

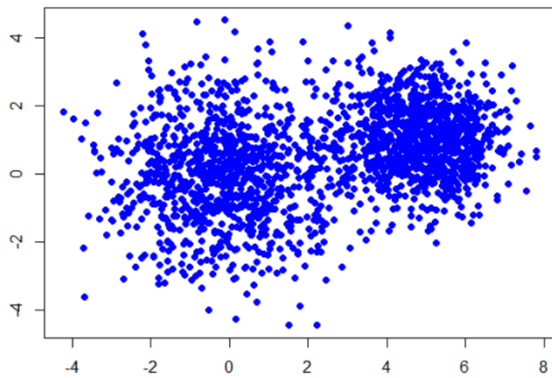
**IE 529**

**Computational Assignment 2**

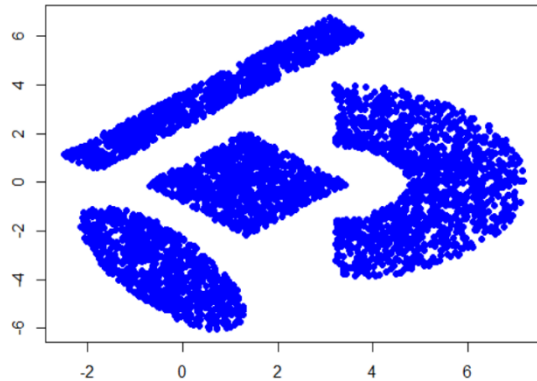
**Shulu Chen(shuluc2)**

**2019/12/18**

1. Include a scatter plot for each set of data, as is, i.e., with no clusters identified yet.



Dataset1: Clustering



Dataset2: ShapedData

2. Submit your pseudocodes for the Lloyds/K-means algorithm, the Greedy K centers and Single Swap algorithms, and the Spectral Clustering algorithm, each on a separate sheet of paper. State clearly whether you are considering normalized or unnormalized Spectral Clustering.

### Lloyd's algorithm

1. Initialization – arbitrarily select  $k$  points as initialize centroids  $m_j$
2. Computer distances (L-2 norm) between each sample  $x_i$  and each centroid  $m_j$  to form a distance matrix  $D$ .
3. Assign each  $x_i$  to its nearest  $m_j$  (partitions data into  $k$  clusters).
4. Update: recompute new “means” for each of the clusters

$$m_j = \frac{1}{|c_j|} \sum_{i=1}^n x_i p_{ij}$$

5. Repeat: until sum of the move-distance between “old centroids” and “new centroids” convergence.

### Greedy k-centers algorithm

1. Initialization: set  $C = \{c_1\}$  for arbitrary  $c_1 \in X$
2. While  $|C| < k$ , do:
  - a) Find  $x_l \in X/C$  that maximize the minimum distance between  $x_l$  and  $C$

$$\max_{x_l} \{ \min_i (\|x_l - c_i\|_2^2) \}$$

- b) Set  $C = C \cup \{x_l\}$ .
  - c) End while.
3. Assign each  $x_i$  to its nearest  $c_i$  (partitions data into k clusters).

### Single-swap heuristic algorithm

1. Initialization: select set of initial medians  $C^0 = \{m_1^0, m_2^0, \dots, m_k^0\}$  and compute the clusters based on  $C^0$ .
2. Compute the cost value:

$$cost^0 = \max_{x_i \in X} (\min_{c_j \in C} \|x_i - c_j\|_2^2)$$

3. Set iteration times  $n$ .
4. While iteration time  $l < n$ , do:
  - a) Arbitrarily select a point  $x_l \in X/C$  and a centroid  $m_j^l \in C$ , swap them and get a temporary centroid set  $C^{l'}$ .
  - b) Compute the new cost value  $cost^{l'}$  based on  $C^{l'}$ .
  - c) IF  $\frac{cost^{l'}}{cost^l} \leq 0.95$ :  
Update the centroid set:  $C^l \leftarrow C^{l'}$
  - d) End while.

