

Class 15

Natalie Avina (PID: A15590695)

11/16/2021

countData and colData

```
head(counts)
```

```
##          SRR1039508 SRR1039509 SRR1039512 SRR1039513 SRR1039516
## ENSG000000000003      723        486       904       445      1170
## ENSG000000000005       0         0         0         0         0
## ENSG00000000419      467       523       616       371      582
## ENSG00000000457      347       258       364       237      318
## ENSG00000000460       96        81        73        66      118
## ENSG00000000938       0         0         1         0         2
##          SRR1039517 SRR1039520 SRR1039521
## ENSG000000000003     1097       806       604
## ENSG000000000005       0         0         0
## ENSG00000000419      781       417       509
## ENSG00000000457      447       330       324
## ENSG00000000460       94        102       74
## ENSG00000000938       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
```

How many genes are in this dataset?

```
nrow(counts)
```

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

How many control cell lines do we have?

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

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

First I need to extract all the control columns. Then I will take the rowwise mean to get the average count values for all genes in these four experiments.

```
control inds <- metadata$dex == "control"  
control.counts <- counts[,control inds]  
head(control.counts)
```

```
##          SRR1039508 SRR1039512 SRR1039516 SRR1039520  
## ENSG000000000003     723      904     1170      806  
## ENSG000000000005      0        0        0        0  
## ENSG000000000419    467      616      582      417  
## ENSG000000000457    347      364      318      330  
## ENSG000000000460     96       73      118      102  
## ENSG000000000938      0        1        2        0
```

```
control.mean <- rowMeans(control.counts)
```

Do the same thing for “treated”.

```
treated inds <- metadata$dex == "treated"  
treated.counts <- counts[,treated inds]  
head(treated.counts)
```

```
##          SRR1039509 SRR1039513 SRR1039517 SRR1039521  
## ENSG000000000003    486      445     1097      604  
## ENSG000000000005      0        0        0        0  
## ENSG000000000419    523      371      781      509  
## ENSG000000000457    258      237      447      324  
## ENSG000000000460     81       66      94       74  
## ENSG000000000938      0        0        0        0
```

```
treated.mean <- rowMeans(treated.counts)
```

We are combining the data from both means into a data frame.

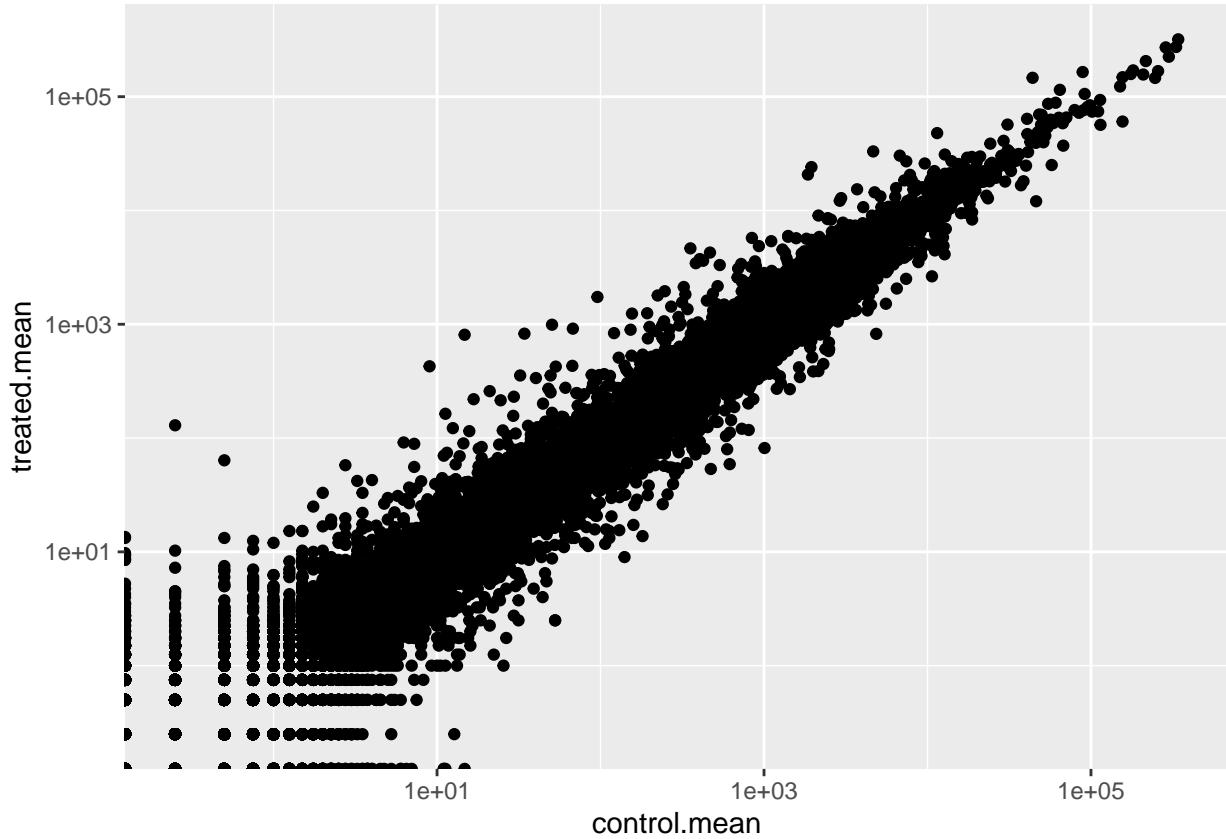
```
meancounts <- data.frame(control.mean, treated.mean)
```

