

Class 15 RNA Seq

Stefanie Hodapp (PID: A53300084)

11/17/2021

Load the contData 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)
```

```
##           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
```

Q1. How many genes are in this dataset? 38694

```
nrow(counts)
```

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

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

```
View(metadata)
```

Check the correspondance of the metadata with the counts data

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

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

Compare control to treated

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

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

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

```
## [1] "SRR1039508" "SRR1039512" "SRR1039516" "SRR1039520"
```

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

```
control.mean <- rowMeans(counts[ , control.inds])  
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 drug treated...

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

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

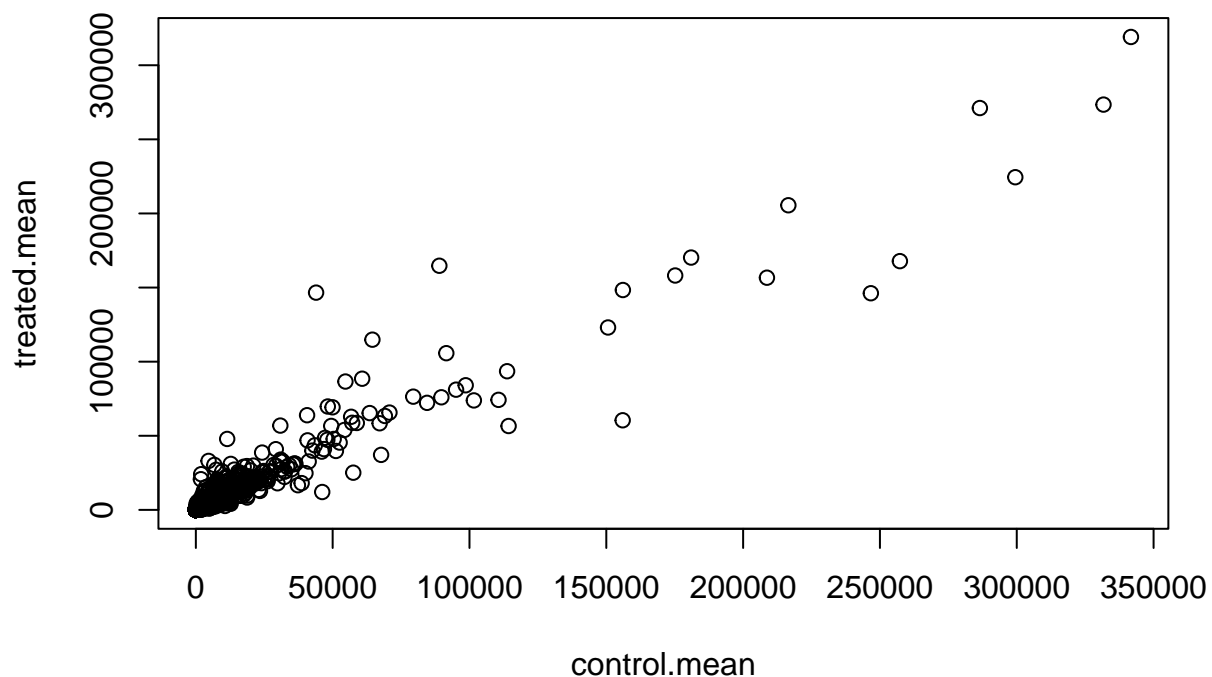
We will combine our meancount data for bookkeeping purposes.

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

Compare the control and treated

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)
```



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?

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

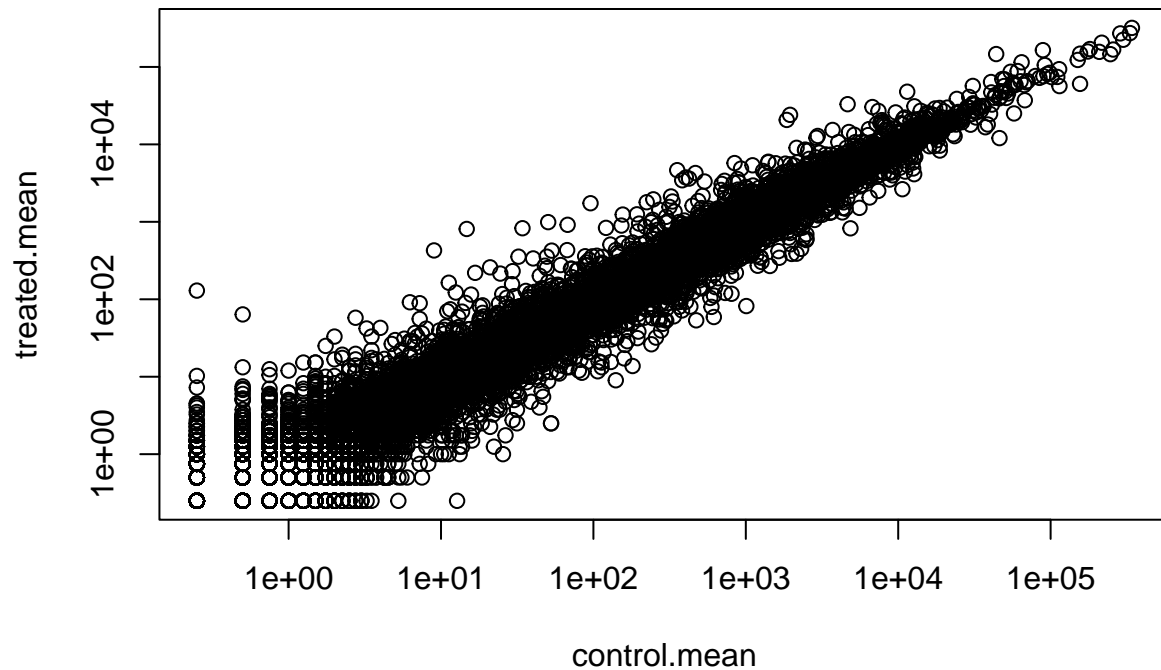
Q6. Try plotting both axes on a log scale. What is the argument to plot() that allows you to do this?

Try plotting both axes 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
```



Let's calculate the log2 fold change between control and dex-treated samples.

```
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
```

```
## ENSG00000000460      FALSE      FALSE
## ENSG00000000938      FALSE      TRUE
```

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? `arr.ind=TRUE` returns the positions in a matrix that are “TRUE” as integers. The `unique()` function will delete duplicate rows (e.g. genes with zero counts in both control and treated samples).

