# class11 : DESeq analysis

Soobin (PID:A15201229)

2/23/2022

This week we are looking at differential expression analysis.

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).

## Import/Read the data from Himes et al.

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

Lets have a peak at this data

```
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
```

Sanity check on correspondence of counts and metadata

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

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

```
dim(counts)
```

```
## [1] 38694     8
```

> Q1. How many genes are in this dataset?

There are 38694 genes in this dataset.

```
n.control <- sum( metadata$dex == "control" )
```

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

There are 4 'control' cell lines.

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

```
##
## FALSE   TRUE
##     4      4
```

**Extract and summarize the control samples**

To find out where the control samples are we need the metadata

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

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

Use rowMeans() instead of rowSums()/4.

**Extract and summarize the treated 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$dex == "treated", ]
treated.count <- counts[ , treated$id ]
treated.mean <- rowMeans(treated.count)
head(treated.mean)
```

```
## ENSG00000000003 ENSG00000000005 ENSG00000000419 ENSG00000000457 ENSG00000000460
##          658.00            0.00          546.00          316.50           78.75
## ENSG00000000938
##            0.00
```
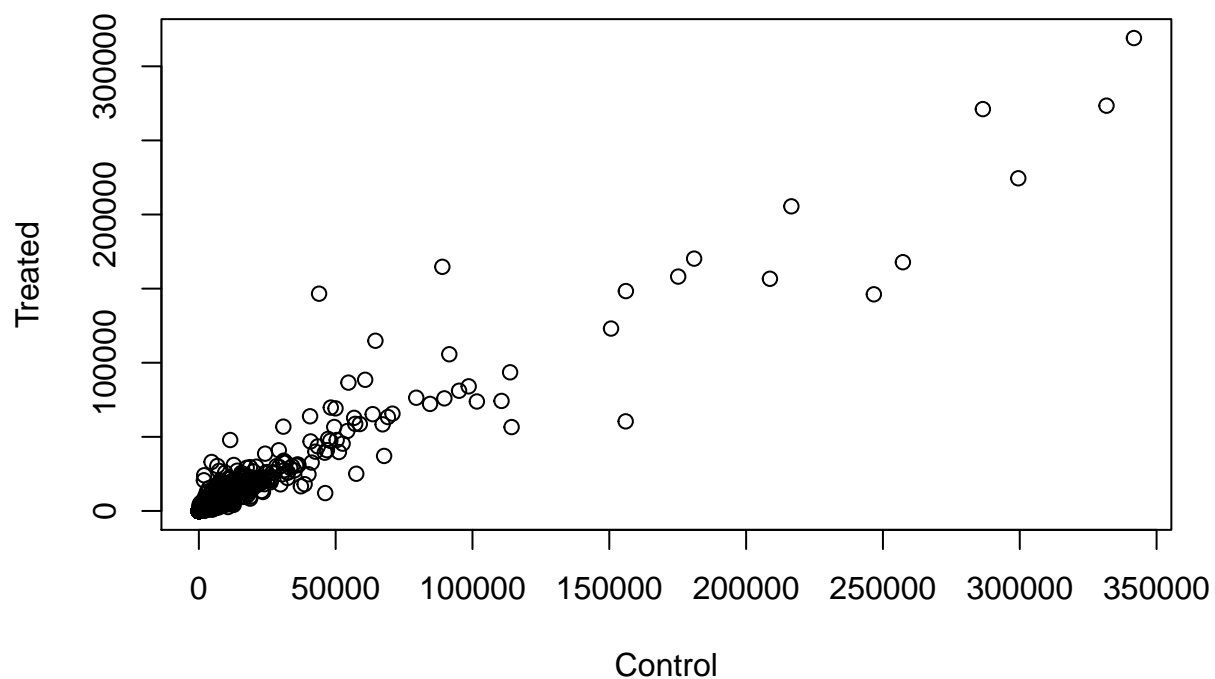
Store these results together in a new data frame called `meancounts`.

