

Clustering

Brandon Kozak

30/10/2019

```
library(tidyverse)
library(BiocManager)
library(here)
```

Goals for this chapter

- Study different types of data that work well with clustering
- Learn the types of measures that determine clusters
- Uncover latent clustering via partitioning the data into tighter sets
- Look at non parametric algorithms such as k-means, k-medoids on real cell data
- Hierarchical clustering
- Use bootstrap to validate clusters

5.2 Why cluster?

5.2.1

It can lead to discoveries, for example in cancer biology.

Clustering is more general than the EM Approach and can be applied to more complex data. This is because many of the clustering techniques do not assume anything about the underlying generating mechanism of the data.

Bioconductor has 212 packages on clustering as of now!

5.3 How do we measure similarity?

We are looking at bird data.

Our first question is what variables are we going to use. Weight and size will give different clustering compared to diet or habitat.

Next we list a couple of ways to define distance.

Assume our data exists in a p -dimensional space, such that point $A = (a_1, \dots, a_p)$ and point $B = (b_1, \dots, b_p)$.

- **Euclidean Distance:** defined to be the square root of the sum of the squared differences between each component of each point. That is, $d(A, B) = \sqrt{(a_1 - b_1)^2 + \dots + (a_p - b_p)^2}$
- **Manhattan Distance:** defined to be the sum of the absolute differences in each component of each point. That is, $d(A, B) = |a_1 - b_1| + \dots + |a_p - b_p|$
- **Maximum Distance:** defined to be the max of the absolute differences between A and B, that is $d_{\infty}(A, B) = \max |a_i - b_i|$

- **Weighted Euclidean Distance:** similar to the euclidean distance, but now we apply certain weight (w_1, \dots, w_p) to each difference.
- **Mahalanobis Distance:** is a special type of weighted euclidean distance that uses the sample covariance matrix of the data, that is $d(A, B) = \sqrt{(A - B)^T S_n^{-1} (A - B)}$, where S_n^{-1} is the inverse of the sample covariance matrix.
- **Malinowski Distance:** is a general form of the Euclidean distance, where we can take the exponent to be m rather than 2. That is, $d(A, B) = ((a_1 - b_1)^m + \dots + (a_p - b_p)^m)^{\frac{1}{p}}$
- **Edit and Hamming Distance:** Simplify the number of differences for each index of a character string. Typical used in DNA sequences.

5.3.1 Computations related to distances in R

We can use the `dist()` function. By default we get the euclidean distance, but can set the “method” parameter

```
# Our points
mx = c(0, 0, 0, 1, 1, 1)
my = c(1, 0, 1, 1, 0, 1)
mz = c(1, 1, 1, 0, 1, 1)
mat = rbind(mx, my, mz)

dist(mat)
```

```
##           mx           my
## my 1.732051
## mz 2.000000 1.732051
```

```
dist(mat, method = "binary")
```

```
##           mx           my
## my 0.6000000
## mz 0.6666667 0.5000000
```

What about the Jaccard distance?

```
mut = read_csv(here("Data", "HIVmutations.csv"))
mut[1:3, 10:16]
```

```
## # A tibble: 3 x 7
##   p32I p33F p34Q p35G p43T p46I p46L
##   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1     0     1     0     0     0     0     0
## 2     0     1     0     0     0     1     0
## 3     0     1     0     0     0     0     0
```

```
library("vegan")
mutJ = vegdist(mut, "jaccard")
mutC = sqrt(2 * (1 - cor(t(mut))))
mutJ
```

```
##           1           2           3           4
## 2 0.8000000
## 3 0.7500000 0.8888889
## 4 0.9000000 0.7777778 0.8461538
## 5 1.0000000 0.8000000 0.8888889 0.9000000
```

5.4 Non parametric mixture detection

There are two k-methods to clustering. k-means and k-medoids.

k-means uses the mean in it's algorithms, while k-medoids uses the true center.

One example of a k-medoid algorithm is the partitioning around medoids (PAM) algorithm. The steps are as follows:

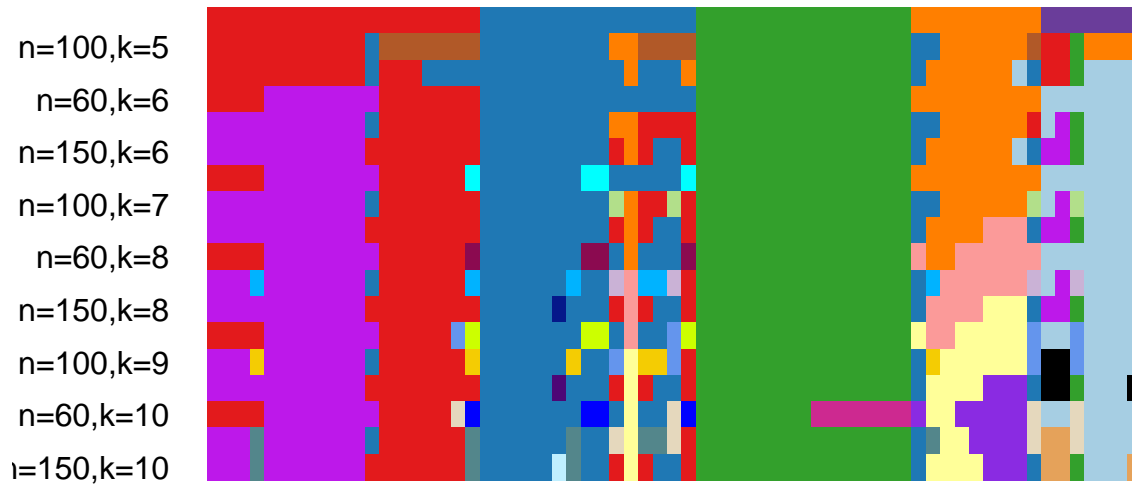
- Start with p features and n observations
- Randomly pick k distinct cluster centers out of the n observations.
- Assign each of the remaining observations to the group whose center is the closest
- For each group, choose a new center from the observations in that group, such that the sum of the distances of the groups members to the new center is minimal.
- Repeat steps 3 and 4 until the groups stabilize.

5.4.2 Tight clusters with re sampling

We can repeat clustering algorithms on the same data but with different starting points. Observations that are almost always grouped together are called tight clusters.

Here is an example using the clusterExperiment package.

```
library("clusterExperiment")
data("fluidigm", package = "scRNAseq")
se = fluidigm[, fluidigm$Coverage_Type == "High"]
assays(se) = list(normalized_counts =
  round(limma::normalizeQuantiles(assay(se))))
ce = clusterMany(se, clusterFunction = "pam", ks = 5:10, run = TRUE,
  isCount = TRUE, reduceMethod = "var", nFilterDims = c(60, 100, 150))
clusterLabels(ce) = sub("FilterDims", "", clusterLabels(ce))
plotClusters(ce, whichClusters = "workflow", axisLine = -1)
```



We can see that the right most red, left most blue, and green clusters appear to be tight clusters.

5.5 examples

5.5.1 Flow cytometry and mass cytometry

```
library("flowCore")
library("flowViz")
fcsB = read.FCS("../data/Bendall_2011.fcs")
slotNames(fcsB)
```

```
## [1] "exprs"      "parameters" "description"
```

Look at the structure of the fcsB object (hint: the colnames function). How many variables were measured?

```
41
```

```
length(colnames(fcsB))
```

```
## [1] 41
```

Subset the data to look at the first few rows (hint: use Biobase::exprs(fcsB)). How many cells were measured?

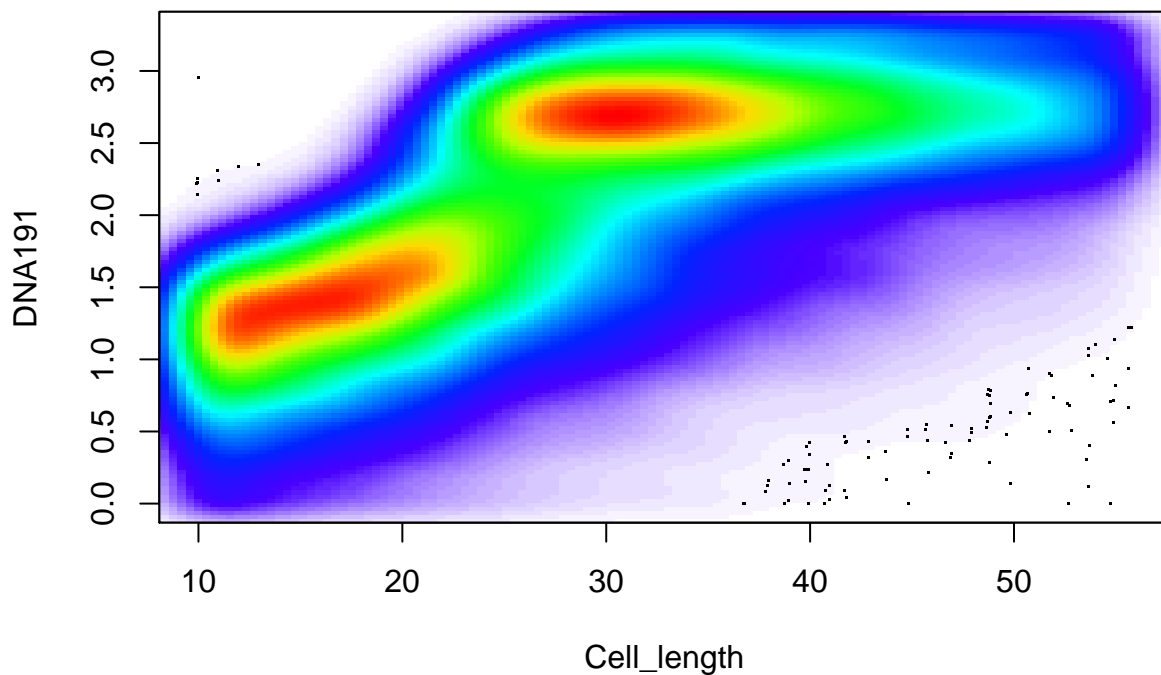
```
nrow(Biobase::exprs(fcsB))
```

```
## [1] 91392
```

5.5.2 Data preprocessing

```
markersB = readr::read_csv("../data/Bendall_2011_markers.csv")
mt = match(markersB$isotope, colnames(fcsB))
stopifnot(!any(is.na(mt)))
colnames(fcsB)[mt] = markersB$marker
```

```
flowPlot(fcsB, plotParameters = colnames(fcsB)[2:3], logy = TRUE)
```



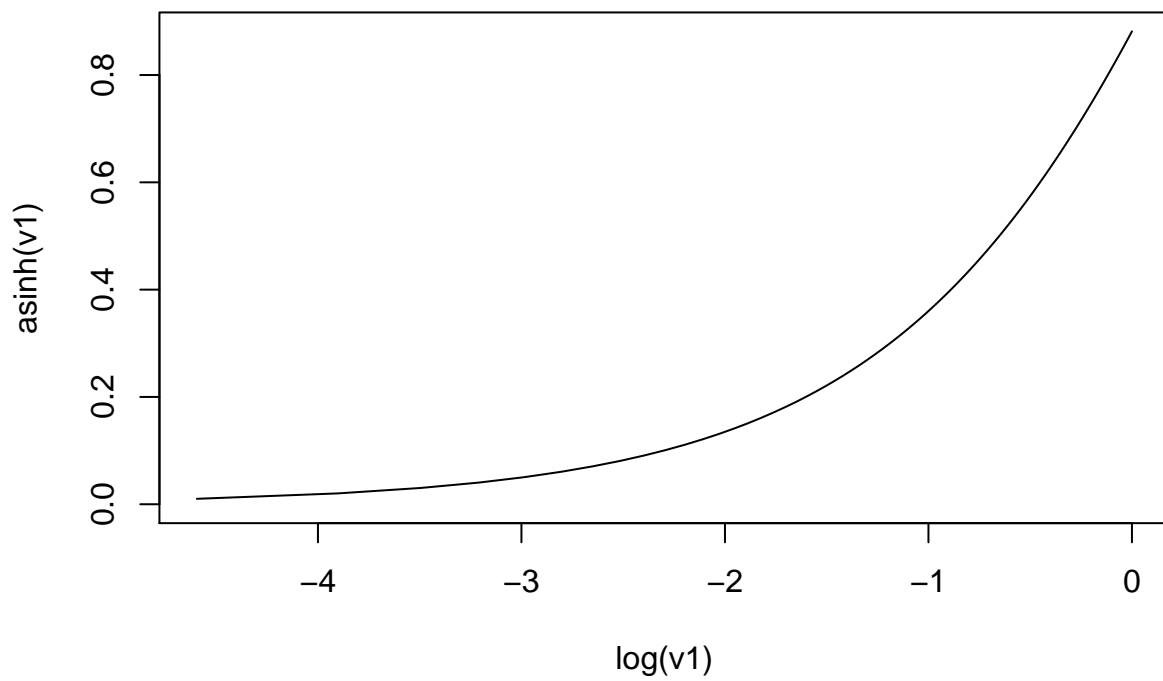
Here we can see clear clustering between cell length and DNA191.

A common data transformation is the hyperbolic arcsin.

