

STAT646 - Homework 2

Daniel Osorio - dcosoriorh@tamu.edu
Department of Veterinary Integrative Biosciences
Texas A&M University

1. For the yeast data, do the following:

```
> yeastFiles <- list.files("yeastData/", full.names = TRUE)
> yeastInfo <- lapply(yeastFiles, function(file){
+   read.csv(file, sep = "\t", comment.char = "!")
+ })
> spots <- table(unlist(lapply(yeastInfo, function(X){X[,1]})))
> spots <- names(spots)[spots == length(yeastInfo)]
> yeastInfo <- lapply(yeastInfo, function(X){
+   X[X[,1] %in% spots, c("Ch1.Net..Mean.", "Ch2.Net..Mean.")]
+ })
> sNames <- gsub(".txt$", "", basename(yeastFiles))
> chanel1 <- sapply(yeastInfo, function(X){X[,1]})
> chanel2 <- sapply(yeastInfo, function(X){X[,2]})
> colnames(chanel1) <- colnames(chanel2) <- sNames
> rownames(chanel1) <- rownames(chanel2) <- spots
```

- (a) For each of average, complete, and single linkage, carry out hierarchical clustering on the samples using Euclidean distance.

```
> distanceMatrix <- dist(t(chanel1))
> avgHclust <- hclust(distanceMatrix, method = "average")
> sglHclust <- hclust(distanceMatrix, method = "single")
> cptHclust <- hclust(distanceMatrix, method = "complete")
```

- i. Which two samples are merged first, and what is the distance between them?

```
> c(sNames[abs(avgHclust$merge[1,])], round(avgHclust$height[1],2))
[1] "69971" "69972" "46589.94"
> c(sNames[abs(sglHclust$merge[1,])], round(sglHclust$height[1],2))
[1] "69971" "69972" "46589.94"
> c(sNames[abs(cptHclust$merge[1,])], round(cptHclust$height[1],2))
[1] "69971" "69972" "46589.94"
```

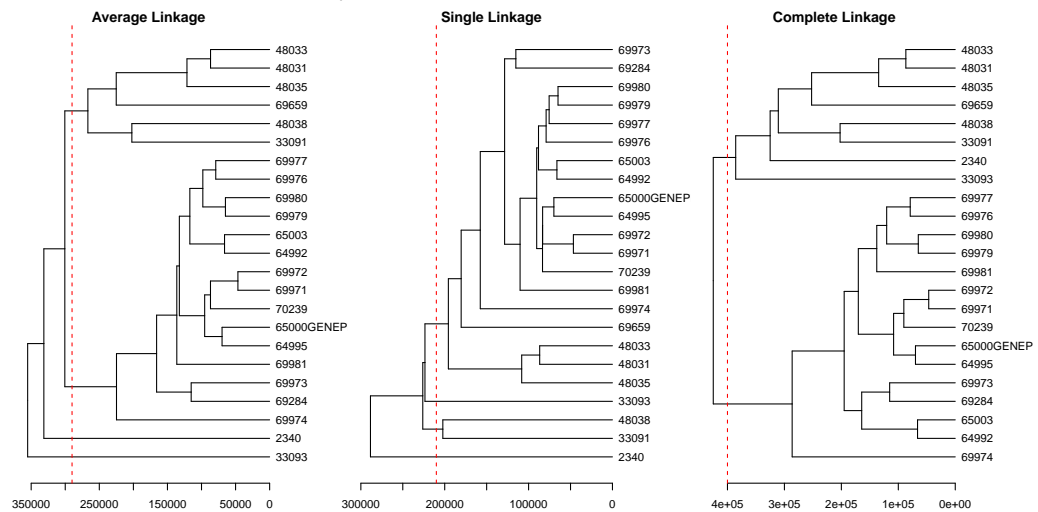
- ii. Report a height plot showing merge distances. How many clusters would you say there are?

```
> par(mfrow=c(1,3), mar=c(3,1,1,5))
> plot(as.dendrogram(avgHclust), horiz = TRUE,
+      main = "Average Linkage")
> abline(v = 2.9e5, lty = 2, col= "red")
> plot(as.dendrogram(sglHclust), horiz = TRUE,
```

```

+     main = "Single Linkage")
> abline(v = 2.1e5, lty = 2, col= "red")
> plot(as.dendrogram(cptHclust), horiz = TRUE,
+     main = "Complete Linkage")
> abline(v = 4e5, lty = 2, col= "red")

