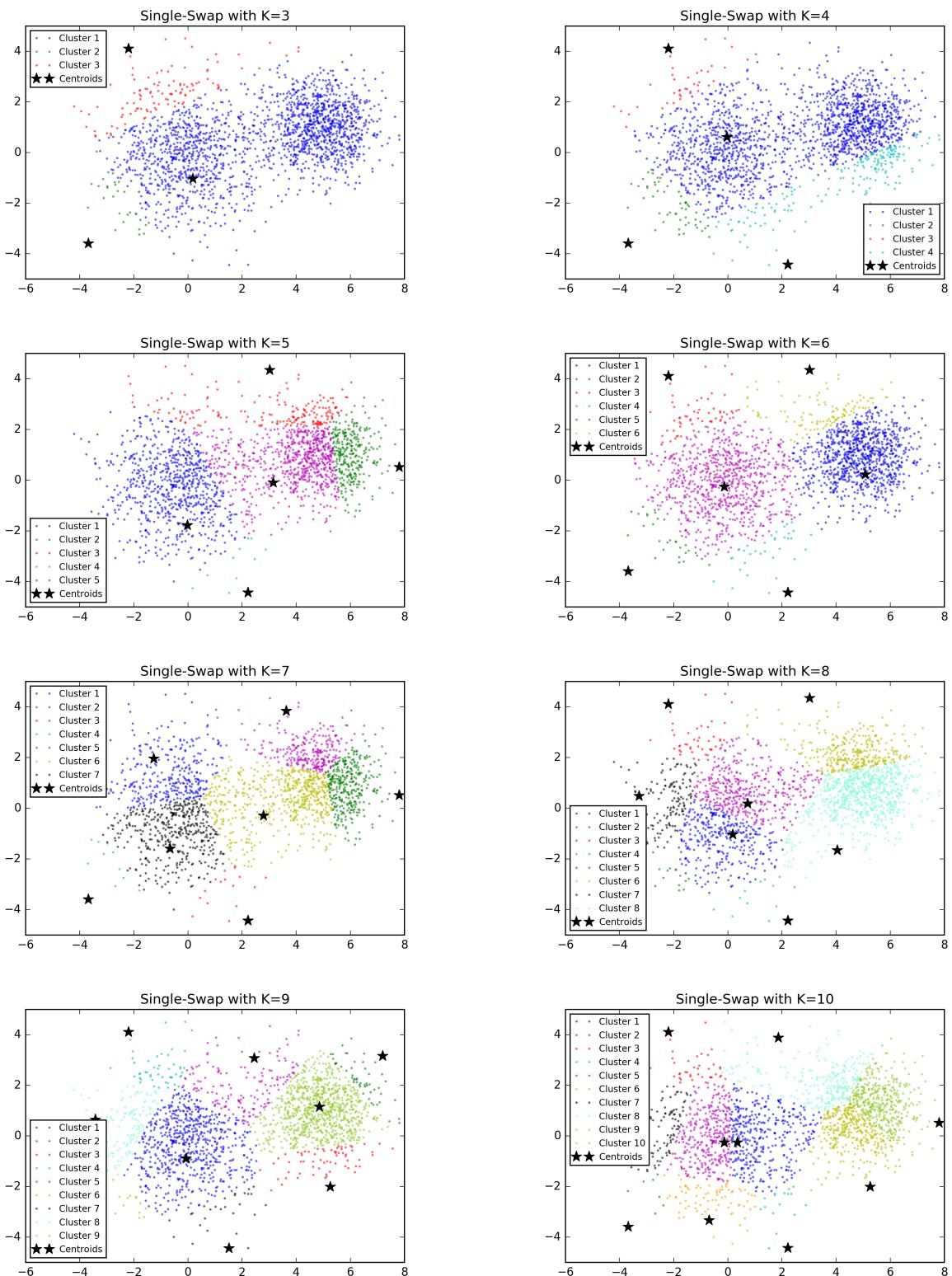


Figure 8: Clustering Result for clustering.txt with Single-Swap Algorithm

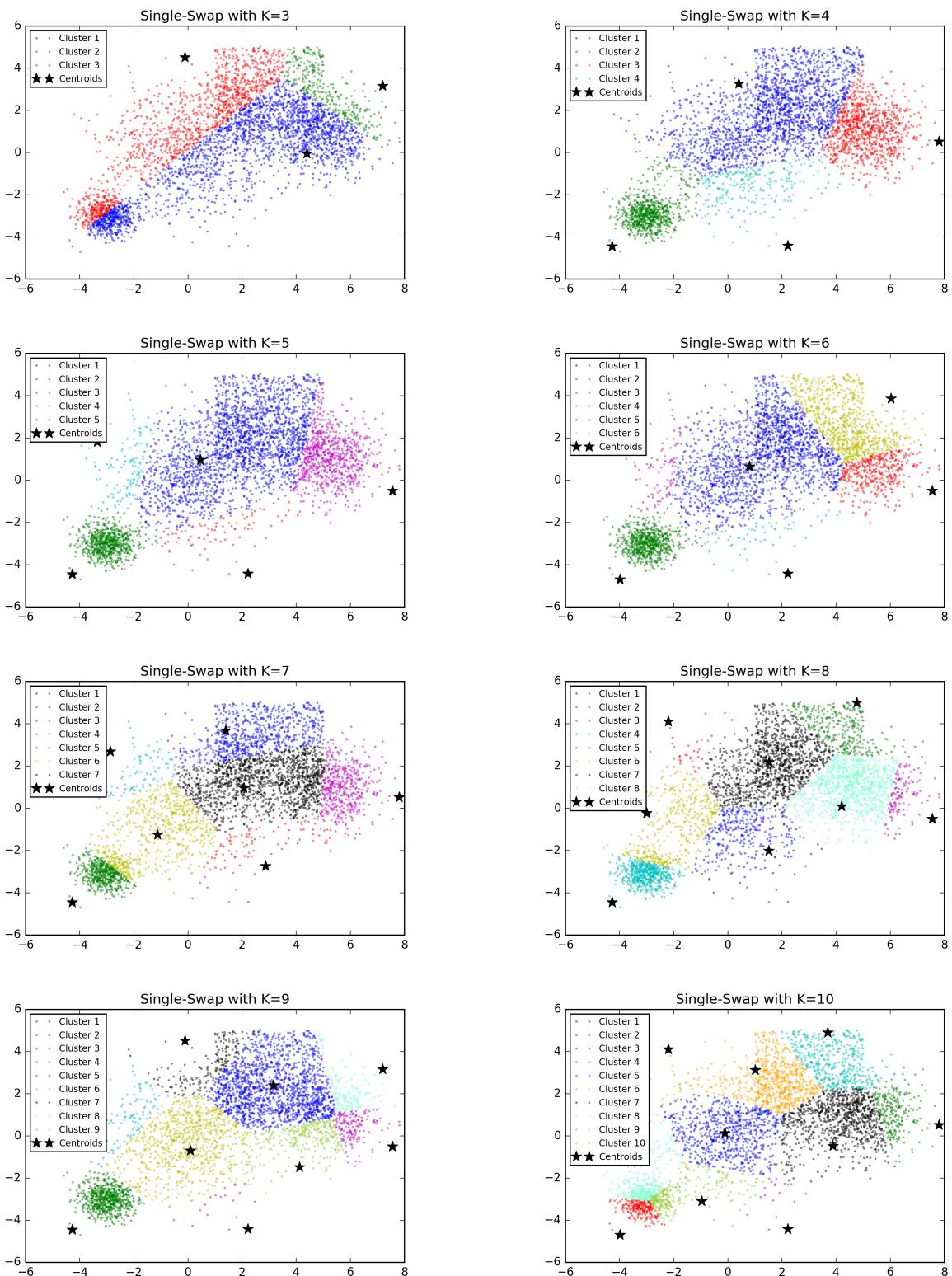


Figure 9: Clustering Result for bigClusteringData.txt with Single-Swap Algorithm

xi

3. The Python code used for K-Centers Algorithm is shown in Listing 3.

```

1 import numpy as np
2 import time
3 from k_centers import kCenters
4
5 def singleSwap(X, K, tau=0.05, random_state=None, verbose=True):
6     """ function to implement the single-swap for k-centers algorithm """
7     t0 = time.time()
8
9     # calculate the initial centers
10    Q, _, _ = kCenters(X, K, random_state=random_state,
11                        verbose=False)
12    N, d = X.shape
13
14    # compute the distance based on current centers
15    distance = np.zeros((N, K))
16    for idx in range(K):
17        distance[:, idx] = np.sqrt(np.sum((X - Q[idx, :]) ** 2, axis=1))
18    cost = np.max(np.min(distance, axis=1)) # calculate cost
19
20    i = 0
21    while i < K:
22        if i == 0:
23            min_dist = np.min(distance[:, 0:], axis=1)
24        elif i == (K - 1):
25            min_dist = np.min(distance[:, :-1], axis=1)
26        else:
27            min_dist = np.minimum(np.min(distance[:, :i], axis=1),
28                                  np.min(distance[:, (i + 1):], axis=1))
29        swap = False # keep recording whether or not swaped
30        for j in range(N):
31            tmp_dist = np.sqrt(np.sum((X - X[j, :]) ** 2, axis=1))
32            new_cost = np.max(np.minimum(min_dist, tmp_dist))
33            if new_cost / cost < (1 - tau):
34                Q[i, :] = X[j, :]
35                distance[:, i] = tmp_dist
36                swap = True
37                cost = new_cost
38        i += 1
39
40        if swap is False:
41            if i == K - 1:
42                break
43            else:
44                i += 1
45        elif (swap is True) and (i == K):
46            i = 0
47
48    C = np.argmin(distance, axis=1)
49
50    if verbose is True:
51        t = np.round(time.time() - t0, 4)
52        print('Single-Swap is finished in ' + str(t) + 's')
53
54    return Q, C, cost

```

Listing 3: Single-Swap Algorithm Python Code

IV. Spectral Clustering Algorithm

1. In this part, with the Lloyd's algorithm for k-means clustering, we choose the best distortion D as the the objection function. The change of D versus cluster number K is shown in Fig. 10, where the result for clustering.txt is shown in Fig. 10a and the result for bigClusteringData.txt is show in Fig. 10b.

