

Class 15: RNASeq Analysis

Julia Napoli

11/17/2021

Background

Today we're examining 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).

We need: 1) count data 2) col data

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

Let's take a look at each.

```
head(counts)
```

```
##          SRR1039508 SRR1039509 SRR1039512 SRR1039513 SRR1039516
## ENSG000000000003      723       486       904       445      1170
## ENSG000000000005       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
## ENSG000000000003     1097       806       604
## ENSG000000000005       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
```

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

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

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

Let's perform some exploratory differential gene expression analysis.

Compare control to treated

First we need to access all the control 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.

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

```
## ENSG0000000003 ENSG0000000005 ENSG00000000419 ENSG00000000457 ENSG00000000460  
##         900.75          0.00         520.50         339.75         97.25  
## ENSG00000000938  
##         0.75
```

Do the same for drug treated...

```
drug inds <- metadata$dex == 'treated'  
drug.ids <- metadata[drug inds,]$id  
drug.mean <- rowMeans(counts[,drug.ids])  
head(drug.mean)
```

```
## ENSG0000000003 ENSG0000000005 ENSG00000000419 ENSG00000000457 ENSG00000000460  
##         658.00          0.00         546.00         316.50         78.75  
## ENSG00000000938  
##         0.00
```

We will combine our meancount data for bookkeeping purposes.

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

```
##           control.mean drug.mean  
## ENSG0000000003      900.75    658.00  
## ENSG0000000005       0.00      0.00  
## ENSG00000000419     520.50    546.00  
## ENSG00000000457     339.75    316.50  
## ENSG00000000460      97.25     78.75  
## ENSG00000000938      0.75      0.00
```

There are 38694 rows/genes in this dataset.

```
nrow(counts)
```

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

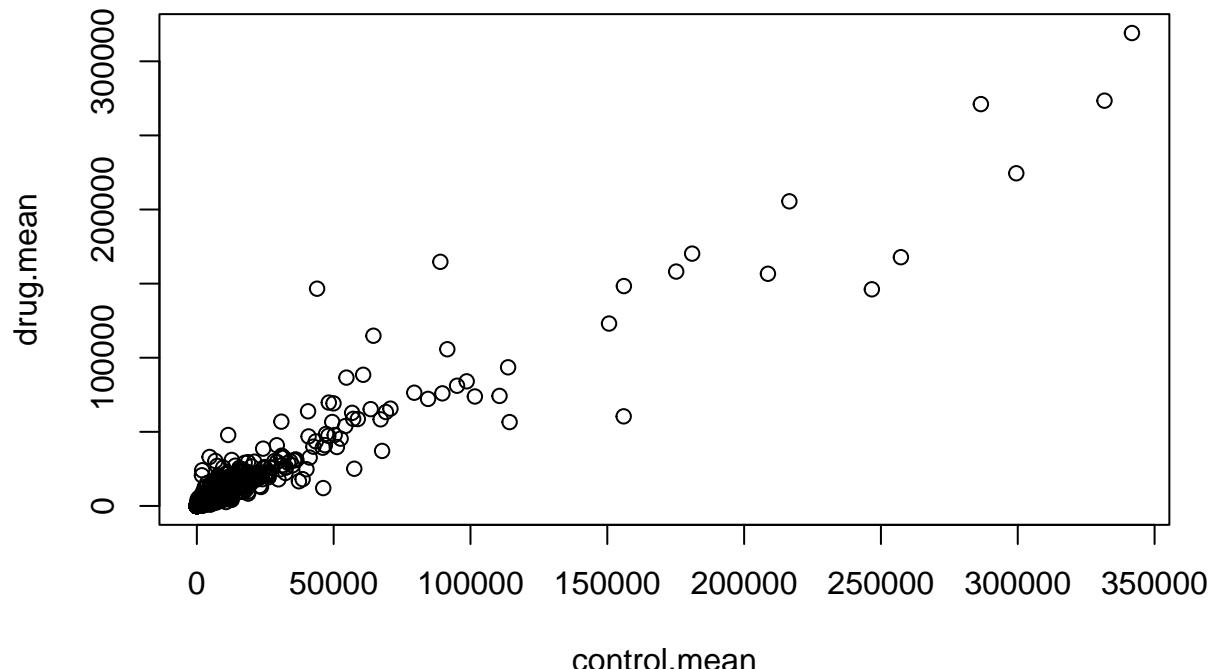
Compare the control and treated

```
colSums(meancounts)
```

