

lab11r

Rachel Kraft

2022-02-24

Section 1

use the following commands to install biocmanager and DESeq2

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

```
library(BiocManager)
library(DESeq2)
library(ggplot2)
```

Section 2: Import countData and colData

read our count data and metadata files

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

take a look

```
head(counts)
```

```
##          SRR1039508 SRR1039509 SRR1039512 SRR1039513 SRR1039516
## ENSG00000000003      723        486       904       445      1170
## ENSG00000000005        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
## ENSG00000000003     1097       806       604
## ENSG00000000005       0         0         0
## ENSG00000000419      781       417       509
## ENSG00000000457      447       330       324
## ENSG00000000460      94        102       74
## ENSG00000000938       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?

```
View(counts)
```

38694

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

```
View(metadata)
```

4

Section 3: Toy differential gene expression

exploratory differential gene expression analysis

find sample id for those labeled control in metadata

```

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

```

```

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

```

or, use dplyr package from tidyverse

```
#install.packages("dplyr")
```

```
library(dplyr)
```

```

control <- metadata %>% filter(dex=="control")
control.counts <- counts %>% select(control$id)
control.mean <- rowSums(control.counts)/4
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?

if more samples were added to metadata these values wouldn't be correct because we know that there are 4 control columns and divide by 4. we could set it to divide by how many control columns there are, not specifically 4

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.mean <- rowSums( counts[, treated$id] )/4  
names(treated.mean) <- counts$ensgene  
head(treated.mean)
```

```
## [1] 658.00 0.00 546.00 316.50 78.75 0.00
```

combine our meancount data

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

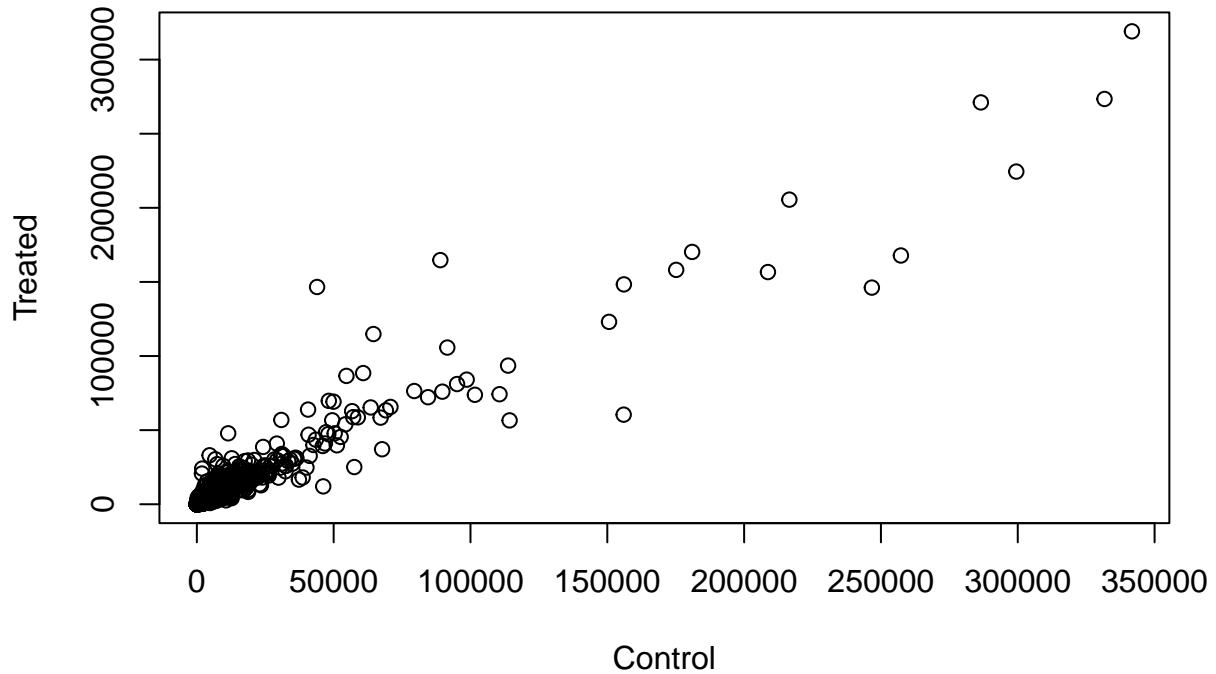
use colSums() to show sum of the mean counts across all genes for each group

