# class13_fr

Cameron Finch (A16734770)

Import csv files with `read.csv()`

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

```r
m <- metadata$id
c <- colnames(counts)
all(m == c)
```

```
[1] TRUE
```

Q1. How many genes in this dataset?

```r
nrow(counts)
```

```
[1] 38694
```

A1. 38694

Q2. How many 'control' cell lines?

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

```
[1] 4
```

A2. 4

```r
# Code approach 1
control <- metadata[metadata[,"dex"]=="control",]
control.counts <- counts[ ,control$id]
control.mean <- rowSums( control.counts )/4
head(control.mean)
```

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

```r
# Code approach 2 (yields same result)
library(dplyr)
```

```
Attaching package: 'dplyr'

The following objects are masked from 'package:stats':

    filter, lag

The following objects are masked from 'package:base':

    intersect, setdiff, setequal, union
```

```r
control <- metadata %>% filter(dex=="control")
control.counts <- counts %>% select(control$id)
control.mean <- rowSums(control.counts)/4
head(control.mean)
```

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

Q3. How would you make the above code in either approach more robust? Is there a function that could help here?

**The code above helps us find the IDs of the control samples, but the `...(counts)/4` function limits the code to function correctly only when there are four control samples. If the sample was different, it would need a different code to function the same. The following approach yields the same results but is more universal.**

A3:

```r
control <- metadata[metadata$dex == "control", ]
cc <- ( counts[ , control$id] )
cm <- rowMeans(cc)
head(cm)
```

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

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)

A4. Repeating the same process for the treatment samples:

```
treated <- metadata[metadata$dex == "treated", ]
tc <- counts[ , treated$id]
tm <- rowMeans(tc)
head(tm)
```
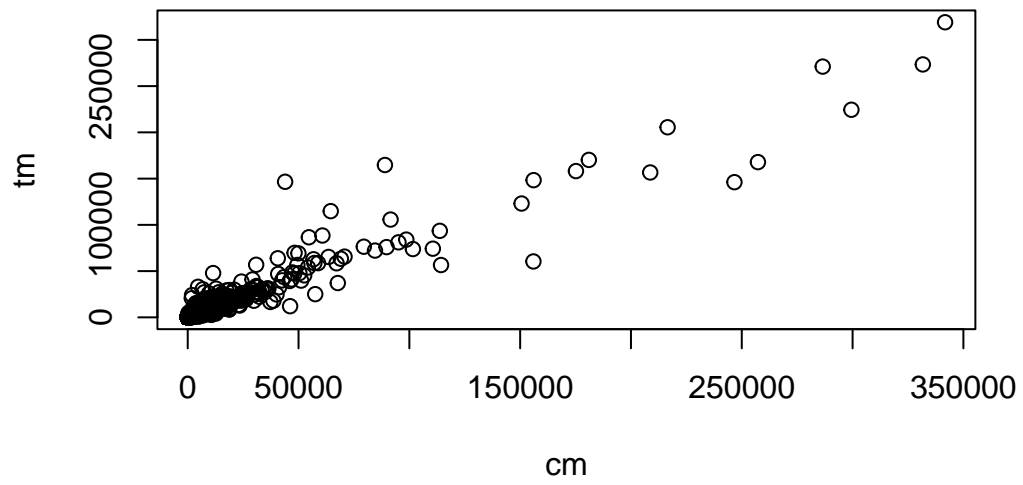
```
ENSG00000000003 ENSG00000000005 ENSG00000000419 ENSG00000000457 ENSG00000000460
         658.00            0.00          546.00          316.50           78.75
ENSG00000000938
           0.00
```

