

Class 13: Transcriptomics and the analysis of RNA-Seq data

Liana Melikian

```
library(BiocManager)
```

```
Bioconductor version '3.16' is out-of-date; the current release version '3.18'  
is available with R version '4.3'; see https://bioconductor.org/install
```

```
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, colAnyNs, 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, rowAnyNs, rowAvgsPerColSet,
rowCollapse, rowCounts, rowCummaxs, rowCummins, rowCumprods,
rowCumsums, rowDiffs, rowIQRDiffs, rowIQRs, rowLogSumExps,
rowMadDiffs, rowMads, rowMaxs, rowMeans2, rowMedians, rowMins,
```

```
rowOrderStats, rowProds, rowQuantiles, rowRanges, rowRanks,  
rowSdDiff, rowSds, rowSums2, rowTabulates, rowVarDiff, 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
```

```
##Import countData and colData
```

```
# Complete the missing code  
counts <- read.csv("airway_scaledcounts.csv", row.names=1)  
metadata <- read.csv("airway_metadata.csv")  
  
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
ENSG000000000419	781	417	509
ENSG000000000457	447	330	324
ENSG000000000460	94	102	74
ENSG000000000938	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

Q1. We have 38694 genes.

```
nrow(counts)
```

```
[1] 38694
```

Q2 We have 4 “control” cell lines.

```
View(metadata)
```

```
##Toy differential gene expression
```

Let's start by calculating the mean counts per genes in the “control” samples. We can then compare this value for each gene to the mean counts in the “treated” samples (i.e. columns).

-Step 1. Find which columns in `counts` correspond to “control” samples. -Step 2. Calculate the mean value per gene in these columns. -Step 3. Store my answer for late in `control.mean`

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

```

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

```

Q3. We can use `rowSums()` to make the code approach more robust.

```

control inds=metadata$dex == "control"

control inds=metadata$dex=="control"

metadata[control inds,]

  id      dex celltype    geo_id
1 SRR1039508 control   N61311 GSM1275862
3 SRR1039512 control   N052611 GSM1275866
5 SRR1039516 control   N080611 GSM1275870
7 SRR1039520 control   N061011 GSM1275874

control counts=counts[,control inds]
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

#apply(control counts,1,mean)
control mean=rowMeans(control counts)

```

Q4. Now do the same steps to get `treated.mean`

```
treated mean=rowMeans(counts[,metadata$dex == "treated"])
```

To keep us tidy lets put `control.mean` and `treated.mean` vectors together as two columns of a new data.frame.

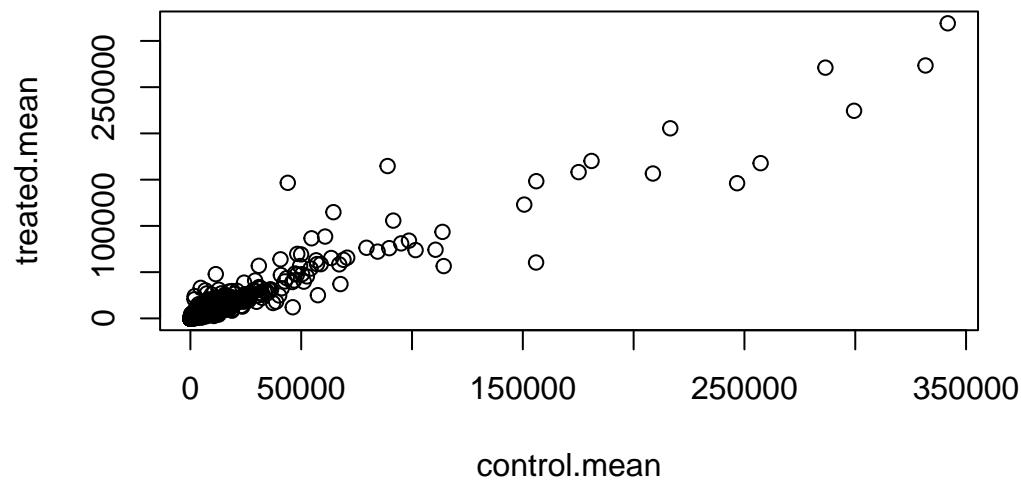
```
meancounts=data.frame(control.mean, treated.mean)
```

