

Class 13: RNASeq with DESeq2

Jessica PID: A15647602

Today we will work with some bulk RNASeq data from Himes et al. where airway smooth muscle (asm) cells were treated with dexamethasone (dex), a synthetic glucocorticoid steroid with anti-inflammatory effects.

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

```
#head(counts)
```

Q1. How many transcripts/ genes are in the `counts` object?

```
nrow(counts)
```

```
[1] 38694
```

Q2. How many “control” samples are there?

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

```
[1] 4
```

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

```
control treated
      4      4
```

I want to compare “control” vs. “treated”

1. Let's split the `counts` into `control.counts` and `treated.counts`

```
metadata$dex == "control"
```

```
[1] TRUE FALSE TRUE FALSE TRUE FALSE TRUE FALSE
```

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

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

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

Syntax with df[ROWS, COLS]

```
control.counts <- counts[ , control.inds]  
#control.counts
```

```
treated.inds <- metadata$dex == "treated"  
treated.counts <- counts[ , treated.inds]  
#treated.counts
```

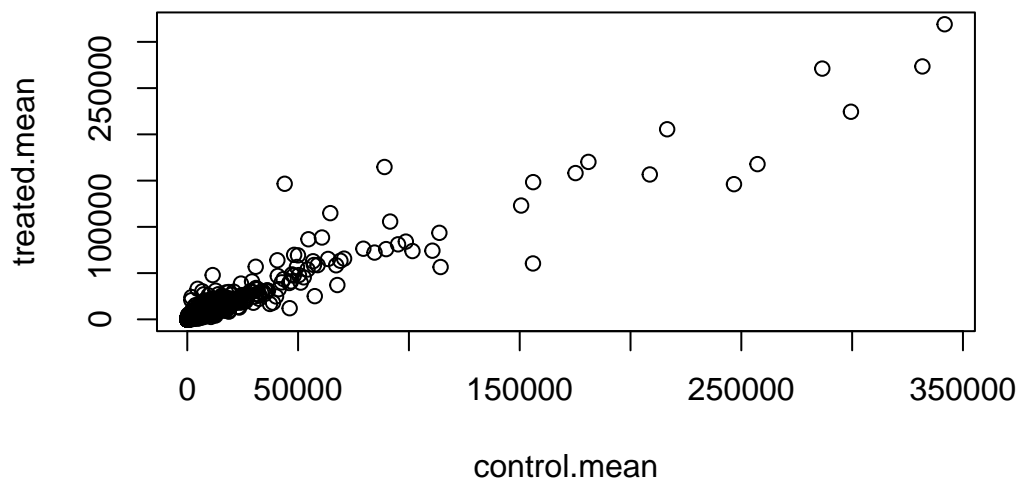
2. Let's find the mean count per gene for "control" and "treated" - then we can compare these :). Let's call it `control.mean` and `treated.mean`.

I can use the `apply()` function to apply `mean()` over the rows or columns of any data frame.