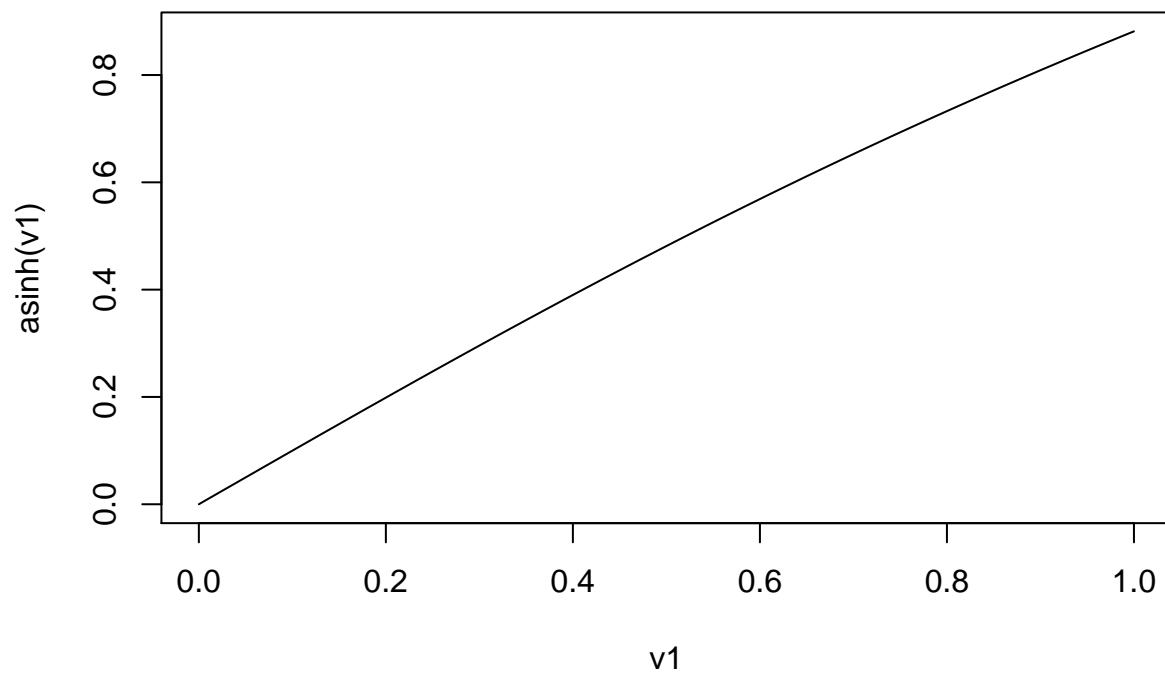
$$\operatorname{asinh}(x) = \log(x + \sqrt{x^2 + 1})$$

We can view this transformation for different values of x.

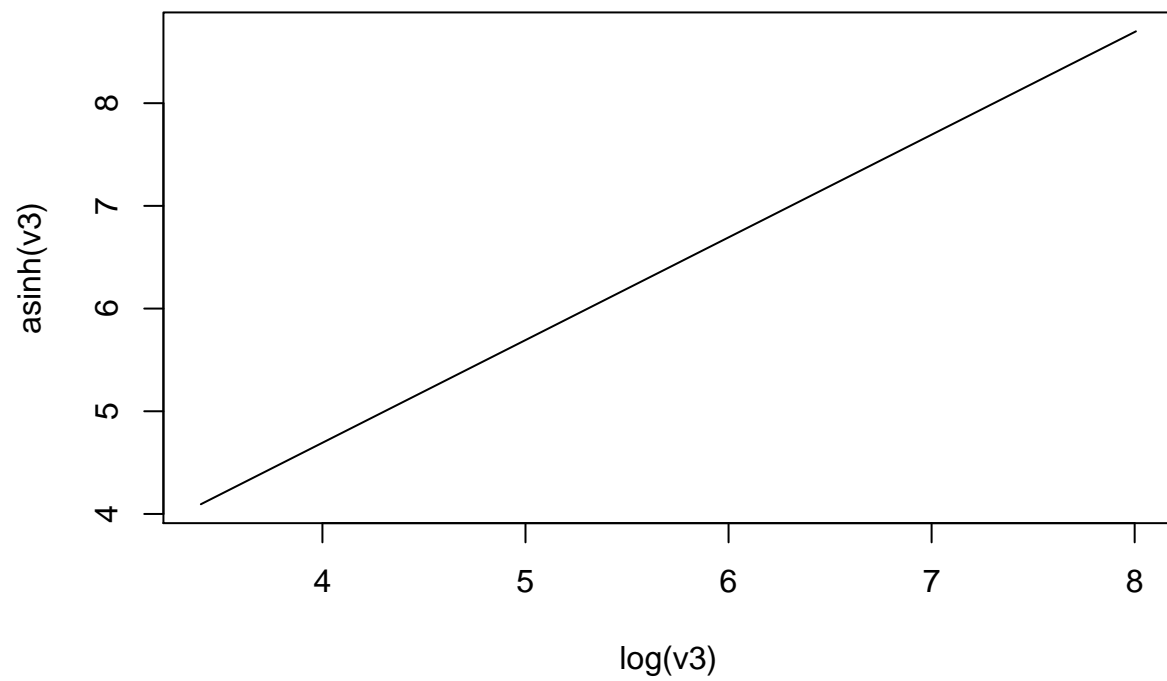
```
v1 = seq(0, 1, length.out = 100)
plot(log(v1), asinh(v1), type = 'l')
```



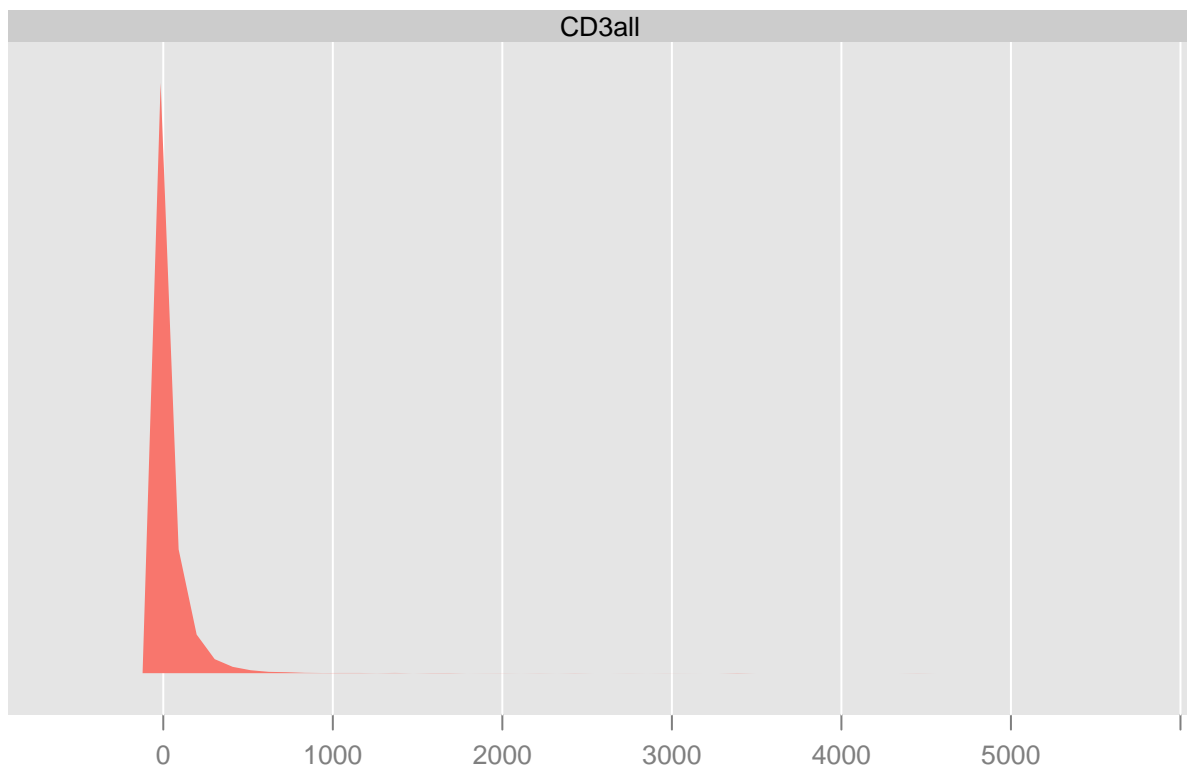
```
plot(v1, asinh(v1), type = 'l')
```



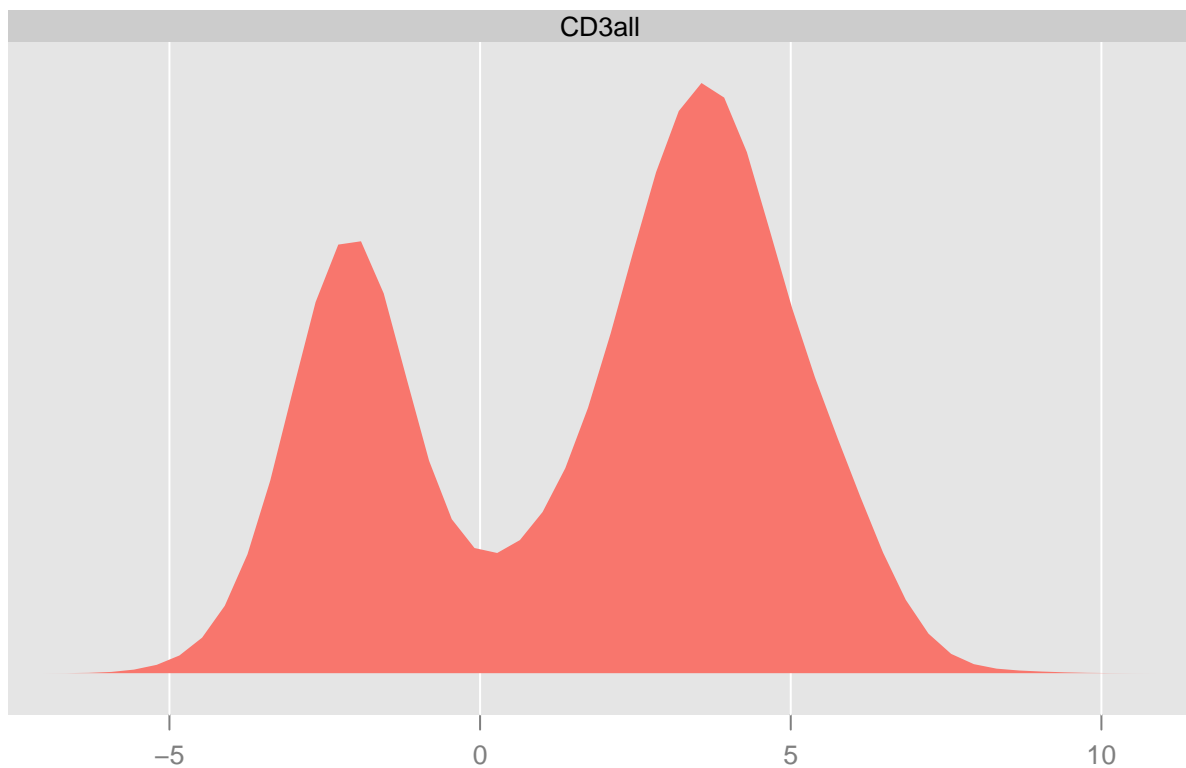
```
v3 = seq(30, 3000, length = 100)
plot(log(v3), asinh(v3), type= 'l')
```



```
asinhtrs = arcsinhTransform(a = 0.1, b = 1)
fcsBT = transform(fcsB,
  transformList(colnames(fcsB)[-c(1, 2, 41)], asinhtrs))
densityplot(~`CD3all`, fcsB)
```

```
densityplot( ~`CD3all`, fcsBT)
```



How many dimensions does the following code use to split the data into 2 groups using k-means ?

```
kf = kmeansFilter("CD3all" = c("Pop1", "Pop2"), filterId="myKmFilter")
fres = flowCore::filter(fcsBT, kf)
summary(fres)
```

```
## Pop1: 33429 of 91392 events (36.58%)
## Pop2: 57963 of 91392 events (63.42%)
```

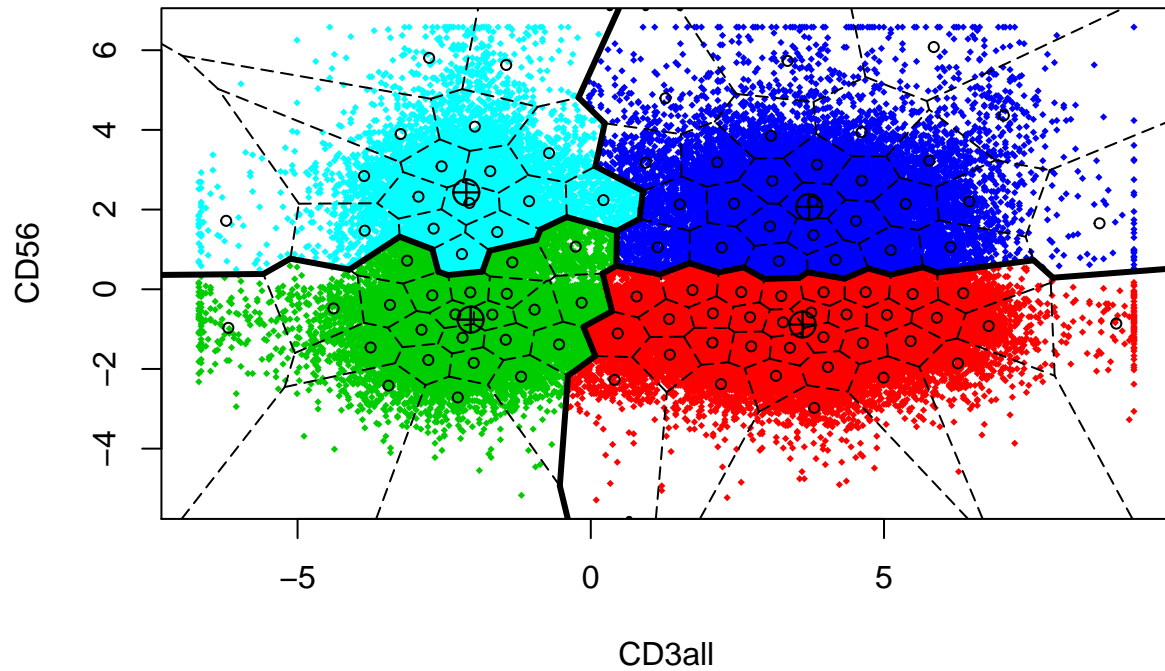
```
fcsBT1 = flowCore::split(fcsBT, fres, population = "Pop1")
fcsBT2 = flowCore::split(fcsBT, fres, population = "Pop2")
```

naive projection of the data into the two dimensions spanned by the CD3 and CD56 markers, clustering was performed using kmeans.

```
library("flowPeaks")
fp = flowPeaks(Biobase::exprs(fcsBT)[, c("CD3all", "CD56")])
```

