

class13

First, install packages

```
install.packages("BiocManager") BiocManager::install() BiocManager::install("DESeq2")
```

In today's class we will explore and 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).

##Data Import We have two input files, so-called "count data" and "coldata"

```
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
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		

```
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(counts)
```

```
[1] 38694
```

There are 38694 genes in the dataset.

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

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

```
control treated
        4         4
```

There are 4 control cells in the dataset.

differential gene expression

Time to do some analysis. We have 4 control and 4 treated samples/experiments/columns. Make sure the metadata id column matches the columns in our countdata.

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

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

```
colnames(counts)
```

```
[1] "SRR1039508" "SRR1039509" "SRR1039512" "SRR1039513" "SRR1039516"
[6] "SRR1039517" "SRR1039520" "SRR1039521"
```

```
metadata$id
```

```
[1] "SRR1039508" "SRR1039509" "SRR1039512" "SRR1039513" "SRR1039516"  
[6] "SRR1039517" "SRR1039520" "SRR1039521"
```

To check that all elements of a vector are TRUE we can use the all() function.

```
all(c(T,T,T,T))
```

```
[1] TRUE
```

```
all(c(T,T,T,F))
```

```
[1] FALSE
```

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

```
[1] TRUE
```

To start I will calculate the ‘control.mean’ and ‘treated.mean’ vlaues and compare them.

- Identify and extract the ‘control’ only columns
- determine the mean value for each gene (i.e. row)
- Do the same for ‘treated’

```
control <- metadata[metadata$dex == "control", ]  
control.counts <- counts[ , control$id]  
control.mean <- apply(control.counts, 1, mean)  
head(control.mean)
```

```
ENSG00000000003 ENSG00000000005 ENSG00000000419 ENSG00000000457 ENSG00000000460  
900.75 0.00 520.50 339.75 97.25  
ENSG00000000938  
0.75
```

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

The issue with given code in the website is that mean is calculated by dividing rowSUs with 4. This method cause issue when there are more than 4 types of cell lines in the dataset. To ensure that all number of cell lines are included, we can use apply function instead of doing “/4” as shown in code above.

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)

