

# Class 13: RNA Seq (pt. 1)

Amanda Wilpitz | A17463962

Today we will analyze data 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 countDATA and colDATA

There are two datasets I need to import/read:

- `countData` the transcript counts per gene (rows) in the different experiments
- `colData` information about the columns (i.e. experiments) in `countData`

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

We can have a peak at these with `head()`

```
head(counts)
```

	SRR1039508	SRR1039509	SRR1039512	SRR1039513	SRR1039516
ENSG000000000003	723	486	904	445	1170
ENSG000000000005	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		
ENSG000000000003	1097	806	604		
ENSG000000000005	0	0	0		
ENSG00000000419	781	417	509		
ENSG00000000457	447	330	324		
ENSG00000000460	94	102	74		
ENSG00000000938	0	0	0		

```
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
7 SRR1039520 control    N061011 GSM1275874
8 SRR1039521 treated    N061011 GSM1275875
```

Q1. How many genes are in this dataset?

```
nrow(counts)
```

```
[1] 38694
```

Q2. How many ‘control’ cell lines do we have?

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

```
[1] 4
```

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

```
control treated
        4       4
```

## Toy Differential Gene Expression

We can find the average (mean) count values per gene for all “control” experiments and compare it to the mean values for “treated”. If there is no difference, the drug didn’t change anything.

- Step 1. Extract all “control” columns from the `counts` data

```
control inds <- metadata$dex == "control"
control counts <- counts[, control inds]
```

```
dim(control.counts)
```

```
[1] 38694     4
```

```
head(control.counts)
```

	SRR1039508	SRR1039512	SRR1039516	SRR1039520
ENSG000000000003	723	904	1170	806
ENSG000000000005	0	0	0	0
ENSG000000000419	467	616	582	417
ENSG000000000457	347	364	318	330
ENSG000000000460	96	73	118	102
ENSG000000000938	0	1	2	0

- Step 2. Find the mean value for each gene in the **control** columns

```
control.mean <- rowSums(control.counts)/ncol(control.counts)
```

```
head(control.mean)
```

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

Q3. How would you make the above code in either approach more robust? Is there a function that could help here?

You can make the code mentioned in the lab sheet be “rowSums(control.counts)/ncol(control.counts)” instead of “rowSums(control.counts)/4”. With it being over 4, the 4 is hard coded in and if any experiments are added, it wouldn’t calculate it correctly. By having a code for the count of how many columns (experiments) there are, it would adjust according to new information.

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**)

- Step 3. Extract all “treated” columns from the **counts** data

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

- Step 4. Find the mean value for each gene in the **treated** columns

```
treated.mean <- rowSums(treated.counts)/ncol(treated.counts)  
  
head(treated.mean)
```

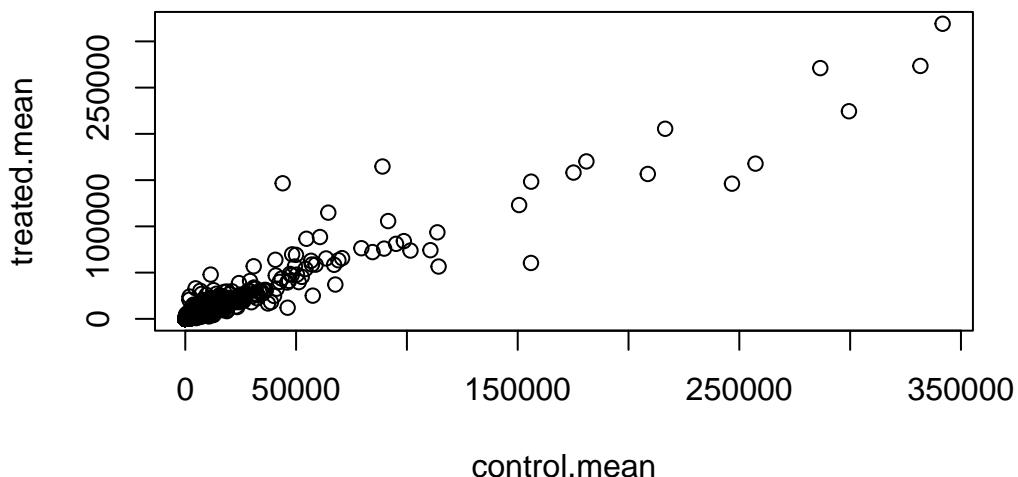
```
ENSG00000000003 ENSG00000000005 ENSG00000000419 ENSG00000000457 ENSG00000000460  
658.00 0.00 546.00 316.50 78.75  
ENSG00000000938  
0.00
```