```
meancounts <- data.frame(control.mean, treated.mean)
```
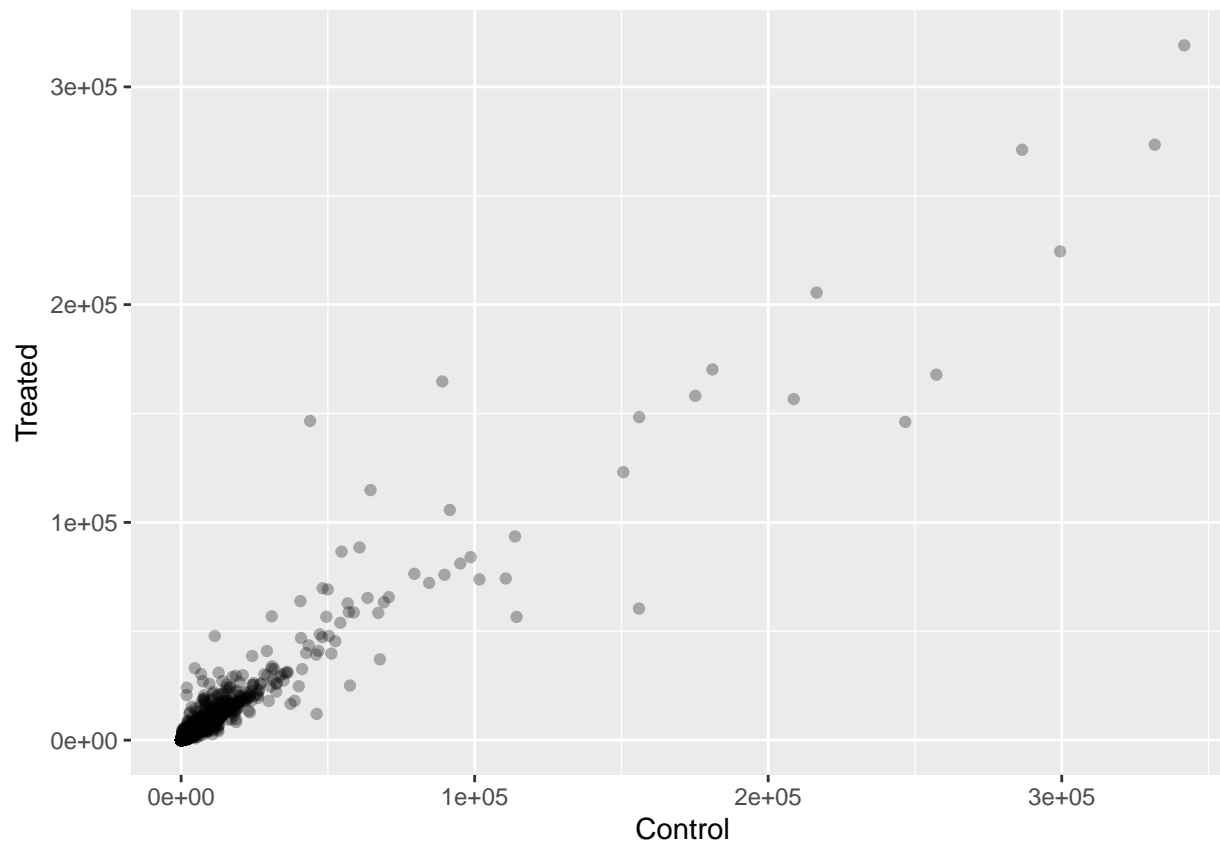
Lets make a plot to explore the results

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[,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(x=control.mean, y=treated.mean) + geom_point(alpha=0.3) + labs(x="Control", y=
```
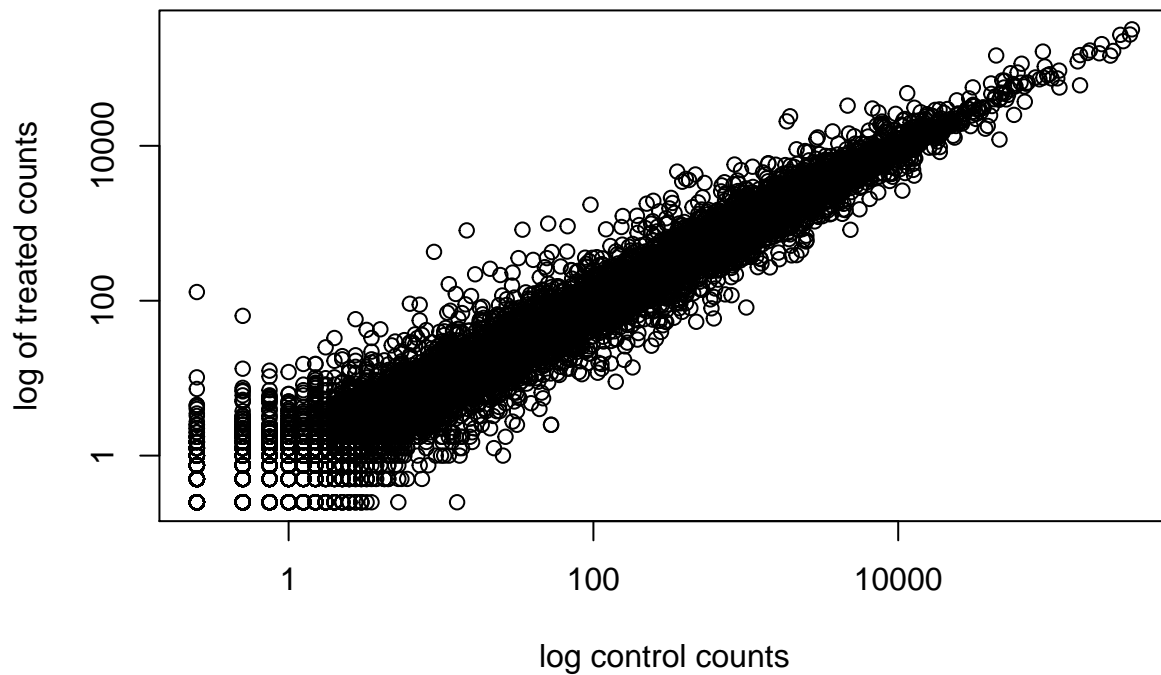
We will make a log-log plot to draw out this skewed data and see what is going on.