```
treated <- metadata[metadata$dex == "treated", ]  
treated.counts <- counts[, treated$id]  
treated.mean <- apply(treated.counts, 1, mean)  
head(treated.mean)
```

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

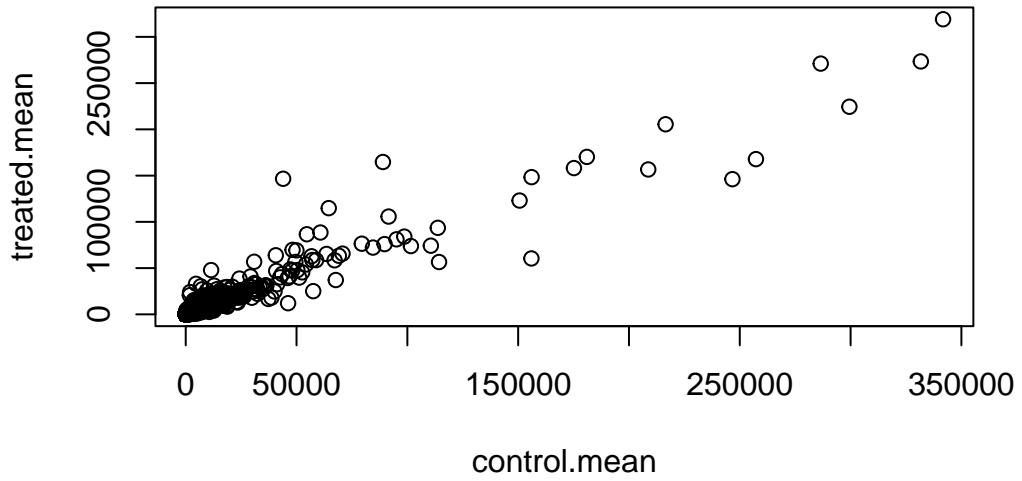
Let's store these togehter for ease of book-keeping

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

have a quick view of this data:

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.

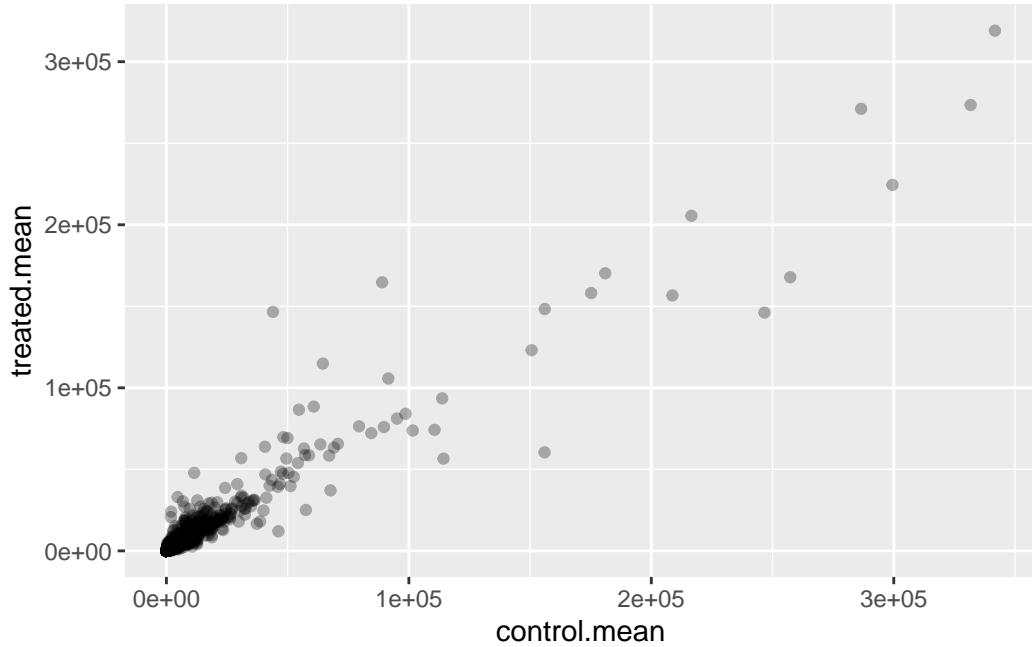
```
plot(meancounts)
```



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?

```
library(ggplot2)

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



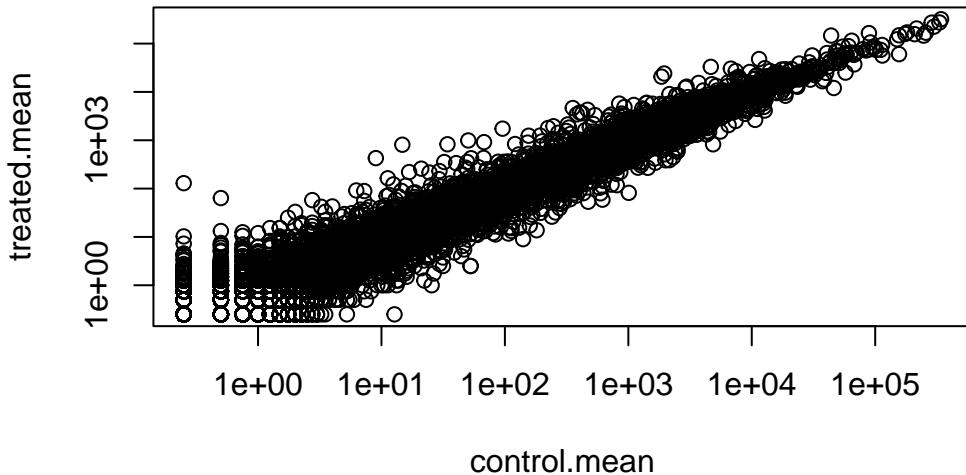
This data is screaming at us to log transform as it is so heavily skewed and over such a wide range

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
```



I want to compare the treated and the control values here and we will use fold change in log2 units to do this. $\log_2(\text{treated}/\text{control})$

```

log2(20/20) #no difference
[1] 0

log2(20/10) #doubling in the treated
[1] 1

log2(5/10) #half in the treated
[1] -1

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

```

A common rule of thumb cut-off for calling a gene “differentially expressed” is a log2 fold change value of either $> +2$ or <-2 for “up regulated” and “down regulated” respectively.