### Spectral Clustering algorithm(unnormalized)

1. Compute weight matrix:

$$s_{ij} = \exp \left( -\frac{\|x_i - x_j\|_2^2}{\sigma} \right)$$

$$W = \{s_{ij}\}$$

2. Compute degree matrix:

$$D = \text{diag}\left\{\sum_j s_{ij}\right\}$$

3. Compute Laplacian matrix:

$$L = D - W$$

4. Compute eigenvalue/eigenvector decomposition for L:

$$L = U^{-1}\Lambda U$$

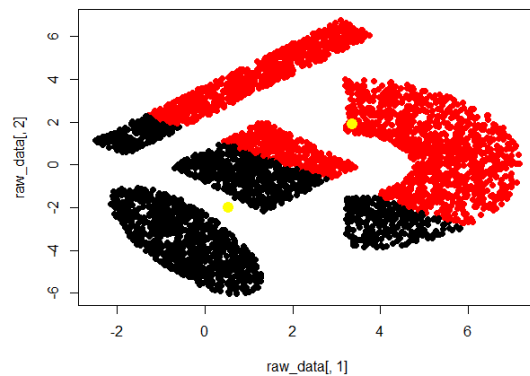
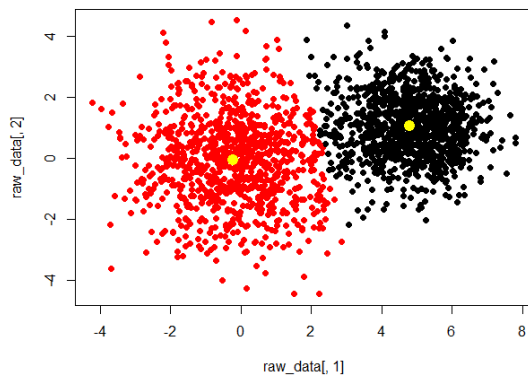
5. Choose least k eigenvalues and take corresponding eigenvectors of  $U$  as  $u_k$ .
6. Denote rows of  $u_k$  by  $y_i$
7. Cluster “points”  $\{y_1, \dots, y_n\}$  using k-means.
8. Define output index sets  $J_l = \{j | y_j \in c_l\}$

- For your K-means algorithm by itself, and for your Spectral Clustering algorithm, present the output(clustering) results for BOTH sets of test data given to you. (See below)
- Evaluate the results for a range of K values (number of clusters), e.g., from  $K = 2$ , to  $K = 11$  (but maybe not all of these; try at least 3 different values, then try to pick 2-3 more that will help you find the best cluster outcome). When implementing the K-means algorithm (by itself and from your Spectral Clustering algorithm), you should try multiple initializations (say 5 to 10) and then select the 'best' result, for each K value. State how many different initializations you tried and why.

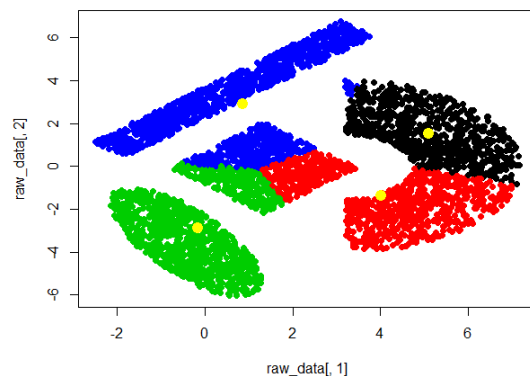
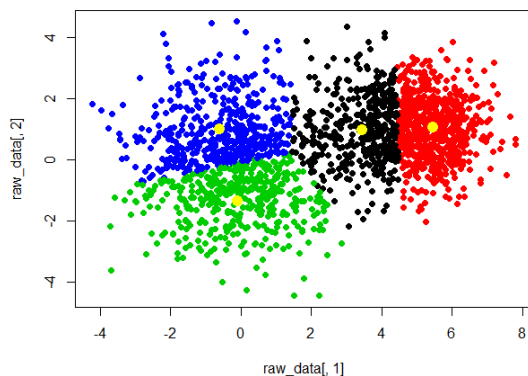
K-means on dataset1

k-means on dataset2

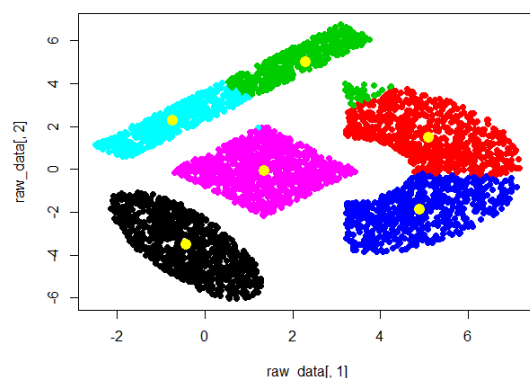
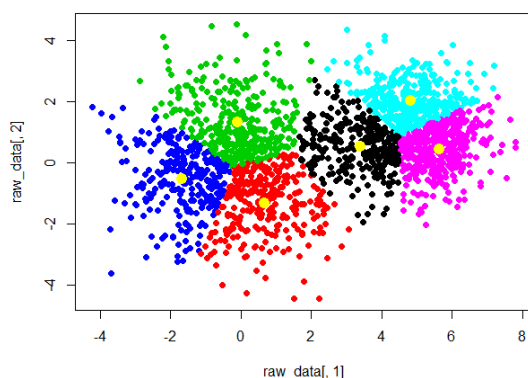
K=2