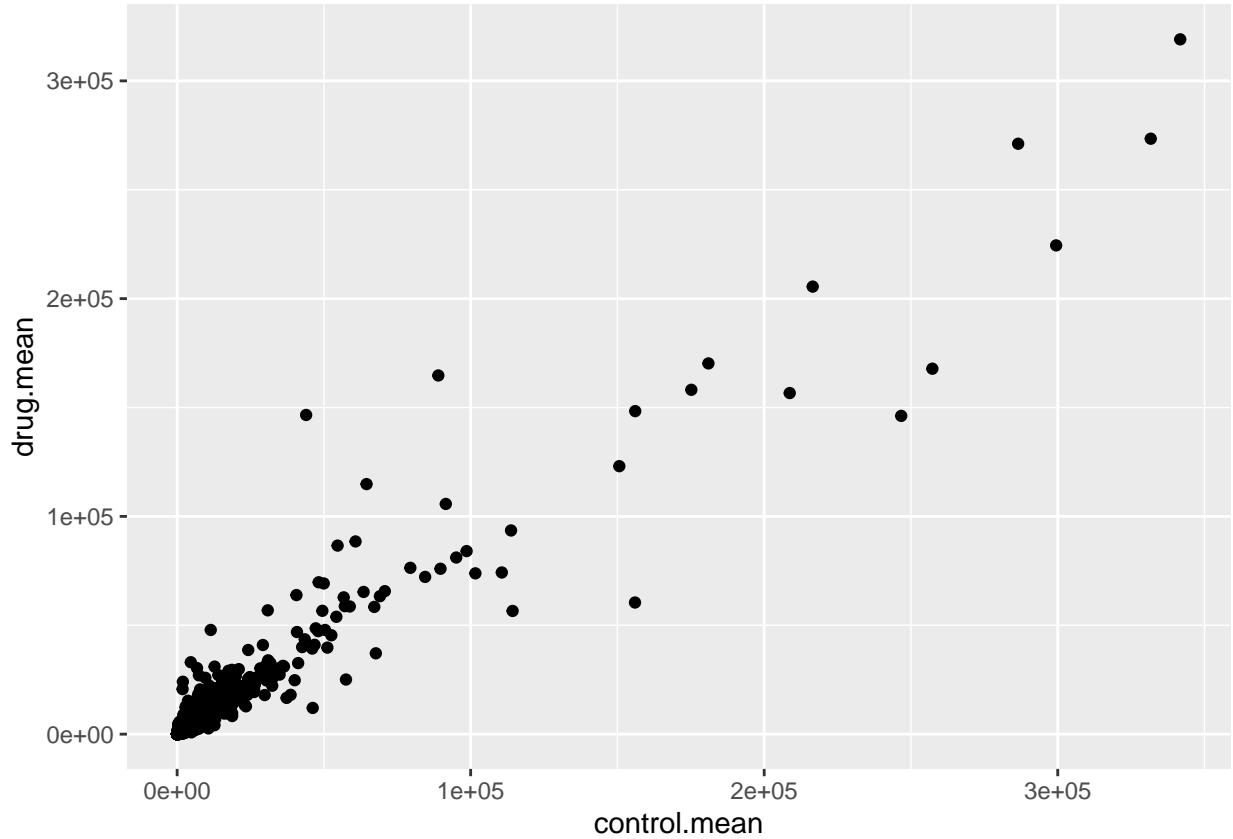
```
## control.mean      drug.mean
##      23005324      22196524
```