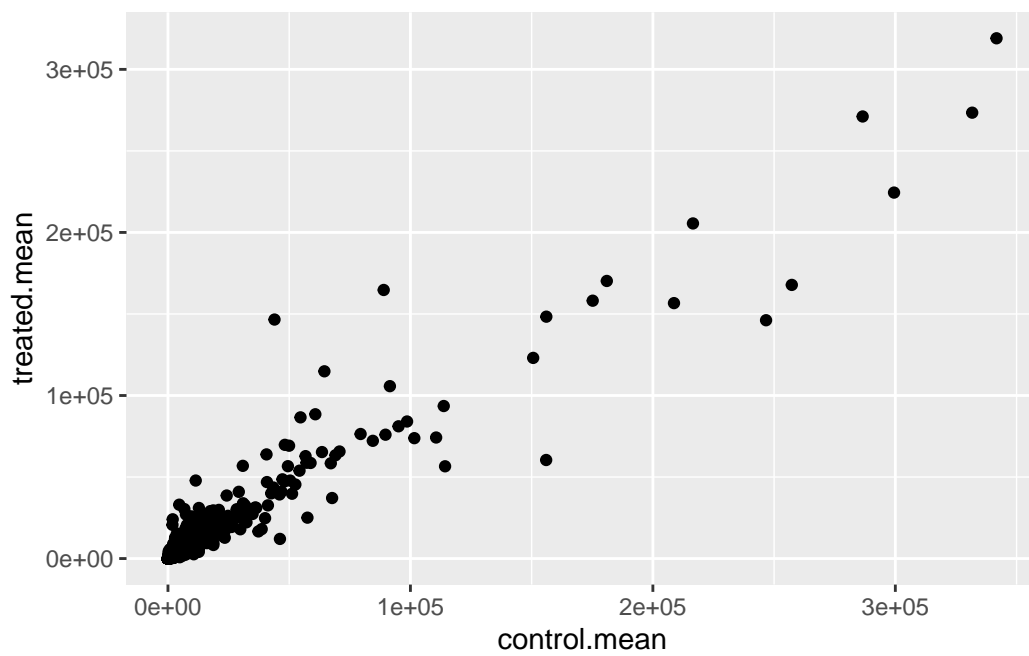
```
control.mean <- apply(control.counts, 1, mean)  
treated.mean <- apply(treated.counts, 1, mean)
```

Put these together for ease of book-keeping

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



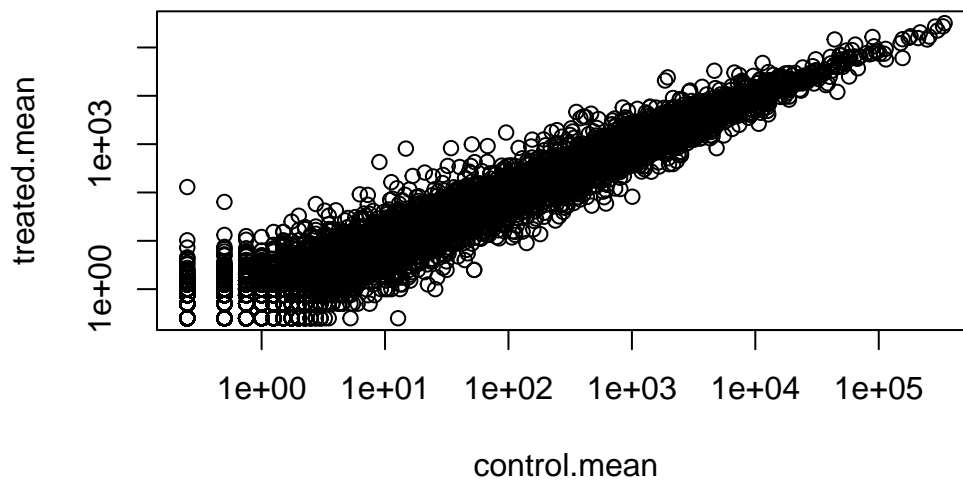
```
library(ggplot2)
ggplot(meancounts, aes(x = control.mean, y= treated.mean)) + geom_point()
```



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



```
log2(40/10)
```

```
[1] 2
```

Let's calculate the log2 fold change and add it to our table `mean.counts`

```
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

FILTER OUT ALL GENES WITH ZERO COUNTS IN EITHER CONTROL OR TREATED

```
to.rm <- rowSums(meancounts[,1:2] == 0) > 0
mycounts <- meancounts[!to.rm, ]
```

```
nrow(mycounts)
```

```
[1] 21817
```

Q. How many “down” regulated genes do we have at the common log2 fold change value of -2...

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

```
[1] 367
```

Q. How many “up” at log2FC > +2

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

```
[1] 250
```

We are missing the stats! ##DESeq Analysis

```
library(DESeq2)
```

DESeq, like many BioConductor packages, wants our input data in a very specific format.

