

Class15

Camryn McCann (PID: A15437387)

11/16/2021

```
#First lets call our programs

library(BiocManager)
library(DESeq2)

## Loading required package: S4Vectors

## Loading required package: stats4

## Loading required package: BiocGenerics

##
## Attaching package: 'BiocGenerics'

## The following objects are masked from 'package:stats':
##       IQR, mad, sd, var, xtabs

## The following objects are masked from 'package:base':
##       anyDuplicated, append, as.data.frame, basename, cbind, colnames,
##       dirname, do.call, duplicated, eval, evalq, Filter, Find, get, grep,
##       grepl, intersect, is.unsorted, lapply, Map, mapply, match, mget,
##       order, paste, pmax, pmax.int, pmin, pmin.int, Position, rank,
##       rbind, Reduce, rownames, sapply, setdiff, sort, table, tapply,
##       union, unique, unsplit, which.max, which.min

##
## Attaching package: 'S4Vectors'

## The following objects are masked from 'package:base':
##       expand.grid, I, unname

## Loading required package: IRanges

## Loading required package: GenomicRanges

## Loading required package: GenomeInfoDb
```

```

## Loading required package: SummarizedExperiment

## Loading required package: MatrixGenerics

## Loading required package: matrixStats

## 
## Attaching package: 'MatrixGenerics'

## The following objects are masked from 'package:matrixStats':
## 
##     colAlls, colAnyNAs, colAnyNs, colAvgsPerRowSet, colCollapse,
##     colCounts, colCummaxs, colCummins, colCumprods, colCumsums,
##     colDiffs, colIQRDiffs, colIQRs, colLogSumExps, colMadDiffs,
##     colMads, colMaxs, colMeans2, colMedians, colMins, colOrderStats,
##     colProds, colQuantiles, colRanges, colRanks, colSdDiffs, colSds,
##     colSums2, colTabulates, colVarDiffs, colVars, colWeightedMads,
##     colWeightedMeans, colWeightedMedians, colWeightedSds,
##     colWeightedVars, rowAlls, rowAnyNAs, rowAnyNs, rowAvgsPerColSet,
##     rowCollapse, rowCounts, rowCummaxs, rowCummins, rowCumprods,
##     rowCumsums, rowDiffs, rowIQRDiffs, rowIQRs, rowLogSumExps,
##     rowMadDiffs, rowMads, rowMaxs, rowMeans2, rowMedians, rowMins,
##     rowOrderStats, rowProds, rowQuantiles, rowRanges, rowRanks,
##     rowSdDiffs, rowSds, rowSums2, rowTabulates, rowVarDiffs, rowVars,
##     rowWeightedMads, rowWeightedMeans, rowWeightedMedians,
##     rowWeightedSds, rowWeightedVars

## Loading required package: Biobase

## Welcome to Bioconductor
## 
## Vignettes contain introductory material; view with
## 'browseVignettes()'. To cite Bioconductor, see
## 'citation("Biobase")', and for packages 'citation("pkgname")'.

## 
## Attaching package: 'Biobase'

## The following object is masked from 'package:MatrixGenerics':
## 
##     rowMedians

## The following objects are masked from 'package:matrixStats':
## 
##     anyMissing, rowMedians

```

#Now let's import our data. Then look at each.

```

counts <- read.csv("airway_scaledcounts.csv",row.names=1)
metadata <- read.csv("airway_metadata.csv")

```

```

head(counts)

##          SRR1039508 SRR1039509 SRR1039512 SRR1039513 SRR1039516
## ENSG000000000003      723        486       904       445      1170
## ENSG000000000005       0         0         0         0         0
## ENSG000000000419     467       523       616       371      582
## ENSG000000000457     347       258       364       237      318
## ENSG000000000460      96        81        73        66      118
## ENSG000000000938      0         0         1         0         2
##          SRR1039517 SRR1039520 SRR1039521
## ENSG000000000003    1097       806       604
## ENSG000000000005       0         0         0
## ENSG000000000419     781       417       509
## ENSG000000000457     447       330       324
## ENSG000000000460      94        102       74
## ENSG000000000938      0         0         0

```

```
head(metadata)
```

```

##      id   dex celltype geo_id
## 1 SRR1039508 control N61311 GSM1275862
## 2 SRR1039509 treated N61311 GSM1275863
## 3 SRR1039512 control N052611 GSM1275866
## 4 SRR1039513 treated N052611 GSM1275867
## 5 SRR1039516 control N080611 GSM1275870
## 6 SRR1039517 treated N080611 GSM1275871

```

Q1. How many genes are in this dataset?

```
nrow(counts)
```

```
## [1] 38694
```

This dataset has 38694 genes.

