Final project for the course HY350.90
**Topic 1: Gene Expression Analysis from microarray data**

**Task 1**
**1.** Find and download the dataset  entitled
"Chronic high-level alcohol consumption effect on brain: post-mortem hippocampus".
Link: http://www.ncbi.nlm.nih.gov/sites/GDSbrowser?acc=GDS4879



**2.** Clean the dataset from the unecessary lines in the beginning and

the end.
These lines start with either of the symbols: # ! ^
Open file with a notepad and search for the characters :

**3.** Read in the dataset in R.
   Seperate with tab character '\t'
   Set the NA string (na.strings)  as either "NA" or "null"

```
dataset <-read.table
(file='C:\\Users\\avon\\Desktop\\hy390\\GDS4879.soft',  sep = '\t',
na.strings = c("NA" , "null" ), header=TRUE,
stringsAsFactors = FALSE)
```

```
R Console

>
>
>
> dataset <- read.table(file='C:\\Users\\avon\\Desktop\\hy390\\GDS4879.soft', sep = '\t', na.strings = c("NA" , "null" ), header=TRUE, stringsAsFactors = FALSE )
> |
```

**4**. Choose only a subset of lines for the final dataset : all lines from 1 - 10000
   that contain no "NA" or "null".
   Remove all NA's with complete cases and then take the first 10000 lines.

```
newdata <- dataset[ complete.cases( dataset ) ,]
newdata <- newdata[ 1:10000,]
```

3

**5.** Save the second column that contains information about the gene names
in a variable called gene.names .

```
gene.names <- newdata[,2:2]
```



4

**6.** Take The final dataset that   consist only of a subset of 24 sample that we need

```
#delete 2 first columns
newdata <- newdata[c(-1,-2)]

#delete unnecessary columns by name
finaldata <- subset(newdata, select =
-c(GSM1085680,GSM1085680,GSM1085682,GSM1085683,GSM1085684,GSM1085687,
GSM1085691,GSM1085697,GSM1085700,GSM1085686,GSM1085688,GSM1085690,GSM
1085692,GSM1085693,GSM1085702,GSM1085703) )

#view data and columns that left
View(finaldata)
colnames(finaldata)
```

**Task 2**

**1.** Make a boxplot of the samples
First convert finaldata to numeric.

```
dataNum <- matrix(data = NA, nrow = dim(finaldata)[1], ncol =
dim(finaldata)[2])

for (i in 1:dim(finaldata)[2])
{
    dataNum[,i] <- c(as.numeric(finaldata[[i]]))
}

boxplot(dataNum , ylim= c(0,15))
```



As shown above the dataset is normalized because one or more arrays don't have intensity levels which are drastically different from the rest of the arrays, this may indicate a problem with these arrays. As normalization refers to the creation of shifted and scaled versions of statistics, where the intention is that these normalized values allow the comparison of corresponding normalized values for different datasets in a way that eliminates the effects of certain gross influences, as in an anomaly time series.

**Task 3**
**1.** Construct a heatmap of all samples

```
#---create table
nba <- finaldata
matrix <- data.matrix(nba)

#---create heatmap
heatmap <- heatmap(matrix)
```



Heatmap reflexes gene expression values . Heatmap has a graphical
representation of data , where the individual values of each sample
contained in a matrix , are represented as colors

**2.** Construct a heatmap of all samples using the function heatmap.2 from the library gplots.
Use scale = 'none' and trace = 'none'.

First Install gplots library :

```
if (!require("gplots")) {
    install.packages("gplots", dependencies = TRUE)
    library(gplots)
    }
```

Then create 2<sup>nd</sup> heatmap  using the function heatmap.2 :

```
heatmap.2 <- heatmap.2(matrix, scale="none",  trace = "none")
```



The differences betwween Heatmap with method 1 heatmap with method 2
is that 2<sup>nd</sup> heatmap's  rows and columns are reordered  according to
some set of values within the restrictions imposed by the dendrogram
is carried out.
2<sup>Nd</sup> heatmap has also a histogram  added to the top  left side.

**Task 4**

```
finaldatanew <- t(matrix)

temp <- prcomp(finaldatanew)
colors <-as.factor
(c(rep('yellow',6),rep('green',6),rep('blue',6),rep('red',6)))

#plotno1
plot(temp$x[,1], temp$x[,2], xlab = 'pc1', ylab = 'pc2')

#plotno2
plot(temp$x[,1], temp$x[,2], xlab = 'pc1', ylab = 'pc2',col=colors)
```

4. Plot the data using the Principal Component 1 and Principal
Component 2.

5. Try to color the data points of the plot according to the class of each sample. e.g. use red for female-alcoholic, green for female non-alcoholic, etc.