```
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

The main function in DESeq2 is called DESeq()

```
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)
```

```
head(res)
```

log2 fold change (MLE): dex treated vs control

Wald test p-value: dex treated vs control

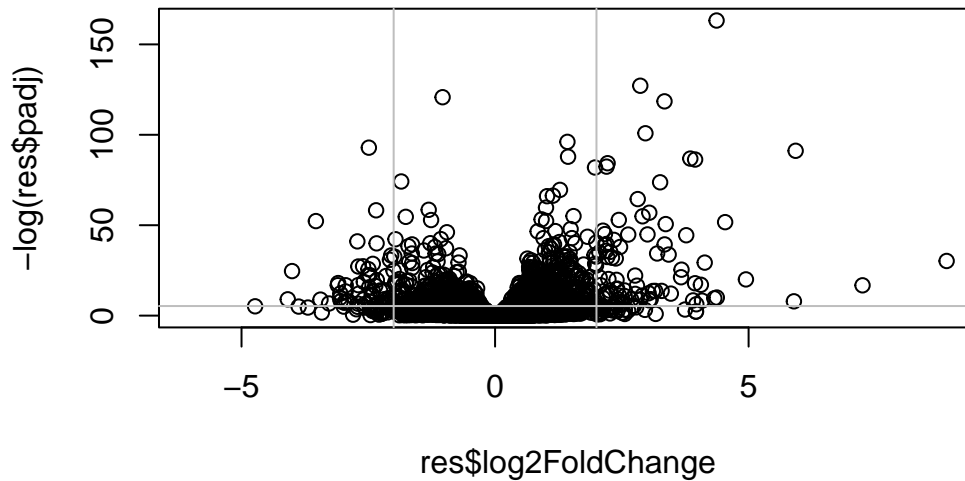
DataFrame with 6 rows and 6 columns

	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
ENSG000000000419	520.134160	0.2061078	0.101059	2.039475	0.0414026
ENSG000000000457	322.664844	0.0245269	0.145145	0.168982	0.8658106

ENSG00000000460	87.682625	-0.1471420	0.257007	-0.572521	0.5669691
ENSG00000000938	0.319167	-1.7322890	3.493601	-0.495846	0.6200029
	padj				
	<numeric>				
ENSG00000000003	0.163035				
ENSG00000000005	NA				
ENSG00000000419	0.176032				
ENSG00000000457	0.961694				
ENSG00000000460	0.815849				
ENSG00000000938	NA				

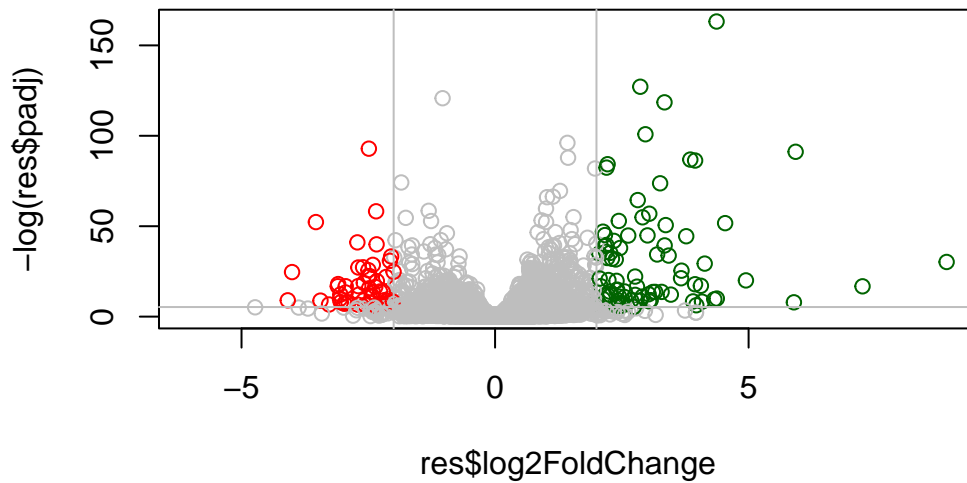
Figure volcano plot for logFC vs. P-value

```
plot(res$log2FoldChange, -log(res$padj))
abline(v=c(-2,2), col = "gray")
abline(h=-log(0.005), col = "gray")
```



```
mycols <- rep("gray", nrow(res))
mycols[res$log2FoldChange > 2] <- "darkgreen"
mycols[res$log2FoldChange < -2] <- "red"
mycols[res$padj > 0.005] <- "gray"
```

```
plot(res$log2FoldChange, -log(res$padj), col = mycols)
abline(v=c(-2,2), col = "gray")
abline(h=-log(0.005), col = "gray")
```



