

class15

Jibin (PID: A53300326)

2021/11/17

load the countData and colData

We need 2 things - countData - colData

```
library(BiocManager)
library(DESeq2)
```

```
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
```

Side-note: Let's check the corepondence of the metadata and count data setup.

```
metadata$id
```

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

```
colnames(counts)
```

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

We can use the `==` thing to see if they are the same

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

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

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

```
## [1] FALSE
```

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

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

Compare control to treated

First we need to access all the control columns in our counts data.

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

```
## [1] "SRR1039508" "SRR1039512" "SRR1039516" "SRR1039520"
```

Use these ids to access just the control columns of our counts data

```
head(counts[, control.inds])
```

```
##           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
```

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

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

Do the same for drug treated

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

```
## [1] "SRR1039509" "SRR1039513" "SRR1039517" "SRR1039521"
```

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

```
## ENSG000000000003 ENSG000000000005 ENSG000000000419 ENSG000000000457 ENSG000000000460  
##           658.00           0.00           546.00           316.50           78.75  
## ENSG0000000000938  
##           0.00
```

we will combine our means count data for bookkeeping purposes

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

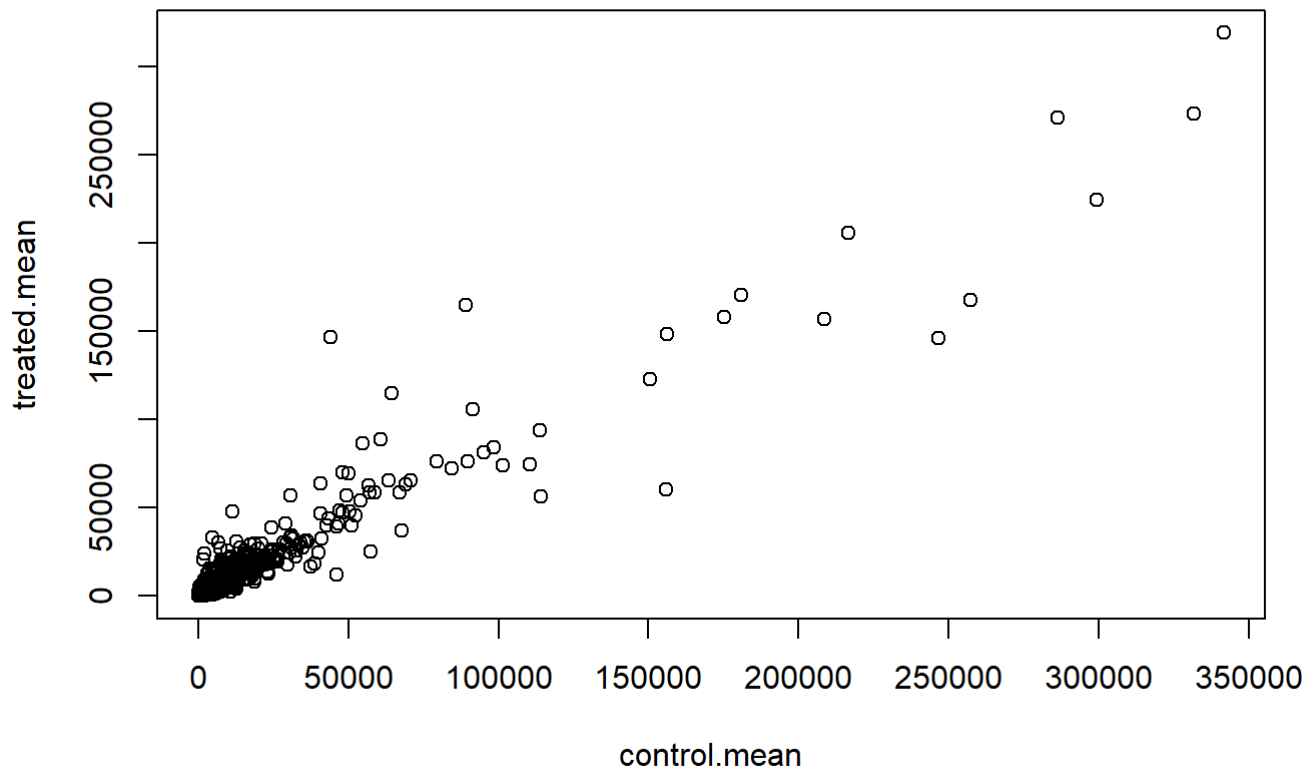
There are 38694 in this dataset

```
nrow(counts)
```

```
## [1] 38694
```