```
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
```

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

```
[1] 1846
```

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

```
[1] 2212
```

We first need to remove zero count genes - as we can't say anything about these genes anyway and their division of log values are messing things up (divide by zero) or the -inf log problem

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

Q7. What is the purpose of the arr.ind argument in the which() function call above? Why would we then take the first column of the output and need to call the unique() function?

While meancounts[,1:2]==0 filters control.mean and treated.mean with value of zero, it returns boolean. arr.ind = TRUE return row and column index where TRUE value is, indicating where the gene with zero value locates. Hence, arr.ind makes easier to summarize genes with zero values as we don't have to look for TRUE throughout entire dataframe.

Q. How many genes do we have left that we can say something about (i.e. they don't have any zero counts)

```
nrow(mycounts)
```

```
[1] 21817
```

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

[1] 250

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

[1] 367

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

There are 250 upregulated genes

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

There are 367 downregulated gene

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

No. Fold change itself doesn't indicate statistically significant difference. To determine that genes expression difference is actually significant, we need to evaluate p-values.

#DESeq analysis Let's do this properly with the help of the DESeq2 package

```
#|message: FALSE
library(DESeq2)

: S4Vectors

: stats4

: BiocGenerics

: '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

: 'S4Vectors'

The following object is masked from 'package:utils':

findMatches

The following objects are masked from 'package:base':

expand.grid, I, unname

: IRanges

: 'IRanges'

The following object is masked from 'package:grDevices':

windows

: GenomicRanges

: GenomeInfoDb

: SummarizedExperiment

```
: MatrixGenerics
```

```
: matrixStats
```

```
: '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,  
rowSdDiffs, rowSds, rowSums2, rowTabulates, rowVarDiffs, rowVars,  
rowWeightedMads, rowWeightedMeans, rowWeightedMedians,  
rowWeightedSds, rowWeightedVars
```

```
: Biobase
```

```
Welcome to Bioconductor
```

```
Vignettes contain introductory material; view with  
'browseVignettes()'. To cite Bioconductor, see  
'citation("Biobase")', and for packages 'citation("pkgname")'.
```

```
: 'Biobase'
```

```
The following object is masked from 'package:MatrixGenerics':
```

```
rowMedians
```

```
The following objects are masked from 'package:matrixStats':
```

```
anyMissing, rowMedians
```

We have to use a specific data object for working with DESeq

```
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
```

Run our main analysis with the ‘DESeq()’ function

```
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 out of our ‘dds’ object we can use DESeq function called ‘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
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

```

```
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 = 0.5)
summary(res05)
```

```

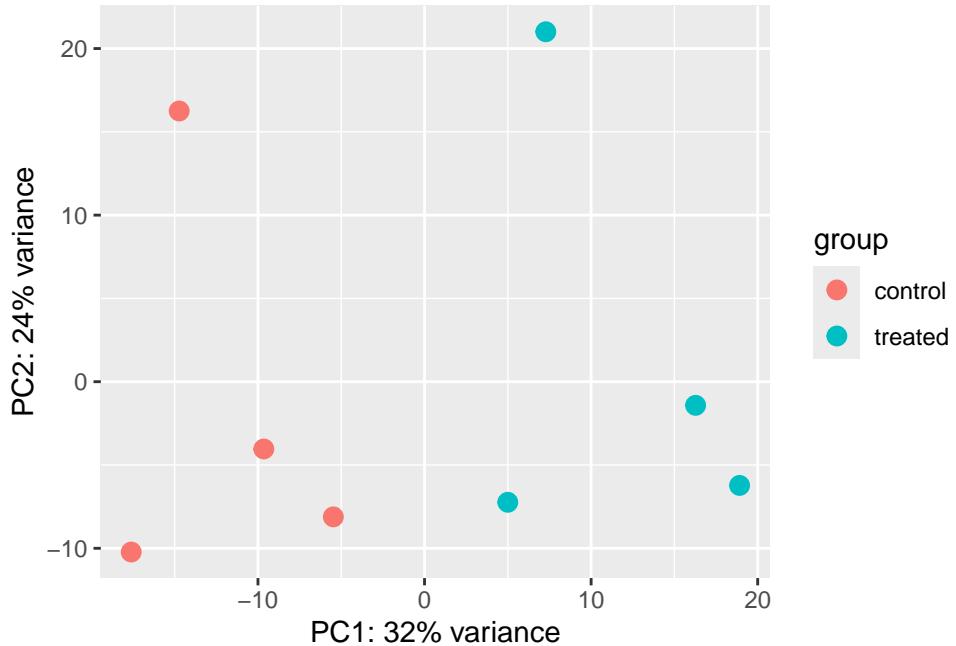
out of 25258 with nonzero total read count
adjusted p-value < 0.5
LFC > 0 (up)      : 3375, 13%