```

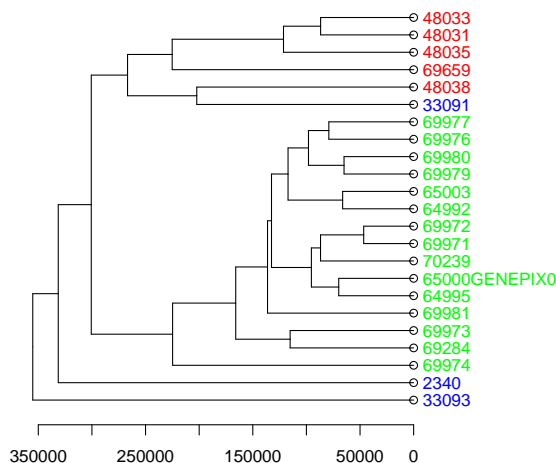


- iii. Interpret the clusters that you found with respect to the “meta” data for this dataset
- (b) Carry out K-means clustering, with $K = 3$. How do the results compare to the hierarchical clustering results you obtained using average linkage?

```

> avgHclust <- as.dendrogram(avgHclust)
> set.seed(11)
> K <- kmeans(distanceMatrix, centers = 3)
> labelCol <- function(x) {
+   if (is.leaf(x)) {
+     label <- attr(x, "label")
+     K1 <- names(K$cluster)[K$cluster == 1]
+     K2 <- names(K$cluster)[K$cluster == 2]
+     K3 <- names(K$cluster)[K$cluster == 3]
+     attr(x, "nodePar") <-
+       list(lab.col=if(label %in% K1){
+         "red"} else {if (label %in% K2){"blue"} else{"green"}})})
+   return(x)
+ }
> avgHclust <- dendrapply(avgHclust, labelCol)
> nodePar <- list(pch = c(NA,NA))
> par(mar=c(3,15,1,15))
> plot(avgHclust, horiz = TRUE, nodePar = nodePar)

```



2. The cardiothoracic data, in file ‘GDS4308.soft’, are described here:

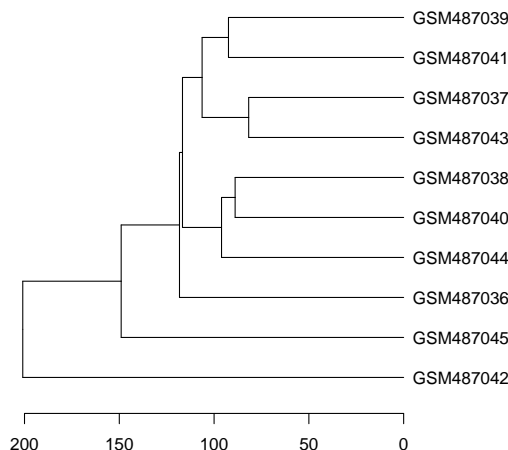
<https://www.ncbi.nlm.nih.gov/sites/GDSbrowser?acc=GDS4308>

Note that the first 97 lines of the data file consist of meta data, so the gene intensities start on line 98. Note also that the last line of the data file contains a table description and should not be included with the gene intensities. In what follows, work with the log₂-transformed intensities. Inspect the meta data in “GDS4308.soft” (you can just open it in a text editor) to figure out what the column names correspond to. Note that these are paired data.

```
> GSE19533 <- read.csv(file = "GDS4308_full.soft",
+                       sep = "\t", skip = 122,
+                       comment.char = "!")
> GSE19533 <- log2(GSE19533[,3:12])
```

