

Math 5365

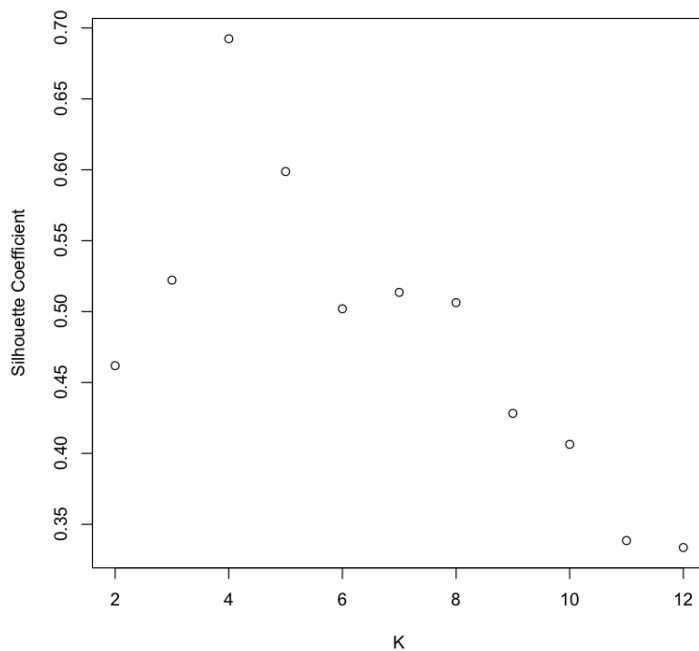
Data Mining 1

Homework 18

Mary Barker

1. Consider the data from problem 1 on Homework 17

- (a) Find the number of clusters that maximizes the silhouette coefficient, and plot the silhouette coefficient, and plot the silhouette coefficient vs K.



The optimal number of clusters for the silhouette coefficient is 4.

- (b) What is the maximum possible value of the silhouette coefficient?

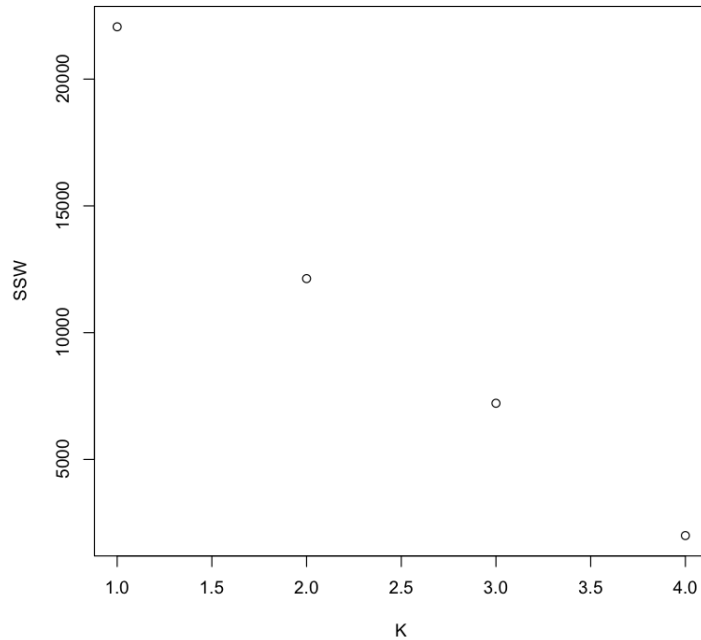
```
$max
```

```
[1] 0.6923732
```

```
$where
```

```
[1] 4
```

- (c) Plot SSW vs K. Does the optimal value of K suggested by this plot agree with the one based on the silhouette coefficient?



