

Class 15 RNA-Seq

Claire Chapman

11/17/2021

Read in Data

DESeq2 expects a dataframe of count data and a second dataframe with information about the samples (metadata)

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

Take a look

```
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
```

There are 38694 genes in this dataset

Column names of counts data MUST equal the IDs in metadata. Check this:

```
metadata$id == colnames(counts)
```

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

Exploratory Diff Gene Analysis - Compare control to treated

For demonstration purposes only, never do actual diff analysis this way

Extract row from metadata with controls

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

```
## ENSG00000000003 ENSG00000000005 ENSG000000000419 ENSG000000000457 ENSG000000000460  
##           900.75           0.00           520.50           339.75           97.25  
## ENSG000000000938  
##           0.75
```

OR use dplyr (I prefer)

```
library(dplyr)
```

```
##  
## Attaching package: 'dplyr'  
  
## The following objects are masked from 'package:stats':  
##  
##     filter, lag  
  
## The following objects are masked from 'package:base':  
##  
##     intersect, setdiff, setequal, union
```

```
control <- metadata %>% filter(dex=="control")  
control.counts <- counts %>% select(control$id)  
control.mean <- rowSums(control.counts)/4  
head(control.mean)
```

```
## ENSG00000000003 ENSG00000000005 ENSG000000000419 ENSG000000000457 ENSG000000000460  
##           900.75           0.00           520.50           339.75           97.25  
## ENSG000000000938  
##           0.75
```

```
treated <- metadata %>% filter(dex=="treated")  
treated.counts <- counts %>% select(treated$id)  
treated.mean <- rowSums(treated.counts)/4  
head(treated.mean)
```

```
## ENSG00000000003 ENSG00000000005 ENSG000000000419 ENSG000000000457 ENSG000000000460
##          658.00          0.00          546.00          316.50          78.75
## ENSG000000000938
##          0.00
```