**Task 5**

```
#building the outcomes
femalesVmales <- as.factor(c(rep('females', 12),rep('males',12)))

#applying on all row genes
myTtest <- function(x,y){
  levs <- unique(y);
  a <- x[y == levs[1]]
  b <- x[y == levs[2]]
  res <- t.test(a, b, var.equal = TRUE)
  res$p.value
}

#save the results p-values in a vector
pvalues <- apply(matrix, 1, myTtest, femalesVmales)
View(pvalues)
```

```
>
> #building the outcomes
> femalesVmales <- as.factor(c(rep('females', 12),rep('males',12)))
>
>
> #applying on all row genes
> myTtest <- function(x,y){
+    levs <- unique(y);
+    a <- x[y == levs[1]]
+    b <- x[y == levs[2]]
+    res <- t.test(a, b, var.equal = TRUE)
+    res$p.value
+ }
> #save the results (p-values) in a vector
> pvalues <- apply(matrix, 1, myTtest, femalesVmales)
> View(pvalues)
> |
```
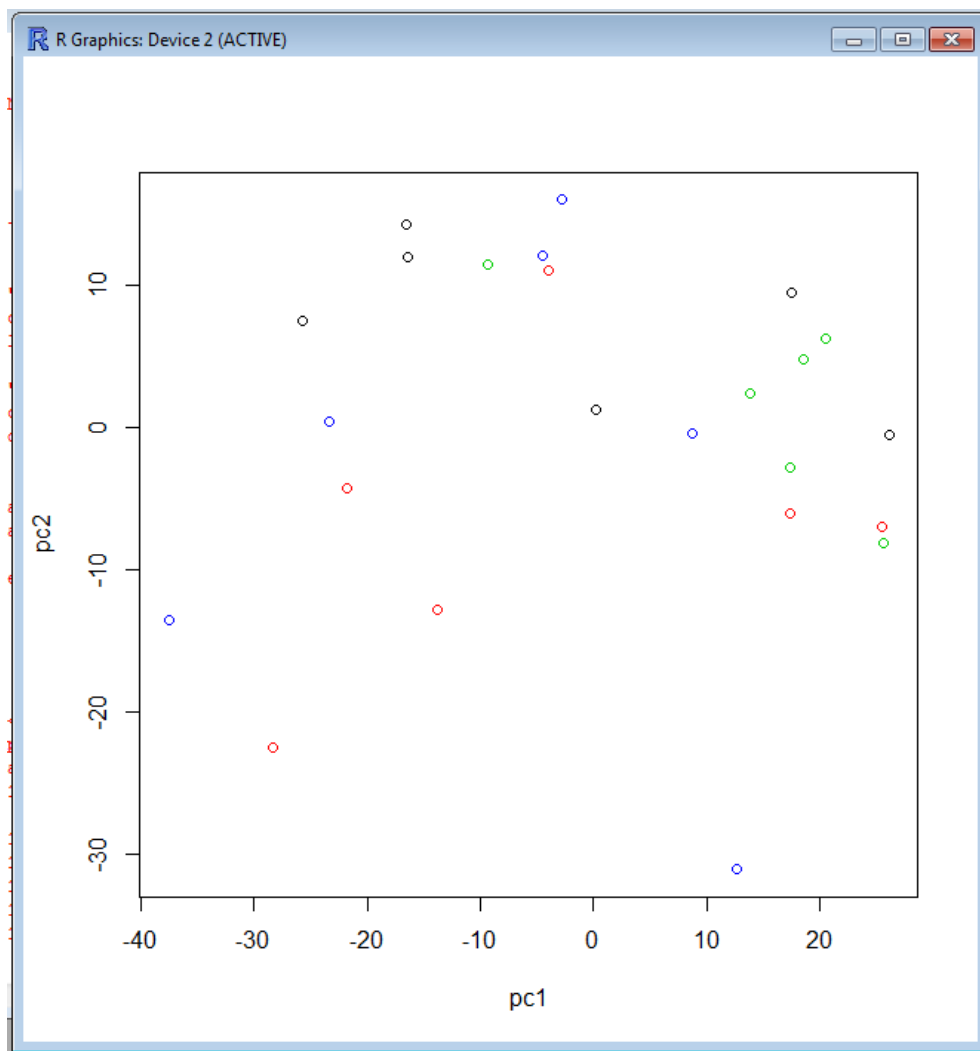
| R Data: pvalues | |
|---|---|
| | **x** |
| 1 | 0.3559557 |
| 2 | 0.3323129 |
| 3 | 0.4436141 |
| 4 | 0.5515539 |
| 5 | 0.8425665 |
| 6 | 0.1381835 |
| 7 | 0.3229616 |
| 8 | 0.306034 |
| 9 | 0.1112976 |
| 10 | 0.5445897 |
| 11 | 0.8069778 |
| 12 | 0.1825728 |
| 13 | 0.526363 |

```
#2.  ascending order of pvalues_femaleVmale
ordered_femaleVmale <-
pvalues_femaleVmale[order(pvalues_femaleVmale)]

#take the first 100 minimum pvalues
list1=head(ordered_femaleVmale ,100)


for (i in 1:length(list1) ) {
     for (y in 1:length(pvalues_femaleVmale) ) {
          if( pvalues_femaleVmale[y] == list1[i])
               write.table (matrix[y] ,
     file='C:\\Users\\avon\\Desktop\\hy390\\tmp.txt' ,
               append=TRUE,sep = "\t")
     }
}
```

## Task 6

```
#1.  applying on all row genes
myTtest <- function(x,y){
  levs <- unique(y);
  a <- x[y == levs[1]]
  b <- x[y == levs[2]]
  res <- t.test(a, b, var.equal = TRUE)
  res$p.value
}

#2.  save the results (p-values) of nonalcoholics and alcoholics in a
vector
pvalues_nonalcoholicsValcoholics <- apply(matrix, 1, myTtest,
nonalcoholicsValcoholics)
#View(pvalues_nonalcoholicsValcoholics)
```

```
#ascending order of pvalues_nonalcoholicsValcoholic
ordered_nonalcoholicsValcoholics <-
pvalues_nonalcoholicsValcoholics[order(pvalues_nonalcoholicsValcoholi
cs)]

#take the first 100 minimum pvalues
list2=head(ordered_nonalcoholicsValcoholics ,100)

for (i in 1:length(list2) ) {
    for (y in 1:length(pvalues_nonalcoholicsValcoholics) ) {
        if( pvalues_nonalcoholicsValcoholics[y] == list2[i])


write.table(matrix[y],file='C:\\Users\\avon\\Desktop\\hy390\\tmp2.txt
',append=TRUE,sep = "\t")
    }
}
```