

Class11 Lab

Morgan Farrell

2/23/2022

```
#RNAseq analysis from Asthma study Himes et. al.
```

```
#install.packages("BiocManager")
BiocManager::install("DESeq2")
library(BiocManager)
library(DESeq2)
```

```
#Import data set, countData and colData
```

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

Look at the data

```
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?

```
nrow(counts)
```

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

There are 38694 genes

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

```
metadata$dex
```

```
## [1] "control" "treated" "control" "treated" "control" "treated" "control"  
## [8] "treated"
```

There are 4 control cell lines

Make sure the control id from the count and metadata match

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

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

```
all(c(T, F, T)) #how the all function works
```

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

Test the ‘all()’ function on the conditional argument we wrote

```
if(all(metadata$id == colnames(counts))){  
  cat("yep")  
}
```

```
## yep
```

Compare control vs treated cell lines

```
control <- metadata[metadata$dex == "control",] #subset control data from the metadata set  
head(control)
```

```

##           id      dex celltype     geo_id
## 1 SRR1039508 control   N61311 GSM1275862
## 3 SRR1039512 control   N052611 GSM1275866
## 5 SRR1039516 control   N080611 GSM1275870
## 7 SRR1039520 control   N061011 GSM1275874

control.counts <- counts[ ,control$id] #extract control ids from the count data
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

Q3. How would you make the above code in either approach more robust? Because we are sub-setting the metadata and the count data by row labels it works for larger datasets too. If we were only working with column locations that may change with different datasets. Also we put the subsetted information into a vector to store all the data we pulled out, opposed to just printing out the entire dataset.

Get treated experiments too, using a similar method to control

```

treated <- metadata[metadata$dex == "treated",]
treated.counts <- counts[ ,treated$id]
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

Find the means of treated and control experiments

```

control.mean <- rowMeans( control.counts )
head(control.mean)

```

```
## ENSG00000000003 ENSG00000000005 ENSG000000000419 ENSG000000000457 ENSG000000000460
##          900.75           0.00        520.50       339.75         97.25
## ENSG000000000938
##          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.mean <- rowMeans( treated.counts )
head(treated.mean)

## ENSG00000000003 ENSG00000000005 ENSG000000000419 ENSG000000000457 ENSG000000000460
##          658.00           0.00        546.00       316.50         78.75
## ENSG000000000938
##          0.00
```

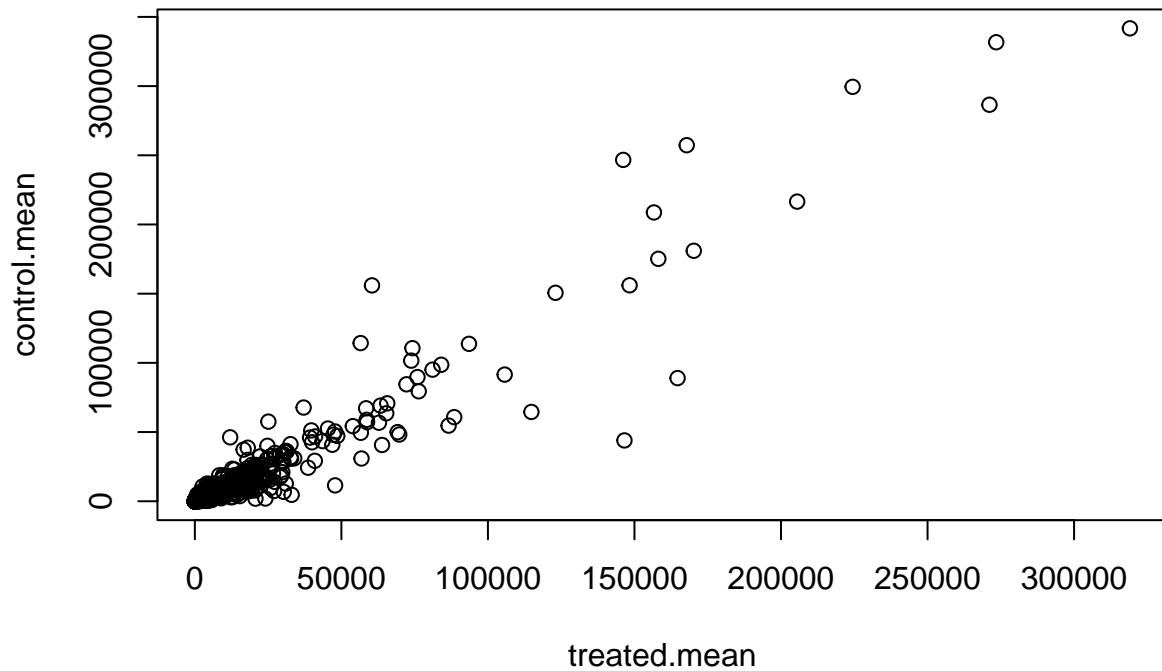
Combine the two datasetts for plotting

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

##           control.mean   treated.mean
## ENSG00000000003      900.75      658.00
## ENSG00000000005       0.00       0.00
## ENSG000000000419     520.50      546.00
## ENSG000000000457     339.75      316.50
## ENSG000000000460      97.25       78.75
## ENSG000000000938      0.75       0.00
```