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

Gene Annotation

```
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


```

ENSG000000000938    0.319167    -1.7322890    3.493601    -0.495846    0.6200029
                    padj
                    <numeric>
ENSG000000000003    0.163035
ENSG000000000005           NA
ENSG000000000419    0.176032
ENSG000000000457    0.961694
ENSG000000000460    0.815849
ENSG000000000938           NA

```

```

library("AnnotationDbi")
library("org.Hs.eg.db")

```

```

columns(org.Hs.eg.db)

```

```

[1] "ACCNUM"      "ALIAS"       "ENSEMBL"     "ENSEMBLPROT" "ENSEMBLTRANS"
[6] "ENTREZID"    "ENZYME"      "EVIDENCE"    "EVIDENCEALL"  "GENENAME"
[11] "GENETYPE"    "GO"          "GOALL"       "IPI"          "MAP"
[16] "OMIM"        "ONTOLOGY"    "ONTOLOGYALL" "PATH"         "PFAM"
[21] "PMID"        "PROSITE"     "REFSEQ"      "SYMBOL"       "UCSCKG"
[26] "UNIPROT"

```

```

res$symbols <- mapIds(org.Hs.eg.db,
                      keys = row.names(res),
                      keytype = "ENSEMBL",
                      column = "SYMBOL",
                      multiVals = "first")

```

'select()' returned 1:many mapping between keys and columns

```

head(res)

```

log2 fold change (MLE): dex treated vs control

Wald test p-value: dex treated vs control

DataFrame with 6 rows and 7 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	symbols
	<numeric>	<character>
ENSG000000000003	0.163035	TSPAN6
ENSG000000000005	NA	TNMD
ENSG0000000000419	0.176032	DPM1
ENSG0000000000457	0.961694	SCYL3
ENSG0000000000460	0.815849	FIRRM
ENSG0000000000938	NA	FGR

##Pathway Analysis

```
library(pathview)
```

```
#####
Pathview is an open source software package distributed under GNU General
Public License version 3 (GPLv3). Details of GPLv3 is available at
http://www.gnu.org/licenses/gpl-3.0.html. Particullary, users are required to
formally cite the original Pathview paper (not just mention it) in publications
or products. For details, do citation("pathview") within R.
```

```
The pathview downloads and uses KEGG data. Non-academic uses may require a KEGG
license agreement (details at http://www.kegg.jp/kegg/legal.html).
```

```
#####
```

```
library(gage)
```

```
library(gageData)
```

```
data(kegg.sets.hs)
```

```
head(kegg.sets.hs, 2)
```

```
$`hsa00232 Caffeine metabolism`
```

```
[1] "10" "1544" "1548" "1549" "1553" "7498" "9"
```

```
$`hsa00983 Drug metabolism - other enzymes`
```

```
[1] "10" "1066" "10720" "10941" "151531" "1548" "1549" "1551"  
[9] "1553" "1576" "1577" "1806" "1807" "1890" "221223" "2990"  
[17] "3251" "3614" "3615" "3704" "51733" "54490" "54575" "54576"  
[25] "54577" "54578" "54579" "54600" "54657" "54658" "54659" "54963"  
[33] "574537" "64816" "7083" "7084" "7172" "7363" "7364" "7365"  
[41] "7366" "7367" "7371" "7372" "7378" "7498" "79799" "83549"  
[49] "8824" "8833" "9" "978"
```

```
res$entrez <- mapIds(org.Hs.eg.db,  
  keys = row.names(res),  
  keytype = "ENSEMBL",  
  column = "ENTREZID",  
  multiVals = "first")
```

'select()' returned 1:many mapping between keys and columns

I can use **gage** to overlap with known KEGG pathways.