```
head(meancounts)
```

	control.mean	treated.mean
ENSG000000000003	900.75	658.00
ENSG000000000005	0.00	0.00
ENSG00000000419	520.50	546.00
ENSG00000000457	339.75	316.50
ENSG00000000460	97.25	78.75
ENSG00000000938	0.75	0.00

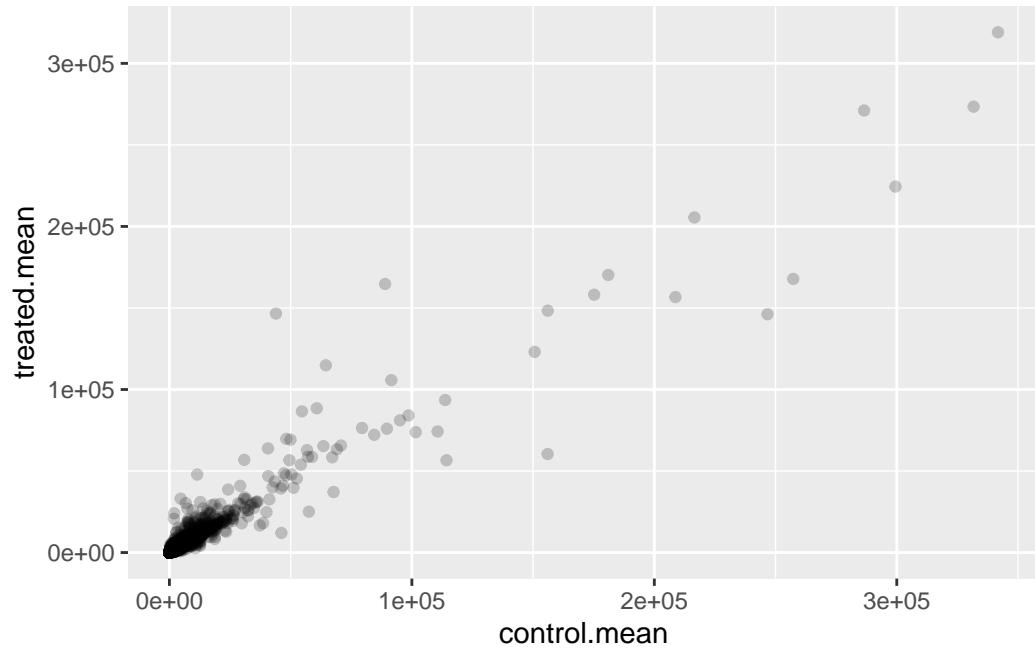
Q5(a).

```
plot(meancounts)
```



And a ggplot version:

```
library(ggplot2)
ggplot(meancounts)+
  aes(control.mean,treated.mean)+
  geom_point(alpha=.2)
```



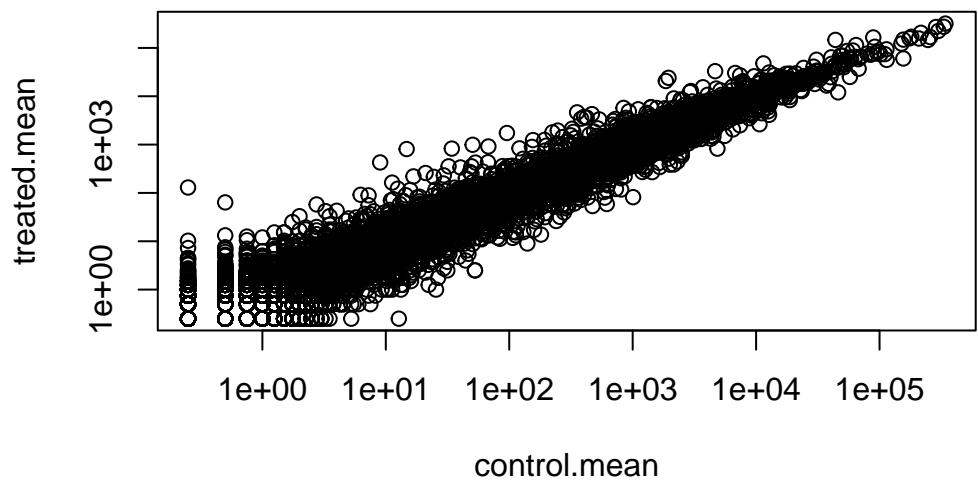
Q5(b). geom_point

Q6.

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

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

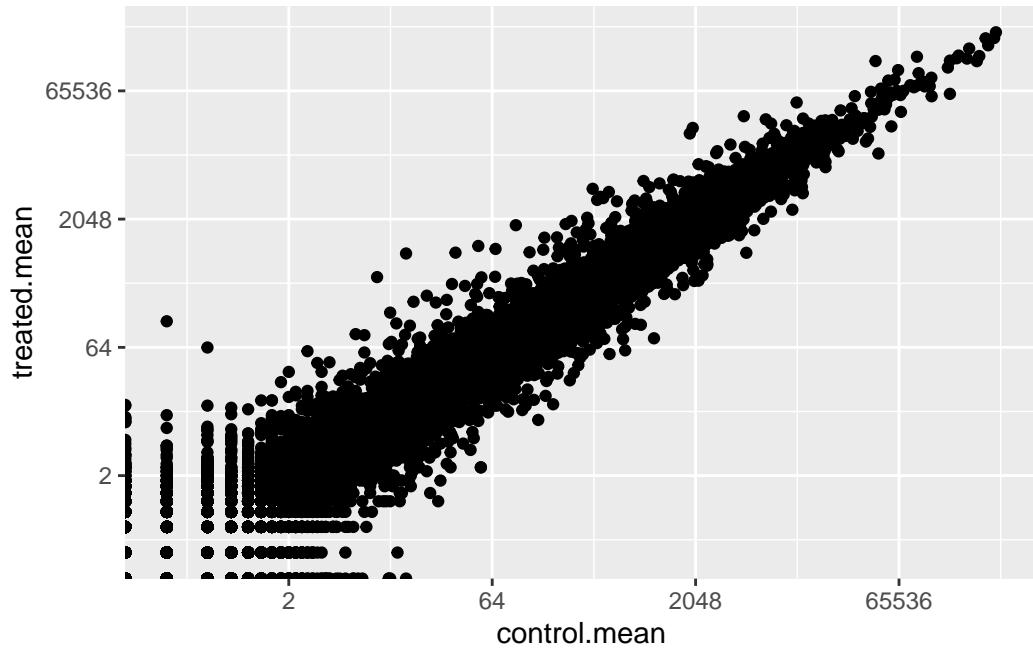
Warning in xy.coords(x, y, xlabel, ylabel, log): 15281 y values <= 0 omitted
from logarithmic plot



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

Warning: Transformation introduced infinite values in continuous x-axis

Warning: Transformation introduced infinite values in continuous y-axis



Log transformations are super useful when our data is skewed and measured over a wide range like this. We can use different log transformations like base10 or natural logs but we most often prefer log2 units.

```
#Control/Treated
log2(10/10)
```

```
[1] 0
```

What if there was a doubling?

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

```
[1] 1
```

Half counts

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

```
[1] -1
```

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

```
[1] 2
```

Let's add a log2 fold-change column to our little `meancounts` data.frame:

```
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

There are a couple of weird results. The ! flips TRUE values to FALSE and vice versa

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

dim(mycounts)
```

```
[1] 21817      3
```

```
head(mycounts)
```

	control.mean	treated.mean	log2fc
ENSG000000000003	900.75	658.00	-0.45303916
ENSG000000000419	520.50	546.00	0.06900279
ENSG000000000457	339.75	316.50	-0.10226805
ENSG000000000460	97.25	78.75	-0.30441833
ENSG000000000971	5219.00	6687.50	0.35769358
ENSG000000001036	2327.00	1785.75	-0.38194109

Q7. The arr.ind = TRUE argument in which() retrieves row and column indices where there are TRUE values. Using unique(zero.vals[, 1]) ensures only unique row indices with zero counts in at least one sample are selected for removal, focusing on rows without zero counts in both specified columns.

