

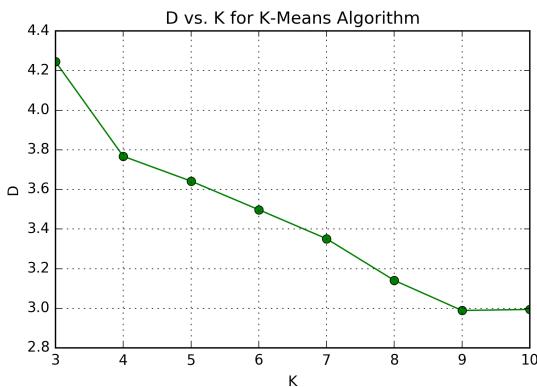
# IE 529 Fall 2016 Computational Assignment 2

Jifu Zhao

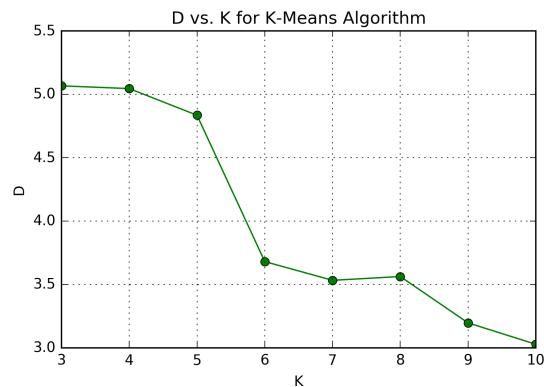
Date: 12/15/2016

## I. Lloyd's (K-means)Algorithm

1. In this part, the polynomial regression is conducted on the given dataset. The derivation is relatively simple, here we only give one simple version.



(a) clustering.txt



(b) bigClusteringData.txt

Figure 3: Change of Distortion versus Cluster Number K for K-Means Algorithm

```

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, :] # initial k centers
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

```

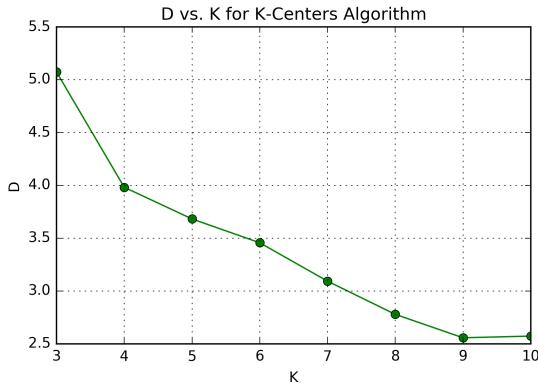
```

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) + ' iters')
39
40     return Y, C, D

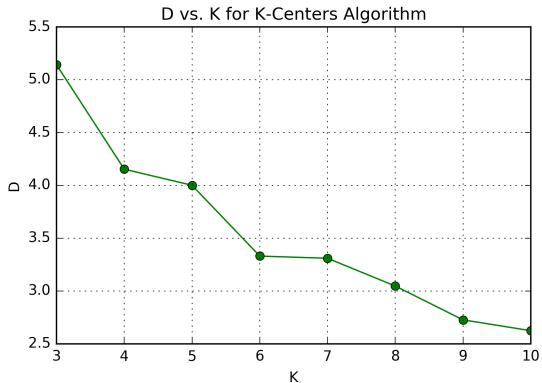
```

Listing 1: K-Means Algorithm Python Code

## II. Greedy K-centers Algorithm



(a) clustering.txt



(b) bigClusteringData.txt

Figure 6: Change of Distortion versus Cluster Number K for K-Center Algorithm

```

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

```

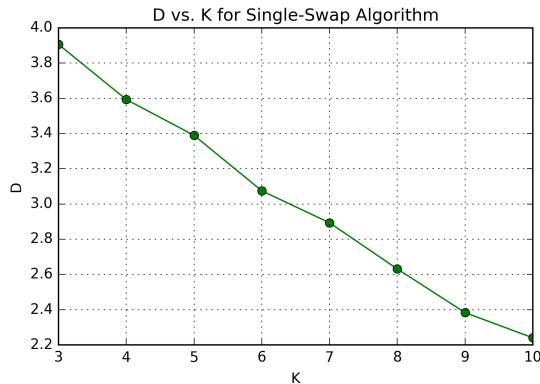
```

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

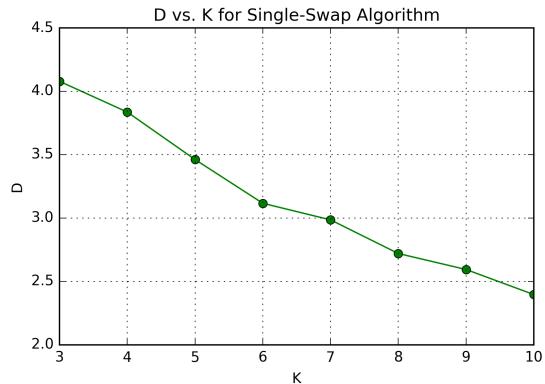
```

Listing 2: K-Centers Algorithm Python Code

### III. Single-Swap Algorithm



(a) clustering.txt



(b) bigClusteringData.txt

Figure 9: Change of Distortion versus Cluster Number K for Single-Swap Algorithm

```

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, _, pre_cost, _ = kCenters(X, K, random_state=random_state,
11                                  verbose=False)
12    N, d = X.shape

```

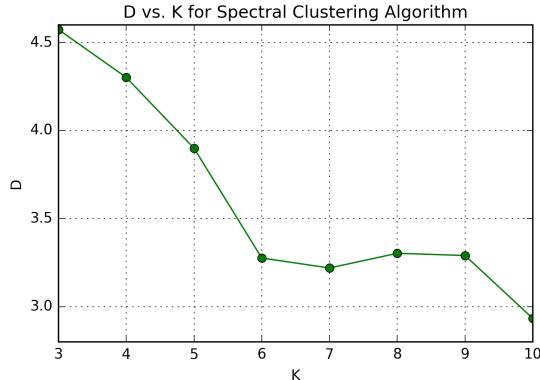
```

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

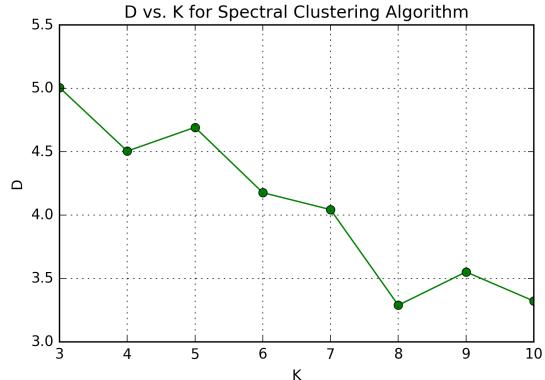
```

