

Class 12: Transcriptomics and the Analysis of RNA-Seq Data

Audrey Nguyen

Important countData and colData

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

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?

```
nrow(counts)
```

```
[1] 38694
```

There are 38694 genes in this dataset.

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

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

```
control treated
        4       4
```

There are 4 ‘control’ cell lines.

Let’s check to see if the id names in metadata match the order of the columns in the count-Data.

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

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

We can use the `all()` function to check that all its inputs are true.

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

```
[1] TRUE
```

Let’s extract only the control columns.

```

control inds <- metadata$dex == "control"
control ids <- metadata$id[control inds]
control counts <- counts[, control ids]
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

I want a single summary counts value for each gene in the control experiments. I will start by taking the average.

```

# you can also use apply(control.counts, 1, mean)
control.mean <- rowMeans(control.counts)

```

We need to do the same thing to get the `treated.mean`.

```

treated inds <- metadata$dex == "treated"
treated ids <- metadata$id[treated inds]
treated counts <- counts[, treated ids]
head(treated counts)

```

	SRR1039509	SRR1039513	SRR1039517	SRR1039521
ENSG000000000003	486	445	1097	604
ENSG000000000005	0	0	0	0
ENSG000000000419	523	371	781	509
ENSG000000000457	258	237	447	324
ENSG000000000460	81	66	94	74
ENSG000000000938	0	0	0	0

```

treated.mean <- rowMeans(treated.counts)

```

Toy differential gene expression

This bit of code will first find the sample id for those labeled control. It will then calculate the mean counts per gene across the samples.

```
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
ENSG000000000938
    0.75
```

An alternative way to do this same thing using the `dplyr` package from the tidyverse is shown below. Which do you prefer and why?

```
library(dplyr)
```

```
Attaching package: 'dplyr'
```

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

```
filter, lag
```

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

```
intersect, setdiff, setequal, union
```

```
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
ENSG000000000938
    0.75
```

I like using the first method better, because I don't like typing the `%>%` in `dplyr`.

Q3. How would you make the above code in either approach more robust?

You can use the function `apply()` to include more functions.

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`)

See above code for `treated.mean`.

We will combine our meancount data for bookkeeping purposes.