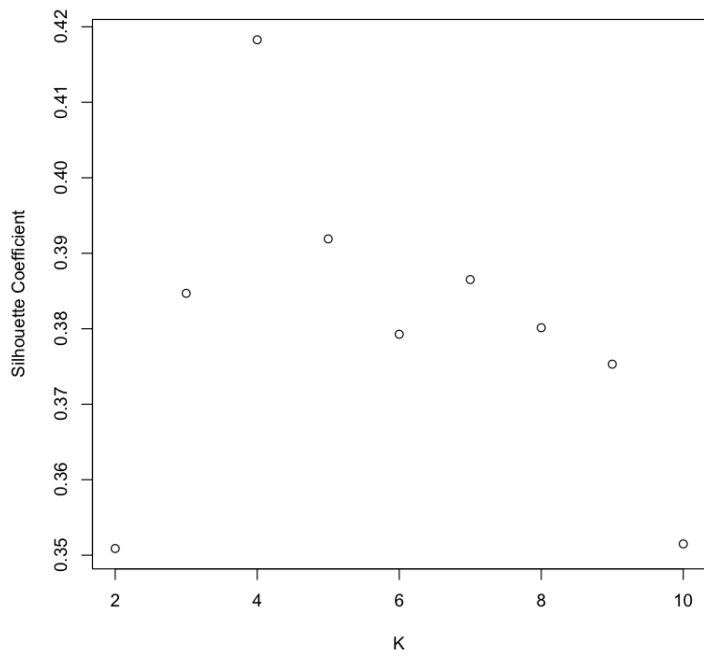
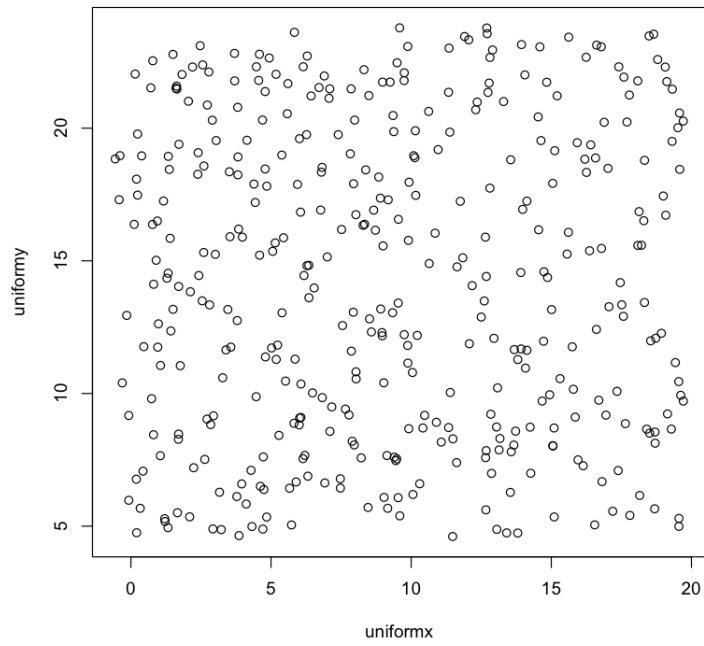
- (d) Optional: Is the silhouette coefficient for this clustering statistically significant? (It may be a good idea to let R run while you're out of the office to do this problem.) Running the same case with uniform data gave a silhouette coefficient maximized at  $k = 4$  also.

```
$max
```

```
[1] 0.4182791
```

```
$where
```

```
[1] 4
```