K=4

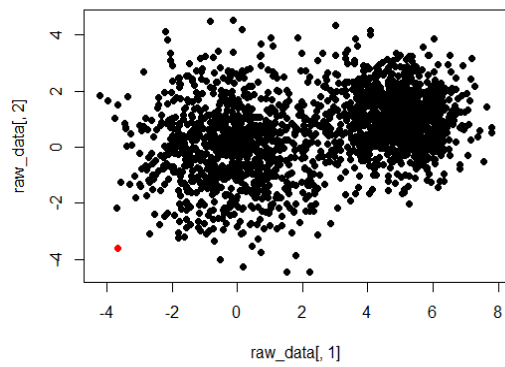


K=6

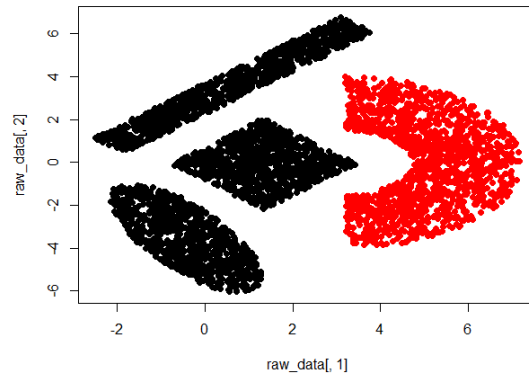


Spectral Clustering on dataset1

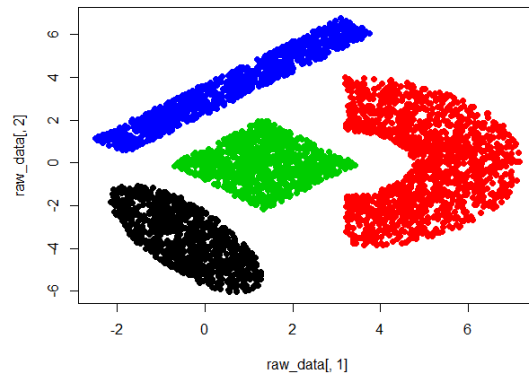
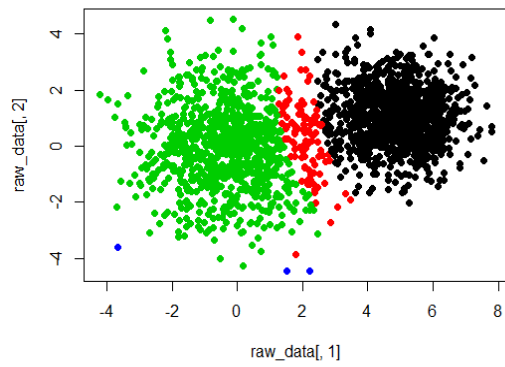
K=2



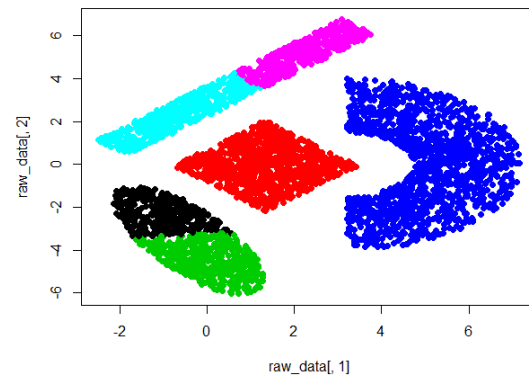
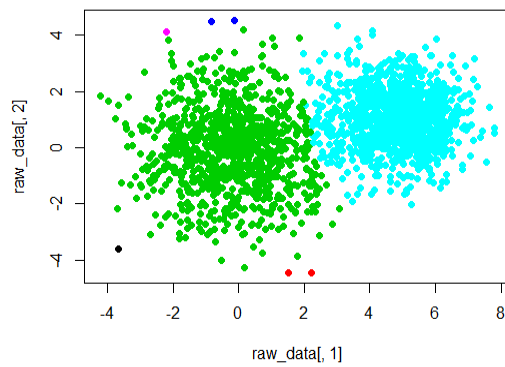
Spectral Clustering on dataset2



