# Class 12: Transcriptomics and the analysis of RNA-Seq data

Pagna Hout

5/12/23

## 1. Bioconductor and DESeq2 setup

```r
library(BiocManager)
```

Bioconductor version '3.16' is out-of-date; the current release version '3.17'
  is available with R version '4.3'; see https://bioconductor.org/install

```r
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 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

## 2. Import countData and colData

```r
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
```

```
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
```

```
View(counts)
View(metadata)
```

**Q1: How many genes are in this dataset?**

```
nrow(counts)
```

```
[1] 38694
```

There are 38694 genes in this dataset.

**Q2: How many 'control' cell lines do we have?**

```
table(metadata$dex)
```

4

```
control treated
      4       4
```

There are 4 "control" cell lines in the dataset.

## 3. Toy differential gene expression

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

```
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
```

**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`)**

```
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
```

Combining meancount data for bookeeping purposes
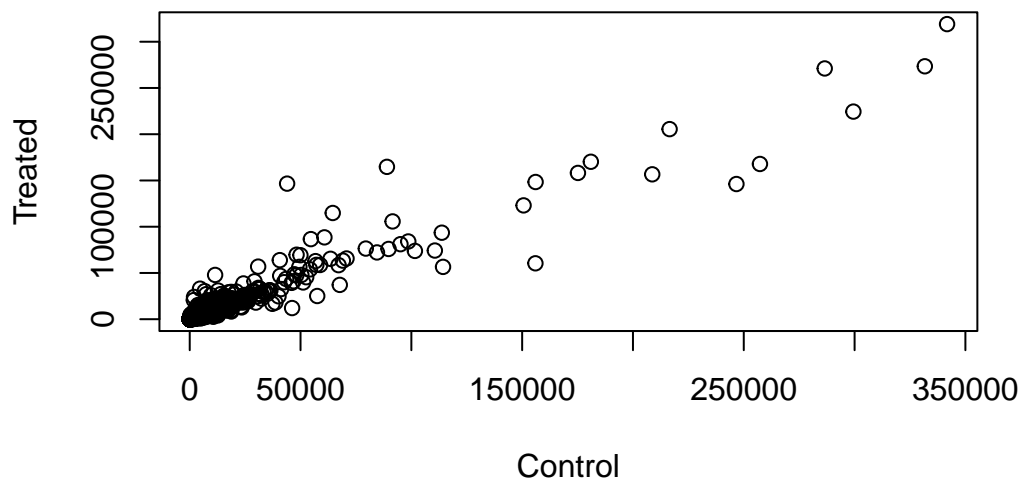
```r
meancounts <- data.frame(control.mean, treated.mean)

colSums(meancounts)
```

```
control.mean treated.mean
    23005324     22196524
```

**Q5 (a): Create a scatter plot showing the mean of the treated samples against the mean of the control samples.**
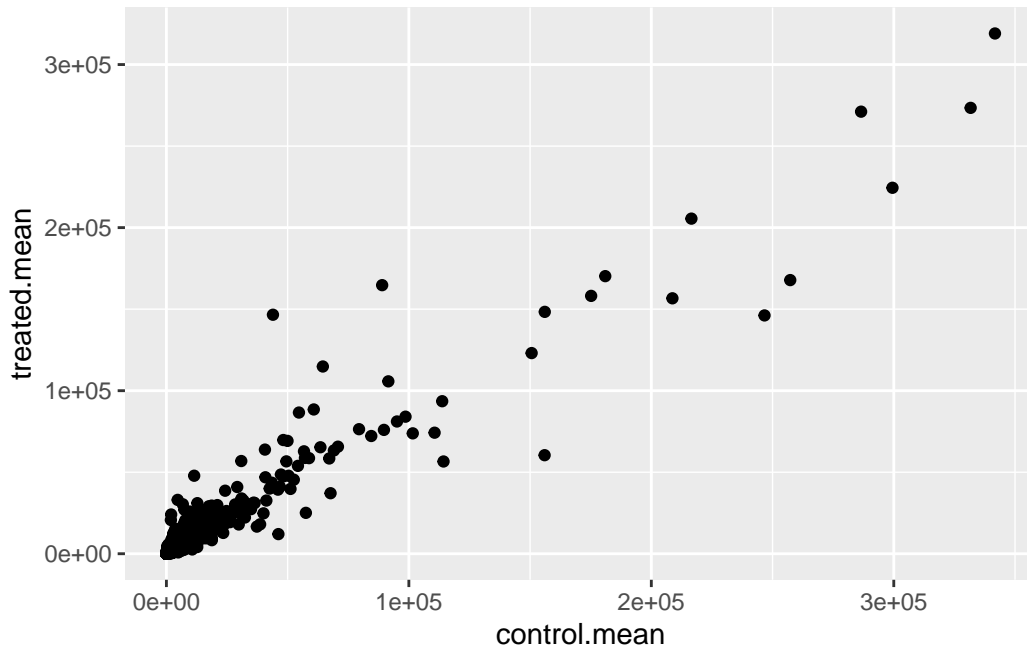
```r
plot(meancounts, xlab="Control", ylab="Treated")
```

**Q5 (b): You could also use the ggplot2 package to make this figure producing the plot below. What geom_() function would you use for this plot?**

