

Class11

Ethan Lai

2/23/2022

Let's call library(DESeq2) and load our two data files

```
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, 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

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

There are 38694 rows, ie “genes” in this dataset

```
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
## 7 SRR1039520 control   N061011 GSM1275874
## 8 SRR1039521 treated   N061011 GSM1275875
```

It looks like we have four treated (ie with drug) and four control (ie no drug). Does the drug do anything?

We want to compare treated vs control.

Fist, let’s make sure the metadata matches the counts data order

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

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

Let’s get going!

First, let’s take a summary statistic of all controls vs all treated:

```
#metadata[metadata$dex=="control", "id"]

controls<- counts[metadata$dex=="control"]
treateds<- counts[metadata$dex=="treated"]
head(controls)
```

```
##          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
```

```
head(treateds)
```

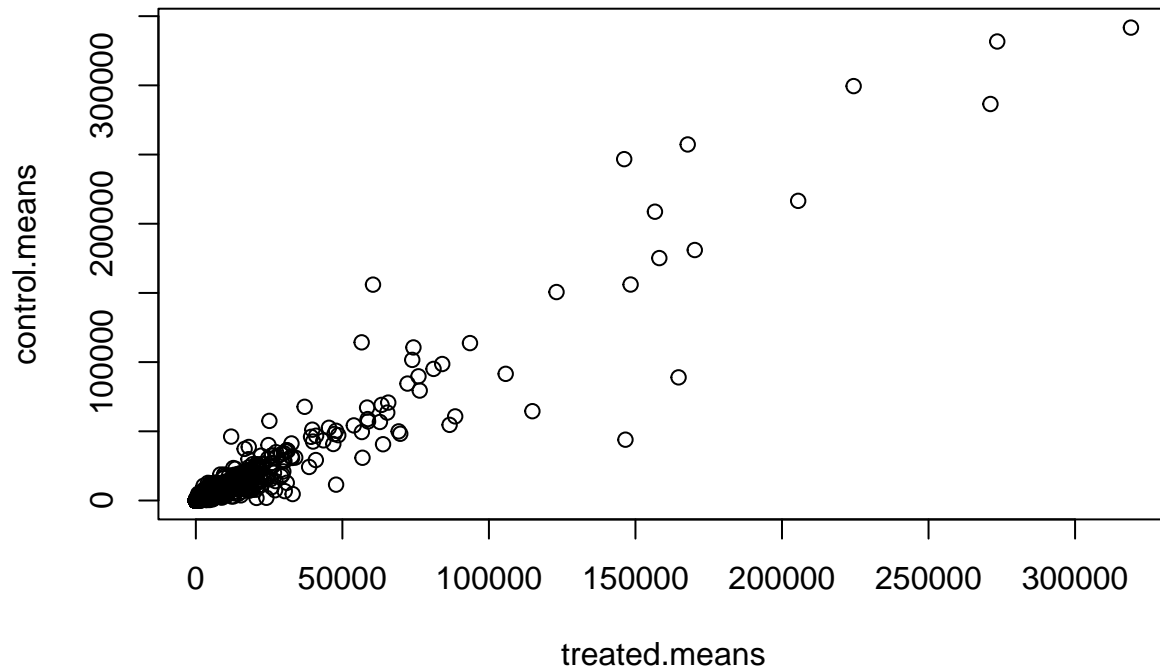
```
##                SRR1039509 SRR1039513 SRR1039517 SRR1039521
## ENSG000000000003         486         445         1097         604
## ENSG000000000005          0          0          0          0
## ENSG000000000419         523         371         781         509
## ENSG000000000457         258  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

Now that we have our control and treated data separated, we'll find the mean count values for each row ie gene. We could use `apply()` or more simply `rowMeans()`

```
control.means<- rowMeans(controls)
treated.means<-rowMeans(treateds)
```

Let's plot control vs treated means:

```
plot( treated.means, control.means,)
```