Let's put these two mean values together for easy book-keeping

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

Q5 (a). Create a scatter plot showing the mean of the treated samples against the mean of the control samples. Your plot should look something like the following.

```
plot(meancounts)
```

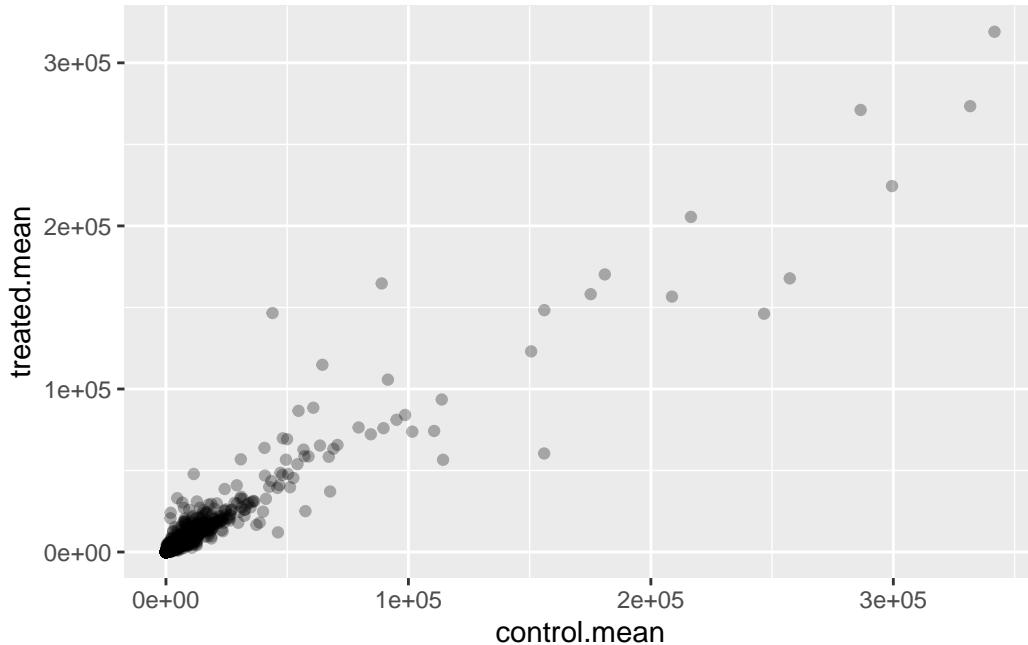


Q5 (b). You could also use the ggplot2 package to make this figure producing the plot below. What geom\_?() function would you use for this plot?

You would use `geom_point()`

```
library(ggplot2)

ggplot(meancounts, aes(control.mean, treated.mean)) +
  geom_point(alpha = 0.3)
```



```
#using alpha allows us to see that there is a lot of overplotting with data points overlapping
```

Any points on the diagonal shows that the drug had no effect.