```
colSums(meancounts)
```

```
## control.mean treated.mean  
## 23005324 22196524
```

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.

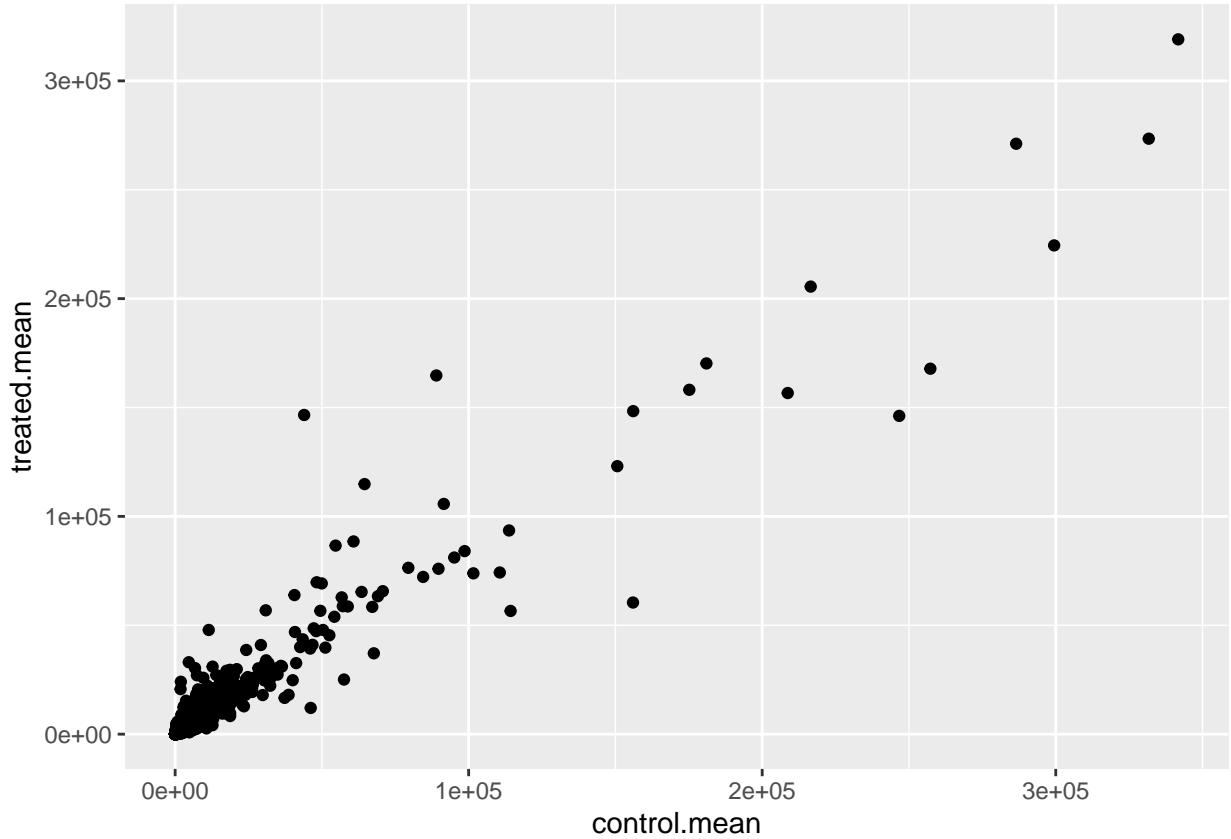
```
plot(meancounts, xlab="Control", ylab="Treated")
```



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?