Q2. How many ‘control’ cell lines do we have?

```
grep("control",metadata[,2])
```

```
## [1] 1 3 5 7
```

```
sum(metadata$dex == "control")
```

```
## [1] 4
```

There are 4 ‘control’ cell lines.

```
#Toy differential gene expression
```

We know row 1,3,5, and 7 are control samples, so we can now use that to find a control mean.

```

control <- metadata[metadata[, "dex"]=="control",]
control.counts <- counts[ ,control$id]
control.mean <- rowSums( control.counts )/4
head(control.mean)

## ENSG000000000003 ENSG000000000005 ENSG000000000419 ENSG000000000457 ENSG000000000460
##          900.75           0.00        520.50       339.75      97.25
## ENSG00000000938
##          0.75

```

Q3. How would you make the above code in either approach more robust?

Instead of rowSums, we can do rowMeans because doing rowSums then dividing it by 4 is less robust, and would not work if the data set changed. A better way to do it is as follows using rowMeans:

```

control.ind <- metadata$dex == "control"
control.counts <- counts[ ,control.ind]
control.mean <- rowMeans( control.counts )
head(control.mean)

## ENSG000000000003 ENSG000000000005 ENSG000000000419 ENSG000000000457 ENSG000000000460
##          900.75           0.00        520.50       339.75      97.25
## ENSG00000000938
##          0.75

```

Q4. Follow the same procedure for the treated samples (i.e. calculate the mean per gene across drug treated samples and assign to a labeled vector called treated.mean)

```

treated.ind <- metadata$dex == "treated"
treated.counts <- counts[ ,treated.ind]
treated.mean <- rowMeans( treated.counts )
head(treated.mean)

## ENSG000000000003 ENSG000000000005 ENSG000000000419 ENSG000000000457 ENSG000000000460
##          658.00           0.00        546.00       316.50      78.75
## ENSG00000000938
##          0.00

```

For easier access we can put our means together.

```

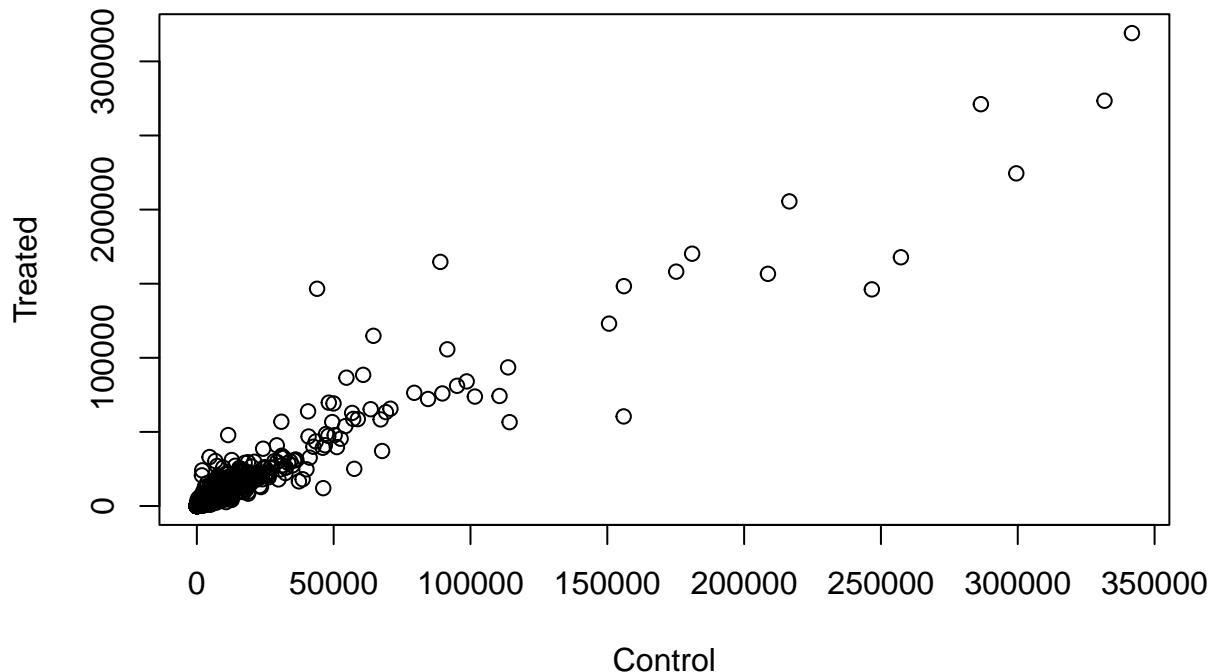
meancounts <- data.frame(control.mean, treated.mean)
colSums(meancounts)

## control.mean treated.mean
##      23005324     22196524

