

# Class 12: Transcriptomics and the analysis of RNA-Seq data

Maria Tavares

## Table of contents

Background . . . . .	1
Data Import . . . . .	2
Toy analysis example . . . . .	3
Filter out zeero count genes . . . . .	6
DESeq analysis . . . . .	8
Volcano Plot . . . . .	9
Save our results . . . . .	10
A nicer ggplot volcano plot . . . . .	10
Add annotation data . . . . .	11
Save my annotated results . . . . .	14
Pathway analysis . . . . .	14

## Background

Today we will analyze some RNAseq data from Himes et al. on the effects of a common steroid (dexamethasone also called “dex”) on airway smooth muscle cells (ASMs).

For this analysis we need to main inputs

- **countData**: a table of **counts** per gene(in rows) accross experiments (in columns) (amount of genes in rows and experiments in columns)
- **colData**: **metadata** about the design of the experiments. The rows match the columns in **countData**

## Data Import

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

	SRR1039508	SRR1039509	SRR1039512	SRR1039513	SRR1039516
ENSG00000000003	723	486	904	445	1170
ENSG00000000005	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
ENSG00000000003	1097	806	604
ENSG00000000005	0	0	0
ENSG000000000419	781	417	509
ENSG000000000457	447	330	324
ENSG000000000460	94	102	74
ENSG000000000938	0	0	0

and the metadata

```
metadata <- read.csv("airway_metadata.csv")
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
```

Q2. How many experiments (i.e columns are in `counts` or rows in `metadata`) are there?

```
ncol(counts)
```

```
[1] 8
```

Q3. How many “control” experiments are there in the dataset?

```
sum(metadata$dex == "control")
```

```
[1] 4
```

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

### Toy analysis example

1. Extract the “control” columns from `counts`
2. Calculate the mean value for each gene in these “control” columns

3-4. Do the same for the “treated” columns 5. Compare these mean values for each gene

Step 1.

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

Step 2.

```
control.mean <- rowMeans(control.counts)
```

Step 3.

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

Step 4.

```
treated.mean <- rowMeans(treated.counts)
```

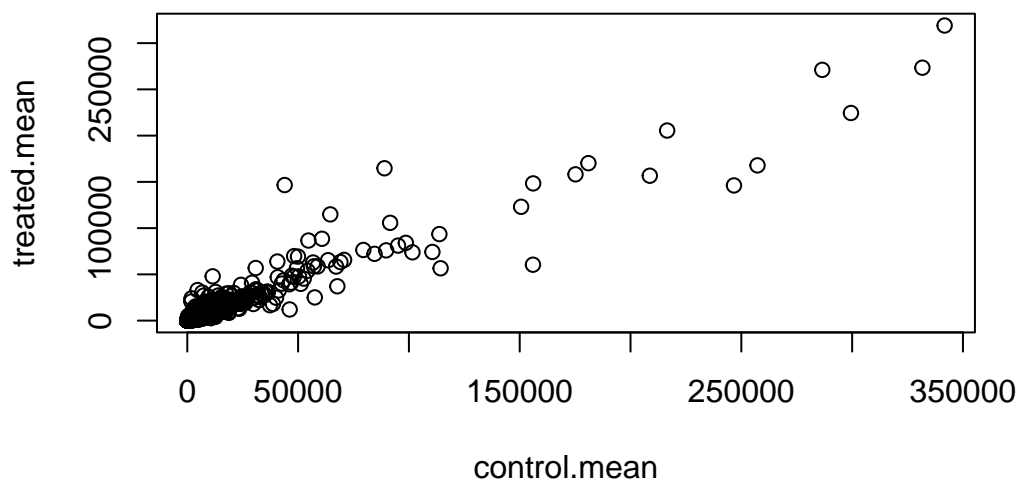
For ease of book-keeping we can store these together in one data frame called `meancounts`

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

	control.mean	treated.mean
ENSG000000000003	900.75	658.00
ENSG000000000005	0.00	0.00
ENSG000000000419	520.50	546.00
ENSG000000000457	339.75	316.50
ENSG000000000460	97.25	78.75
ENSG000000000938	0.75	0.00

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?

