

Class 15

Saumya Ranjan

11/16/2021

```
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.

```
nrow(counts)
```

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

Q2.

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

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

First I need to extract all the “control” columns. Then I will take the row wise mean to get the average count values for all genes in these four experiments

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

```

head(control.counts)

##          SRR1039508 SRR1039512 SRR1039516 SRR1039520
## ENSG000000000003     723      904     1170      806
## ENSG000000000005      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(control.counts)

```

Q3. Other method was prone to data shifts, which could result in getting the wrong answer everytime with no way of knowing the answer was wrong. This method isn't as likely to result in shifts in the data that could change the answers.

Q.4 (Doing the same for the treated experiments)

```

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

##          SRR1039509 SRR1039513 SRR1039517 SRR1039521
## ENSG000000000003     486      445     1097      604
## ENSG000000000005      0         0         0         0
## ENSG00000000419     523      371      781      509
## ENSG00000000457     258      237      447      324
## ENSG00000000460      81       66      94       74
## ENSG00000000938      0         0         0         0

treated.mean <- rowMeans(treated.counts)

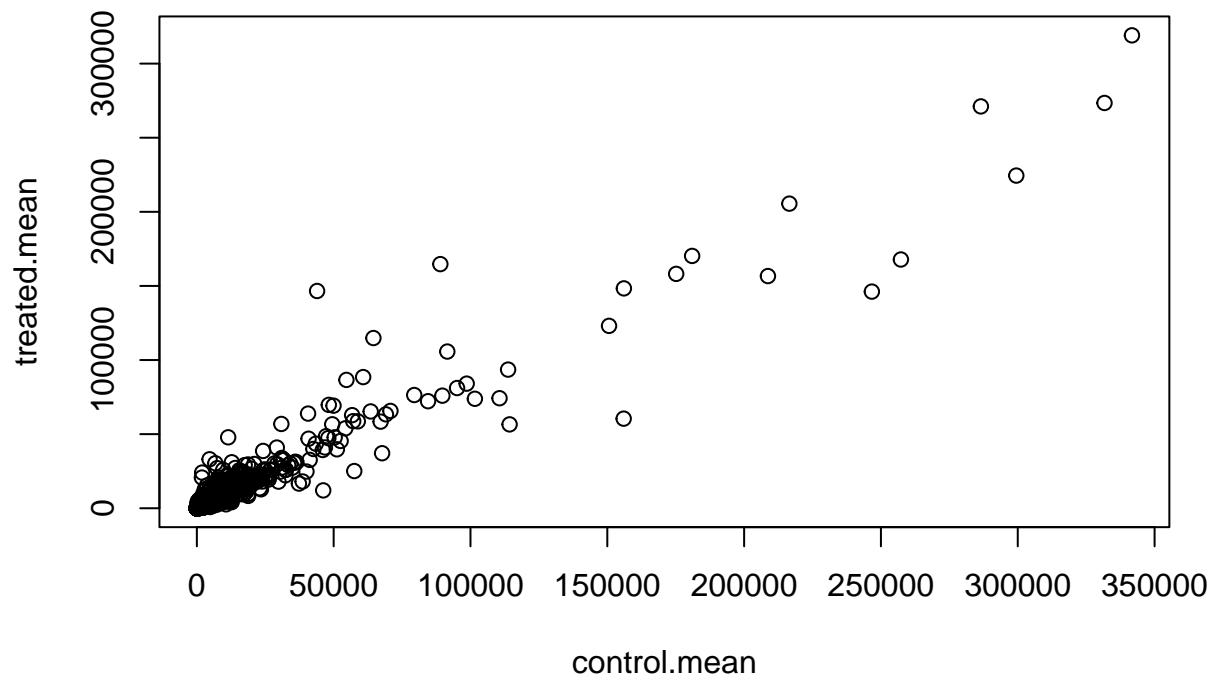
```

We combine our meancount data for bookkeeping purposes

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

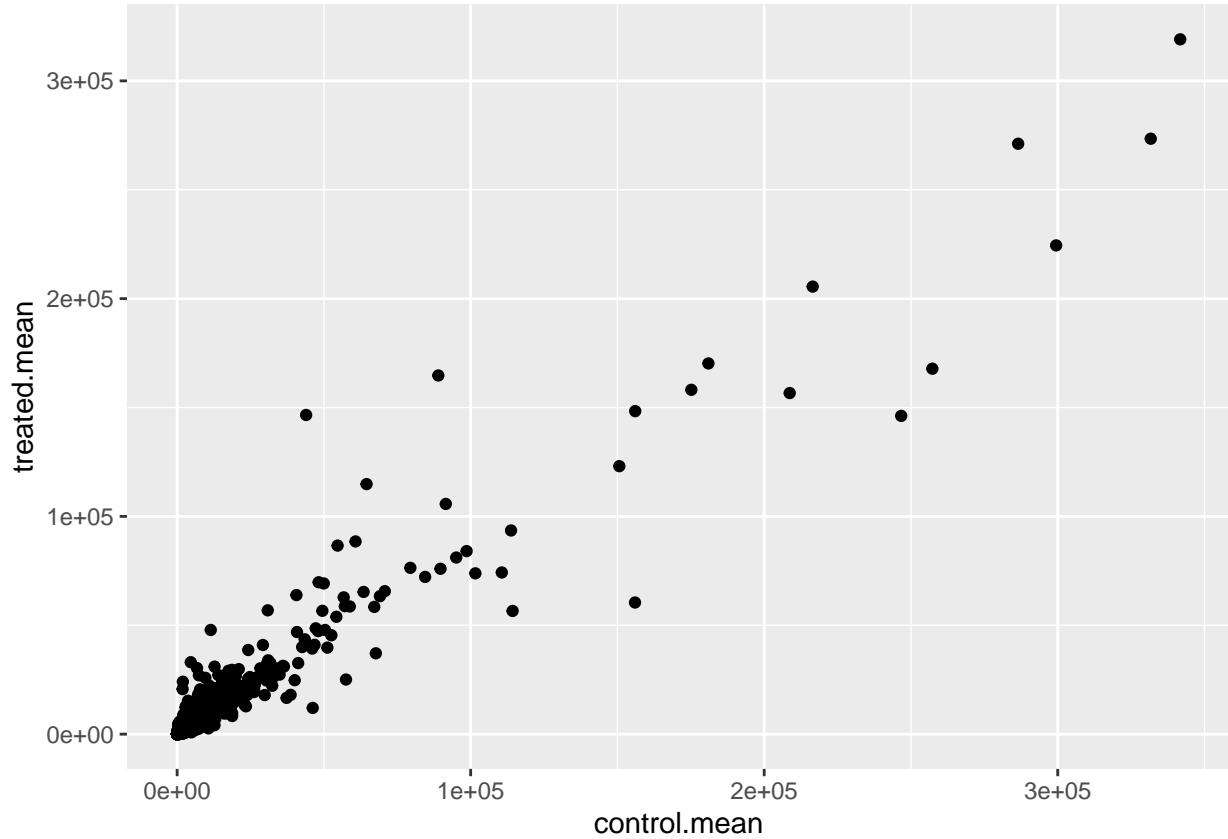
Q5.a Quick plot

