

# class15

Hyeonseok Jang (PID# A59011126)

11/17/2021

```
#install.packages("BiocManager")  
#BiocManager::install()
```

```
#For this class, you'll also need DESeq2:  
#BiocManager::install("DESeq2")
```

```
library(BiocManager)  
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
```

```

##
## Attaching package: 'IRanges'

## 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: '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
```

```
## Import countData and colData
```

```
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
## ENSG000000000419        467          523          616          371          582
## ENSG000000000457        347          258          364          237          318
## ENSG000000000460         96           81           73           66          118
## ENSG000000000938          0           0           1           0           2
##                SRR1039517 SRR1039520 SRR1039521
## ENSG00000000003        1097          806          604
## ENSG00000000005           0           0           0
## ENSG000000000419        781          417          509
## ENSG000000000457        447          330          324
## ENSG000000000460         94          102           74
## ENSG000000000938          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
```

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

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

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

Q1. How many genes are in this dataset?

```
nrow(counts)
```

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

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

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

```
## [1] 4
```

```
## Compare control to treated
```

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

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

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

```
##           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
```

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

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

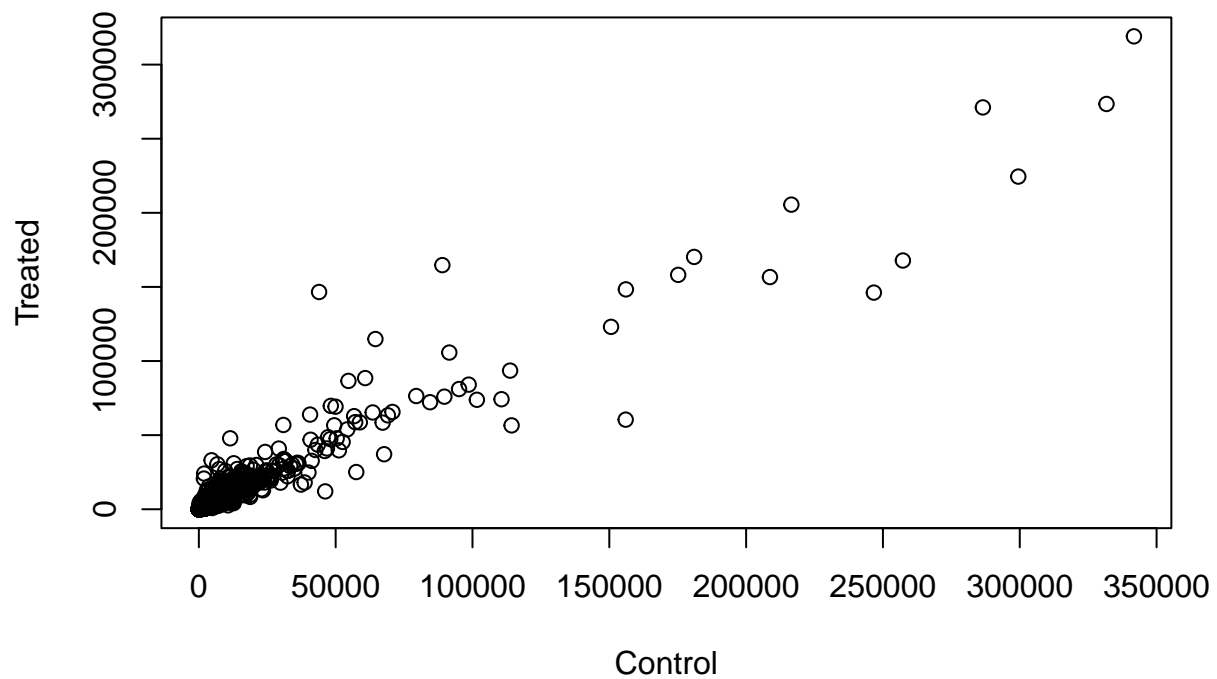
Q4. Follow the same procedure for the treated samples (i.e. calculate the mean per gene across drug treated samples and assign to a labeled vector called treated.mean)

