

class13

Rosa Chavez

2024-06-07

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

head(metadata)

##           id      dex celltype     geo_id
## 1 SRR1039508 control   N61311 GSM1275862
## 2 SRR1039509 treated    N61311 GSM1275863
## 3 SRR1039512 control   NO52611 GSM1275866
## 4 SRR1039513 treated    NO52611 GSM1275867
## 5 SRR1039516 control   NO80611 GSM1275870
## 6 SRR1039517 treated    NO80611 GSM1275871

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

## [1] TRUE
```

Q1. How many genes are in this dataset?

There are 38694 genes in this dataset

Q2. How many ‘control’ cell lines do we have?

There are 4 control cell lines in the dataset

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

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

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

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

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

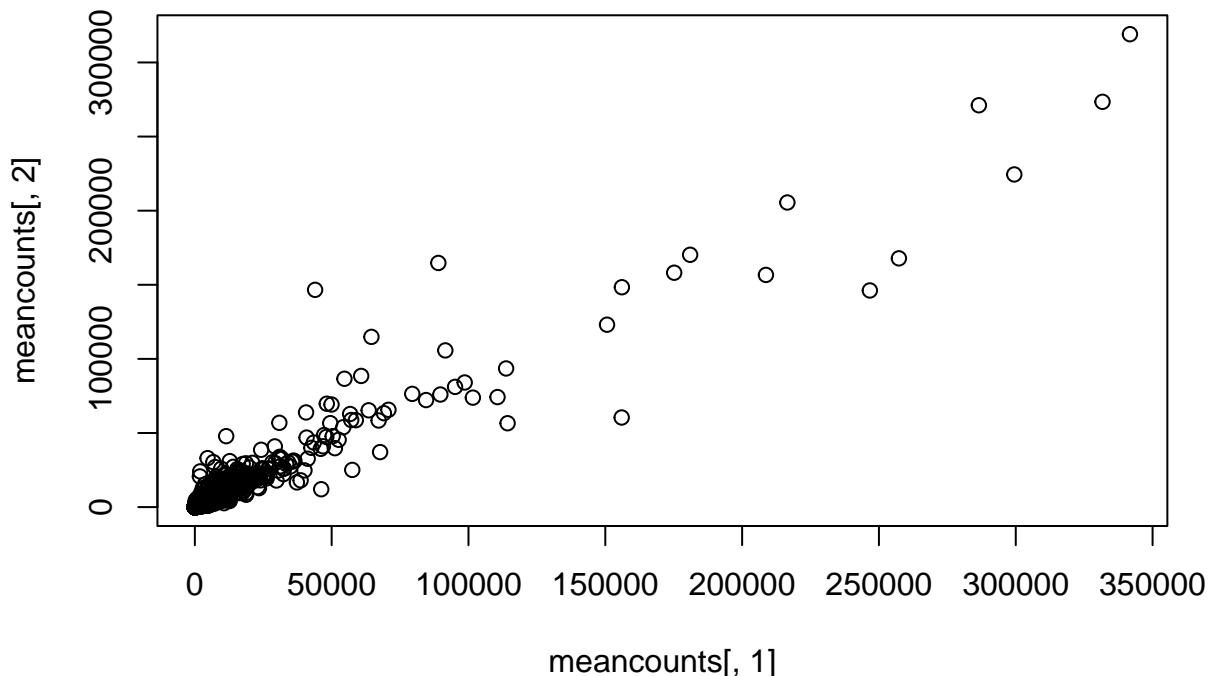
```