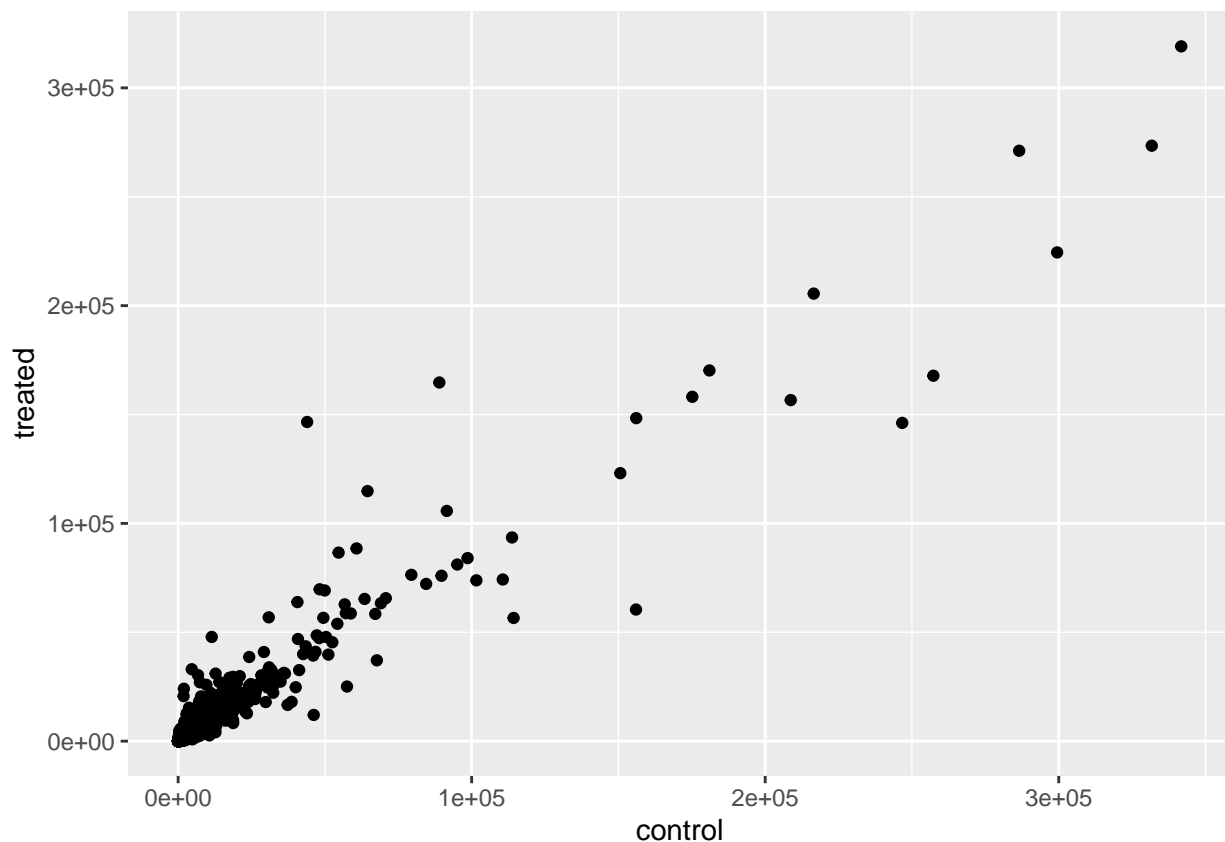
Combine mean count data

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

Compare the control and treated

Quick scatterplot to check work

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



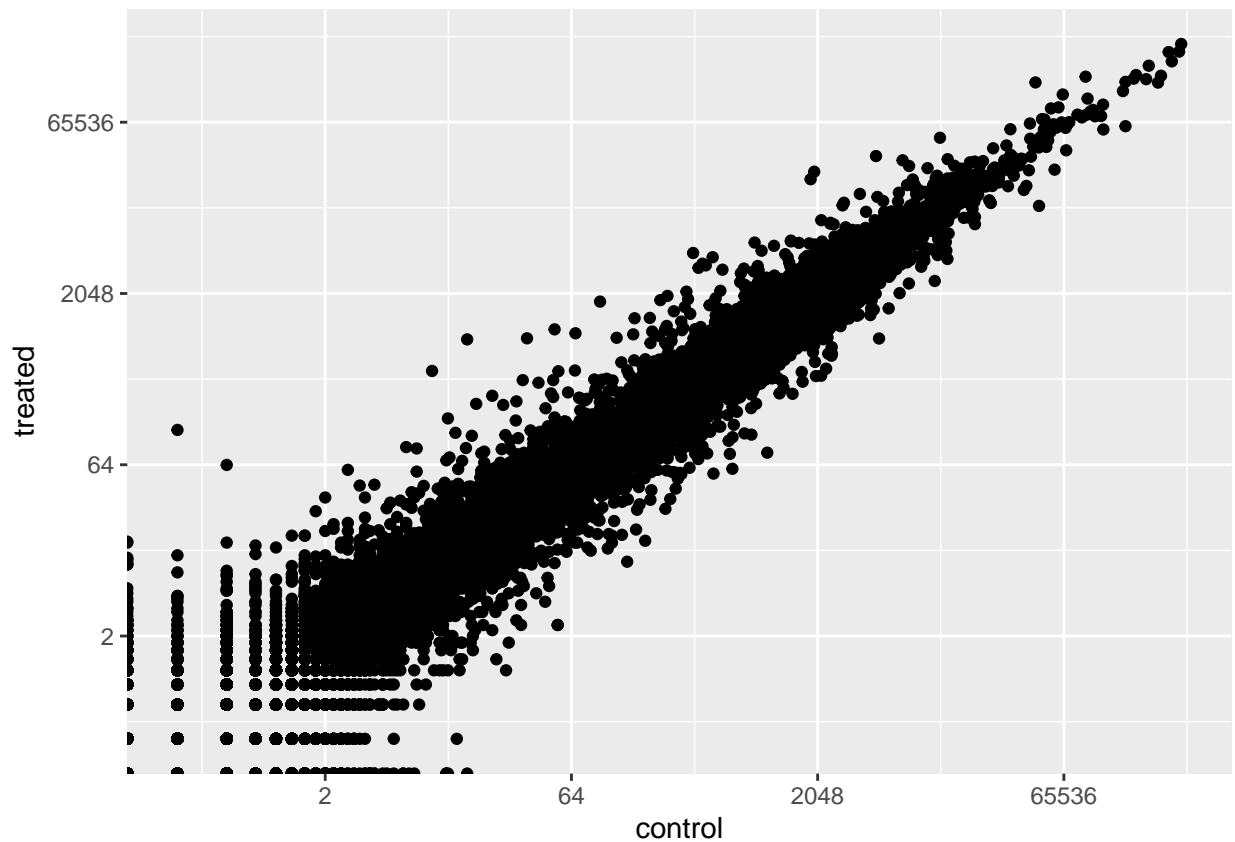
Not truly seeing all the 60,000 data

```
meancounts %>%
  ggplot(aes(control.mean, treated.mean)) +
  scale_y_continuous(trans = "log2")+
  scale_x_continuous(trans = "log2")+
  geom_point()
```

```
geom_point()+  
labs(x = "control", y = "treated")
```

