# Class 13

Abzael Verduzco

```r
library(BiocManager)
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

Loading required package: GenomicRanges

Loading required package: GenomeInfoDb

Loading required package: SummarizedExperiment

Loading required package: MatrixGenerics

Loading required package: matrixStats


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
```

This week we are looking at differential expression analysis.

```r
# Complete the missing code
counts <- read.csv("airway_scaledcounts.csv", stringsAsFactors = FALSE, row.names=1)
metadata <-  read.csv("airway_metadata.csv", stringsAsFactors = FALSE)

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
```

```
View(metadata)
```

```
nrow(counts)
```

```
[1] 38694
```

Q1. How many genes are in this dataset?

There are 38694 genes in this data set.

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

There are 4 control cell lines.

#4. Toy differential gene expression

**Extract and summarize the control 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
ENSG00000000938
           0.75
```

**Extract and summarize the treated (i.e. drug) samples**

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)
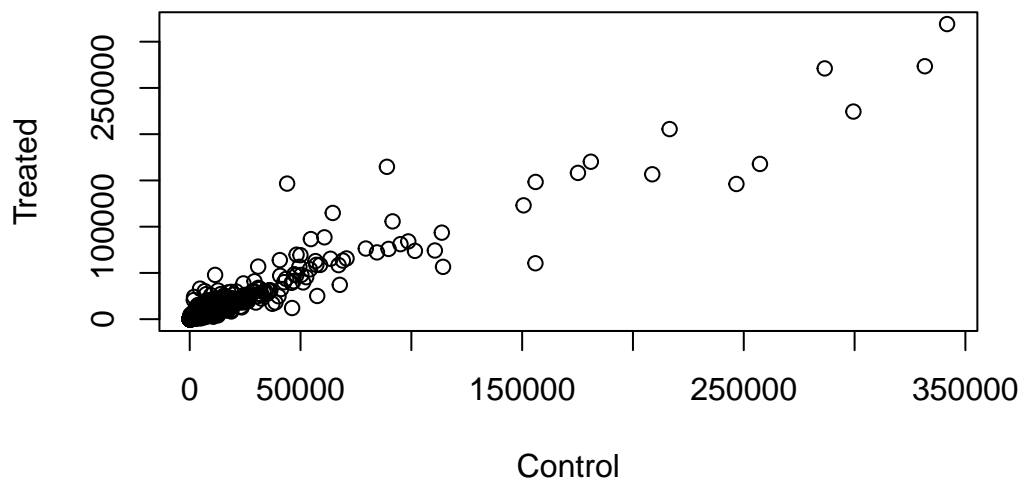
```
treated <- metadata[metadata[,"dex"]=="treated",]
treated.counts <- counts[,treated$id]
treated.mean <- rowMeans(treated.counts)
```

Store these results together in a new data frame called `meancounts`

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

Lets make a plot to explore the results a little

```
plot(meancounts[,1], meancounts[,2], xlab="Control", ylab="Treated")
```
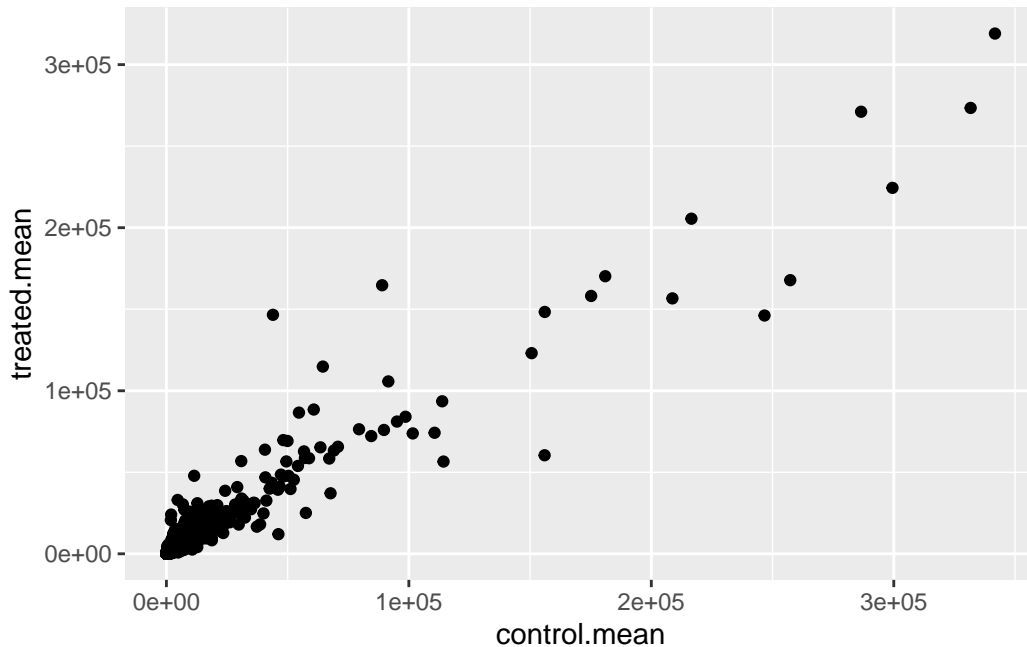


Q5 (b).You could also use the ggplot2 package to make this figure producing the plot below. What geom_?() function would you use for this plot?

geom_point