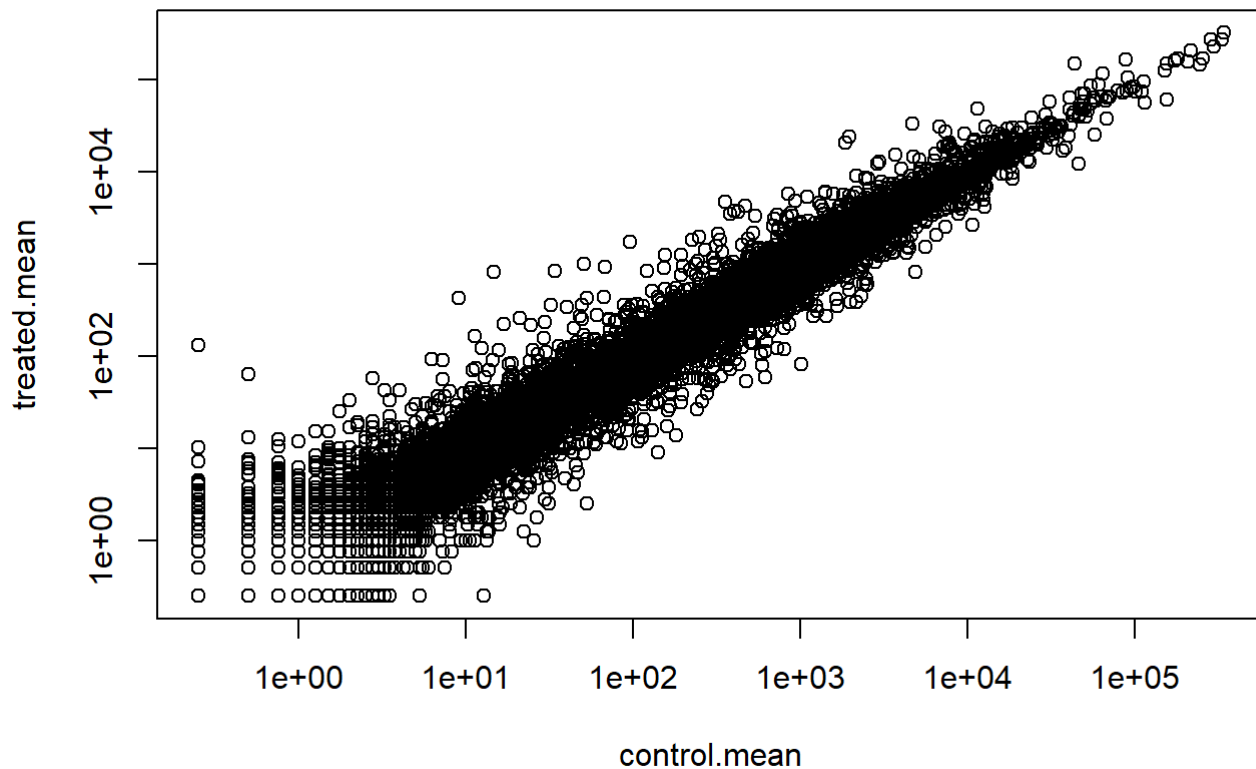
Compare the control and treated

```
plot(meancounts)
```



This would benefit from a log transform! Let's plot on a log scale

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



We often use log trasforamtions as they make life much nicer in this world...

```
log2(40/20)
```

```
## [1] 1
```

Cool. I like log2!

```
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
```