```

Q5 (a). Create a scatter plot showing the mean of the treated samples against the mean of the control samples. Your plot should look something like the following.

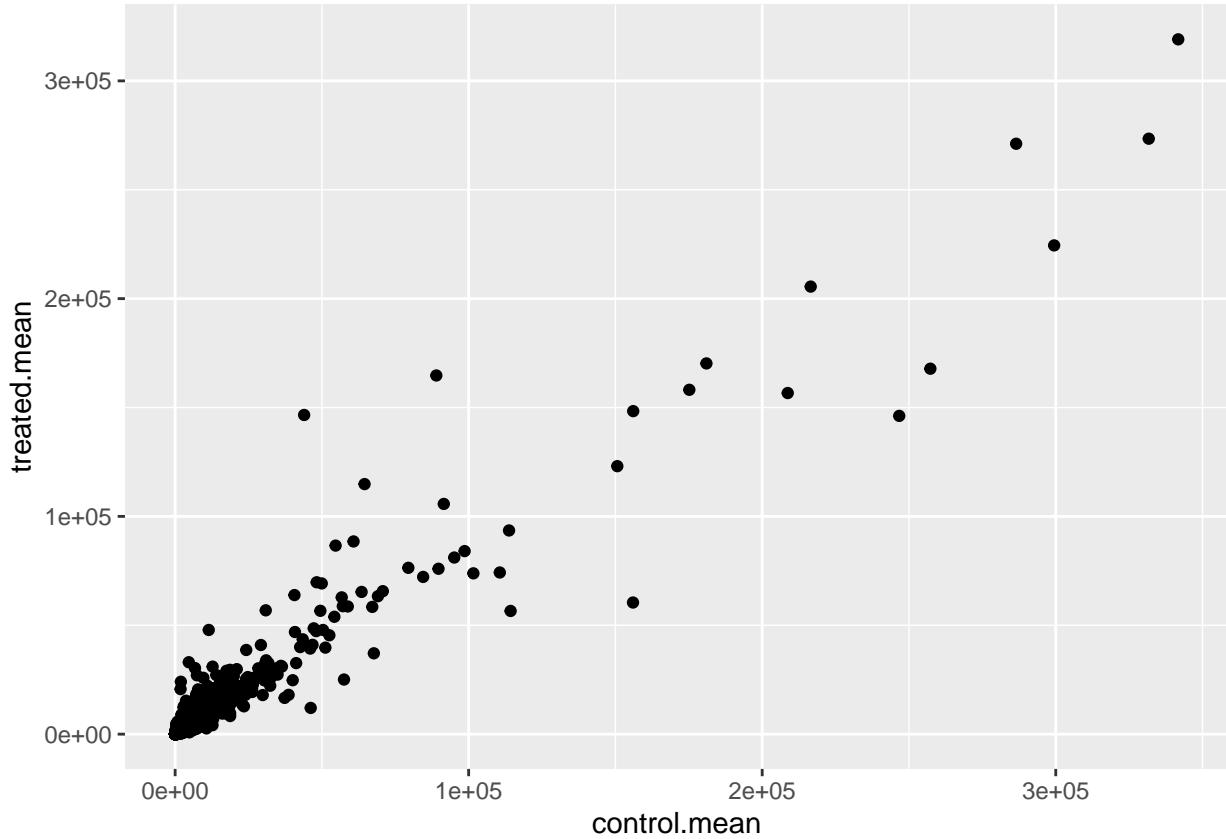
```
plot(control.mean, treated.mean, xlab= "Control", ylab= "Treated")
```



Q5 (b). You could also use the ggplot2 package to make this figure producing the plot below. What geom_?() function would you use for this plot?

