# Class 13: RNA seq analysis with DESeq2

Hsiang-Ying Lu (PID: A15608316)

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).

```
library(DESeq2)
```

## Data import

```
# 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         |

```r
dim(counts)
```

```
[1] 38694     8
```

```r
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
```

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

```
[1] 4
```

```r
table(metadata$dex)
```

```
control treated
      4       4
```

Q1. How many genes are in this dataset? 38694

Q2. How many 'control' cell lines do we have? 4

I want to compare the control to the treated columns. To do this I will

- Step 1. Identify and extract the "control" columns.
- Step 2. Calculate the mean value per gene for all these "control" columns and save as control.mean.
- Step 3. Do the same for treated
- Step 4. Compare the control.mean and treated.mean values.

Step 1:

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

```r
head(counts[,control.inds])
```

```
                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
```

```r
control.means <- rowMeans(counts[,control.inds])
head(control.means)
```

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

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

```r
head(counts[,treated.inds])
```

```
                SRR1039509 SRR1039513 SRR1039517 SRR1039521
ENSG00000000003        486        445       1097        604
ENSG00000000005          0          0          0          0
ENSG00000000419        523        371        781        509
ENSG00000000457        258        237        447        324
ENSG00000000460         81         66         94         74
ENSG00000000938          0          0          0          0
```

```r
treated.means <- rowMeans(counts[,treated.inds])
head(treated.means)
```

```
ENSG00000000003 ENSG00000000005 ENSG00000000419 ENSG00000000457 ENSG00000000460
         658.00            0.00          546.00          316.50           78.75
ENSG00000000938
           0.00
```
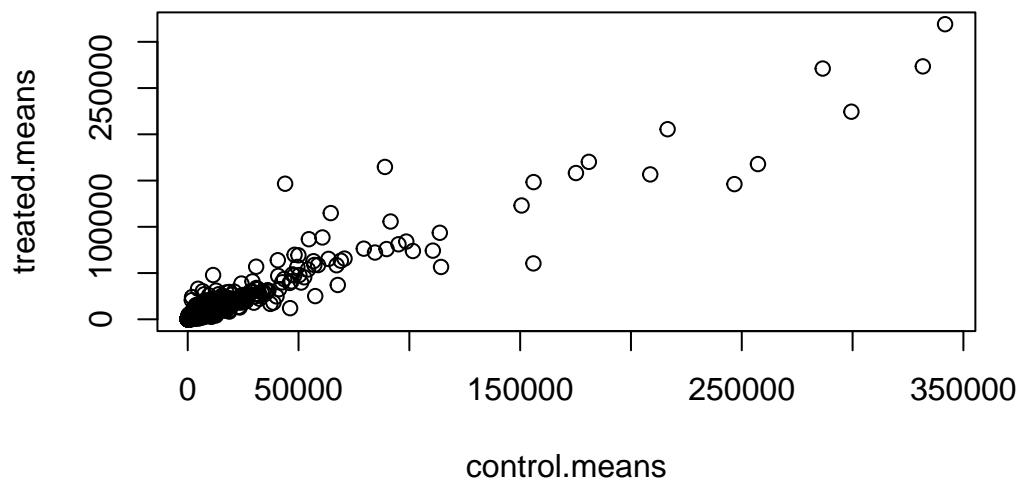
We will combine our meancount data for bookkeeping purpose

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

```
control.means treated.means
     23005324      22196524
```
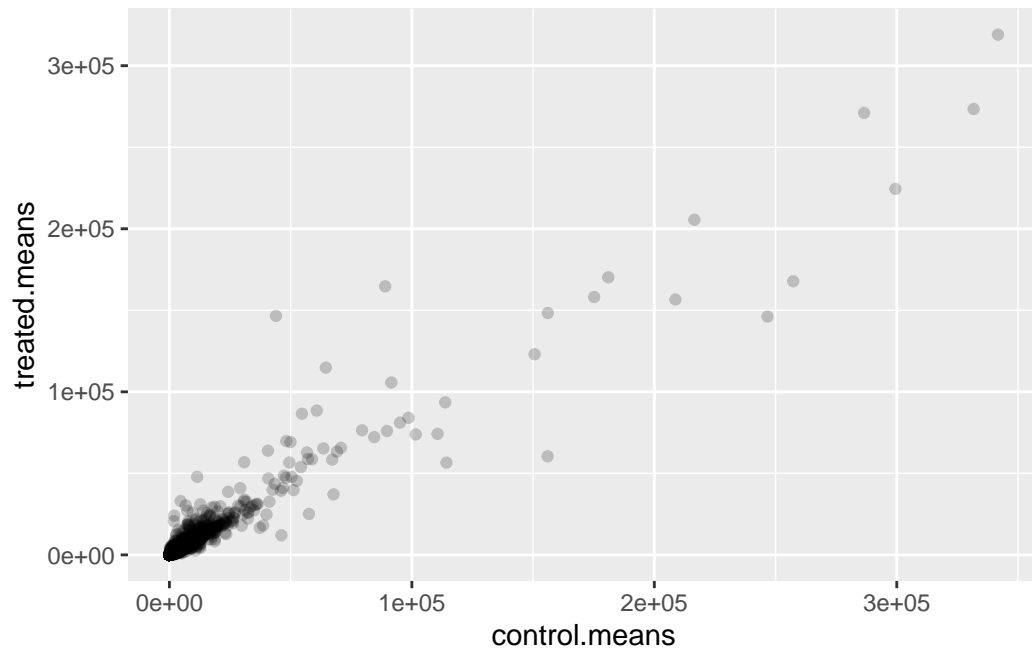
Let's see what these count values look like...