K=4

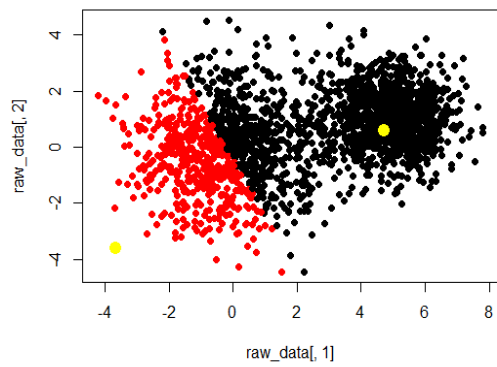


K=6

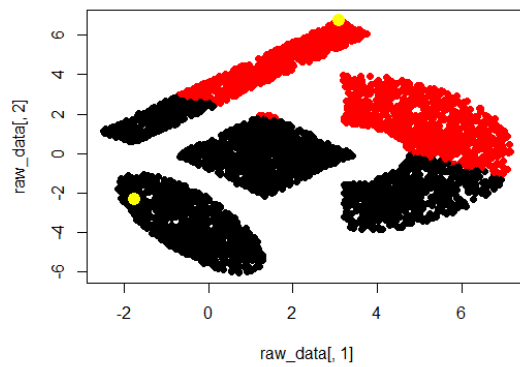


Greedy K-Centers on dataset1

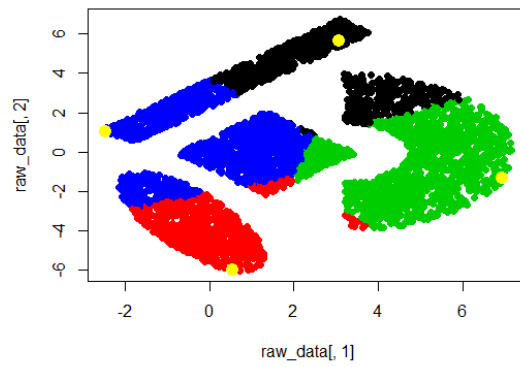
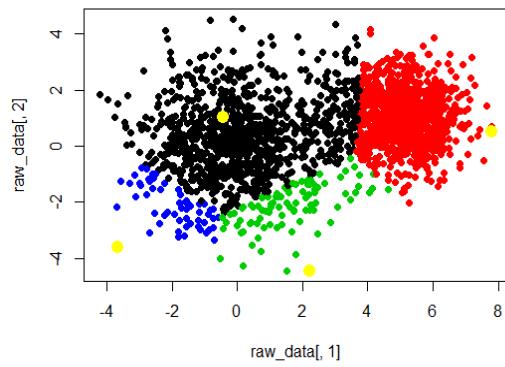
K=2



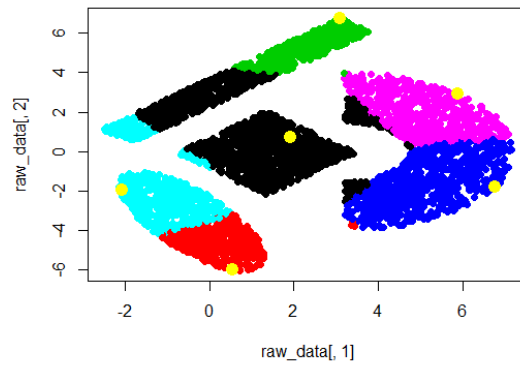
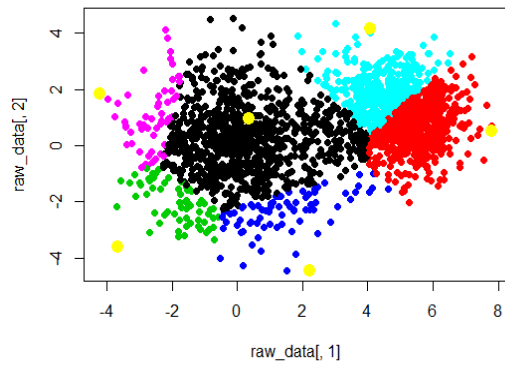
Greedy K-Centers on dataset2