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

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



We often use log transformations to make life easier...

If there is no difference, log change is 0

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

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

If there is twice as much, expression goes up, log change is 1

```
log2(40/20)
```

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

If there is half as much, expression goes down, log change is -1

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

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

Transform our data Adding a column called log2 fold change representing $\log_2(\text{treated}/\text{control})$

```
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
## ENSG000000000419      520.50      546.00  0.06900279
## ENSG000000000457      339.75      316.50 -0.10226805
## ENSG000000000460       97.25       78.75 -0.30441833
## ENSG000000000938        0.75        0.00        -Inf
```

You get some weird values “NaN” when you try to divide by zero “-inf” Need to exclude the zeros

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

```
##           control.mean treated.mean
## ENSG000000000003      FALSE      FALSE
## ENSG000000000005       TRUE       TRUE
## ENSG000000000419      FALSE      FALSE
## ENSG000000000457      FALSE      FALSE
## ENSG000000000460      FALSE      FALSE
## ENSG000000000938      FALSE       TRUE
```

the `which()` function tells us the indices of TRUE entries in a logical vector. It is not useful in default mode, doesn't break between columns. Need to argue with it.

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

```
##           row col
## ENSG000000000005     2  1
## ENSG000000004848    65  1
## ENSG000000004948    70  1
## ENSG000000005001    73  1
## ENSG000000006059   121  1
## ENSG000000006071   123  1
```

Some genes have 0 value in both columns, so select unique values

```
to.rm <- unique(inds[, 1])
mycounts <- meancounts[-to.rm, ]
head(mycounts)
```

##	control.mean	treated.mean	log2fc
## ENSG000000000003	900.75	658.00	-0.45303916
## ENSG000000000419	520.50	546.00	0.06900279
## ENSG000000000457	339.75	316.50	-0.10226805
## ENSG000000000460	97.25	78.75	-0.30441833
## ENSG000000000971	5219.00	6687.50	0.35769358
## ENSG00000001036	2327.00	1785.75	-0.38194109

We now have 21817 genes remaining.

