

Exercises 2-4

Giovanni Zurlo

12/10/2021

Exercise 1

1. Use the gap statistic method to estimate the number of clusters for the olive oil data, and for Artificial Dataset2. Comment on the solutions.
2. Generate a dataset from a two-dimensional uniform distribution on the rectangle $[min\ x1; max\ x1] \times [min\ x2; max\ x2]$, where $x1, x2$ are the values of the first and second variable of Artificial Dataset2. Compute K-means clusterings for K from 1 to 10 and show at least three scatterplots of the dataset with clusterings with different K . Compare the values of $\log Sk$ from clustering Artificial Dataset2 with those from clustering the uniformly distributed dataset

Olive Oil Data

```
# Loading and standardizing oliveoil data
library(cluster)
data("oliveoil")
diag(var(oliveoil))
```

```
## macro.area      region    palmitic palmitoleic    stearic     oleic
##      NA          NA       28423.3515   2755.6583 1350.1903 164681.9365
##      linoleic    linolenic arachidic eicosenoic
##      58951.4616 168.1871   485.3319   198.3392
```

```
solve=scale(oliveoil[,3:10])
library(cluster)
set.seed(1234)
cg.solve=clusGap(solve,kmeans, K.max=10, B=150, d.power=2,
                 spaceH0 = "scaledPCA", nstart=100, iter.max=100)
plot(cg.solve, main="Gap Statistic Plot for std Oliveoil Data")
```

Gap Statistic Plot for std Oliveoil Data



```
# Extracting the best K according to different criteria
maxSE(cg.solve$Tab[,3],cg.solve$Tab[,4],"globalSEmax",SE.factor=2)
```

```
## [1] 9
```

```
maxSE(cg.solve$Tab[,3],cg.solve$Tab[,4],"firstSEmax",SE.factor=2)
```

```
## [1] 9
```

```
maxSE(cg.solve$Tab[,3],cg.solve$Tab[,4],"Tibs2001SEmax",SE.factor=3)
```

