

Lab 13

```
library(DESeq2)
```

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

```
[1] 38694
```

```
#Q1 there are 38694 genes
metadata$dex == "control"
```

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

```
#Q2 there are 4 control cell lines
```

compare ctrl to treated cols 1. identify and extract “control” columns 2. calculate the mean value per gene for all these “control” columns 3. do the same for treated 4. compare the “control.mean” and “treated.mean” values

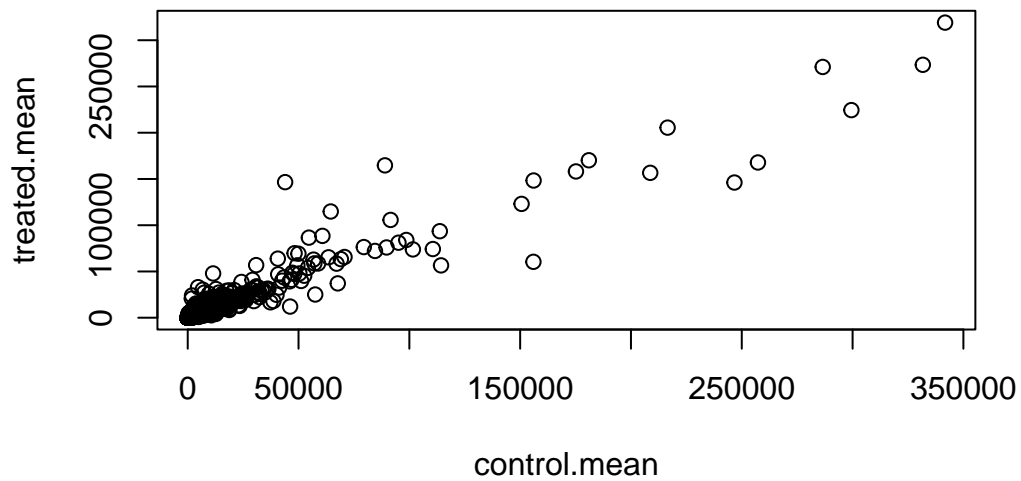
```
#step 1
control.inds <- metadata$dex == "control"
#metadata[control.inds, ]
control.mean <- rowMeans(counts[, control.inds])
#Q2 You would need to add a function that allows you to consider the mean when more samples
#Q4
treated.inds <- metadata$dex == "treated"
#metadata[treated.inds, ]
```

```
treated.mean <- rowMeans(counts[, treated.inds])
```

```
#Q5a
```

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

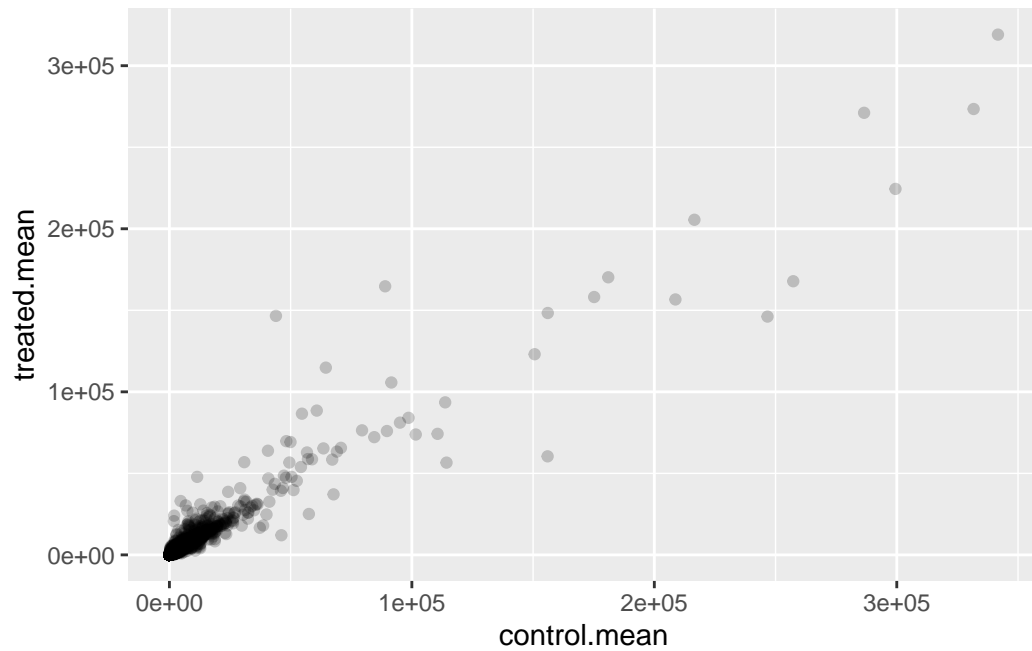
```
plot(meancounts)
```



```
#Q5b - geom_point()
```

```
library(ggplot2)
```