We would use geom_point after calling the ggplot2 package to make a similar figure.

```
library(ggplot2)
ggplot(meancounts, aes(x=control.mean, y=treated.mean)) + geom_point()
```



These plots are not great, and we need a log transformation in order to see more details from our data.

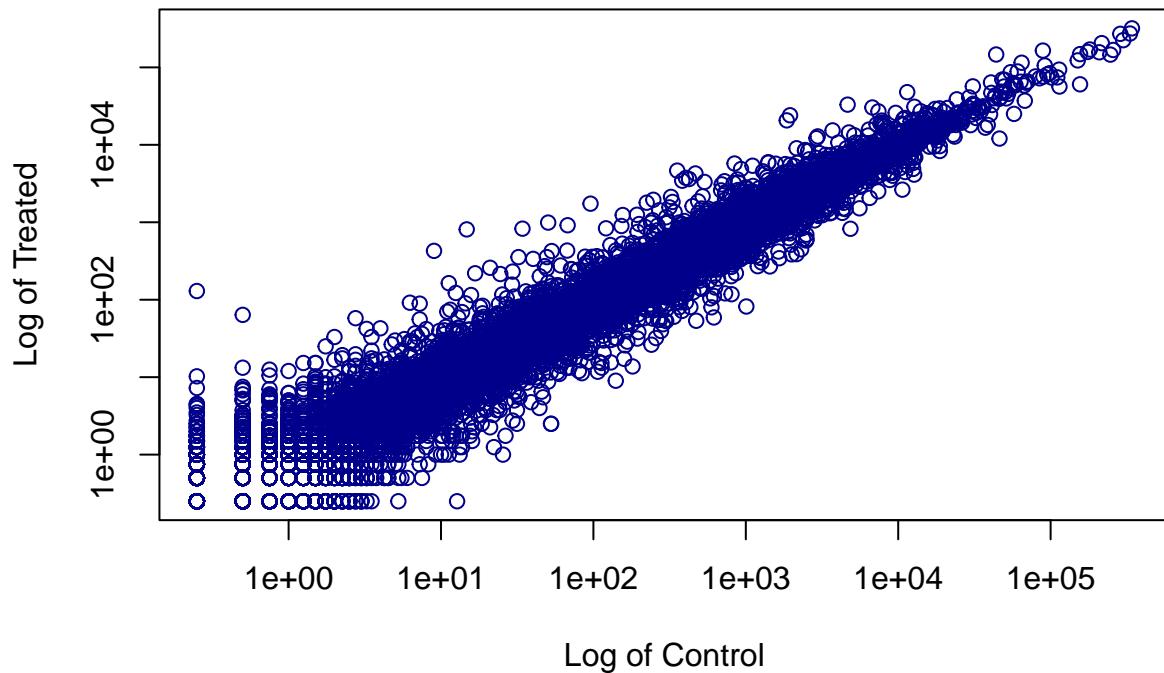
Q6. Try plotting both axes on a log scale. What is the argument to plot() that allows you to do this?

Adding the argument log to the plot() argument will allow for the axes to be on a log scale.

```
plot(meancounts, log = "xy", xlab= "Log of Control", ylab= "Log of Treated", col = "dark blue")

## Warning in xy.coords(x, y, xlabel, ylabel, log): 15032 x values <= 0 omitted
## from logarithmic plot

## Warning in xy.coords(x, y, xlabel, ylabel, log): 15281 y values <= 0 omitted
## from logarithmic plot
```



We often use log2 in the biology field because of its mathematical properties that make interpretation more straight forward.

```
#lets look at a math example to understand more about log2
log2(10/10)
```

```
## [1] 0
```

```
log2(20/10)
```

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

```
log2(40/10)
```

```
## [1] 2
```

```
log2(80/10)
```

```
## [1] 3
```

```
# we can also look at decreasing
log2(5/10)
```

```
## [1] -1
log2(2.5/10)
```

```
## [1] -2
```