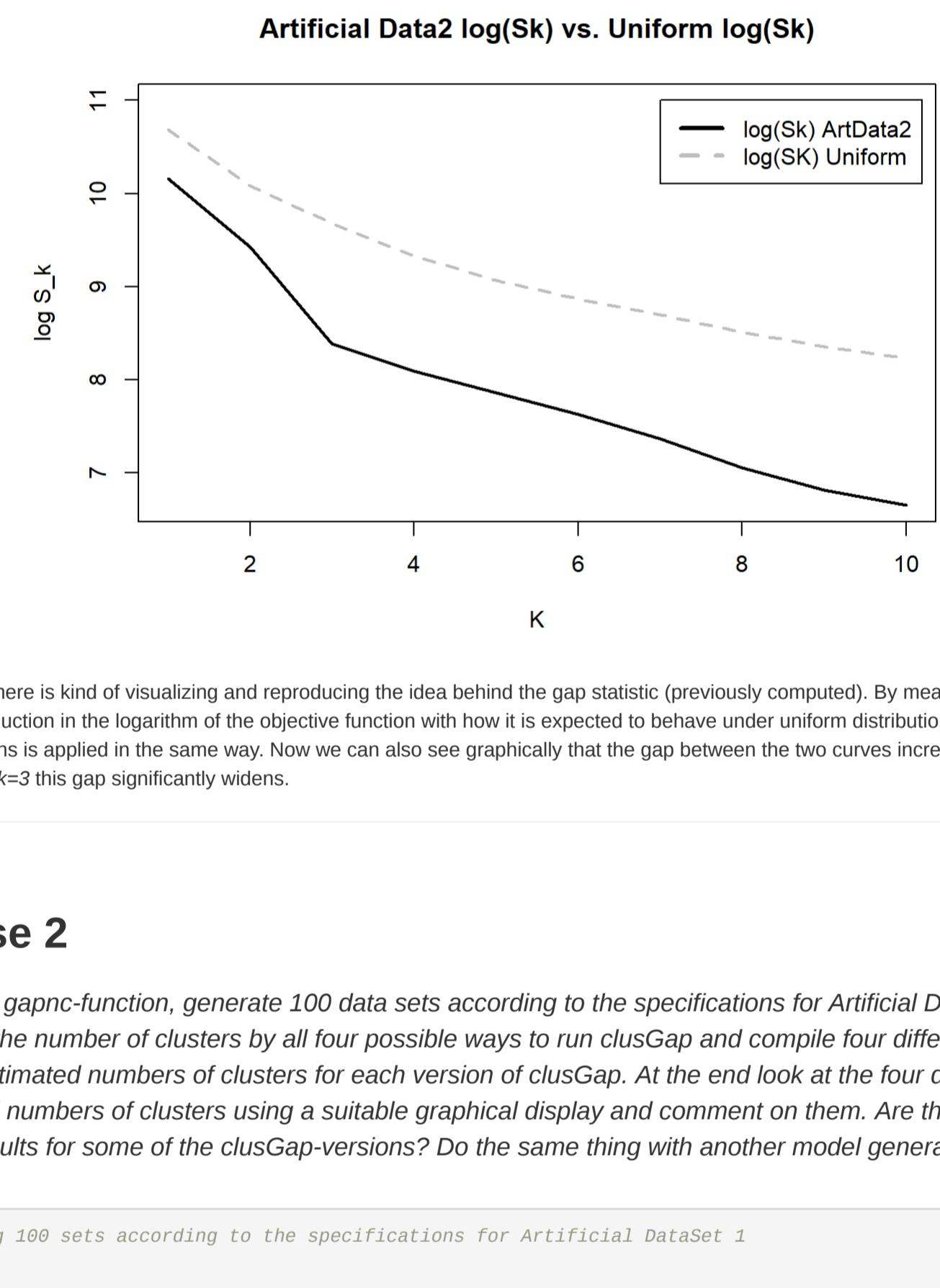
```
## [1] 5
```

The gap statistic plot shows the elbow at $k=5$, however the index grows monotonically up to $k=10$. Both the "globalSEmax" and "firstSEmax" criterions suggest choosing $k=9$ clusters, for which solution we saw a clear correspondence with the "region" covariate. The only criterion suggesting the elbow value $k=5$ is "Tibs2001SEmax" with 3 as standard error factor. I would choose between $k=5$ or 9 depending on the different usage of these results, maybe 5 leads to insights on the data that are more catching and concise.

Artificial Dataset 2

```
# Loading and standardizing the Artificial Dataset 2
artdata=read.table(file.choose())
sartdata=scale(artdata)
library(cluster)
set.seed(1234)
cg.sart=clusGap(sartdata,kmeans, K.max=10, B=130, d.power=2,
                 spaceH0 = "scaledPCA", nstart=100, iter.max=100)
plot(cg.sart, main="Gap Statistic Plot for std Artificial Data 2")
```

Gap Statistic Plot for std Artificial Data 2



```
# Extracting the best K according to different criteria
maxSE(cg.sart$Tab[,3],cg.sart$Tab[,4],"globalSEmax",SE.factor=2)
```

```
## [1] 9
```

```
maxSE(cg.sart$Tab[,3],cg.sart$Tab[,4],"firstSEmax",SE.factor=2)
```

```
## [1] 3
```

```
maxSE(cg.sart$Tab[,3],cg.sart$Tab[,4],"Tibs2001SEmax",SE.factor=1)
```

```
## [1] 3
```

The gap statistic plot shows a local maximum at $k=3$, which we know to be the most natural choice (based on the data generating process). "GlobalSEmax" suggests 9 as the optimal number of groups, but both "Tibs2001SEmax" and "firstSEmax" point towards the local maximum of the graph $k=3$, which seems the most reasonable value.

Comparison with a uniform distribution

```
set.seed(1234)
# Generating a dataset from a two-dimensional uniform distribution
U1=runif(140, max = max(artdata$V1), min = min(artdata$V1))
U2=runif(140, max = max(artdata$V2), min = min(artdata$V2))
U=data.frame(U1,U2)

# Collecting the outputs of the k-means function in a list, from 1 to K.max
km.out=list()
set.seed(1234)
for (i in 1:10) {km.out[[i]]=kmeans(U,iter.max = 100,nstart = 100, centers = i)}
```