Plot control vs treated

```
library(ggplot2)  
ggplot(meancounts, aes(x = control.mean, y = treated.mean)) + geom_point() + scale_x_log10() + scale_y_
```



```
## Warning: Transformation introduced infinite values in continuous x-axis  
## Warning: Transformation introduced infinite values in continuous y-axis
```



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

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

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

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

```
log2(5/10)
```

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

We see 0 values for no change, positive values for increased expression, and negative values for decreased expression. **log2(fold-change)** is used a lot in genomics/proteomics.

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

	control.mean	treated.mean	log2fc
## ENSG000000000003	900.75	658.00	-0.45303916
## ENSG000000000005	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

We have to get rid of the genes (rows) with 0.00 counts because we can't say anything about these; they have no data.

```
head(meancounts[,1:2]==0)
```

```
## control.mean treated.mean
## ENSG000000000003 FALSE FALSE
## ENSG000000000005 TRUE TRUE
## ENSG00000000419 FALSE FALSE
## ENSG00000000457 FALSE FALSE
## ENSG00000000460 FALSE FALSE
## ENSG00000000938 FALSE TRUE
```

```
which(c(F,F,T,T))
```

```
## [1] 3 4
```

We can use the **which()** function with ‘arr.ind=TRUE’ argument to get the columns and rows where the TRUE (zero counts) values are.

```
zero.vals <- which(meancounts[,1:2] == 0, arr.ind = TRUE)
head(zero.vals)
```

```
## row col
## ENSG000000000005 2 1
## ENSG00000004848 65 1
## ENSG00000004948 70 1
## ENSG00000005001 73 1
## ENSG00000006059 121 1
## ENSG00000006071 123 1
```

```
to.rm <- unique(zero.vals[,"row"])
head(sort(to.rm))
```

```
## [1] 2 6 65 70 73 81
```

```
mycounts <- meancounts[-to.rm,]
head(mycounts)
```

```
## control.mean treated.mean log2fc
## ENSG000000000003 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
```

```
nrow(mycounts)
```

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

```

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, colAnys, 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, rowAnys, 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

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},
##    }

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

dds <- DESeq(dds)

## estimating size factors

## estimating dispersions

## gene-wise dispersion estimates

## mean-dispersion relationship

## final dispersion estimates

## fitting model and testing

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