Listing 3: Single-Swap Algorithm Python Code

## IV. Spectral Clustering Algorithm



(a) clustering.txt



(b) bigClusteringData.txt

Figure 12: Change of Distortion versus Cluster Number K for Spectral Clustering

```

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

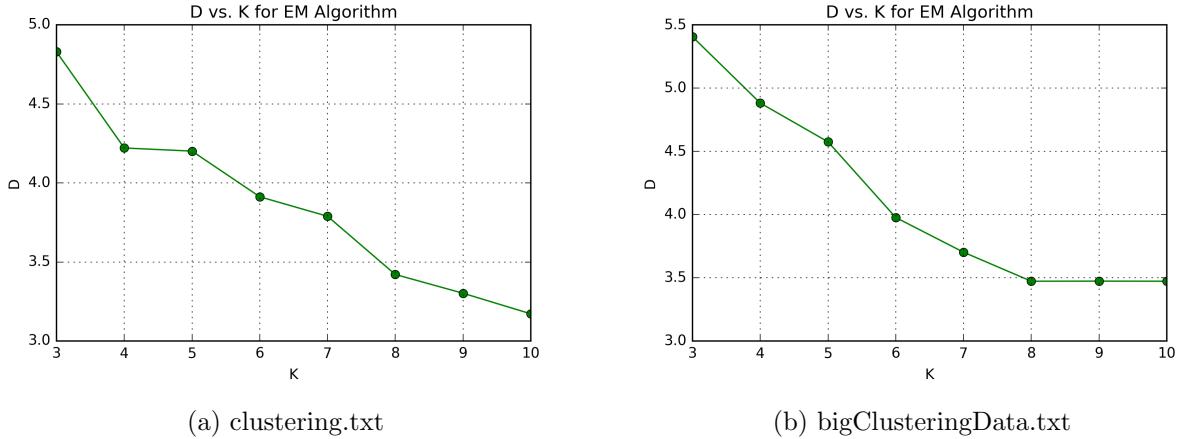


Figure 15: Change of Distortion versus Cluster Number K for EM Algorithm

```
1 import numpy as np
2 import time
3
4 class EM(object):
5     """ self-defined calss for EM algorithm """
6
7     def __init__(self, m, threshold=0.01, random_state=None, maxIter=500):
8         """ initialize the EM algorithm """
9         self.m = m
10        self.threshold = threshold
11        self.random_state = random_state
12        self.maxIter = maxIter
13        self.w = None
14        self.gamma = None
15        self.mu = None
16        self.sigma = None
17        self.gaussianProb = None
18        self.logLikelihood = None
19        self.distance = None
20        self.D = None
21
22    def train(self, x, verbose=False):
23        """ function to perform EM algorithm on X """
24        t0 = time.time()
25        np.random.seed(self.random_state)
26
27        # initialize the mean and covariance matrix
28        self.initialize(x)
29
30        # iterate through E and M steps
31        for i in range(1, self.maxIter + 1):
32            self.estep(x)
33            self.mstep(x)
34            if abs(self.logLikelihood[-1] - self.logLikelihood[-2]) \
35            / abs(self.logLikelihood[-2]) < self.threshold:
36                for i in range(self.m):
```

```

37         self.distance[:, i] = np.sqrt(np.sum((x - self.mu[i])**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=1))
48     self.D = np.max(np.min(self.distance, axis=1))
49     if verbose is True:
50         t = np.round(time.time() - t0, 4)
51         print('Stopped, reach the maximum iteration ' + str(t) + 's')
52
53     def initialize(self, x):
54         """ function to initialize the parameters """
55         n, dim = x.shape # find the dimensions
56         self.distance = np.zeros((n, self.m))
57         self.w = np.ones(self.m) * (1 / self.m)
58         self.gamma = np.zeros((n, self.m))
59         self.gaussianProb = np.zeros((n, self.m))
60         self.mu = [None] * self.m
61         self.sigma = [None] * self.m
62         self.logLikelihood = []
63
64         cov = np.cov(x.T)
65         mean = np.mean(x, axis=0)
66         for k in range(self.m):
67             self.mu[k] = mean + np.random.uniform(-0.5, 0.5, dim)
68             self.sigma[k] = cov
69
70         # update gamma
71         self.gamma = self.gammaprob(x, self.w, self.mu, self.sigma)
72
73         # calculate the expectation of log-likelihood
74         self.logLikelihood.append(self.likelihood())
75
76     def estep(self, x):
77         """ function to conduct E-Step for EM algorithm """
78         self.gamma = self.gammaprob(x, self.w, self.mu, self.sigma)
79
80     def mstep(self, x):
81         """ function to conduct M-Step for EM algorithm """
82         n, dim = x.shape
83         sumGamma = np.sum(self.gamma, axis=0)
84         self.w = sumGamma / n
85
86         for k in range(self.m):
87             self.mu[k] = np.sum(x.T * self.gamma[:, k], axis=1) / sumGamma[k]
88             diff = x - self.mu[k]
89             weightedDiff = diff.T * self.gamma[:, k]
90             self.sigma[k] = np.dot(weightedDiff, diff) / sumGamma[k]
91             if np.linalg.matrix_rank(self.sigma[k]) != 3:
92                 randv = np.random.random(dim) / 10000
93                 self.sigma[k] = self.sigma[k] + np.diag(randv)
94
95         # calculate the expectation of log-likelihood

```