```
plot(meancounts)
```



Q.5b

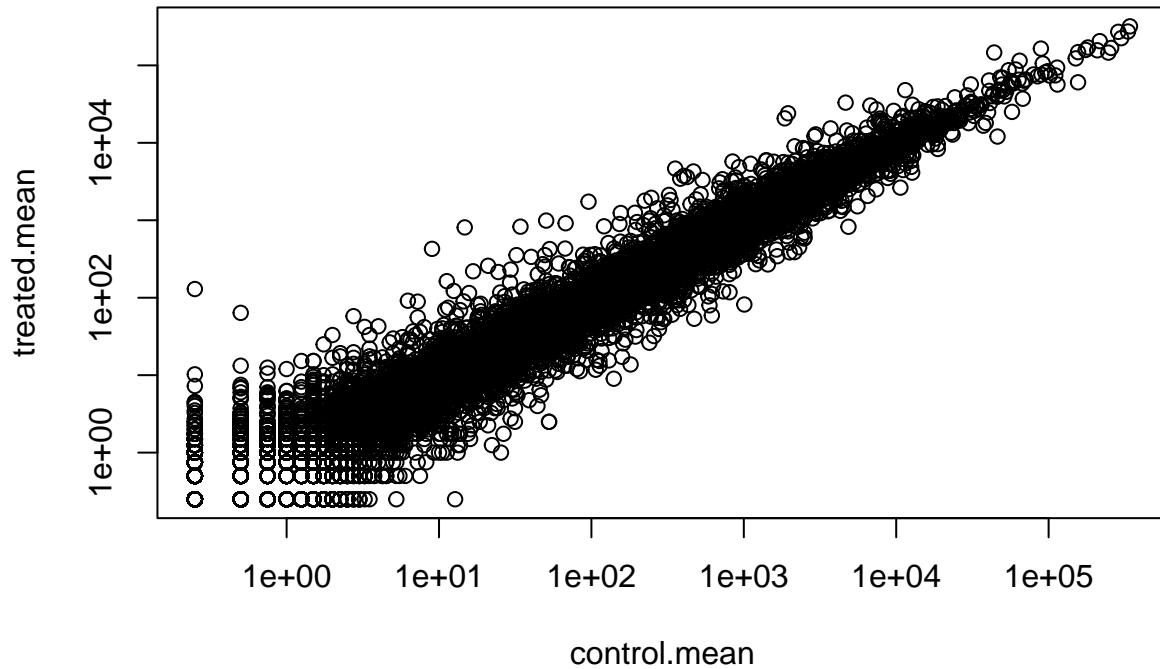
```
library(ggplot2)
ggplot(meancounts, aes(control.mean , treated.mean)) +geom_point()
```



Q.6 We're going to have to replot on a log-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
```