`geom_point`

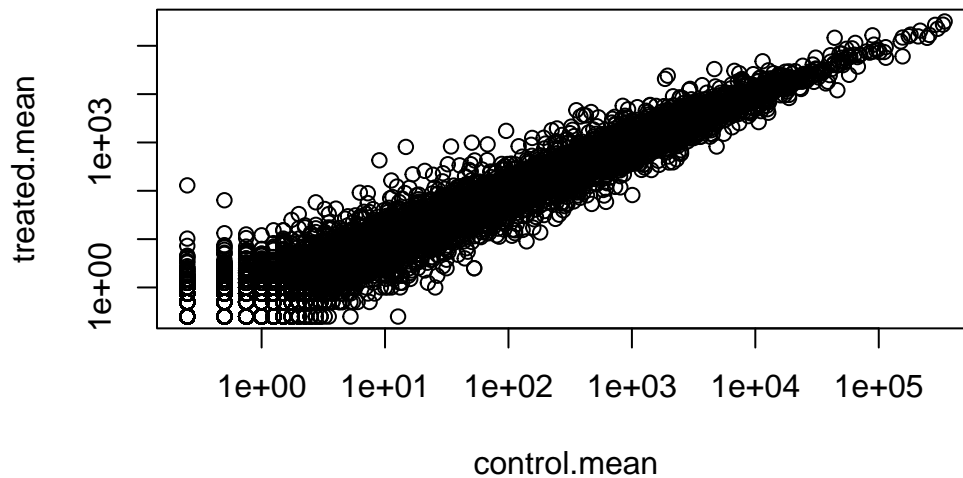
This is screaming at me to log transform this data!

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  $\leq 0$  omitted from logarithmic plot

Warning in `xy.coords(x, y, xlabel, ylabel, log)`: 15281 y values  $\leq 0$  omitted from logarithmic plot



We use  $\log_2$  “fold-change” as a way to compare

```
#treated/control  
log2(10/10) #no change
```

```
[1] 0
```

```
log2(20/10) #up regulated
```

```
[1] 1
```

```
log2(10/20) #down regulated
```

```
[1] -1
```

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

```
[1] 2
```

```
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
ENSG0000000000419	520.50	546.00	0.06900279
ENSG0000000000457	339.75	316.50	-0.10226805
ENSG0000000000460	97.25	78.75	-0.30441833
ENSG0000000000938	0.75	0.00	-Inf

### Filter out zero count genes

```
nonzero.inds <- rowSums(counts) != 0
mycounts <- meancounts [nonzero.inds,]
```

```
zero.inds <- which(meancounts[,1:2] == 0, arr.ind = T)[,1]
mygenes <- meancounts[-zero.inds,]
```

```
y <- data.frame(a=c(1,5,0,5), b=c(1,0,5,5))
y
```

	a	b
1	1	1
2	5	0
3	0	5
4	5	5

```
which(y==0, arr.ind=TRUE)
```

```
      row col  
[1,]    3  1  
[2,]    2  2
```

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 = TRUE` option makes `which()` return a matrix of (row, column) positions of zeros instead of a single index vector. We take the first column to get the row numbers of those zeros, and use `unique()` so each row is only listed once before removing it.

A common “rule-of-thumb” threshold for calling something “up” regulated is a log2-fold-change of +2 or greater. For “down” regulated -2 or less.

Q8. How many genes are “up” regulated at the +2 log2FC threshold?

```
sum(mygenes$log2fc >= 2)
```

```
[1] 314
```

Q9. How many genes are “down” regulated at the -2 log2FC threshold?

```
sum(mygenes$log2fc <= -2)
```

```
[1] 485
```

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

Can’t trust it yet until we make sure the differences are statistically significant with further analysis.

## DESeq analysis

Let's do this with DESeq2 and put some stats behind these numbers.

```
library(DESeq2)
```

Warning: package 'IRanges' was built under R version 4.4.2

Warning: package 'GenomeInfoDb' was built under R version 4.4.2

Warning: package 'MatrixGenerics' was built under R version 4.4.2

DESeq wants 3 things for analysis, countData, colData and design.

```
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

The main function in the DESeq package to run analysis is called DESeq().

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

estimating size factors

estimating dispersions

gene-wise dispersion estimates

mean-dispersion relationship

final dispersion estimates

fitting model and testing



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

Get the results of this DESeq object with the function `results()`

```
res <- results(dds)
head(res)
```