```
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

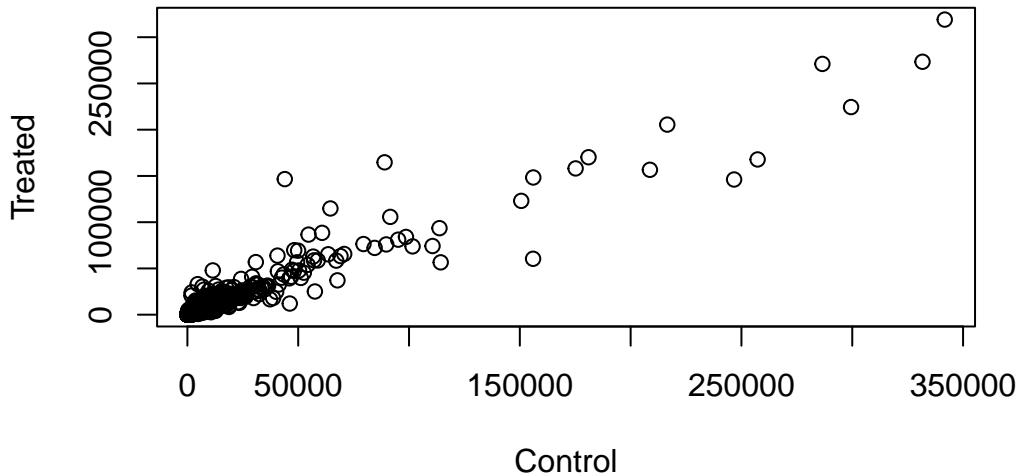
`colSums()` the data to show the sum of the mean counts across all genes for each group.

```
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.

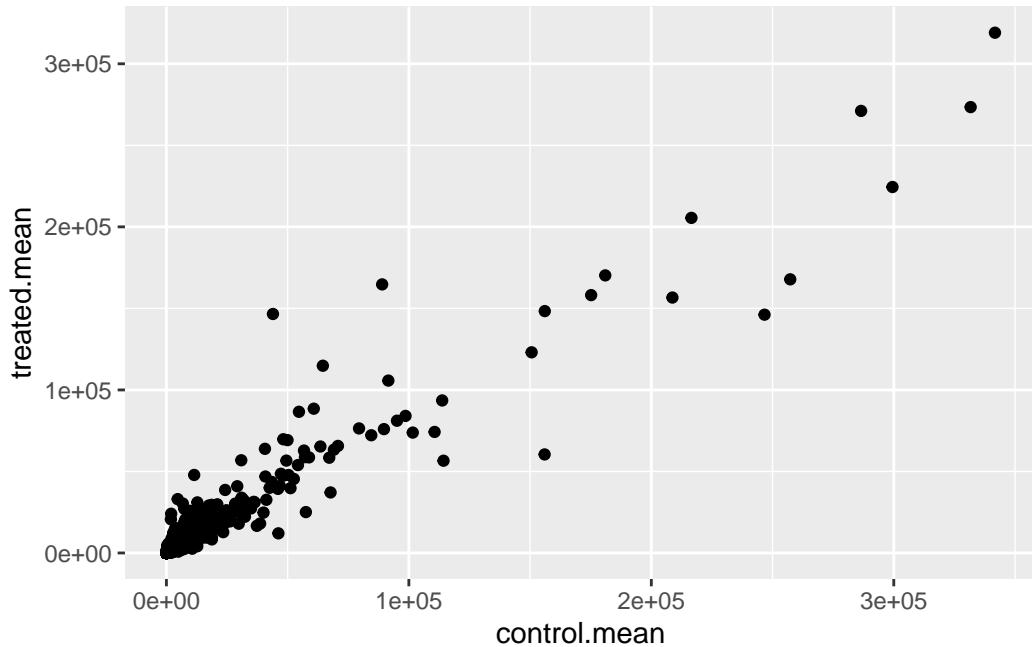
```
plot(meancounts$control.mean, meancounts$treated.mean, xlab = "Control", ylab = "Treated")
```



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?

You would use `geom_point()`.

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



There are too many points around the origin, even though there are 60,000 points. Use a log

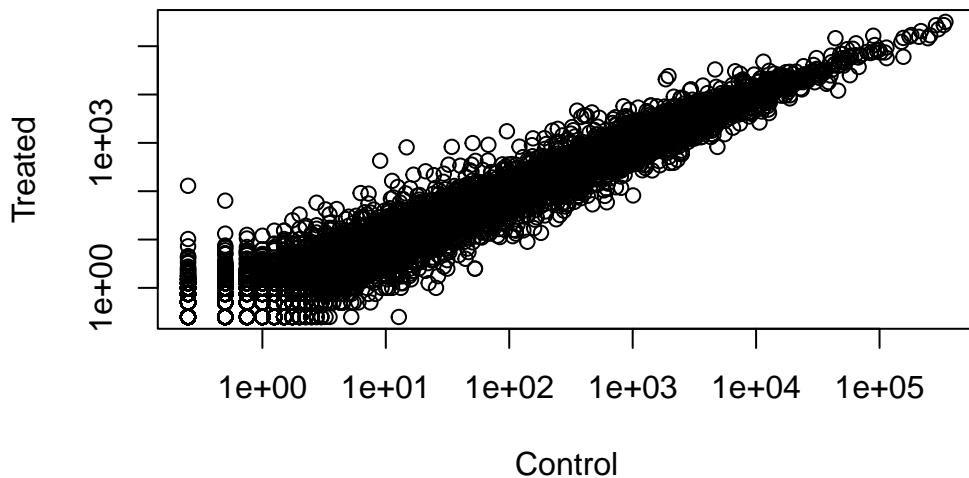
function to visualize more points.

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

```
plot(meancounts$control.mean, meancounts$treated.mean, xlab = "Control", ylab = "Treated",
```

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



The most useful and most straightforward to understand is log2 transform.

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

```
[1] 0
```

0 means no change.

Doubling

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

```
[1] 1
```

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

```
[1] -1
```

```
log2(80/20)
```

```
[1] 2
```

We can calculate the log2fold change, add it to our `meancounts` data.frame and inspect the results either with the `head()` or the `View()` function.