As we see above, when there is no change, there is a 0 value. When there is increases there are positive values, and then for decreases there are negative values. This property leads us to work with $\log_2(foldchange)$ all the time in the genomics and proteomic fields.

```
meancounts$log2fc <- log2(meancounts[, "treated.mean"] / meancounts[, "control.mean"])
head(meancounts)
```

	control.mean	treated.mean	log2fc
## ENSG00000000003	900.75	658.00	-0.45303916
## ENSG00000000005	0.00	0.00	NaN
## ENSG00000000419	520.50	546.00	0.06900279
## ENSG00000000457	339.75	316.50	-0.10226805
## ENSG00000000460	97.25	78.75	-0.30441833
## ENSG00000000938	0.75	0.00	-Inf

Some of the results are not helpful, specifically NaN (not a number) and -Inf (negative infinity). These are due to zeros in our gene set, and are a result of dividing by or doing the log of 0. So let's edit our data to remove those genes.

```
#first, we find the zeros.
zero.vals <- which(meancounts[, 1:2]==0, arr.ind=TRUE)
to.rm <- unique(zero.vals[, 1])
#Now we remove them from our meancounts dataframe.
mycounts <- meancounts[-to.rm,]
head(mycounts)
```

	control.mean	treated.mean	log2fc
## ENSG00000000003	900.75	658.00	-0.45303916
## ENSG00000000419	520.50	546.00	0.06900279
## ENSG00000000457	339.75	316.50	-0.10226805
## ENSG00000000460	97.25	78.75	-0.30441833
## ENSG00000000971	5219.00	6687.50	0.35769358
## ENSG00000001036	2327.00	1785.75	-0.38194109

How many genes are left now?

```
nrow(mycounts)
```

```
## [1] 21817
```

Q7. What is the purpose of the arr.ind argument in the which() function call above?
 Why would we then take the first column of the output and need to call the unique() function?

The arr.ind argument within the which() function is for array indices, and it indicates when marked as TRUE, that the array indices should be returned when both positions are arrays. So now we know which positions have zeros. We use the unique() function to extract unique elements. Specifically in this case, we use it to not count any row with a zero in both samples more than once.

A common threshold used for calling something differentially expressed is a log2(FoldChange) of greater than 2 or less than -2. Let's filter the dataset both ways to see how many genes are up or down-regulated.

```
up.ind <- mycounts$log2fc > 2  
down.ind <- mycounts$log2fc < (-2)
```

Q8. Using the up.ind vector above can you determine how many up regulated genes we have at the greater than 2 fc level?

```
count(up.ind)
```

```
## [1] 250
```

Q9. Using the down.ind vector above can you determine how many down regulated genes we have at the greater than 2 fc level?

```
count(down.ind)
```

```
## [1] 367
```

Q10. Do you trust these results? Why or why not?

While this is a good preliminary analysis, we do not have any indicators of statistical significance that can lead us to trust a conclusion of sorts, so we should not trust them. We must continue to analyze in order to actually have results we can trust.

```
#DESeq2 analysis
```

We can call our DESeq2 package to better analyze count data like ours.

```
library(DESeq2)  
citation("DESeq2")  
  
##  
## Love, M.I., Huber, W., Anders, S. Moderated estimation of fold change  
## and dispersion for RNA-seq data with DESeq2 Genome Biology 15(12):550  
## (2014)  
##  
## A BibTeX entry for LaTeX users is  
##  
## @Article{,  
##   title = {Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2},  
##   author = {Michael I. Love and Wolfgang Huber and Simon Anders},  
##   year = {2014},  
##   journal = {Genome Biology},  
##   doi = {10.1186/s13059-014-0550-8},  
##   volume = {15},  
##   issue = {12},  
##   pages = {550},  
## }
```

Now lets set up the input.

```
dds <- DESeqDataSetFromMatrix(countData=counts,
                               colData=metadata,
                               design=~dex)

## converting counts to integer mode

## Warning in DESeqDataSet(se, design = design, ignoreRank): some variables in
## design formula are characters, converting to factors

dds

## class: DESeqDataSet
## dim: 38694 8
## metadata(1): version
## assays(1): counts
## rownames(38694): ENSG00000000003 ENSG00000000005 ... ENSG00000283120
##   ENSG00000283123
## rowData names(0):
## colnames(8): SRR1039508 SRR1039509 ... SRR1039520 SRR1039521
## colData names(4): id dex celltype geo_id
```