## ENSG00000000938
##          0.00
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[,1], meancounts[,2])
```



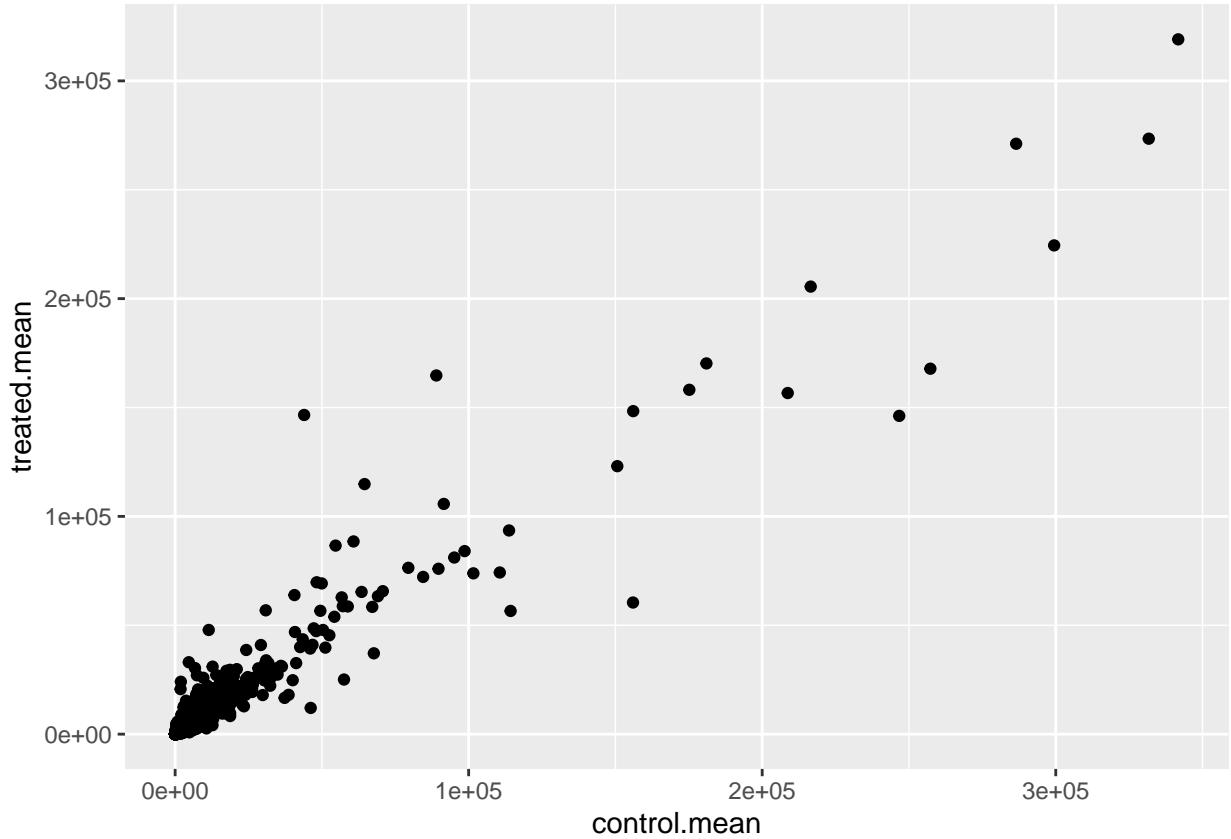
Q5 (b). You could also use the ggplot2 package to make this figure producing the plot below. What geom_?() function would you use for this plot?

```

library(ggplot2)