Lets add the **log2(fold-change)** to our ‘meancounts’ df

```
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
## 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 the genes (rows) with zero counts

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

```
to.rm <- unique(zero.vals[,1])
mycounts <- meancounts[-to.rm,]
head(mycounts)
```

```
##                  control.mean treated.mean      log2fc
## ENSG000000000003      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
nrow(mycounts)

## [1] 2181

```

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?

- which() when used in combination of arr.ind=TRUE tells us which positions are true elements. In this case, TRUE means the expression=0 and is something we need to exclude.

- Calling unique() will ensure we don't count any row twice if it has zero entries in both samples

How many genes are upregulated/downregulated upon drug treatment? We will use log2 threshold of +2 for this.

```

up.ind <- mycounts$log2fc > 2
down.ind <- mycounts$log2fc < (-2)

```

Q8.

```
sum(up.ind)
```

```
## [1] 250
```

Q9.

```
sum(down.ind)
```

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

Q10. I trust the counts to be correct, but I have no reason to believe that the results are significant

```
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, colAnyNs, 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, rowAnyNs, 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("pkename")'.

##
## Attaching package: 'Biobase'
## The following object is masked from 'package:MatrixGenerics':
## 
##     rowMedians

## The following objects are masked from 'package:matrixStats':
## 
##     anyMissing, rowMedians
citation("DESeq2")

##
## Love, M.I., Huber, W., Anders, S. Moderated estimation of fold change
## and dispersion for RNA-seq data with DESeq2 Genome Biology 15(12):550
## (2014)
##
```