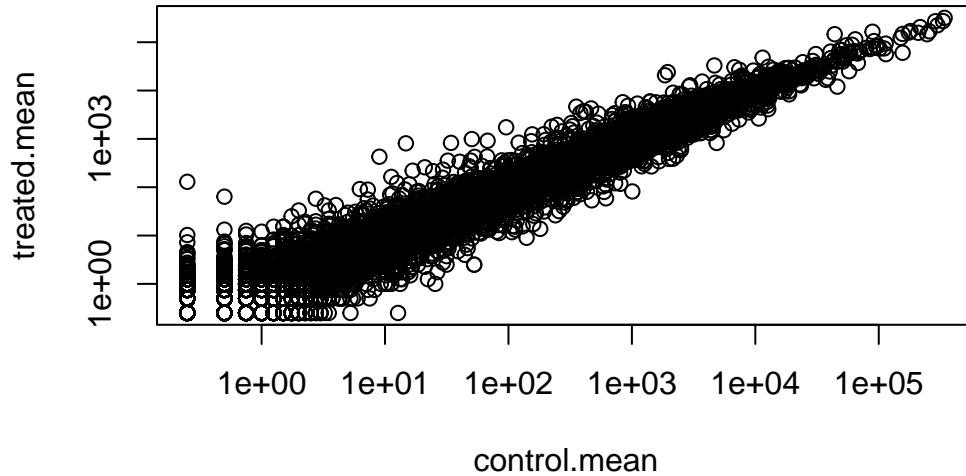
Whenever we see data that is so heavily skewed like this, we often log transform it so we can see what is going on more easily.

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

```
plot(meancounts, log="xy")
```

```
Warning in xy.coords(x, y, xlabel, ylabel, log): 15032 x values <= 0 omitted
from logarithmic plot
```

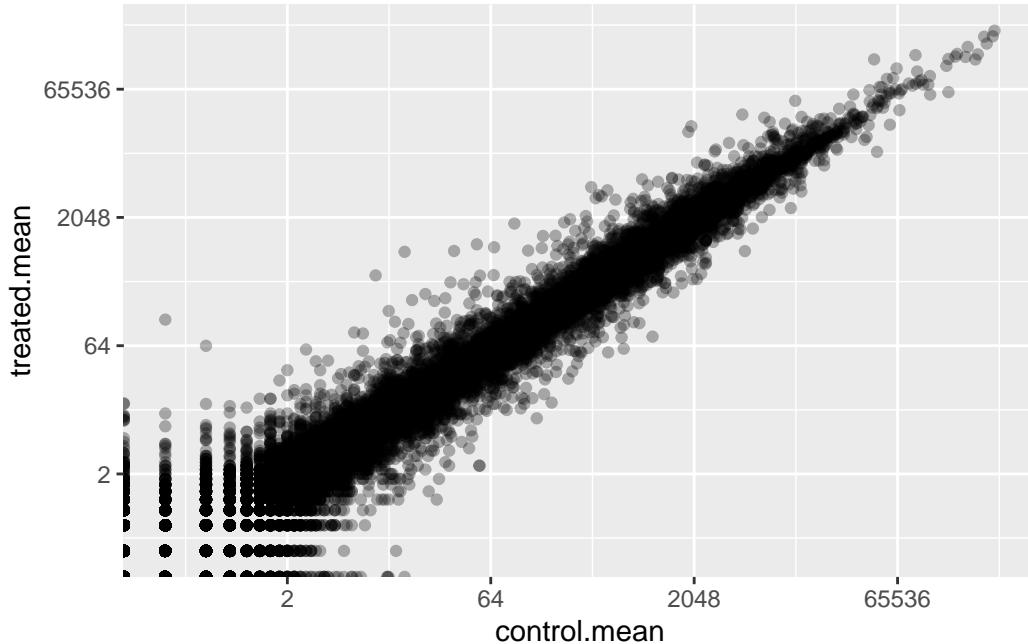
```
Warning in xy.coords(x, y, xlabel, ylabel, log): 15281 y values <= 0 omitted  
from logarithmic plot
```



```
ggplot(meancounts, aes(control.mean, treated.mean)) +  
  geom_point(alpha = 0.3) +  
  scale_x_continuous(trans="log2") +  
  scale_y_continuous(trans="log2")
```

```
Warning in scale_x_continuous(trans = "log2"): log-2 transformation introduced  
infinite values.
```

```
Warning in scale_y_continuous(trans = "log2"): log-2 transformation introduced  
infinite values.
```



We most often work in log<sub>2</sub> units as this makes the math easier. Let's have a play to see this.

```
#treated/control
log2(20/10)
```

```
[1] 1
```

```
#control/treated
log2(10/20)
```

```
[1] -1
```

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

```
[1] 0
```

We can now add “log<sub>2</sub> fold-change” values to our `meancounts` dataset.

```

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

NaN means not a number (dividing a zero by a zero and finding the log of it) -Inf means infinity

We need to filter out zero count genes - i.e. remove the rows (genes) that have a 0 value in either control or treated means.