```
geom_point()
```

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



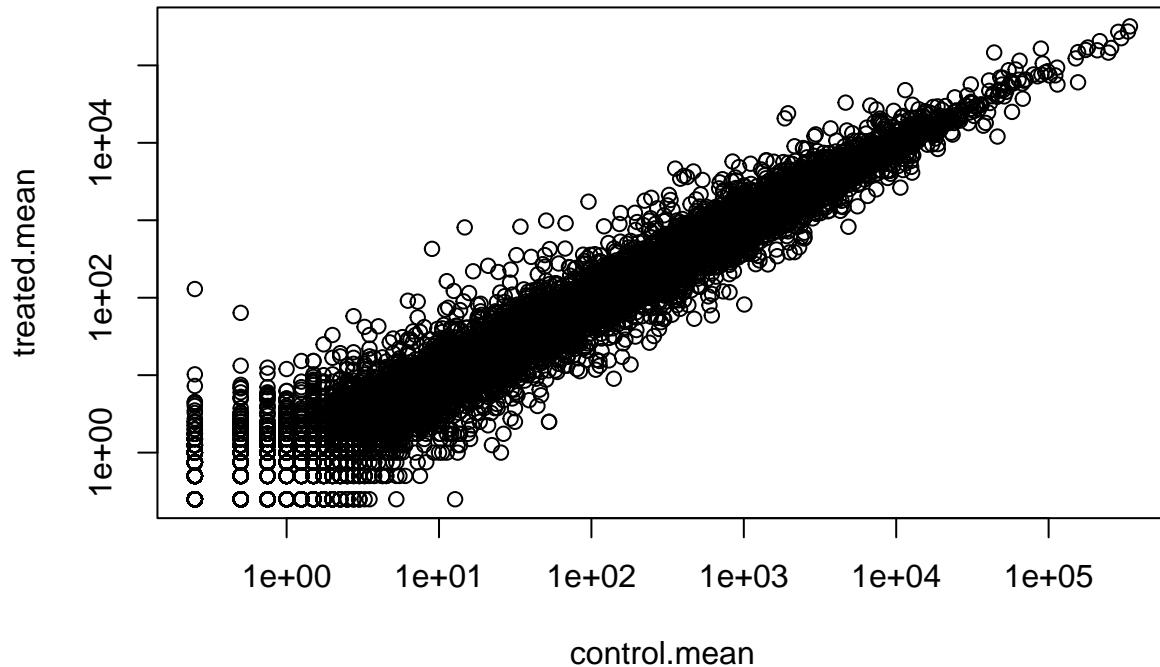
Q6. Try plotting both axes on a log scale. What is the argument to `plot()` that allows you to do this?

```
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
```



find candidate differentially expressed genes by looking for genes with a large change between control and dex-treated samples. calculate log2foldchange, add it to our meancounts data.frame and inspect

```
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
## ENSG00000000419     520.50    546.00  0.06900279
## ENSG00000000457     339.75    316.50 -0.10226805
## ENSG00000000460      97.25     78.75 -0.30441833
## ENSG00000000938      0.75     0.00       -Inf
```

get rid of NaN and -Inf results

```
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
## ENSG00000000419     520.50    546.00  0.06900279
```

```

## ENSG00000000457      339.75      316.50 -0.10226805
## ENSG00000000460      97.25       78.75 -0.30441833
## ENSG00000000971     5219.00     6687.50  0.35769358
## ENSG00000001036     2327.00     1785.75 -0.38194109

```

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?

arr.ind returns the row and column positions where the value is TRUE, or in this case have zero counts. we do this so we can ignore any genes with zero counts. unique() ensures that we don't count any row twice if it has zero entries in both samples

filter our data set with log2(FoldChange) of greater than 2 or less than -2