Let's visualize the results...

```
plot(meancounts)
library(ggplot2)
```



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

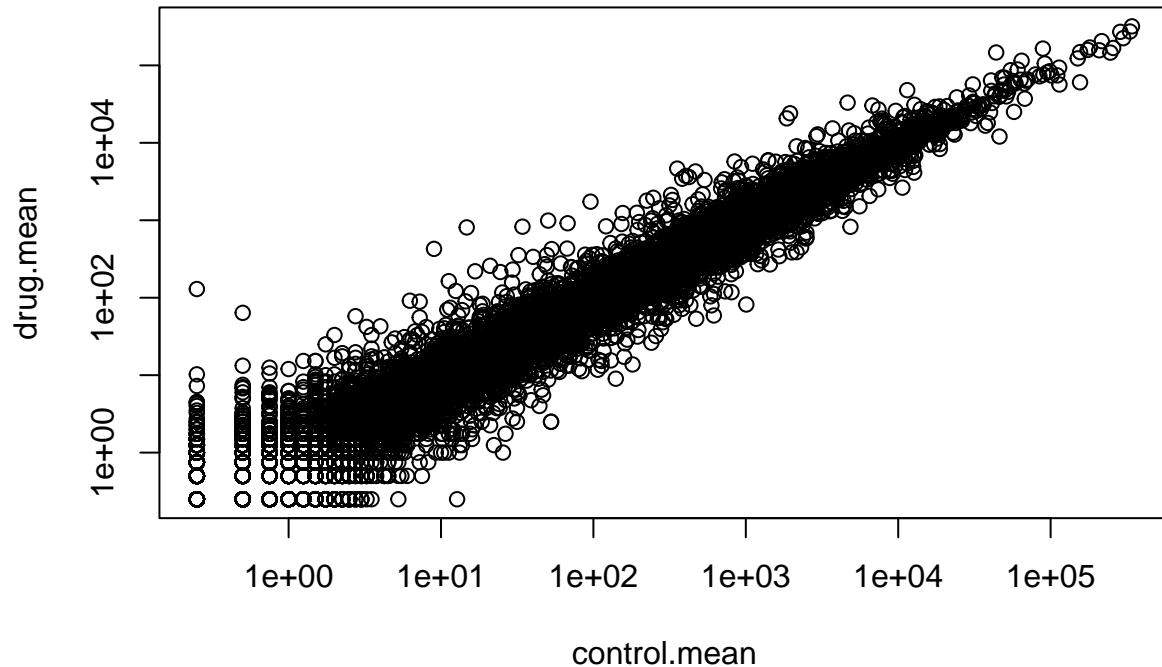


This would benefit from a log transform! Let's do that.

```
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 transforms, bc they make life easier sometimes :)

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

```
## [1] 0
```

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

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

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

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

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

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

```
meancounts$log2fc <- log2(meancounts[, "drug.mean"] / meancounts[, "control.mean"])
head(meancounts)
```

	control.mean	drug.mean	log2fc
## ENSG000000000003	900.75	658.00	-0.45303916

```

## ENSG000000000005      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 drop the zero count genes/rows!

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

```

##                  control.mean drug.mean
## ENSG000000000003     FALSE     FALSE
## ENSG000000000005     TRUE      TRUE
## ENSG00000000419     FALSE     FALSE
## ENSG00000000457     FALSE     FALSE
## ENSG00000000460     FALSE     FALSE
## ENSG00000000938     FALSE      TRUE

```

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

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

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

However, it is not that useful in default mode on our type of multi-column input.

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

```

##                  row col
## ENSG000000000005   2   1
## ENSG00000004848  65   1
## ENSG00000004948  70   1
## ENSG00000005001  73   1
## ENSG00000006059 121   1
## ENSG00000006071 123   1

```

I only care about rows here (if there is a zero in any column I will exclude this row eventually).