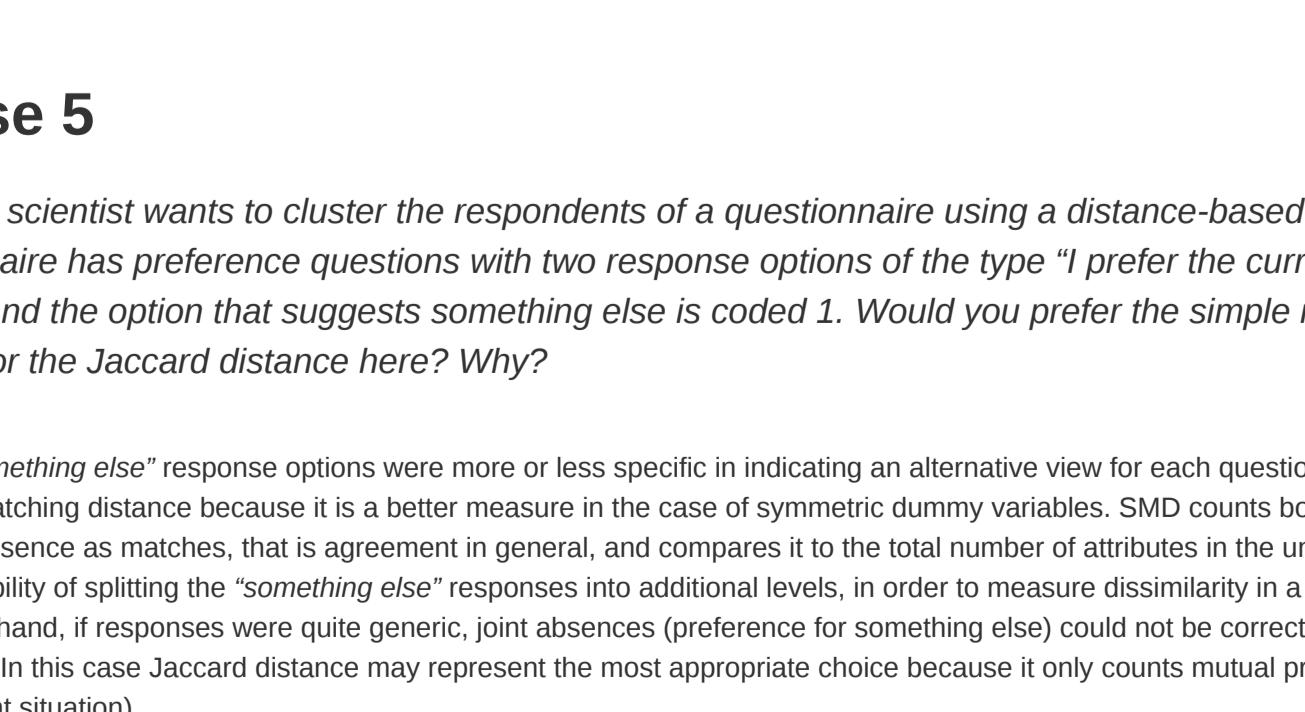
```
# Three scatterplots of the uniform data with different K
plot(U,col=(km.out[[2]]$cluster),cex=1.2,pch=0,lwd=2,
     main='Uniform Data Clustering, K=2')
```

Uniform Data Clustering, K=2



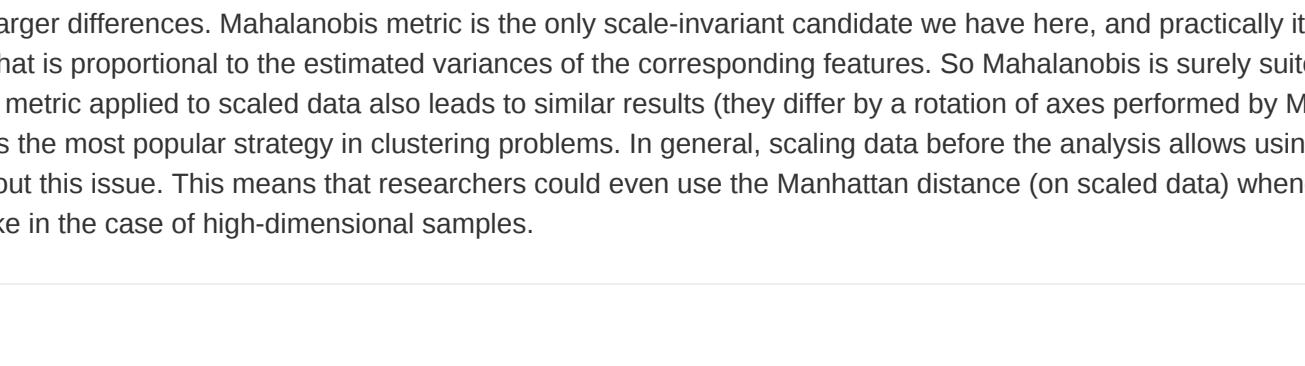
```
plot(U,col=(km.out[[4]]$cluster),cex=1.4,pch=0,lwd=2,
     main='Uniform Data Clustering, K=4')
```

Uniform Data Clustering, K=4



```
plot(U,col=(km.out[[6]]$cluster),cex=1.2,pch=0,lwd=2,
     main='Uniform Data Clustering, K=6')
```