```

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

```

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

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

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

sum the true values in the vector; there are 250 up regulated genes at the greater than 2 fc level

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

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

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

367 down regulated genes

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

We should not trust these results because there has not been any analysis done if these changes are statistically significant.

DESeq2 analysis

```
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)
##
```

```

## 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},
## }

```

we can construct a DESeqDataSet from (1) a count matrix, (2) a metadata file, and (3) a formula indicating the design of the experiment

Take a look at metadata again. The thing we're interested in is the dex column, which tells us which samples are treated with dexamethasone versus which samples are untreated controls. We'll specify the design with a tilde, like this: design=~dex.

We will use the DESeqDataSetFromMatrix() function to build the required DESeqDataSet object and call it dds, short for our DESeqDataSet

```
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): ENSG00000000003 ENSG00000000005 ... ENSG00000283120
##   ENSG00000283123
## rowData names(0):
## colnames(8): SRR1039508 SRR1039509 ... SRR1039520 SRR1039521
## colData names(4): id dex celltype geo_id

```

Next, let's run the DESeq analysis pipeline on the dataset, and reassign the resulting object back to the same variable

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

```

## estimating size factors

## estimating dispersions

```

```
## gene-wise dispersion estimates  
## mean-dispersion relationship  
## final dispersion estimates
```

getting the results

```
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

```

we can adjust the p-value cut off by resetting alpha

```

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

```

Section 5: adding annotation data

for interpreting results, we need to understand differentially expressed genes. get gene names (gene SYMBOLs)

```
install "AnnotationDbi" install "org.Hs.eg.db"
```

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

```

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

```

What DB identifiers can be looked up

```
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"

```

Use mapIds() function to add indiv columns

```

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
```

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$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 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
## 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      entrez      genename
##           <numeric> <character> <character> <character>
## ENSG000000000003  0.163035    TSPAN6      7105    tetraspanin 6
## ENSG000000000005   NA         TNMD      64102    tenomodulin
## ENSG00000000419   0.176032    DPM1       8813    dolichyl-phosphate m..
## ENSG00000000457   0.961694    SCYL3      57147    SCY1 like pseudokina..
## ENSG00000000460   0.815849    C1orf112    55732    chromosome 1 open re..
## ENSG00000000938   NA         FGR       2268    FGR proto-oncogene, ..