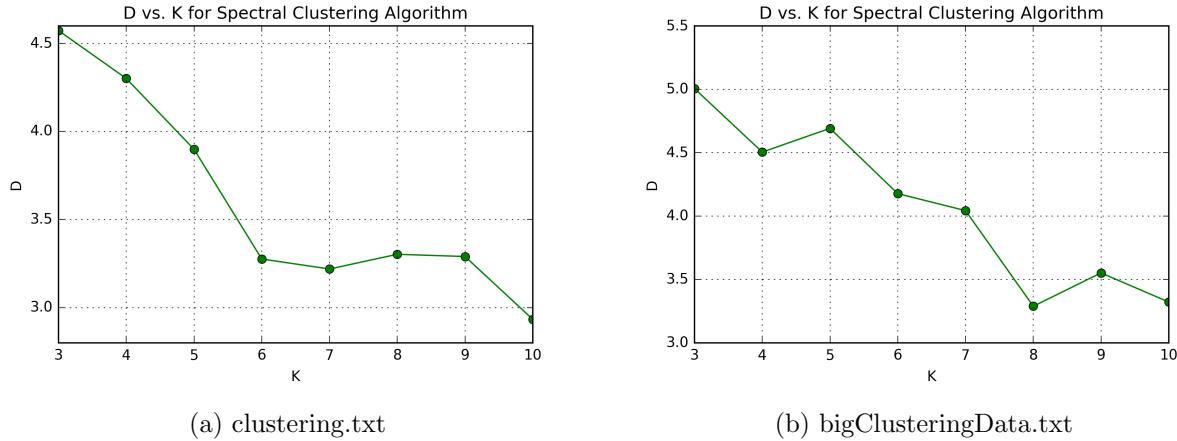


Figure 10: Change of Distortion versus Cluster Number K for Spectral Clustering

2. The scatter plot of the clustering result for clustering.txt is shown in Fig. 11 and the scatter plot of the clustering result for bigClusteringData.txt is shown in Fig. 12. The cluster centroids are clearly marked and different clusters are denoted by different colors.

