

Week 9 DESeq analysis

Ralph Goguanco

3/4/2022

This week we are looking at differential expression analysis.

The data for this hands-on session comes from a published RNA-seq experiment where airway smooth muscle cells were treated with dexamethasone, a synthetic glucocorticoid steroid with anti-inflammatory effects (Himes et al. 2014).

Import/Read the data from Himes et al.

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

Lets have a wee peak at this data

```
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
```

Sanity check on corespondence of counts and metadata

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

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

Q1. How many genes are in this dataset?

There are 38694 genes in this dataset.

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

There are 4 control cell lines in this dataset.

Extract and summarize the control samples

To find out where the control samples are we need the metadata

```
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
```

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

Instead of using rowSums/4, we can just use rowMeans

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)

Extract and summarize the treated (i.e. drug) samples

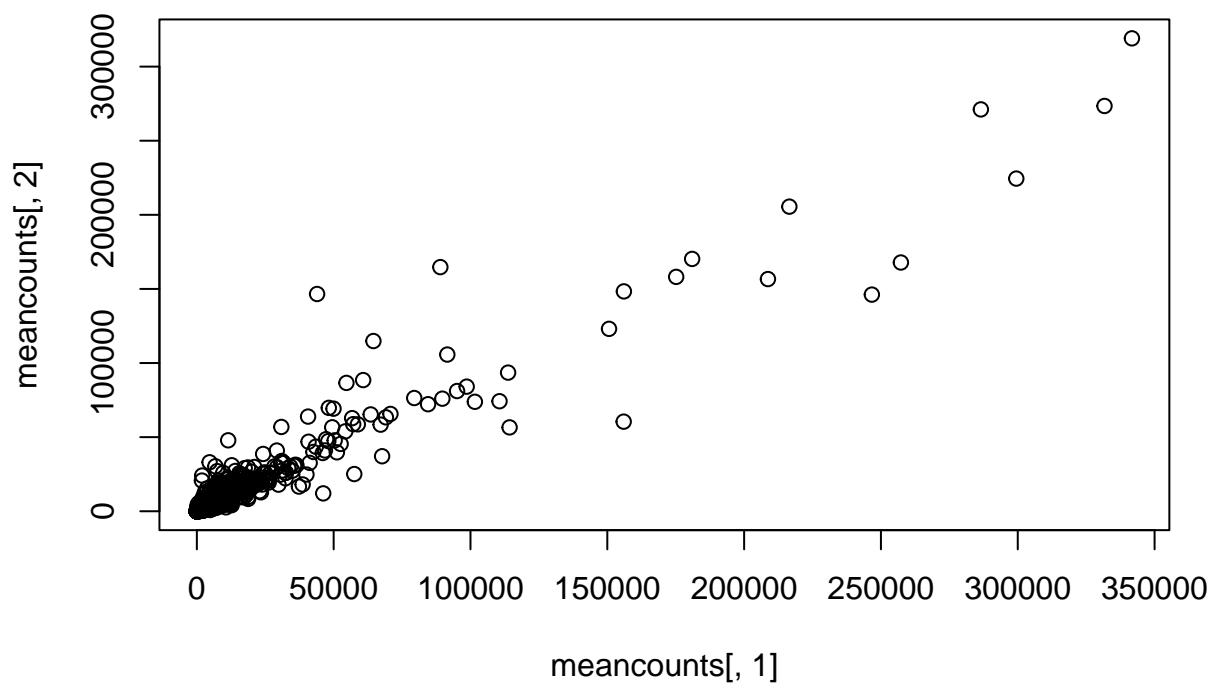
```
treated <- metadata[metadata$dex == "treated", ]  
treated.counts <- counts[, treated$id]  
treated.mean <- rowMeans(treated.counts)
```

Store these results together in a new data frame called `meancounts` Lets make a plot to explore the results a little