```

96         self.logLikelihood.append(self.likelihood())
97
98     def gammaprob(self, x, w, mu, sigma):
99         """ function to calculate the gamma probability """
100        for k in range(self.m):
101            self.gaussianProb[:, k] = self.gaussian(x, mu[k], sigma[k])
102
103        weightedSum = np.sum(w * self.gaussianProb, axis=1)
104        gamma = ((w * self.gaussianProb).T / weightedSum).T
105
106        return gamma
107
108    def gaussian(self, x, mu, sigma):
109        """ function to calculate the multivariate gaussian probability """
110        inversion = np.linalg.inv(sigma)
111        part1 = (-0.5 * np.sum(np.dot(x - mu, inversion) * (x - mu), axis=1))
112        part2 = 1 / ((2 * np.pi) ** (len(mu) / 2) *
113                      (np.linalg.det(sigma) ** 0.5))
114
115        pdf = part2 * np.exp(part1)
116
117        return pdf
118
119    def likelihood(self):
120        """ function to calculate the log likelihood """
121        log = np.log(np.sum(self.w * self.gaussianProb, axis=1))
122        logLikelihood = np.sum(log)
123
124        return logLikelihood
125
126    def get_label(self):
127        """ function to predict the classes using calculated parameters """
128        label = np.argmax(self.w * self.gaussianProb, axis=1)