```
library(ggplot2)

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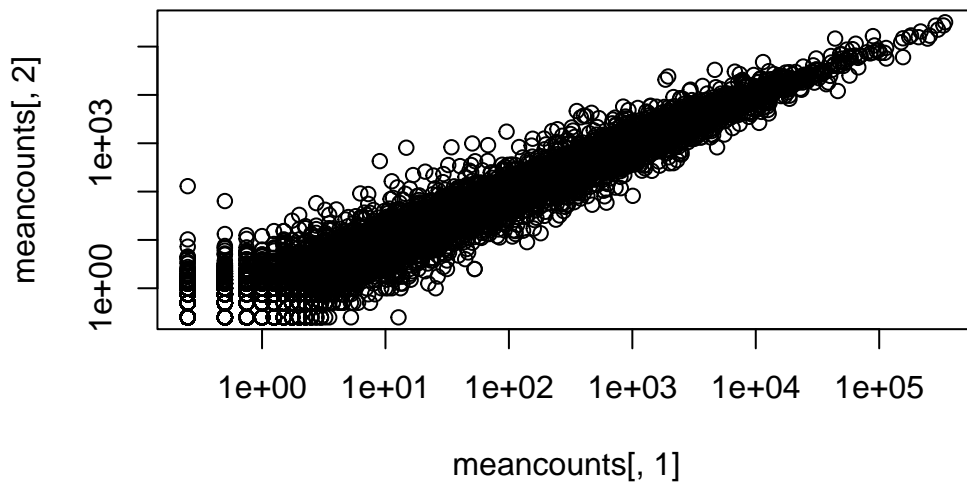


For this plot, I would use the geom_point () function.

**Q6: Try plotting both axes on a log scale. What is the argument to plot() that allows you to do this?**

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

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

The log argument allows us to plot both axes on a log scale.

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

```
                control.mean treated.mean        log2fc
ENSG00000000003       900.75       658.00 -0.45303916
ENSG00000000005         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
```

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

to.rm <- unique(zero.vals[,1])
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
```

**Q7: What is the purpose of the `arr.ind` argument in the which() function call above? Why would we then take the first column of the output and need to call the unique() function?**

The purpose of the arr.ind argument is to tell us which genes (rows) and samples (columns) have zero counts.

The purpose of the unique function is to prevent R to count any row twice if it has zer entries in both samples.

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

**Q8: Using the `up.ind` vector above can you determine how many up regulated genes we have at the greater than 2 fc level?**

```
sum(up.ind)
```

```
[1] 250
```

There are 250 up regulated genes greater than 2 fc level.

**Q9: Using the `down.ind` vector above can you determine how many down regulated genes we have at the greater than 2 fc level?**

```
sum(down.ind)
```

```
[1] 367
```

There are 367 down regulated genes greater than 2 fc level.

**Q10: Do you trust these results? Why or why not?**

These results might not be reliable because we have not determined whether or not the differences are statistically significant (for example, using p-values).