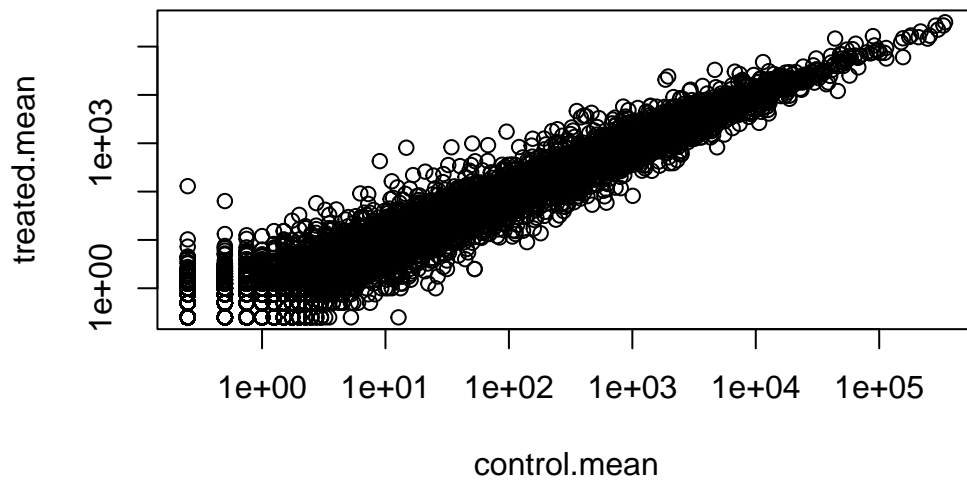
```
ggplot(meancounts) + aes(control.mean, treated.mean) + geom_point(alpha=0.2)
```



```
#Q6 - log  
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(10/10)
```

```
[1] 0
```

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

```
[1] 1
```

```
log2(5/10)
```

```
[1] -1
```

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

need to exclude any genes with 0 counts! ! !

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

```
nrow(mycounts)
```

```
[1] 21817
```

#Q7 arr.ind=TRUE returns the row and #columns indices where there are true values
#unique() will ensure no rows are counted twice if there are no entries in both samples

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

```
[1] 250
```

```
#Q8 250 genes are upregulated  
sum(mycounts$log2fc > -2)
```

```
[1] 21332
```

```
#Q9 21332 genes are down regulated  
#Q10 no, fold change can be large without  
# being statistically significant + it  
# depends on the p-value
```

```
library(DESeq2)  
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>
ENSG00000000003	747.1942	-0.3507030	0.168246	-2.084470	0.0371175
ENSG00000000005	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>				
ENSG00000000003	0.163035				
ENSG00000000005	NA				
ENSG000000000419	0.176032				
ENSG000000000457	0.961694				
ENSG000000000460	0.815849				
...	...				
ENSG00000283115	NA				
ENSG00000283116	NA				
ENSG00000283119	NA				
ENSG00000283120	NA				
ENSG00000283123	NA				

```

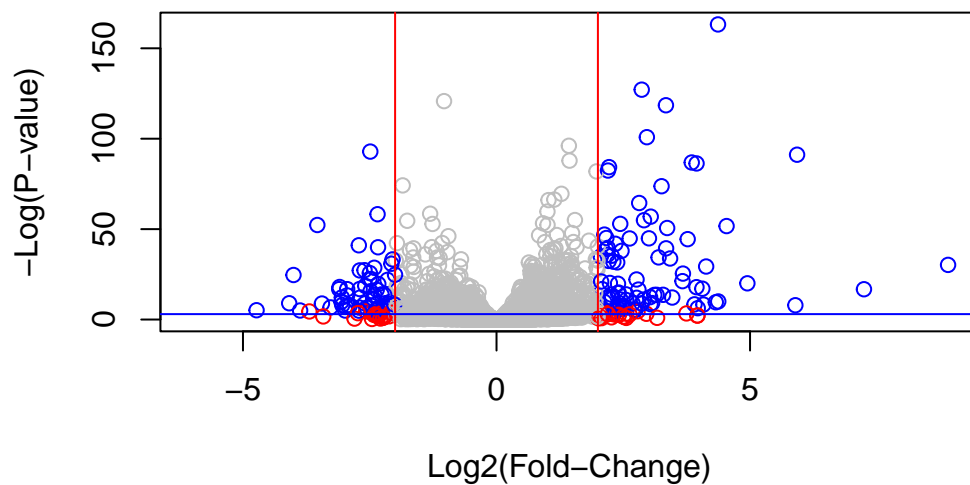
mycols <- rep("gray", nrow(res))
mycols[ abs(res$log2FoldChange) > 2 ] <- "red"