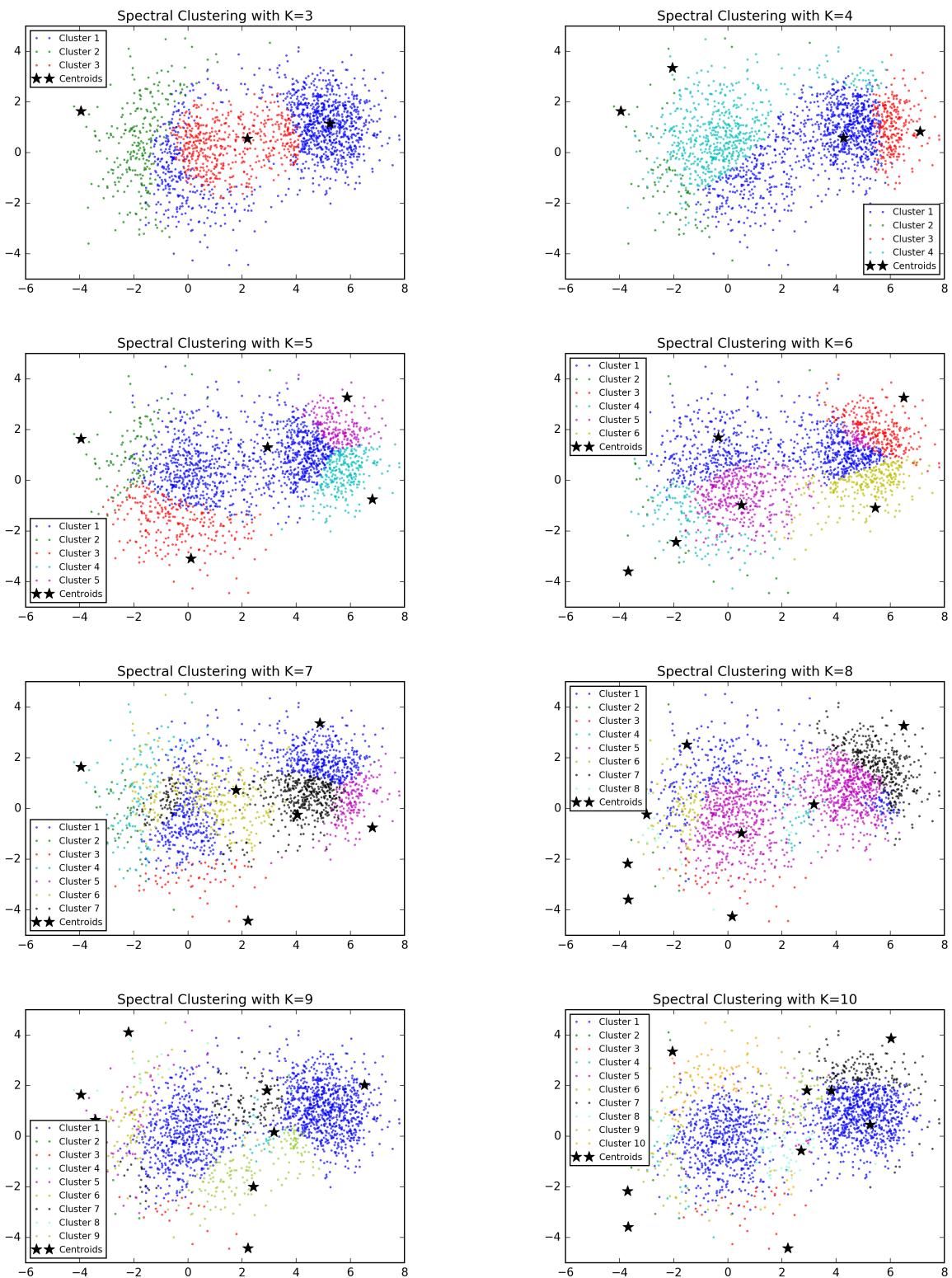


Figure 11: Clustering Result for clustering.txt with Spectral Clustering

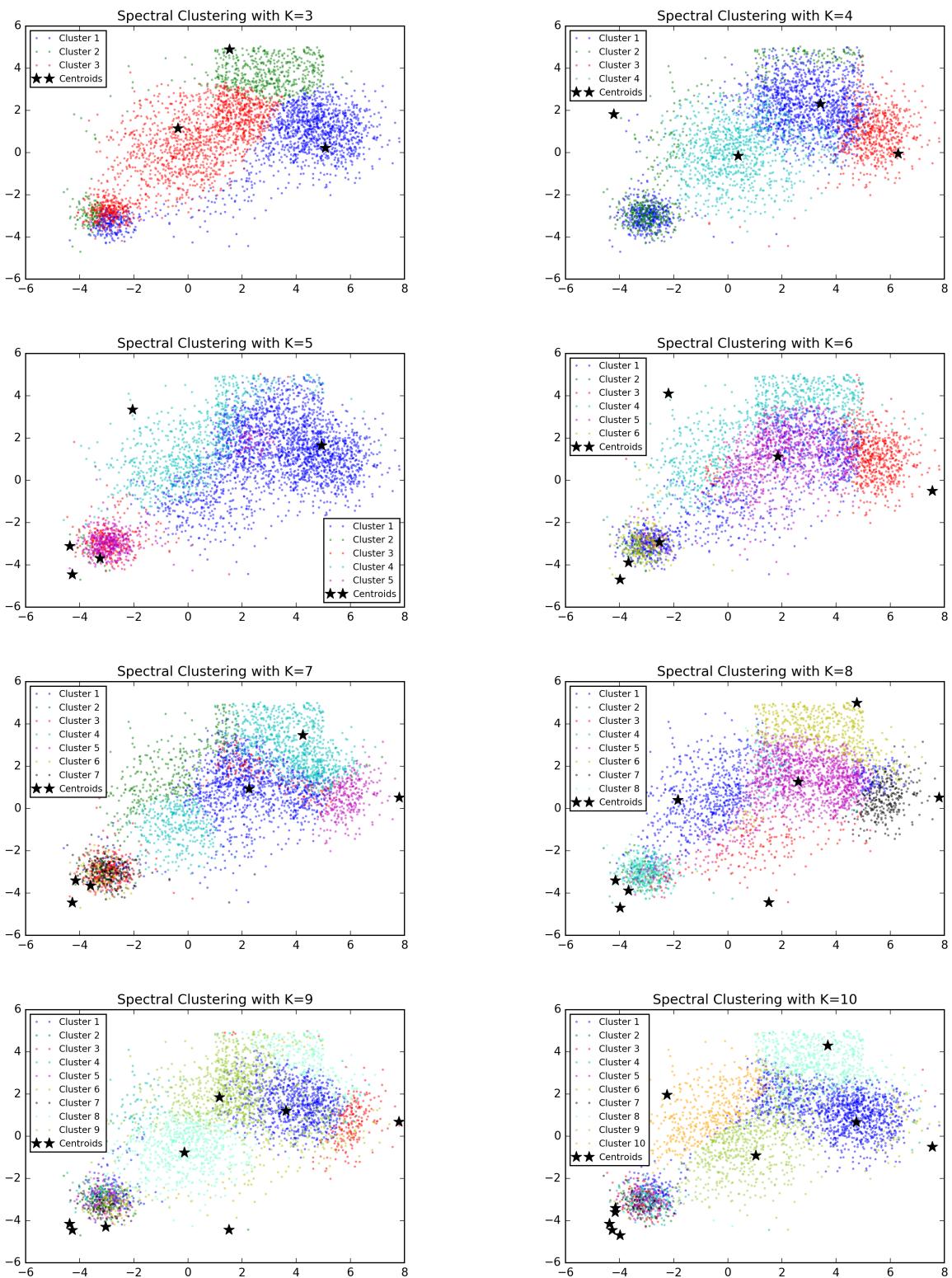


Figure 12: Clustering Result for bigClusteringData.txt with Spectral Clustering

3. The Python code used for K-Centers Algorithm is shown in Listing 4.

```
1 import numpy as np
2 import time
3 from k_centers import kCenters
4
5 def spectralClustering(X, K, random_state=None, verbose=True):
6     """ function to implement the spectral clustering algorithm """
7     t0 = time.time()
8
9     N, d = X.shape
10    W = np.zeros((N, N)) # adjacency matrix W
11    for i in range(N):
12        distance = np.sqrt(np.sum((X - X[i, :]) **2, axis=1))
13        W[:, i] = distance
14
15    diag = np.sum(W, axis=1)
16    D = np.diag(diag) # diagonal matrix D
17    L = D - W # Laplacian matrix L
18    L = np.identity(N) - np.dot(np.linalg.inv(D), W)
19    eigvals, U = np.linalg.eigh(L)
20    U = U[:, -K:] # first K eigenvectors
21
22    # call k-means for clustering
23    _, C, _, idx = kCenters(U, K, random_state=random_state, verbose=False)
24
25    Q = X[idx, :]