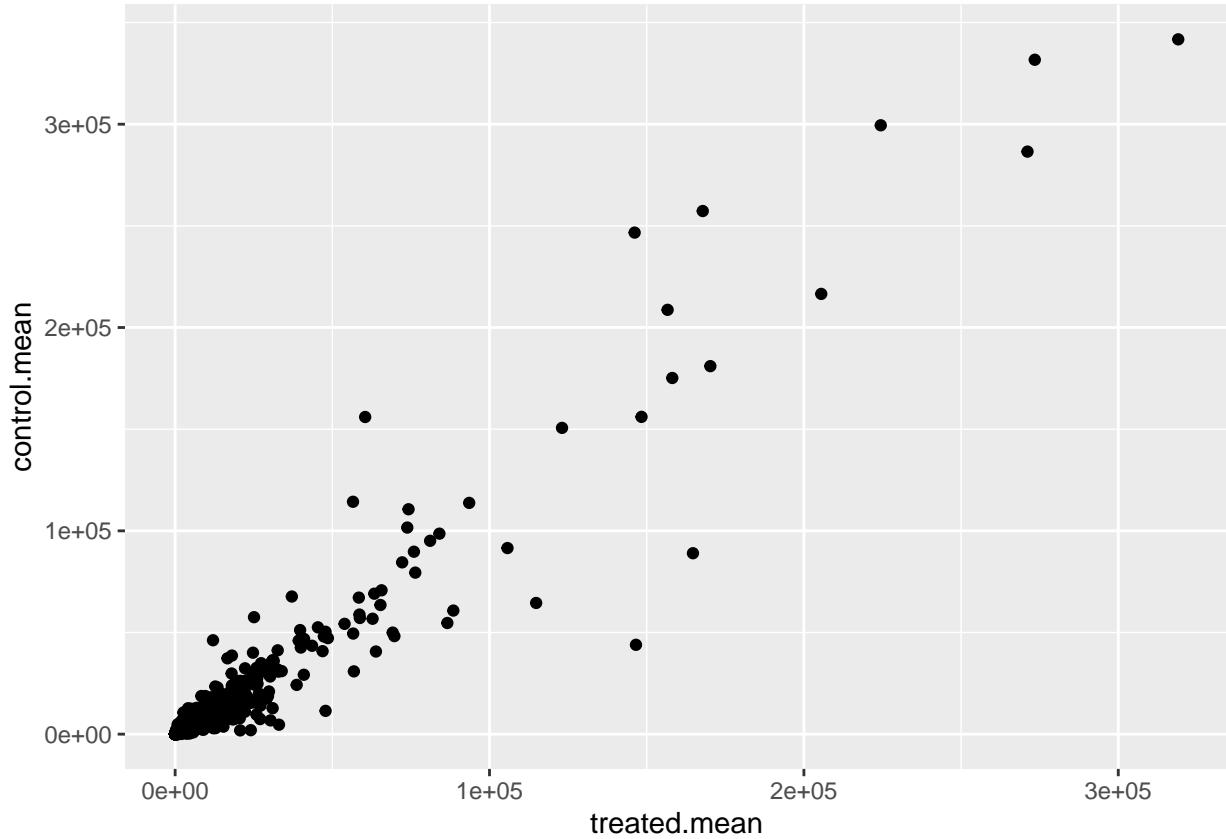
Q5 (a). Create a scatter plot showing the mean of the treated samples against the mean of the control samples. Your plot should look something like the following.

```
plot(treated.mean, control.mean)
```