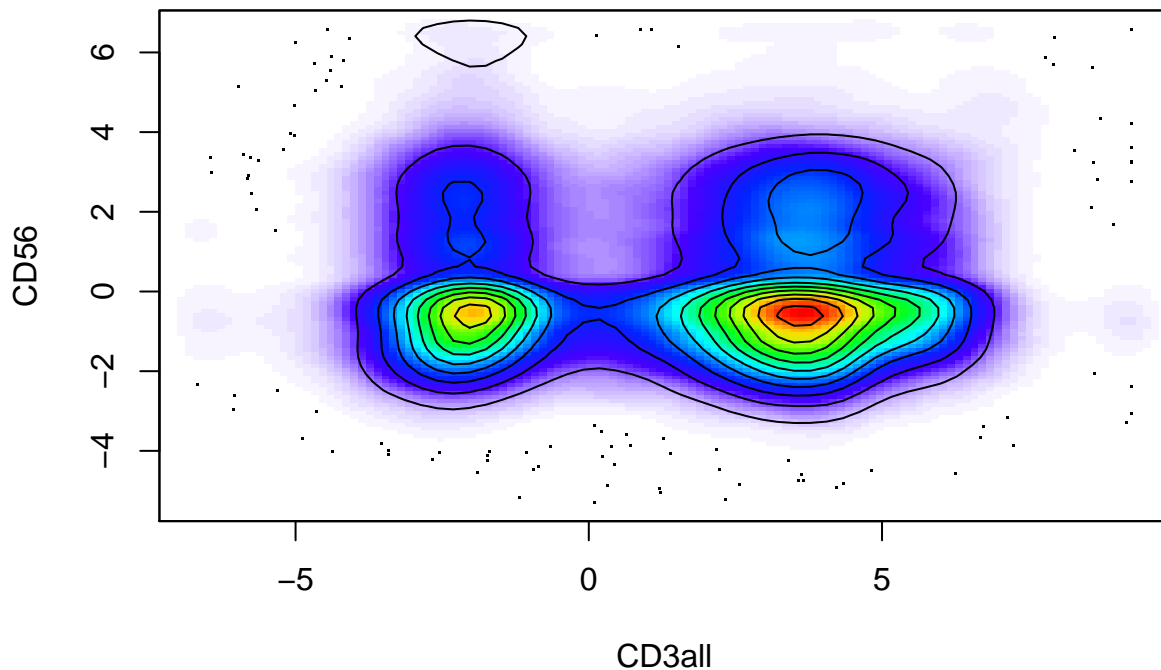
```
##          step 0, set the initial seeds, tot.wss=17597.8
##          step 1, do the rough EM, tot.wss=11900.7 at 0.35 sec
##          step 2, do the fine transfer of Hartigan-Wong Algorithm
##          tot.wss=11846.3 at 0.691 sec
```

```
plot(fp)
```



To avoid overplotting, we can use a 2d kernel density plot.

```
flowPlot(fcsBT, plotParameters = c("CD3all", "CD56"), logy = FALSE)  
contour(fcsBT[, c(40, 19)], add = TRUE)
```

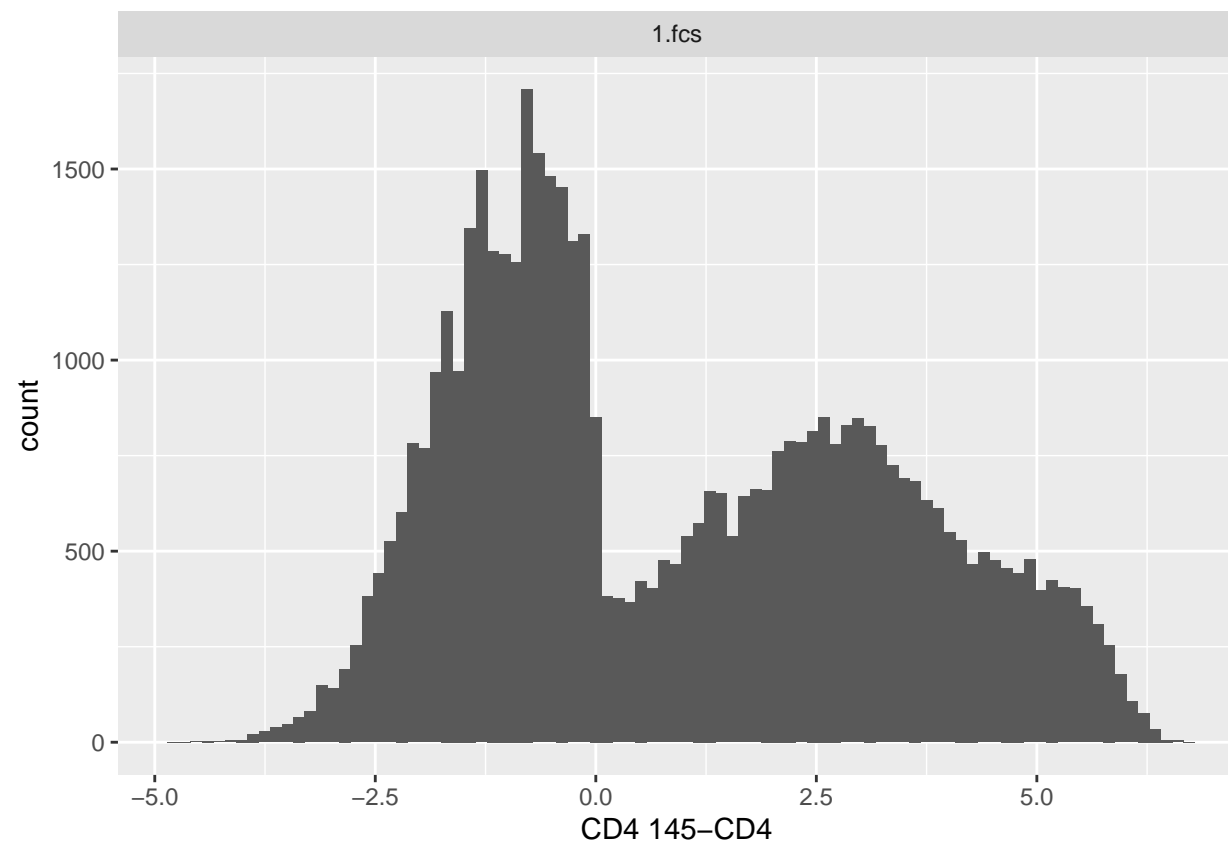


A more recent Bioconductor package, `ggcyto`, has been designed to enable the plotting of each patient in a different facet using `ggplot`.

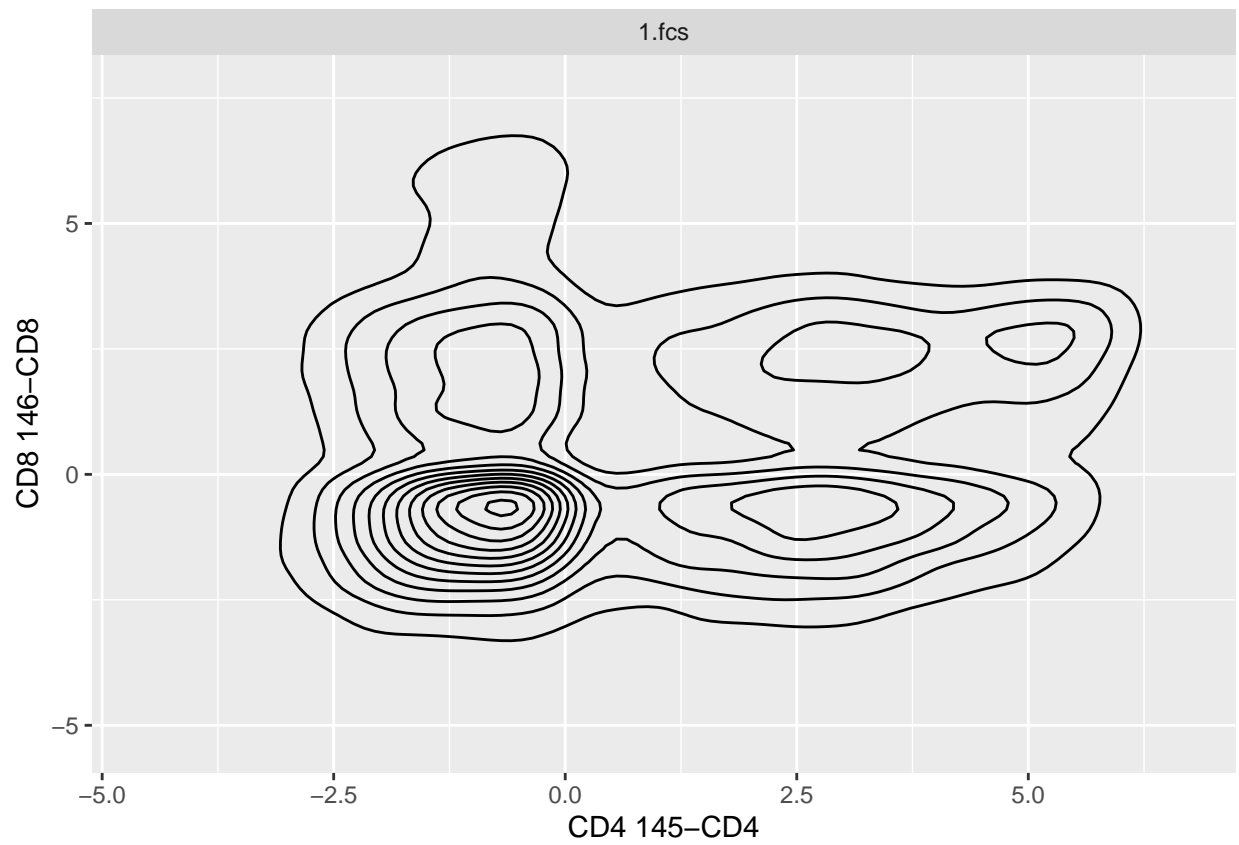
Try comparing the output using this approach to what we did above using the following:

```
library("ggcyto")
library("labeling")
ggcd4cd8=ggcyto(fcsB,aes(x=CD4,y=CD8))
ggcd4=ggcyto(fcsB,aes(x=CD4))
ggcd8=ggcyto(fcsB,aes(x=CD8))
p1=ggcd4+geom_histogram(bins=60)
p1b=ggcd8+geom_histogram(bins=60)
asinhT = arcsinhTransform(a=0,b=1)
transl = transformList(colnames(fcsB)[-c(1,2,41)], asinhT)
fcsBT = transform(fcsB, transl)
p1t=ggcyto(fcsBT,aes(x=CD4))+geom_histogram(bins=90)
p2t=ggcyto(fcsBT,aes(x=CD4,y=CD8))+geom_density2d(colour="black")
p3t=ggcyto(fcsBT,aes(x=CD45RA,y=CD20))+geom_density2d(colour="black")

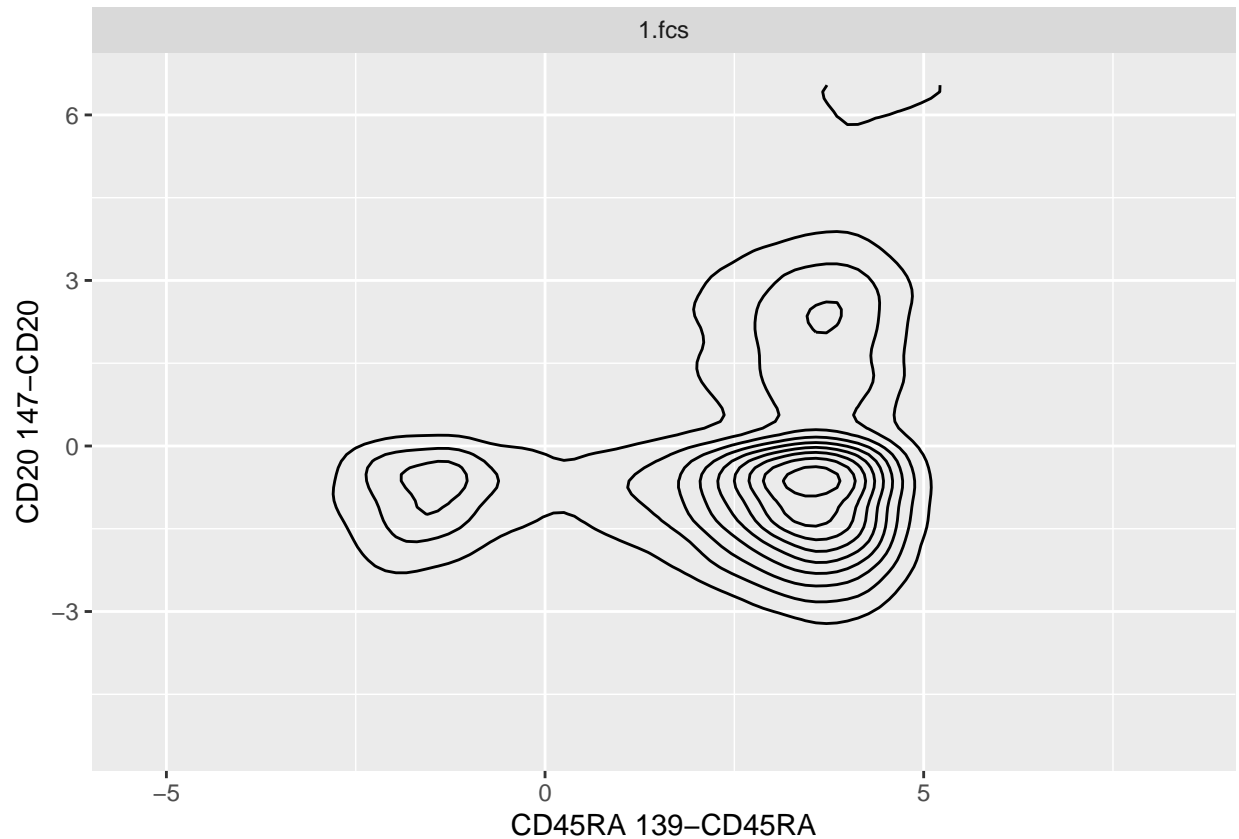
p1t
```



p2t



p3t



Output looks similar, but this method appears to be less sensitive to clusters.