26    loss = np.zeros((N, K))
27    for i in range(K):
28        loss[:, i] = np.sqrt(np.sum((X - Q[i, :]) **2, axis=1))
29    D = np.max(np.min(loss, axis=1))
30
31    if verbose is True:
32        t = np.round(time.time() - t0, 4)
33        print('Spectral Clustering finished in ' + str(t) + 's')
34
35    return W, U, Q, C, D
```

Listing 4: Spectral Clustering Algorithm Python Code

V. Expectation Maximization (EM) Algorithm

1. In this part, with the Lloyd's algorithm for k-means clustering, we choose the best distortion D as the the objection function. The change of D versus cluster number K is shown in Fig. 13, where the result for clustering.txt is shown in Fig. 13a and the result for bigClusteringData.txt is show in Fig. 13b.

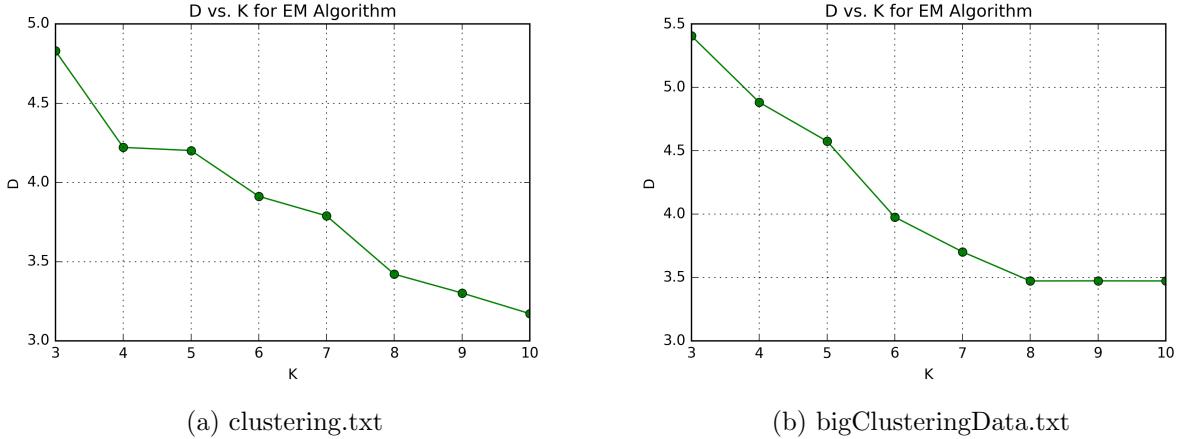


Figure 13: Change of Distortion versus Cluster Number K for EM Algorithm

2. The scatter plot of the clustering result for clustering.txt is shown in Fig. 14 and the scatter plot of the clustering result for bigClusteringData.txt is shown in Fig. 15. The cluster centroids are clearly marked and different clusters are denoted by different colors.