```

```
LFC < 0 (down)      : 3108, 12%
outliers [1]        : 142, 0.56%
low counts [2]       : 8564, 34%
(mean count < 4)
[1] see 'cooksCutoff' argument of ?results
[2] see 'independentFiltering' argument of ?results
```

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

using ntop=500 top features by variance



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

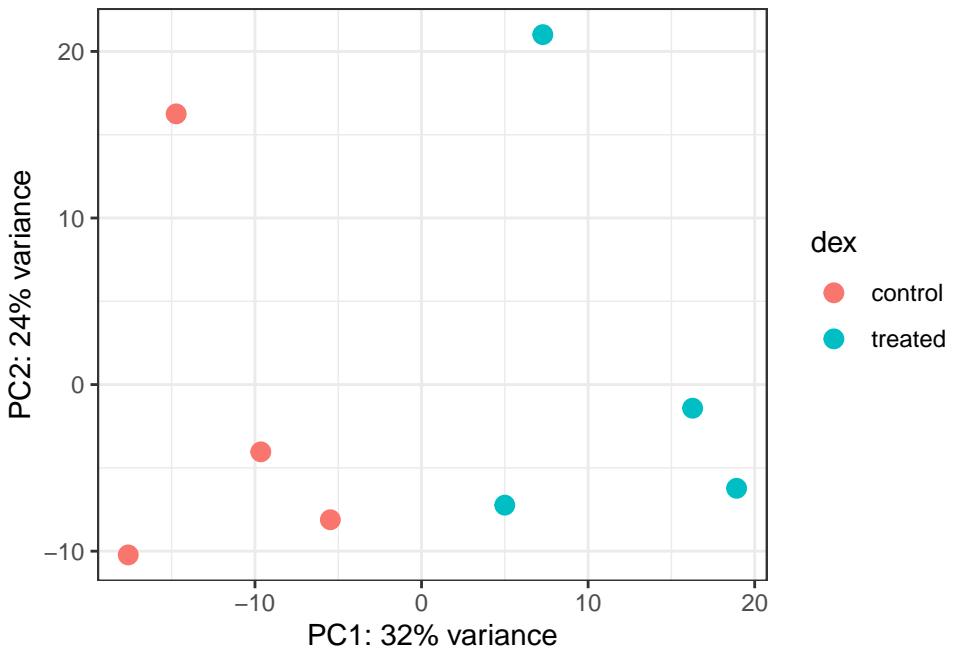
using ntop=500 top features by variance

```
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()
```



```
#BiocManager::install("AnnotationDbi")
#BiocManager::install("org.Hs.eg.db")
```

```

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

columns(org.Hs.eg.db)

[1] "ACCCNUM"        "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"

res

og2 fold change (MLE): dex treated vs control
ald test p-value: dex treated vs control
ataFrame with 38694 rows and 6 columns
  baseMean log2FoldChange      lfcSE      stat      pvalue
  <numeric>      <numeric> <numeric> <numeric> <numeric>
NSG000000000003  747.1942   -0.3507030  0.168246 -2.084470  0.0371175
NSG000000000005   0.0000      NA        NA        NA        NA
NSG000000000419  520.1342   0.2061078  0.101059  2.039475  0.0414026
NSG000000000457  322.6648   0.0245269  0.145145  0.168982  0.8658106
NSG000000000460   87.6826   -0.1471420  0.257007 -0.572521  0.5669691
...
          ...
  padj
  <numeric>
NSG000000000003  0.163035
NSG000000000005   NA
NSG000000000419  0.176032
NSG000000000457  0.961694
NSG000000000460  0.815849