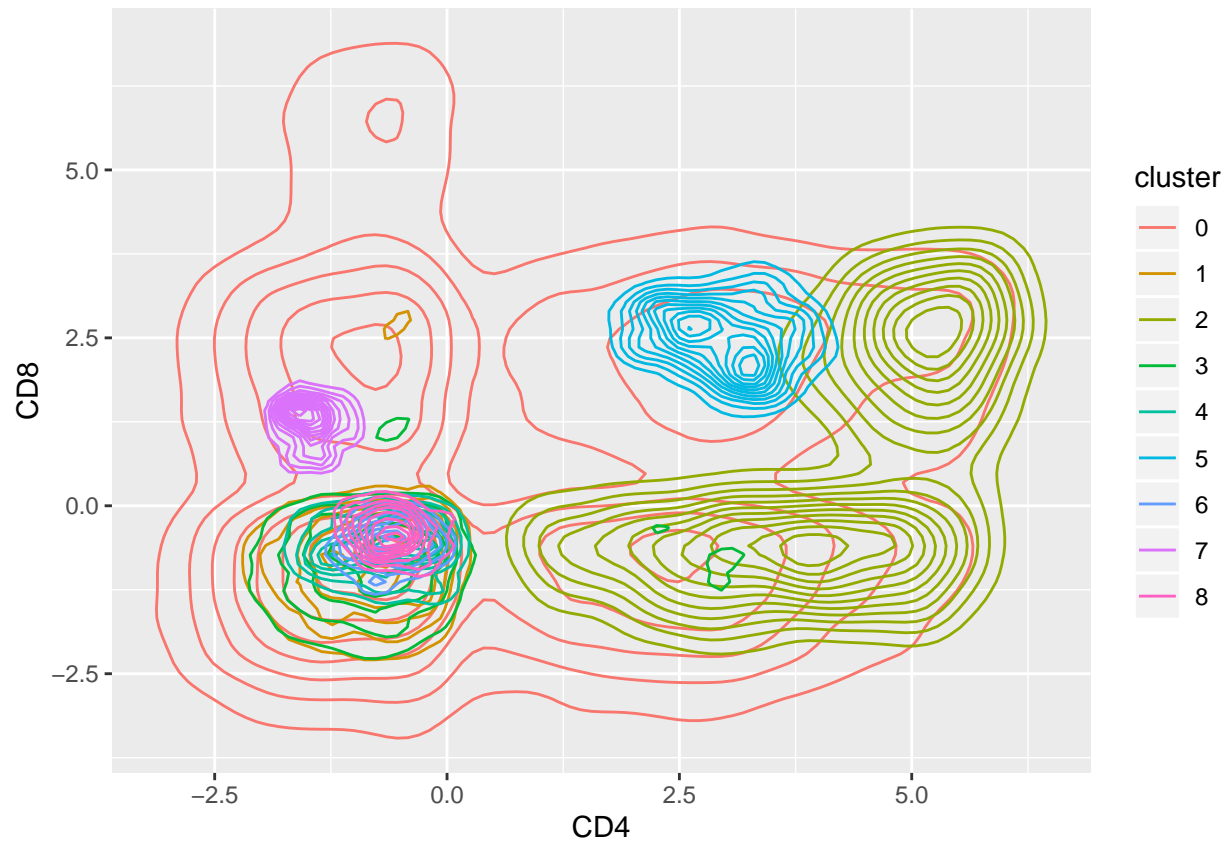
5.5.3 Density-based clustering

Used when we have only a few markers but a large amount of data (cells).

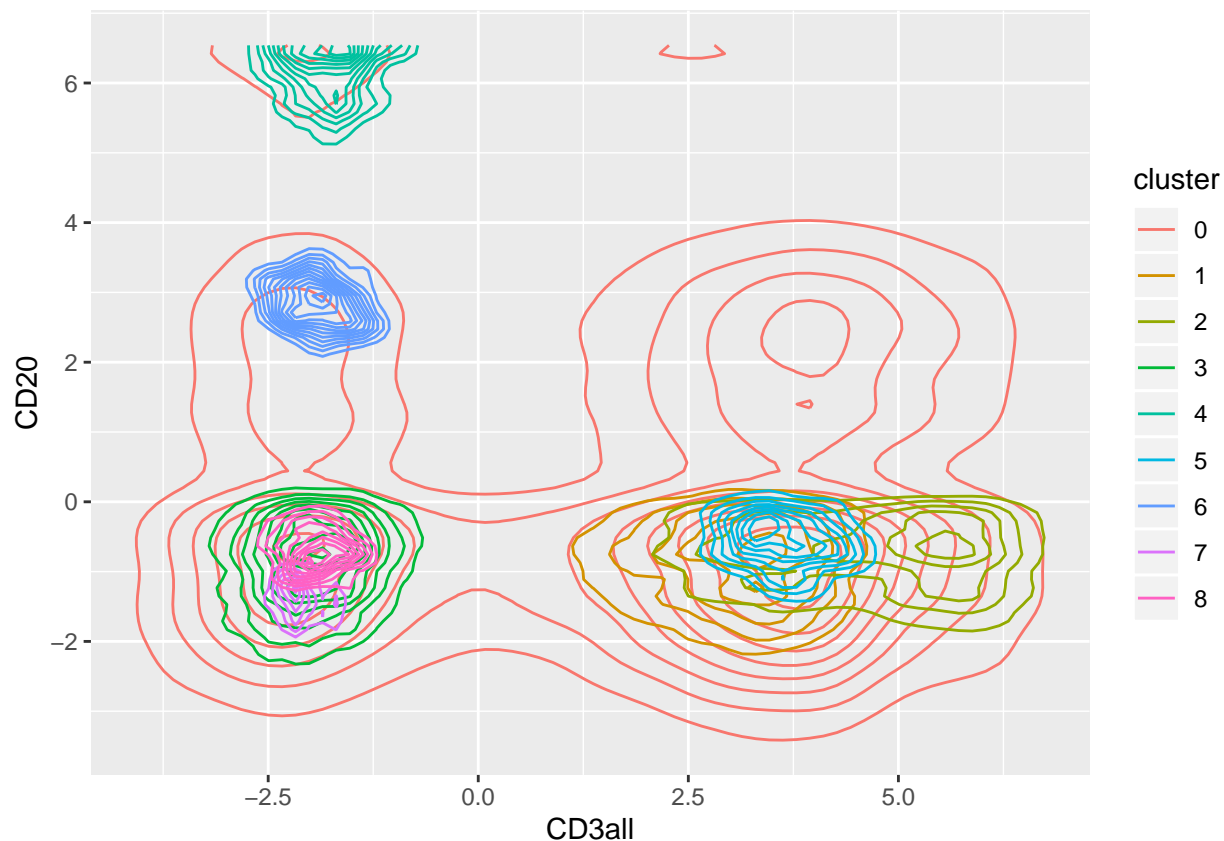
```
library("dbscan")
mc5 = Biobase::exprs(fcsBT)[, c(15,16,19,40,33)]
res5 = dbscan::dbscan(mc5, eps = 0.65, minPts = 30)
mc5df = data.frame(mc5, cluster = as.factor(res5$cluster))
table(mc5df$cluster)
```

```
##
##      0      1      2      3      4      5      6      7      8
## 75954 4031 5450 5310  259  257   63   25   43
```

```
ggplot(mc5df, aes(x=CD4, y=CD8, col=cluster))+geom_density2d()
```



```
ggplot(mc5df, aes(x=CD3all, y=CD20, col=cluster))+geom_density2d()
```

Try increasing the dimension to 6 by adding one CD marker-variables from the input data.

Then vary eps, and try to find four clusters such that at least two of them have more than 100 points.

Repeat this with 7 CD marker-variables, what do you notice?

```
mc6 = Biobase::exprs(fcsBT)[, c(15, 16, 19, 33, 25, 40)]
res = dbscan::dbscan(mc6, eps = 0.65, minPts = 20)
mc6df = data.frame(mc6, cluster = as.factor(res$cluster))
table(mc6df$cluster)
```

```
##
##      0      1      2      3      4      5      6
## 91068    34    61    20    67   121    21
```

```
mc7 = Biobase::exprs(fcsBT)[, c(11, 15, 16, 19, 25, 33, 40)]
res = dbscan::dbscan(mc7, eps = 0.95, minPts = 20)
mc7df = data.frame(mc7, cluster = as.factor(res$cluster))
table(mc7df$cluster)
```

```
##
##      0      1      2      3      4      5      6      7      8      9     10
## 90249    21   102   445   158   119    19   224    17    20    18
```

In general, the higher the dimension, the larger we need to set eps in order to find large enough clusters.

How does density-based clustering work?

Uses something called the density-connectedness criterion. That is, it looks at small neighborhood spheres of radius ϵ to see if points are connected.

The rest of this sounds very abstract and mathematical...

5.6 Hierarchical clustering

Bottom-up approach, where similar observations and subclasses are assembled iteratively.

Note that the order of the labels does not matter within sibling pairs and the horizontal distances are usually meaningless, while the vertical distances do encode some information.

There is also a top-down approach that takes all the objects and splits them sequentially according to a chosen criterion.

5.6.1 How to compute (dis)similarities between aggregated clusters?

We will need more than just the distances between all pairs of individual objects. We also need a way to calculate distances between the aggregates.

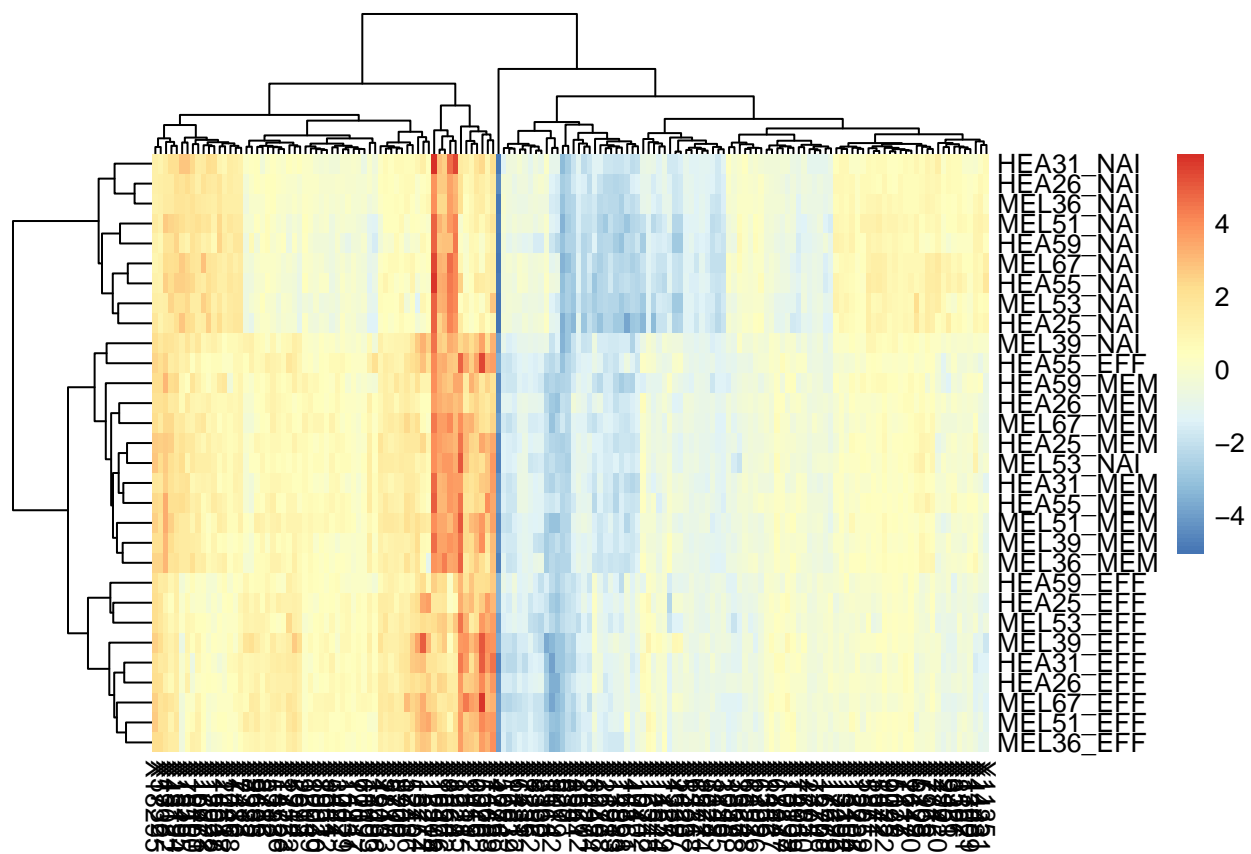
Hierarchical clustering for cell populations The Morder data are gene expression measurements for 156 genes on T cells of 3 types (naïve, effector, memory) from 10 patients (Holmes et al. 2005).

Using the pheatmap package, make two simple heatmaps, without dendrogram or reordering, for Euclidean and Manhattan distances of these data.

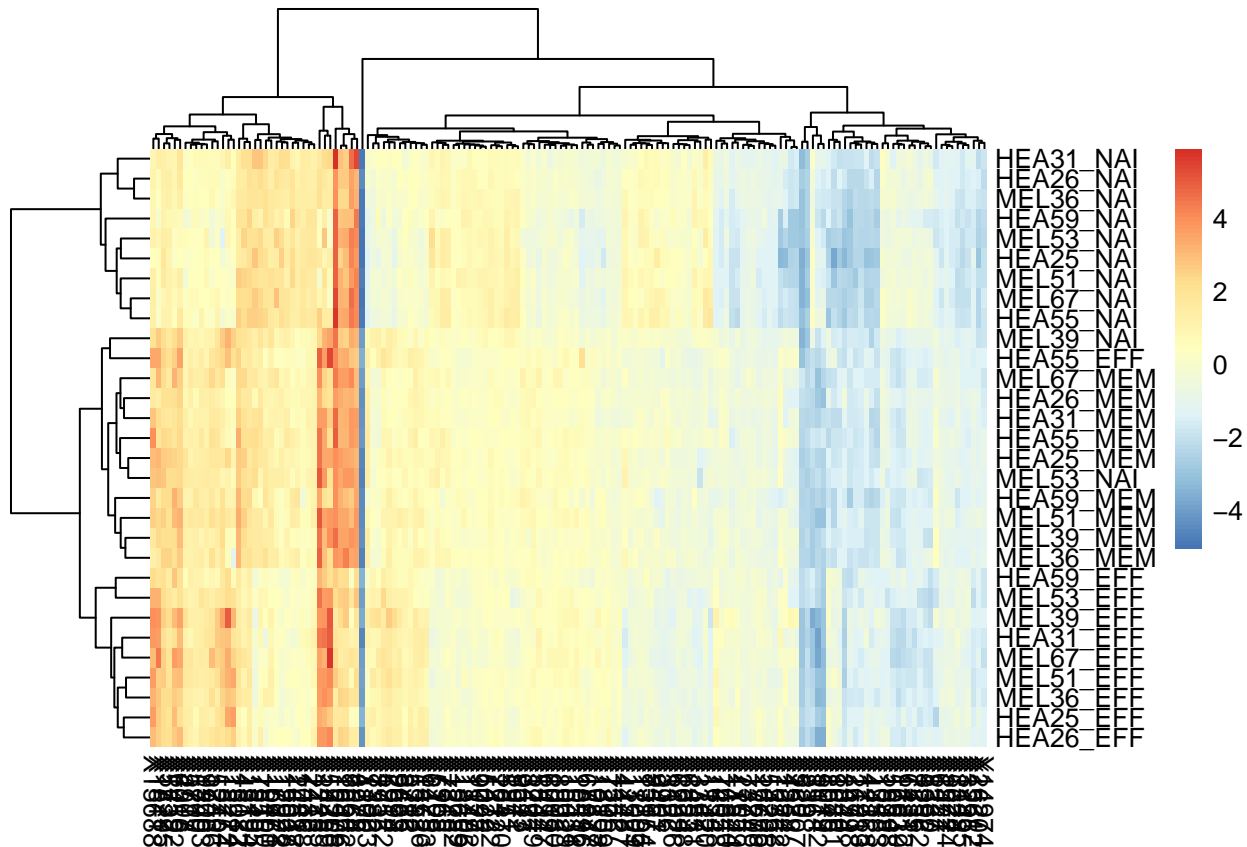
```
library(pheatmap)

load(here("Data", "Morder.RData"))

pheatmap(Morder)
```



```
pheatmap(Morder,clustering_distance_rows="manhattan",clustering_distance_cols="manhattan")
```



5.7 Validating and choosing the number of clusters

Need a way to validate our clusters with a number.