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

```
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
```

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
```

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

```
## [1] 2
```

This log2 transformation has this nice property where if there is no change the log2 value will be zero and if it double the log2 value will be 1 and if halved it will be -1.

If the drug had no effect, the log of treated.mean vs. control.mean would just have a straight line. However, the log2 fold change have some up and down from 0, which indicates the possibility that the drug have an effect.

So lets add a log2 fold change column to our results.

```
# add a column
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
```

We need to get rid of zero count genes that we can not say anything about.

```
zero.values <- which( meancounts[ , 1:2] == 0, arr.ind=TRUE )
to.rm <- unique( zero.values[ , 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?

which() tells us which elements are true in the vector. The arr.ind=TRUE argument will tell which() to return both row and columns that have a TRUE value. unique() will prevent us from counting the zero twice from both column 1 and 2.

How many genes are remaining?

```
nrow(mycounts)
```

```
## [1] 21817
```

> 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
sum(up.ind)
```

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

There are 250 up-regulated genes that have greater fold change than 2.

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

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

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

There are 367 down-regulated genes with fold changes smaller than -2.

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

No, not yet. We do not know whether the change is statistically significant.

# DESeq2 Analysis

Let's do this the right way. DEseq2 is an R package specifically for analyzing count-based NGS data like RNA-seq.

