

Lab #2

Data visualization

In this lab, we are going to work with a very popular time course cDNA data set from Paul Spellman's lab at Stanford. This microarray was designed with probes from the yeast *Saccharomyces cerevisiae* genome. The data set includes 3 different experiments, each with its own time course (each array is a different time point) for measuring transcript levels that are induced by various cyclins. The transcripts that respond to this stimulus are seen to be regulated at the different stages of the cell cycle. The 3 experiments differ by the method that the yeast cultures were synchronized: α factor arrest, elutriation, and arrest of a *cdc15* temperature-sensitive mutant. The *cdc15* time course experiment is the one that we will use in this lab to conduct some simple mathematical manipulations and plots.

The paper, "Comprehensive Identification of Cell Cycle-regulated Genes of the Yeast *Saccharomyces cerevisiae* by Microarray Hybridization" is available on the course website.

- 1.) Go to class website under Course Documents > Data Sets and download the Spellman yeast cell cycle dataset zip file.
- 2.) Unzip the text file, and read into R (Hint: using the `read.table()` function with a "header=T" argument and "row.names=1" argument is one method to do this).
- 3.) Look at the dimensions of the data frame and make sure that there are 6,178 genes and 77 arrays/sample.

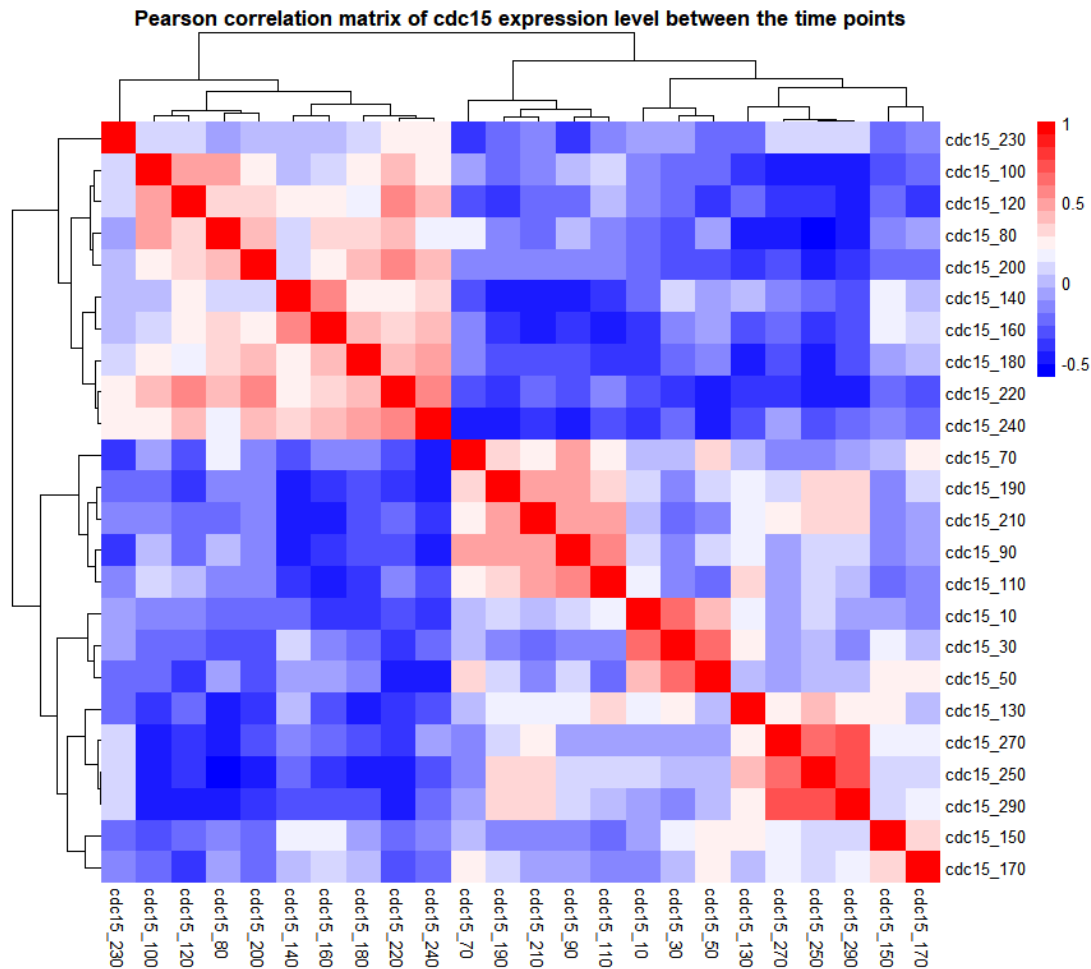
```
> yeastdata=data.frame( read.table("C:/Users/dell/Desktop/spellman.txt", header=T,row.name=1))
> dim(yeastdata)
[1] 6178  77
```