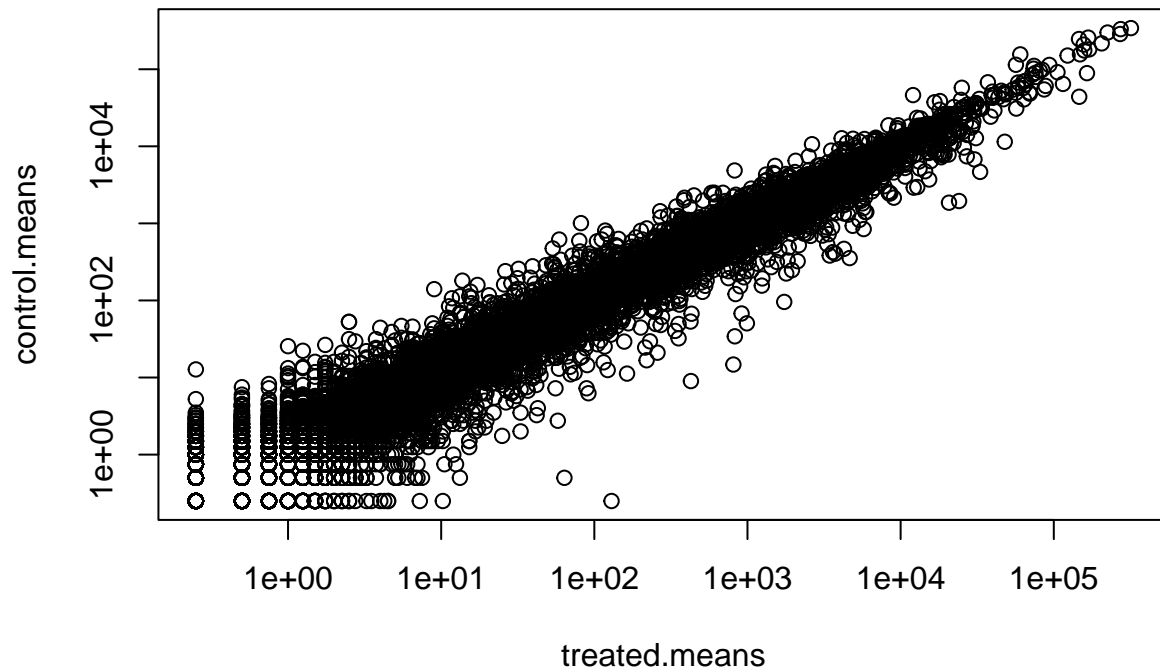


Most of the data is in the bottom left. Let's make it a log-log plot

```
plot( treated.means, control.means, log="xy")
```

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

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



We often use log transformation because it is more intuitive:

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

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

Zero fold change

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

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

1 fold log change, meaning doubling.

Let's compute the log2 fold change

```
log2fc <- log2(treated.means/control.means)
```

Let's make a data.frame to store our results to data:

```
meancounts<-data.frame(control.means, treated.means, log2fc)
head(meancounts)
```

```
##               control.means treated.means      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
```

We need a way to remove NaN and -Inf: We use the which() function in a complicated one-liner.

```
to.rm<- unique(which(meancounts[,1:2]==0, arr.ind=TRUE)[,"row"])
mycounts<- meancounts[-to.rm,]
```

Q1. How many genes do we have left? Q2. How many have a log2fc more than 2?

```
remainingGenes<- nrow(mycounts)
```

There are 21817 genes left.

```
greater<- sum(mycounts$log2fc>2)
```

There are 250 genes with log2fc >2

DESeq2

```
library("DESeq2")
```

First we need to setup the object that DESeq needs with

```
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
```

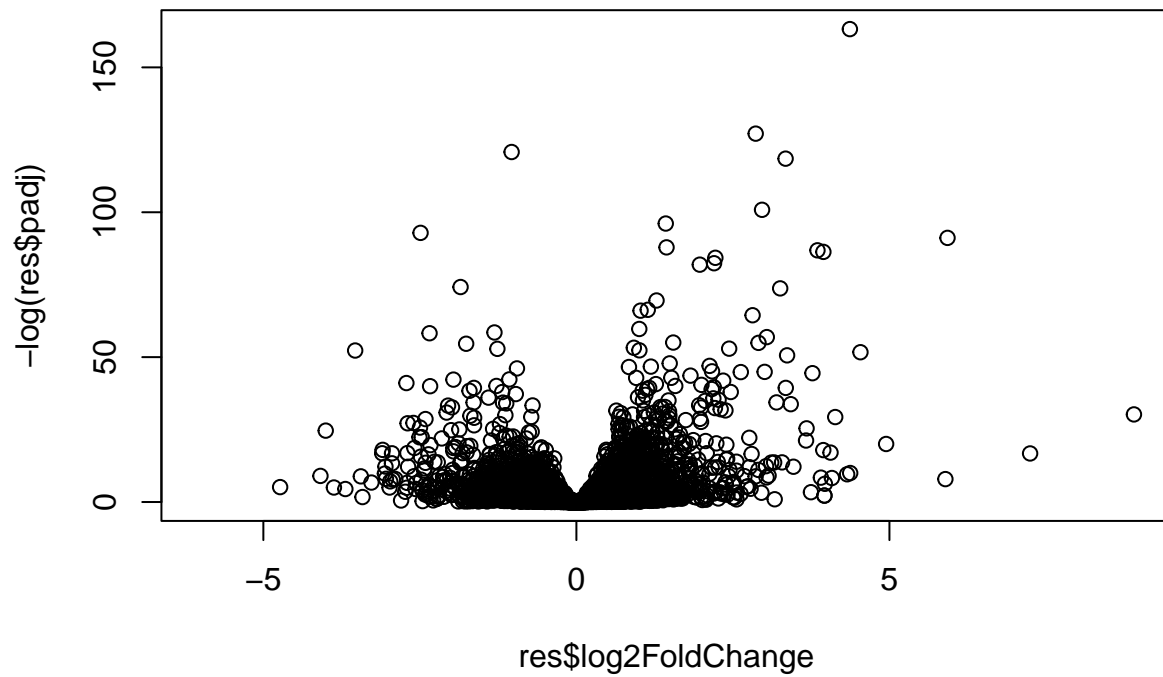
```
## DataFrame with 38694 rows and 6 columns
```

```
##           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
## ENSG000000000460  0.815849
## ...           ...
## ENSG00000283115    NA
## ENSG00000283116    NA
## ENSG00000283119    NA
## ENSG00000283120    NA
## ENSG00000283123    NA
```

A main result figure

A common main result figure from this type of analysis is a volcano plot. This is a plot of log2 fold change vs P value

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



Let's color code our significant hits

```
plot(res$log2FoldChange, -log(res$padj), col=ifelse(res$padj <= 0.05 & (res$log2FoldChange>2 | res$log2
```