```
# load up DESeq2
library(DESeq2)
```

```
## Loading required package: S4Vectors
```

```
## Loading required package: stats4
```

```
## Loading required package: BiocGenerics
```

```
## Loading required package: parallel
```

```
##
## Attaching package: 'BiocGenerics'
```

```
## The following objects are masked from 'package:parallel':
##
##     clusterApply, clusterApplyLB, clusterCall, clusterEvalQ,
##     clusterExport, clusterMap, parApply, parCapply, parLapply,
##     parLapplyLB, parRapply, parSapply, parSapplyLB
```

```
## The following objects are masked from 'package:stats':
##
##     IQR, mad, sd, var, xtabs
```

```
## The following objects are masked from 'package:base':
##
##      anyDuplicated, 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:base':
##
##      expand.grid


## 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
```

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

```
## estimating size factors

## estimating dispersions

## gene-wise dispersion estimates

## mean-dispersion relationship

## final dispersion estimates

## fitting model and testing
```

```r
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
```

We can get some basic summary tallies using the `summary()` function.
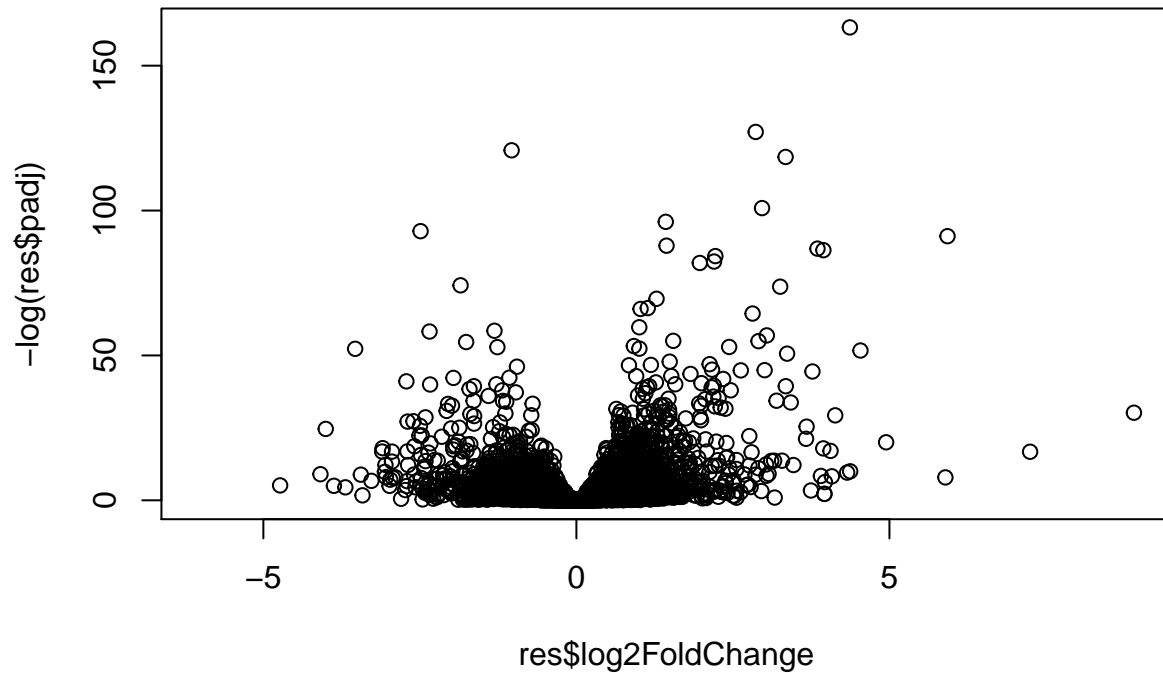
```r
summary(res, alpha=0.05)
```

```
##
## out of 25258 with nonzero total read count
## adjusted p-value < 0.05
## LFC > 0 (up)       : 1242, 4.9%
## LFC < 0 (down)     : 939, 3.7%
## 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
```

# Volcano plot

Make a summary plot of our results.

```r
plot( res$log2FoldChange, -log(res$padj) )
```



```r
log(0.1)
```

```
## [1] -2.302585
```

```r
log(0.005)
```

```
## [1] -5.298317
```

Finish for today by saving our results

```r
write.csv(res, file="DESeq2_results.csv")
```

---

## Adding annotation data

To help interpret our results, we need to understand what the differentially expressed genes are. A first step here is to get the gene names (i.e. gene SYMBOLS).

For this I will install: - BiocManager::install("AnnotationDbi") - BiocManager::install("org.Hs.eg.db")

```
# BiocManager main annotation packages
library("AnnotationDbi")
library("org.Hs.eg.db")
```

##

What DB identifiers can I look up?