K=4



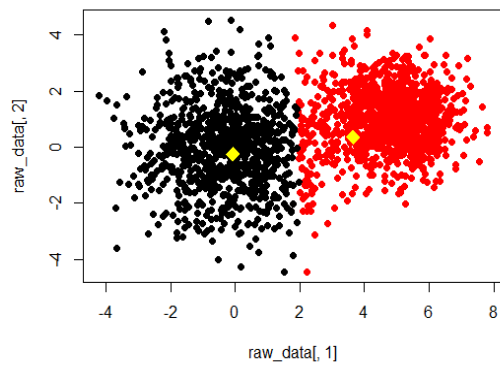
K=6



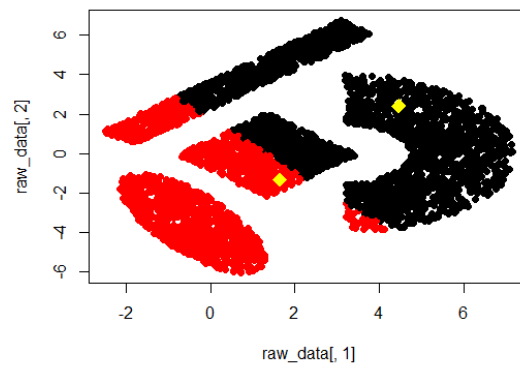


Single-Swap Heuristic on dataset1

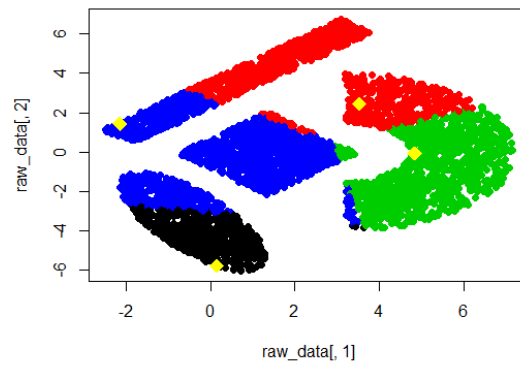
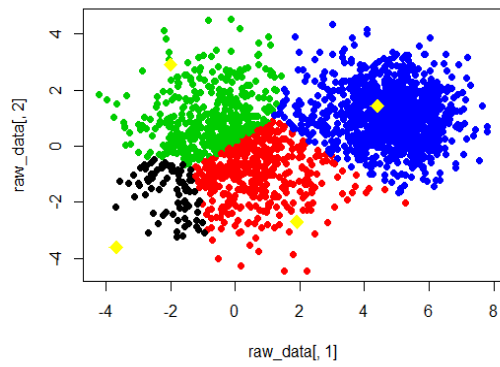
K=2



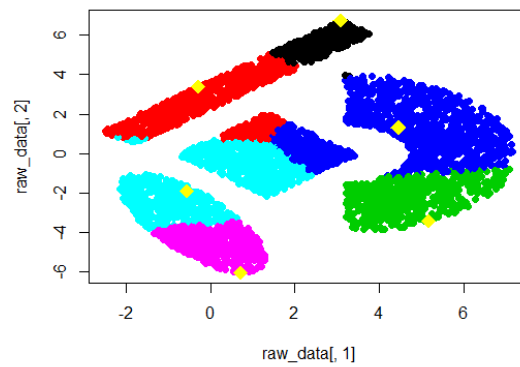
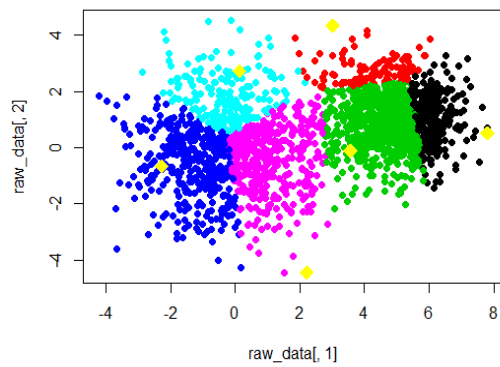
Single-Swap Heuristic on dataset2



K=4



K=6



5. **For the K-means algorithm, describe how you tested for convergence and why you choose this as a test.**

For k-means algorithm, I computer the distance between the “old centroids” and “new centroids” as the iteration gains. When the iteration gains is less than 0.0001, which shows that the centroids do not change greatly anymore, and then we can regard it as convergence.

6. **Discuss how many “natural” clusters you think each data set has, and why.**

Dataset1: Two clusters

Dataset2: Four clusters

Visually, sample data in dataset1 mixed together, and the left side and the right side seems cluster at each center, so two clusters make sense. For dataset2, it is clearly to see that there are four shaped sample data, so four clusters would work better.

7. **Provide a short analysis of the computational effort you found was required for each of the algorithms. State what measure you used to evaluate computational effort (e.g., running time or flops, or ??), and how you observed this. Briefly discuss how these computational efforts compare to your expectations.**

For Lloyd's algorithm, time complexity is  $O(nkd)$  at each iteration.

For Spectral Clustering, time complexity  $O(n^3)$  because it needs to computer the similarity matrix and eigenvalue decomposition. When it comes to a large dataset, it would be time-consuming.

