# Lab Class 13 (DESeq2)

## Zainub Darsot (A16294217)

The data for this hands-on session comes 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).

Today we will examine this RNASeq data.

### Section 3

```
# Complete the missing code
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 |

Q1.How many genes are in this dataset?

```r
nrow(counts)
```

```
[1] 38694
```

There are 38694 genes.

> Q2. How many 'control' cell lines do we have?

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

```
[1] 4
```

There are 4 'control' cell lines

## Section 4

Start by counting the mean counts per gene in the 'control' samples, then compare this to mean counts in the 'treated' column.

Step 1: Find the counts for "control" samples Step 2: Calculate the mean coutns per gene in the "control" sample and store this in 'control.mean'.

Step 1:

```r
control.inds <- metadata$dex == "control"

metadata[control.inds, ]
```

```
          id     dex celltype    geo_id
1 SRR1039508 control   N61311 GSM1275862
3 SRR1039512 control  N052611 GSM1275866
5 SRR1039516 control  N080611 GSM1275870
7 SRR1039520 control  N061011 GSM1275874
```

```r
control.counts <- counts[,control.inds]

head(control.counts)
```

2

```
              SRR1039508 SRR1039512 SRR1039516 SRR1039520
ENSG00000000003        723         904        1170         806
ENSG00000000005          0           0           0           0
ENSG00000000419        467         616         582         417
ENSG00000000457        347         364         318         330
ENSG00000000460         96          73         118         102
ENSG00000000938          0           1           2           0
```

Step 2:

```
#apply(control.counts,1, mean)
```

OR

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

```
cont.inds <- rowMeans( counts[, metadata$dex == "control"])

head(cont.inds)
```

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

> 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)

For Treated:

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

```
          id       dex celltype     geo_id
2 SRR1039509 treated   N61311 GSM1275863
4 SRR1039513 treated  N052611 GSM1275867
6 SRR1039517 treated  N080611 GSM1275871
8 SRR1039521 treated  N061011 GSM1275875
```

```
treated.counts <- counts[,treated.inds]
```

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

```
ENSG00000000003 ENSG00000000005 ENSG00000000419 ENSG00000000457 ENSG00000000460
         658.00            0.00          546.00          316.50           78.75
ENSG00000000938
           0.00
```

To keep things tidy, we will store `control.mean` and `treated.mean` together as two columns in a data frame

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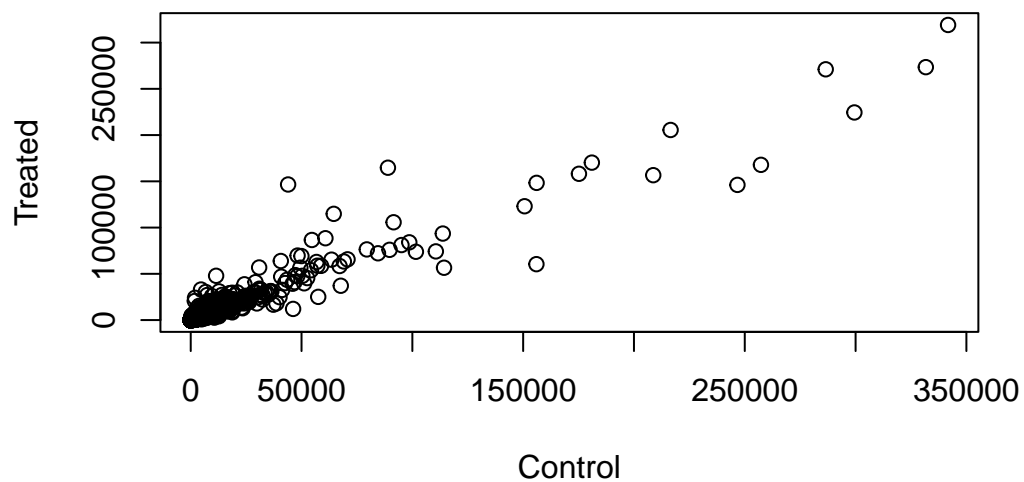
```
head(meancounts)
```

