

Class15

Jaquish (PID: A59010386)

11/17/2021

Background

Today we examine 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)

Load the `countData` and `colData`

We need 2 things - 1: count Data - 2: colData (the metadata that tells us about the design of the experiment)

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

```
#head(counts)
```

```
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
```

Side note: Let's check the correspondence of the metadata and count data setup.

```
metadata$id
```

```
## [1] "SRR1039508" "SRR1039509" "SRR1039512" "SRR1039513" "SRR1039516"
## [6] "SRR1039517" "SRR1039520" "SRR1039521"
```

```
colnames(counts)
```

```
## [1] "SRR1039508" "SRR1039509" "SRR1039512" "SRR1039513" "SRR1039516"
## [6] "SRR1039517" "SRR1039520" "SRR1039521"
```

We can use `'=='` thing to see if they are the same

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

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

Checks to see if everything is what you thought, in this case confirming there are no falses.

```
all( c(T,T,T,F))
```

```
## [1] FALSE
```

Compare control to treated

First we need to access all the controls columns in our counts data.

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

Use these ids to access just the control columns of our 'counts' data

```
head(counts[ , control.ids])
```

```
##                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
```

```
control.mean <- rowMeans(counts[ , control.ids])  
head(control.mean)
```

```
## ENSG000000000003 ENSG000000000005 ENSG000000000419 ENSG000000000457 ENSG000000000460  
##           900.75           0.00           520.50           339.75           97.25  
## ENSG000000000938  
##           0.75
```

Do the same for the treated.

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

```
head(counts[ , treated.ids])
```

```
##                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(counts[ , treated.ids])
head(treated.mean)
```

```
## ENSG000000000003 ENSG000000000005 ENSG000000000419 ENSG000000000457 ENSG000000000460
##           658.00           0.00           546.00           316.50           78.75
## ENSG000000000938
##           0.00
```