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

```



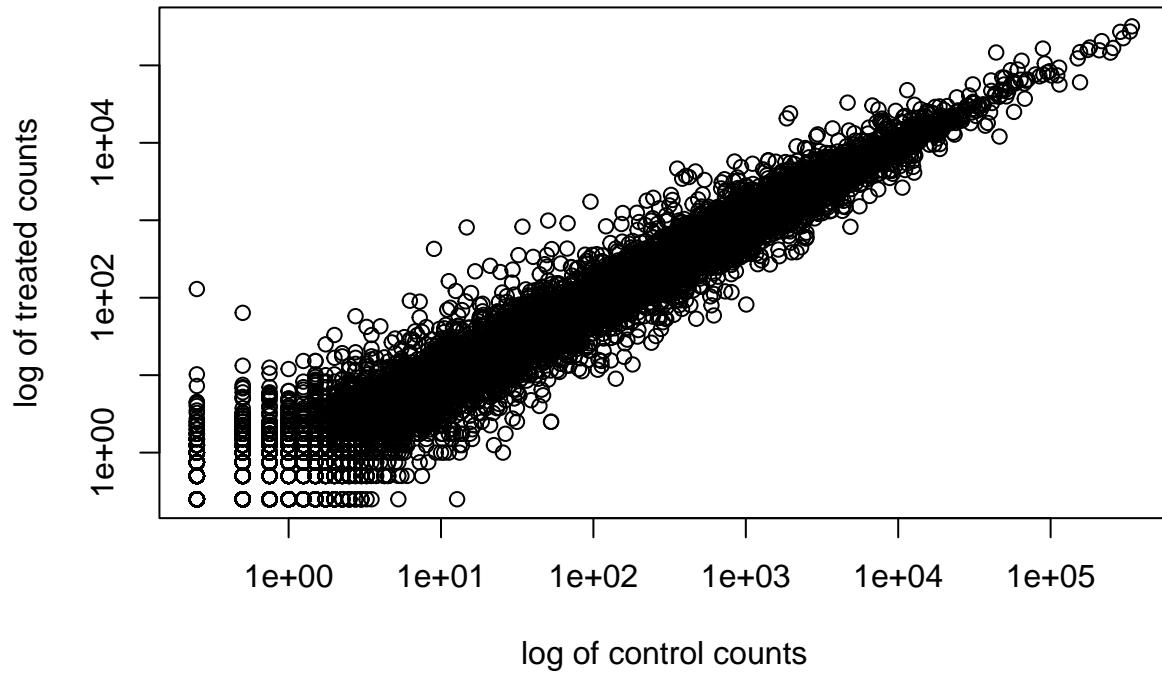
Q6. Try plotting both axes on a log scale. What is the argument to plot() that allows you to do this?

We will make a log-log plot to draw this skewed data and see what is going on.

```
plot(meancounts[,1], meancounts[,2], log="xy",
      xlab= "log of control counts",
      ylab= "log of 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
```



We often used log2 transformations when dealing with this sort of data

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

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

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

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

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

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

This log2 transformation has this nice property where if there is no change the log2 value will be zero and if it double the log2 value will be 1 and if halved it will be -1.

```
meancounts$log2fc <- log2(meancounts$treated.mean /
  meancounts$control.mean)
```

```
head(meancounts)
```

	control.mean	treated.mean	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

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?

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

How many 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
```

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
```

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

Where are missing the statistics, which gives us a tunnel vision approach on the data. Further analysis must be done to truly trust these results.

```
# load up DESeq2
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, 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 object is masked from 'package:utils':
##
##     findMatches