- 4.) Isolate only the *cdc15* experiment (samples 23-46).

```
> cdcyeast=yeastdata[grepl("cdc15",colnames(yeastdata))]
> dim(cdcyeast)
[1] 6178  24
```

- 5.) Now calculate a correlation matrix between the time points (use Pearson's correlation). Make sure to title the plot, label the axes, and provide a legend of the color gradient. In the correlation calculation, make sure to use the argument 'use' and value=`pairwise.complete.obs` since all of these arrays have at least one missing value.

```
> library(corrplot)
> library(heatmap)
> personmatrix=cor(cdcyeast, method = "pearson", use = "pairwise.complete.obs")
> col<- colorRampPalette(c("blue", "white", "red"))(20)
> p <- heatmap(personmatrix, col = col,
+             clustering_distance_rows = "correlation",
+             clustering_distance_cols = "correlation",
+             border = FALSE,
+             main = "Pearson correlation matrix of cdc15 expression level between the time points")
```



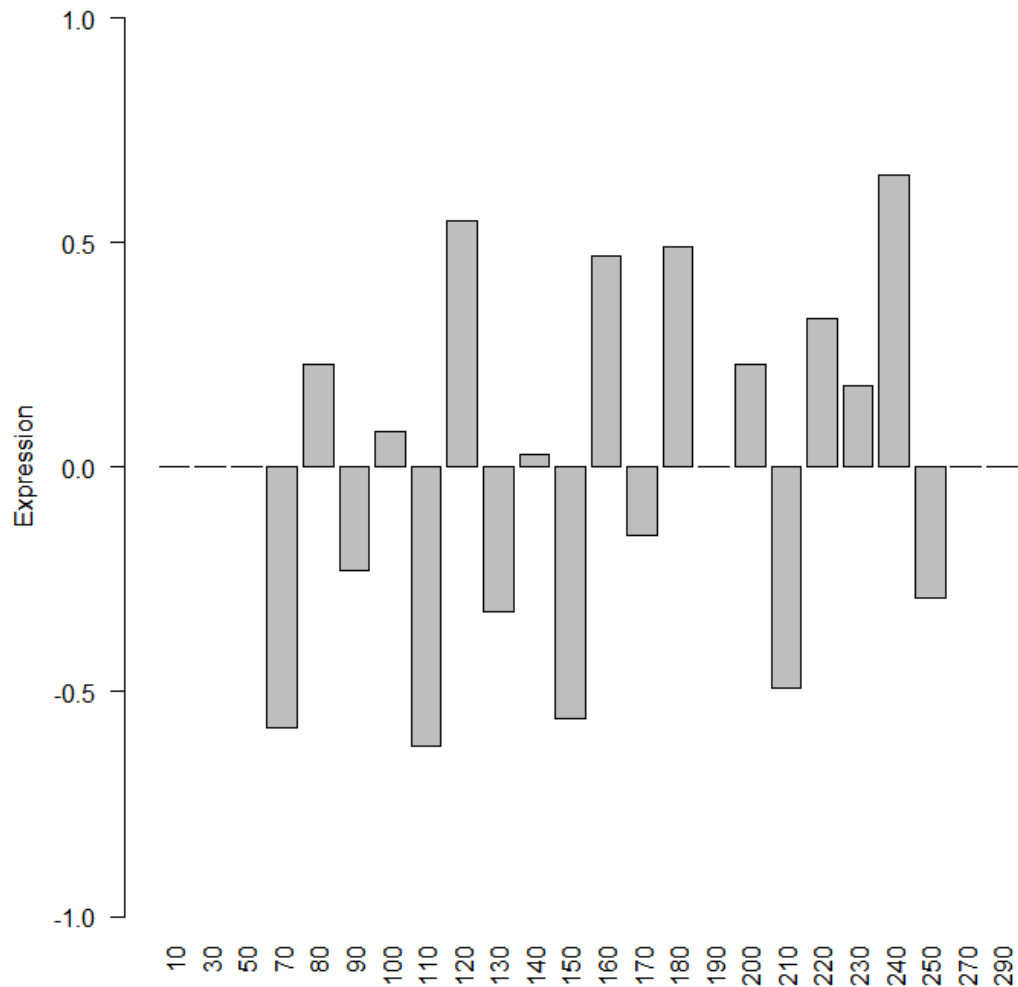
6.) Select the gene YAL002W, which is VPS8, a membrane-binding component of the CORVET complex; involved in endosomal vesicle tethering and fusion in the endosome to vacuole protein targeting pathway. Impute the missing values with the row mean (save as something). Make sure to cast the gene to numeric.

```
> data_YAL=as.numeric(cdcyeast['YAL002W',])
> data_YAL[is.na(data_YAL)]=mean(data_YAL,na.rm=T)
> data_YAL
[1] 3.857953e-18 3.857953e-18 3.857953e-18 -5.800000e-01 2.300000e-01 -2.300000e-01 8.000000e-02 -6.200000e-01
[9] 5.500000e-01 -3.200000e-01 3.000000e-02 -5.600000e-01 4.700000e-01 -1.500000e-01 4.900000e-01 3.857953e-18
[17] 2.300000e-01 -4.900000e-01 3.300000e-01 1.800000e-01 6.500000e-01 -2.900000e-01 3.857953e-18 3.857953e-18
```

