

Cluster Analysis

In this lab, you will be working with an R data set that was run on the Affymetrix human HGU95Av2 array. The microarray data are from genomic primary fibroblast cell lines and were generated for 46 samples: 23 human (*Homo sapien*), 11 bonobo (*Pan paniscus*), and 12 gorilla (*Gorilla gorilla*) donors. This is a publicly available dataset within the 'fibroEset' package in R. It should be noted that two identical human donor arrays are in this dataset. This data set is good for clustering and classification problems since there is a large difference in transcript profiles between all 3 species.

The analysis that you will conduct is based on clustering methods. The first problems require hierarchical clustering, while the last problems use spectral k -means clustering. We denote this as 'spectral' because instead of using the genes/probes as input into the clustering algorithm like the hierarchical clustering method, some form of spectral decomposition (e.g. PCA) is first computed and these eigenfunctions are used in the clustering algorithm. This method can be more useful than using the genes/probes in some cases where the variability is best summarized in a few components (or eigenfunctions).

- 1.) Load the fibroEset library and data set. Obtain the classifications for the samples.

```
> if (!requireNamespace("BiocManager", quietly = TRUE))
  install.packages("BiocManager")
> BiocManager::install("fibroEset")

> library(fibroEset)
> data(fibroEset)

> dat<-exprs(fibroEset)
> anno<- fibroEset@phenoData@data[["species"]]
> anno
[1] b b b b b b b b b b g g g g g g g g g g h h h h h h h h h h h h h h h h h h
h h h
Levels: b g h
```

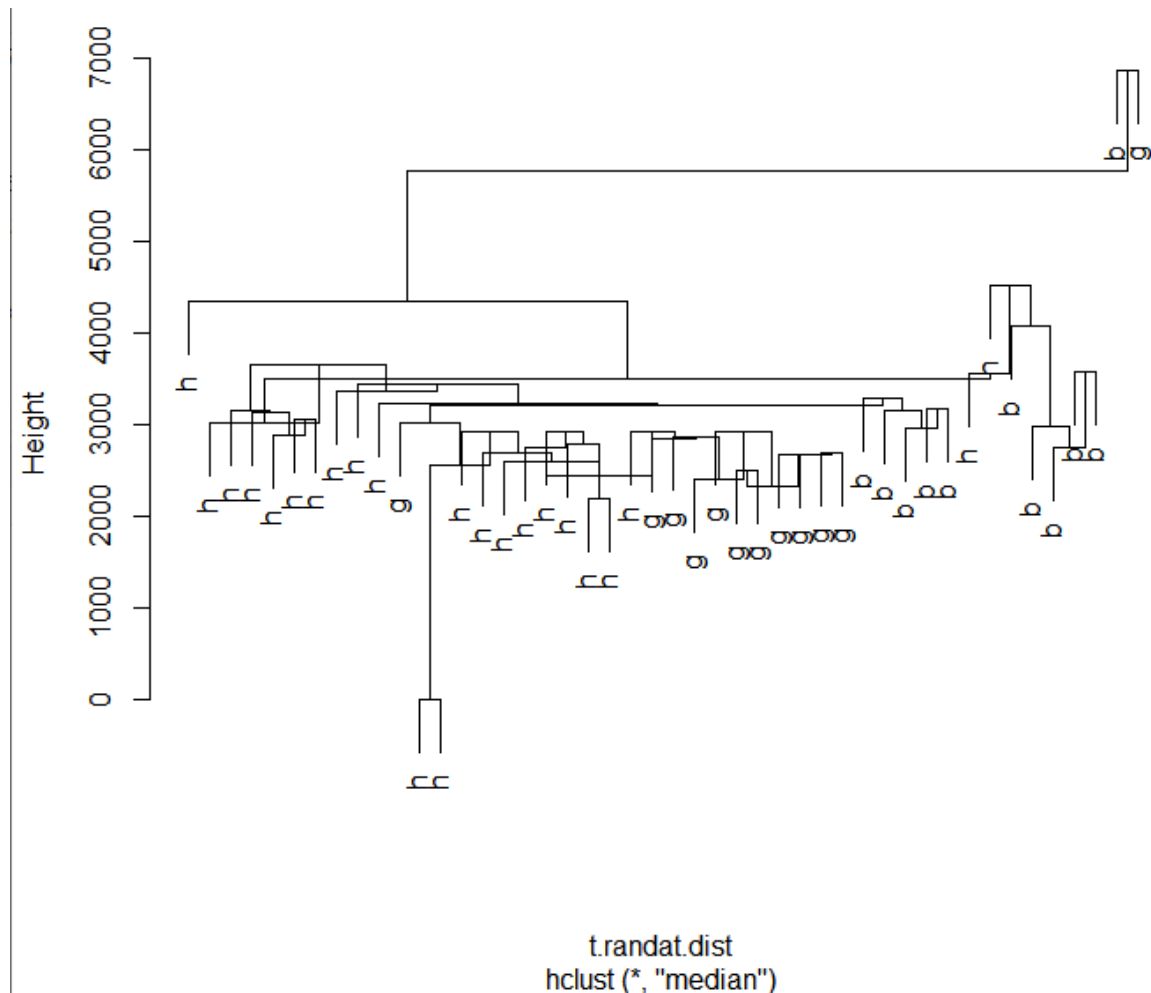
- 2.) Select a random set of 50 genes from the data frame, and subset the data frame.

```
> colnames(dat)<-anno
> set.seed(46)
> randat<-dat[sample(nrow(dat), 50),]
```