```

```

...
ENSG00000283115      NA
ENSG00000283116      NA
ENSG00000283119      NA
ENSG00000283120      NA
ENSG00000283123      NA

res$symbol <- mapIds(org.Hs.eg.db,
                      keys=row.names(res),          # Gene name
                      keytype="ENSEMBL",            # Format of gene name
                      column="SYMBOL",             # New format to add
                      multiVals="first")

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

```

Q11. Run the mapIds() function two more times to add the Entrez ID and UniProt accession and GENENAME as new columns called res\$entrez, res\$uniprot, and res\$genename.

```

#res$entrez
res$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

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

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

#res$genename
res$genename <- mapIds(org.Hs.eg.db,

```

```

    keys=row.names(res),
    keytype="ENSEMBL",
    column="GENENAME",
    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      uniport
  <numeric> <character> <character> <character>
ENSG000000000003 0.163035      TSPAN6       7105 AOA024RC10
ENSG000000000005   NA          TNMD        64102 Q9H2S6
ENSG000000000419 0.176032      DPM1        8813 060762
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, ..

```

```

ord <- order(res$padj)
#View(res[ord,])

```

```

head(res[ord,])

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>
ENSG00000152583    954.771      4.36836  0.2371268   18.4220 8.74490e-76
ENSG00000179094    743.253      2.86389  0.1755693   16.3120 8.10784e-60
ENSG00000116584   2277.913     -1.03470  0.0650984  -15.8944 6.92855e-57
ENSG00000189221   2383.754      3.34154  0.2124058   15.7319 9.14433e-56
ENSG00000120129   3440.704      2.96521  0.2036951   14.5571 5.26424e-48
ENSG00000148175   13493.920     1.42717  0.1003890   14.2164 7.25128e-46
  padj      symbol      entrez      uniprot
  <numeric> <character> <character> <character>
ENSG00000152583 1.32441e-71      SPARCL1      8404  AOA024RDE1
ENSG00000179094 6.13966e-56       PER1        5187  015534
ENSG00000116584 3.49776e-53      ARHGEF2      9181  Q92974
ENSG00000189221 3.46227e-52       MAOA        4128  P21397
ENSG00000120129 1.59454e-44      DUSP1        1843  B4DU40
ENSG00000148175 1.83034e-42       STOM        2040  F8VSL7
  genename
  <character>
ENSG00000152583           SPARC like 1
ENSG00000179094           period circadian reg..
ENSG00000116584           Rho/Rac guanine nucl..
ENSG00000189221           monoamine oxidase A
ENSG00000120129           dual specificity pho..
ENSG00000148175           stomatin

```

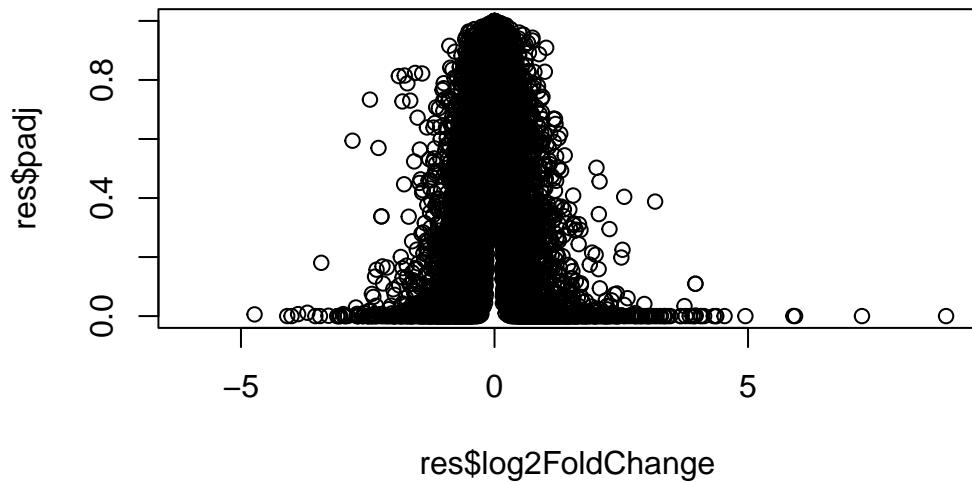
write.csv(res[ord,], "deseq_results.csv")

Adjusted p value ensure low false positive. While 5% is widely used, 5% of 25000 is hella lots of inaccuracy.