129
130        return label

```

Listing 5: EM Algorithm Python Code

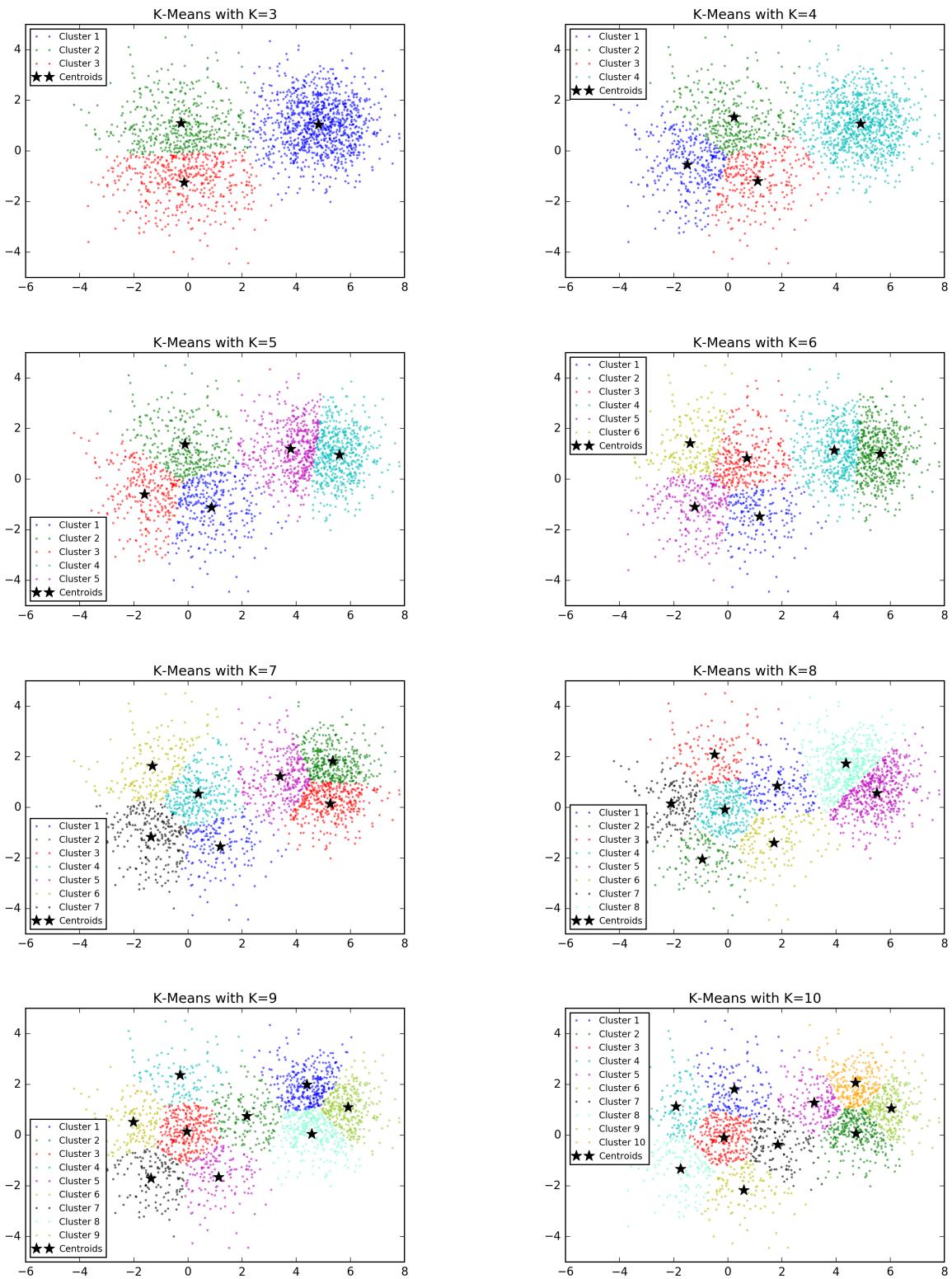


Figure 1: Clustering Result for clustering.txt with K-Means Algorithm

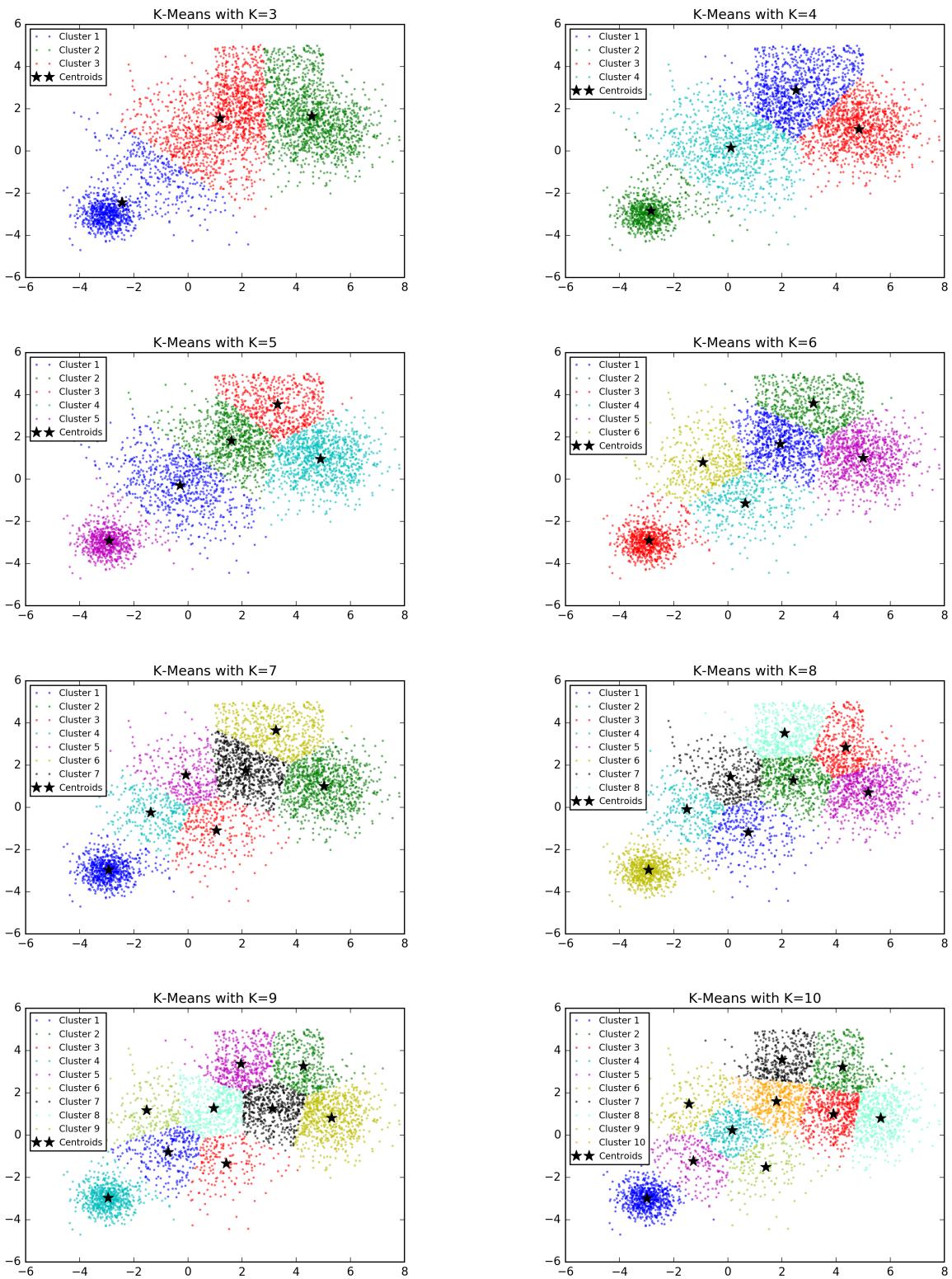


Figure 2: Clustering Result for `bigClusteringData.txt` with K-Means Algorithm

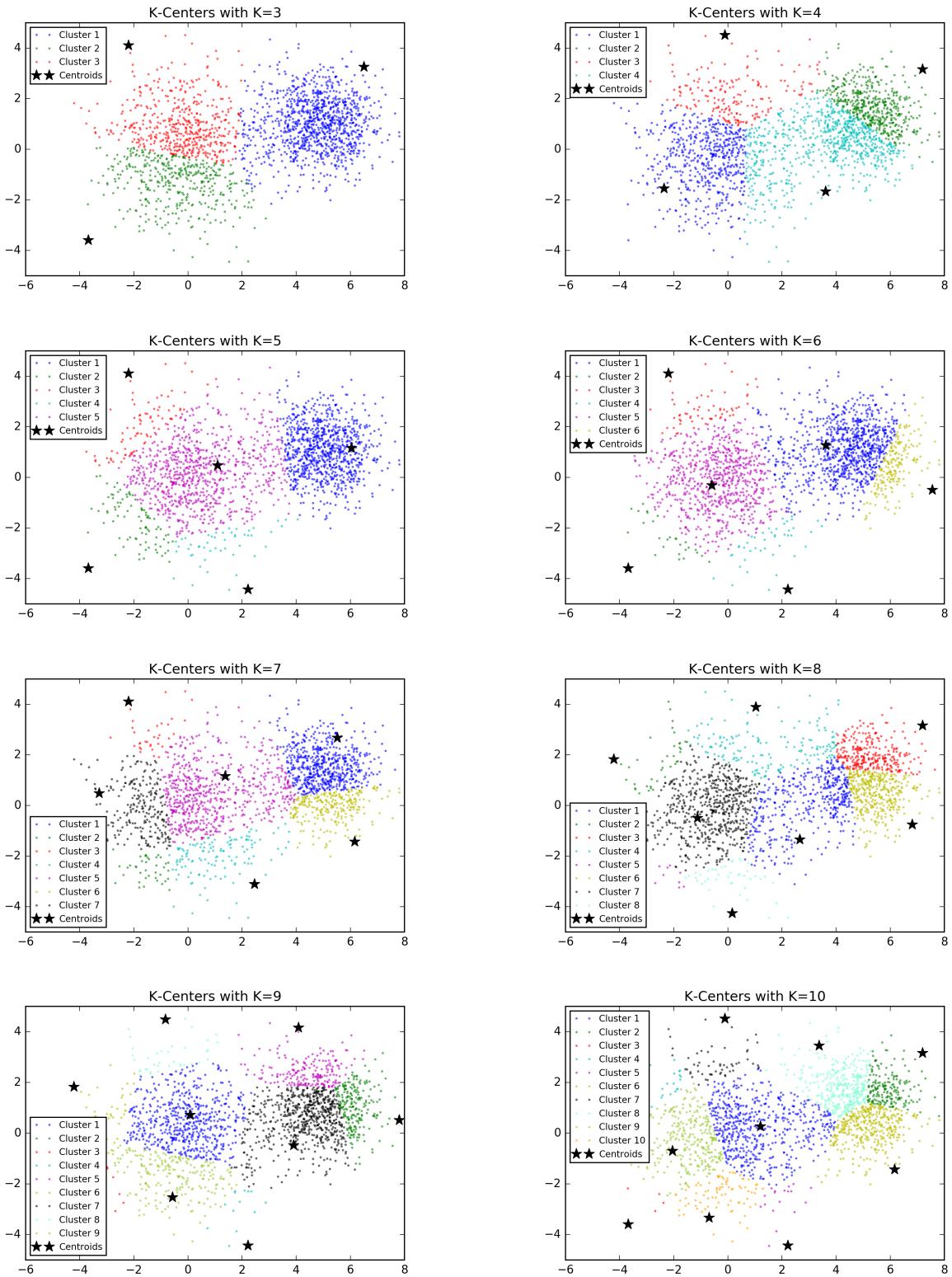


Figure 4: Clustering Result for clustering.txt with K-Center Algorithm