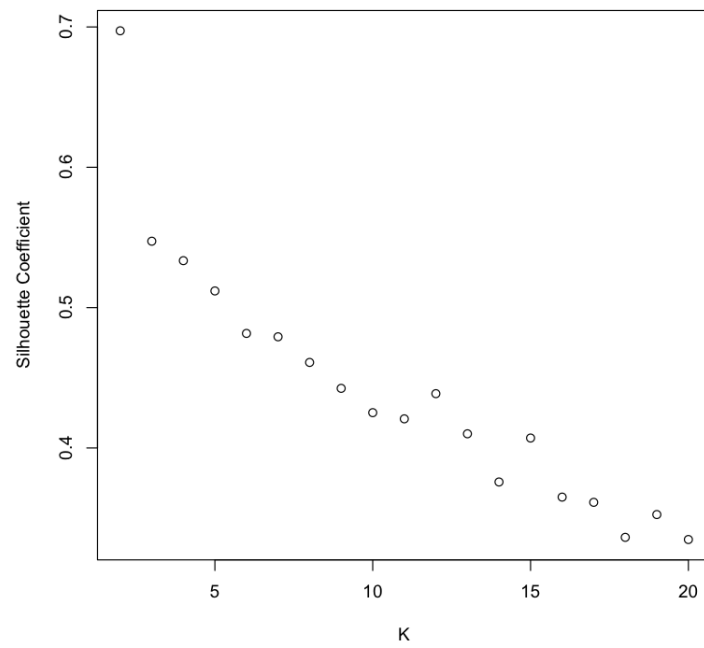
2. Repeat problem 1 for the wdbc data set. Also, find the weighted entropy and purity for the optimal clustering.

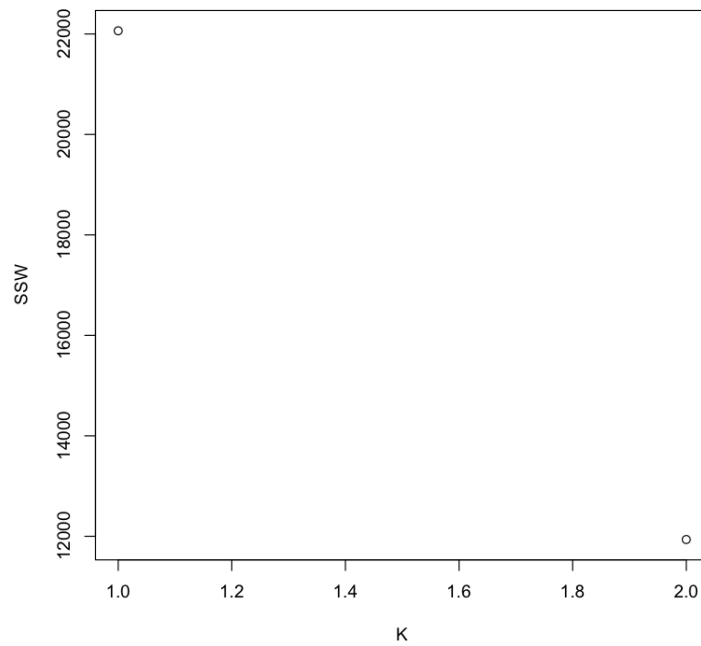
\$max

[1] 0.6972646

\$where

[1] 2





entropy = 0.5503462

purity = 0.8541301

```

1 #Data Mining hw 18
2 library(stats)
3 library(cluster)
4 library(fields)
5
6 wdbc <- read.csv('~/.Dropbox/Tarleton/data_mining/dfiles/wdbc.data',
7                 header=F,sep=',')
8 wdbc <- wdbc[,-1]
9 source('~/.Dropbox/Tarleton/data_mining/generic_functions/dataset_ops.R')
10 source('~/.Dropbox/Tarleton/data_mining/generic_functions/measures.R')
11
12 x <- c(rnorm(100, 5, 1.5), rnorm(100, 15, 1.5),
13        rnorm(100, 5, 1.5), rnorm(100, 15, 1.5))
14 y <- c(rnorm(100, 10, 1.5), rnorm(100, 10, 1.5),
15        rnorm(100, 20, 1.5), rnorm(100, 20, 1.5))
16 plot(x, y)
17 points <- data.frame(x = x, y = y)