Next, we can run DESeq analysis.

```
dds <- DESeq(dds)

## estimating size factors

## estimating dispersions

## gene-wise dispersion estimates

## mean-dispersion relationship

## final dispersion estimates

## fitting model and testing
```

Then, we can get the results.

```
res <- results(dds)
res

## log2 fold change (MLE): dex treated vs control
## Wald test p-value: dex treated vs control
## DataFrame with 38694 rows and 6 columns
##           baseMean log2FoldChange      lfcSE       stat     pvalue
## <numeric>      <numeric> <numeric> <numeric> <numeric>
## ENSG00000000003  747.1942    -0.3507030  0.168246 -2.084470 0.0371175
```

```

## ENSG00000000005    0.0000      NA      NA      NA      NA
## ENSG00000000419  520.1342   0.2061078  0.101059  2.039475  0.0414026
## ENSG00000000457  322.6648   0.0245269  0.145145  0.168982  0.8658106
## ENSG00000000460   87.6826  -0.1471420  0.257007 -0.572521  0.5669691
## ...
## ENSG00000283115  0.000000      NA      NA      NA      NA
## ENSG00000283116  0.000000      NA      NA      NA      NA
## ENSG00000283119  0.000000      NA      NA      NA      NA
## ENSG00000283120  0.974916  -0.668258   1.69456 -0.394354  0.693319
## ENSG00000283123  0.000000      NA      NA      NA      NA
##          padj
## <numeric>
## ENSG00000000003  0.163035
## ENSG00000000005      NA
## ENSG00000000419  0.176032
## ENSG00000000457  0.961694
## ENSG00000000460  0.815849
## ...
## ENSG00000283115      NA
## ENSG00000283116      NA
## ENSG00000283119      NA
## ENSG00000283120      NA
## ENSG00000283123      NA

```

We need to convert the object to a data frame as well to view it in full.

```
#View(as.data.frame(res))
```

Some of the results have NA, and their baseMean is 0. This is because DESeq does independent filtering that filters our results with little or no chance of significance.

```
summary(res)
```

```

##
## out of 25258 with nonzero total read count
## adjusted p-value < 0.1
## LFC > 0 (up)      : 1563, 6.2%
## LFC < 0 (down)     : 1188, 4.7%
## outliers [1]       : 142, 0.56%
## low counts [2]     : 9971, 39%
## (mean count < 10)
## [1] see 'cooksCutoff' argument of ?results
## [2] see 'independentFiltering' argument of ?results

```

```
res05 <- results(dds, alpha=0.05)
summary(res05)
```

```

##
## out of 25258 with nonzero total read count
## adjusted p-value < 0.05
## LFC > 0 (up)      : 1236, 4.9%
## LFC < 0 (down)     : 933, 3.7%

```

```

## outliers [1]      : 142, 0.56%
## low counts [2]    : 9033, 36%
## (mean count < 6)
## [1] see 'cooksCutoff' argument of ?results
## [2] see 'independentFiltering' argument of ?results