```
head(mycounts)
```

	control.mean	treated.mean	log2fc
ENSG000000000003	900.75	658.00	-0.45303916
ENSG000000000419	520.50	546.00	0.06900279
ENSG000000000457	339.75	316.50	-0.10226805
ENSG000000000460	97.25	78.75	-0.30441833
ENSG000000000971	5219.00	6687.50	0.35769358
ENSG00000001036	2327.00	1785.75	-0.38194109

A common threshold used for calling something differentially expressed is a log2(FoldChange) of greater than 2 or less than -2. Let's filter the dataset both ways to see how many genes are up or down-regulated.

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

Q8. 250 up-regulated genes

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

```
[1] 250
```

Q9. 367 down-regulated genes

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

```
[1] 367
```

Q10. No, we haven't considered the statistical significance of the differences in fold changes.

We will use the DESeq2 package to do this analysis properly.

```
##Setting up for DESeq We must load it up with a library() call.
```

```
library(DESeq2)
citation("DESeq2")
```

To cite package 'DESeq2' in publications use:

Love, M.I., Huber, W., Anders, S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2 Genome Biology 15(12):550 (2014)

A BibTeX entry for LaTeX users is

```
@Article{,
  title = {Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2},
  author = {Michael I. Love and Wolfgang Huber and Simon Anders},
  year = {2014},
  journal = {Genome Biology},
  doi = {10.1186/s13059-014-0550-8},
  volume = {15},
  issue = {12},
  pages = {550},
}
```

```
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
```

```
class: DESeqDataSet
dim: 38694 8
metadata(1): version
```

```
assays(1): counts
rownames(38694): ENSG000000000003 ENSG000000000005 ... ENSG00000283120
  ENSG00000283123
rowData names(0):
colnames(8): SRR1039508 SRR1039509 ... SRR1039520 SRR1039521
colData names(4): id dex celltype geo_id
```

Now we can run our DESeq analysis

```
dds=DESeq(dds)
```

```
estimating size factors
```

```
estimating dispersions
```

```
gene-wise dispersion estimates
```

```
mean-dispersion relationship
```

```
final dispersion estimates
```

```
fitting model and testing
```

Get our results back from the dds object.

```
dds=DESeq(dds)
```

```
using pre-existing size factors
```

```
estimating dispersions
```

```
found already estimated dispersions, replacing these
```

```
gene-wise dispersion estimates
```

```
mean-dispersion relationship
```

```
final dispersion estimates
```

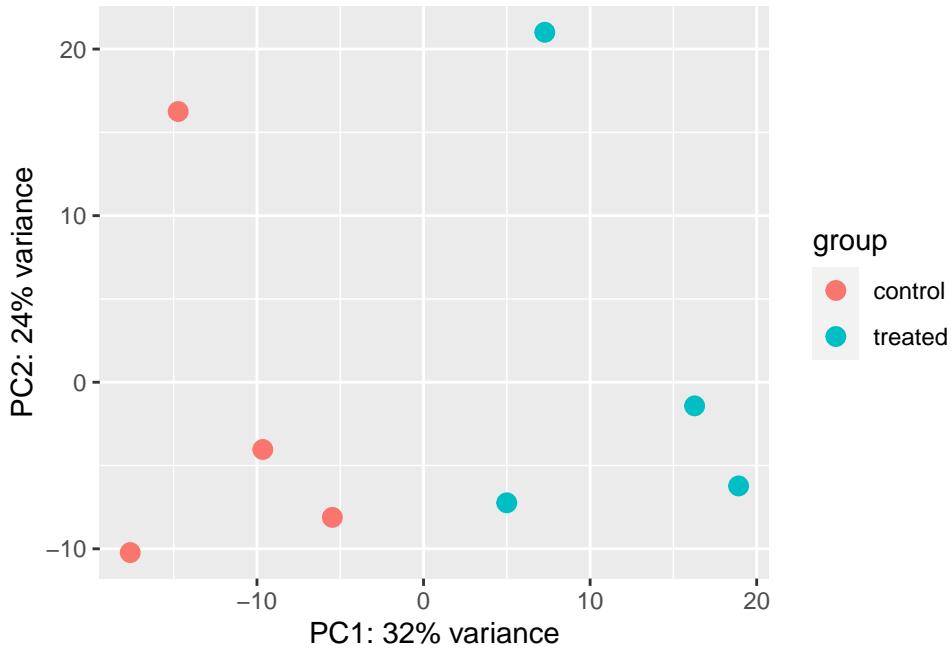
```
fitting model and testing
```