```

```

18
19
20 kmeans_reps <- function(data, centers, reps){
21   w_ss = Inf
22   for(i in 1:reps){
23     k_cluster <- kmeans(x = data, centers = centers)
24     if((k_cluster$tot.withinss) < w_ss){
25       ssw = k_cluster$tot.withinss
26       my_k_cluster <- k_cluster
27     }
28   }
29   return(my_k_cluster)
30 }
31
32 min_rep <- function(K, eps){
33   ceiling(log(eps) / log(1 - factorial(K)/K^K))
34 }
35
36 # a. Find the number of clusters that maximizes the silhouette coefficient,
37 # and plot the silhouette coefficient, and plot the silhouette
38 # coefficient vs K.
39
40 mysil <- function(x, dmat){
41   return(mean(silhouette(x = x, dmat = dmat)[,3]))
42 }
43
44 find_sil <- function(data, kmax, niter, eps){
45   dmat <- rdist(data)
46   sil_v <- 1:kmax
47   for(K in 2:kmax){
48     iter <- min(niter, min_rep(K, eps))
49     kmeans_tmp <- kmeans_reps(data, K, iter)
50     sil_v[K] <- mysil(kmeans_tmp$cluster, dmat)
51   }
52   sil_v <- sil_v[2:kmax]
53   plot(2:kmax, sil_v, xlab='K', ylab='Silhouette Coefficient')
54   return(list(max = max(sil_v), where = which.max(sil_v) + 1))
55 }
56 max_k <- find_sil(points, 12, 1000, 0.01)
57
58
59 # b. What is the maximum possible value of the silhouette coefficient?
60 max_k$max
61

```

```

62 # c. Plot SSW vs K. Does the optimal value of K suggested by this plot agree
63 #     with the one based on the silhouette coefficient?
64
65 plot_ssw <- function(data, kmax, niter, eps){
66     ssw_v <- 1:kmax
67
68     for(K in 1:kmax){
69         iter <- min(niter, min_rep(K, eps))
70         kmeans_tmp <- kmeans_reps(data, K, iter)
71         ssw_v[K] <- kmeans_tmp$tot.withinss
72     }
73     plot(1:kmax, ssw_v, xlab='K', ylab='SSW')
74 }
75 plot_ssw(points, max_k$where, 1000, 0.01)
76
77 # d. Optional: Is the silhouette coefficient for this clustering statistically
78 #     significant? (It may be a good idea to let R run while you're out of the
79 #     office to do this problem.)
80
81 rmat <- apply(points, 2, range)
82 uniformx <- runif(nrow(points), rmat[1,1], rmat[2,1])
83 uniformy <- runif(nrow(points), rmat[1,2], rmat[2,2])
84 upoints <- data.frame(uniformx, uniformy)
85
86 find_sil(upoints, 10, 1000, 0.01)
87
88 # 2. Repeat problem 1 for the wdbc data set. Also, find the weighted entropy
89 #     and purity for the optimal clustering.
90 # a.
91 max_k <- find_sil(wdbc, 20, 1000, 0.01)
92
93 # b.
94 max_k$max
95
96 # c.
97 plot_ssw(points, max_k$where, 1000, 0.01)
98
99 table_ent <- function(table){
100     col_sums <- apply(table, 2, sum)
101     col_props <- col_sums / sum(col_sums)
102     for(j in 1:ncol(table)){
103         if(sum(table[,j] != 0)){
104             table[,j] <- table[,j] / sum(table[,j])
105         }

```

```

106     }
107     table_entropies <- apply(table, 2, entropy_eval)
108     return(col_props %*% table_entropies)
109 }
110
111 best_wdbc <- kmeans_reps(wdbc[,2:ncol(wdbc)], 2, 1000)
112 predicted <- rep('',nrow(wdbc))
113 predicted[1 * (best_wdbc$cluster == 2) == 1] <- 'B'
114 predicted[1 * (best_wdbc$cluster == 1) == 1] <- 'M'
115 wdbc_tab <- table(wdbc$V2, predicted)
116 entropy <- table_ent(wdbc_tab)
117
118 purity <- sum(apply(wdbc_tab, 2, max)) / sum(wdbc_tab)

```