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

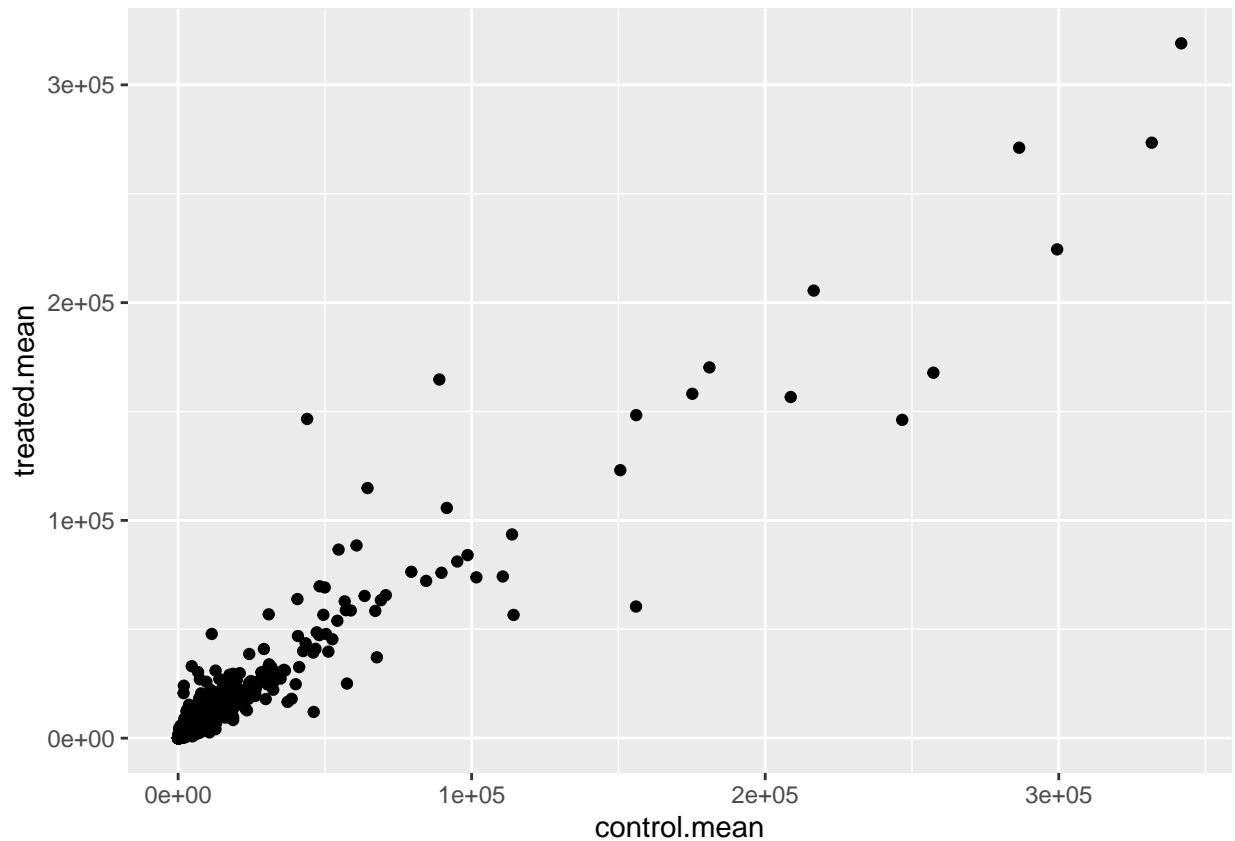
Lets make a plot to explore the results a little

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



```
library(ggplot2)

ggplot(meancounts) +
  aes(control.mean, treated.mean) +
  geom_point()
```