- 3.) Run and plot hierarchical clustering of the samples using manhattan distance metric and median linkage method. Make sure that the sample classification labels are along the x-axis. Title the plot.

```
> randat.dist<-dist(randat, method="manhattan")
> dat.hca<-hclust(randat.dist, method="median")
> plot(dat.hca, main="fibroEset Cluster Dendotram")
#did not give sample classification on x axis, so will transpose data.

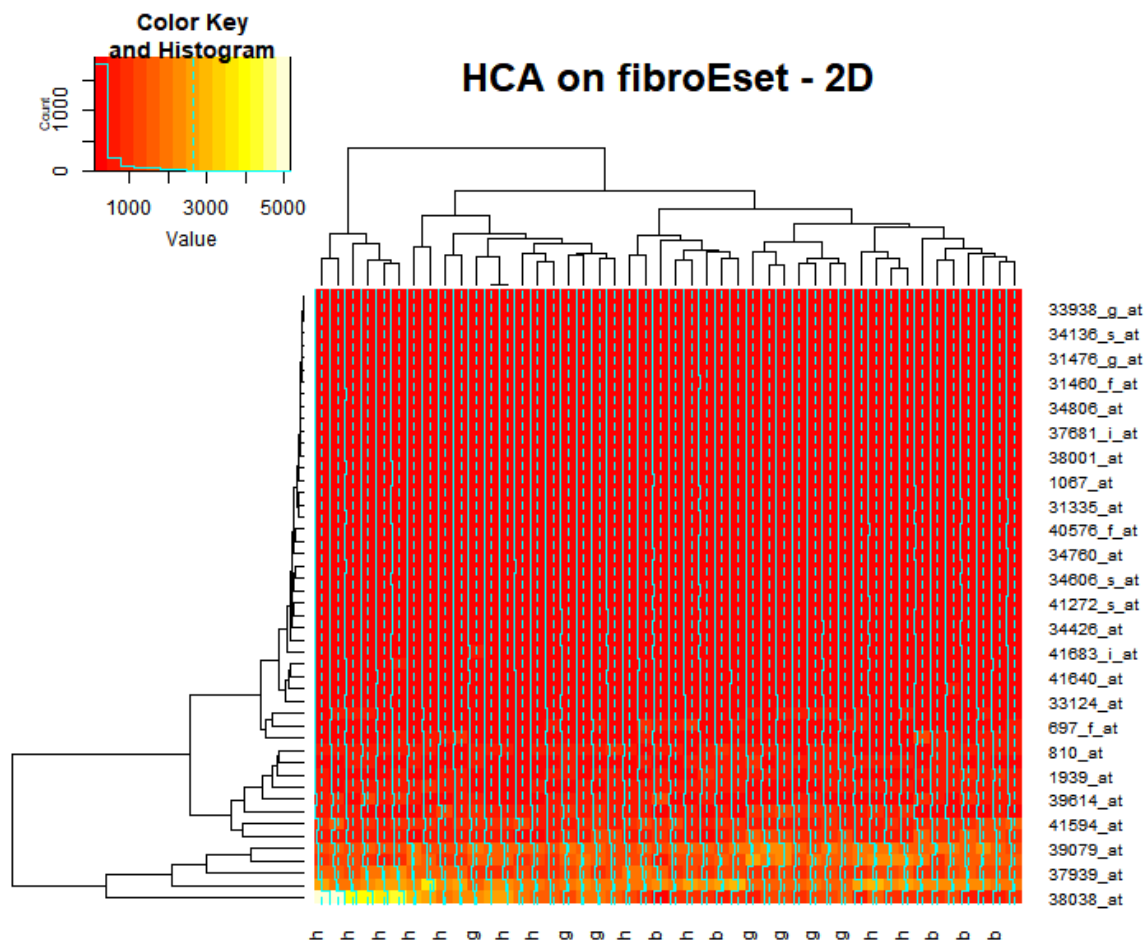
> t.randat<-t(randat)
> t.randat.dist<-dist(t.randat, method="manhattan")
> t.dat.hca<-hclust(t.randat.dist, method="median")
> plot(t.dat.hca, main="fibroEset Cluster Dendotram")
```



- 4.) Now both run hierarchical clustering and plot the results in two dimensions (on samples and genes). Plot a heatmap with the genes on the y-axis and samples on

the x-axis. Once again, make sure that the sample and genes labels are present.
Title the plot. -could use manhattan and median again

```
> heatmap.2(randat, main="HCA on fibroEset - 2D")
```



5.) Calculate PCA on the samples and retain the first two components vectors (eigenfunctions). Calculate k -means clustering on these first two components with $k=3$.

```
> pc.randat<-prcomp(randat)
> pc2<-pc.randat[["rotation"]][,1:2]
> cl<-kmeans(pc2, centers=3)
```

K-means clustering with 3 clusters of sizes 6, 26, 14

Cluster means:

	PC1	PC2
1	0.1907579	0.2946022
2	0.1333776	-0.1212075
3	0.1478109	0.0349093