

Class 12

Emmanuel Robles

Import the data

We need two things for this analysis: - countdata (counts for every transcript/gene in each experiment) - coldata (metadata that describes the experimental setup)

```
countData <- read.csv("airway_scaledcounts.csv", row.names=1)
head(countData)
```

	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

```
metadata <- read.csv("airway_metadata.csv")
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?

```
nrow(countData)
```

```
[1] 38694
```

Q2. How many “control” cell lines do we have?

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

```

control treated
      4      4

```

Another method

```
sum(metadata$dex == "control")
```

```
[1] 4
```

- Step 1 Calculate the mean of the control sample Calculate the mean of the treated samples

(a) We need to find which columns in countData are “control” samples.

- Look in the metadata (a.k.a. colData), \$dex column

```
control.ind <- metadata$dex == "control"
```

(b) Extract all the control columns from countData and call it control.counts

```
control.counts <- countData[,control.ind]
```

(c) Calculate the mean values across the rows of control and calculate the mean count values

```
control.means <- rowMeans(control.counts)
head(control.means)
```

```
ENSG000000000003 ENSG000000000005 ENSG000000000419 ENSG000000000457 ENSG000000000460
          900.75           0.00           520.50           339.75           97.25
ENSG000000000938
          0.75
```

- Step 2 Calculate the mean of the treated samples...

```
treated.ind <- metadata$dex == "treated"
treated.counts <- countData[,treated.ind]
treated.means <- rowMeans(treated.counts)
head(treated.means)
```

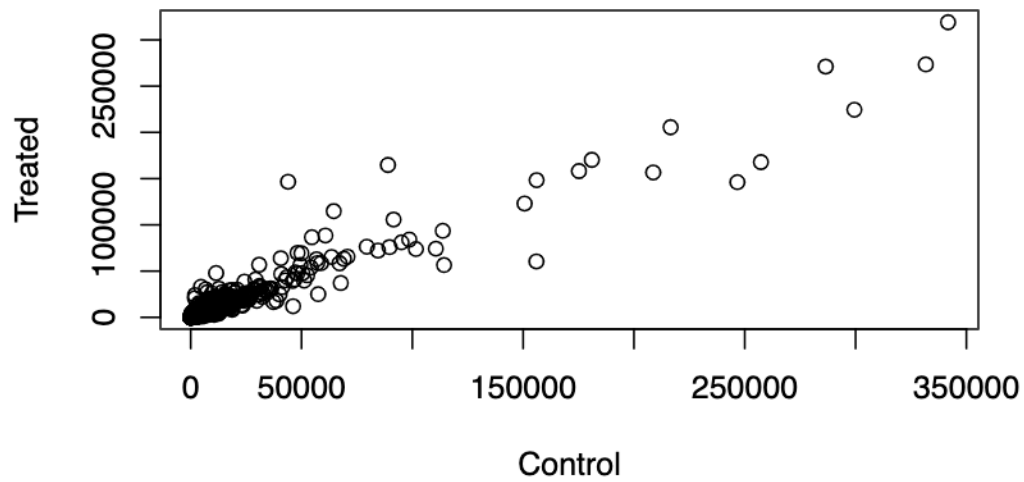
```
ENSG000000000003 ENSG000000000005 ENSG000000000419 ENSG000000000457 ENSG000000000460
          658.00           0.00           546.00           316.50           78.75
ENSG000000000938
          0.00
```

We now have control and treated mean count values, For ease of keeping I will combine these vectors into a new data.frame `meancounts`.

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

	control.means	treated.means
ENSG000000000003	900.75	658.00
ENSG000000000005	0.00	0.00
ENSG000000000419	520.50	546.00
ENSG000000000457	339.75	316.50
ENSG000000000460	97.25	78.75
ENSG000000000938	0.75	0.00

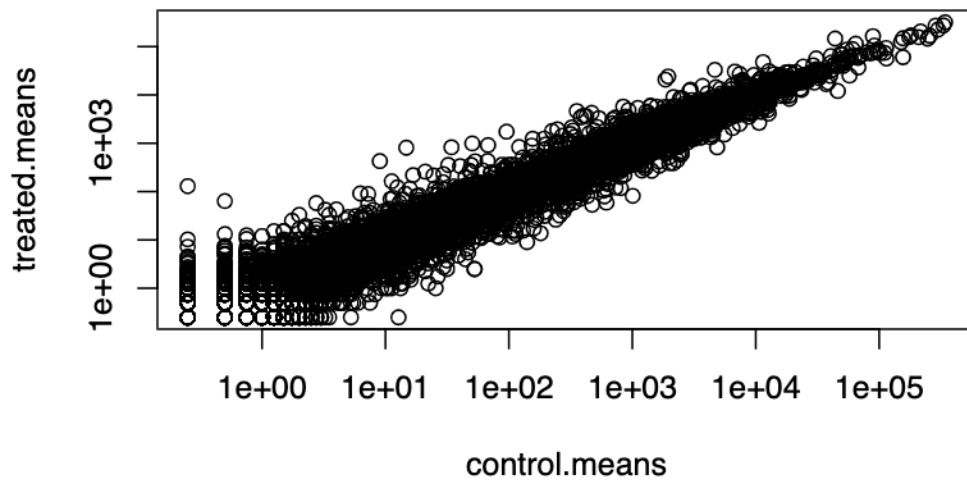
```
plot(meancounts[,1],meancounts[,2], xlab="Control", ylab="Treated")
```



```
plot(meancounts, 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



We use log transforms for skewed data such as this ans because we really care most about relative changes in magnitude.