```
# There should be some similar data bases
columns(org.Hs.eg.db)
```

```
##  [1] "ACCNUM"       "ALIAS"        "ENSEMBL"      "ENSEMBLPROT"  "ENSEMBLTRANS"
##  [6] "ENTREZID"     "ENZYME"       "EVIDENCE"     "EVIDENCEALL"  "GENENAME"
## [11] "GO"           "GOALL"        "IPI"          "MAP"          "OMIM"
## [16] "ONTOLOGY"     "ONTOLOGYALL"  "PATH"         "PFAM"         "PMID"
## [21] "PROSITE"      "REFSEQ"       "SYMBOL"       "UCSCKG"       "UNIGENE"
## [26] "UNIPROT"
```

We will use `mapIds()` function to translate between different ids.

```
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>
## 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    C1orf112
## ENSG00000000938        NA         FGR
```

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

12

```r
# entrez : NCBI database
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

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

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

```r
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

```r
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  A0A024RCI0
## ENSG00000000005        NA        TNMD       64102      Q9H2S6
## ENSG00000000419  0.176032        DPM1        8813      O60762
## ENSG00000000457  0.961694       SCYL3       57147      Q8IZE3
## ENSG00000000460  0.815849    C1orf112       55732  A0A024R922
## 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, ..
```

# Pathway analysis with R and Bioconductor

Here we play with just one, the GAGE package (which stands for Generally Applicable Gene set Enrichment), to do KEGG pathway enrichment analysis on our RNA-seq based differential expression results.

I need to install the gage package along with the pathview package for generating pathway figures from my results.

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

```
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)

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"
```

We need a vector of fold-change labeled with the names of our genes in ENTREZ format.

```
# Make a separate vector for res$log2FoldChange
foldchanges = res$log2FoldChange

# Assign "literally" names to this vector that we can identify the foldchanges with entrez format
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 the GAGE analysis passing in our foldchanges vector and the KEGG genesets we are interested in.

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

Let's have a look at what is contained in this `keggres` results object (i.e. it's attributes).

```
attributes(keggres)
```

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

```
# Look at the first three down (less) pathways
head(keggres$less, 3)
```
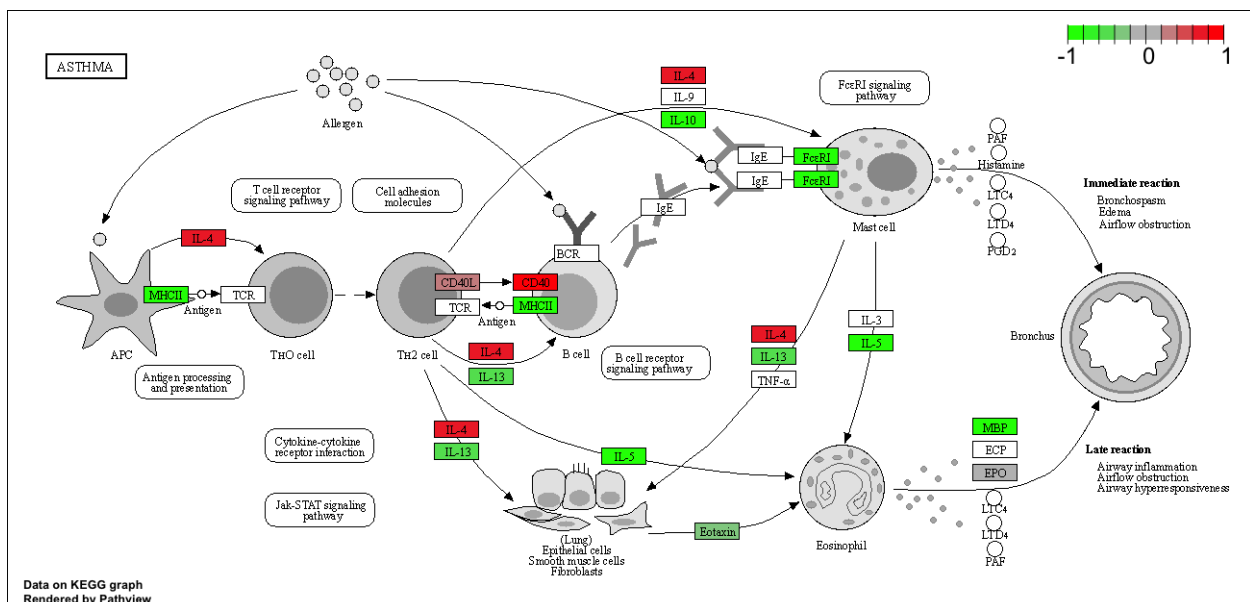
```
##                                    p.geomean stat.mean       p.val
## hsa05332 Graft-versus-host disease 0.0004250461 -3.473346 0.0004250461
## hsa04940 Type I diabetes mellitus  0.0017820293 -3.002352 0.0017820293
## hsa05310 Asthma                    0.0020045888 -3.009050 0.0020045888
##                                       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
```