```
library(ggplot2)

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



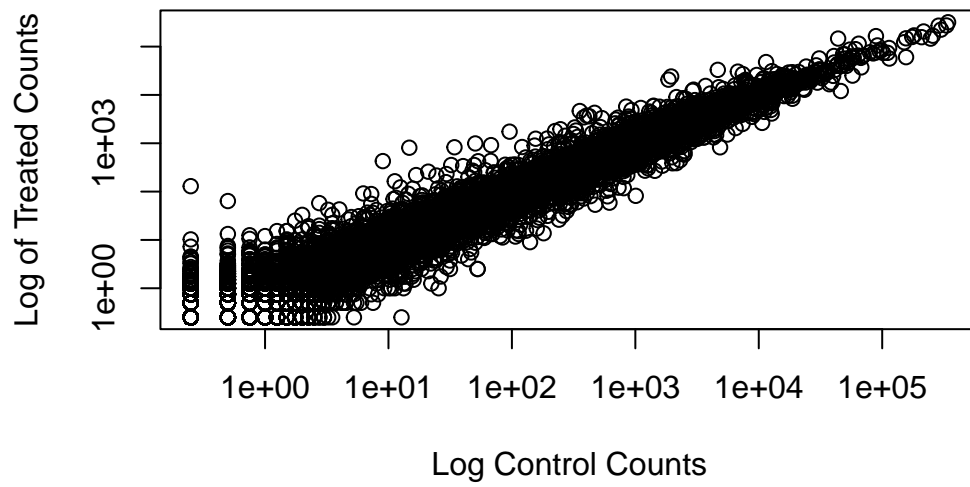We will make a log-log plot to draw out this skewed data and see what is going on.

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

log

```
plot(meancounts[,1], meancounts[,2], log="xy",
     xlab="Log Control Counts",
     ylab="Log of Treated Counts")
```

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

6

We often log2 transformations when dealing with this sort of data.

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

```
[1] 0
```

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

```
[1] 1
```

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

```
[1] -1
```

```
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
```

```r
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
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?

It will return both the rows and colums where there are TRUE values. Unique() will make sure they are not present twice.

```r
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?

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

```
[1] 250
```

There are 250 up regulated 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?

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

```
[1] 367
```

There are 367 down regulated genes.

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

No, as we haven't made sure they are significant we will need to use DESeq2 to make sure of this.