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

```
[1] 0
```

If I have half the amount I will have log2 fold-change of -1

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

```
[1] -1
```

If I have double the amount (20 compared to 10), change of +1

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

```
[1] 1
```

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

Q. How many genes are upregulated at the common threshold of +2 log2FC values?

```
sum(meancounts$log2fc >= 2, na.rm=TRUE )
```

```
[1] 1910
```

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, 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, 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 = countData,  
                              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

To run the analysis I can now use the main DESeq2 function called DESeq() with dds as
input.

```
dds <- DESeq(dds)
```

estimating size factors

estimating dispersions

gene-wise dispersion estimates

mean-dispersion relationship

final dispersion estimates

fitting model and testing

To get the results back from this object we can use the `results()` function from the package.

```
res <- results(dds)
head(res)
```

log2 fold change (MLE): dex treated vs control

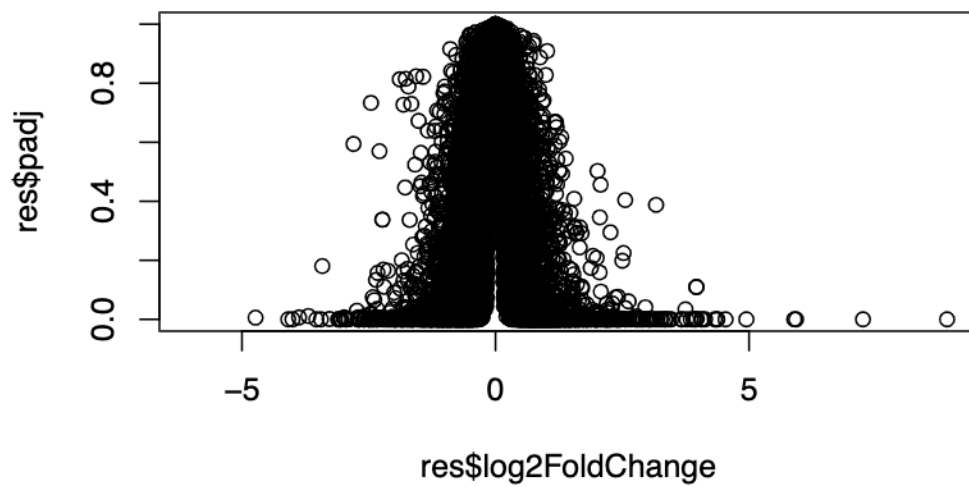
Wald test p-value: dex treated vs control

DataFrame with 6 rows and 6 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
ENSG0000000000419	520.134160	0.2061078	0.101059	2.039475	0.0414026
ENSG0000000000457	322.664844	0.0245269	0.145145	0.168982	0.8658106
ENSG0000000000460	87.682625	-0.1471420	0.257007	-0.572521	0.5669691
ENSG0000000000938	0.319167	-1.7322890	3.493601	-0.495846	0.6200029
	padj				
	<numeric>				
ENSG0000000000003	0.163035				
ENSG0000000000005	NA				
ENSG00000000000419	0.176032				
ENSG00000000000457	0.961694				
ENSG00000000000460	0.815849				
ENSG00000000000938	NA				

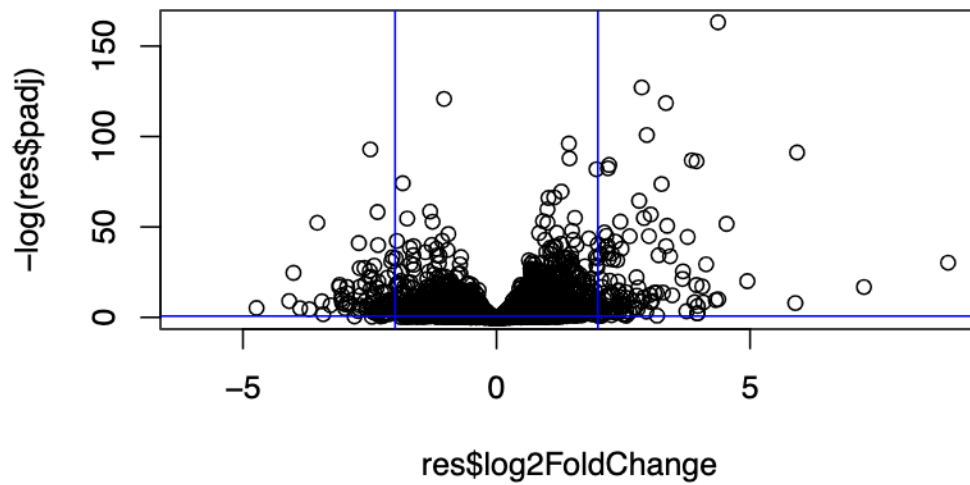
Let's make a final (for today) plot of log2 fold-change vs the adjusted p-value.

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



It is the low p-values that we care about and these are lost in to our skewed plot above.

```
plot(res$log2FoldChange, -log(res$padj))
abline(v=c(+2, -2), col="blue")
abline(h= -log(0.5), col="blue")
```



Finally we can make a color vector to use in the plot to highlight the genes we care about.

```
mycols <- rep("grey", nrow(res))
mycols[abs(res$log2FoldChange) >= 2] <- "red"
mycols[res$padj] <- "grey"
#mycols

plot(res$log2FoldChange, -log(res$padj), col= mycols)
abline(v=c(+2, -2), col="blue")
abline(h= -log(0.5), col="blue")
```