```

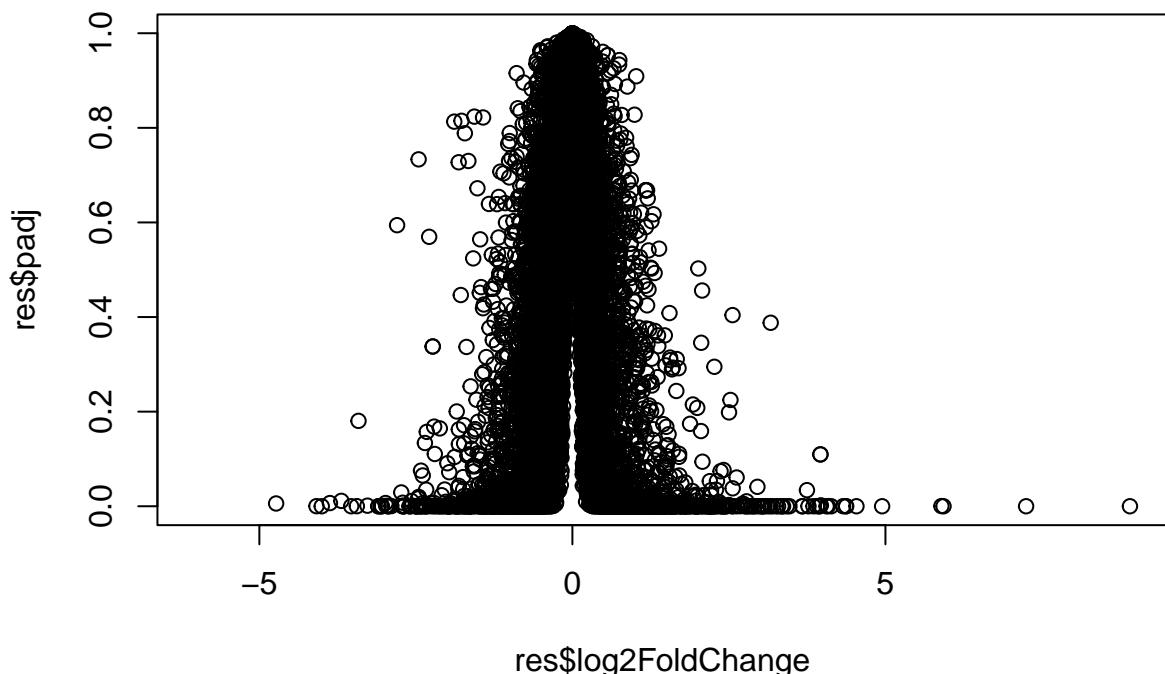
```

## 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>
## ENSG0000000003 0.163035
## ENSG0000000005 NA
## ENSG00000000419 0.176032
## ENSG00000000457 0.961694
## ENSG00000000460 0.815849
## ...
## ...
## ENSG00000283115 NA
## ENSG00000283116 NA
## ENSG00000283119 NA
## ENSG00000283120 NA
## ENSG00000283123 NA