```

We will come back to this section next class

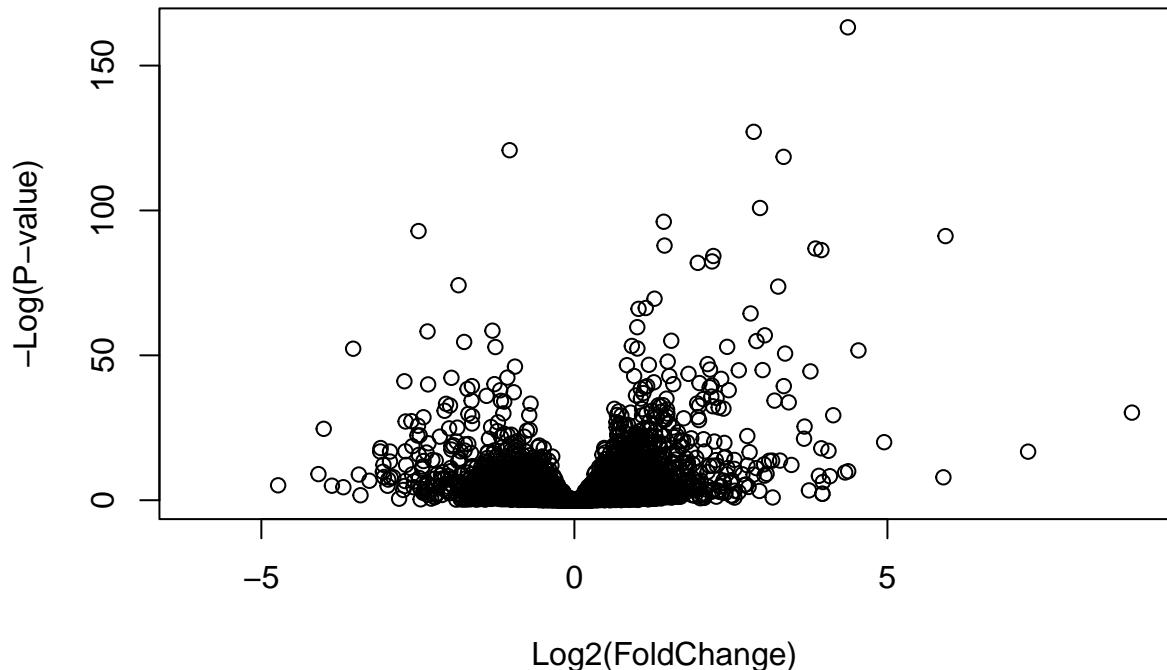
Data Visualization

Let's make a volcano plot. These summary figures are frequently used to highlight the proportion of genes that are both significantly regulated and display a high fold change.

```

plot( res$log2FoldChange, -log(res$padj),
      xlab="Log2(FoldChange)",
      ylab="-Log(P-value)")

```



This plot is not very useful because all the p-values of significance are clustered and hidden in the bottom. We can transform this with log and add lines to make this visualization more useful.

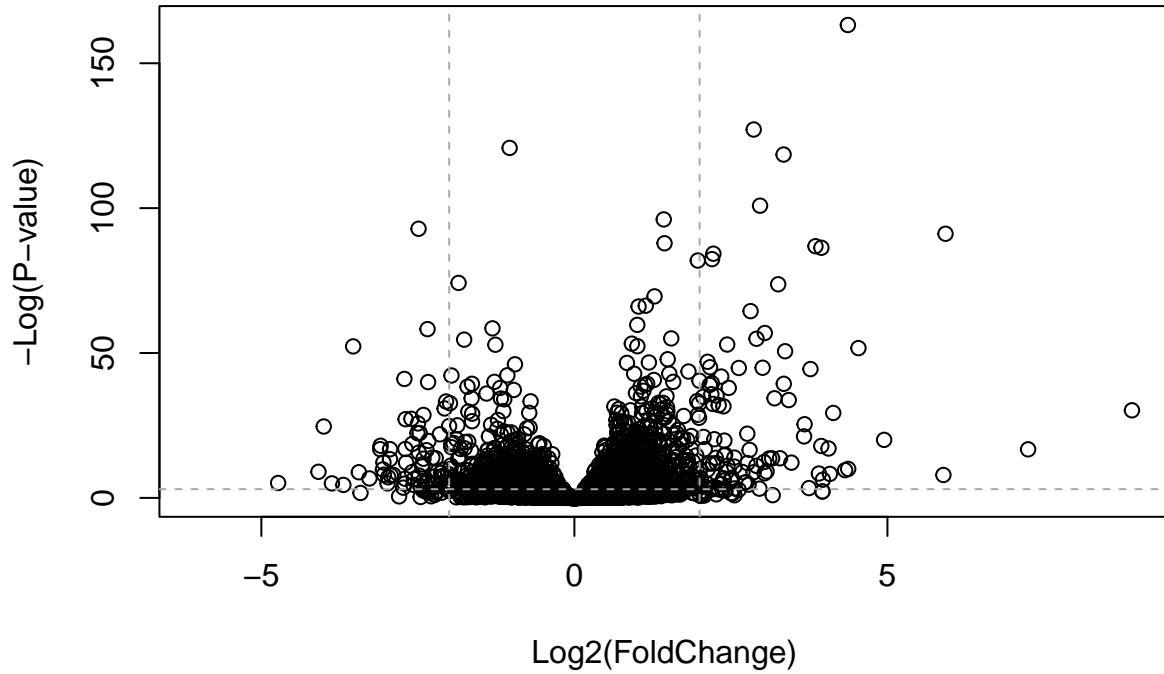
```

plot( res$log2FoldChange, -log(res$padj),
      ylab="-Log(P-value)", xlab="Log2(FoldChange)")

# Add some cut-off lines