```

## A BibTeX entry for LaTeX users is
##
## @Article{,
##   title = {Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2},
##   author = {Michael I. Love and Wolfgang Huber and Simon Anders},
##   year = {2014},
##   journal = {Genome Biology},
##   doi = {10.1186/s13059-014-0550-8},
##   volume = {15},
##   issue = {12},
##   pages = {550},
## }
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): ENSG00000000003 ENSG00000000005 ... ENSG00000283120
##   ENSG00000283123
## rowData names(0):
## colnames(8): SRR1039508 SRR1039509 ... SRR1039520 SRR1039521
## colData names(4): id dex celltype geo_id
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

```

Write out whole dataset (including genes that don't change)

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

Focus in on those genes with a small p-value

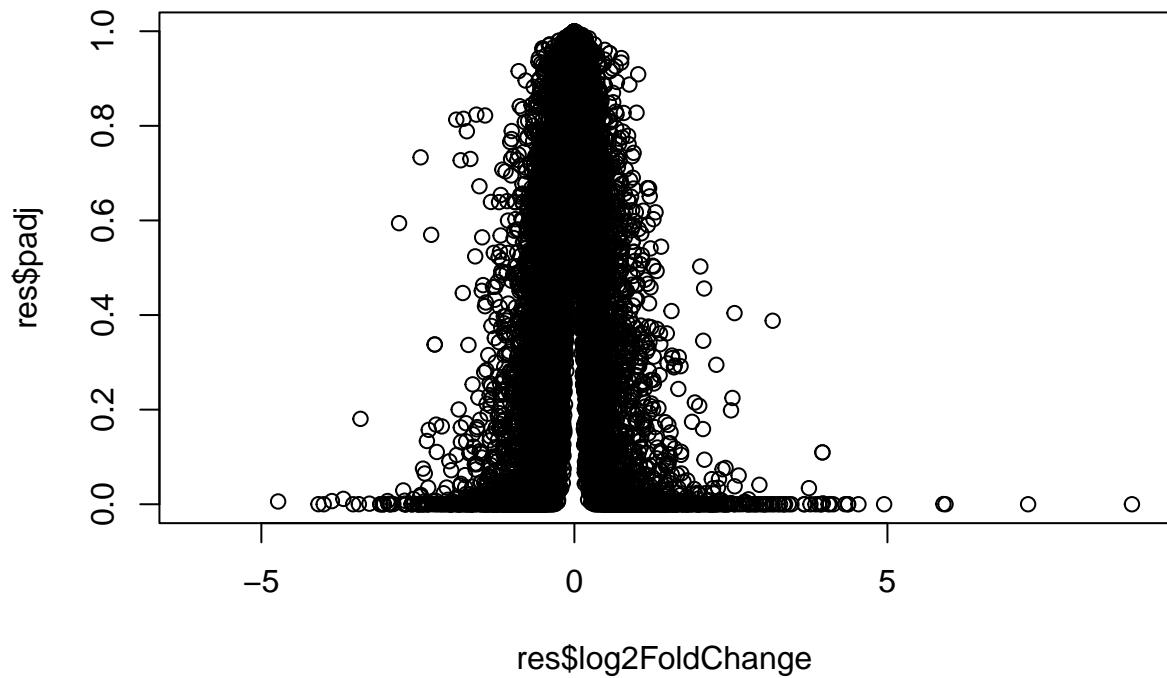
```
res05 <- results(dds, alpha=0.05)
```

```
summary(res05)
```

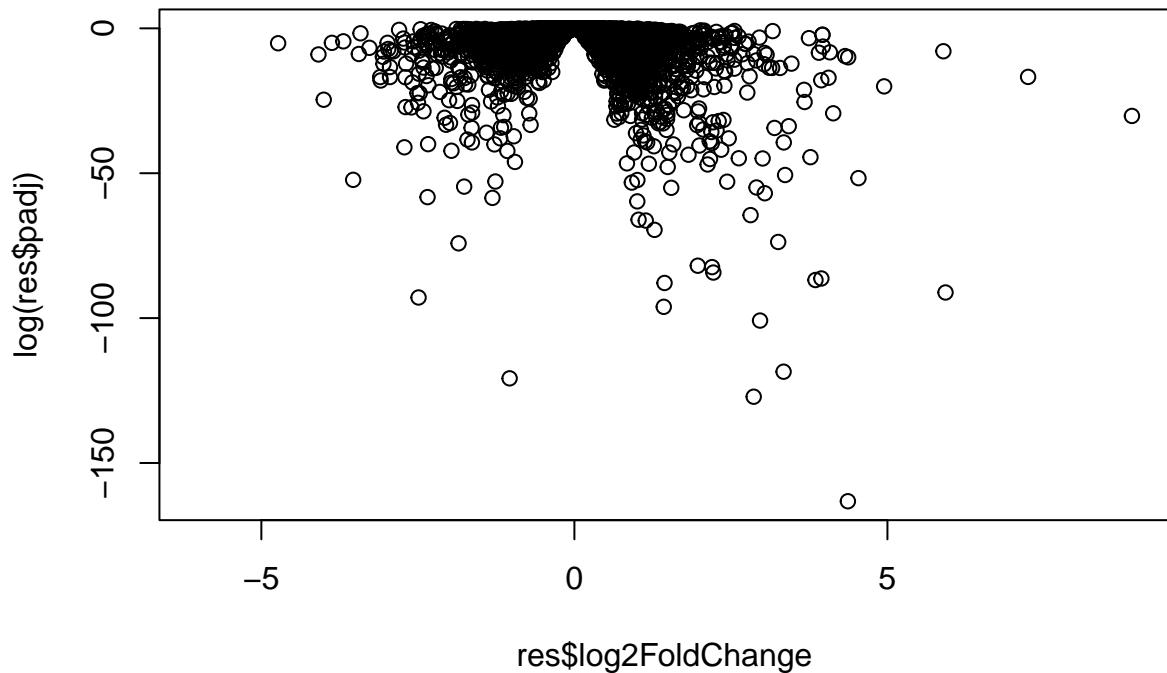
```

##
## out of 25258 with nonzero total read count
## adjusted p-value < 0.05
## LFC > 0 (up)      : 1236, 4.9%
## LFC < 0 (down)    : 933, 3.7%
## outliers [1]       : 142, 0.56%
## low counts [2]     : 9033, 36%
## (mean count < 6)
## [1] see 'cooksCutoff' argument of ?results
## [2] see 'independentFiltering' argument of ?results
plot(res$log2FoldChange, res$padj)