Up-regulation and Down-regulation

A common threshold for log2fc is up-regulated if > 2 and down-regulated if < -2 .

```
up_reg <- mycounts$log2fc > 2
down_reg <- mycounts$log2fc < -2
```

```
sum(up_reg)
```

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

There are 250 genes up-regulated, or 1.1458954 %

There are 367 genes down-regulated or 1.6821745 %

DESeq2 Analysis

```
library(DESeq2)
```

```
## Loading required package: S4Vectors
```

```
## Loading required package: stats4
```

```
## Loading required package: BiocGenerics
```

```
##
```

```
## Attaching package: 'BiocGenerics'
```

```
## The following objects are masked from 'package:dplyr':
```

```
##
```

```
##      combine, intersect, setdiff, union
```

```
## 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:dplyr':
##
##   first, rename

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

## Loading required package: IRanges

##
## Attaching package: 'IRanges'

## The following objects are masked from 'package:dplyr':
##
##   collapse, desc, slice

## The following object is masked from 'package:grDevices':
##
##   windows

## Loading required package: GenomicRanges

## Loading required package: GenomeInfoDb

## Loading required package: SummarizedExperiment

## Loading required package: MatrixGenerics

## Loading required package: matrixStats

##
## Attaching package: 'matrixStats'

## The following object is masked from 'package:dplyr':
##
##   count

```

```

##
## Attaching package: 'MatrixGenerics'

## The following objects are masked from 'package:matrixStats':
##
##   colAlls, colAnyNAs, colAnys, colAvgPerRowSet, 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, rowAvgPerColSet,
##   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

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): ENSG000000000003 ENSG000000000005 ... 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
```

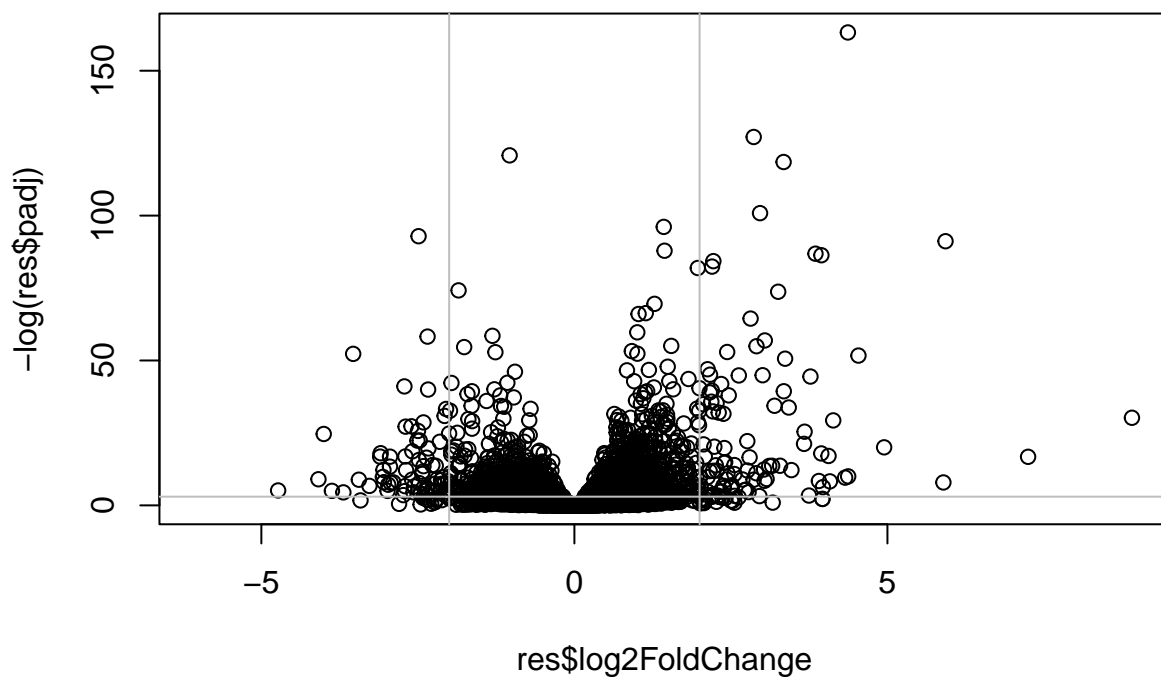
```
## log2 fold change (MLE): dex treated vs control
## Wald test p-value: dex treated vs control
## DataFrame with 38694 rows and 6 columns
##      baseMean log2FoldChange      lfcSE      stat      pvalue
##      <numeric>      <numeric> <numeric> <numeric> <numeric>
## ENSG000000000003  747.1942   -0.3507030  0.168246 -2.084470  0.0371175
## ENSG000000000005    0.0000           NA           NA           NA           NA
## ENSG000000000419  520.1342    0.2061078  0.101059  2.039475  0.0414026
## ENSG000000000457  322.6648    0.0245269  0.145145  0.168982  0.8658106
## ENSG000000000460   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>
## ENSG000000000003  0.163035
## ENSG000000000005    NA
## ENSG000000000419  0.176032
## ENSG000000000457  0.961694
```

```
## ENSG00000000460 0.815849
## ...
## ENSG00000283115 NA
## ENSG00000283116 NA
## ENSG00000283119 NA
## ENSG00000283120 NA
## ENSG00000283123 NA
```