```

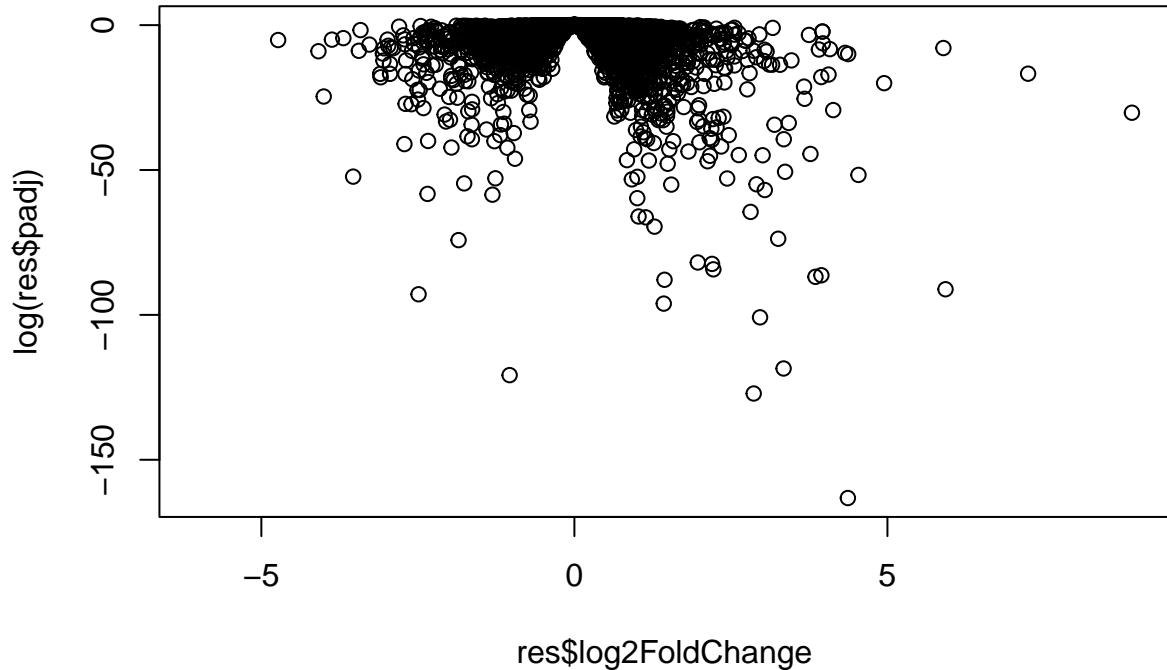
VOLCANO PLOTSSSSS

```
plot(res$log2FoldChange, res$padj)
```



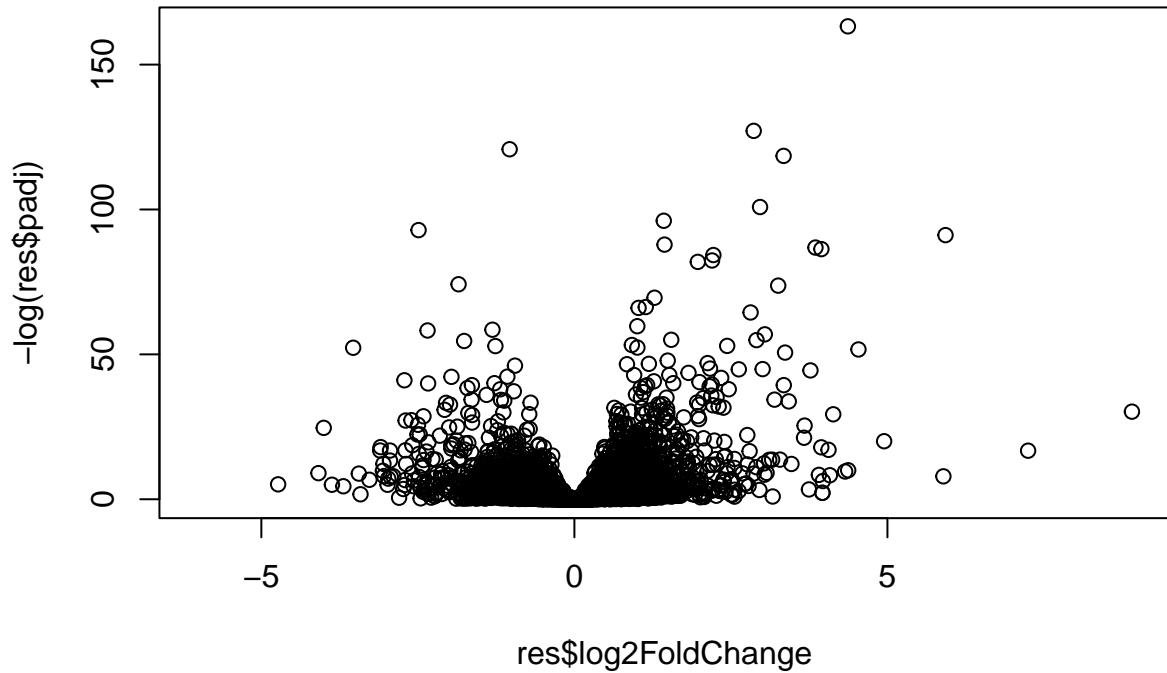
That was not a useful plot because all the small p-values are hidden at the bottom of the plot and we can't see them very well. Log can help

```
plot(res$log2FoldChange, log(res$padj))
```