$x_i$

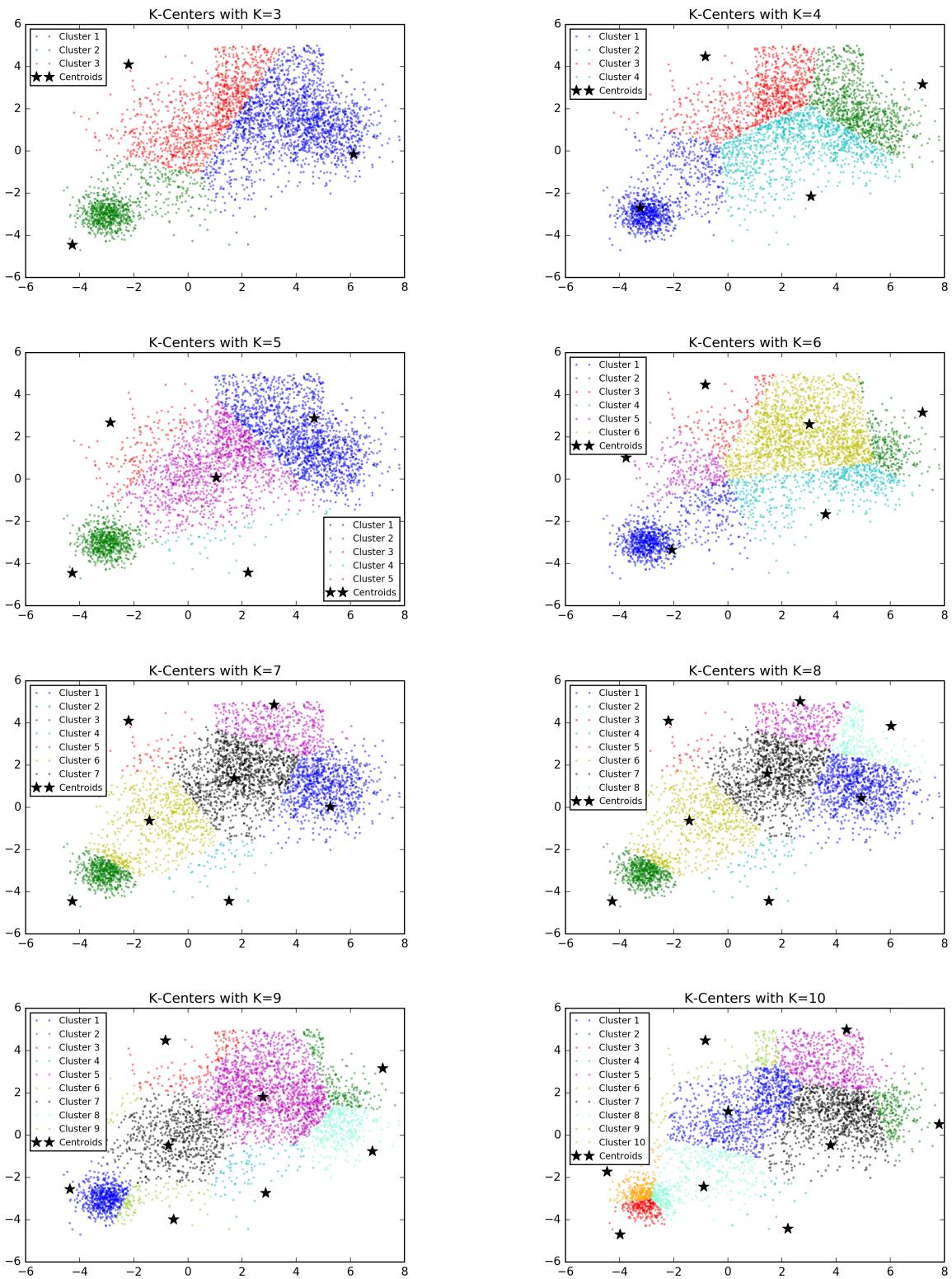


Figure 5: Clustering Result for bigClusteringData.txt with K-Center Algorithm

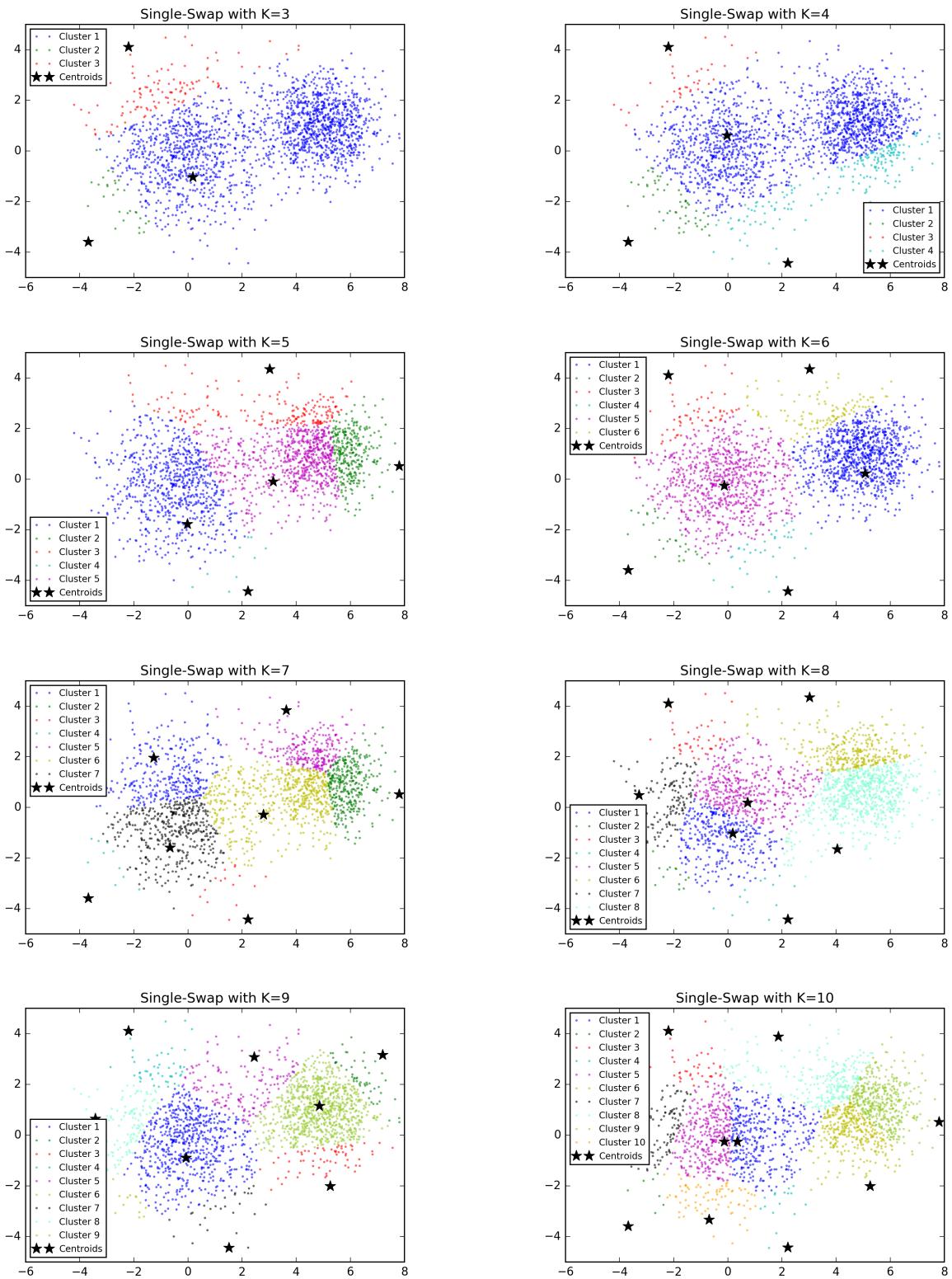


Figure 7: Clustering Result for clustering.txt with Single-Swap Algorithm

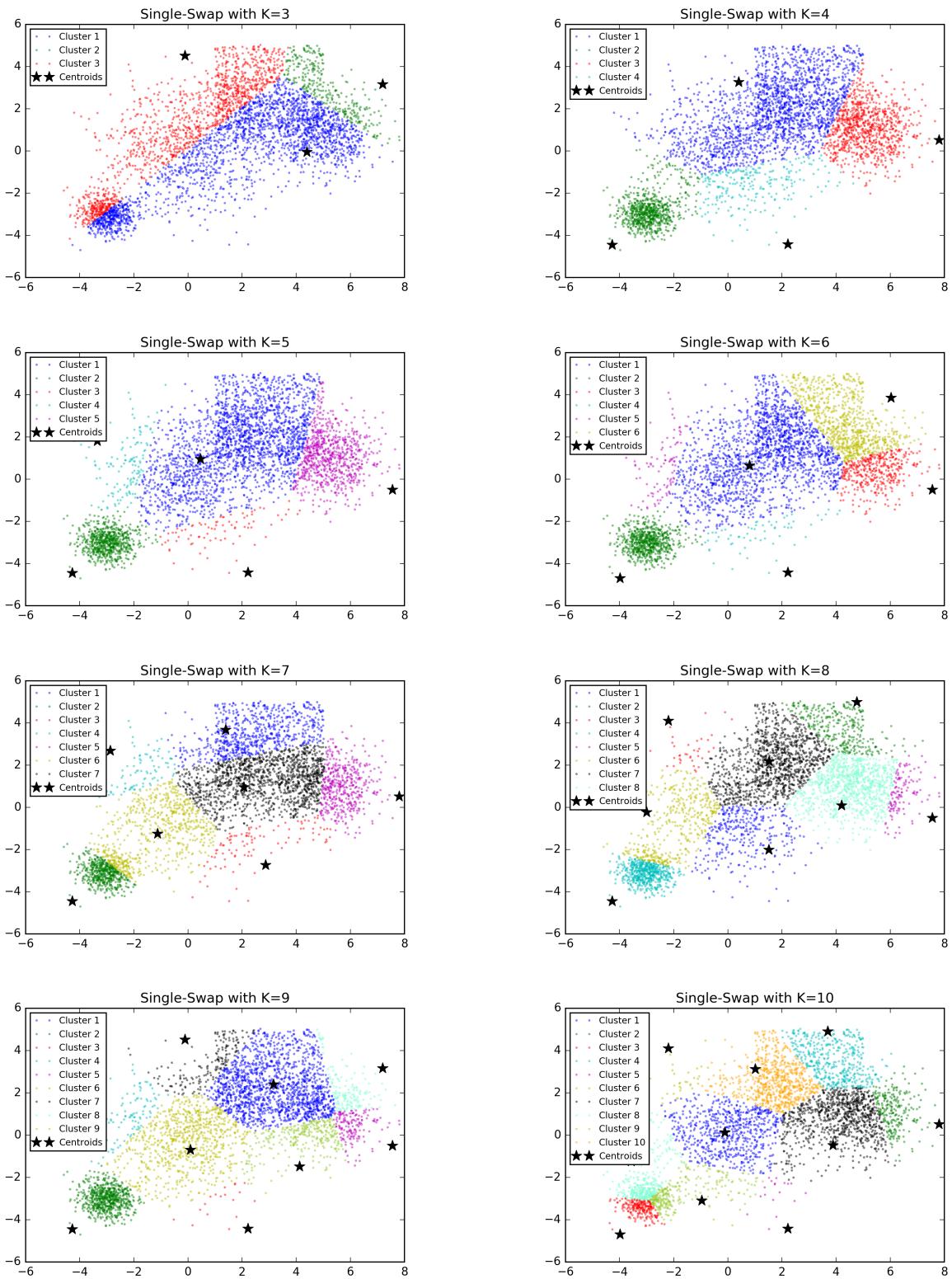


Figure 8: Clustering Result for bigClusteringData.txt with Single-Swap Algorithm

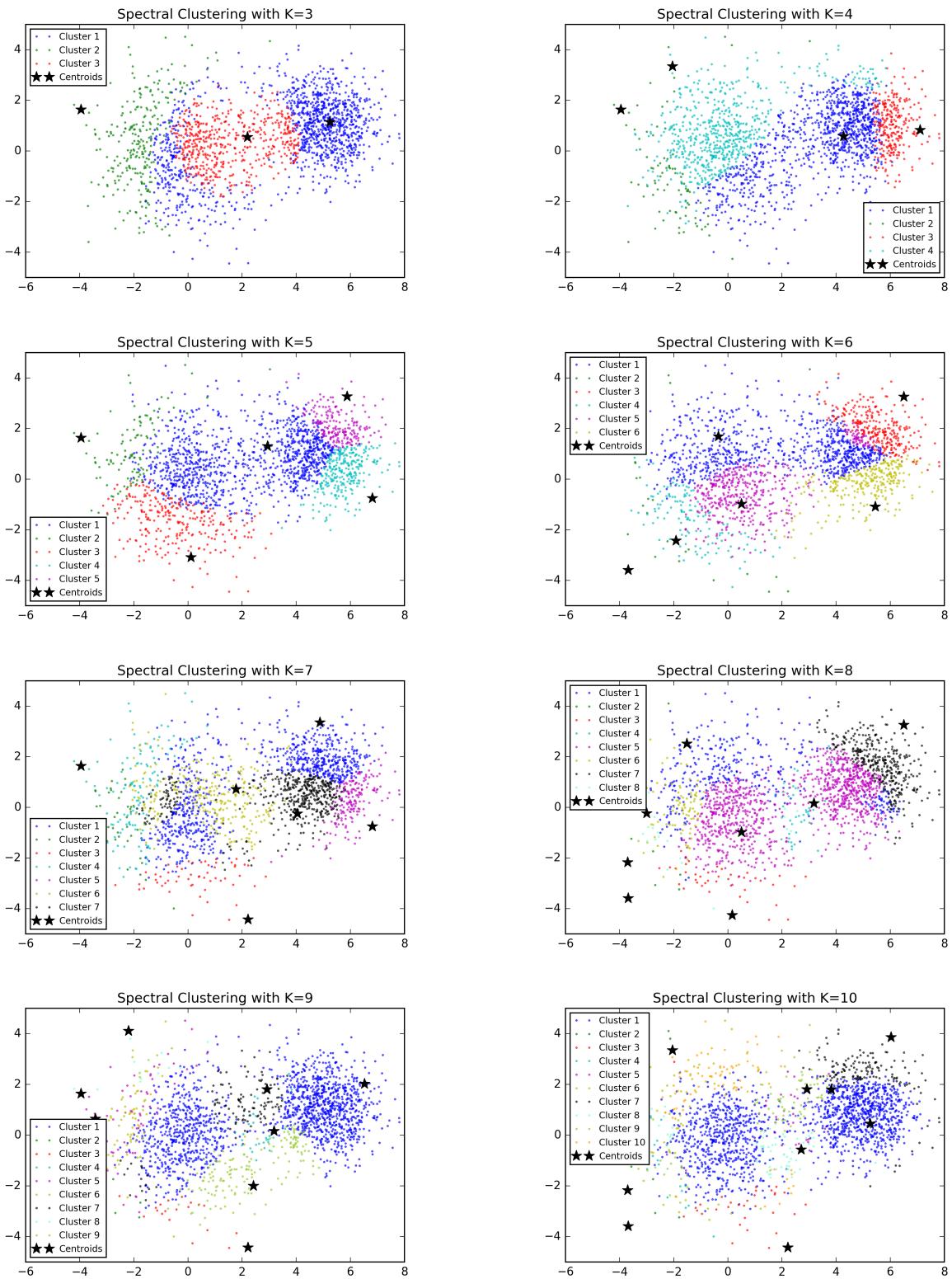