## The following objects are masked from 'package:base':
##
##     expand.grid, I, uname

## 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, 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

```

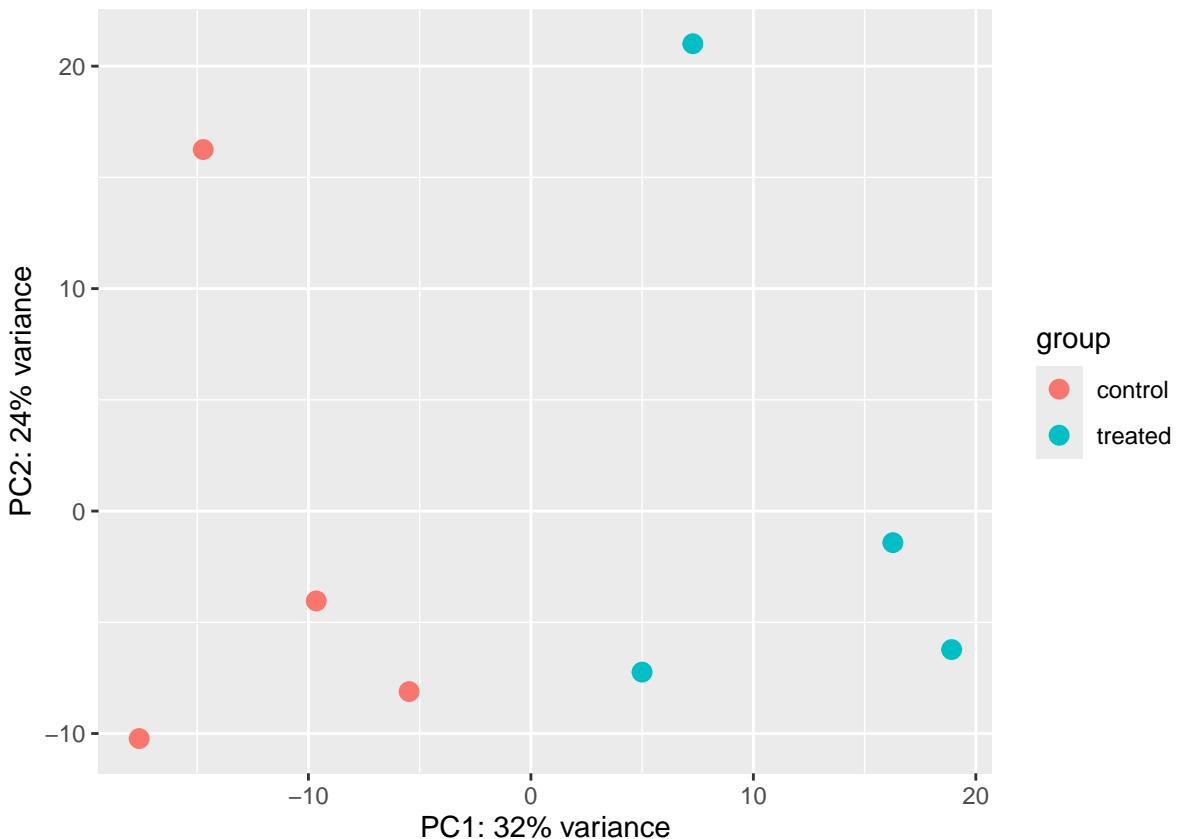
```

##      'citation("Biobase")', and for packages 'citation("pkgname")'.
##
## Attaching package: 'Biobase'
## The following object is masked from 'package:MatrixGenerics':
##
##     rowMedians
## The following objects are masked from 'package:matrixStats':
##
##     anyMissing, rowMedians
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
vsd <- vst(dds, blind = FALSE)
plotPCA(vsd, intgroup = c("dex"))