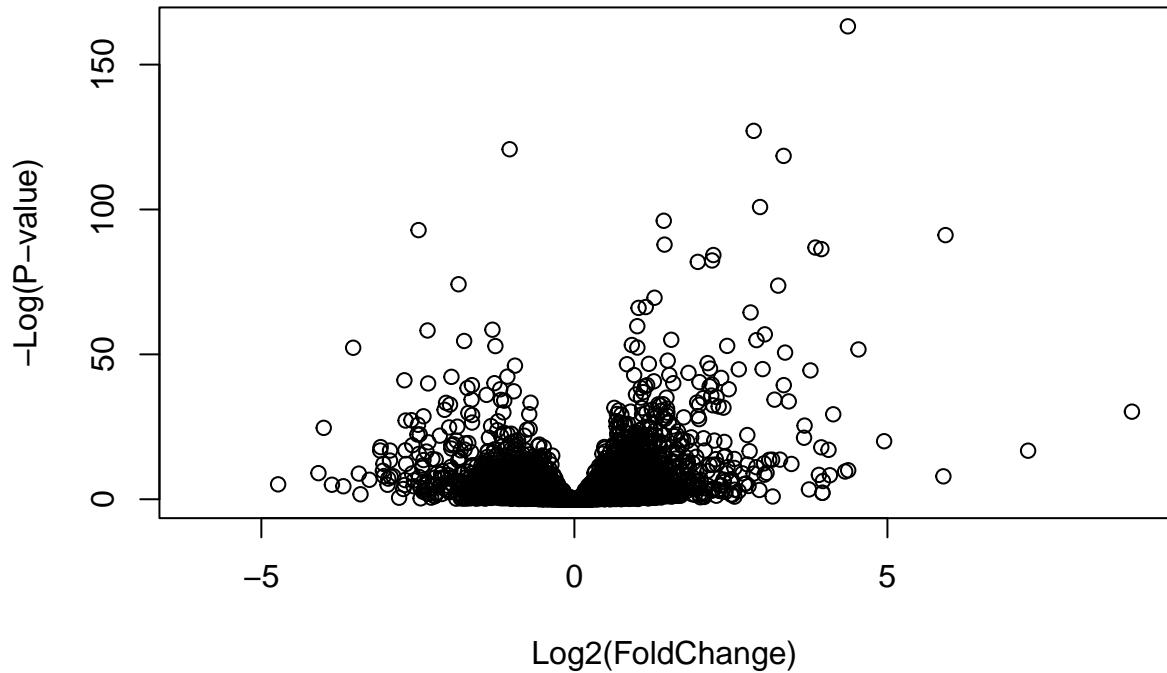
```



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



```
plot( res$log2FoldChange, -log(res$padj),  
      xlab="Log2(FoldChange)",  
      ylab="-Log(P-value)")
```