8. **Write a brief summary paragraph of your findings, that is, briefly summarize your KEY findings.**

For Lloyd's algorithm, I find that the initialize points is not so important compared with Greedy K-center algorithm. The result is similar with the package “kmeans” in R language. So, we can make a conclusion that the Lloyd’s algorithm is robust, and the accuracy is good.

For Greedy K-Centers algorithm, the result would be greatly influenced by the initialized points, the results changes greatly at each iteration. Also, I do not get a good result(compared with Lloyd’s algorithm).

The Singular Swap algorithm is a good way to update the centroids after Greedy K-Centers. However, when the clusters k is small, e.g. k=2, the update policy may not work well.

The Spectral Clustering works really well when the dataset is divided clearly with different shapes, but the result is not well when the dataset is traditional. So, Spectral Clustering may be used in some special area.

## Appendix: Personal codes

### Lloyd's algorithm

```
1  ###Lloyd's algorithm###
2  setwd("C:/Users/Lenovo/Desktop")
3  raw_data<-read.csv("clustering.csv",header = F)
4  k=4
5  n=nrow(raw_data)
6  plot(x=raw_data[,1],y=raw_data[,2],pch=16,col="blue")
7
8  ##cluster algorithm
9  clusters<-function(ct,data){
10     distance<-matrix(data = NA,nrow=n,ncol=k)
11     cluster<-c()
12     for (j in 1:k){
13         distance[,j]<-(data[,1]-ct[j,1])^2+(data[,2]-ct[j,2])^2
14     }
15     for (i in 1:n){
16         cluster[i]<-which.min(distance[i,])
17     }
18     cluster
19 }
20
21 count_cent<-function(cluster_data){
22     count_table<-table(cluster_data[,3])
23     x<-c()
24     y<-c()
25     for (i in 1:k){
26         x[i]<-(1/count_table[i])*sum(data[data[,3]==i,1])
27         y[i]<-(1/count_table[i])*sum(data[data[,3]==i,2])
28     }
29     cent_point<-cbind(x,y)
30     cent_point
31 }
32 #initialization:select k points randomly
33 initial_point<-as.matrix(raw_data[sample(row.names(raw_data), k),])
34 points(initial_point,col="red",pch=16)
35 #points(initial_point,col="red",pch=16)
36 cluster_int<-clusters(initial_point,raw_data)
37 data<-cbind(raw_data,cluster_int)
38 cent_points<-count_cent(data)
39 iteration_distance<-(cent_points-initial_point)
40 iteartion_rate<-sum(iteration_distance[,1]^2+iteration_distance[,2]^2)
41 record_iteation<-c()
42 while (iteartion_rate>=0.0001) {
43     cluster<-clusters(cent_points,raw_data)
44     data<-cbind(raw_data,cluster)
45     cent_points_new<-count_cent(data)
46     iteration_distance<-(cent_points_new-cent_points)
47     iteartion_rate<-sum(iteration_distance[,1]^2+iteration_distance[,2]^2)
48     print(iteartion_rate)
49     cent_points<-cent_points_new
50     append(record_iteation, iteartion_rate)
51 }
52
53 plot(x=raw_data[,1],y=raw_data[,2],col=data[,3],pch=16)
54 points(cent_points,col="blue",pch=18,cex=1.5)
55
56 #km<-kmeans(raw_data,k)
57 #plot(x=raw_data[,1],y=raw_data[,2],col=km$cluster,pch=16)
58 #sum(data[,3]-km$cluster)
```

## Greedy k-centers algorithm

```
1  ###GreedyKCenters###
2  setwd("C:/Users/Lenovo/Desktop")
3  raw_data<-read.csv("clustering.csv",header = F)
4  k=6
5  n=nrow(raw_data)
6  plot(x=raw_data[,1],y=raw_data[,2],col="blue",pch=16)
7  #initialization:select a point randomly
8  initial_point<-as.matrix(raw_data[sample(row.names(raw_data), 1),])
9  distance_matrix<-matrix(data=NA,nrow=n,ncol=k)
10 point_matrix<-matrix(data=NA,nrow=k,ncol=2)
11 point_matrix[1,]<-initial_point
12 points(x=point_matrix[1,1],y=point_matrix[1,2],col="red",pch=16,cex=2)
13 i=1
14 distance_matrix[i,]<-(raw_data[,1]-point_matrix[1,1])^2+(raw_data[,2]-point_matrix[1,2])^2
15 y<-which.max(distance_matrix[,1])
16 for (i in 2:k) {
17   point_new<-as.matrix(raw_data[y,])
18   point_matrix[i,]<-point_new
19   distance_matrix[i,]<-(raw_data[,1]-point_matrix[i,1])^2+(raw_data[,2]-point_matrix[i,2])^2
20   min_dis=c()
21   for (j in 1:n){
22     min_dis[j]<-min(distance_matrix[j,1:i])
23   }
24   y<-which.max(min_dis)
25   points(x=point_matrix[i,1],y=point_matrix[i,2],col="red",pch=16,cex=2)
26 }
27
28 clusters<-function(ct,data){
29   distance<-matrix(data = NA,nrow=n,ncol=k)
30   cluster<-c()
31   for (j in 1:k){
32     distance[,j]<-(data[,1]-ct[j,1])^2+(data[,2]-ct[j,2])^2
33   }
34   for (i in 1:n){
35     cluster[i]<-which.min(distance[i,])
36   }
37   cluster
38 }
39
40 cluster<-clusters(point_matrix,raw_data)
41 data<-cbind(raw_data,cluster)
42 plot(x=raw_data[,1],y=raw_data[,2],col=data[,3],pch=16)
```