```
                control.mean treated.mean
ENSG00000000003       900.75       658.00
ENSG00000000005         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(meancounts)
```

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

Plot: > 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.
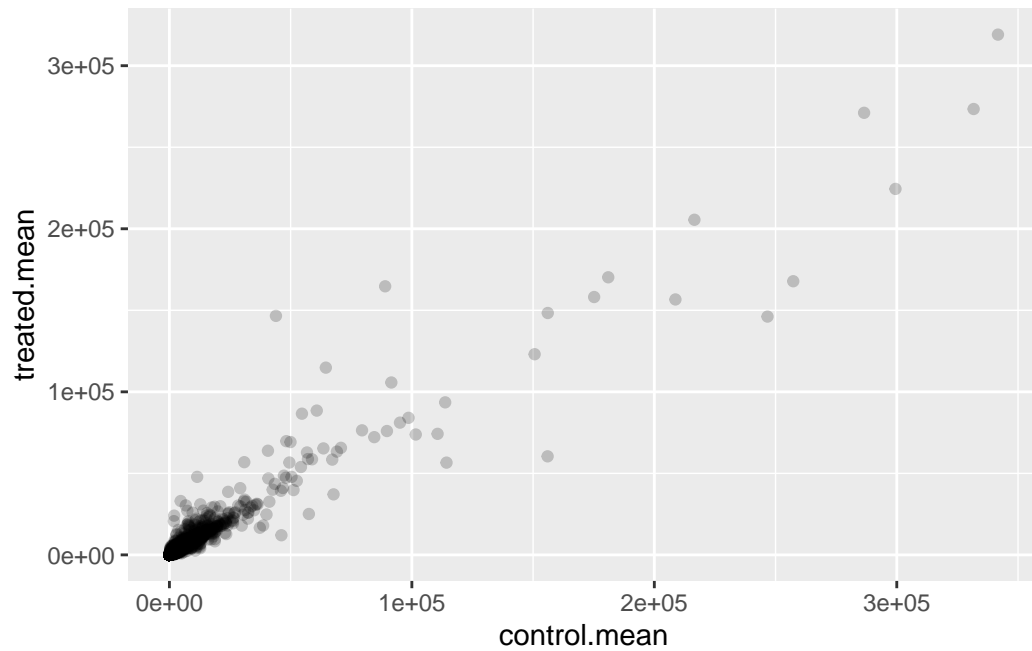
4

```
plot(meancounts[,1],meancounts[,2], 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?

```
library(ggplot2)

ggplot(meancounts) +
  aes(control.mean, treated.mean) +
  geom_point(alpha = 0.2)
```
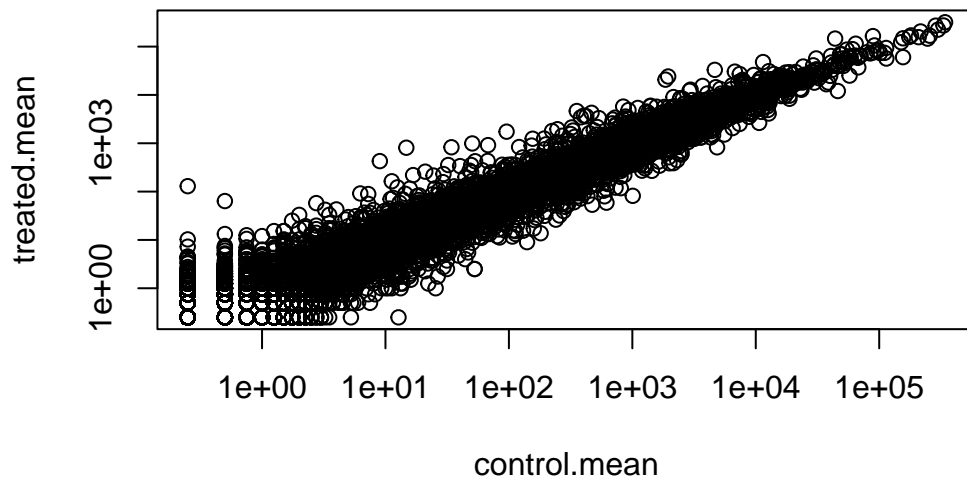
Q.6 Try plotting both axes on a log scale. What is the argument to plot() that allows you to do this?

```
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
```

We often use Log Transformations for when the data is skewed and measured over a large range. Base10 and natural logs are all valid, but Log2 units is preferred because they are much easier to understand

Add a Log2 Fold-change column to `meancounts` data.frame:

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


head(meancounts)
```

```
                control.mean treated.mean        log2fc
ENSG00000000003       900.75       658.00  -0.45303916
ENSG00000000005         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
```

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

7

The ! flips TRUE values to False

```
x <- c(T, F, T)
!x
```

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

```
dim(mycounts)
```

```
[1] 21817     3
```

```
head(mycounts)
```

```
                control.mean treated.mean      log2fc
ENSG00000000003       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?