We define the within-groups sum of squared distances (WSS)

$$WSS_k = \sum_{i=1}^k \sum_{x_i \in C_i} d^2(x_i, \bar{x}_i)$$

where:

- k is the number of clusters
- C_i is the set of objects in the i^{th} cluster
- \bar{x}_i is the center of mass (average point) of the i^{th} cluster.

This is a good starting point, but we can always minimize this value by taking the number of clusters to be equal to the number of data points.

```
library("dplyr")
simdat = lapply(c(0, 8), function(mx) {
  lapply(c(0,8), function(my) {
    tibble(x = rnorm(100, mean = mx, sd = 2),
           y = rnorm(100, mean = my, sd = 2),
           class = paste(mx, my, sep = ":"))
  }) %>% bind_rows
}) %>% bind_rows
simdat
```

```
## # A tibble: 400 x 3
##       x     y class
##   <dbl> <dbl> <chr>
## 1  1.76  1.72 0:0
## 2 -2.34  2.48 0:0
## 3 -0.168 -2.90 0:0
## 4  1.29  1.17 0:0
## 5  0.724  0.547 0:0
## 6 -1.87  2.17 0:0
## 7  2.71  0.864 0:0
## 8 -0.321 -2.66 0:0
## 9 -0.795 -1.44 0:0
## 10 -2.39 -1.06 0:0
## # ... with 390 more rows
```

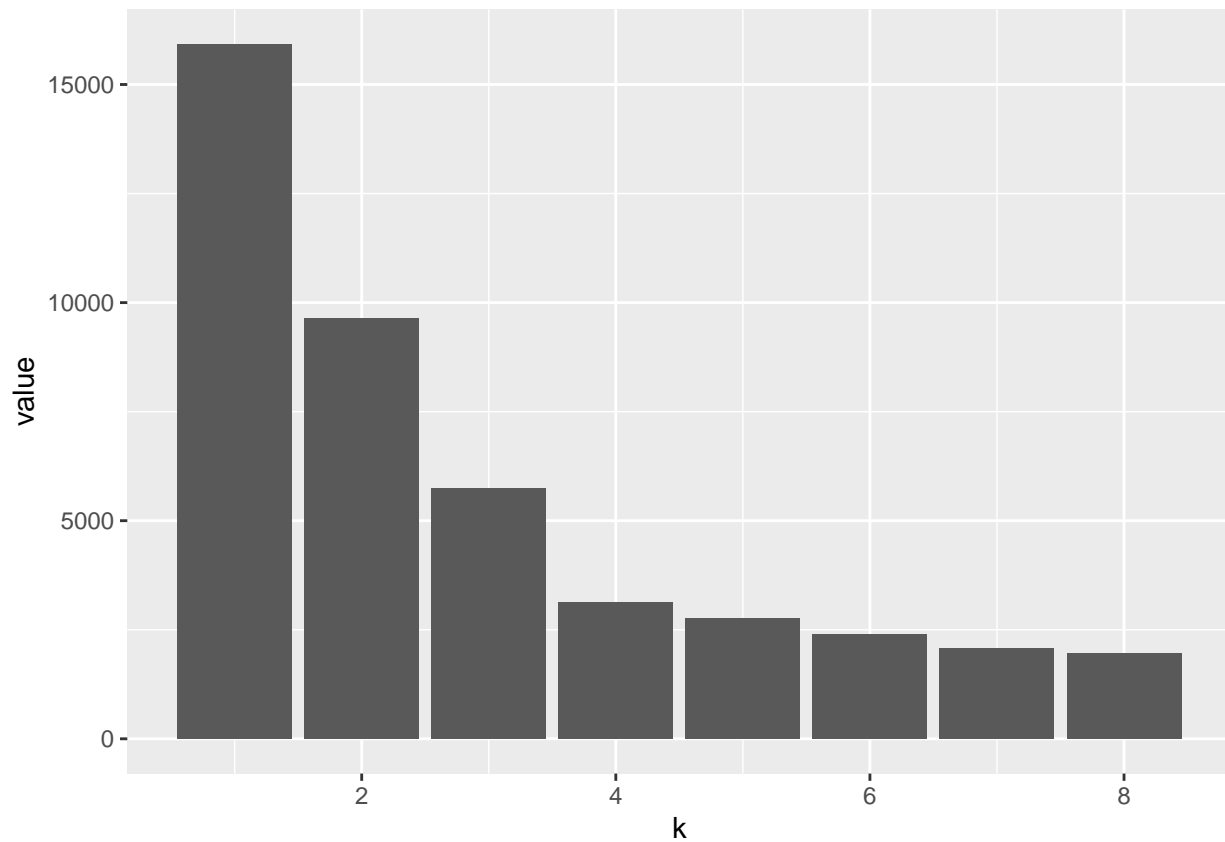
```
simdatxy = simdat[, c("x", "y")] # without class label
```

```
ggplot(simdat, aes(x = x, y = y, col = class)) + geom_point() +
  coord_fixed()
```



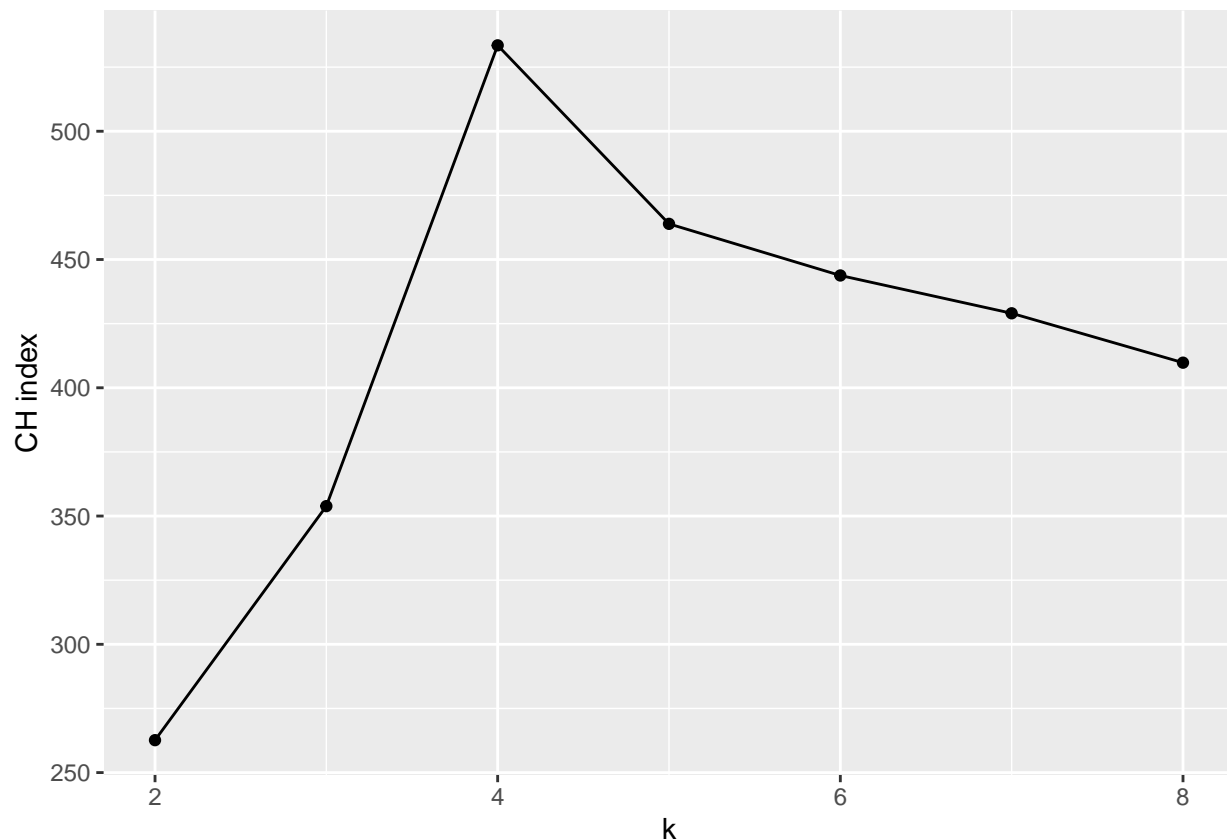
```
wss = tibble(k = 1:8, value = NA_real_)
wss$value[1] = sum(scale(simdatxy, scale = FALSE)^2)
for (i in 2:nrow(wss)) {
  km = kmeans(simdatxy, centers = wss$k[i])
  wss$value[i] = sum(km$withinss)
```

```
}
ggplot(wss, aes(x = k, y = value)) + geom_col()
```



We can also use the calinski-Harabasz index, it is the ratio between Between Sum of Squares and Within Sum of squares errors distances.

```
library("fpc")
library("cluster")
CH = tibble(
  k = 2:8,
  value = sapply(k, function(i) {
    p = pam(simdatxy, i)
    calinhara(simdatxy, p$cluster)
  })
)
ggplot(CH, aes(x = k, y = value)) + geom_line() + geom_point() +
  ylab("CH index")
```



5.7.1 Using the gap statistic

The basic idea is to take some transformation of the within-sum-of-squares (typically the log) and compare it to the averages from simulated data with less structure.

The general steps for computing the gap statistic are:

- Cluster the data with k clusters and compute WSS_k for various choices of k
- Generate B plausible reference data sets, using Monte Carlo sampling from a homogeneous distribution and redo Step 1 above for these new simulated data.
- Compute:

$$gap(k) = \bar{l}_k - \log WSS_k \text{ with}$$

$$\bar{l}_k = \frac{1}{B} \sum_{b=1}^B \log W_{k,b}^*$$

- Then we can use the standard deviation, defined as:

$$sd_k^2 = \frac{1}{B-1} \sum_{b=1}^B (\log W_{k,b}^* - \bar{l}_k)^2$$

in order to determine the best value for k .

The packages `cluster` and `clusterCrit` provide implementations.

Make a function that plots the gap statistic as in Figure 5.27.

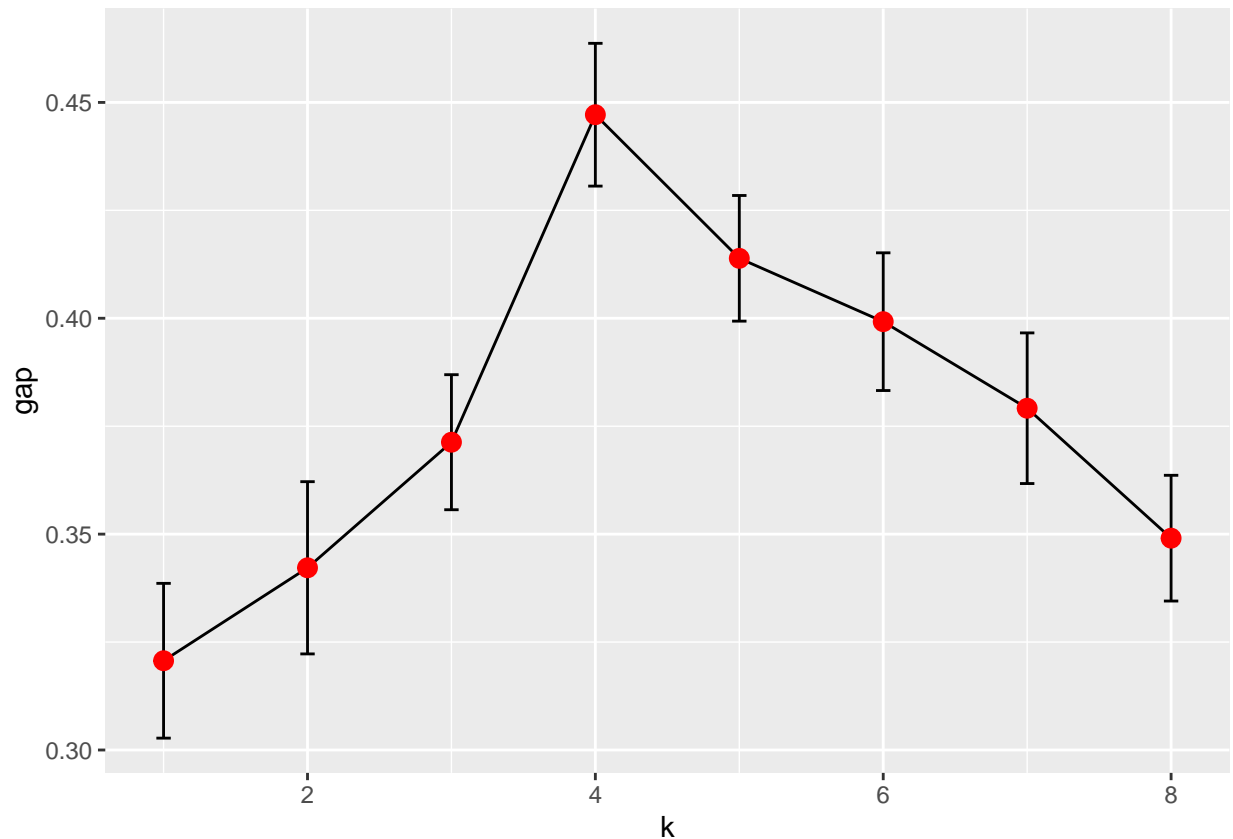
Show the output for the `simdat` example dataset clustered with the `pam` function.

```

library("cluster")
library("ggplot2")
pamfun = function(x, k)
  list(cluster = pam(x, k, cluster.only = TRUE))

gss = clusGap(simdatxy, FUN = pamfun, K.max = 8, B = 50,
  verbose = FALSE)
plot_gap = function(x) {
  gstab = data.frame(x$Tab, k = seq_len(nrow(x$Tab)))
  ggplot(gstab, aes(k, gap)) + geom_line() +
    geom_errorbar(aes(ymax = gap + SE.sim,
      ymin = gap - SE.sim), width=0.1) +
    geom_point(size = 3, col= "red")
}
plot_gap(gss)