```
results(dds)
```

```
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
ENSG00000000419   520.1342    0.2061078  0.101059  2.039475 0.0414026
ENSG00000000457   322.6648    0.0245269  0.145145  0.168982 0.8658106
ENSG00000000460   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
ENSG00000000419   0.176032
ENSG00000000457   0.961694
ENSG00000000460   0.815849
...
...
ENSG00000283115       NA
ENSG00000283116       NA
ENSG00000283119       NA
ENSG00000283120       NA
ENSG00000283123       NA
```

```
##PCA
```

```
vsd <- vst(dds, blind = FALSE)
plotPCA(vsd, intgroup = c("dex"))
```

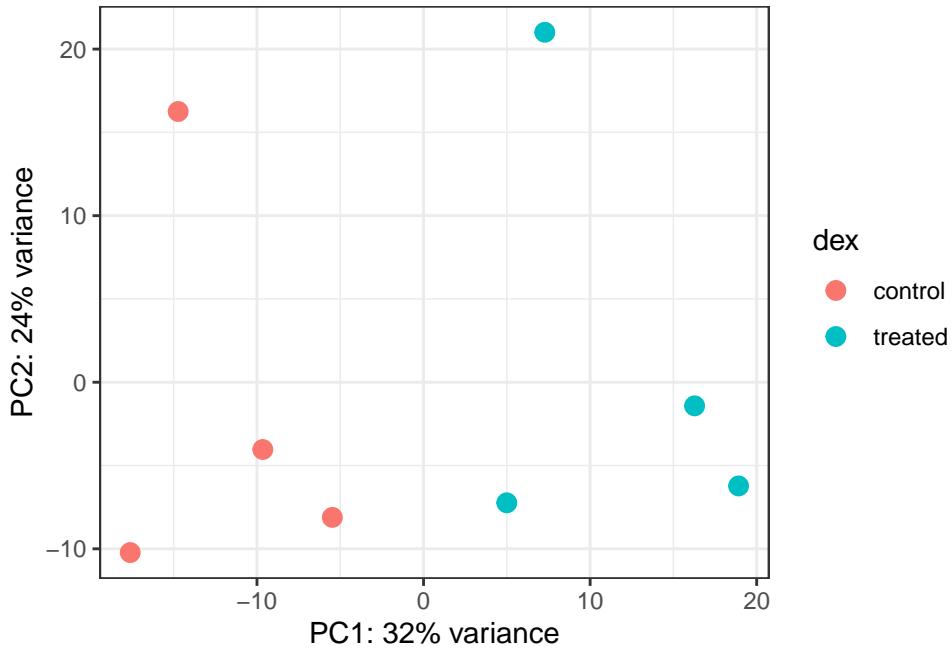


```
pcaData <- plotPCA(vsd, intgroup=c("dex"), returnData=TRUE)
head(pcaData)
```

	PC1	PC2	group	dex	name
SRR1039508	-17.607922	-10.225252	control	control	SRR1039508
SRR1039509	4.996738	-7.238117	treated	treated	SRR1039509
SRR1039512	-5.474456	-8.113993	control	control	SRR1039512
SRR1039513	18.912974	-6.226041	treated	treated	SRR1039513
SRR1039516	-14.729173	16.252000	control	control	SRR1039516
SRR1039517	7.279863	21.008034	treated	treated	SRR1039517

```
# Calculate percent variance per PC for the plot axis labels
percentVar <- round(100 * attr(pcaData, "percentVar"))
```

```
ggplot(pcaData) +
  aes(x = PC1, y = PC2, color = dex) +
  geom_point(size = 3) +
  xlab(paste0("PC1: ", percentVar[1], "% variance")) +
  ylab(paste0("PC2: ", percentVar[2], "% variance")) +
  coord_fixed() +
  theme_bw()
```