The arr.ind=TRUE functions returns the row and column for `TRUE` values. The unique() function makes sure that any row is not counted twice when it has zero entries in both samples.

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

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

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

```
[1] 250
```

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? 367

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

[1] 367

There are 367 downregulated genes.

We forgot about statistical significance of these differences...

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

The main limitation is that statistical significance is unknown. We will use DESeq2 package to do this analysis properly.

## Section 5: Setting up for DESeq

We must load DESeq2 with `library()` function.

```
library(DESeq2)
```

Loading required package: S4Vectors

Loading required package: stats4

Loading required package: BiocGenerics

Attaching package: '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


Attaching package: 'S4Vectors'

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

    findMatches

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

    expand.grid, I, unname

Loading required package: IRanges


Attaching package: 'IRanges'

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

    windows

Loading required package: GenomicRanges

Loading required package: GenomeInfoDb

Loading required package: SummarizedExperiment

Loading required package: MatrixGenerics

Loading required package: matrixStats

```
Warning: package 'matrixStats' was built under R version 4.3.2


Attaching package: 'MatrixGenerics'

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

    colAlls, colAnyNAs, colAnys, 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, rowAnys, 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

Loading required package: Biobase

Welcome to Bioconductor

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


Attaching package: 'Biobase'

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

    rowMedians

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

    anyMissing, rowMedians
```

Setting up DESeq:

```r
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

Now we can run our DESeq analysis

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

estimating size factors

estimating dispersions

gene-wise dispersion estimates

mean-dispersion relationship

final dispersion estimates

fitting model and testing

Getting results back from **dds** object:

```r
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>
ENSG00000000003 747.194195     -0.3507030  0.168246 -2.084470 0.0371175
```
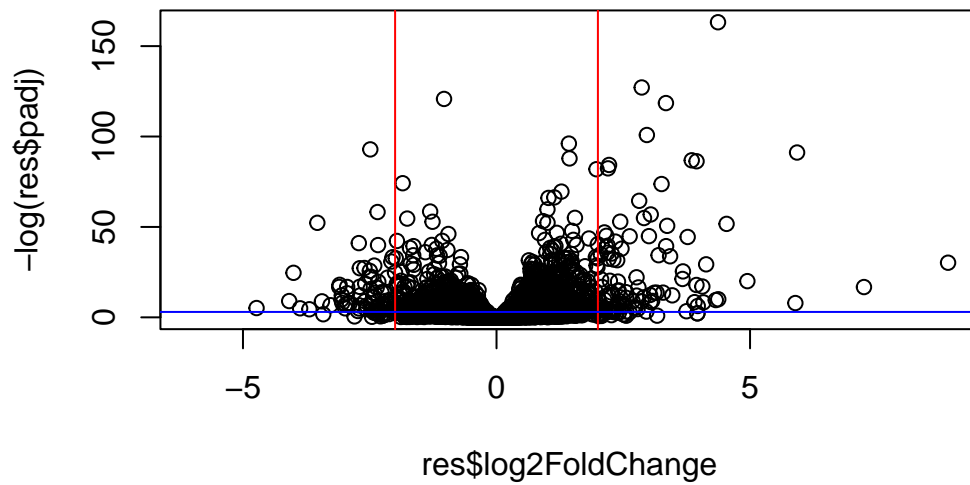
```
ENSG00000000005    0.000000                NA        NA        NA        NA
ENSG00000000419 520.134160         0.2061078  0.101059  2.039475 0.0414026
ENSG00000000457 322.664844         0.0245269  0.145145  0.168982 0.8658106
ENSG00000000460  87.682625        -0.1471420  0.257007 -0.572521 0.5669691
ENSG00000000938   0.319167        -1.7322890  3.493601 -0.495846 0.6200029
                     padj
                <numeric>
ENSG00000000003  0.163035
ENSG00000000005        NA
ENSG00000000419  0.176032
ENSG00000000457  0.961694
ENSG00000000460  0.815849
ENSG00000000938        NA
```

## A summary results plot:

Volcano plot. This is a common type of summary figure that keeps both our inner biologist and inner statistician happy because it shows both P-values and Log2(Fold-Changes).

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

Save our result to date.

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

## Section 8: Adding annotation Data:

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

Warning: package 'AnnotationDbi' was built under R version 4.3.2

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

Avalaible key types:

```
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"
```

The main function we will use here is called `mapIds()`

```
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>
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
                <numeric>
ENSG00000000003  0.163035
ENSG00000000005        NA
ENSG00000000419  0.176032
ENSG00000000457  0.961694
ENSG00000000460  0.815849
ENSG00000000938        NA
```

```
res$symbol <- mapIds(org.Hs.eg.db,
                    keys = row.names(res),
                    keytype = "ENSEMBL",
                    column = "SYMBOL",
                    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>
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
                <numeric> <character>
ENSG00000000003  0.163035      TSPAN6
ENSG00000000005        NA        TNMD
ENSG00000000419  0.176032        DPM1
ENSG00000000457  0.961694       SCYL3
ENSG00000000460  0.815849       FIRRM
ENSG00000000938        NA         FGR
```

Adding 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 8 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              genename
                 <numeric> <character>           <character>
ENSG00000000003   0.163035      TSPAN6          tetraspanin 6
ENSG00000000005         NA        TNMD            tenomodulin
ENSG00000000419   0.176032        DPM1 dolichyl-phosphate m..
ENSG00000000457   0.961694       SCYL3 SCY1 like pseudokina..
ENSG00000000460   0.815849       FIRRM FIGNL1 interacting r..
ENSG00000000938         NA         FGR FGR proto-oncogene, ..
```

```r
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
```

## Pathway Analysis

We will use **gage** package along with **pathview** here to do geneset enrichment (a.k.a. pathway analysis) and figure generation respectively

```r
#1 message
library(pathview)
```

```
#############################################################################
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).
#############################################################################
```

```r
library(gageData)
library(gage)
```

```r
data(kegg.sets.hs)

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 we need for `ggplot()` is our genes in ENTREZ id format with a measure of their importance.

It wants a vector of e.g. fold-changes

```r
foldchanges <- res$log2FoldChange
head(foldchanges)
```

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

```r
x <- c(100, 80, 100)
names(x) <- c("destiny", "barry", "chris")
x
```

```
destiny   barry   chris
    100      80     100
```

Add ENTREZ ids as `names()` to my `foldchanges` vector

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

```
      7105       64102       8813      57147      55732       2268
-0.35070302          NA  0.20610777  0.02452695 -0.14714205 -1.73228897
```

Now we can run `gage` with this input vector and the geneset we want to examine for overlap/enrichment...

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

```
attributes(keggres)
```

```
$names
[1] "greater" "less"    "stats"
```

```
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
```

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

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

```
Info: Working in directory C:/Users/zidar/OneDrive/Desktop/BIMM 143/class 13
```
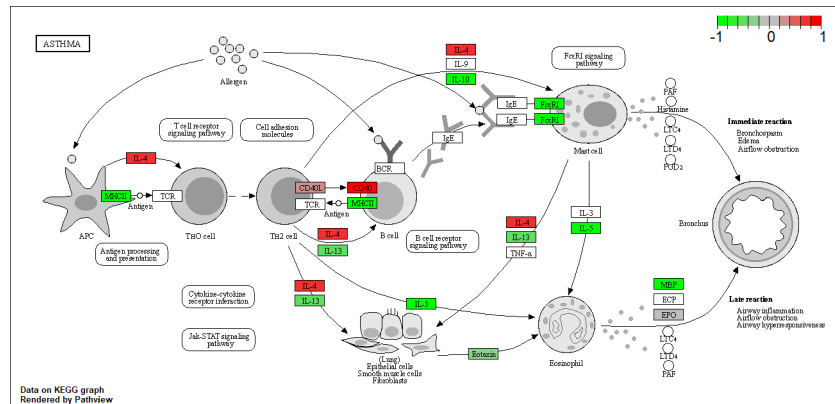
```
Info: Writing image file hsa05310.pathview.png
```

Figure 1: My genes involved in Asthma pathway