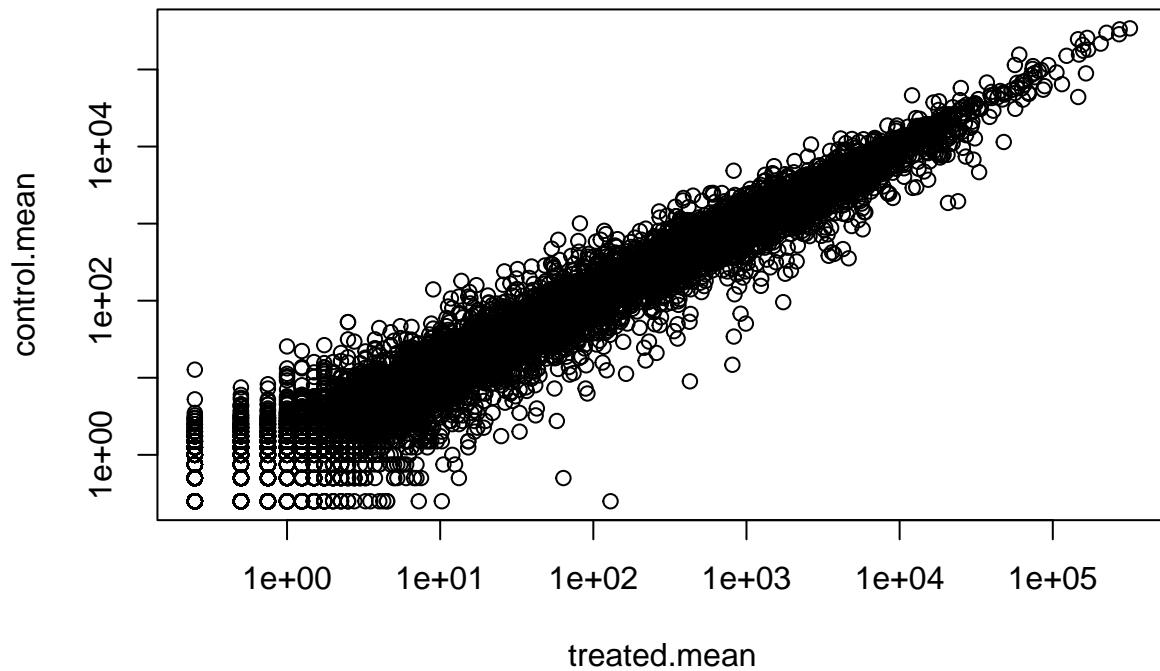
Q5 (b). You could also use the ggplot2 package to make this figure producing the plot below. What geom_?() function would you use for this plot?

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



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

```
plot(treated.mean, control.mean, log="xy")  
  
## Warning in xy.coords(x, y, xlabel, ylabel, log): 15281 x values <= 0 omitted  
## from logarithmic plot  
  
## Warning in xy.coords(x, y, xlabel, ylabel, log): 15032 y values <= 0 omitted  
## from logarithmic plot
```



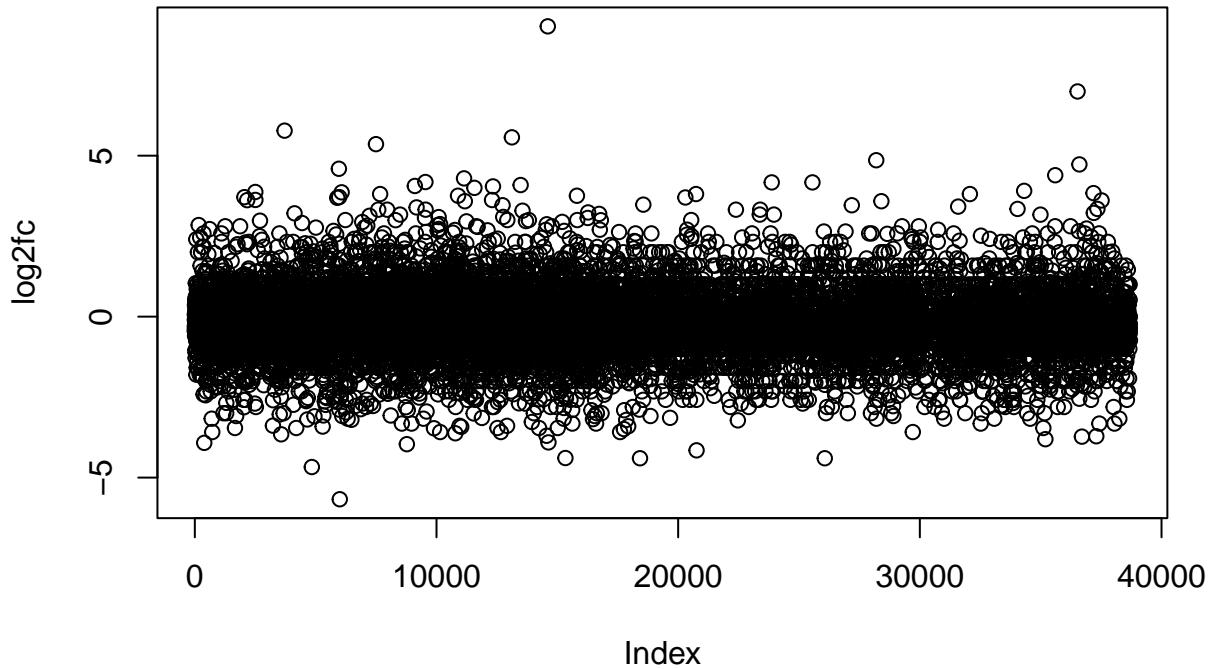
We often use log 2 transformation, because it indicates fold change