## using ntop=500 top features by variance

```



```
pcaData <- plotPCA(vsd, intgroup=c("dex"), returnData=TRUE)
```

```
## using ntop=500 top features by variance
```

```

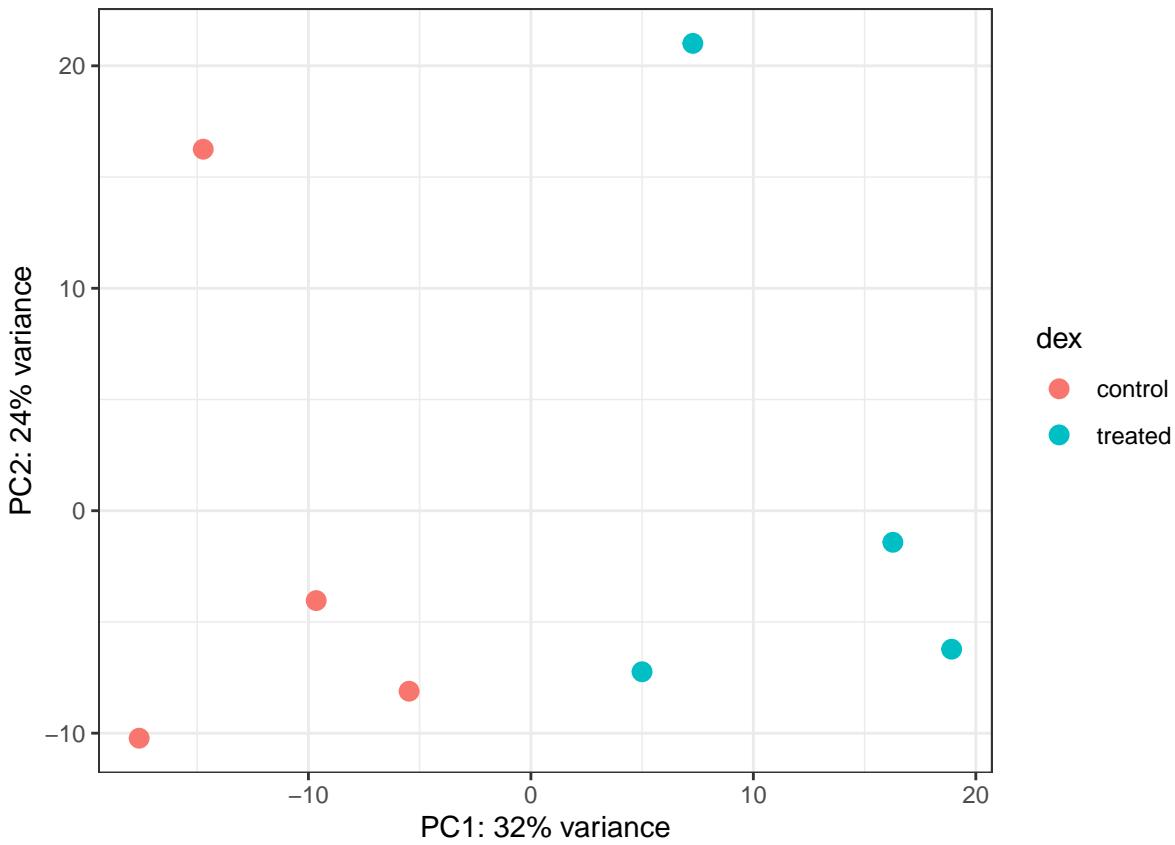
head(pcaData)

##          PC1      PC2   group    dex     name
## SRR1039508 -17.607922 -10.225252 control control SRR1039508
## SRR1039509  4.996738  -7.238117 treated treated SRR1039509
## SRR1039512 -5.474456 -8.113993 control control SRR1039512
## SRR1039513 18.912974 -6.226041 treated treated SRR1039513
## SRR1039516 -14.729173 16.252000 control control SRR1039516
## SRR1039517  7.279863 21.008034 treated treated SRR1039517

# Calculate percent variance per PC for the plot axis labels
percentVar <- round(100 * attr(pcaData, "percentVar"))

ggplot(pcaData) +
  aes(x = PC1, y = PC2, color = dex) +
  geom_point(size = 3) +
  xlab(paste0("PC1: ", percentVar[1], "% variance")) +
  ylab(paste0("PC2: ", percentVar[2], "% variance")) +
  coord_fixed() +
  theme_bw()

```



```

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>
## ENSG000000000003  747.1942    -0.3507030  0.168246 -2.084470 0.0371175
## ENSG000000000005   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>
## ENSG000000000003   0.163035
## ENSG000000000005   NA
## ENSG00000000419   0.176032
## ENSG00000000457   0.961694
## ENSG00000000460   0.815849
## ...
##           ...
## ENSG00000283115   NA
## ENSG00000283116   NA
## ENSG00000283119   NA
## ENSG00000283120   NA
## ENSG00000283123   NA

```

We can get some basic summary tallies using `summary()` function

```

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

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

```

```

##
```

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

```
## [1] "ACCCNUM"      "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"
```

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.