Uniform Data Clustering, K=6



```
# Extracting log(Sk) from the list
U.logSk=rep(NA,10)
for (i in 1:10) {U.logSk[i]=log(km.out[[i]]$tot.withinss)}
```

```
# Extracting log(Sk) for the Artificial Dataset 2 output
set.seed(1234)
cg.sart=clusGap(artdata,kmeans, K.max=10, B=300, d.power=2,
                 spaceH0 = "scaledPCA", nstart=100, iter.max=100)
AD2.logSk=cg.sart$Tab[1,]
```

```
# Plotting the two sequences of indexes
plot(1:10,U.logSk,type='l',col='grey',lty=2,ylim=c(min(AD2.logSk),11),
     ylab="log Sk",main="Artificial Data2 log(Sk) vs. Uniform log(Sk)",xlab='K')
lines(1:10,AD2.logSk,lwd=2)
```

```
legend(7,11,legend=c("log(Sk) ArtData2","log(Sk) Uniform"), lty=1, lwd=3, col=c(1, "grey"))
```

Artificial Data2 log(Sk) vs. Uniform log(Sk)

What I'm doing here is kind of visualizing and reproducing the idea behind the gap statistic (previously computed). By means of the last plot I can compare the reduction in the logarithm of the objective function with how it is expected to behave under uniform distribution on the same "domain", for which k-means is applied in the same way. Now we can also see graphically that the gap between the two curves increases with the number of clusters, and at $k=3$ this gap significantly widens.

Exercise 2

Exercise 3

Exercise 4

- Using the gapnc-function, generate 100 data sets according to the specifications for Artificial Data Set 1, estimate the number of clusters by all four possible ways to run clusGap and compile four different vectors with all 100 estimated numbers of clusters for each version of clusGap. At the end look at the four distributions of estimated numbers of clusters using a suitable graphical display and comment on them. Are there better or worse results for some of the clusGap-versions? Do the same thing with another model generating data

```
# Generating 100 sets according to the specifications for Artificial DataSet 1
data=list()
simruns=100
```

```
kmmax=10
```

```
library(cluster)
```

```
set.seed(123)
```

```
for (i in 1:simruns) {
```

```
  v1 <- c(rnorm(50,0,1), rsn(70,5,1,8), rnorm(30,6,1))
```

```
  v2 <- c(rnorm(50,0,1), rsn(70,6,1,8), 8+rtr(30,5))
```

```
  data[[i]] <- cbind(v1,v2)
```

```
}
```

```
out.orig=list()
```

```
set.seed(1)
```

```
for (i in 1:100) {
```

```
  out.orig[[i]]=gapnc(data[[i]], K.max = kmmax ,spaceH0='original',SE.factor=1)
```

```
  iter.max=100, nstart=3)
```

```
# Collecting the outputs of the k-means function in a list, from 1 to K.max
```

```
km.out=list()
```

```
set.seed(1234)
```

```
for (i in 1:10) {km.out[[i]]=kmeans(U,iter.max = 100,nstart = 100, centers = i)}
```

```
# Three scatterplots of the uniform data with different K
```

```
plot(U,col=(km.out[[2]]$cluster),cex=1.2,pch=0,lwd=2,
```

```
     main='Uniform Data Clustering, K=2')
```

Uniform Data Clustering, K=2


```
plot(U,col=(km.out[[4]]$cluster),cex=1.4,pch=0,lwd=2,
```

```
     main='Uniform Data Clustering, K=4')
```

Uniform Data Clustering, K=4


```
plot(U,col=(km.out[[6]]$cluster),cex=1.2,pch=0,lwd=2,
```

```
     main='Uniform Data Clustering, K=6')
```

Uniform Data Clustering, K=6


```
# Extracting log(Sk) from the list
```

```
U.logSk=rep(NA,10)
```

```
for (i in 1:10) {U.logSk[i]=log(km.out[[i]]$tot.withinss)}
```

```
# Extracting log(Sk) for the Artificial Dataset 2 output
```

```
set.seed(1234)
```

```
cg.sart=clusGap(artdata,kmeans, K.max=10, B=300, d.power=2,
```

```
                 spaceH0 = "scaledPCA", nstart=100, iter.max=100)
```

```
AD2.logSk=cg.sart$Tab[1,]
```

```
# Plotting the two sequences of indexes
```

```
plot(1:10,U.logSk,type='l',col='grey',lty=2,ylim=c(min(AD2.logSk),11),
```

```
     ylab="log Sk",main="Artificial Data2 log(Sk) vs. Uniform log(Sk)",xlab='K')
```

```
lines(1:10,AD2.logSk,lwd=2)
```

```
legend(7,11,legend=c("log(Sk) ArtData2","log(Sk) Uniform"), lty=1, lwd=3, col=c(1, "grey"))
```