```

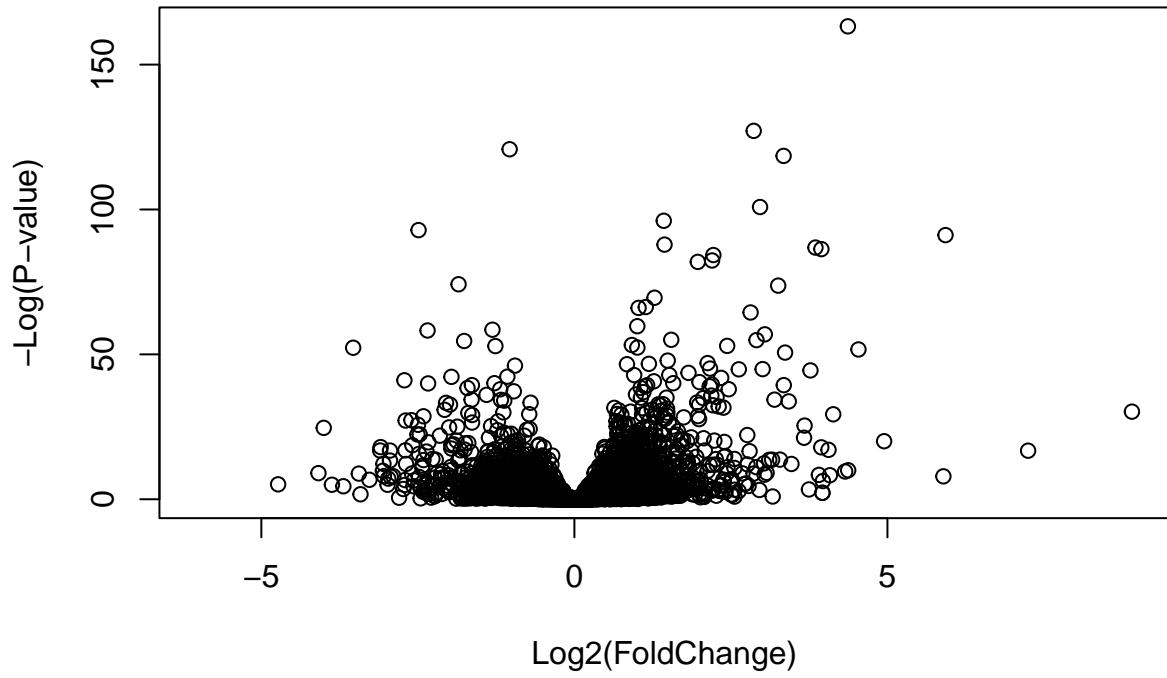
Section 6: Data Visualization

let's make a basic volcano plot

```

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

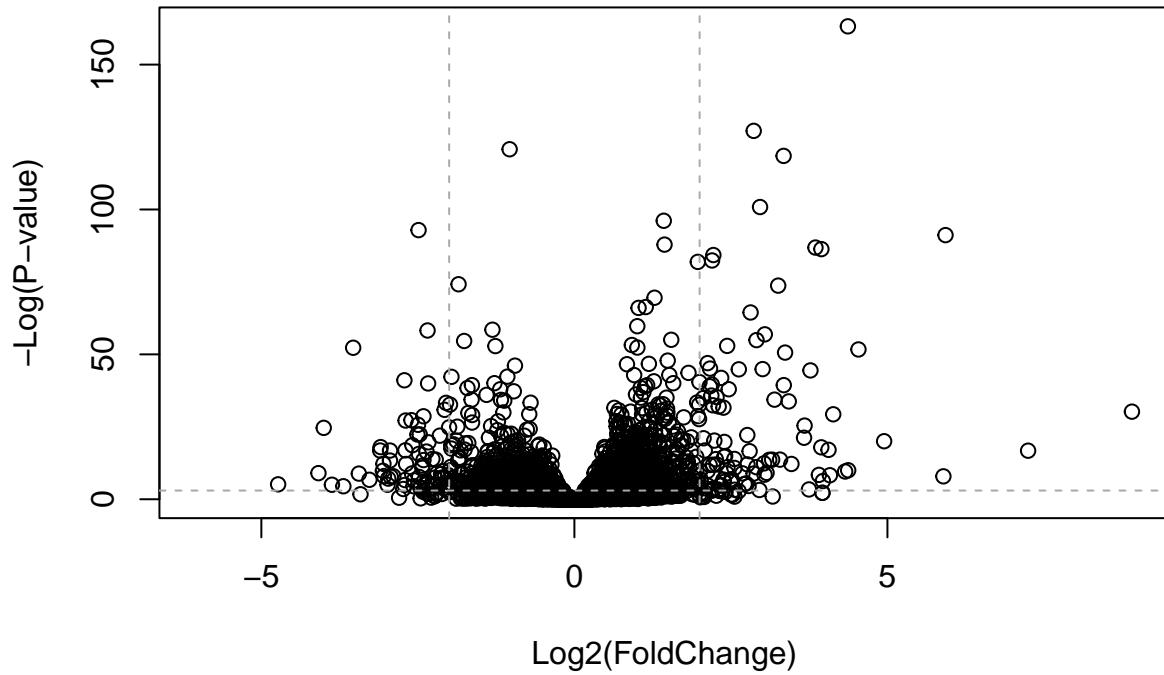
```



To make this more useful we can add some guidelines (with the abline() function) and color (with a custom color vector) highlighting genes that have $\text{padj} < 0.05$ and the absolute $\text{log2FoldChange} > 2$.

```
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)
```