```
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")
```

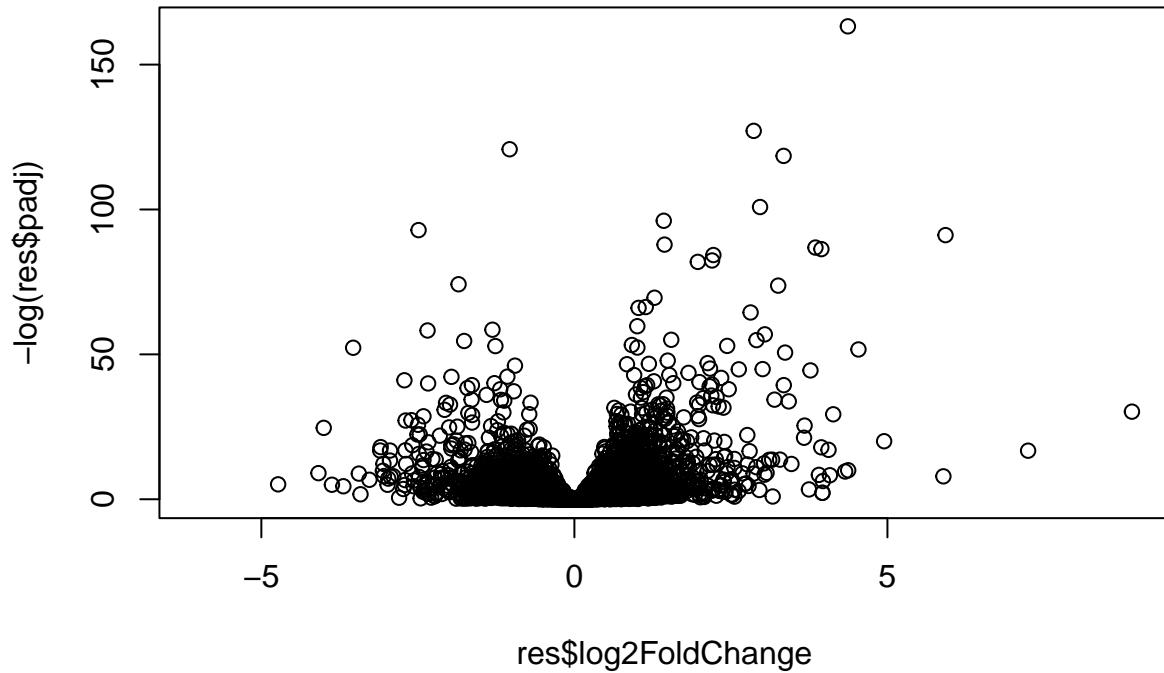
```
## 'select()' returned 1:many mapping between keys and columns
```

```
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    FIRRM
## ENSG00000000938   NA          FGR
```

Make a summary plot of our results.

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



```

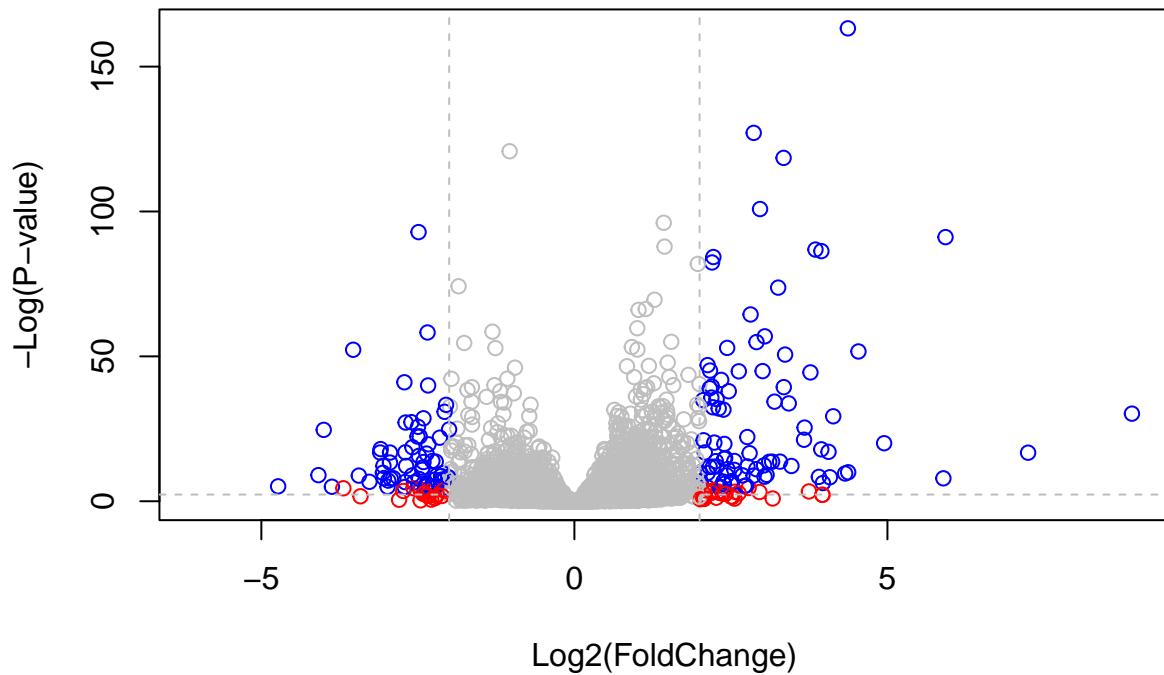
# Setup our custom point color vector
mycols <- rep("gray", nrow(res))
mycols[ abs(res$log2FoldChange) > 2 ] <- "red"