```
##DESeq Analysis
```

```
results(dds)
```

```
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
ENSG00000000419	520.1342	0.2061078	0.101059	2.039475	0.0414026
ENSG00000000457	322.6648	0.0245269	0.145145	0.168982	0.8658106
ENSG00000000460	87.6826	-0.1471420	0.257007	-0.572521	0.5669691
...
ENSG0000283115	0.000000		NA	NA	NA
ENSG0000283116	0.000000		NA	NA	NA
ENSG0000283119	0.000000		NA	NA	NA
ENSG0000283120	0.974916	-0.668258	1.69456	-0.394354	0.693319
ENSG0000283123	0.000000		NA	NA	NA
		padj			
		<numeric>			

```
ENSG000000000003 0.163035
ENSG000000000005 NA
ENSG000000000419 0.176032
ENSG000000000457 0.961694
ENSG000000000460 0.815849
...
ENSG0000283115 NA
ENSG0000283116 NA
ENSG0000283119 NA
ENSG0000283120 NA
ENSG0000283123 NA
```

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

using pre-existing size factors

estimating dispersions

found already estimated dispersions, replacing these

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>
ENSG000000000003	747.1942	-0.3507030	0.168246	-2.084470	0.0371175

```

ENSG000000000005  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>
ENSG000000000003 0.163035
ENSG000000000005 NA
ENSG000000000419 0.176032
ENSG000000000457 0.961694
ENSG000000000460 0.815849
...
...      ...
ENSG00000283115  NA
ENSG00000283116  NA
ENSG00000283119  NA
ENSG00000283120  NA
ENSG00000283123  NA

```

```
summary(res)
```

```

out of 25258 with nonzero total read count
adjusted p-value < 0.1
LFC > 0 (up)      : 1563, 6.2%
LFC < 0 (down)     : 1188, 4.7%
outliers [1]       : 142, 0.56%
low counts [2]     : 9971, 39%
(mean count < 10)
[1] see 'cooksCutoff' argument of ?results
[2] see 'independentFiltering' argument of ?results

```

```
res05 <- results(dds, alpha=0.05)
summary(res05)
```

```

out of 25258 with nonzero total read count
adjusted p-value < 0.05
LFC > 0 (up)      : 1236, 4.9%
LFC < 0 (down)    : 933, 3.7%
outliers [1]       : 142, 0.56%
low counts [2]     : 9033, 36%
(mean count < 6)
[1] see 'cooksCutoff' argument of ?results
[2] see 'independentFiltering' argument of ?results