```
foldchanges <- res$log2FoldChange  
names(foldchanges) <- res$entrez  
head(foldchanges)
```

```
      7105      64102      8813      57147      55732      2268  
-0.35070302      NA  0.20610777  0.02452695 -0.14714205 -1.73228897
```

```
# Get the results  
keggres = gage(foldchanges, gsets=kegg.sets.hs)
```

```
attributes(keggres)
```

```
$names
```

```
[1] "greater" "less" "stats"
```

```
# Look at the first three down (less) pathways
head(keggres$less, 3)
```

		p.geomean	stat.mean	p.val
hsa05332	Graft-versus-host disease	0.0004250461	-3.473346	0.0004250461
hsa04940	Type I diabetes mellitus	0.0017820293	-3.002352	0.0017820293
hsa05310	Asthma	0.0020045888	-3.009050	0.0020045888

		q.val	set.size	exp1
hsa05332	Graft-versus-host disease	0.09053483	40	0.0004250461
hsa04940	Type I diabetes mellitus	0.14232581	42	0.0017820293
hsa05310	Asthma	0.14232581	29	0.0020045888

```
pathview(gene.data=foldchanges, pathway.id="hsa05310")
```

'select()' returned 1:1 mapping between keys and columns

Info: Working in directory /Users/jessicaraygoza/Desktop/Graduate School/BGGN 213/Class 13

Info: Writing image file hsa05310.pathview.png

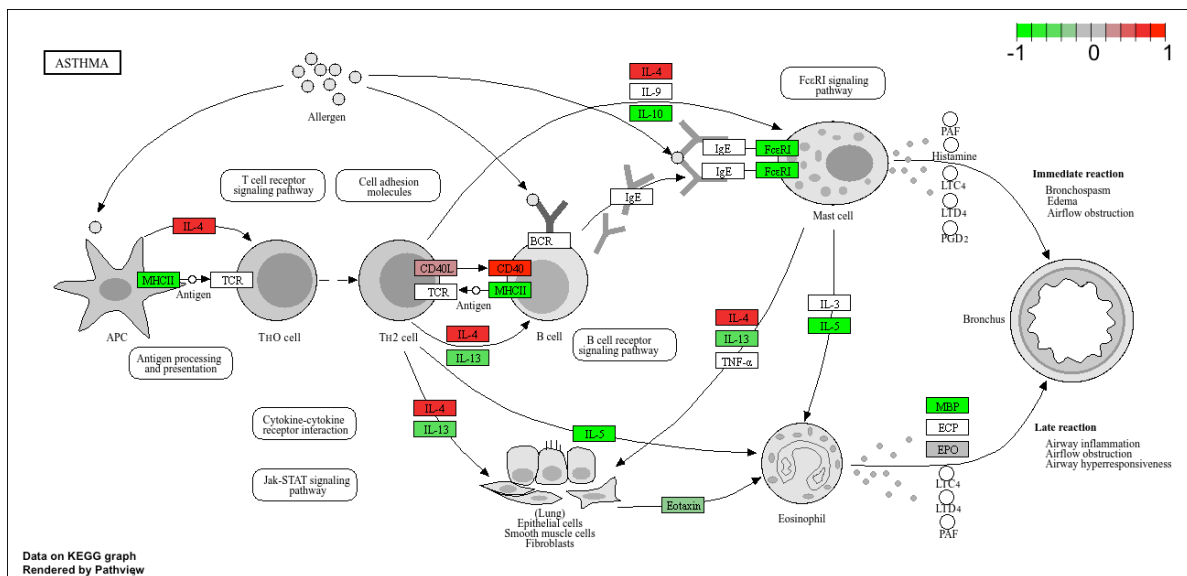


Figure 1: A pathway figure