The ‘`which()`’ functions tell us the indices of TRUE entries in a logical vector

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

```
##           row col
## ENSG00000000005    2  1
## ENSG00000004848   65  1
## ENSG00000004948   70  1
## ENSG00000005001   73  1
## ENSG00000006059  121  1
## ENSG00000006071  123  1
```

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

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

Remove the genes with zero expression from our `meancounts` data

```
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
```

We now have “`nrow(mycounts)`” genes 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
```

What percentage is this?

```
round((sum(up.ind) / nrow(mycounts))*100, 2)
```

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

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
```

```
round((sum(down.ind) / nrow(mycounts))*100, 2)
```

```
## [1] 1.68
```

Q10. Do you trust these results? Why or why not? I don't trust that these results are significant as they only take into account fold change and not p-values (i.e. significance). Genes exhibiting a 2-fold change may not be significant.

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, colAvgPerRowSet, 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, rowAvgPerColSet,
##     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
```

```
dds
```

```
## class: DESeqDataSet  
## dim: 38694 8  
## metadata(1): version  
## assays(1): counts  
## rownames(38694): ENSG000000000003 ENSG000000000005 ... ENSG00000283120  
## ENSG00000283123  
## rowData names(0):  
## colnames(8): SRR1039508 SRR1039509 ... SRR1039520 SRR1039521  
## colData names(4): id dex celltype geo_id
```

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
```



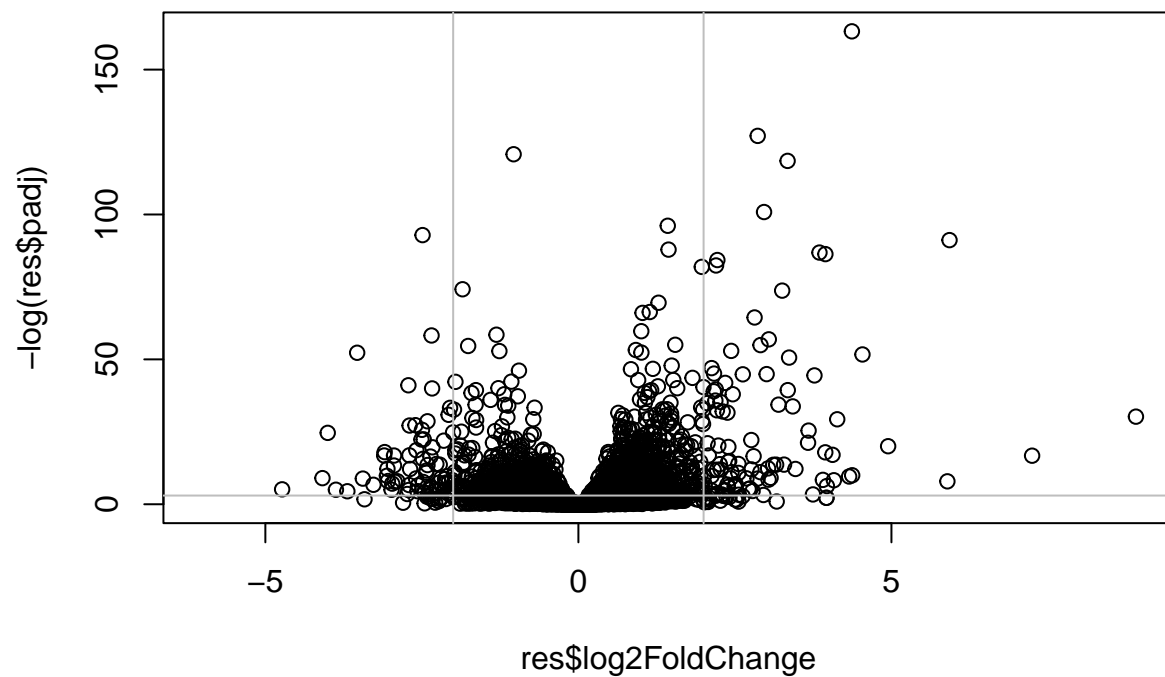
```
res <- results(dds)
head(res)
```

```
## 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
```

A Volcano plot

Volcano plots highlight the proportion of genes that are both significantly regulated and display a high fold change

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



Adding annotation data

```
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"       "UCSCKG"
## [26] "UNIPROT"
```

Here, we map to "SYMBOL", the common gene names

```
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$symbol)

## ENSG00000000003 ENSG00000000005 ENSG00000000419 ENSG00000000457 ENSG00000000460
##      "TSPAN6"      "TNMD"      "DPM1"      "SCYL3"      "Clorf112"
## ENSG00000000938
##      "FGR"
```

Let's save our results to date

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

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.

```
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 in res\$log2FoldChange.

```
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, let's run the gage pathway analysis.

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

We can look at the attributes() of this