```

```
abline(v=c(-2,2), col="darkgray", lty=2)
abline(h=-log(0.05), col="darkgray", lty=2)
```



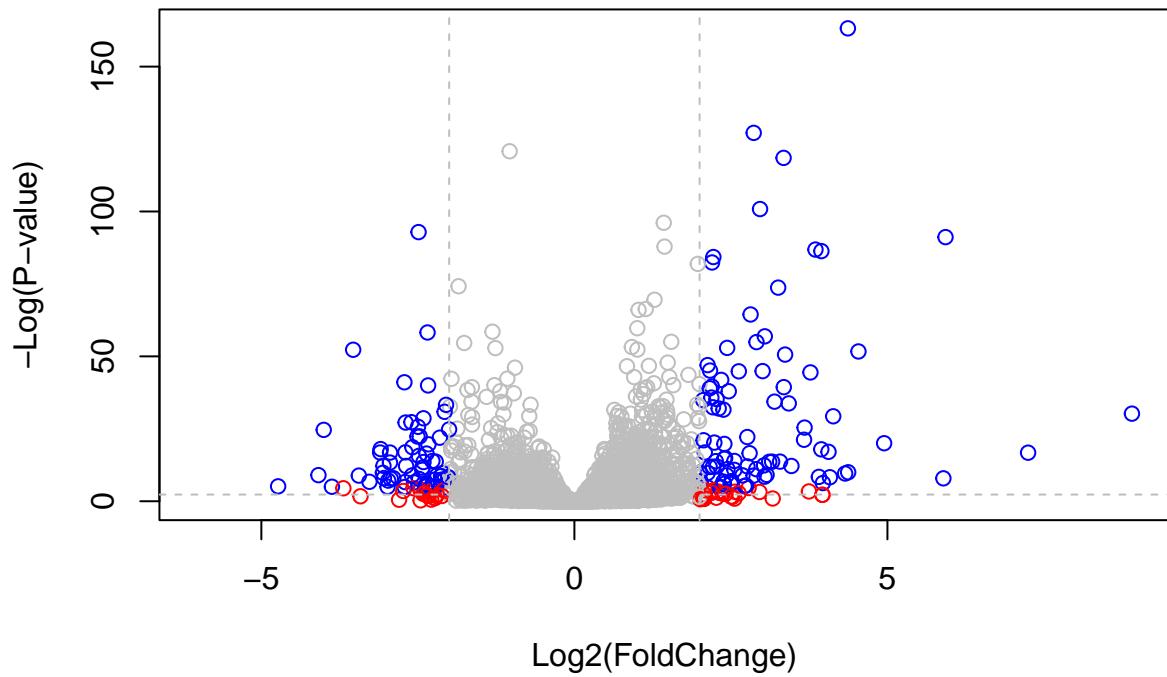
Now we can plot it once again but re-color the plot based on the genes (i.e points) that we care about – that is those with large fold-change and low p-values.

```
# Setup our custom point color vector
mycols <- rep("gray", nrow(res))
mycols[ abs(res$log2FoldChange) > 2 ] <- "red"

inds <- (res$padj < 0.01) & (abs(res$log2FoldChange) > 2 )
mycols[ inds ] <- "blue"

# Volcano plot with custom colors
plot( res$log2FoldChange, -log(res$padj),
      col=mycols, ylab="-Log(P-value)", xlab="Log2(FoldChange)" )

# Cut-off lines
abline(v=c(-2,2), col="gray", lty=2)
abline(h=-log(0.1), col="gray", lty=2)
```



Instead of this, we can also use the “EnhancedVolcano” feature of BiocManager that utilizes ggplot under the hood, which we will do next class.