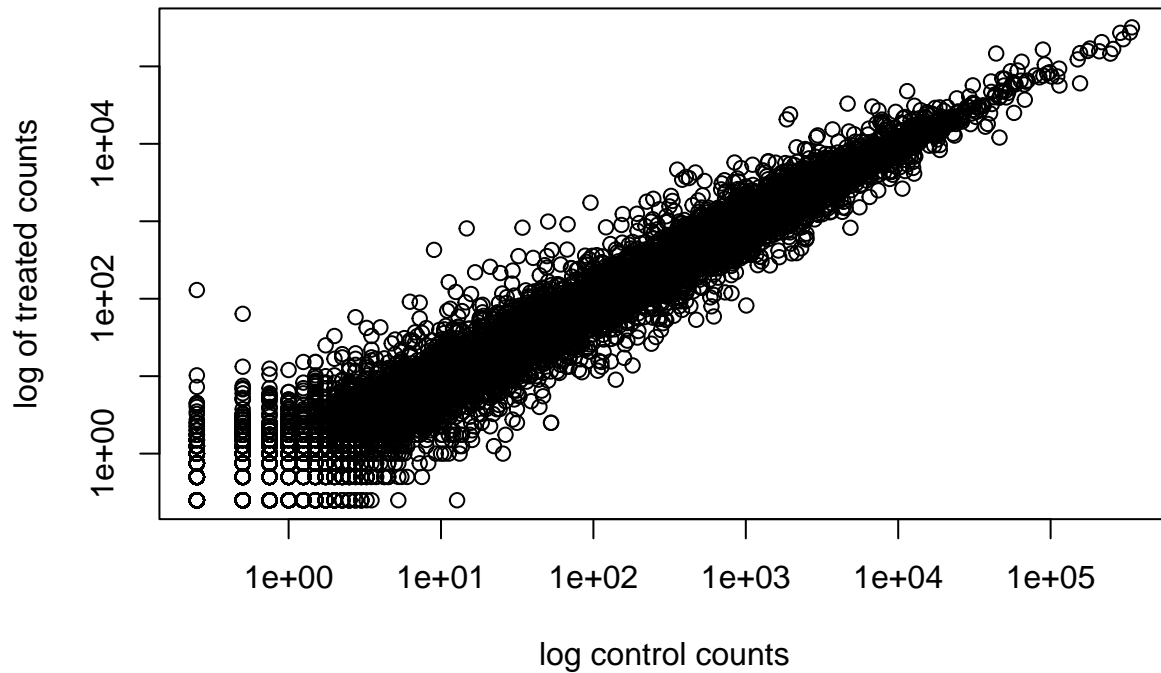


We will make a log-log plot to draw out this skewed data and see what is going on

```
plot(meancounts[,1], meancounts[,2], log="xy",
     xlab="log 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
```



We often log2 transformations when dealing with this sort of data.

```
log2(80/20)
```

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

This log2 transformation has this nice property where if there is no change the log2 value will be zero and if it double the log2 value will be 1 and if halved it will be -1.

So lets add a log2 fold change column to our results so far

```
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
```

We need to get rid of zero count genes that we can not say anything about

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

How many genes are remaining?

```
nrow(mycounts)
```

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

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

Acommon 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

```
sum(mycounts$log2fc > 2)
```

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

and d own regulated

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

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

DESeq2 analysis

Let's do this the right way. DESeq2 is an R package specifically for analyzing count-based NGS data like RNA-seq.

```
# load up DESeq2
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

##
## Attaching package: 'IRanges'

## 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: '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, 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

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 <- 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
```

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

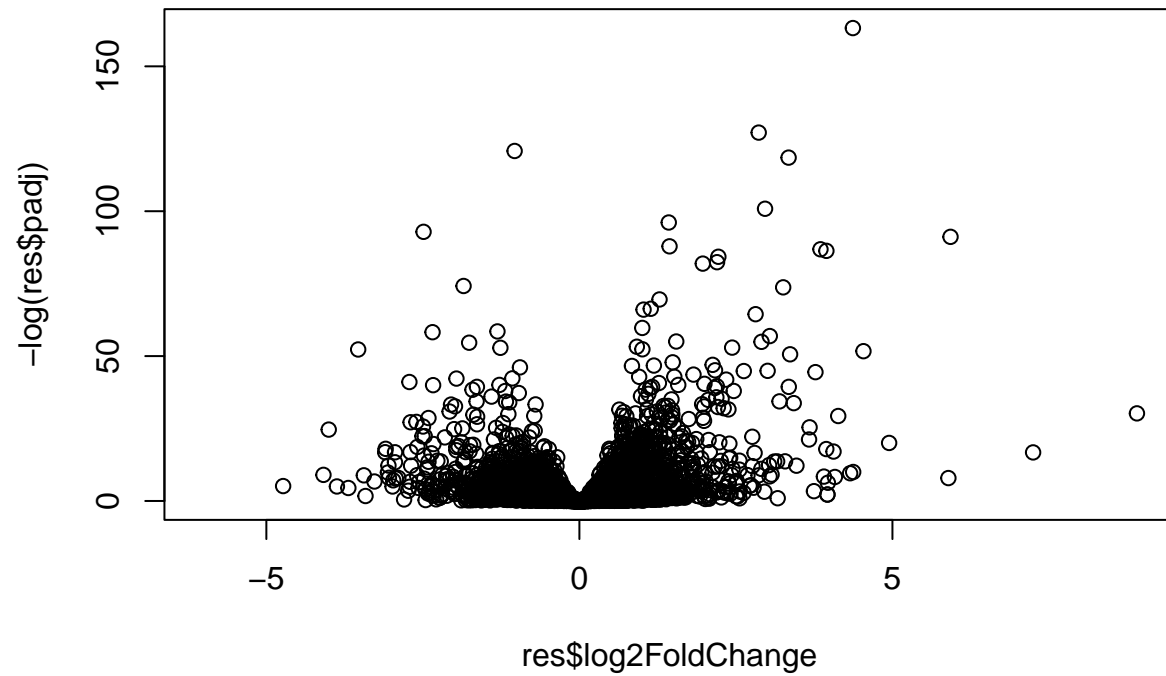
```
summary(res, alpha=0.05)
```

```
##
## 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
```

Make a summary plot of our results.

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



Finish for today by saving our results

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