```

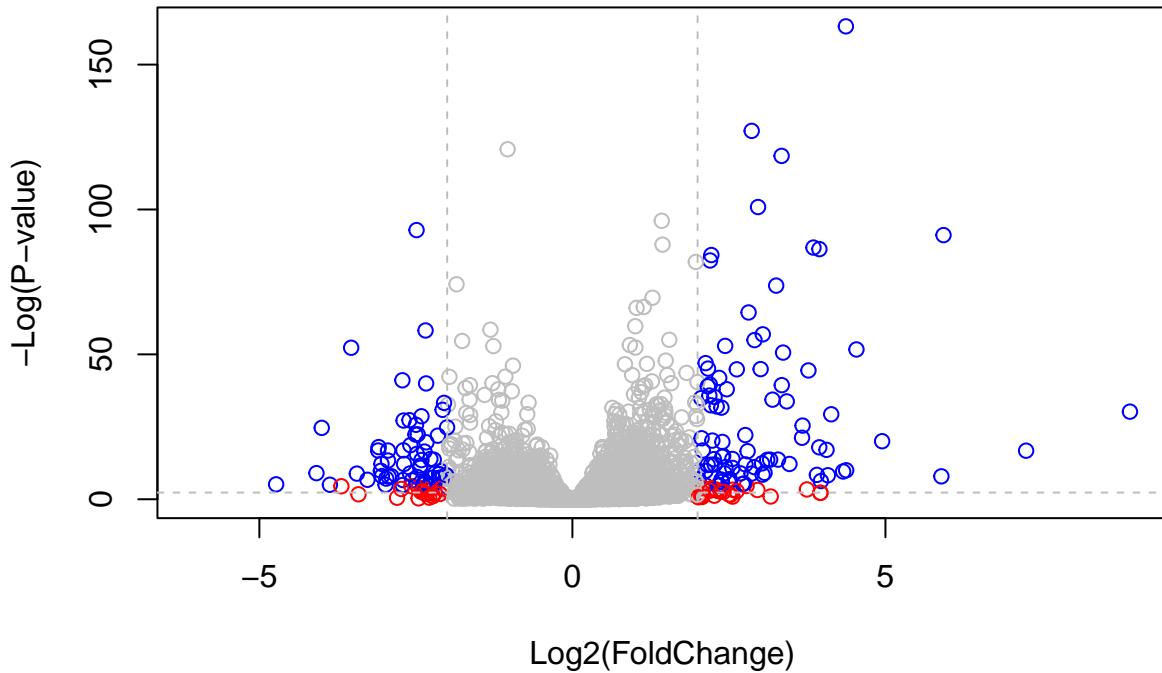
# Setup our custom point color vector
mycols <- rep("gray", nrow(res))
mycols[ abs(res$log2FoldChange) > 2 ] <- "red"

inds <- (res$padj < 0.01) & (abs(res$log2FoldChange) > 2 )
mycols[ inds ] <- "blue"

# Volcano plot with custom colors
plot( res$log2FoldChange, -log(res$padj),
      col=mycols, ylab="-Log(P-value)", xlab="Log2(FoldChange)" )

# Cut-off lines
abline(v=c(-2,2), col="gray", lty=2)
abline(h=-log(0.1), col="gray", lty=2)

```



adding annotation data for our genes

For this we need to biocond. packages

```

library("AnnotationDbi")
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"

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
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>
## ENSG000000000003 747.194195 -0.3507030  0.168246 -2.084470 0.0371175
## ENSG000000000005  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>
## ENSG000000000003  0.163035    TSPAN6
## ENSG000000000005  NA          TNMD
## ENSG00000000419   0.176032    DPM1
## ENSG00000000457   0.961694    SCYL3
## ENSG00000000460   0.815849    C1orf112
## ENSG00000000938   NA          FGR

#Enhanced volcano unable to install

# x <- as.data.frame(res)

# EnhancedVolcano(x,
#   lab = x$symbol,
#   x = 'log2FoldChange',
#   y = 'pvalue')

res$entrez <- 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
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"

```

The main `gage()` function requires a names vector of fold changes, where the names of the values are Entrez genes IDs

```

foldchanges <- res$log2FoldChange
names(foldchanges) <- res$entrez
head(foldchanges)

```

```
##      TSPAN6          TNMD         DPM1        SCYL3       C1orf112        FGR
## -0.35070302          NA  0.20610777  0.02452695 -0.14714205 -1.73228897
```

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

```
attributes(keggres)
```

```

## $names
## [1] "greater" "less"     "stats"

```

This separates outs the results by “greater” and “less” ie those that are up vs down regulated.

Look at the first 3 down pathways

```
head(keggres$less, 3)
```

	p.geomean	stat.mean	p.val	q.val
## hsa00232 Caffeine metabolism	NA	NaN	NA	NA
## hsa00983 Drug metabolism - other enzymes	NA	NaN	NA	NA
## hsa01100 Metabolic pathways	NA	NaN	NA	NA
	set.size	exp1		
## hsa00232 Caffeine metabolism	0	NA		
## hsa00983 Drug metabolism - other enzymes	0	NA		
## hsa01100 Metabolic pathways	0	NA		

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

```
## Warning: None of the genes or compounds mapped to the pathway!
```

```
## Argument gene.idtype or cpd.idtype may be wrong.
```

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

```
## Info: Working in directory /home/saranjan
```

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