log2 fold change (MLE): dex treated vs control

Wald test p-value: dex treated vs control

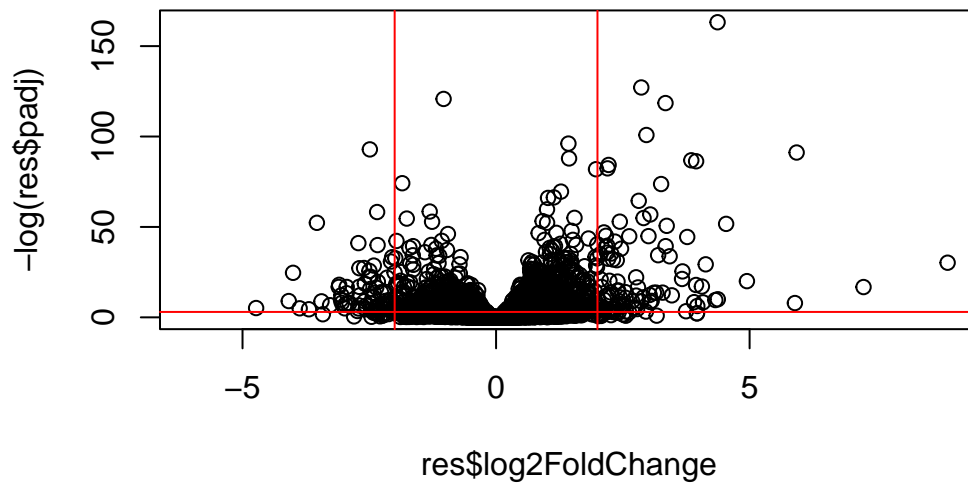
DataFrame with 6 rows and 6 columns

	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
ENSG0000000000419	520.134160	0.2061078	0.101059	2.039475	0.0414026
ENSG0000000000457	322.664844	0.0245269	0.145145	0.168982	0.8658106
ENSG0000000000460	87.682625	-0.1471420	0.257007	-0.572521	0.5669691
ENSG0000000000938	0.319167	-1.7322890	3.493601	-0.495846	0.6200029
	padj				
	<numeric>				
ENSG000000000003	0.163035				
ENSG000000000005	NA				
ENSG0000000000419	0.176032				
ENSG0000000000457	0.961694				
ENSG0000000000460	0.815849				
ENSG0000000000938	NA				

## Volcano Plot

This is a plot of log2FC vs adjusted p-value

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



### Save our results

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

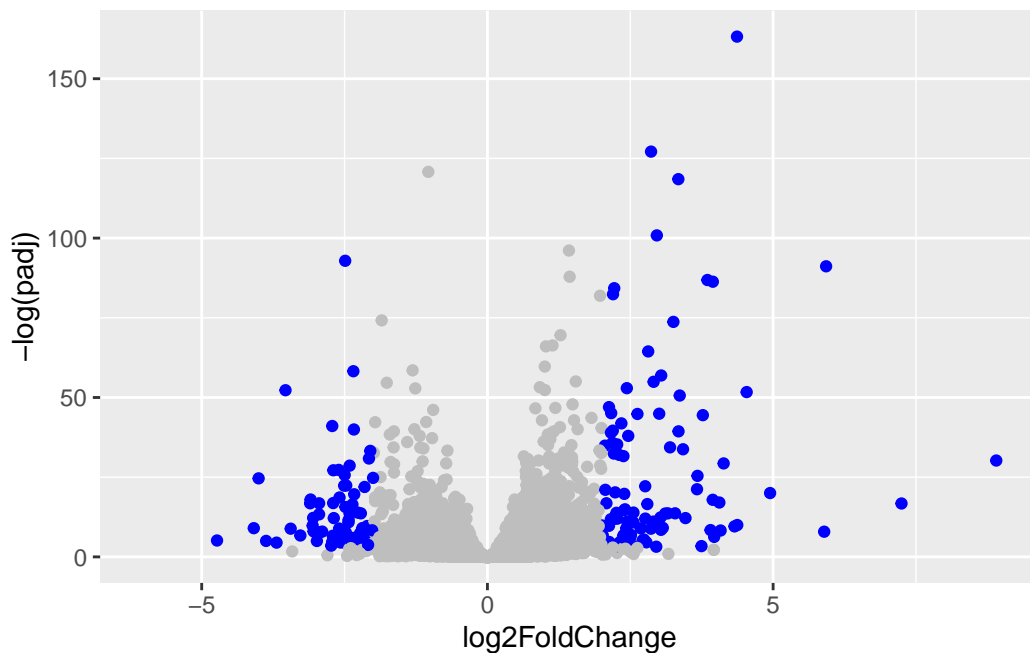
### A nicer ggplot volcano plot

```
library(ggplot2)

mycols <- rep ("gray", nrow(res))
mycols[ abs(res$log2FoldChange) > 2 ] <- "blue"
mycols[ res$padj >= 0.05 ] <- "gray"

ggplot(res) +
  aes(log2FoldChange, -log(padj)) +
  geom_point(col = mycols)
```

Warning: Removed 23549 rows containing missing values or values outside the scale range (`geom\_point()`).