## Combine results in data frame

```
meancounts <- data.frame(cm, tm)
View(meancounts)
```
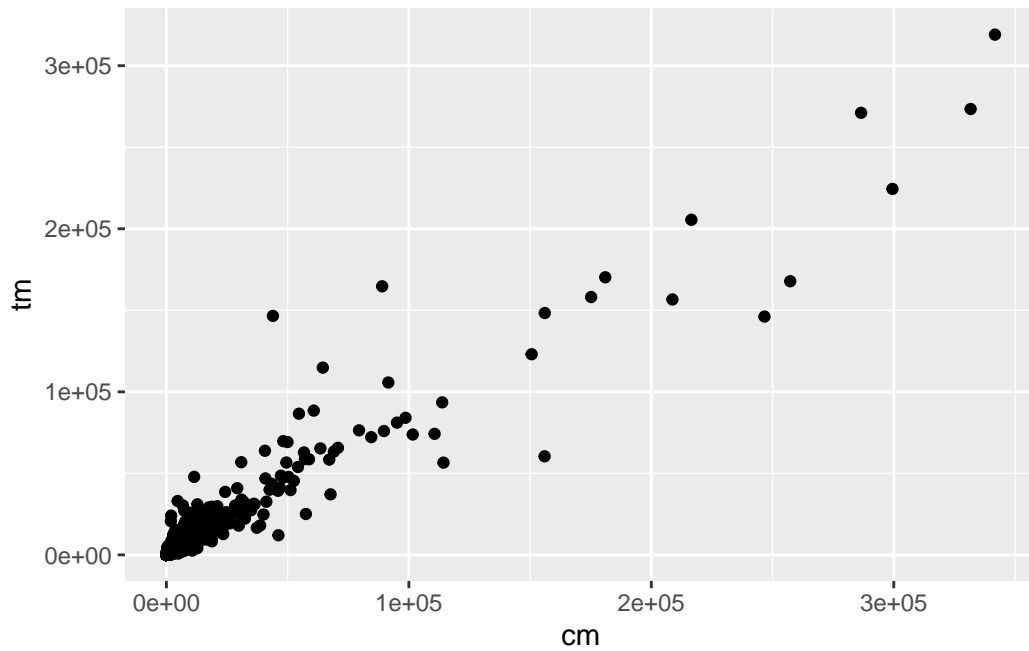
Q5a. Plot results showing correlation between treated and control data using base
R plots A5a:

```
plot(cm, tm)
library(ggplot2)
```

4

Q5b. Recreate plot using ggplot2 A5b:

```
ggplot(meancounts) +
  aes(cm, tm) +
  geom_point()
```
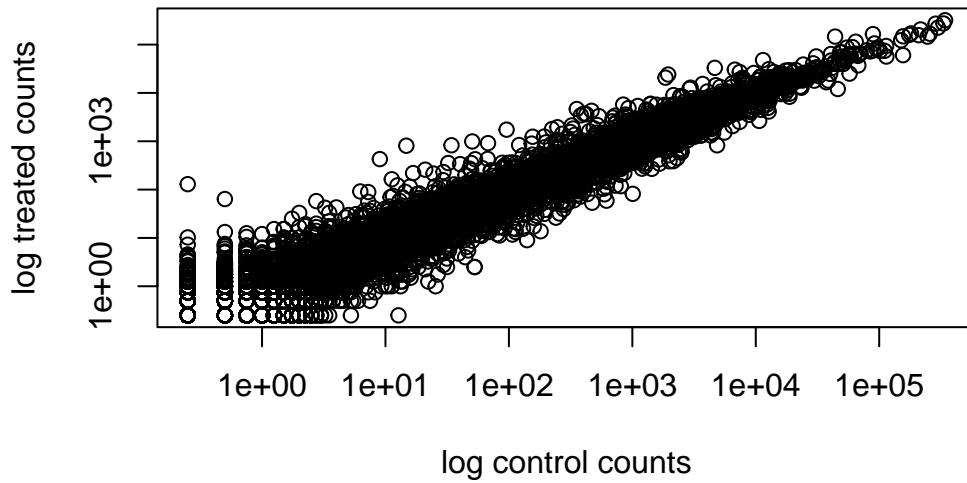
Q6. Make a logarithmic plot to get a better view on the data close to the origin

A6:

```
plot(cm, tm, log="xy", xlab="log control counts", ylab="log 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

6

```r
meancounts$log2fc <- log2(meancounts[,"tm"]/meancounts[,"cm"])
head(meancounts)
```

```
                   cm     tm       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
```

```r
zero.vals <- which(meancounts[,1:2]==0, arr.ind=T)
to.rm <- unique(zero.vals[,1])
mycounts <- meancounts[-to.rm,]
head(mycounts)
```

```
                    cm      tm       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=TRUE` argument? A7. `arr.ind = T` will help organize the data into rows and columns so the elements of the data frame can be separated

Q8. Using the up.ind vector above can you determine how many up regulated genes we have at the greater than 2 fc level?