## Setting up for DESeq

```
library(DESeq2)
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},
  }
```

```r
dds <- DESeqDataSetFromMatrix(countData=counts,
                              colData=metadata,
                              design=~dex)
```
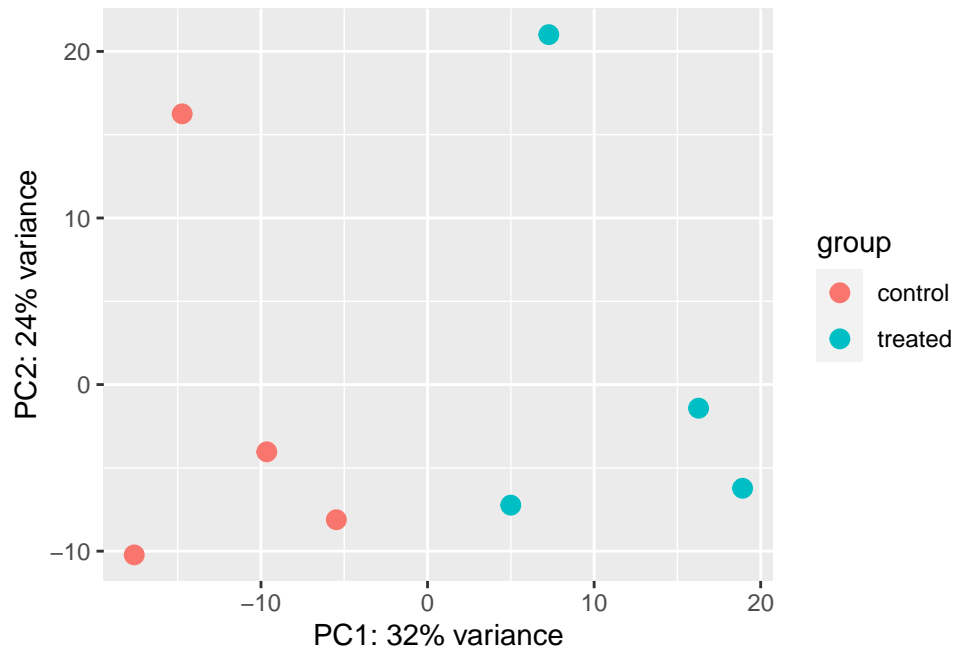
converting counts to integer mode

Warning in DESeqDataSet(se, design = design, ignoreRank): some variables in
design formula are characters, converting to factors

```r
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
```

```r
vsd <- vst(dds, blind = FALSE)
plotPCA(vsd, intgroup = c("dex"))
```

using ntop=500 top features by variance

```
pcaData <- plotPCA(vsd, intgroup=c("dex"), returnData=TRUE)
```

using ntop=500 top features by variance

```
head(pcaData)
```

```
                  PC1        PC2    group      dex       name
SRR1039508 -17.607922 -10.225252 control control SRR1039508
SRR1039509    4.996738  -7.238117 treated treated SRR1039509
SRR1039512   -5.474456  -8.113993 control control SRR1039512
SRR1039513   18.912974  -6.226041 treated treated SRR1039513
SRR1039516  -14.729173  16.252000 control control SRR1039516
SRR1039517    7.279863  21.008034 treated treated SRR1039517
```

```
percentVar <- round(100 * attr(pcaData, "percentVar"))
```

```
ggplot(pcaData) +
  aes(x = PC1, y = PC2, color = dex) +
  geom_point(size =3) +
```

```
    xlab(paste0("PC1: ", percentVar[1], "% variance")) +
    ylab(paste0("PC2: ", percentVar[2], "% variance")) +
    coord_fixed() +
    theme_bw()
```



## DESeq Analysis

```
#results(dds)
```

```
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>
ENSG00000000003  747.1942     -0.3507030  0.168246 -2.084470 0.0371175
ENSG00000000005    0.0000             NA        NA        NA         NA
ENSG00000000419  520.1342      0.2061078  0.101059  2.039475 0.0414026
ENSG00000000457  322.6648      0.0245269  0.145145  0.168982 0.8658106
ENSG00000000460   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>
ENSG00000000003   0.163035
ENSG00000000005         NA
ENSG00000000419   0.176032
ENSG00000000457   0.961694
ENSG00000000460   0.815849
...                    ...
ENSG00000283115         NA
ENSG00000283116         NA
ENSG00000283119         NA
ENSG00000283120         NA
ENSG00000283123         NA
```

For statistical significance

```r
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**

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

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

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