## Add annotation data

We need to add gene symbols, gene names and other database ids to make my results useful for further analysis.

```
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

```

```
head(rownames(res))
```

```

[1] "ENSG000000000003" "ENSG000000000005" "ENSG000000000419" "ENSG000000000457"
[5] "ENSG000000000460" "ENSG000000000938"

```

We can use the `mapIds()` function from `bioconductor` to help us.

```

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

```

Let's see what database id formats we can translate between

```
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$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

```

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

```
head(res$symbol)
```

```

ENSG000000000003 ENSG000000000005 ENSG000000000419 ENSG000000000457 ENSG000000000460
      "TSPAN6"           "TNMD"           "DPM1"           "SCYL3"           "FIRRM"
ENSG0000000000938
      "FGR"

```

Add GENENAME then ENTREZID

```

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

```

'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	tetraspanin 6
ENSG000000000005	NA	TNMD	tenomodulin
ENSG000000000419	0.176032	DPM1	dolichyl-phosphate m..
ENSG000000000457	0.961694	SCYL3	SCY1 like pseudokina..
ENSG000000000460	0.815849	FIRRM	FIGNL1 interacting r..
ENSG000000000938	NA	FGR	FGR proto-oncogene, ..

```

res$entrezid <- 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

```

'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	entrezid	
	<numeric>	<character>	<character>	<character>	
ENSG000000000003	0.163035	TSPAN6	tetraspanin 6	7105	
ENSG000000000005	NA	TNMD	tenomodulin	64102	
ENSG000000000419	0.176032	DPM1	dolichyl-phosphate m..	8813	
ENSG000000000457	0.961694	SCYL3	SCY1 like pseudokina..	57147	
ENSG000000000460	0.815849	FIRRM	FIGNL1 interacting r..	55732	
ENSG000000000938	NA	FGR	FGR proto-oncogene, ..	2268	

## Save my annotated results

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

## Pathway analysis

We will use the **gage** function from bioconductor.

```
library(gage)
```

```
library(gageData)
```

```
library(pathview)
```

Warning: package 'pathview' was built under R version 4.4.2

```
#####  
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>).

```
#####
```

```
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 **gage** wants as input is a named vector of importance i.e. a vector with labeled fold-changes.

```
foldchanges <- res$log2FoldChange  
names(foldchanges) <- res$entrez  
head(foldchanges)
```

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

```
data(kegg.sets.hs)
keggres = gage(foldchanges, gsets=kegg.sets.hs)
```

```
head( keggres$less, 5)
```

	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
hsa00230 Purine metabolism	NA	NaN	NA	NA
hsa05340 Primary immunodeficiency	NA	NaN	NA	NA

	set.size	expl
hsa00232 Caffeine metabolism	0	NA
hsa00983 Drug metabolism - other enzymes	0	NA
hsa01100 Metabolic pathways	0	NA
hsa00230 Purine metabolism	0	NA
hsa05340 Primary immunodeficiency	0	NA

Let's look at just one of these hsa05310

```
library(pathview)
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/mariatavares/Desktop/BGGN213/class12

Info: Writing image file hsa05310.pathview.png

Insert figure for this pathway



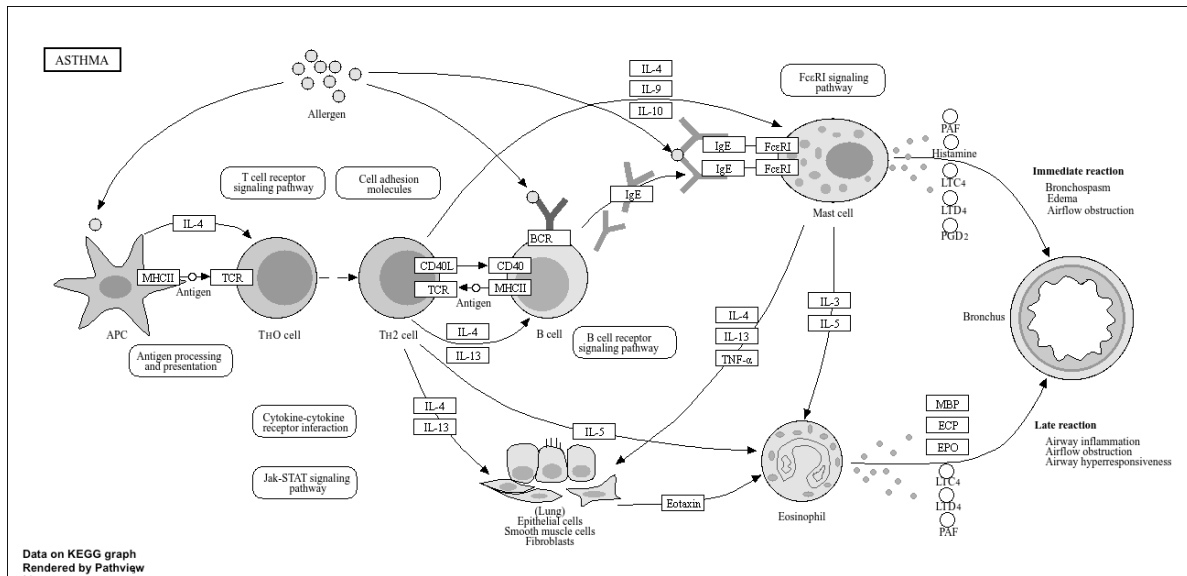


Figure 1: Asthma pathway from KEGG with my differentially expressed genes highlighted