# Lab 13

Bernice Lozada (A16297973)

```
library(BiocManager)
```

Bioconductor version '3.18' is out-of-date; the current release version '3.19'
  is available with R version '4.4'; see https://bioconductor.org/install

```
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, aperm, 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 object is masked from 'package:utils':

    findMatches

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

## Import countData and colData

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

head(counts)
```

|                 | SRR1039508 | SRR1039509 | SRR1039512 | SRR1039513 | SRR1039516 |
|-----------------|-----------:|-----------:|-----------:|-----------:|-----------:|
| ENSG00000000003 | 723        | 486        | 904        | 445        | 1170       |
| ENSG00000000005 | 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
ENSG00000000003       1097         806         604
ENSG00000000005          0           0           0
ENSG00000000419        781         417         509
ENSG00000000457        447         330         324
ENSG00000000460         94         102          74
ENSG00000000938          0           0           0
```

```r
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?

There are 38694 genes in this data set.

Q2. How many 'control' cell lines do we have?

There are 4 control cell lines.

## Extract and summarize the control samples

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

```
ENSG00000000003 ENSG00000000005 ENSG00000000419 ENSG00000000457 ENSG00000000460
         900.75            0.00          520.50          339.75           97.25
ENSG00000000938
           0.75
```

**Extract and summarize the treated samples**

```r
treated <- metadata[metadata$dex == "treated",]
treated.counts <- counts[,treated$id]
treated.mean <- rowMeans(treated.counts)
head(treated.mean)
```
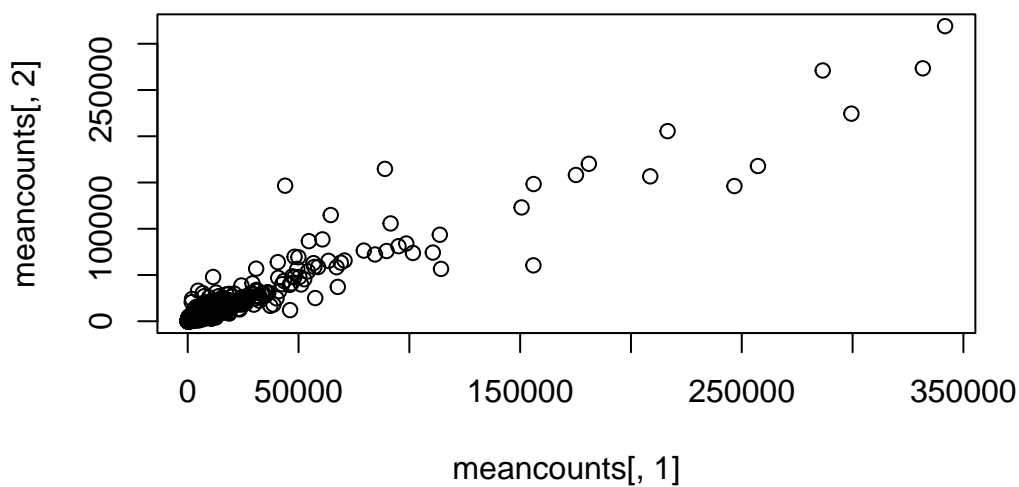
```
ENSG00000000003 ENSG00000000005 ENSG00000000419 ENSG00000000457 ENSG00000000460
        658.00            0.00          546.00          316.50           78.75
ENSG00000000938
          0.00
```

Store these results together in a dataframe called mean counts.

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

Lets make a plot to explore the results a little.

```r
plot(meancounts[,1], meancounts[,2])
```
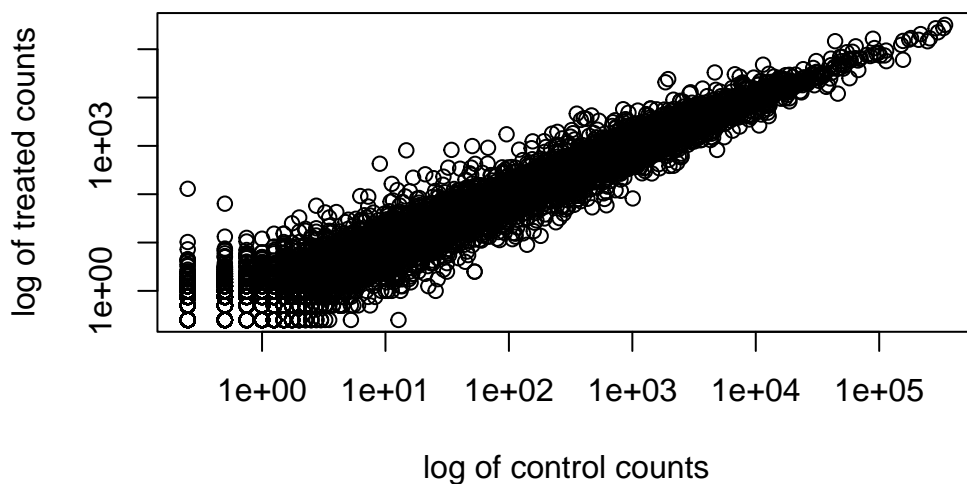


Make log-log plot to draw out this skewed data and see what is going on.

```r
plot(meancounts[,1], meancounts[,2], log="xy", xlab = "log of control counts",
     ylab = "log of treated counts")
```

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



Log2 transfirmation has a nice property, where no change will make the log2 value zero, doubling will lead log2 to be 1 and halving will lead it to be -1.