Volcano plot

Common way to visualize the results Add lines at our -2, 2 fc thresholds and at the significant p value threshold. We are interested in the top left and top right quadrants.

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



Adding annotation data

Use Bioconductor's main annotation packages. Must first install these packages in the console.

```
library("AnnotationDbi")
```

```
##
## Attaching package: 'AnnotationDbi'

## The following object is masked from 'package:dplyr':
##
##      select

library("org.Hs.eg.db")

##

columns(org.Hs.eg.db)

##      [1] "ACCNUM"      "ALIAS"      "ENSEMBL"      "ENSEMBLPROT" "ENSEMBLTRANS"
##      [6] "ENTREZID"    "ENZYME"     "EVIDENCE"     "EVIDENCEALL"  "GENENAME"
##     [11] "GENETYPE"    "GO"         "GOALL"        "IPI"          "MAP"
##     [16] "OMIM"        "ONTOLOGY"   "ONTOLOGYALL"  "PATH"         "PFAM"
##     [21] "PMID"        "PROSITE"    "REFSEQ"       "SYMBOL"       "UCSCKG"
##     [26] "UNIPROT"
```

We want to make a new value “SYMBOL” that includes the common gene name that will be widely recognized and used.

```
res$symbol <- mapIds(org.Hs.eg.db,
                     keys=row.names(res), # Our genenames
                     keytype="ENSEMBL",   # The format of our genenames
                     column="SYMBOL",     # The new format we want to add
                     multiVals="first")
```

'select()' returned 1:many mapping between keys and columns

```
head(res)

## log2 fold change (MLE): dex treated vs control
## Wald test p-value: dex treated vs control
## DataFrame with 6 rows and 7 columns
##      baseMean log2FoldChange      lfcSE      stat      pvalue
##      <numeric>      <numeric> <numeric> <numeric> <numeric>
## ENSG000000000003 747.194195   -0.3507030 0.168246 -2.084470 0.0371175
## ENSG000000000005  0.000000         NA         NA         NA         NA
## ENSG000000000419 520.134160    0.2061078 0.101059  2.039475 0.0414026
## ENSG000000000457 322.664844    0.0245269 0.145145  0.168982 0.8658106
## ENSG000000000460  87.682625   -0.1471420 0.257007 -0.572521 0.5669691
## ENSG000000000938  0.319167   -1.7322890 3.493601 -0.495846 0.6200029
##      padj      symbol
##      <numeric> <character>
## ENSG000000000003 0.163035      TSPAN6
## ENSG000000000005      NA      TNMD
## ENSG000000000419 0.176032      DPM1
## ENSG000000000457 0.961694      SCYL3
## ENSG000000000460 0.815849      C1orf112
## ENSG000000000938      NA      FGR
```

Now let's save it for next time!

```
write.csv(res, file = "DESeq_results.csv")
```