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

There are 38694 rows/genes in this dataset.

```
nrow(counts)
```

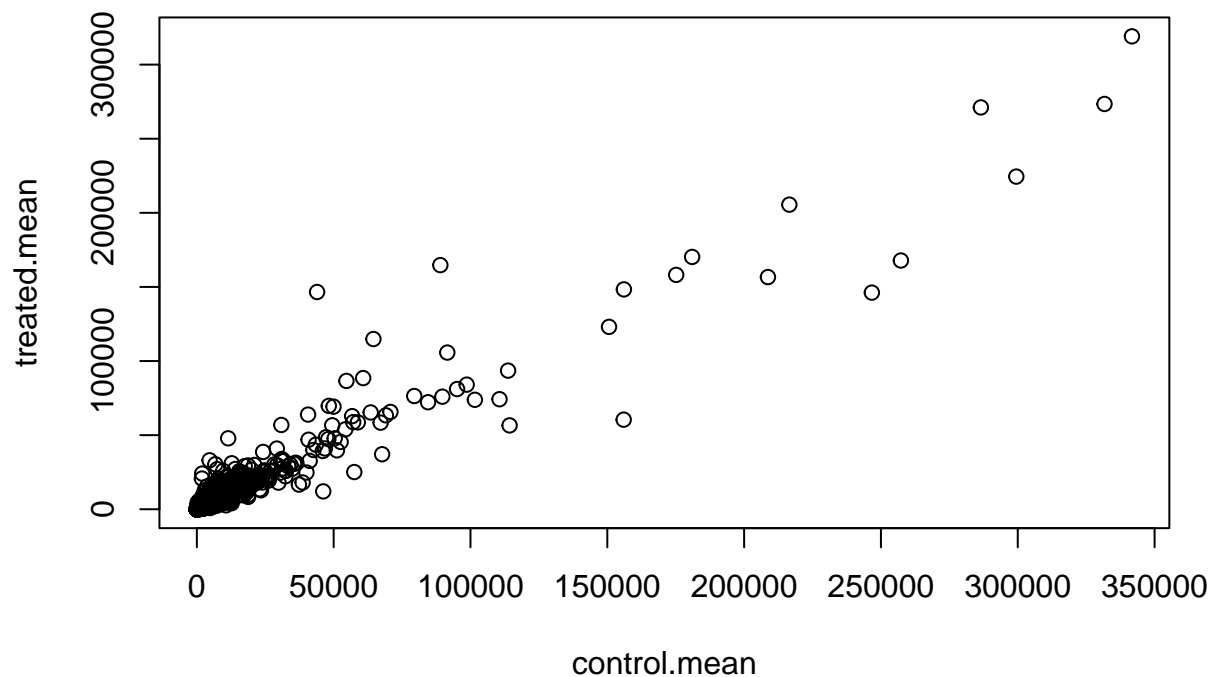
```
## [1] 38694
```

Q1. 38,694 genes

Compare the control v. treated

A quick plot of our progress so far.

```
plot(meancounts)
```



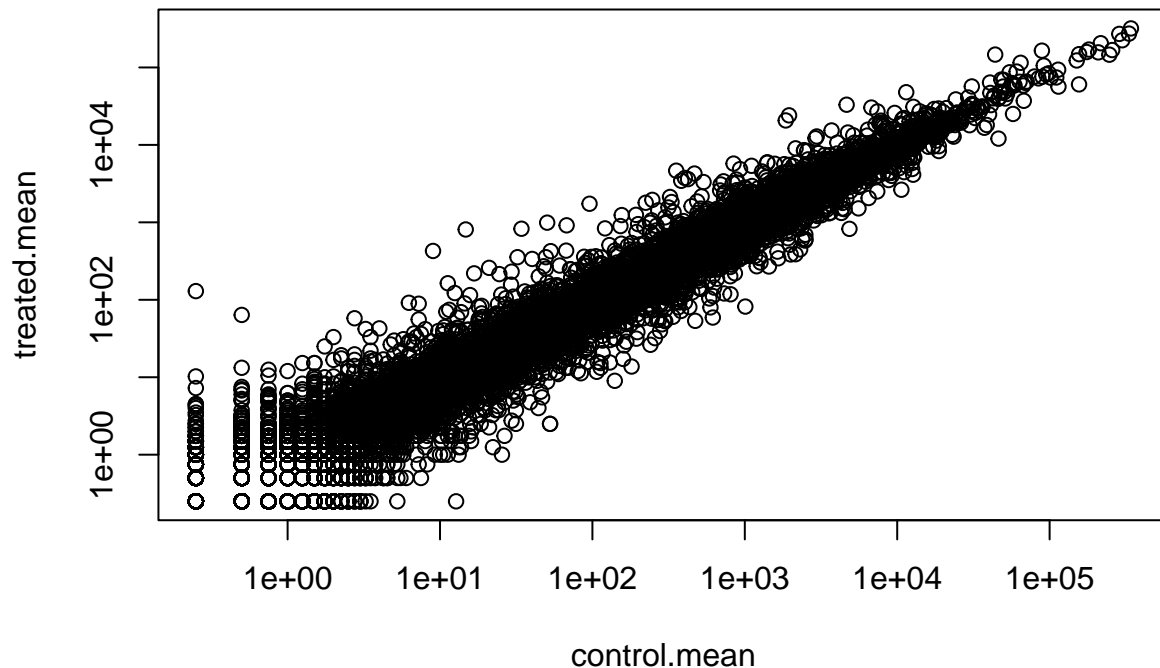
Each dot is a gene. Dots along the linear diagonal see no expression change. Below the diagonal is negative, above positive/increased in treated.

This would benefit from a log transform! Let's plot on a log scale

```
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 to make data more comprehensible/sensible.

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

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

Fold change of 1 if the treatment doubled in expression.

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

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

Fold change of -1 if the treatment expression decreased to half of control.

Cool, lets calculate the fold change of our control v. treated.

```
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
```

We need to drop the zero count genes/rows..

```
head(meancounts[,1:2] == 0)
```

```
##               control.mean treated.mean
## ENSG000000000003      FALSE      FALSE
## ENSG000000000005       TRUE       TRUE
## ENSG000000000419      FALSE      FALSE
## ENSG000000000457      FALSE      FALSE
## ENSG000000000460      FALSE      FALSE
## ENSG000000000938      FALSE       TRUE
```

The which() function tells us the indices of TRUE entries in a logical vector.

```
which( c(T,F,T))
```

```
## [1] 1 3
```

For our case we don't want to know the True positions within the vector, we want to know which rows have True.

```
inds <- which(meancounts[,1:2] == 0, arr.ind=TRUE)
head(inds)
```

```
##           row col
## ENSG000000000005    2  1
## ENSG000000004848   65  1
## ENSG000000004948   70  1
## ENSG000000005001   73  1
## ENSG000000006059  121  1
## ENSG000000006071  123  1
```

I only care about the rows here (if there is a 0 in any column I will exclude this eventually)

```
to.rm <- unique(sort(inds[, "row"]))
head(meancounts[-to.rm])
```

```
##               control.mean      log2fc
## ENSG000000000003      900.75 -0.45303916
## ENSG000000000005         0.00         NaN
## ENSG000000000419      520.50  0.06900279
## ENSG000000000457      339.75 -0.10226805
## ENSG000000000460       97.25 -0.30441833
## ENSG000000000938        0.75      -Inf
```

```
mycounts <- meancounts[-to.rm,]
```

We now have 21817 genes remaining.

```
nrow(mycounts)
```

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

Fold-change threshold of +2 or greater?

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

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

What percentage is this?

```
round((sum(mycounts$log2fc > +2) / nrow(mycounts)) *100, 2)
```

```
## [1] 1.15
```

```
sum(mycounts < -2)
```

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

```
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, 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 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, colAveragesPerRowSet, 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, rowAveragesPerColSet,
##   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
```

We first need to setup the DESeq input object.