```
plot(meancounts)
```
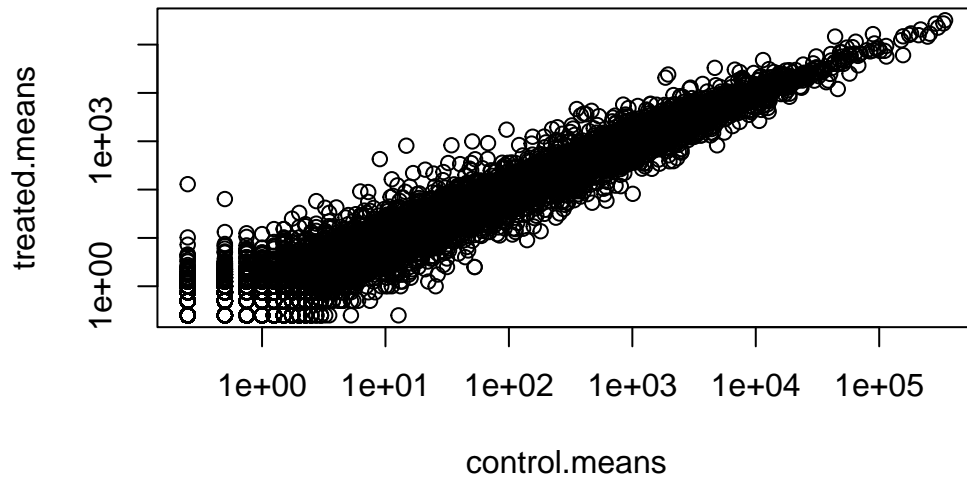


```
library(ggplot2)

ggplot(meancounts)+
  aes(control.means, treated.means)+
  geom_point(alpha=0.2)
```

```
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

Logs are super useful when we have such skewed data

```r
#Treated/control
log2(20/10)
```

```
[1] 1
```

Add log2(Fold-change) values to our wee results table.

```r
meancounts$log2fc <- log2(meancounts$treated.means/meancounts$control.means)
head(meancounts)
```

```
                control.means treated.means       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
```

I need to exclude any genes with zero counts as we can't say anything about them anyway from this experiment and it causes me math pain.

```
# What values in the first two cols are zero

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

Q. How many genes do I have left?

```
nrow(mycounts)
```

[1] 21817

Q. How many genes are "up-degulated" i.e. have a log2(fold-change) greater than +2?

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

[1] 250

Q. How many are "down-regulated"with a log2(fold-change) less than -2?

```
sum(mycounts$log2fc < -2)
```

[1] 367

Q10. Do you trust these results? Why or why not? No, because there's no information on statistical significant.

## Running DESeq

Like many bioconductor analysis packages DESeq wants it's input in a very particular way.

```
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

To run DESeq analysis we call the main function from the package called DESeq(dds)