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

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

```
## Info: Working in directory /Users/sbhwang/Desktop/BIMM 143/class11
```

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



15

Q12. Can you do the same procedure as above to plot the pathview figures for the top 2 down-reguled pathways?

```
pathview(gene.data=foldchanges, pathway.id="hsa05332")
```

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

```
## Info: Working in directory /Users/sbhwang/Desktop/BIMM 143/class11
```
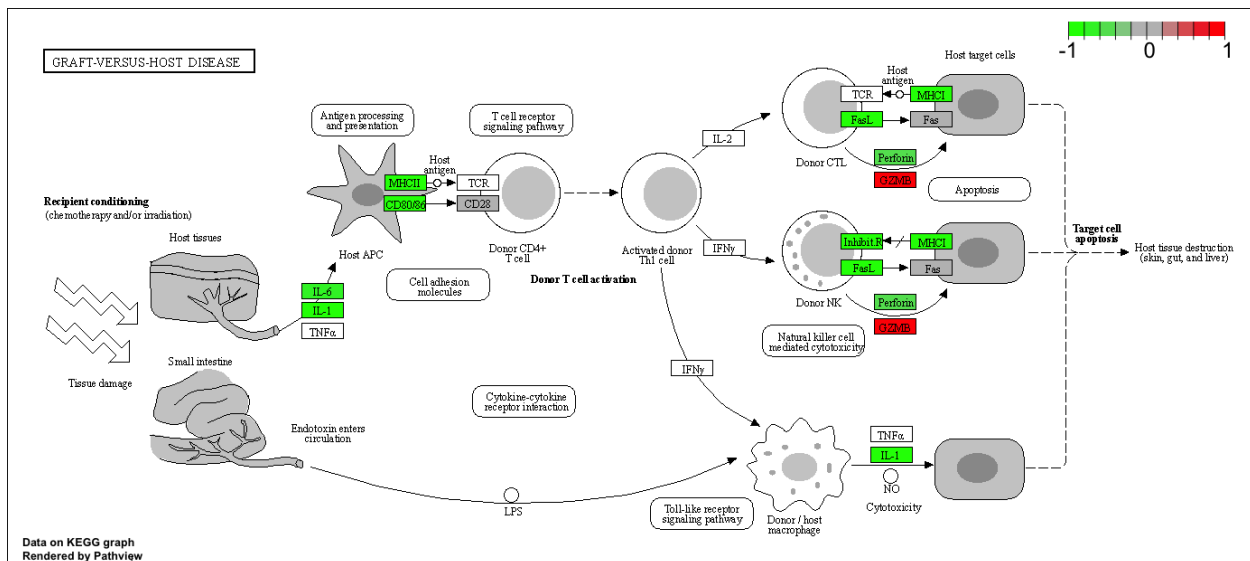
```
## Info: Writing image file hsa05332.pathview.png
```
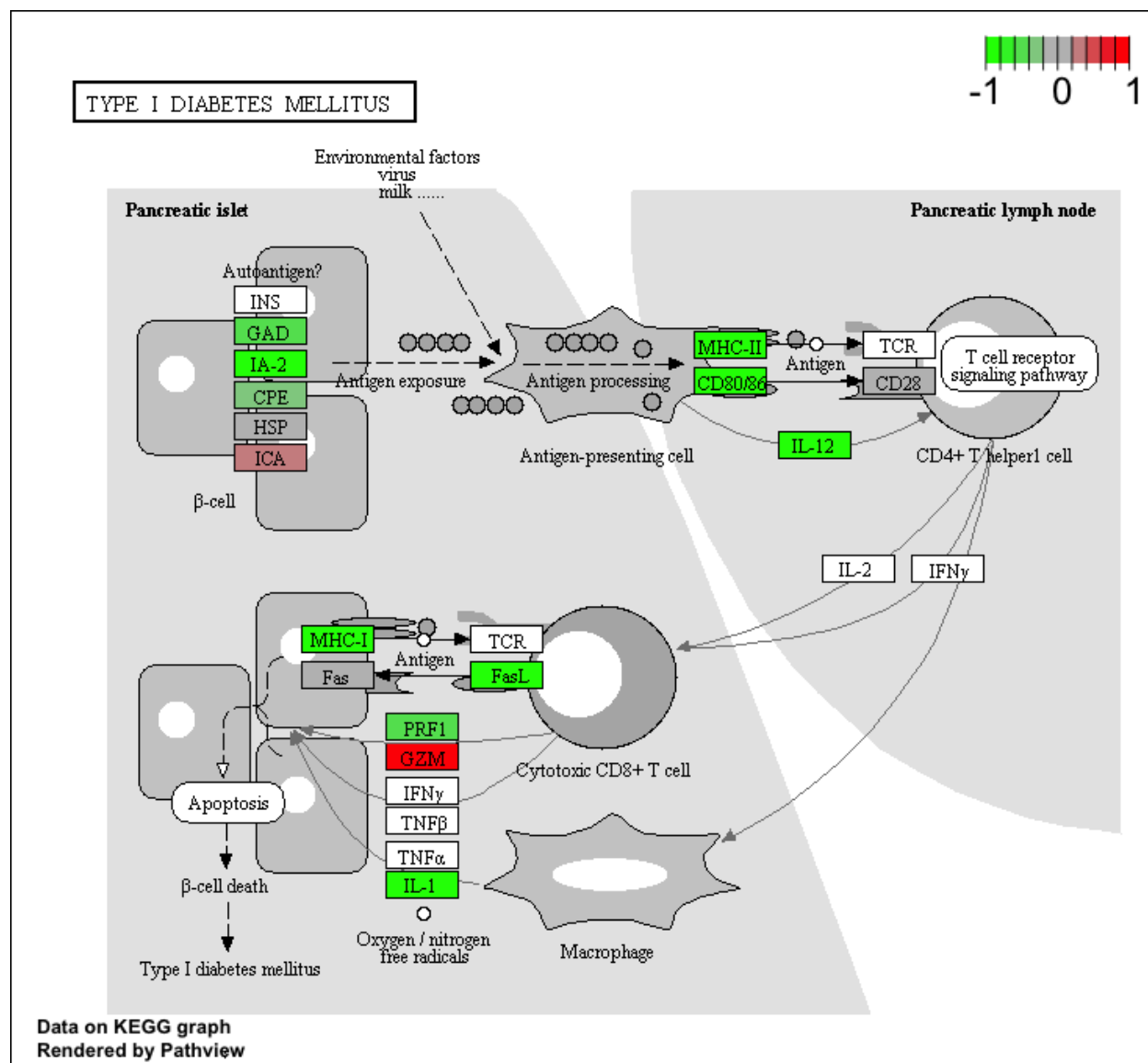
```
pathview(gene.data=foldchanges, pathway.id="hsa04940")
```

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

```
## Info: Working in directory /Users/sbhwang/Desktop/BIMM 143/class11
```

```
## Info: Writing image file hsa04940.pathview.png
```

## Final step save our results

```
write.csv(res, file="deseq_results2.csv")
```