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)

```



```

log(0.1)

## [1] -2.302585

log(0.005)

## [1] -5.298317

write.csv(res, file="DESeq2_results.csv")

library(EnhancedVolcano)

## Loading required package: ggrepel

x <- as.data.frame(res)

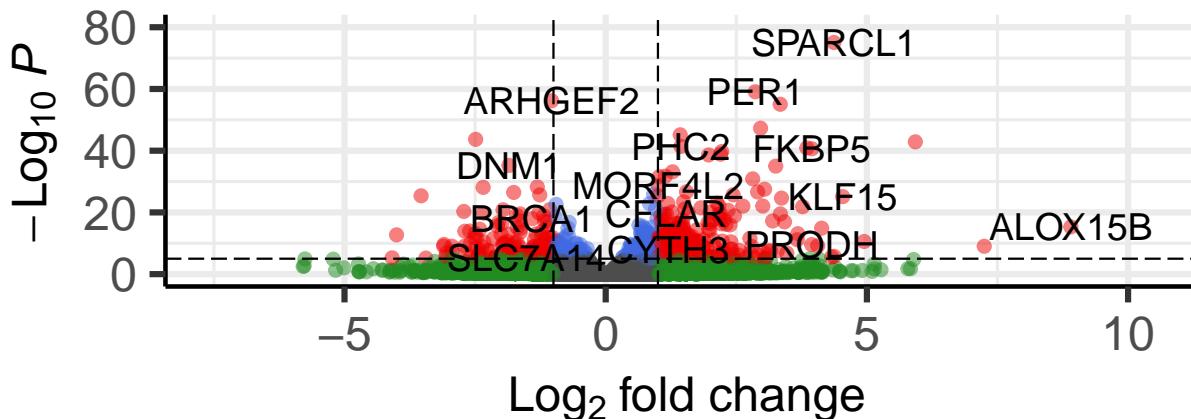
EnhancedVolcano(x,
  lab = x$symbol,
  x = 'log2FoldChange',
  y = 'pvalue')

```

Volcano plot

EnhancedVolcano

● NS ● Log₂ FC ● p-value ● p-value and log₂ FC



```
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"

foldchanges = res$log2FoldChange
names(foldchanges) = res$entrez
head(foldchanges)

## [1] -0.35070302          NA  0.20610777  0.02452695 -0.14714205 -1.73228897

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

## $names
## [1] "greater" "less"     "stats"
head(keggres$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 C:/Users/ITSloaner/Desktop-I89K3M9/OneDrive - UC San Diego/BIMM 143/class
## Info: Writing image file hsa05310.pathview.png
pathview(gene.data=foldchanges, pathway.id="hsa05310", kegg.native=FALSE)

## 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 C:/Users/ITSloaner/Desktop-I89K3M9/OneDrive - UC San Diego/BIMM 143/class
## Info: Writing image file hsa05310.pathview.pdf

```