```



```

library("Hiiragi2013")
data("x")

selFeats = order(rowVars(Biobase::exprs(x)), decreasing = TRUE)[1:50]
embmat = t(Biobase::exprs(x)[selFeats, ])
embgap = clusGap(embmat, FUN = pamfun, K.max = 24
  , verbose = FALSE)
k1 = maxSE(embgap$Tab[, "gap"], embgap$Tab[, "SE.sim"])
k2 = maxSE(embgap$Tab[, "gap"], embgap$Tab[, "SE.sim"],

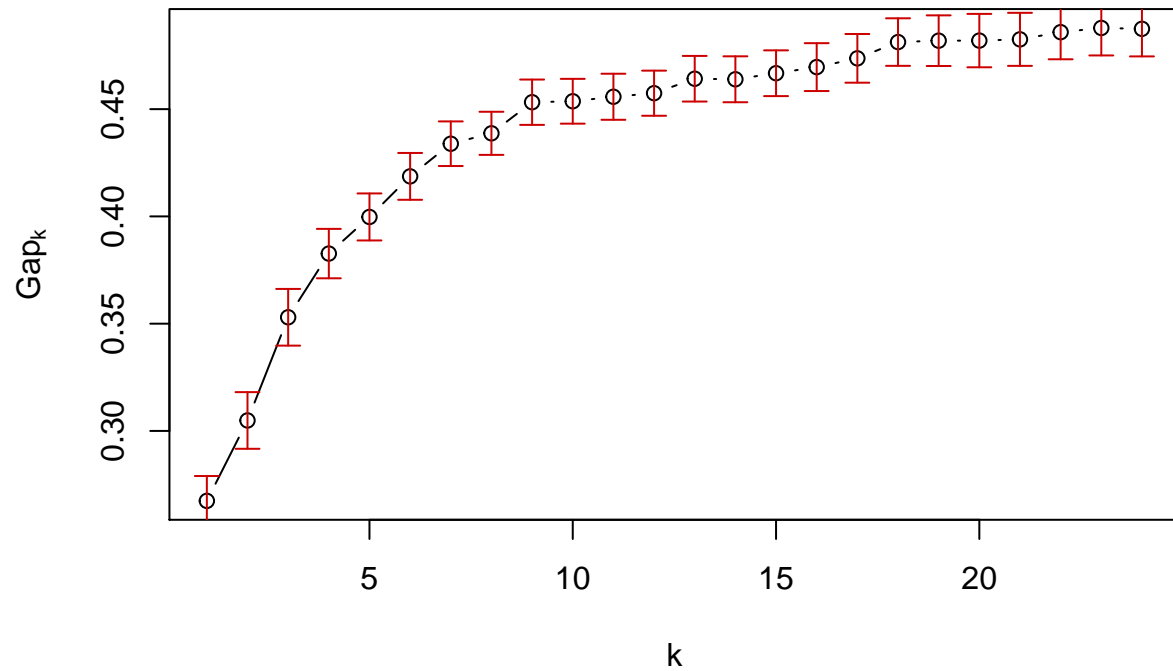
```



```
method = "Tibs2001SEmax")
c(k1, k2)
```

```
## [1] 10 7
```

```
plot(embgap, main = "")
```



```
c1 = pamfun(embmat, k = k1)$cluster
table(pData(x)[names(c1), "sampleGroup"], c1)
```

```
##          c1
##          1  2  3  4  5  6  7  8  9 10
## E3.25      23 11  1  1  0  0  0  0  0  0
## E3.25 (FGF4-KO) 0  0  1 16  0  0  0  0  0  0
## E3.5 (EPI)   2  1  0  0  0  8  0  0  0  0
## E3.5 (FGF4-KO) 0  0  8  0  0  0  0  0  0  0
## E3.5 (PE)    0  0  0  0  6  1  4  0  0  0
## E4.5 (EPI)   0  0  0  0  0  0  0  0  4  0
## E4.5 (FGF4-KO) 0  0  0  0  0  0  0  0  0 10
## E4.5 (PE)    0  0  0  0  0  0  0  4  0  0
```

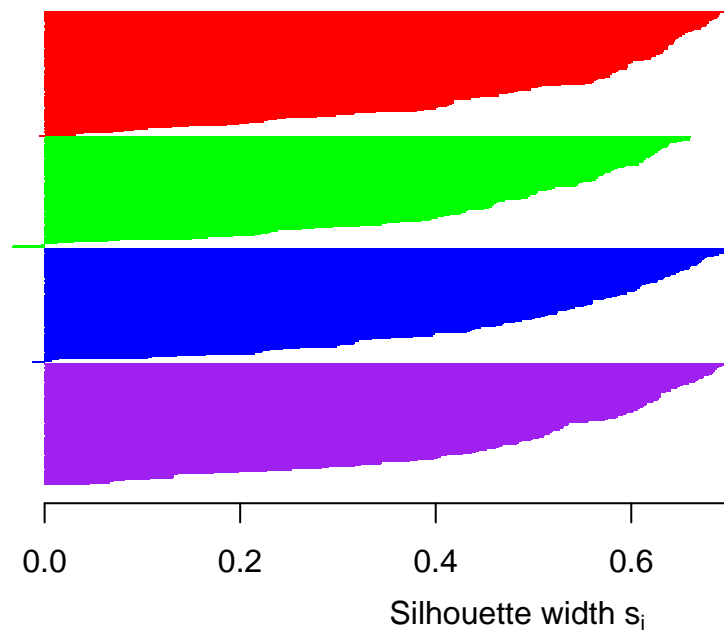
Exercises

5.1

```
library("cluster")
pam4 = pam(simdatxy, 4)
sil = silhouette(pam4, 4)
plot(sil, col=c("red","green","blue","purple"), main="Silhouette")
```

Silhouette

n = 400



4 clusters C_j

$j : n_j \mid \text{ave}_{i \in C_j} s_i$

1 : 106 | 0.50

2 : 94 | 0.46

3 : 97 | 0.50

4 : 103 | 0.50

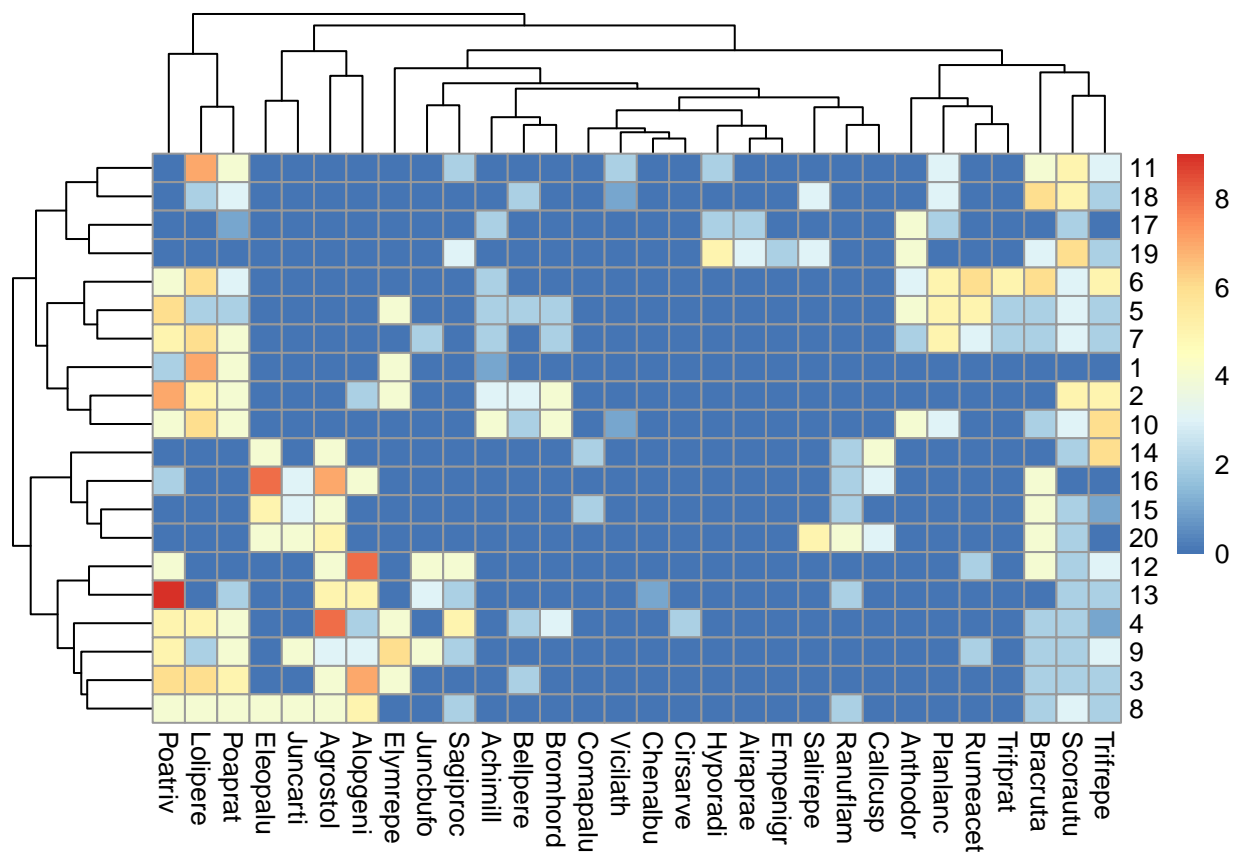
Average silhouette width : 0.49

I believe 4 gives the best silhouette index.

c.) Not sure about this

5.2

```
data(dune)
pheatmap(dune)
```



5.3

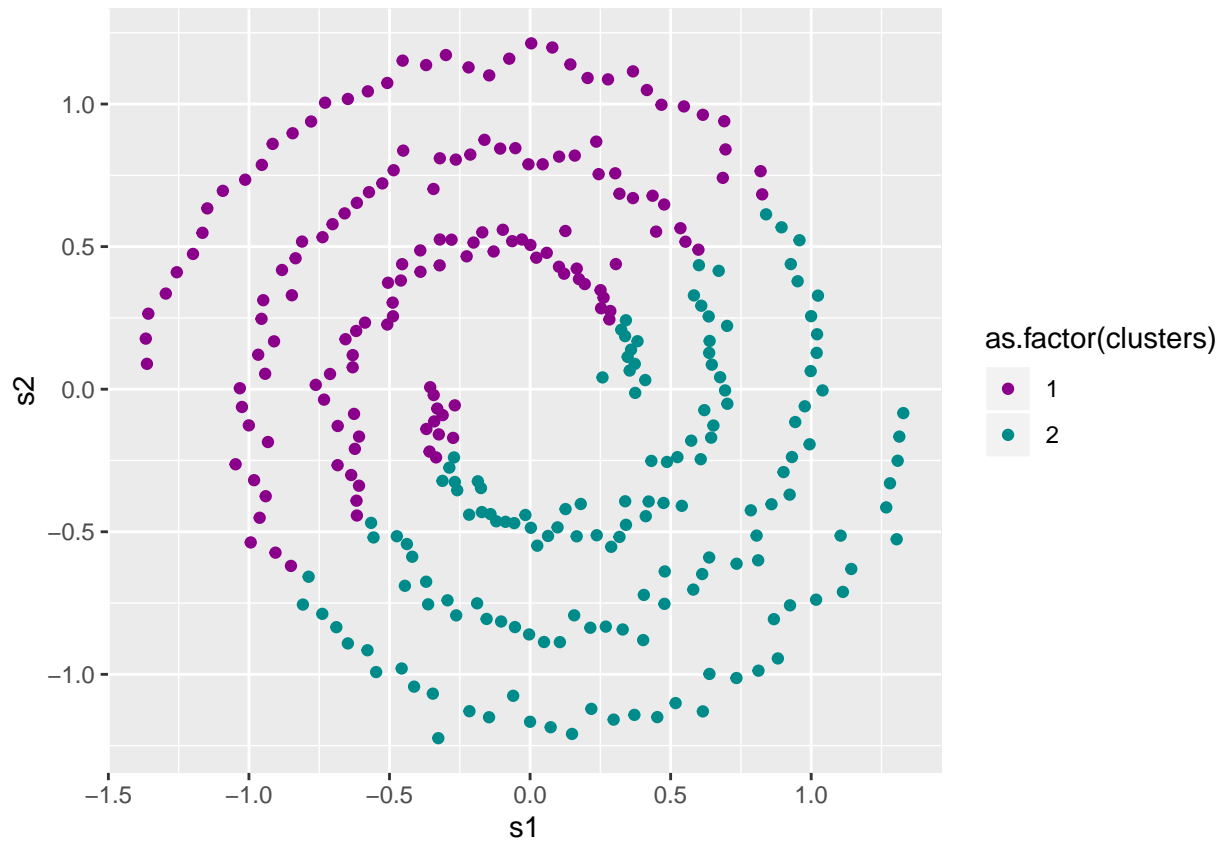
```
library(kernlab)
library(tidyverse)
data(spirals)

kmeans = kmeans(spirals,2)

clusters = kmeans$cluster

data = tibble(s1=spirals[,1],s2=spirals[,2],clusters)

data %>% ggplot(aes(x=s1,y=s2,color=as.factor(clusters))) +
  geom_point() +
  scale_color_manual(values=c("darkmagenta","cyan4"))
```