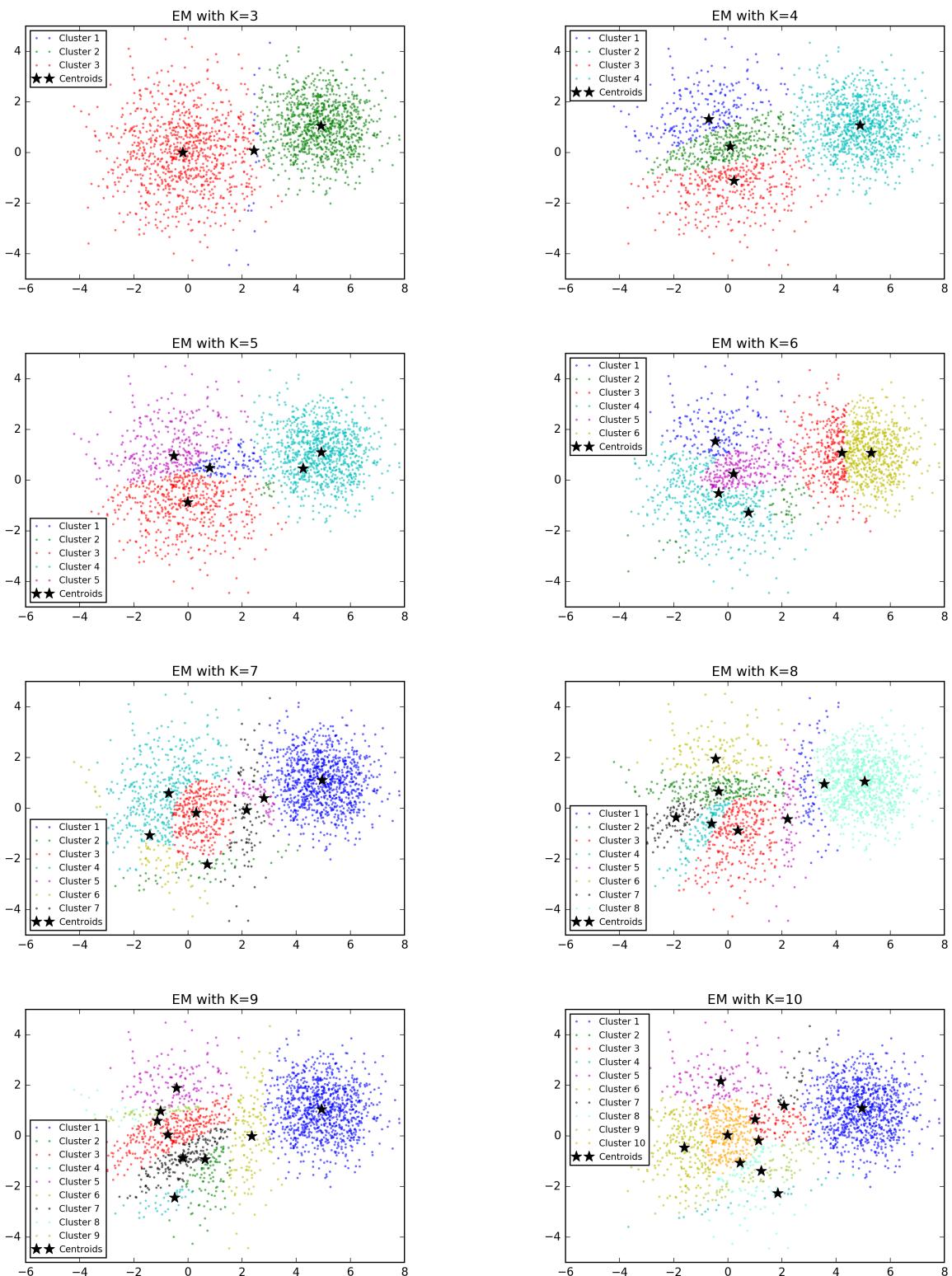


Figure 14: Clustering Result for clustering.txt with EM Algorithm

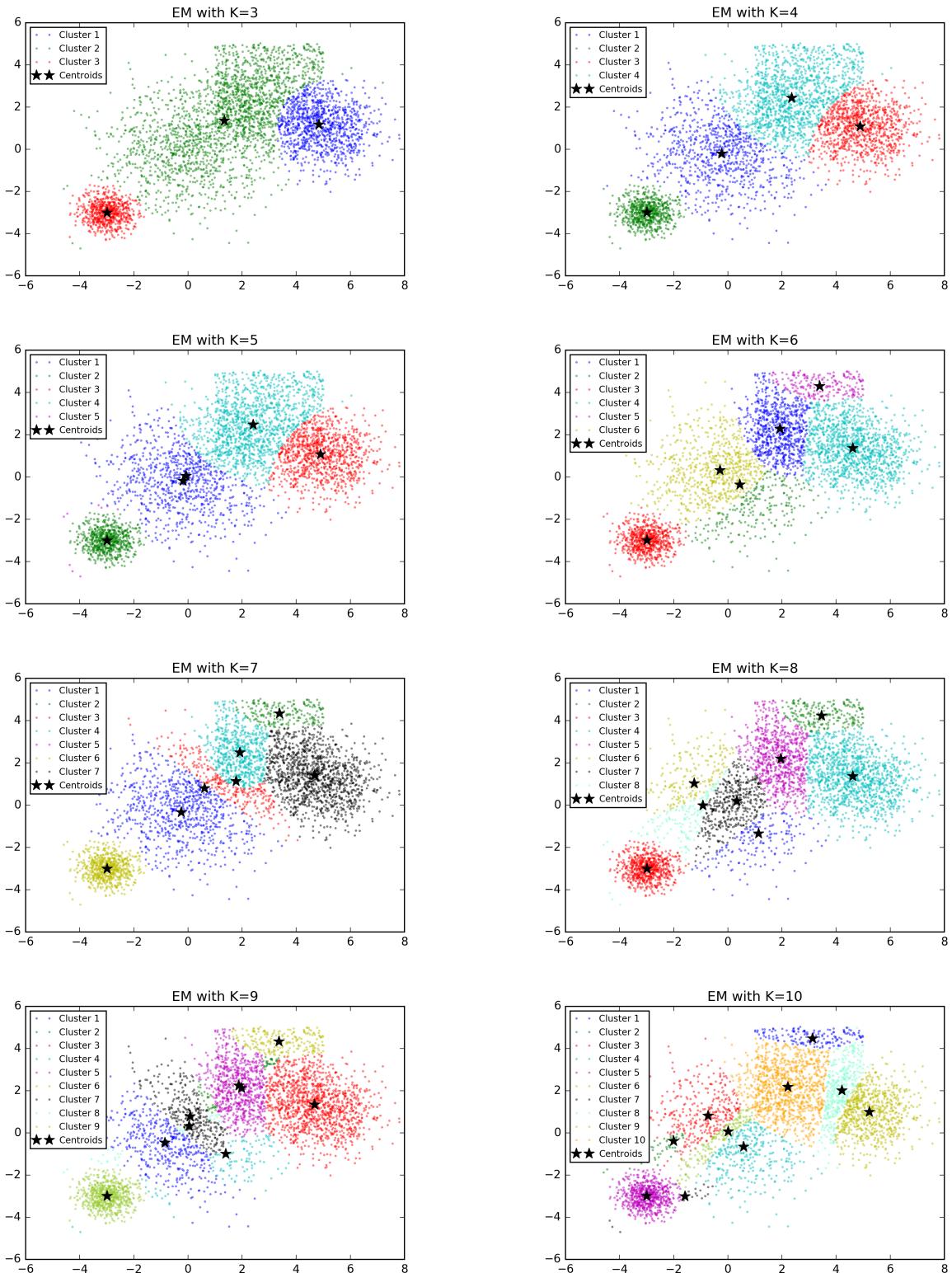


Figure 15: Clustering Result for bigClusteringData.txt with EM Algorithm

3. The Python code used for K-Centers Algorithm is shown in Listing 5.

```

1 import numpy as np
2 import time
3
4 class EM(object):
5     """ self-defined class for EM algorithm """
6     def __init__(self, m, threshold=0.01, random_state=None, maxIter=500):
7         """ initialize the EM algorithm """
8         self.m = m
9         self.threshold = threshold
10        self.random_state = random_state
11        self.maxIter = maxIter
12        self.w = None
13        self.gamma = None
14        self.mu = None
15        self.sigma = None
16        self.gaussianProb = None
17        self.logLikelihood = None
18        self.distance = None
19        self.D = None
20
21    def train(self, x, verbose=False):
22        """ function to perform EM algorithm on X """
23        t0 = time.time()
24        np.random.seed(self.random_state)
25
26        # initialize the mean and covariance matrix
27        self.initialize(x)
28
29        # iterate through E and M steps
30        for i in range(1, self.maxIter + 1):
31            self.estep(x)
32            self.mstep(x)
33            if abs(self.logLikelihood[-1] - self.logLikelihood[-2]) \
34                / abs(self.logLikelihood[-2]) < self.threshold:
35                for i in range(self.m):
36                    self.distance[:, i] = np.sqrt(np.sum((x - self.mu[i]) \
37                        **2,
38                                         axis=1))
39                    self.D = np.max(np.min(self.distance, axis=1))
40                    if verbose is True:
41                        t = np.round(time.time() - t0, 4)
42                        print('Reach threshold at', i,
43                              'th iters in ' + str(t) + 's')
44                    return
45
46                    for i in range(self.m):
47                        self.distance[:, i] = np.sqrt(np.sum((x - self.mu[i]) **2, axis
48                        =1))
49                    self.D = np.max(np.min(self.distance, axis=1))
50                    if verbose is True:
51                        t = np.round(time.time() - t0, 4)
52                        print('Stopped, reach the maximum iteration ' + str(t) + 's')
53