```

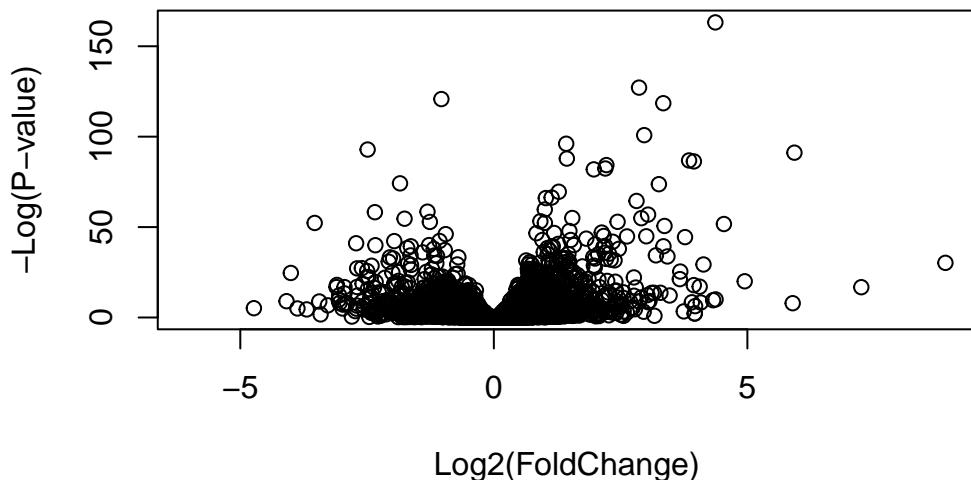
##Data Visualization

Volcano plot. This is a common type of summary figure that keep both our inner biologist and inner stats nerd happy because it shows both P-values and Log2(Fold-Changes).

```

plot( res$log2FoldChange, -log(res$padj),
      xlab="Log2(FoldChange)",
      ylab="-Log(P-value)")

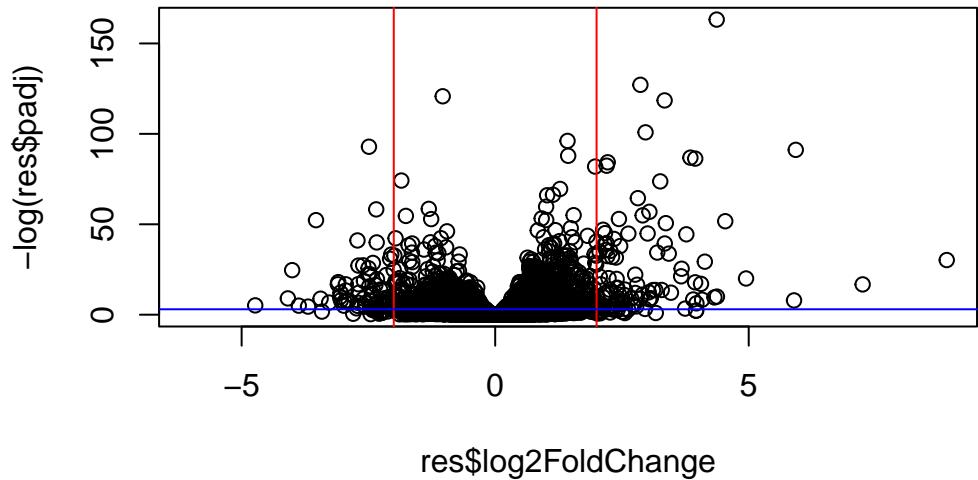
```



```

plot(res$log2FoldChange, -log(res$padj))
abline(v=2,col="red")
abline(v=-2,col="red")
abline(h=-log(0.05),col="blue")

```



```

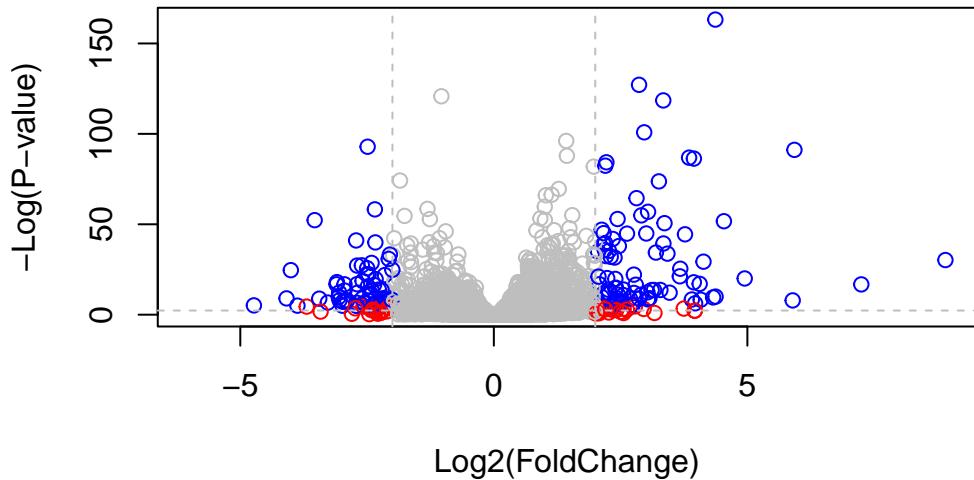
# Setup our custom point color vector
mycols <- rep("gray", nrow(res))
mycols[ abs(res$log2FoldChange) > 2 ] <- "red"

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

# Volcano plot with custom colors
plot( res$log2FoldChange, -log(res$padj),
col=mycols, ylab="-Log(P-value)", xlab="Log2(FoldChange)" )

# Cut-off lines
abline(v=c(-2,2), col="gray", lty=2)
abline(h=-log(0.1), col="gray", lty=2)

```



```
library(EnhancedVolcano)
```

Loading required package: ggrepel

```
res=results(dds)

x <- as.data.frame(res)

#EnhancedVolcano(x,
#                 #lab = x$symbol,
#                 # x = 'log2FoldChange',
#                 # y = 'pvalue')
```