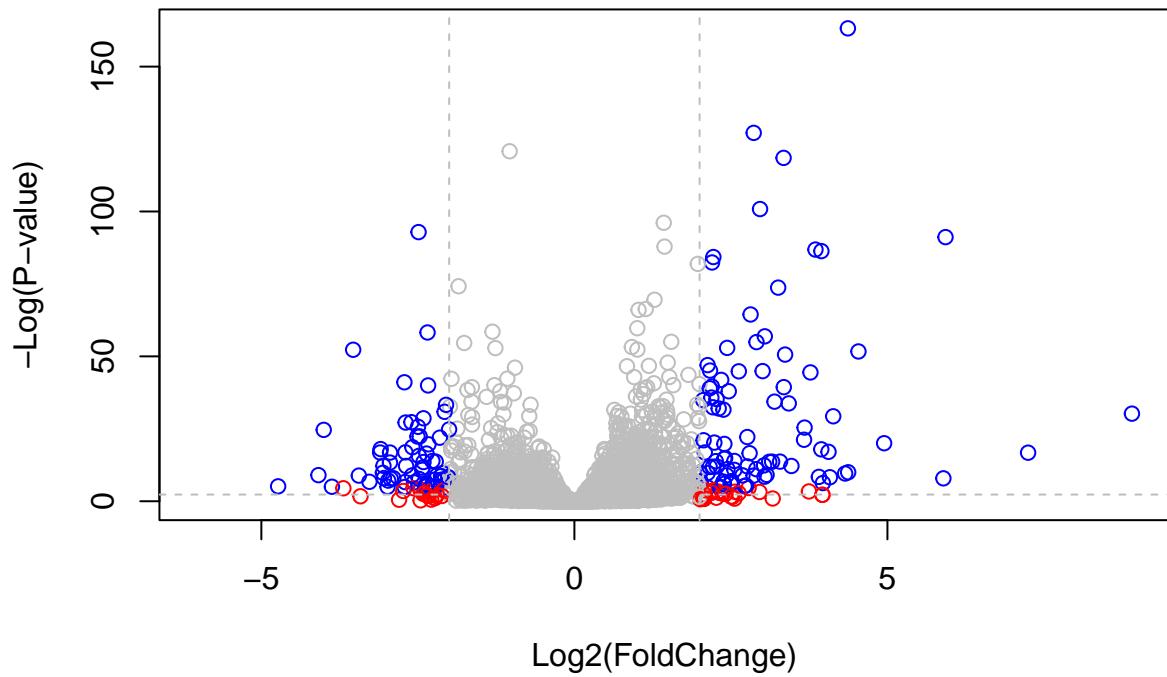
to color the points, we can setup a custom color vector indicating transcripts with large fold change and significant differences between conditions

```
# 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)
```



try enhanced volcano

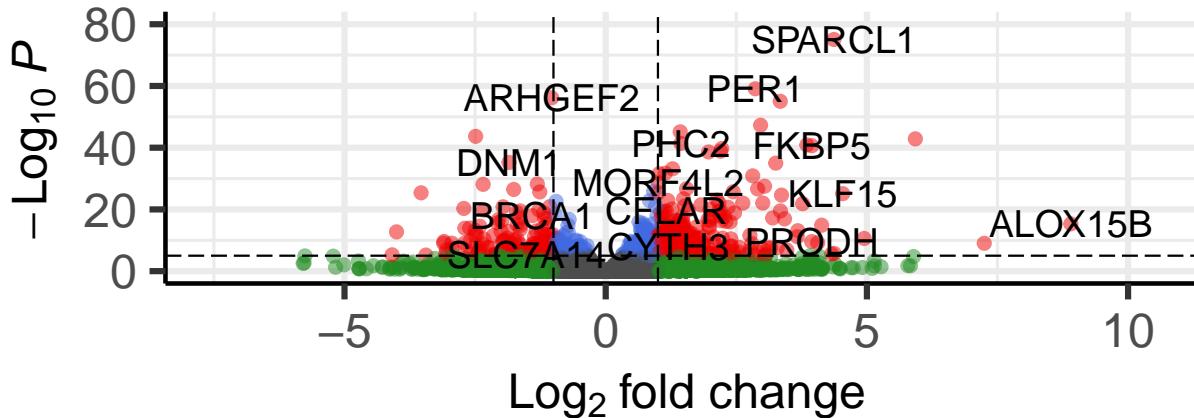
```
library(EnhancedVolcano)

x <- as.data.frame(res)
EnhancedVolcano(x, lab = x$symbol, x = 'log2FoldChange', y = 'pvalue')
```