The `which()` function tells us the indices of TRUE netries in a logical vector.

```
which(c(T, F, T))
```

```
## [1] 1 3
```

```
zero.vals <- which(meancounts[,1:2]==0, arr.ind=TRUE)

to.rm <- unique(zero.vals[,1])
mycounts <- meancounts[-to.rm,]
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
```

```
nrow(mycounts)
```

```
## [1] 21817
```

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

```
## [1] 250
```

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

```
## [1] 367
```

What the percentage is this?

```
round(sum(mycounts$log2fc > 2)/nrow(mycounts)*100, 2)
```

```
## [1] 1.15
```

DESeq2 analysis

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

```
##
## 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)
##
## LaTeX的用户的BibTeX条目是
##
## @Article{,
##   title = {Moderated estimation of fold change and dispersion for RNA-seq data with DESeq
2},
##   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)

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
```

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

```
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
```

We can summarize some basic tallies using the summary function.

```
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
```

If the adjusted p value cutoff will be a value other than 0.1, alpha should be set to that value:

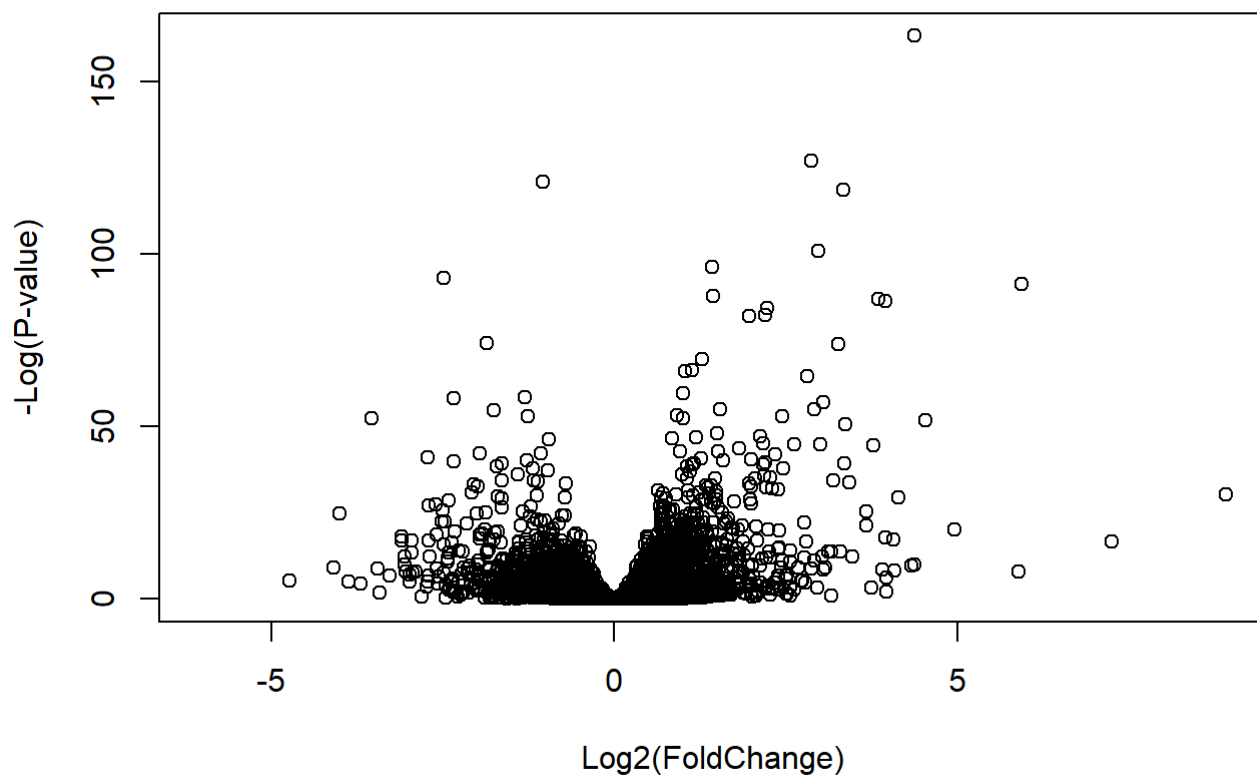
```
res05 <- results(dds, alpha=0.05)
summary(res05)
```

```
##
## 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
```