```r
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()`

Our current IDs are here:

```
#mapIds()
head(row.names(res))
```

[1] "ENSG00000000003" "ENSG00000000005" "ENSG00000000419" "ENSG00000000457"
[5] "ENSG00000000460" "ENSG00000000938"

These are in ENSEMBLE format. I want "SYMBOL" ids:

```
res$symbol <-mapIds(org.Hs.eg.db,
                    keys = row.names(res),
                    keytype="ENSEMBL",        #Our genenames
                    column = "SYMBOL",        #Format our our new genenames
                    multiVals = "first")      #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 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 |

Let's add GENENAME

```

```r
res$genename <- mapIds(org.Hs.eg.db,
                       keys = row.names(res),
                       keytype= "ENSEMBL",       #Our genenames
                       column = "GENENAME",    #Format our our new genenames
                       multiVals = "first")    #New format we want to add
```

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

```r
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", #Our genenames
                     column = "ENTREZID",    #Format our our new genenames
                     multiVals = "first")    #New format we want to add
```

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

```
res$uniprot <- mapIds(org.Hs.eg.db,
                      keys = row.names(res),
                      keytype =  "ENSEMBL",        #Our genenames
                      column = "UNIPROT",      #Format our our new genenames
                      multiVals = "first")    #New format we want to add
```

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

Part 10: Pathway Analysis

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

```
library(gage)
```

```
library(gageData)
```

Let's have a peak at the first two pathways in KEGG

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

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

```
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
```

Now we can run `gage()` with this input vector and the genset we want to examine for overlap/enrichment…

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

Look at the results

```
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
```

We can view these pathways with our geneset genes highlighted using the `pathview()` function. E.g. for "Asthma" I will use the pathway.id hsa05310 as seen above.

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

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

```
Info: Working in directory /Users/abzael/Desktop/BIMM 143/Class 13 - BIMM 143
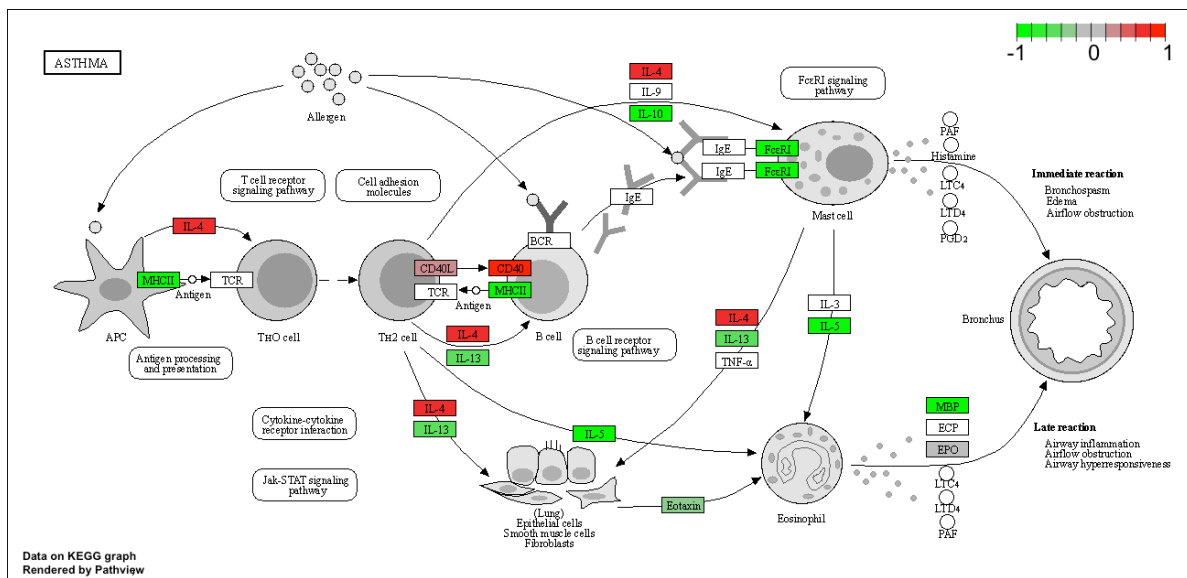```

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



Figure 1: My genes involved in Asthma pathway