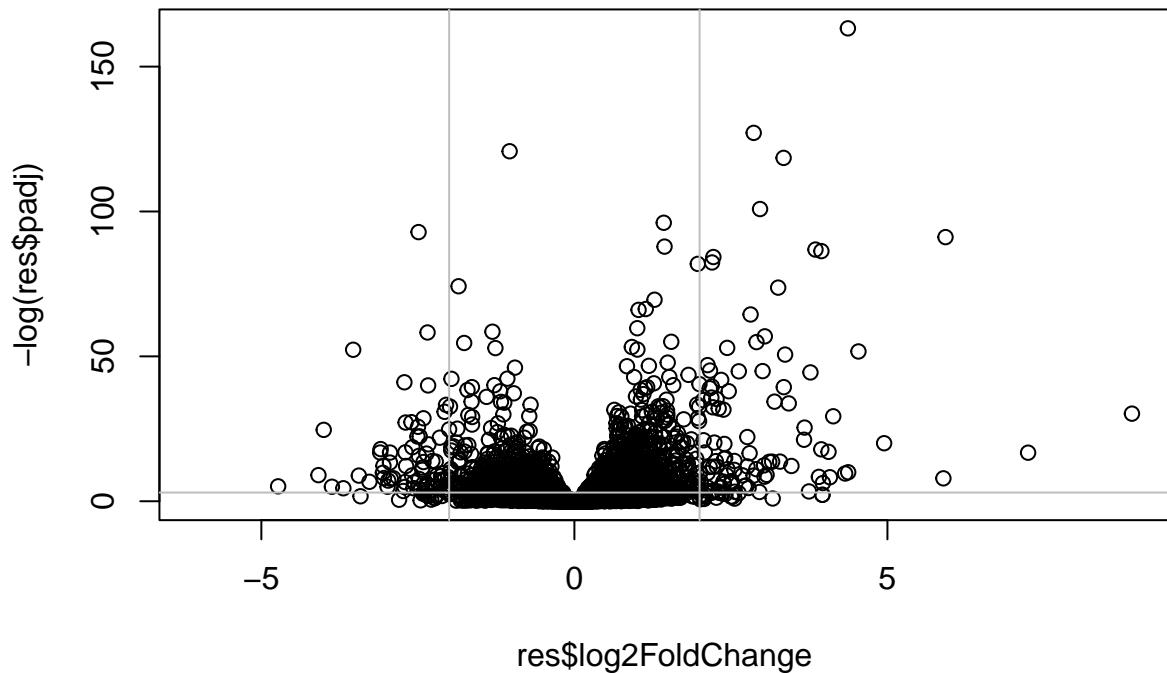
We can flip this p-value axis by putting a minus sign on it, giving us a classic volcano plot.

```
plot(res$log2FoldChange, -log(res$padj))
```



Let's make our plot pretty by adding color to draw attention to the genes we care about- the points with large fold-change and small p-values (high -log values)

```
plot(res$log2FoldChange, -log(res$padj))
abline(v=c(-2, +2), col= "gray")
abline(h=-log(0.05), col="gray")
```



```
mycols <- rep("gray", nrow(res))
mycols[abs(res$log2FoldChange)>2] <- "red"
mycols[res$padj>0.05] <- "gray"
```

```
plot(res$log2FoldChange, -log(res$padj), col=mycols)
abline(v=c(-2, +2), col= "gray")
abline(h=-log(0.05), col="gray")
```