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

estimating size factors

estimating dispersions

gene-wise dispersion estimates

mean-dispersion relationship

final dispersion estimates

fitting model and testing

To get the results out of this dds object we can use the DESeq results() function.

```
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
```
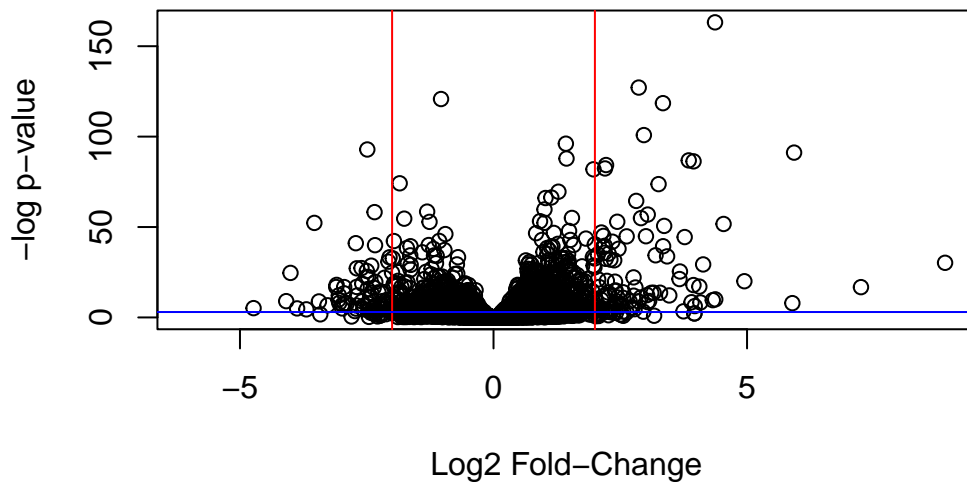
A common summary visualization is callsed a Volcano plot.

```r
plot(res$log2FoldChange, -log(res$padj),
     xlab="Log2 Fold-Change",
     ylab="-log p-value")
abline(v=c(-2, 2), col="red")
abline(h=-log(0.05), col="blue")
```



```r
mycols <- rep("gray", nrow(res))
mycols[res$log2FoldChange > 2] <- "black"
mycols[res$log2FoldChange < -2] <- "black"
mycols[res$padj>0.05] <- "gray"


plot(res$log2FoldChange, -log(res$padj), col=mycols,
     xlab="Log2 Fold-Change",
     ylab="-log p-value")
```
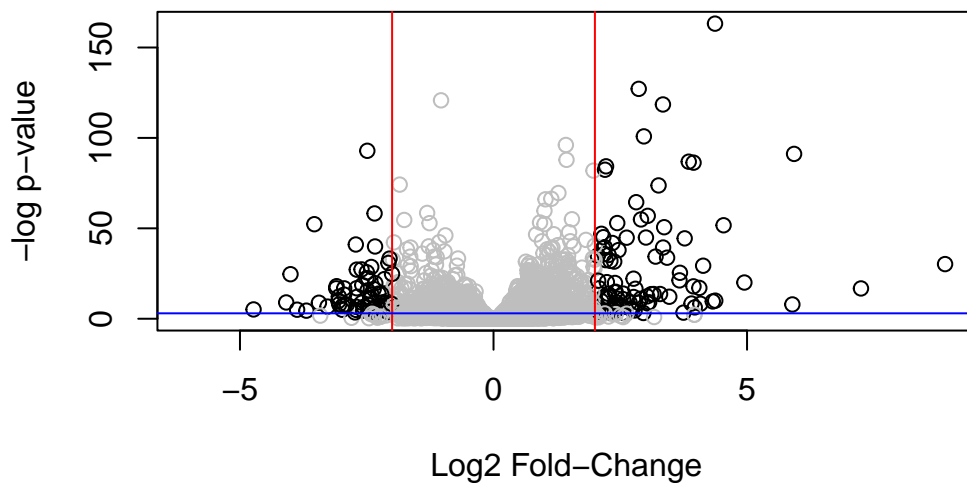
```
abline(v=c(-2, 2), col="red")
abline(h=-log(0.05), col="blue")
```



## Save our results to date

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

## adding annotation data

We need to translate or "map" our ensemble IDs into more understandable gene names

```
library("AnnotationDbi")
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"       "UCSCKG"
[26] "UNIPROT"
```

```r
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
```

```r
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
```

```r
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

```
  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      entrez
                 <numeric> <character> <character>
ENSG00000000003   0.163035      TSPAN6        7105
ENSG00000000005         NA        TNMD       64102
ENSG00000000419   0.176032        DPM1        8813
ENSG00000000457   0.961694       SCYL3       57147
ENSG00000000460   0.815849       FIRRM       55732
ENSG00000000938         NA         FGR        2268
```

```
  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

```
  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>
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       FIRRM       55732  A0A024R922
ENSG00000000938        NA         FGR        2268      P09769
```

```r
res$genenames <- 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       FIRRM       55732  A0A024R922
ENSG00000000938        NA         FGR        2268      P09769
                                genenames
                              <character>
ENSG00000000003          tetraspanin 6
ENSG00000000005            tenomodulin
ENSG00000000419 dolichyl-phosphate m..
ENSG00000000457 SCY1 like pseudokina..
ENSG00000000460 FIGNL1 interacting r..
ENSG00000000938 FGR proto-oncogene, ..
```

## 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)
```

```r
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"
```

```r
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
```

Run gage:

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

```r
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
```

Let's have a look at one of these pathways

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

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

```
Info: Working in directory /Users/sophialu1999/Desktop/UCSD Biological Sciences Ph.D./BGGN213
```

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