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

```
[1] 250
```

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

```
[1] 367
```

A8. 250

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

A9. 367

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

A10. Not entirely, because although we know each of these numbers indicate at least a 2fold change, we can't ensure that this indicates a significant difference

```
library(DESeq2)
```

```
Loading required package: S4Vectors
```

```
Loading required package: stats4
```

```
Loading required package: BiocGenerics
```

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

8

The following objects are masked from 'package:dplyr':

    combine, intersect, setdiff, union

The following objects are masked from 'package:stats':

    IQR, mad, sd, var, xtabs

The following objects are masked from 'package:base':

    anyDuplicated, aperm, 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:dplyr':

    first, rename

The following object is masked from 'package:utils':

    findMatches

The following objects are masked from 'package:base':

    expand.grid, I, unname

Loading required package: IRanges


Attaching package: 'IRanges'

The following objects are masked from 'package:dplyr':

    collapse, desc, slice

The following object is masked from 'package:grDevices':

    windows

Loading required package: GenomicRanges

Loading required package: GenomeInfoDb

Loading required package: SummarizedExperiment

Loading required package: MatrixGenerics

Loading required package: matrixStats


Attaching package: 'matrixStats'

The following object is masked from 'package:dplyr':

    count


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
```

11

```
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)
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
```

```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
```

> 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. # Making Volcano Plots

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

```
Attaching package: 'AnnotationDbi'

The following object is masked from 'package:dplyr':

    select
```

```r
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
res$entrez <- mapIds(org.Hs.eg.db,
                     keys=row.names(res),
                     column="ENTREZID",
                     keytype="ENSEMBL",
                     multiVals="first")
```

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

```r
res$uniprot <- mapIds(org.Hs.eg.db,
                      keys=row.names(res),
                      column="UNIPROT",
                      keytype="ENSEMBL",
                      multiVals="first")
```

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

```r
res$genename <- mapIds(org.Hs.eg.db,
                       keys=row.names(res),
                       column="GENENAME",
```

```
                          keytype="ENSEMBL",
                          multiVals="first")
```

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

```
  head(res)
```

```
log2 fold change (MLE): dex treated vs control
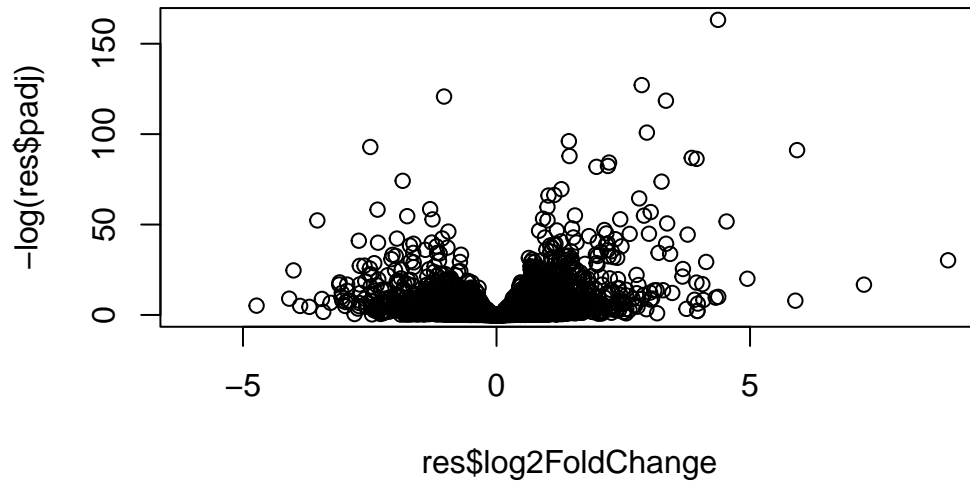Wald test p-value: dex treated vs control
DataFrame with 6 rows and 10 columns
                  baseMean log2FoldChange     lfcSE      stat    pvalue
                 <numeric>      <numeric> <numeric> <numeric> <numeric>
ENSG00000000003 747.194195     -0.3507030  0.168246 -2.084470 0.0371175
ENSG00000000005   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
                      genename
                   <character>
ENSG00000000003     tetraspanin 6
ENSG00000000005      tenomodulin
ENSG00000000419 dolichyl-phosphate m..
ENSG00000000457 SCY1 like pseudokina..
ENSG00000000460 FIGNL1 interacting r..
ENSG00000000938 FGR proto-oncogene, ..
```

## Volcano Plot

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



The genes far from the 'volcano spout' are the ones that have changed the most significantly.