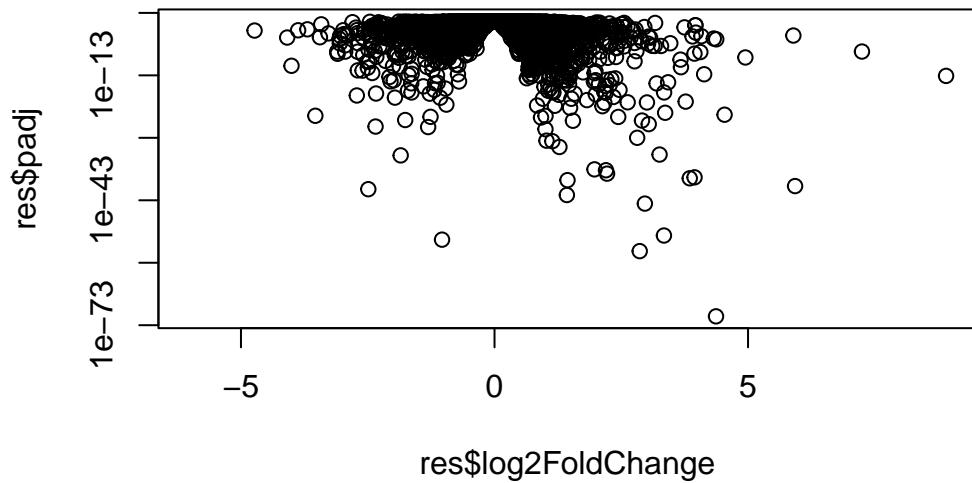
Volcano Plot

A very common and useful summary results figure from this type of analysis is called a volcano plot - a plot of log2FC vs p-value. We use the ‘padj’ the adjusted p-value for multiple testing.

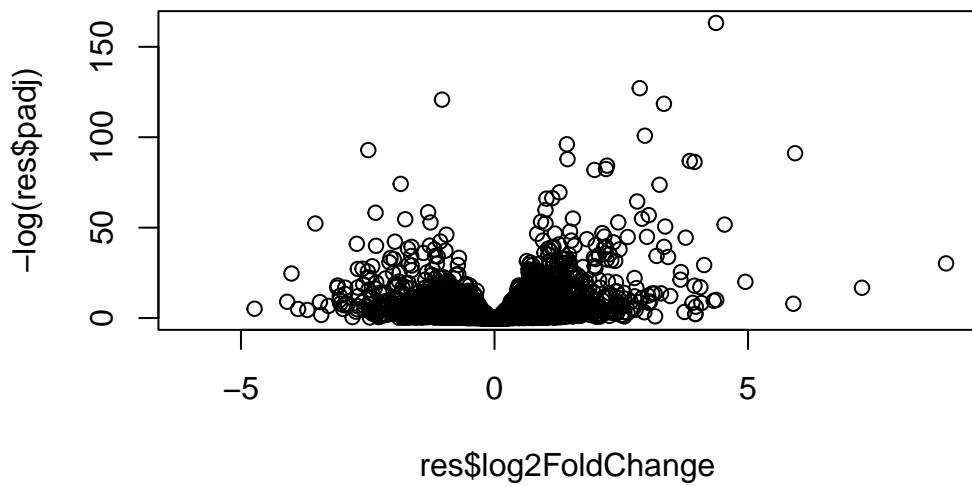
```
plot(res$log2FoldChange, res$padj)
```



```
plot(res$log2FoldChange, res$padj, log = "y")
```



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



Add some color and nice labels for this plot

```
# 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(0.05), col="gray", lty=2)
```