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

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

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, xlab="Control", ylab="Treated")
```

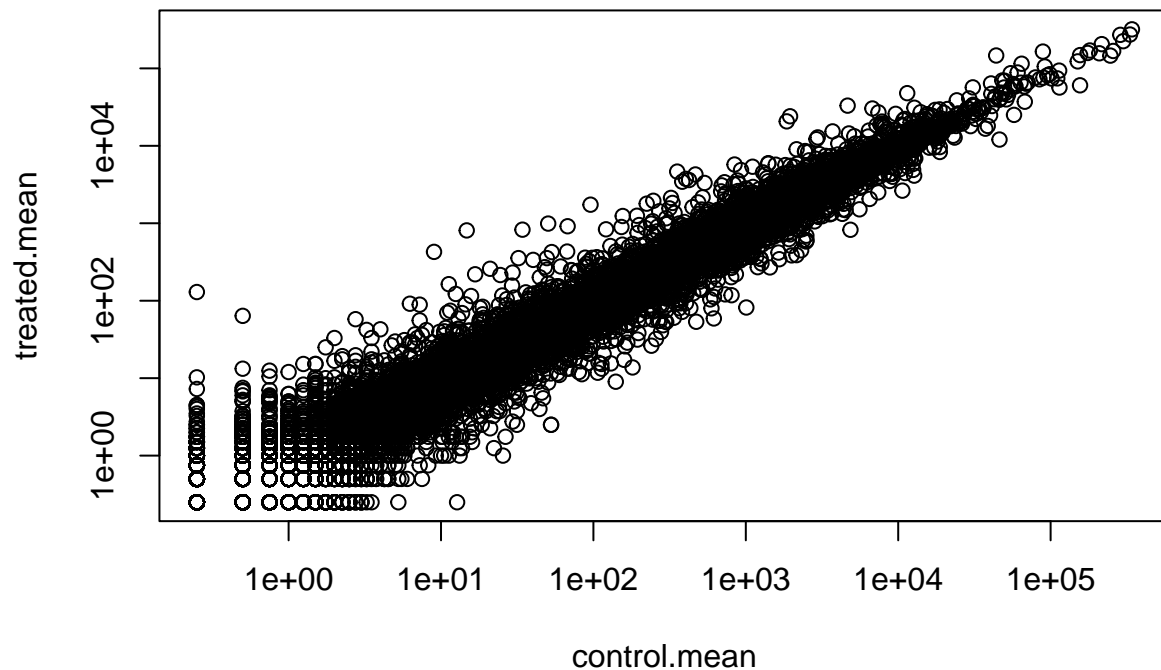


Q6. Try plotting both axes on a log scale. What is the argument to `plot()` that allows you to do this?

```
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
```



```
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
```

```
inds <- which(meancounts[,1:2]==0, arr.ind=TRUE)
to.rm <- unique(sort(inds[, "row"]))

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

What percentage of genes are above the fold-change threshold of +2 or greater?

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

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

How about down?

```
round(sum(mycounts$log2fc < -2)/nrow(mycounts) * 100, 2)
```

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

## DESeq2 analysis

We first need to setup the data for DESeq2

```
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
```

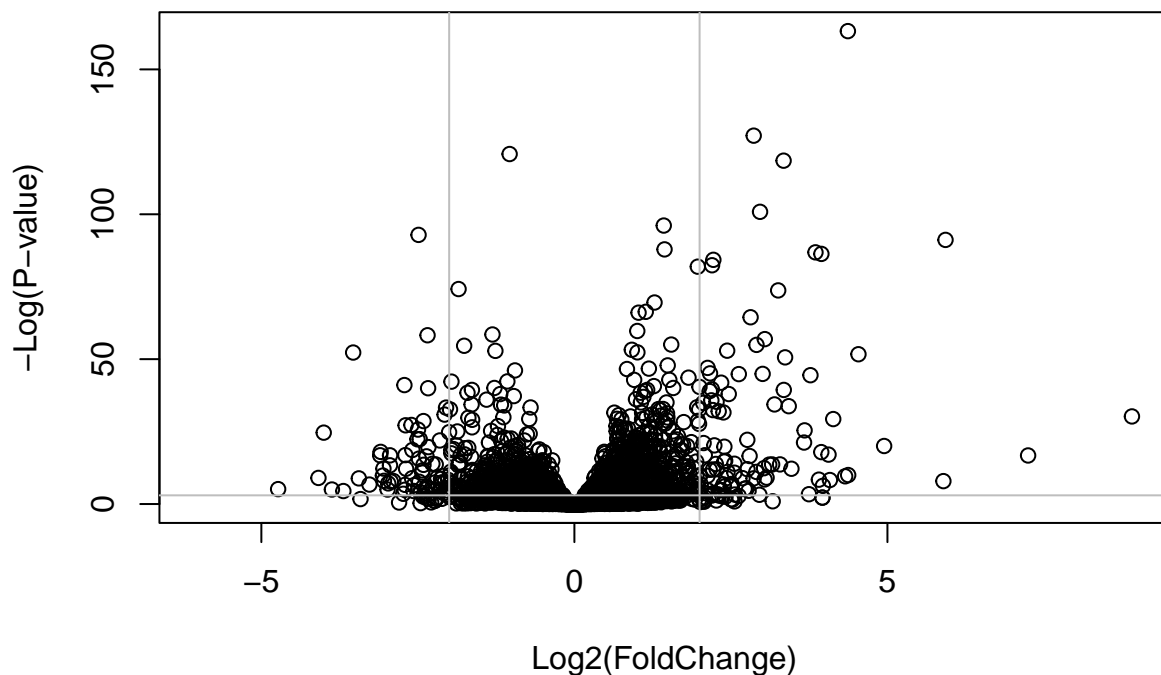
```
## DataFrame with 6 rows and 6 columns
```

```
##           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
```

## Volcano Plot

```
plot(res$log2FoldChange, -log(res$padj), ylab="-Log(P-value)", xlab="Log2(FoldChange)")
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 so we can make some sense of what is going on here...

For this we will use two bioconductor packages, one does the work and is called **AnnotationDbi** the other is to contain the data and called **org.Hs.eg.db**

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



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

```
##
```

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$symbol)
```

```
## ENSG000000000003 ENSG000000000005 ENSG000000000419 ENSG000000000457 ENSG000000000460
##      "TSPAN6"      "TNMD"      "DPM1"      "SCYL3"      "Clorf112"
## ENSG0000000000938
##      "FGR"
```

Lets finally save the result to data

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