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

We now have 21817 genes remaining.

```
nrow(mycounts)
```

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

How many of these genes are up regulated at the log 2-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
```

How about down-regulated genes?

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

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

DESeq2 Analysis

```
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, 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

```

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

```

Look at the results.

```

res <- results(dds)
head(res)

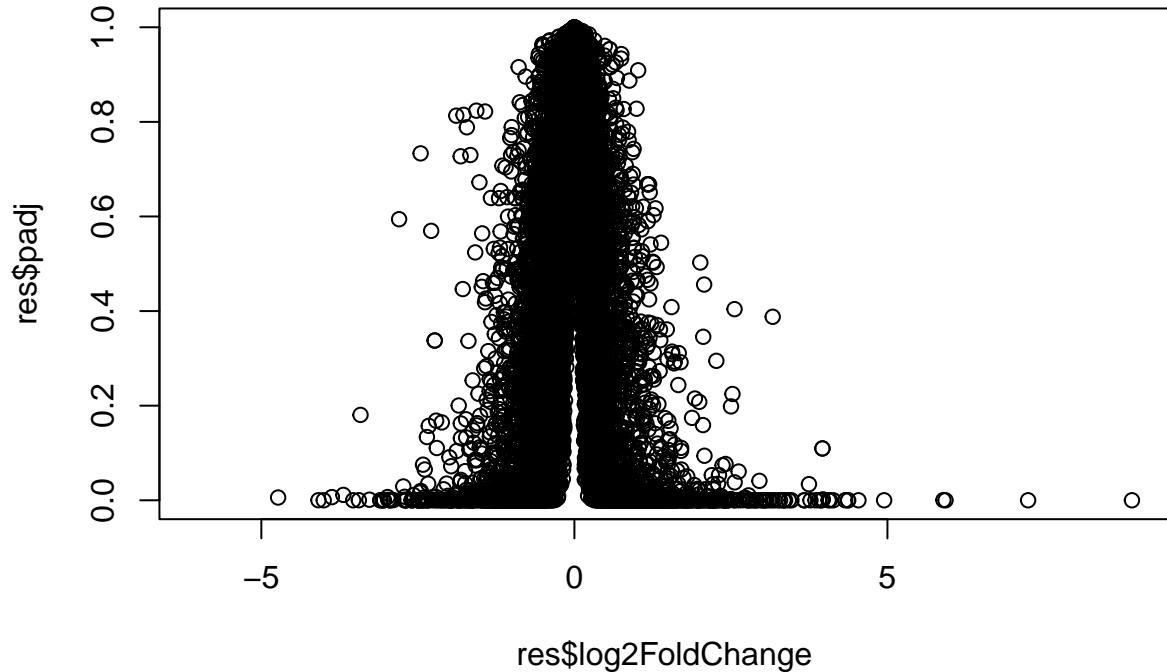
## log2 fold change (MLE): dex treated vs control
## Wald test p-value: dex treated vs control
## DataFrame with 6 rows and 6 columns
##           baseMean log2FoldChange      lfcSE      stat     pvalue
##           <numeric>      <numeric> <numeric> <numeric> <numeric>
## ENSG00000000003 747.194195     -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