Artificial Data2 log(Sk) vs. Uniform log(Sk)

I then extracted the optimal number of clusters for each simulated dataset and each variant of the clusGap function, obtaining 4 vectors of length 100. These have been represented by means of a heatmap in which orange bars represent deviations from the "correct" number of groups. As you can see, methods based on the rotation of principal axes perform slightly worse than those which keep the original space of the sample, but these findings cannot be generalized. Results are also tabulated under here in 2-by-2 tables.

```
# Compile 4 different vectors with all 100 est numbers of clusters for each clusGap run
```

```
clu.orig=rep(NA,simruns)
```

```
for (i in 1:simruns) {clu.orig[i]=clusGap(data[[i]], K.max = kmmax ,spaceH0='original',SE.factor=1)
```

```
  iter.max=100, nstart=3)}
```

```
clu.orig2=rep(NA,simruns)
```

```
for (i in 1:simruns) {clu.orig2[i]=clusGap(data[[i]], K.max = kmmax ,spaceH0='PCA',SE.factor=1)
```

```
  iter.max=100, nstart=3)}
```

```
clu.pca1=rep(NA,simruns)
```

```
for (i in 1:simruns) {clu.pca1[i]=clus
```

EXERCISE 3



The Mahalanobis distance is defined as :

$$d_M(x, y) = \sqrt{(x-y)^T S^{-1} (x-y)}$$

where S is known to be SYMMETRIC and POSITIVE SEMI-DEFINITE. By the SPECTRAL THEOREM, there exists always a decomposition.

$$S = U^T D U = (U^T D^{1/2})(D^{1/2} U) = Q^T Q$$

where $D = \text{diag}(\lambda_1, \dots, \lambda_p)$ is a diagonal $p \times p$ matrix and U is an orthonormal $p \times p$ matrix of eigenvectors.

Then, let $\tilde{x} = D^{-1/2}(U^{-1})^T x$ and $\tilde{y} = D^{-1/2}(U^{-1})^T y$
we can write the Mahalanobis distance between x and y as the euclidean distance between \tilde{x} and \tilde{y}

$$d_M(x, y) = \sqrt{(x-y)^T S^{-1} (x-y)} = \sqrt{(\tilde{x}-\tilde{y})^T (\tilde{x}-\tilde{y})} = d_E(\tilde{x}, \tilde{y})$$

WHICH IS KNOWN TO SATISFY THE TRIANGULAR INEQUALITY

$$d_E(\tilde{x}, \tilde{y}) \leq d_E(\tilde{x}, w) + d_E(w, \tilde{y}) \quad w \in \mathbb{R}^p$$

$$\text{since } \|\tilde{x} - \tilde{y}\|_E = \|\tilde{x} - w + w - \tilde{y}\|_E \leq \|\tilde{x} - w\|_E + \|w - \tilde{y}\|_E$$

The SIMPLE MATCHING DISTANCE is defined as:



$$d_{SM}(x, y) = \frac{1}{P} \sum_{j=1}^P 1(x_j \neq y_j)$$

For binary data (*) it can be interpreted as PROPORTIONAL TO THE MANHATTAN DISTANCE :

$$\frac{1}{P} \sum_{j=1}^P 1(x_j \neq y_j) = \frac{1}{P} \sum_{j=1}^P |x_j - y_j| = \frac{1}{P} d_{L_1}(x, y)$$

Then, $\frac{1}{P} d_{L_1}(x, y) \leq \frac{1}{P} d_{L_1}(x, z) + \frac{1}{P} d_{L_1}(z, y)$

$$d_{L_1}(x, y) \leq d_{L_1}(x, z) + d_{L_1}(z, y)$$

without loss of generality, assume that $z = 0_P$

$$\sum_{j=1}^P |x_j - y_j| \leq \sum_{j=1}^P |x_j| + \sum_{j=1}^P |y_j|$$

which holds using the fact that $|Q+b| \leq |Q| + |b|$ for any reals Q, b .

- (*) In general, one could recode any categorical variable by splitting the attributes into a set of dummy variables; in fact, triangular inequality holds in general for the SMD.

EXERCISE 4

Increasing the value of the S.E. factor q may lead to smaller optimal number K_0 of clusters
(Keeping constant the dataset and the collection of B randomly generated uniform samples)

so $K_{0,1} \geq K_{0,2}$ because

$$GAP(K^*) - S_{K^*} > GAP(K^*) - 2S_{K^*}$$



$$S_{K^*} < 2S_{K^*}$$

ALWAYS TRUE SINCE S_{K^*} IS ALWAYS POSITIVE BEING THE ESTIMATED S.E. OF $GAP(K^*)$