```

to.keep <- rowSums(meancounts[, 1:2] == 0) == 0
mycounts <- meancounts[to.keep,]
nrow(mycounts)

```

[1] 21817

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

How many genes are “up” regulated at the common log2 fold-change threshold of +2?

```

up inds <- mycounts$log2fc >= 2
sum(up inds, na.rm = T)

```

[1] 314

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

How many genes are “down” regulated at the common log2 fold-change threshold of -2?

```
down.ind <- mycounts$log2fc <= -2  
sum(down.ind, na.rm = T)
```

[1] 485

Q10. Do you trust these results? Why or why not?

I do not trust these results to be different because they may be downregulated or upregulated in their magnitude of their differences, but according to a p-value may not be significant.

## DESeq2 Analysis

To do this the right way we need to consider the significance of the differences not just their magnitude.

```
library(DESeq2)
```

To use this package, it wants countData and colData in a 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

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

estimating size factors

estimating dispersions

gene-wise dispersion estimates

mean-dispersion relationship

```
final dispersion estimates
```

```
fitting model and testing
```

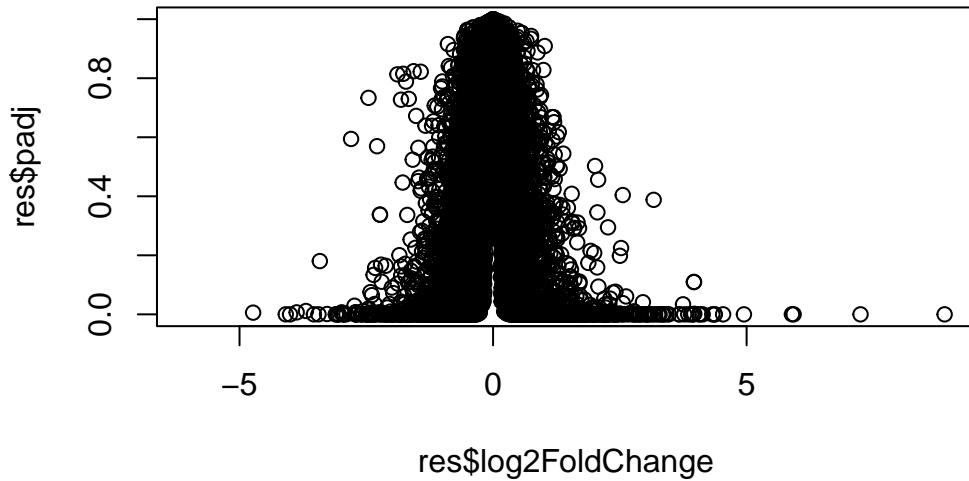
Extract my results

```
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
ENSG00000000419   520.134160  0.2061078  0.101059  2.039475 0.0414026
ENSG00000000457   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>
ENSG000000000003  0.163035
ENSG000000000005  NA
ENSG00000000419   0.176032
ENSG00000000457   0.961694
ENSG00000000460   0.815849
ENSG00000000938   NA
```

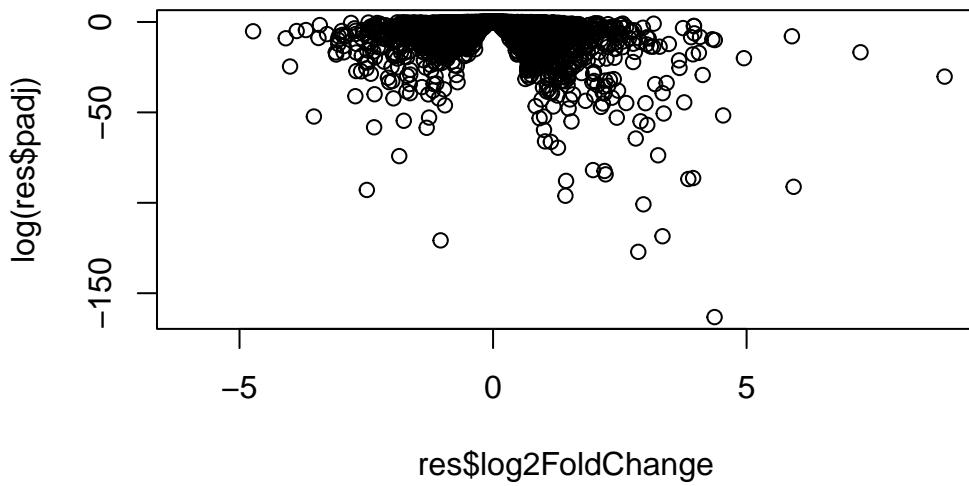
Plot of fold-change vs P-value (adjusted for multiple testing):

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



Take the log of the p-value

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