```

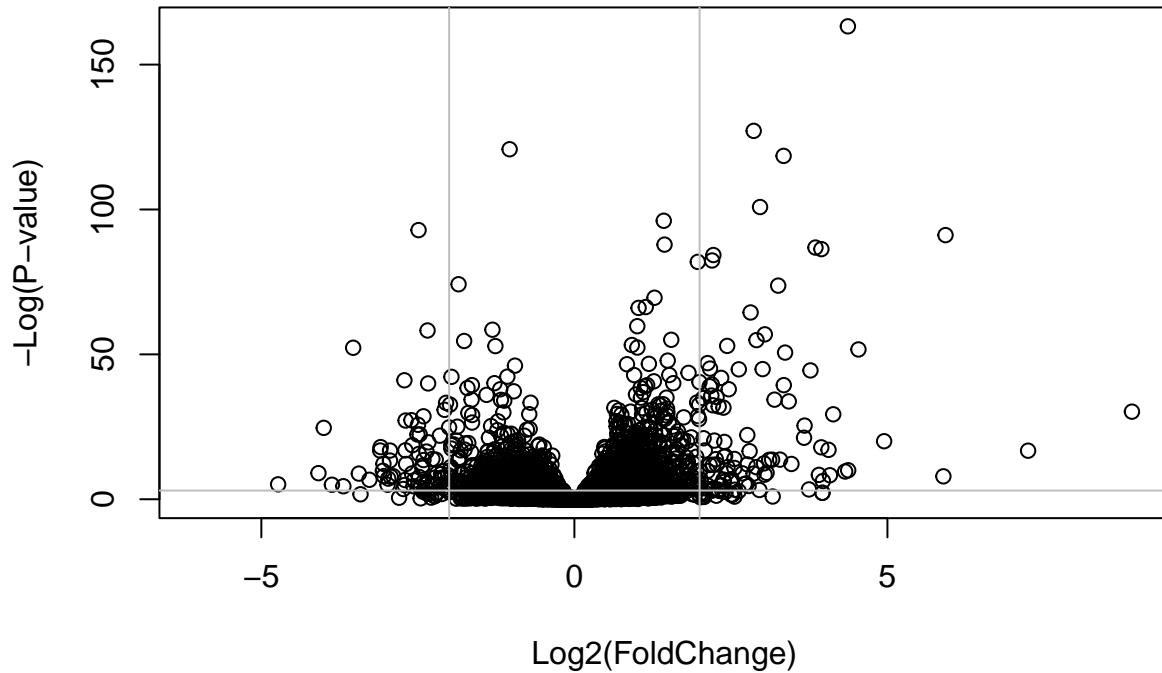
Volcano plot

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

```
plot(res$log2FoldChange,res$padj)
```



```
plot(res$log2FoldChange,-log(res$padj), xlab="Log2(FoldChange)",  
     ylab="-Log(P-value)")  
abline(v=c(-2,2),col = "gray")  
abline(h = -log(0.05),col = "gray")
```



Adding annotation data

Let's add some meaningful gene names to our dataset so we can make sense of what is going on here!

```
library("AnnotationDbi")

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

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

##
columns(org.Hs.eg.db)

## [1] "ACNUM"          "ALIAS"           "ENSEMBL"         "ENSEMLPROT"      "ENSEMLTRANS"
## [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"
```

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),           # Our gene names
                      keytype="ENSEMBL",            # The format of our gene names
                      column="SYMBOL",             # The new format we want to add
                      multiVals="first")

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

head(res$symbol)

## ENSG00000000003 ENSG00000000005 ENSG00000000419 ENSG00000000457 ENSG00000000460
##      "TSPAN6"          "TNMD"          "DPM1"          "SCYL3"          "C1orf112"
## ENSG00000000938
##      "FGR"

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

```

Let's finally save our data to a file

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

Pathway analysis

Let's try to bring some biology insight back into this work. For 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"

```

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" "ENSG00000000419" "ENSG00000000457"
## [5] "ENSG00000000460" "ENSG00000000938"

columns(org.Hs.eg.db)

## [1] "ACNUM"      "ALIAS"       "ENSEMBL"     "ENSEMLPROT"  "ENSEMLTRANS"
## [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

```

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>
## ENSG0000000003 747.194195 -0.3507030  0.168246 -2.084470 0.0371175
## ENSG0000000005  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      genename
##          <numeric> <character> <character> <character>
## ENSG0000000003  0.163035   TSPAN6      7105      tetraspanin 6
## ENSG0000000005   NA        TNMD       64102      tenomodulin
## ENSG00000000419  0.176032   DPM1       8813      dolichyl-phosphate m..
## ENSG00000000457  0.961694   SCYL3      57147      SCY1 like pseudokina..
## ENSG00000000460  0.815849   C1orf112    55732      chromosome 1 open re..
## ENSG00000000938   NA        FGR        2268      FGR proto-oncogene, ..

```

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 foldchange 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

```

```
# Get the results

keggres = gage(foldchanges, gsets=kegg.sets.hs)
```

We can look at the attributes() of this or indeed any R object

```
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 pathwview() function will add our 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/julianapolli/Library/Mobile Documents/com~apple~CloudDocs/Year1_Q1/
```

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