## Single-swap heuristic algorithm

```
1  ###single-swap heuristic###
2  setwd("C:/Users/Lenovo/Desktop")
3  raw_data<-read.csv("clustering.csv",header = F)
4  k=6
5  n=nrow(raw_data)
6  #plot(x=raw_data[,1],y=raw_data[,2],col="blue",pch=16)
7  #initialization:select a point randomly
8
9  greedy_centers<-function(data){
10   initial_point<-as.matrix(data[sample(row.names(data), 1),])
11   y1<-which(data==initial_point)
12   distance_matrix<-matrix(data=NA,nrow=n,ncol=k)
13   point_matrix<-matrix(data=NA,nrow=k,ncol=3)
14   point_matrix[1,1:2]<-initial_point
15   point_matrix[1,3]<-y1
16   i=1
17   distance_matrix[i,<-](data[,1]-point_matrix[1,1])^2+(data[,2]-point_matrix[1,2])^2
18   y<-which.max(distance_matrix[,1])
19
20   for (i in 2:k) {
21     point_new<-as.matrix(data[y,])
22     point_matrix[i,1:2]<-point_new
23     distance_matrix[i,<-](data[,1]-point_matrix[i,1])^2+(data[,2]-point_matrix[i,2])^2
24     min_dis=c()
25     for (j in 1:n){
26       min_dis[j]<-min(distance_matrix[j,1:i])
27     }
28     y<-which.max(min_dis)
29     #points(x=point_matrix[i,1],y=point_matrix[i,2],col="red",pch=16,cex=2)
30     point_matrix[i,3]<-y
31   }
32   point_matrix
33 }
34
35 count_cluster<-function(ct,data){
36   distance<-matrix(data = NA,nrow=n,ncol=k)
37   min_distance<-c()
38   cluster<-c()
39   for (j in 1:k){
40     distance[,j]<-c((data[,1]-ct[j,1])^2+(data[,2]-ct[j,2])^2)
41   }
42   for (i in 1:n){
43     cluster[i]<-which.min(distance[i,])
44     min_distance[i]<-min(distance[i,])
45   }
46   new_distance<-cbind(cluster,min_distance)
47   new_distance
48 }
49
50 point_matrix<-as.data.frame(greedy_centers(raw_data))
51 colnames(point_matrix)<-c("x","y","index")
52 plot(x=raw_data[,1],y=raw_data[,2],col="blue",pch=16)
53 points(point_matrix,col="red",pch=16)
54 cluster_data<-count_cluster(point_matrix[,1:2],raw_data)
55 data<-cbind(raw_data,cluster)
56 data_index<-as.data.frame(row.names(data))
57 data<-cbind(data,data_index)
58 colnames(data)<-c("x","y","cluster","min_distance","index")
59 data_swap<-data[-c(point_matrix[,3]),]
60 cost<-max(data_swap[,4])
61 ##choose and swap point
62 for (i in 1:1000){
63   m_new<-as.data.frame(data_swap[sample(row.names(data_swap), 1),c(1,2,5)])
64   c_old<-as.data.frame(point_matrix[sample(row.names(point_matrix), 1),])
65   point_matrix_new<-rbind(point_matrix[-which(point_matrix[,3]==c_old[1,3]),],m_new)
66   cluster_data_new<-count_cluster(point_matrix_new[,1:2],raw_data)
67   data_new<-cbind(raw_data,cluster_data_new,data_index)
68   data_swap_new<-data[-as.numeric(c(point_matrix_new[,3])),]
69   cost_new<-max(data_new[,4])
70   if (cost_new/cost<=0.95){
71     point_matrix<-point_matrix_new
72     cost<-cost_new
73     data_swap<-data_swap_new
74     print(cost)
75   }
76 }
77 points(point_matrix,col="green",pch=16)
78 plot(x=raw_data[,1],y=raw_data[,2],col=cluster_data_new[,1],pch=16)
```