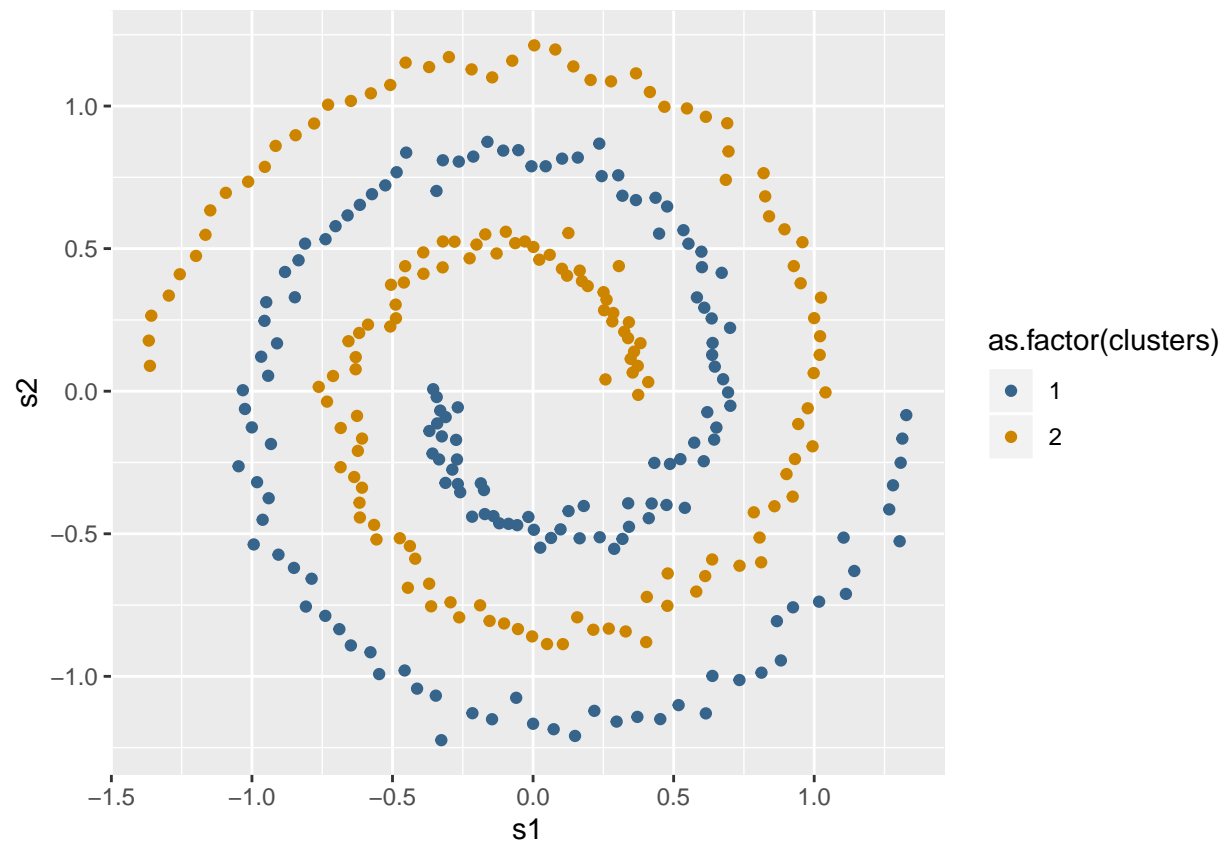
Let's try specc

```
spec = specc(spirals,2)

clusters = spec@.Data

data = tibble(s1=spirals[,1],s2=spirals[,2],clusters)

data %>% ggplot(aes(x=s1,y=s2,color=as.factor(clusters))) +
  geom_point() +
  scale_color_manual(values=c("steelblue4","orange3"))
```

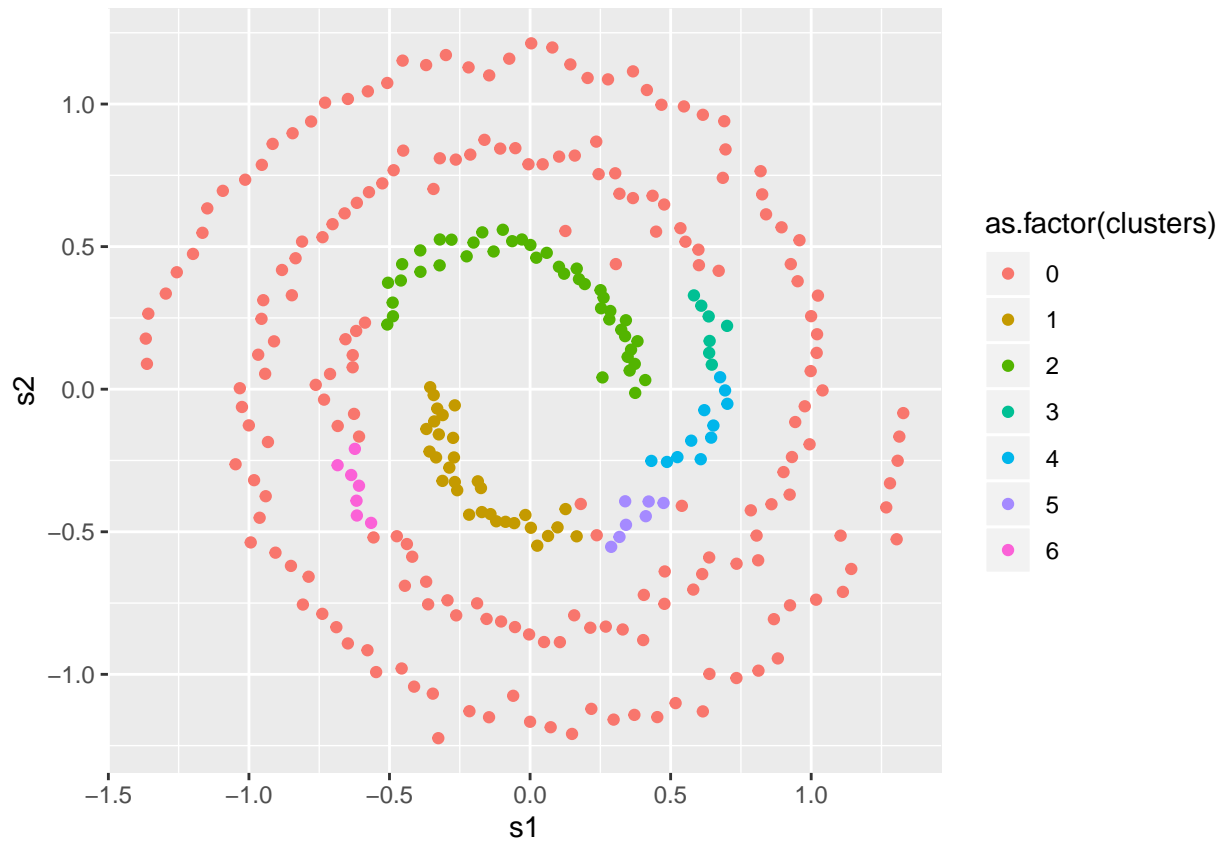


```
library(dbscan)
db = dbscan(spirals,.1)

clusters = db$cluster

data = tibble(s1=spirals[,1],s2=spirals[,2],clusters)

data %>% ggplot(aes(x=s1,y=s2,color=as.factor(clusters))) +
  geom_point()
```



5.4

Not sure why this data is clustered.

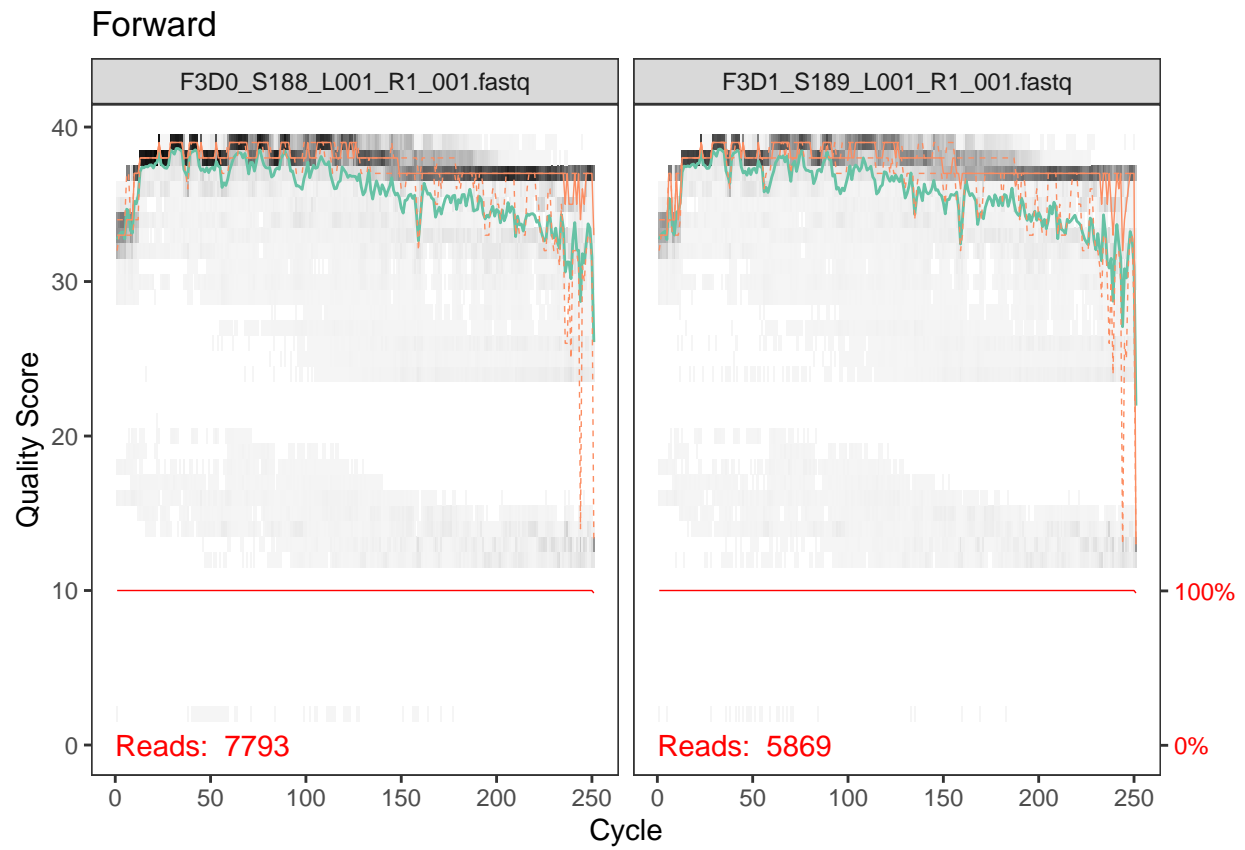
5.5

```
library(dada2)

base_dir = "../data"
miseq_path = file.path(base_dir, "MiSeq_SOP")
filt_path = file.path(miseq_path, "filtered")
fnFs = sort(list.files(miseq_path, pattern="_R1_001.fastq"))
fnRs = sort(list.files(miseq_path, pattern="_R2_001.fastq"))
sampleNames = sapply(strsplit(fnFs, "_"), `[`, 1)
if (!file_test("-d", filt_path)) dir.create(filt_path)
filtFs = file.path(filt_path, paste0(sampleNames, "_F_filt.fastq.gz"))
filtRs = file.path(filt_path, paste0(sampleNames, "_R_filt.fastq.gz"))
fnFs = file.path(miseq_path, fnFs)
fnRs = file.path(miseq_path, fnRs)
print(length(fnFs))
```

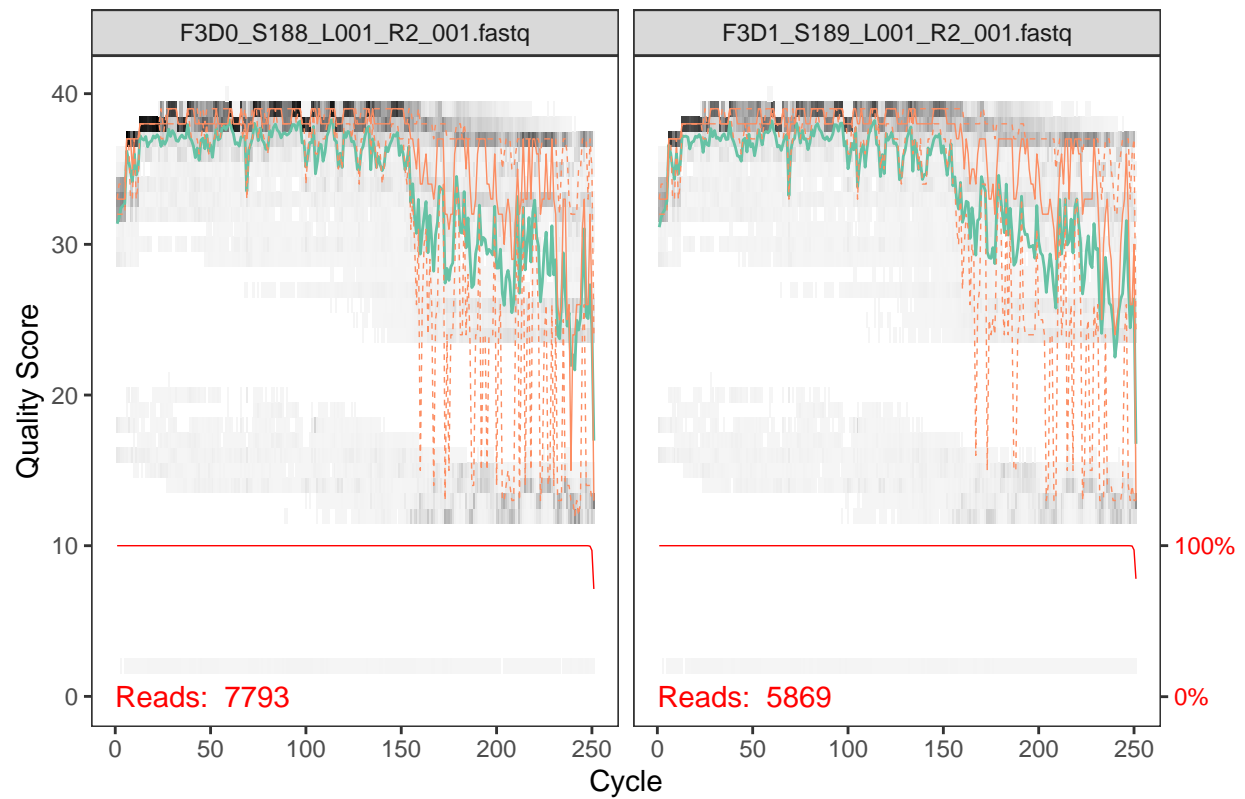
```
## [1] 20
```

```
plotQualityProfile(fnFs[1:2]) + ggtitle("Forward")
```



```
plotQualityProfile(fnRs[1:2]) + ggtitle("Reverse")
```