Volcano plot

EnhancedVolcano

● NS ● Log₂ FC ● p-value ● p-value and log₂ FC



Section 7: Pathway analysis with R and bioconductor

use GAGE package for KEGG pathway enrichment analysis on RNA-seq based differential expression results

install gage package and pathview package

```
BiocManager::install( c("pathview", "gage", "gageData") )
```

```
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"

```

recall that vectors can have names attribute

need vector of fold-change labeled with names of genes in ENTREZ format

```

foldchange <- res$log2FoldChange
names(foldchange) <- res$entrez

```

now can run GAGE analysis passing in our foldchange vector and KEGG genesets we are interested in

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

```
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	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

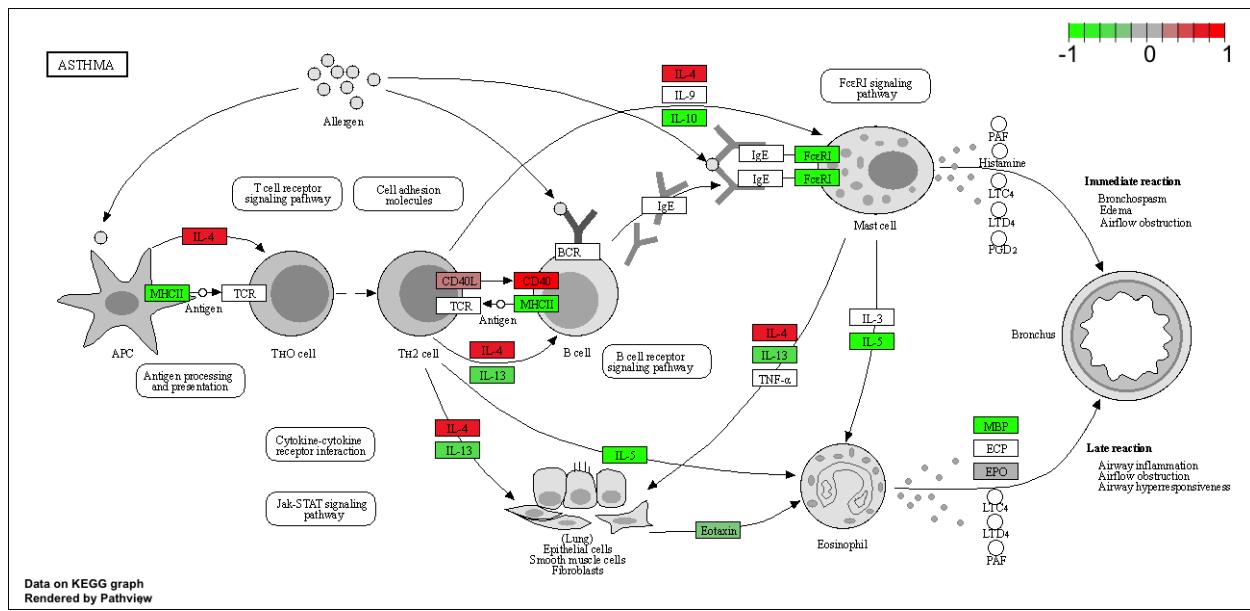
map foldchange results onto any KEGG pathway

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

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

```
## Info: Working in directory /Users/rachel/lab11
```

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



Final step- save results

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