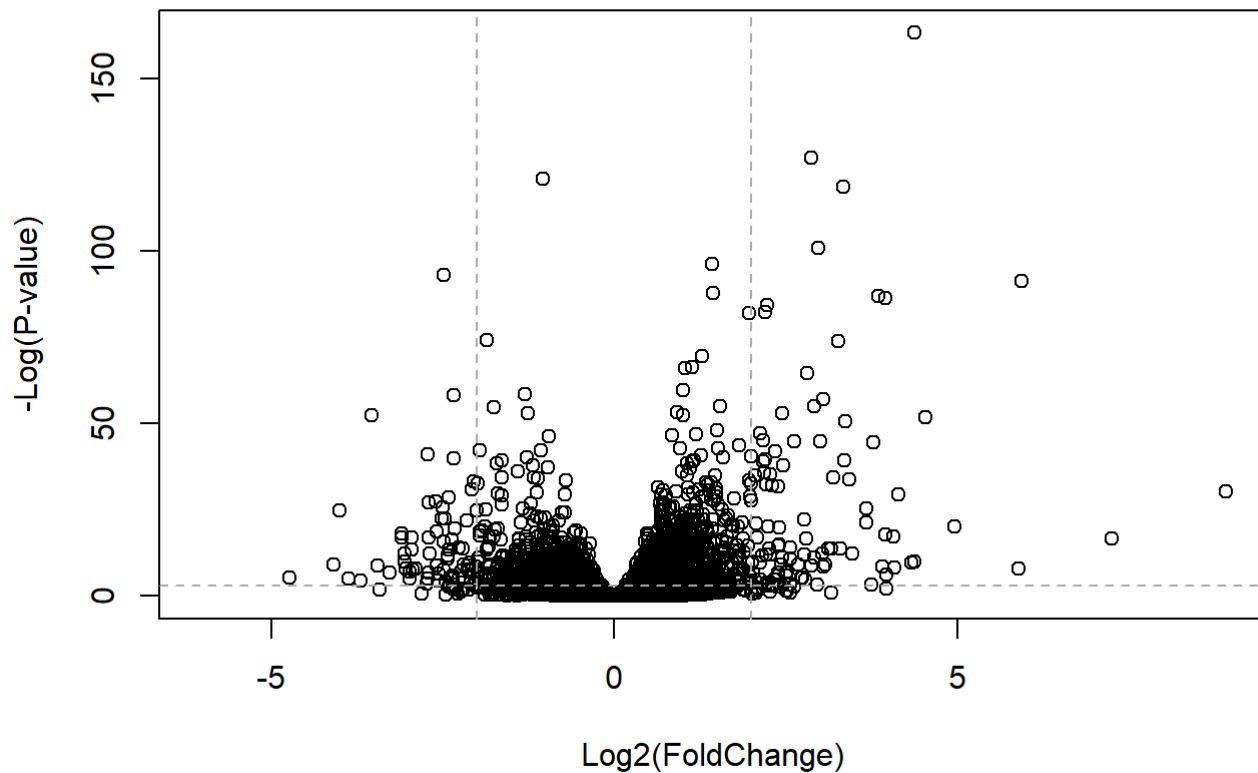
A volcano plot

this is a very common data viz of this

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



```
plot( res$log2FoldChange, -log(res$padj),  
      ylab="-Log(P-value) ", xlab="Log2(FoldChange) ")  
  
# Add some cut-off lines  
abline(v=c(-2,2), col="darkgray", lty=2)  
abline(h=-log(0.05), col="darkgray", lty=2)
```