54    def initialize(self, x):
55        """ function to initialize the parameters """
56        n, dim = x.shape # find the dimensions
57        self.distance = np.zeros((n, self.m))

```

```

56     self.w = np.ones(self.m) * (1 / self.m)
57     self.gamma = np.zeros((n, self.m))
58     self.gaussianProb = np.zeros((n, self.m))
59     self.mu = [None] * self.m
60     self.sigma = [None] * self.m
61     self.logLikelihood = []
62
63     cov = np.cov(x.T)
64     mean = np.mean(x, axis=0)
65     for k in range(self.m):
66         self.mu[k] = mean + np.random.uniform(-0.5, 0.5, dim)
67         self.sigma[k] = cov
68
69     # update gamma
70     self.gamma = self.gammaprob(x, self.w, self.mu, self.sigma)
71
72     # calculate the expectation of log-likelihood
73     self.logLikelihood.append(self.likelihood())
74
75 def estep(self, x):
76     """ function to conduct E-Step for EM algorithm """
77     self.gamma = self.gammaprob(x, self.w, self.mu, self.sigma)
78
79 def mstep(self, x):
80     """ function to conduct M-Step for EM algorithm """
81     n, dim = x.shape
82     sumGamma = np.sum(self.gamma, axis=0)
83     self.w = sumGamma / n
84
85     for k in range(self.m):
86         self.mu[k] = np.sum(x.T * self.gamma[:, k], axis=1) / sumGamma[k]
87         diff = x - self.mu[k]
88         weightedDiff = diff.T * self.gamma[:, k]
89         self.sigma[k] = np.dot(weightedDiff, diff) / sumGamma[k]
90         if np.linalg.matrix_rank(self.sigma[k]) != 3:
91             randv = np.random.random(dim) / 10000
92             self.sigma[k] = self.sigma[k] + np.diag(randv)
93
94     # calculate the expectation of log-likelihood
95     self.logLikelihood.append(self.likelihood())
96
97 def gammaprob(self, x, w, mu, sigma):
98     """ function to calculate the gamma probability """
99     for k in range(self.m):
100         self.gaussianProb[:, k] = self.gaussian(x, mu[k], sigma[k])
101
102     weightedSum = np.sum(w * self.gaussianProb, axis=1)
103     gamma = ((w * self.gaussianProb).T / weightedSum).T
104
105     return gamma
106
107 def gaussian(self, x, mu, sigma):
108     """ function to calculate the multivariate gaussian probability """
109     inversion = np.linalg.inv(sigma)
110     part1 = (-0.5 * np.sum(np.dot(x - mu, inversion) * (x - mu), axis
111 =1))
112     part2 = 1 / ((2 * np.pi) ** (len(mu) / 2) *
113                  (np.linalg.det(sigma) ** 0.5))