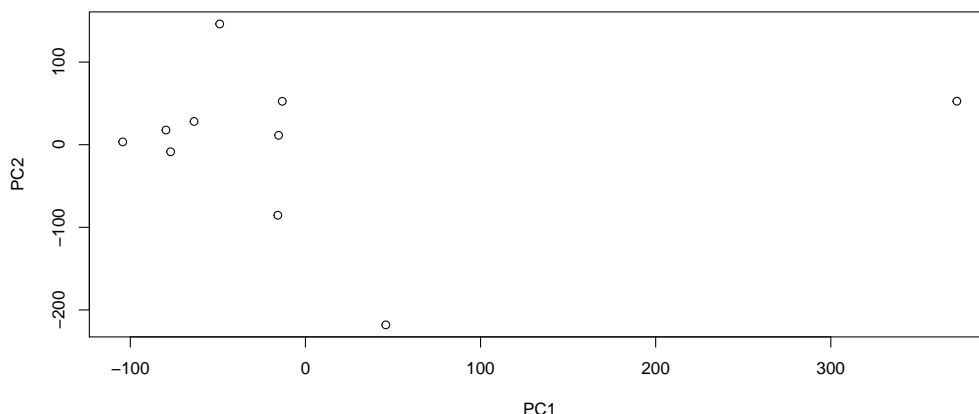
- (a) Carry out hierarchical clustering of the individuals, using Euclidean distance and complete linkage. Interpret the results.

```
> dGSE19533 <- dist(t(GSE19533))
> hc <- hclust(dGSE19533, method = "complete")
> par(mar=c(3,15,1,15))
> plot(as.dendrogram(hc), horiz = TRUE)
```



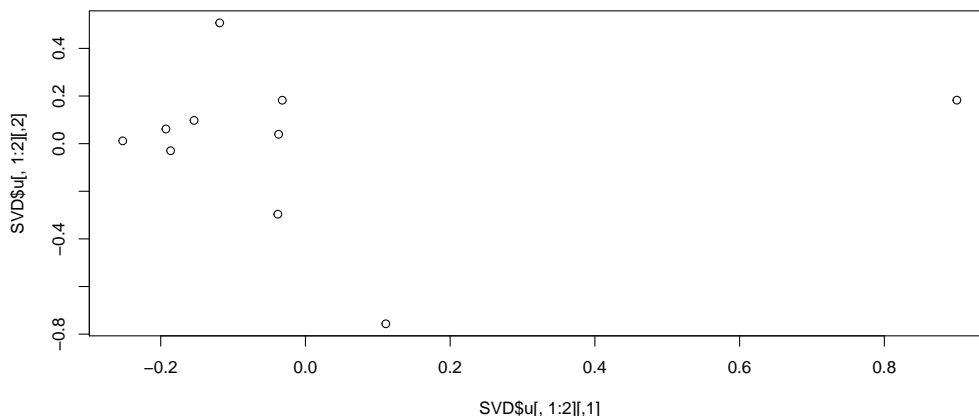
- (b) Carry out a principal component analysis of the intensities, treating individuals (columns) as variables and features (rows) as replicates. Interpret the results.

```
> cGSE19533 <- t(scale(t(GSE19533), center = TRUE, scale = TRUE))
> PC <- prcomp(t(cGSE19533))
> plot(PC$x[,1:2])
```



- (c) Carry out singular value decomposition of the column-centered intensities. Interpret the results.

```
> SVD <- svd(t(cGSE19533))
> plot(SVD$u[,1:2])
```



- (d) Compute one-sample t-statistics for each gene to search for features for which the mean paired difference is not 0. Use the bootstrap to obtain p-values. Use `p.adjust` to translate the p-values to FDR estimates (specify argument `method = 'fdr'`). How many features are significant at an estimated FDR of 0.05? I will give tips on R code for this problem in class and Q&A.

```
> pre <- paste0("GSM4870",seq(37,45,2))
> pos <- paste0("GSM4870",seq(36,44,2))
> diff <- GSE19533[,pos] - GSE19533[,pre]
> tStat <- apply(diff,1,function(X){t.test(X)$statistic})
> pre <- t(scale(t(GSE19533[,pre]), center = TRUE, scale = FALSE))
> pos <- t(scale(t(GSE19533[,pos]), center = TRUE, scale = FALSE))
>
```

```
> # boot <- sapply(1:100, function(b){
> #   print(b)
> #   sColumns <- sample(1:5, replace = TRUE)
> #   if(length(unique(sColumns)) > 1){
> #     dColumns <- pos[,sColumns] - pre[,sColumns]
> #     apply(dColumns,1,function(X){t.test(X)$statistic})
> #   }
> # })
> #
> # pValues <- rowMeans(boot > tStat)
> # sum(p.adjust(pValues, method = "fdr") < 0.05)
```