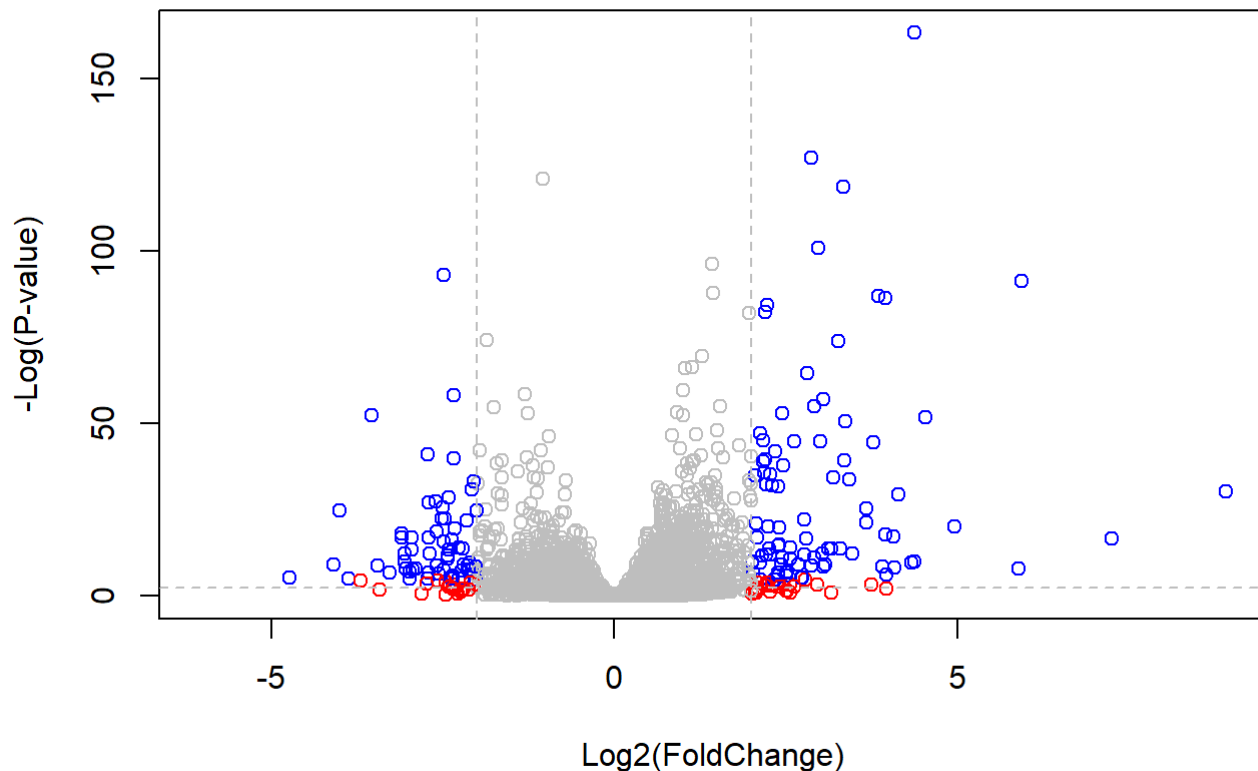


```
# Setup our custom point color vector
mycols <- rep("gray", nrow(res))
mycols[ abs(res$log2FoldChange) > 2 ] <- "red"

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)
```

Adding annotation data

We want to add meaningful gene names to our dataset so we can make some sense of what is going on here

For this we will use two bioconductor packages, one does the work and is called **AnnotationDbi** and the other contains the data we are going to map between and is called **org.Hs.eg.db**

```
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"       "UCSCCKG"
## [26] "UNIPROT"
```

```
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")
```

```
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    Clorf112
## ENSG000000000938         NA      FGR
```

```
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 7 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
##           <numeric> <character>
## ENSG00000152583 1.32441e-71    SPARCL1
## ENSG00000179094 6.13966e-56      PER1
## ENSG00000116584 3.49776e-53    ARHGEF2
## ENSG00000189221 3.46227e-52      MAOA
## ENSG00000120129 1.59454e-44      DUSP1
## ENSG00000148175 1.83034e-42      STOM
```