```

```

113
114     pdf = part2 * np.exp(part1)
115
116     return pdf
117
118 def likelihood(self):
119     """ function to calculate the log likelihood """
120     log = np.log(np.sum(self.w * self.gaussianProb, axis=1))
121     logLikelihood = np.sum(log)
122
123     return logLikelihood
124
125 def get_label(self):
126     """ function to predict the classes using calculated parameters """
127     label = np.argmax(self.w * self.gaussianProb, axis=1)
128
129     return label

```

Listing 5: EM Algorithm Python Code

2 Natural Clusters Discussion

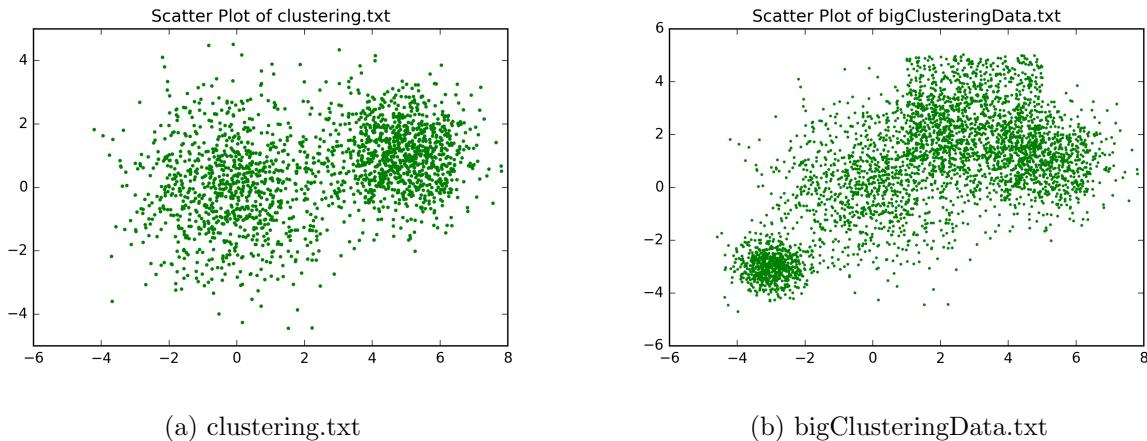


Figure 16: Scatter Plot of Original Data

3 K-Means Convergence Discussion

4 Computational Effort Discussion

Table 1: Computation Time Comparison for clustering.txt

K	K-Means	K-Centers	Single-Swap	Spectral Clustering	EM
3	0.2673	0.0067	0.1582	1.5716	0.4412
4	0.2543	0.0010	0.1971	1.4587	0.5078
5	0.3624	0.0012	0.1895	1.4132	0.4881
6	0.4325	0.0013	0.2668	1.4840	0.7462
7	0.2837	0.0015	0.2693	1.4877	0.8227
8	0.1923	0.0019	0.3306	1.4708	0.9216
9	0.3481	0.0024	0.3012	1.4669	1.0106
10	0.1989	0.0026	0.3790	1.5053	1.1277

Note: all measured time is in unit of seconds (s)

Table 2: Computation Time Comparison for bigClusteringData.txt

K	K-Means	K-Centers	Single-Swap	Spectral Clustering	EM
3	0.4725	0.0020	0.4511	11.1860	0.3412
4	0.8326	0.0019	0.5692	10.2542	0.2553
5	0.4722	0.0021	0.6445	10.5166	0.5763
6	1.2580	0.0024	0.8809	10.2167	0.6196
7	0.9619	0.0025	0.7638	9.8986	1.1591
8	0.7767	0.0034	0.7815	9.4616	1.3177
9	0.7177	0.0037	0.9010	9.8733	1.4977
10	1.2369	0.0041	1.0744	9.8636	1.7560

Note: all measured time is in unit of seconds (s)

5 Single-Swap Algorithm Discussion