inds <- (res$padj < 0.01) & (abs(res$log2FoldChange) > 2 )
mycols[ inds ] <- "blue"

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

abline(v=c(-2,2), col="red")
abline(h=-log(0.05), col="blue")

```



```

mycols <- rep("gray",nrow(res))
mycols[ res$log2FoldChange >2 ] <- "black"
mycols[ res$log2FoldChange < -2 ] <- "black"
mycols[ res$padj > 0.05 ] <- "gray"

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

```

```
library(AnnotationDbi)
```

Warning: package 'AnnotationDbi' was built under R version 4.3.2

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

```
res$symbol <- mapIds(org.Hs.eg.db,  
                     keys=row.names(res), # Our genenames  
                     keytype="ENSEMBL",   # The format of our genenames  
                     column="SYMBOL",     # The new format we want to add  
                     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
ENSG000000000419	520.134160	0.2061078	0.101059	2.039475	0.0414026
ENSG000000000457	322.664844	0.0245269	0.145145	0.168982	0.8658106
ENSG000000000460	87.682625	-0.1471420	0.257007	-0.572521	0.5669691
ENSG000000000938	0.319167	-1.7322890	3.493601	-0.495846	0.6200029
	padj	symbol			
	<numeric>	<character>			
ENSG000000000003	0.163035	TSPAN6			
ENSG000000000005	NA	TNMD			
ENSG000000000419	0.176032	DPM1			
ENSG000000000457	0.961694	SCYL3			
ENSG000000000460	0.815849	FIRRM			
ENSG000000000938	NA	FGR			


```
#Q11
res$entrez <- mapIds(org.Hs.eg.db,
  keys=row.names(res),
  # Our genenames
  keytype="ENSEMBL",
  # The format of our genenames
  column="ENTREZID",
  # The new format we want to add
  multiVals="first")
```

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

```
res$uniprot <- mapIds(org.Hs.eg.db,
  keys=row.names(res),
  # Our genenames
  keytype="ENSEMBL",
  # The format of our genenames
  column="UNIPROT",
  # The new format we want to add
  multiVals="first")
```

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

```
res$genename <- mapIds(org.Hs.eg.db,
  keys=row.names(res),
  # Our genenames
  keytype="ENSEMBL",
  # The format of our genenames
  column="GENENAME",
  # The new format we want to add
  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 10 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
ENSG000000000460	87.682625	-0.1471420	0.257007	-0.572521	0.5669691
ENSG000000000938	0.319167	-1.7322890	3.493601	-0.495846	0.6200029
	padj	symbol	entrez	uniprot	
	<numeric>	<character>	<character>	<character>	
ENSG000000000003	0.163035	TSPAN6	7105	AOA024RCIO	
ENSG000000000005	NA	TNMD	64102	Q9H2S6	
ENSG000000000419	0.176032	DPM1	8813	O60762	
ENSG000000000457	0.961694	SCYL3	57147	Q8IZE3	
ENSG000000000460	0.815849	FIRRM	55732	AOA024R922	
ENSG000000000938	NA	FGR	2268	P09769	
		genename			
		<character>			
ENSG000000000003		tetraspanin 6			
ENSG000000000005		tenomodulin			
ENSG000000000419		dolichyl-phosphate m..			
ENSG000000000457		SCY1 like pseudokina..			
ENSG000000000460		FIGNL1 interacting r..			
ENSG000000000938		FGR proto-oncogene, ..			

```
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)
```

```
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
```

```
keggres = gage(foldchanges, gsets=kegg.sets.hs)
attributes (keggres)
```

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

```
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 C:/Users/isbel/Documents/BGGN 213/Class 13

Info: Writing image file hsa05310.pathview.png