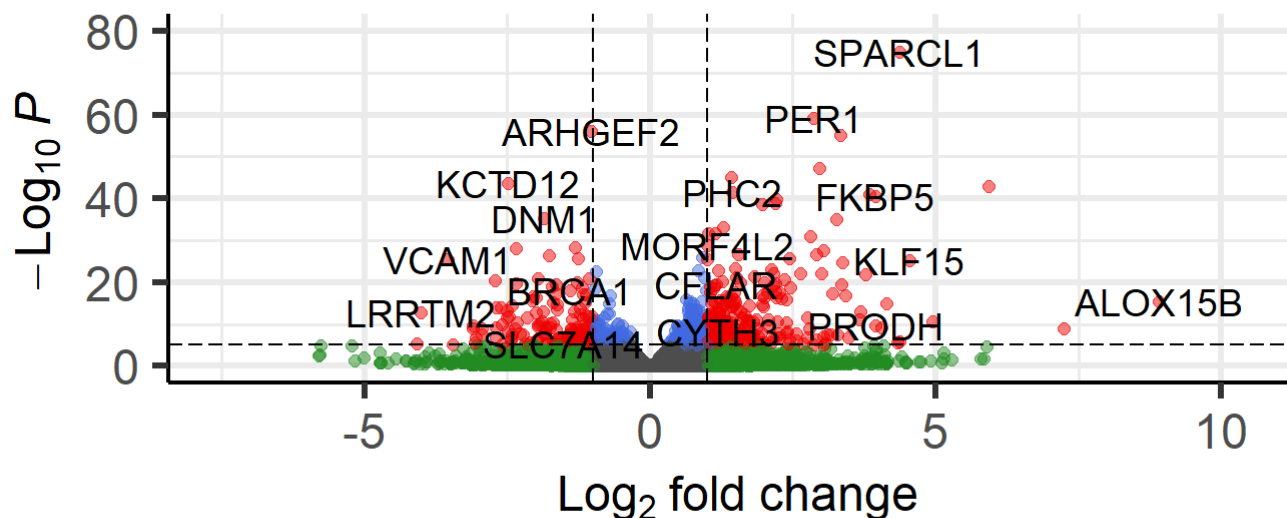
```
library(EnhancedVolcano)
x <- as.data.frame(res)

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

Volcano plot

EnhancedVolcano

● NS ● $\text{Log}_2 \text{FC}$ ● p-value ● p – value and $\text{log}_2 \text{FC}$



total = 38694 variables

Let's finally save our results to data

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

#Pathway Analysis

Let's try to bring some biology insights back into this work

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

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"
```

Before we can use KEGG we need to get our gene identifiers in the correct format for KEGG, which is ENTREZ format in this case.

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

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

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

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      060762
## ENSG000000000457  0.961694      SCYL3      57147      Q8IZE3
## ENSG000000000460  0.815849      Clorf112     55732      AOA024R922
## ENSG000000000938         NA      FGR      2268      P09769
##           genename
##           <character>
## ENSG000000000003      tetraspanin 6
## ENSG000000000005      tenomodulin
## ENSG000000000419 dolichyl-phosphate m..
## ENSG000000000457 SCY1 like pseudokina..
## ENSG000000000460 chromosome 1 open re..
## ENSG000000000938 FGR proto-oncogene, ..
```

Assign names to this vector that are the gene IDs that KEGG wants.

```
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)
```

We can look at the attributes() of this or indeed any R object.