Save our results to date

```
#write.csv(res,fule="deseq-results.csv")
```

```
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
ENSG000000000460 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

```

##Adding annotation data

```
library("AnnotationDbi")
```

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

```
columns(org.Hs.eg.db)
```

```
[1] "ACNUM"      "ALIAS"       "ENSEMBL"     "ENSEMLPROT"  "ENSEMLTRANS"
[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"
```

The main function we will use here is called `mapIDs()`

Our current IDs are here:

```
#mapIDs()
head(rownames(res))
```

```
[1] "ENSG000000000003" "ENSG000000000005" "ENSG00000000419" "ENSG00000000457"  
[5] "ENSG00000000460" "ENSG00000000938"
```

Q11. These are in ENSEMBLE format. I want “SYMBOL” ids.

```
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  
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      symbol  
  <numeric> <character>  
ENSG000000000003 0.163035    TSPAN6  
ENSG000000000005  NA         TNMD  
ENSG00000000419  0.176032    DPM1  
ENSG00000000457  0.961694    SCYL3  
ENSG00000000460  0.815849    C1orf112  
ENSG00000000938  NA         FGR
```

Let's add GENENAME

```
res$genename = 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 8 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    genename
      <numeric> <character> <character>
ENSG000000000003  0.163035    TSPAN6    TSPAN6
ENSG000000000005   NA         TNMD     TNMD
ENSG000000000419  0.176032    DPM1     DPM1
ENSG000000000457  0.961694    SCYL3    SCYL3
ENSG000000000460  0.815849    C1orf112 C1orf112
ENSG000000000938   NA         FGR      FGR

res$entrez=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 9 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  genename      entrez
      <numeric> <character> <character> <character>
ENSG000000000003 0.163035    TSPAN6    TSPAN6    TSPAN6
ENSG000000000005  NA        TNMD      TNMD      TNMD
ENSG000000000419 0.176032    DPM1      DPM1      DPM1
ENSG000000000457 0.961694    SCYL3    SCYL3    SCYL3
ENSG000000000460 0.815849    C1orf112  C1orf112  C1orf112
ENSG000000000938  NA        FGR       FGR       FGR

```

##Pathway analysis

We will use the **gage** package along with **pathview** here to do geneset enrichment (aka pathway analysis) and figure generation respectively.

```

library(pathview)
library(gage)
library(gageData)

```

Let's have a peak at the first two pathways in KEGG.

```

data(kegg.sets.hs)

# Examine the first 2 pathways in this kegg set for humans
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"
```

What we need for `gage()` is our genes in ENTREZ id format with a measure of their importance.

It wants a vector of e.g. fold-changes.

```
foldchanges = res$log2FoldChange
head(foldchanges)
```

```
[1] -0.35070302           NA  0.20610777  0.02452695 -0.14714205 -1.73228897
```

Add ENTREZ ids as `names()` to my `foldchanges` vector.

```
names(foldchanges)=res$entrez
head(foldchanges)
```

```
TSPAN6          TNMD         DPM1        SCYL3       C1orf112        FGR
-0.35070302      NA  0.20610777  0.02452695 -0.14714205 -1.73228897
```

Now we can run `gage()` with this input vector and the gensec we want to examine for overlap/enrichment.

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

Look at the results.

```
attributes(keggres)
```

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

```
head(keggres$less,3)
```

	p.geomean	stat.mean	p.val	q.val
hsa00232 Caffeine metabolism	NA	NaN	NA	NA
hsa00983 Drug metabolism - other enzymes	NA	NaN	NA	NA
hsa01100 Metabolic pathways	NA	NaN	NA	NA
	set.size	exp1		
hsa00232 Caffeine metabolism	0	NA		
hsa00983 Drug metabolism - other enzymes	0	NA		
hsa01100 Metabolic pathways	0	NA		

We can view these pathways with our geneset genes highlighted using the `pathview()` function.
E.g. for "Asthma" I will use the pathway.id hs05310 as seen above.

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

Warning: None of the genes or compounds mapped to the pathway!
Argument gene.idtype or cpd.idtype may be wrong.

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

```
Info: Working in directory /Users/lhmel/Desktop/BIMM 143/class13
```

```
Info: Writing image file hsa05310.pathview.png
```