7.) Generate a profile plot of the same gene. Title the plot, label the axes, and on the x-axis, provide the time points only for each array (no "cdc15_" prefix) so we can visualize the transcript pattern over time. Use lwd in the plot command (lwd=line width).

```
> barplot(data_YAL, ylim=c(-1, 1), lwd=1,
+         ylab='Expression', las = 2, cex.names = 1,
+         names.arg = c("10","30","50","70","80","90","100","110","120",
+                       "130","140","150","160","170","180","190","200",
+                       "210","220","230","240","250","270","290"),
+         main = 'Expression Levels for Gene YAL002W Across 3 Experiments Over Time')
```

Expression Levels for Gene YAL002W Across 3 Experiments Over Time



8.) Now let's create a simple shiny app which allows the user to select and correlate any time point verse another time point across all genes. To do this, we can create a server and ui function within the same file, paste both into the R session, then call them with:

```
>shinyApp(ui = ui, server = server)
```

Use the Iris dataset example from the lecture as a model. You can remove the kmeans clustering code and just focus on plotting the columns (time points) of the CDC15 data matrix against each other.

```
yeastdata=read.table("C:/Users/97481/Downloads/spellman.txt", header=T, row.names = 1)
```

```
cdcyeast=yeastdata[grepl("cdc15",colnames(yeastdata))]
```

```
ui<- fluidPage(
  sidebarLayout(
    sidebarPanel(selectInput('xcol', 'X Variable', dimnames(cdcyeast)[[2]]),
```

```

        selectInput('ycol','y Variable', dimnames(cdcyeast)[[2]]),
        selectInput('color','Point color',list("Red" = "Red", "Blue" = "Blue",
                                                "Green" = "#70AD47"),
selected = 1)

    ),
    mainPanel(plotOutput("plot1"))
  )

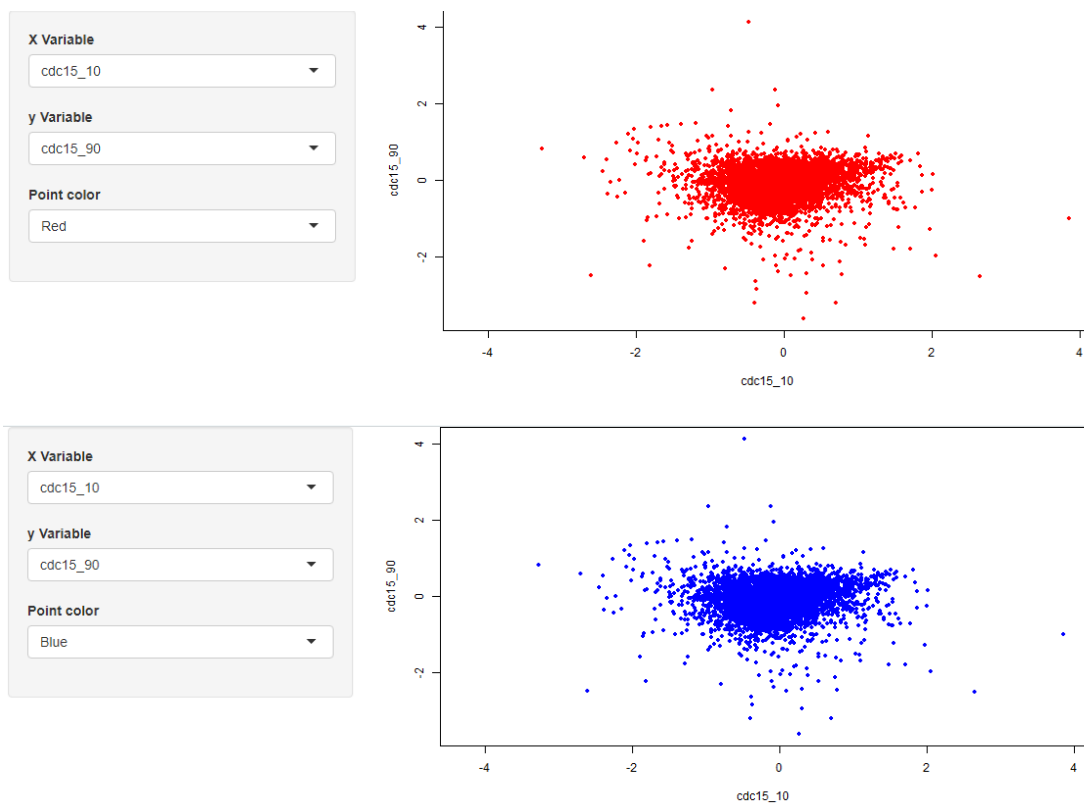
)
server <- function(input, output,session) {

  selectedData <- reactive({cdcyeast[, c(input$xcoll, input$ycol)]})

  output$plot1 <- renderPlot({ par(mar=c(5.1, 4.1, 0, 1))
    plot(selectedData(), col = input$color, pch = 20, cex
      = 1)
  })
}

```

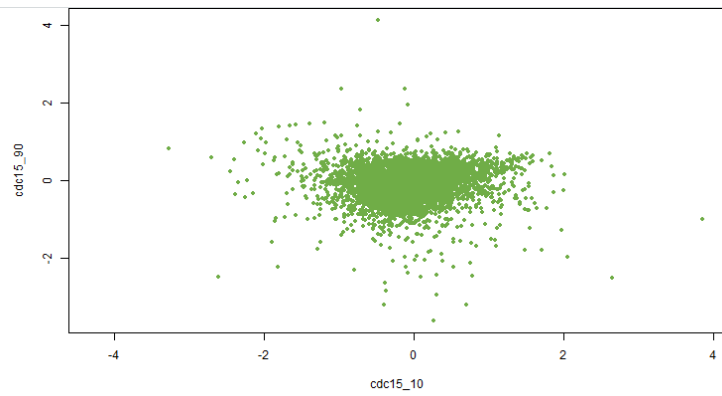
```
shinyApp(ui=ui,server = server)
```



X Variable
cdc15_10

y Variable
cdc15_90

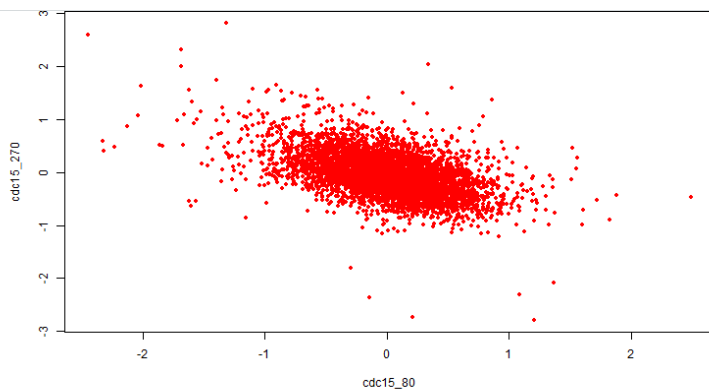
Point color
Green



X Variable
cdc15_80

y Variable
cdc15_270

Point color
Red



X Variable
cdc15_10

y Variable
cdc15_290

Point color
Red