## Spectral Clustering algorithm

```
1  ###Spectral Clustering algorithm,###
2
3  setwd("C:/Users/Lenovo/Desktop")
4  raw_data<-read.csv("ShapedData.csv",header = F)
5  k=4
6  #raw_data<-raw_data[1:100,]
7  n=nrow(raw_data)
8  sigma<-0.9
9  plot(raw_data,col="blue",pch=16)
10 weight_matrix<-matrix(data=NA,nrow=n,ncol=n)
11 #compute weight matrix#
12 for (i in 1:n){
13   weight_matrix[,i]<-(raw_data[,1]-raw_data[i,1])^2+(raw_data[,2]-raw_data[i,2])^2
14 }
15 weight_matrix_new<-exp(-weight_matrix/sigma)
16 sum_w<-c()
17 for (i in 1:n){
18   sum_w[i]<-sum(weight_matrix_new[i,])
19 }
20 D<-diag(sum_w)
21 I<-diag(1,nrow=n)
22 L<-D-weight_matrix_new
23 #L<-I-D^(-.5)*weight_matrix_new*D^(-.5)
24 L_e<-eigen(L)
25
26 uk<-L_e$vector[, (n-k+1):n]
27
28 #km<-kmeans(uk,k)
29 #data<-cbind(raw_data,km$cluster)
30 #plot(x=data[,1],y=data[,2],col=data[,3],pch=16,colnames("x"),rownames("y"))
31
32 clusters<-function(ct,data){
33   d<-ncol(ct)
34   distance<-matrix(data = 0,nrow=n,ncol=k)
35   cluster<-c()
36   for (j in 1:k){
37     for (i in 1:d){
38       distance[,j]<-(data[,i]-ct[j,i])^2+distance[,j]
39     }
40   }
41   for (i in 1:n){
42     cluster[i]<-which.min(distance[i,])
43   }
44   cluster
45 }
46
47 count_cent<-function(cluster_data){
48   d<-ncol(cluster_data)
49   print(d)
50   count_table<-table(cluster_data[,d])
51   cent_point<-matrix(data=0,ncol=(d-1),nrow=k)
52   for (i in 1:k){
53     for (j in 1:(d-1)){
54       cent_point[i,j]<-(1/count_table[i])*sum(cluster_data[cluster_data[,d]==i,j])
55     }
56   }
57   cent_point
58 }
59
60 #initialization:select k points randomly
61 uk<-as.data.frame(uk)
62 initial_point<-as.matrix(uk[sample(row.names(uk), k),])
63 #points(initial_point,col="red",pch=16)
64 cluster_int<-clusters(initial_point,uk)
65 data<-cbind(uk,cluster_int)
66 cent_points<-count_cent(data)
67 iteration_distance<-(cent_points-initial_point)
68 iteartion_rate<-0
69 for (i in 1:k){
70   iteartion_rate<-iteartion_rate+sum(iteration_distance[,i]^2)
71 }
72
```

```

60 #initialization:select k points randomly
61 uk<-as.data.frame(uk)
62 initial_point<-as.matrix(uk[sample(row.names(uk), k),])
63 #points(initial_point,col="red",pch=16)
64 cluster_int<-clusters(initial_point,uk)
65 data<-cbind(uk,cluster_int)
66 cent_points<-count_cent(data)
67 iteration_distance<-(cent_points-initial_point)
68 iteartion_rate<-0
69 for (i in 1:k){
70   iteartion_rate<-iteartion_rate+sum(iteration_distance[,i]^2)
71 }
72
73 record_iteation<-c()
74 while (iteartion_rate>=0.0001) {
75   cluster<-clusters(cent_points,uk)
76   data<-cbind(uk,cluster)
77   cent_points_new<-count_cent(data)
78   iteration_distance<-(cent_points_new-cent_points)
79   iteartion_rate<-0
80   for (i in 1:k){
81     iteartion_rate<-iteartion_rate+sum(iteration_distance[,i]^2)
82   }
83   print(iteartion_rate)
84   cent_points<-cent_points_new
85   append(record_iteation, iteartion_rate)
86 }
87
88 plot(x=raw_data[,1],y=raw_data[,2],col=data[,k+1],pch=16)
89 #points(cent_points,col="blue",pch=18,cex=1.5)
90

```