```
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

Let's filter our data to remove the genes with “NaN” (not a number) and -Inf (negative infinity) 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
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

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?

The `arr.ind` argument returns the array indices if `x` is an array. In this case, it returns the zero values. The `unique()` function is used to eliminate or delete the duplicate values or rows present. In this case, it deletes the zero values if it has zero entries in both samples.

A common threshold used for calling something differentially expressed is a $\log_2(\text{FoldChange})$ 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. Using the `up.ind` vector above can you determine how many upregulated genes we have at the greater than 2 fc level?

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

```
[1] 250
```

There are 250 genes that are 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(mycounts$log2fc < -2)
```

```
[1] 367
```

There are 367 downregulated genes that are greater than 2 fc level.

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

Not really, because we haven't really seen if these results are statistically significant, based on p-values.

DESeq2 analysis

```
library(DESeq2)
```

```
Warning: package 'DESeq2' was built under R version 4.2.2
```

```
Warning: package 'S4Vectors' was built under R version 4.2.2
```

```
Warning: package 'BiocGenerics' was built under R version 4.2.1
```

```
Warning: package 'IRanges' was built under R version 4.2.1
```

```
Warning: package 'GenomicRanges' was built under R version 4.2.2
```

```
Warning: package 'GenomeInfoDb' was built under R version 4.2.2
```

```
Warning: package 'SummarizedExperiment' was built under R version 4.2.1
```

```
Warning: package 'MatrixGenerics' was built under R version 4.2.1
```

```
Warning: package 'Biobase' was built under R version 4.2.1
```

```
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},  
}
```

Importing data

```
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
```

DESeq analysis

You have to run **DESeq()** before looking at results.

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

```
estimating size factors
```

```
estimating dispersions
```

```
gene-wise dispersion estimates
```

```
mean-dispersion relationship
```

```
final dispersion estimates
```

```
fitting model and testing
```

Getting 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
```

Adjusting the p-value to 0.05. Its default value is 0.1.

```
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
```

Adding annotation data

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

```
Warning: package 'AnnotationDbi' was built under R version 4.2.1
```

```
Attaching package: 'AnnotationDbi'
```

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

```
select
```



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



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



```
[1] "ACCCNUM"          "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"
```

The main function we will use from the **AnnotationDbi** package is called **mapIds()**. We can use the **mapIds()** function to add individual columns to our results table.

```
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	C1orf112			
ENSG000000000938	NA	FGR			

Q11. Run the `mapIds()` function 2 more times to add the Entrez ID and UniProt accession and GENENAME a new columns called `res$entrez`, `res$uniprot`, and `res$genename`.

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

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

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

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

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

```

    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>
ENSG00000000003 747.194195 -0.3507030  0.168246 -2.084470 0.0371175
ENSG00000000005  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      uniprot
  <numeric> <character> <character> <character>
ENSG00000000003 0.163035    TSPAN6      7105  AOA024RCI0
ENSG00000000005  NA        TNMD       64102  Q9H2S6
ENSG00000000419  0.176032    DPM1       8813  O60762
ENSG00000000457  0.961694    SCYL3      57147  Q8IZE3
ENSG00000000460  0.815849    C1orf112   55732  AOA024R922
ENSG00000000938  NA        FGR        2268   P09769
  genename
  <character>
ENSG00000000003      tetraspanin 6
ENSG00000000005      tenomodulin
ENSG00000000419      dolichyl-phosphate m..
ENSG00000000457      SCY1 like pseudokina..
ENSG00000000460      chromosome 1 open re..
ENSG00000000938      FGR proto-oncogene, ..

```

You can arrange and view the results by the adjusted p-value.