```
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
```

Run the DESeq analysis pipeline.

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

```
## estimating size factors
```

```
## estimating dispersions
```

```
## gene-wise dispersion estimates
```

```
## mean-dispersion relationship
```

```
## final dispersion estimates
```

```
## fitting model and testing
```

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



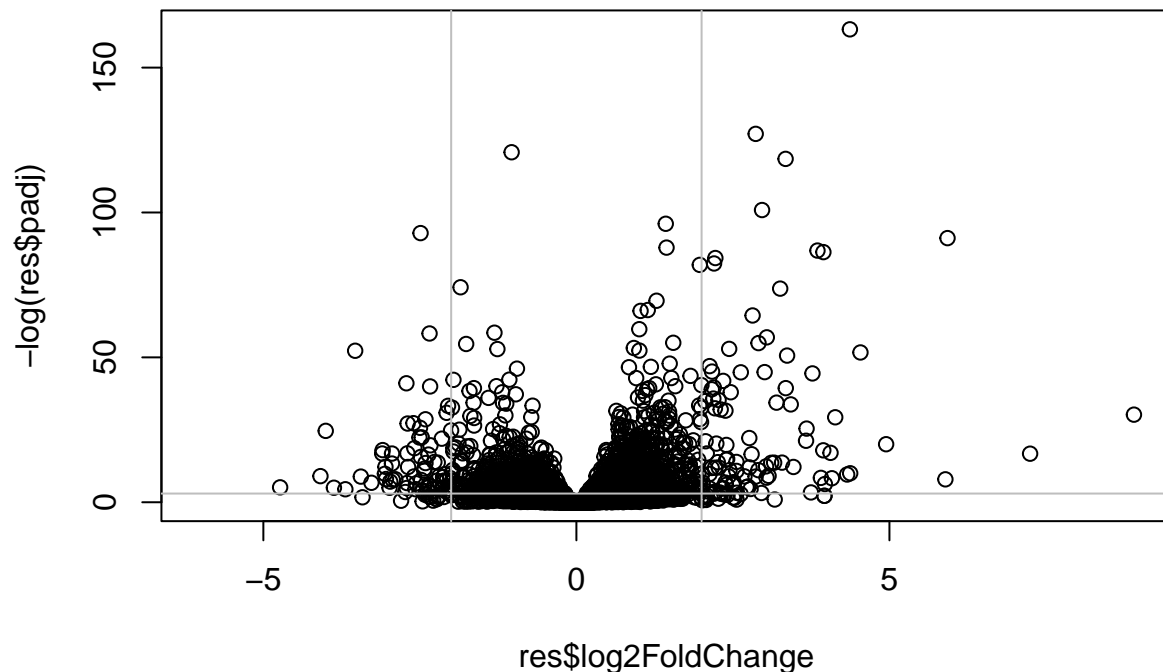
```
## log2 fold change (MLE): dex treated vs control
## Wald test p-value: dex treated vs control
## DataFrame with 6 rows and 6 columns
##           baseMean log2FoldChange      lfcSE      stat      pvalue
##           <numeric>      <numeric> <numeric> <numeric> <numeric>
## ENSG000000000003 747.194195    -0.3507030  0.168246 -2.084470 0.0371175
## ENSG000000000005   0.000000         NA         NA         NA         NA
## 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
```

Focus on the genes with a good p-value(low): plot the p-value against the log2 with volcano plot

A Volcano Plot

This is a very common data viz of this type of data that does not really look like a volcano.

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



Adding annotation data

We want to add meaningful gene names to our dataset to determine which are up/down most significantly. For this we will use 2 bioconductor packages, one does the work and is called AnnotationDbi the other contains the data we are going to map between and is called “org.Hs.eg.db”

```
BiocManager::install("org.Hs.eg.db") BiocManager::install("AnnotationDbi")
```

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

```
## Warning: package 'AnnotationDbi' was built under R version 4.1.2
```

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

```
##
```

```
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"       "UCSCCKG"
## [26] "UNIPROT"
```

Here we map to “SYMBOL” the common gene name that the world understands and wants.

```
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>  
## 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
```

Lets finally save our results to date.

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

Pathway Analysis

Let's try to bring some biology insight back into this. We will start with KEGG.

```
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"
```

The above genes are in Entrez format.

Before we can use KEGG we need to get our gene identifiers in the correct format for KEGG, which is ENTREZ format in this case.

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

```
## [1] "ENSG00000000003" "ENSG00000000005" "ENSG000000000419" "ENSG000000000457"
## [5] "ENSG000000000460" "ENSG000000000938"
```

```
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$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
```

```
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
```

The main `gage()` function requires a named vector of fold changes, where the names of the values are the Entrez gene IDs.

Note that we used the `mapIds()` function above to obtain Entrez gene IDs (stored in `res$entrez`) and we have the fold change results

```
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
```

```
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
```

The `pathway()` function will add out genes to a KEGG pathway as colored entries:

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

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

```
## Info: Working in directory /Users/abigailjaquish/Documents/Bioinformatics/R/BGGN213_github/class15
```

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

1. Data import/read countdata coldata(metadata)
2. PCA (QC) if happy then...
3. Deseq analysis (DESeq() function)
4. Figures: Volcano Plot (ggplot or plot function)
5. Annotation (use biology to understand the up/down regulated) KEGG, GO, etc.
6. Pathway Analysis (gage() function)