```
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      expl
## 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")
```

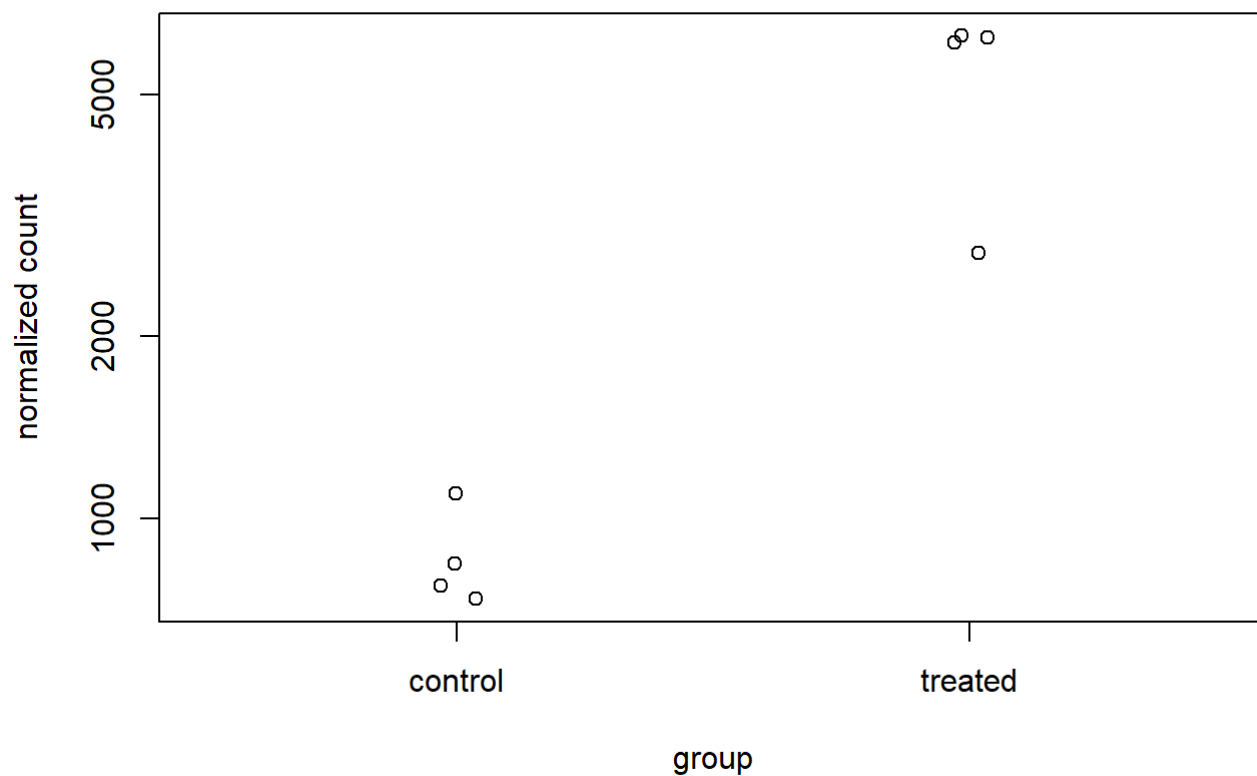
! (hsa05310.pathview.png)

Plotting counts for genes of interest

```
i <- grep("CRISPLD2", res$symbol)
res[i,]
```

```
## log2 fold change (MLE): dex treated vs control
## Wald test p-value: dex treated vs control
## DataFrame with 1 row and 10 columns
##           baseMean log2FoldChange    lfcSE      stat      pvalue
##           <numeric>      <numeric> <numeric> <numeric>  <numeric>
## ENSG00000103196    3096.16         2.62603  0.267444   9.81899 9.32747e-23
##           padj      symbol      entrez      uniprot
##           <numeric> <character> <character> <character>
## ENSG00000103196 3.36344e-20    CRISPLD2      83716    A0A140VK80
##           genename
##           <character>
## ENSG00000103196 cysteine rich secret..
```

```
plotCounts(dds, gene="ENSG00000103196", intgroup="dex")
```


ENSG00000103196

```
d <- plotCounts(dds, gene="ENSG00000103196", intgroup="dex", returnData=TRUE)
head(d)
```

```
##           count      dex
## SRR1039508  774.5002 control
## SRR1039509 6258.7915 treated
## SRR1039512 1100.2741 control
## SRR1039513 6093.0324 treated
## SRR1039516  736.9483 control
## SRR1039517 2742.1908 treated
```

```
library(ggplot2)
ggplot(d, aes(dex, count, fill=dex)) +
  geom_boxplot() +
  scale_y_log10() +
  ggtitle("CRISPLD2")
```