Figure 10: Clustering Result for clustering.txt with Spectral Clustering

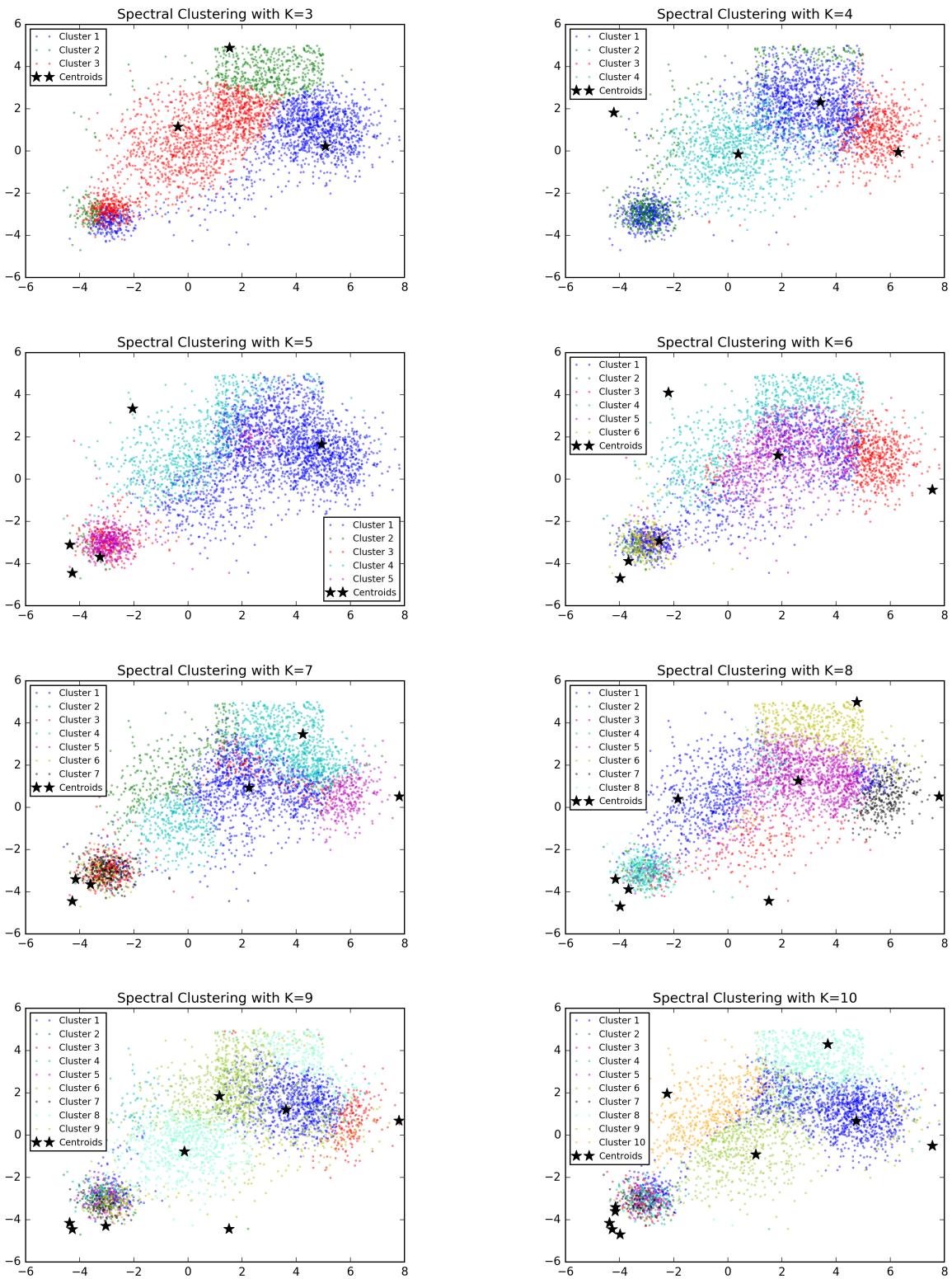


Figure 11: Clustering Result for bigClusteringData.txt with Spectral Clustering

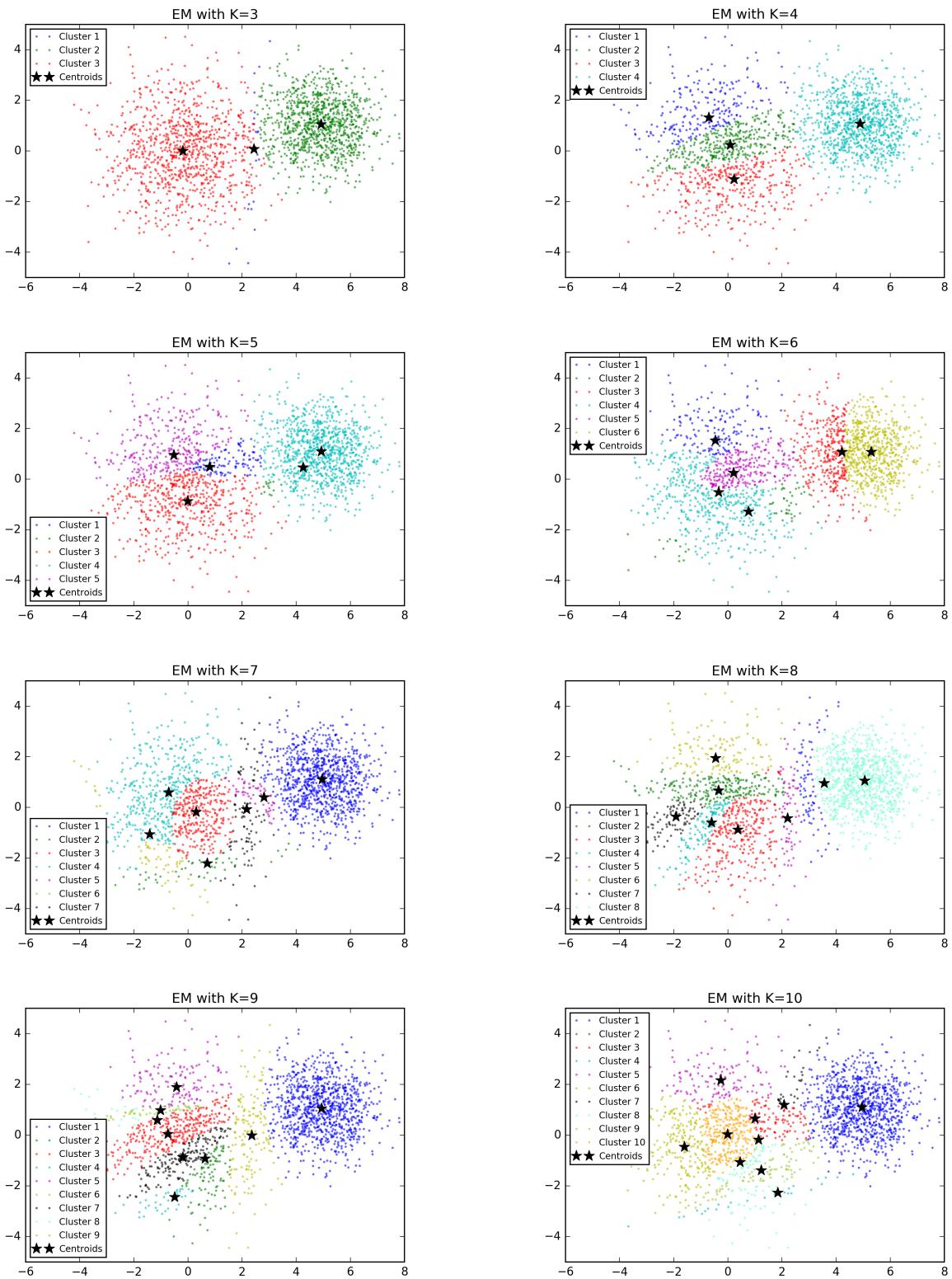


Figure 13: Clustering Result for clustering.txt with EM Algorithm

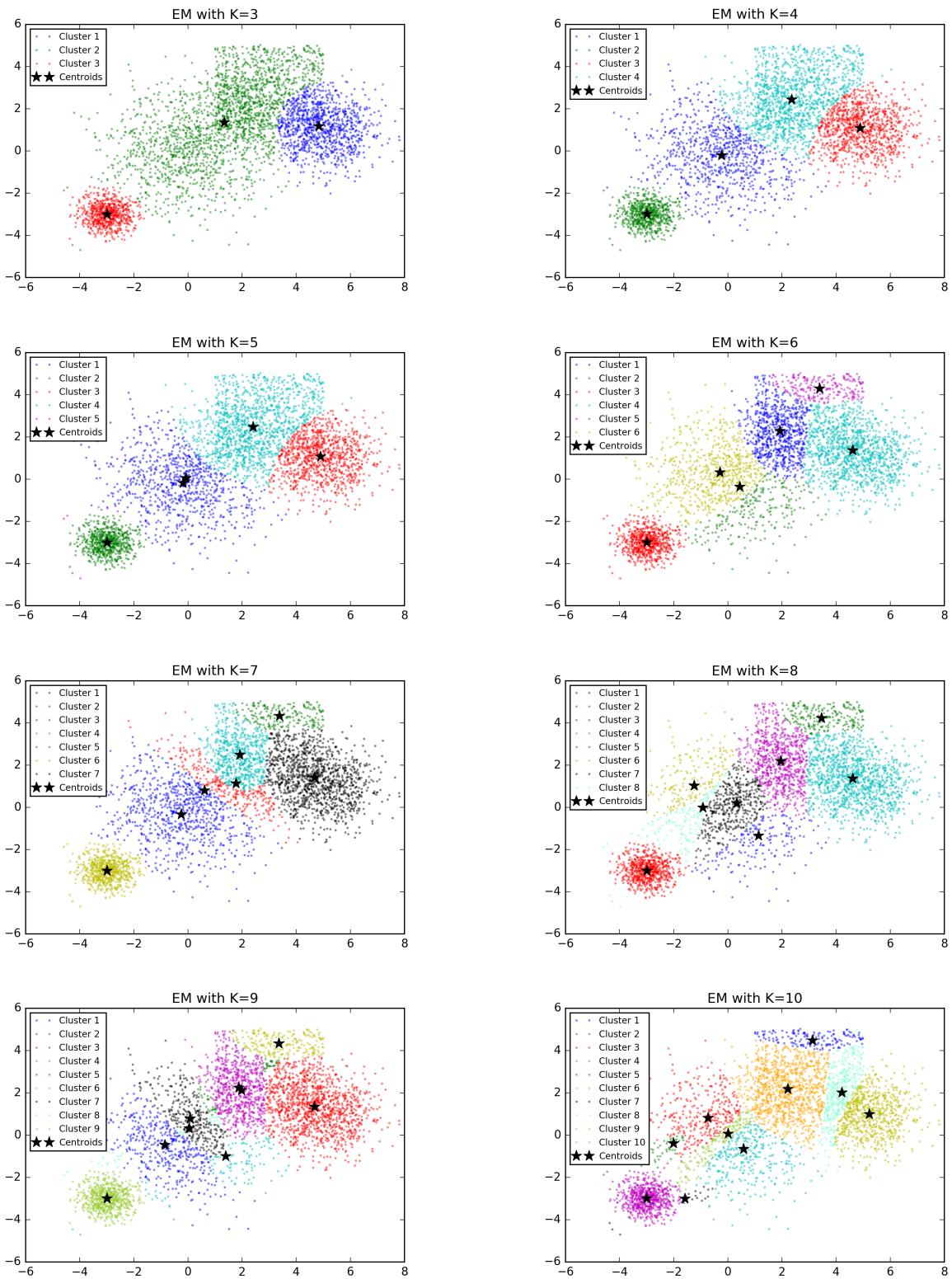


Figure 14: Clustering Result for bigClusteringData.txt with EM Algorithm