```
log(0.01)
```

```
[1] -4.60517
```

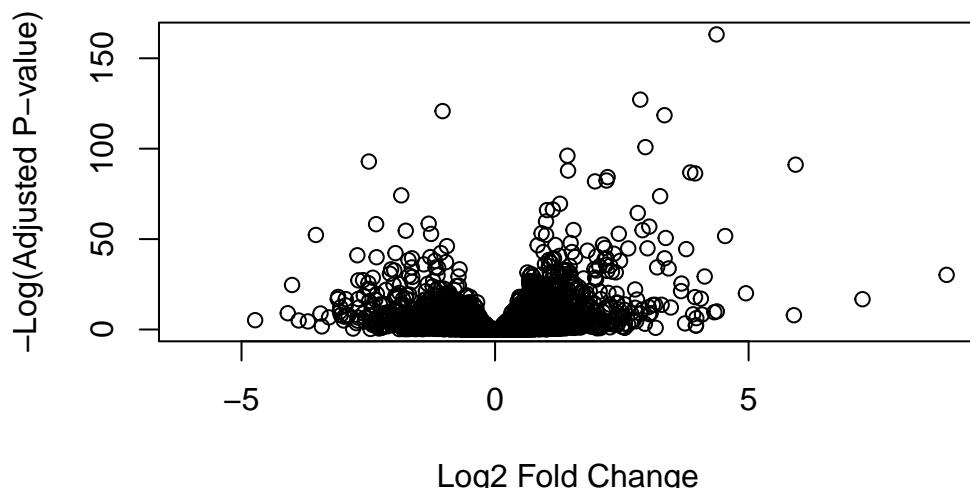
```
log(0.0000000001)
```

```
[1] -23.02585
```

A lower p-value has a lower log value so we need to look down the axis.

We can just flip that y-axis by putting a minus sign on it.

```
plot(res$log2FoldChange, -log(res$padj),  
      xlab = "Log2 Fold Change",  
      ylab = "-Log(Adjusted P-value)")
```



Upregulated is to the right, downregulated genes are to the left

Let's save our work to date

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

To finish off let's make a nicer volcano plot.

- Add the log2 threshold at  $+2/-2$
- Add p-value threshold at 0.05
- Add color to highlight the subset of genes that meet both of the above thresholds

Make it with ggplot2

```
# Setup our custom point color vector
mycols <- rep("grey", nrow(res))
mycols[res$log2FoldChange >= 2] <- "red"
mycols[res$log2FoldChange <= -2] <- "blue"
mycols[res$padj > 0.05] <- "grey"
```

```
ggplot(res) +
  aes(log2FoldChange, -log(res$padj)) +
  geom_point(col = mycols) +
  geom_vline(xintercept = c(-2, 2), col = "grey", linetype = 2) +
  geom_hline(xintercept = 0, yintercept = 0.05, col = "grey", linetype = 2) +
  labs(x = "Log2(FoldChange)", y = "-Log(P-Value)")
```

```
Warning in geom_hline(xintercept = 0, yintercept = 0.05, col = "grey", linetype = 2): Ignoring unknown parameters: `xintercept`
```

```
Warning: Removed 23549 rows containing missing values or values outside the scale range
(`geom_point()`).
```