Add log2 fold change column to our results so far.

```r
meancounts$log2fc <- log2(meancounts$treated.mean/meancounts$control.mean)

# To get rid of NaN:

# says where the count is 0
zero.vals <- which(meancounts[,1:2]==0, arr.ind=TRUE)
```

```
to.rm <- unique(zero.vals[,1])
# removes genes with 0 counts
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 remaining?

There are 21817 genes remaining.

## Use fold change to see up and down regulated genes.

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

## DESeq2 analysis

```
#load up DESeq2
library(DESeq2)

dds <- DESeqDataSetFromMatrix(countData=counts,
                              colData=metadata,
                              design=~dex) # design - which col to look at
```

```
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>
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 can get some basic summary tallies using the `summary()` function.

```
  summary(res, alpha = 0.05)
```
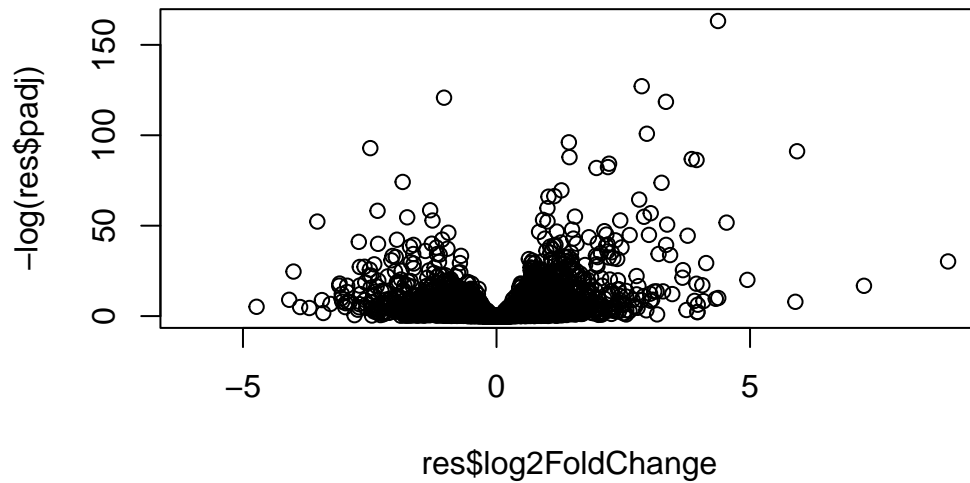
```
out of 25258 with nonzero total read count
adjusted p-value < 0.05
LFC > 0 (up)       : 1242, 4.9%
LFC < 0 (down)     : 939, 3.7%
outliers [1]       : 142, 0.56%
low counts [2]     : 9971, 39%
(mean count < 10)
[1] see 'cooksCutoff' argument of ?results
[2] see 'independentFiltering' argument of ?results
```

## Volcano Plot

Let's make a summary plot of our results.

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



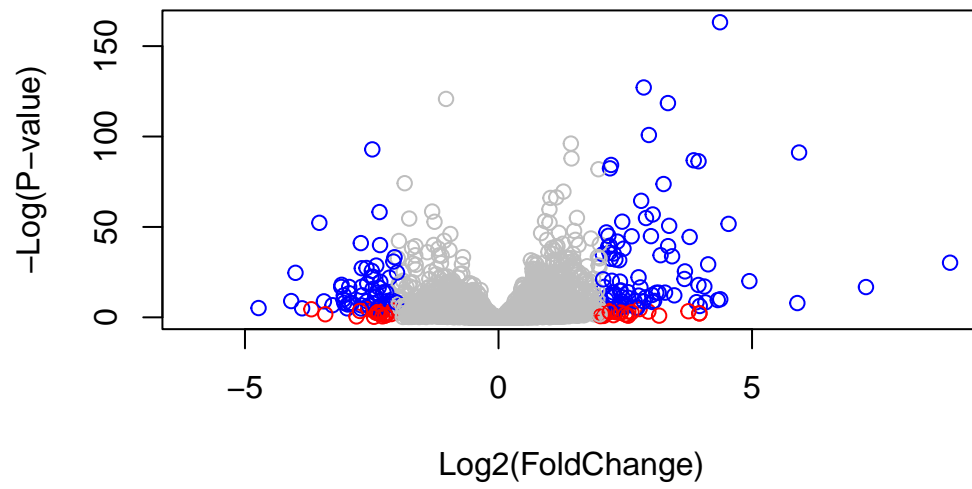Let's add colors:

```
# 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)" )
```

Finish for today by saving our results.

```
# write.csv(res, file = "DESeq2_results.csv")
```