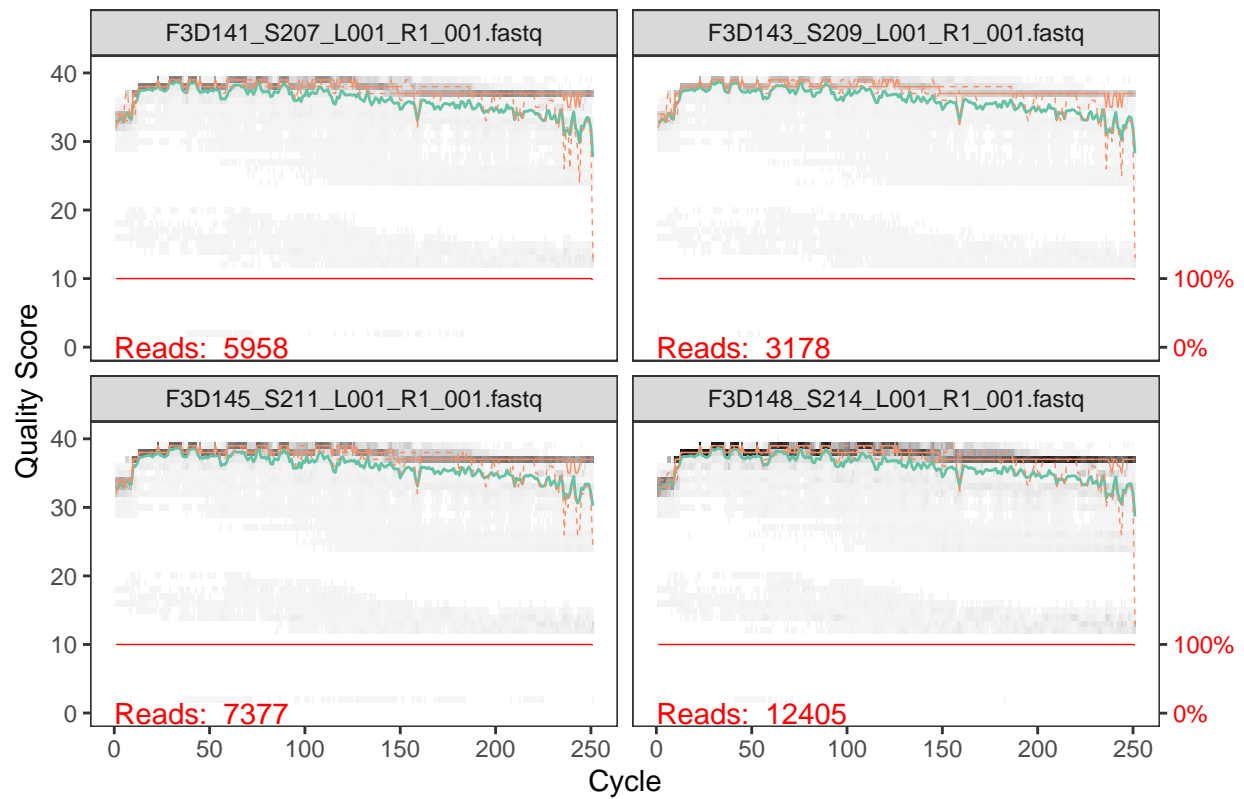
Reverse



5.6

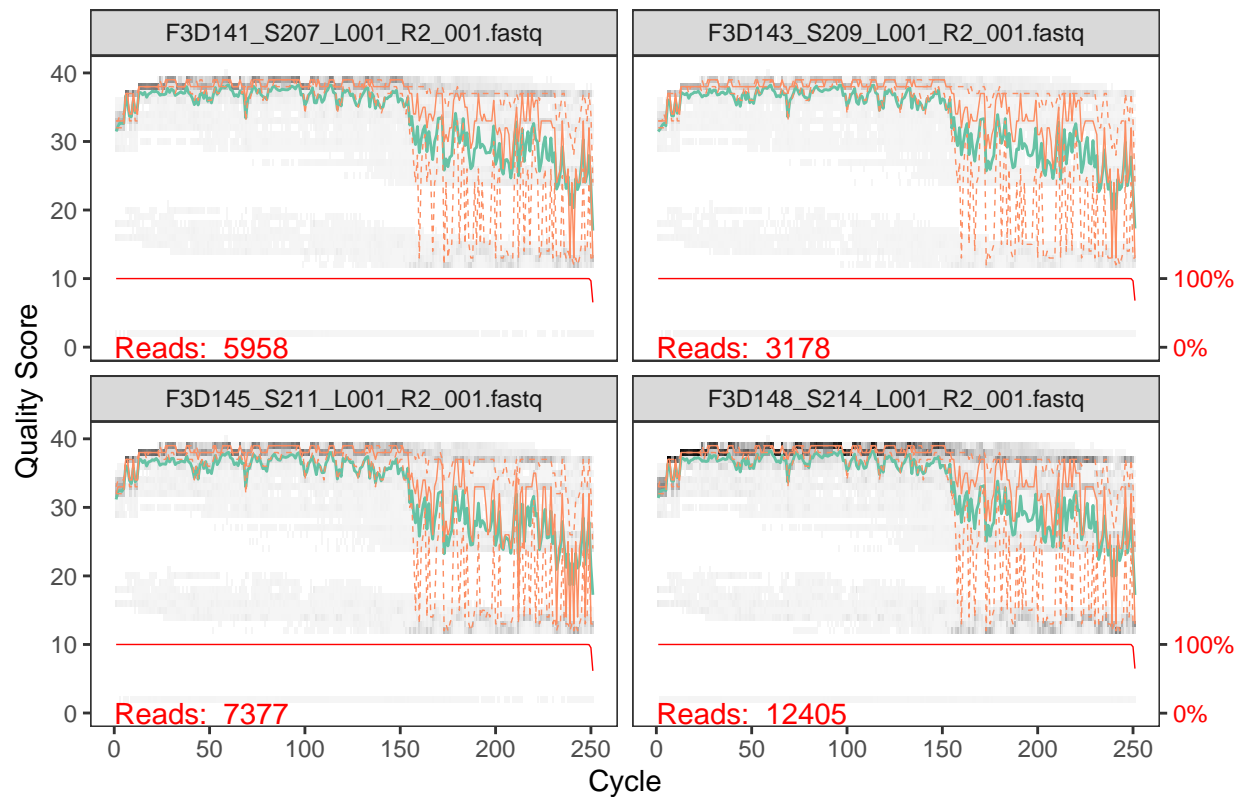
```
plotQualityProfile(fnFs[c(3,5,7,10)]) + ggtitle("Forward")
```


Forward



```
plotQualityProfile(fnRs[c(3,5,7,10)]) + ggtitle("Reverse")
```

Reverse



Reverse plots have unstable quality scores after 150 cycles.

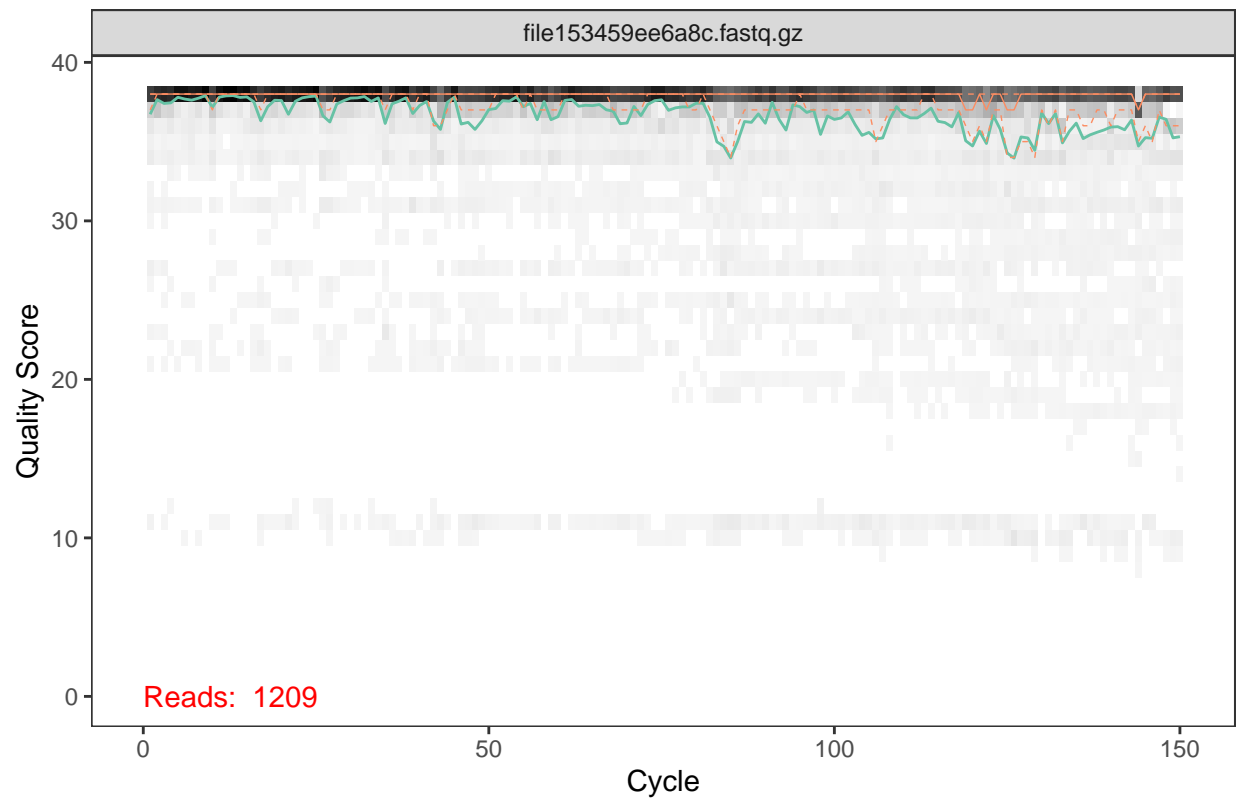
5.7

```
fnF1 <- system.file("extdata", "sam1F.fastq.gz", package="dada2")
fnR1 <- system.file("extdata", "sam1R.fastq.gz", package="dada2")
filtF1 <- tempfile(fileext=".fastq.gz")
filtR1 <- tempfile(fileext=".fastq.gz")

filterAndTrim(fwd=fnF1, filt=filtF1, rev=fnR1, filt.rev=filtR1,
              trimLeft=10, truncLen=160,
              maxN=0, maxEE=2,
              compress=TRUE, verbose=TRUE)

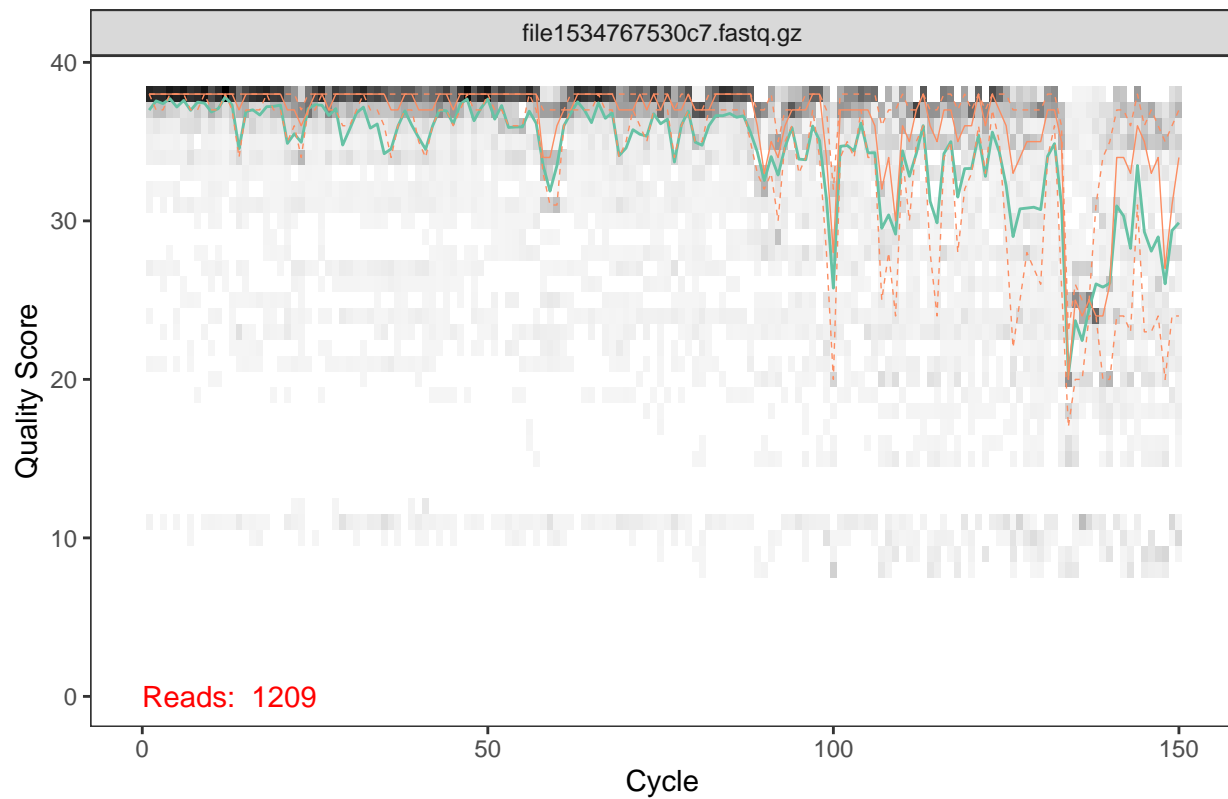
plotQualityProfile(filtF1) + ggtitle("Forward")
```

Forward



```
plotQualityProfile(filtR1) + ggtitle("Reverse")
```

Reverse



5.8

Cant seem to do this, since I need to put my credit card info into a google account.

```
# library(readxl)
# library(ggmap)
#
# bomb_raw = read_xlsx(here("Data", "sep_7_bomb.xlsx"))
#
# geocode(bomb_raw$Location)
```