```

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  Q15534
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

```

Finally, let's write out the ordered significant results with annotations.

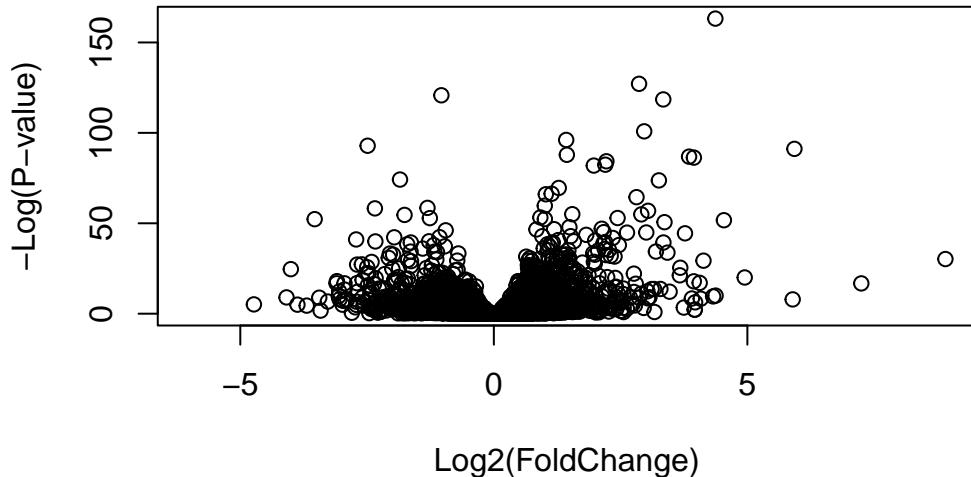
```
write.csv(res[ord, ], "deseq_results.csv")
```

Data Visualization

Volcano plots

Volcano plots are frequently used to highlight the proportion of genes that are both significantly regulated and display a high fold change.

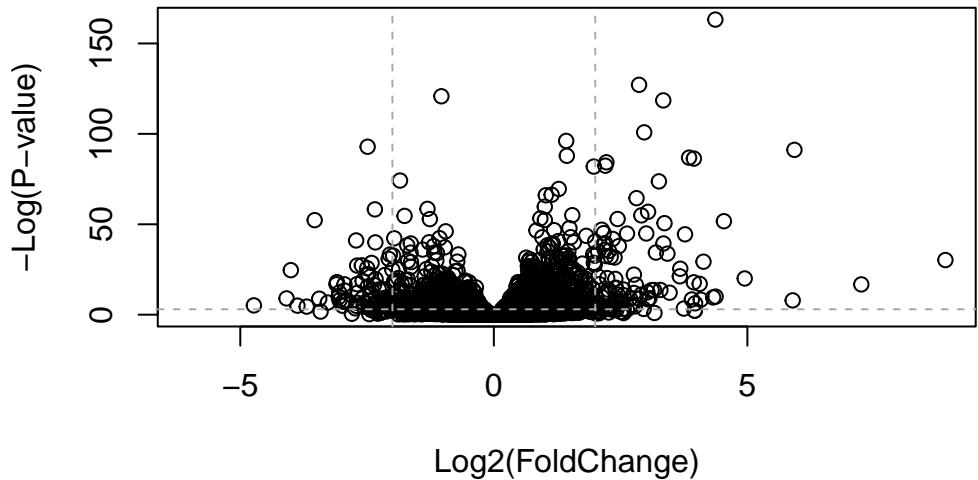
```
plot( res$log2FoldChange, -log(res$padj),
      xlab="Log2(FoldChange)",
      ylab="-Log(P-value)")
```



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 will setup a custom color vector indicating transcripts with large fold change and significant differences between conditions:

```

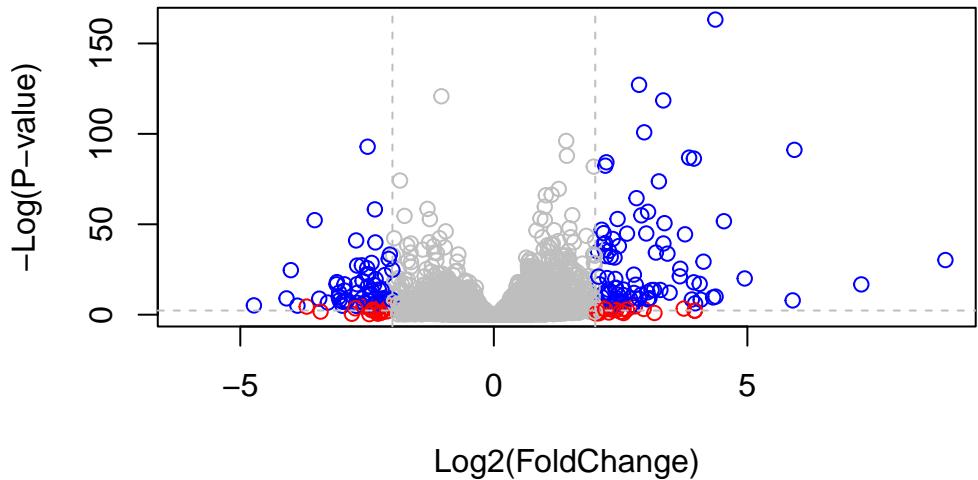
# Setup our custom point color vector
mycols <- rep("gray", nrow(res))
mycols[ abs(res$log2FoldChange) > 2 ] <- "red"

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)

```



You can use the **EnhancedVolcano** bioconductor package for more customization.

```
BiocManager::install("Enhanced Volcano")
```

```
Bioconductor version 3.16 (BiocManager 1.30.19), R 4.2.0 (2022-04-22)
```

```
Installing package(s) 'Enhanced Volcano'
```

```
Warning: package 'Enhanced Volcano' is not available for Bioconductor version '3.16'
```

```
A version of this package for your version of R might be available elsewhere,  
see the ideas at
```

```
https://cran.r-project.org/doc/manuals/r-patched/R-admin.html#Installing-packages
```

```
Installation paths not writeable, unable to update packages
```

```
path: /Library/Frameworks/R.framework/Versions/4.2/Resources/library  
packages:
```

```
boot, brew, bslib, callr, class, cli, cluster, codetools, colorspace,  
commonmark, cpp11, crayon, curl, desc, devtools, digest, evaluate, fansi,  
farver, foreign, fs, gert, ggplot2, gh, gitcreds, gtable, highr, htmltools,  
httr, isoband, jsonlite, knitr, lifecycle, markdown, MASS, Matrix, mgcv,  
nlme, nnet, openssl, pillar, pkgbuild, pkgload, processx, ps, purrr, rlang,  
rmarkdown, roxygen2, rpart, rstudioapi, rversions, sass, scales, spatial,  
stringi, stringr, survival, sys, testthat, tibble, tinytex, utf8, vctrs,  
viridisLite, whisker, xfun, yaml, zip
```

```
Old packages: 'data.table', 'RSQLite'
```

```
library(EnhancedVolcano)
```

Warning: package 'EnhancedVolcano' was built under R version 4.2.1

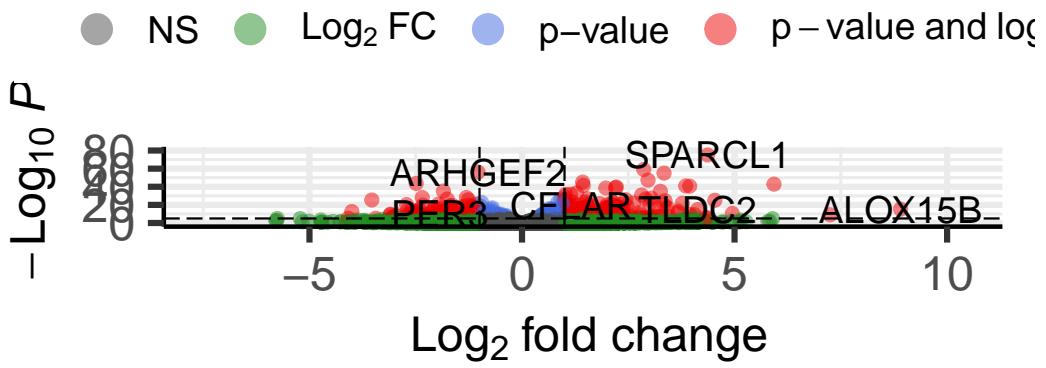
Loading required package: ggrepel

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

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

Volcano plot

EnhancedVolcano



total = 38694 variables

Pathway analysis

```
library(pathview)
```

Warning: package 'pathview' was built under R version 4.2.1

```
#####
# 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)
```

```
Warning: package 'gage' was built under R version 4.2.1
```

```
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"
```

```
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

Running the **gage** pathway analysis:

```
# get the results
keggres = gage(foldchanges, 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
```

Let's try out the **pathview()** function from the **pathview** package to make a pathway plot with our RNA-Seq expression results shown in color.

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

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

Info: Working in directory /Users/audreynguyen/Desktop/BIMM 143/class12

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

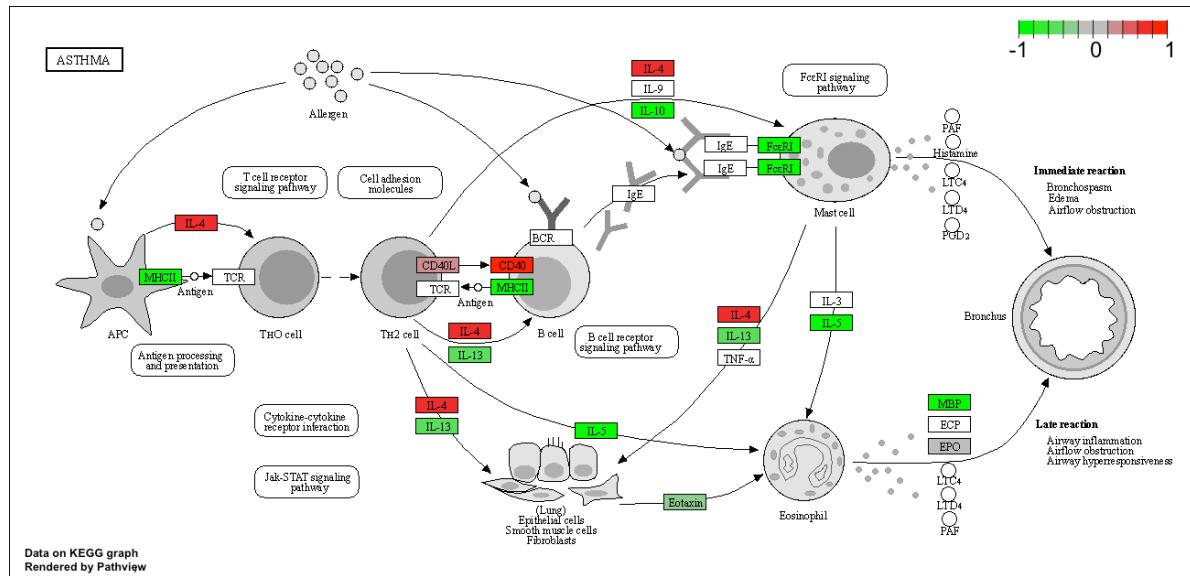


Figure 1: The Asthma pathway with my highlighted differentially expressed genes in color