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

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

log2 fold change on data

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



Put all our data in a data.frame

```
meancounts <- data.frame(control.mean, treated.mean, log2fc)
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?

arr.ind is logical; should array indices be returned when x is an array? So it sorts the matching info into an array with the row and column locations

Get rid of 0s and infinities

```
zero.vals <- which(meancounts[,1:2]== "0", arr.ind = TRUE) #use which function to find the 0s and return  
head(zero.vals)
```

```
##           row col  
## ENSG000000000005   2   1  
## ENSG00000004848  65   1  
## ENSG00000004948  70   1  
## ENSG00000005001  73   1  
## ENSG00000006059 121   1  
## ENSG00000006071 123   1
```

Find the unique rows with these 0s

```
to.rm <- unique(which(meancounts[,1:2]== "0", arr.ind = TRUE)[,"row"]) #genes to exclude sorted by the 0s  
mycounts <- meancounts[-to.rm,] #remove the vector with the 0s  
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
```

Genes left after removing the ones with 0s

```
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?

Not entirely. These differences in the mean may not be statistically significant based on many factors. The control could be skewed by low sample size for example. We should probably do more summary statistics for both groups before determining significance.

Use DESeq2 for these summary statistics on this dataset

```
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, 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

```

```
#citation("DESeq2")
```

Create a DESeq data matrix with our data

```
dds <- DESeqDataSetFromMatrix(countData=counts,
                               colData=metadata,
                               design=~dex) #how are we comparing the data look at the dex column for th

## 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

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
```

```

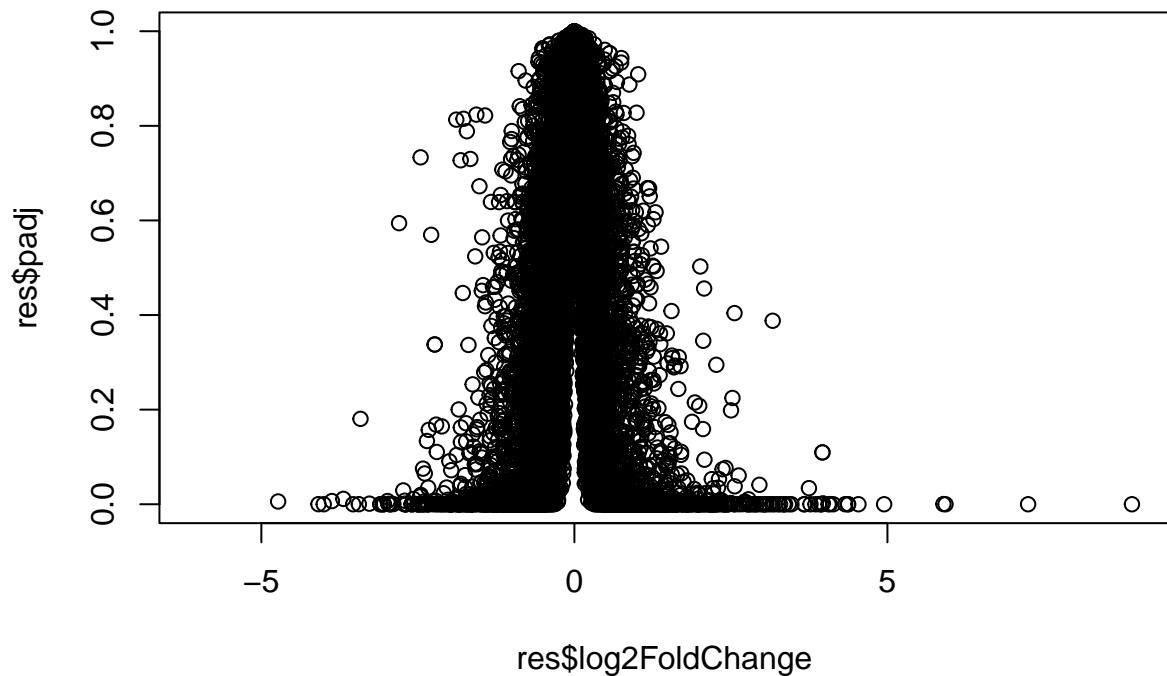
## ENSG00000000419 520.1342      0.2061078  0.101059  2.039475 0.0414026
## ENSG00000000457 322.6648      0.0245269  0.145145  0.168982 0.8658106
## ENSG00000000460 87.6826      -0.1471420  0.257007 -0.572521 0.5669691
## ...
## ...
## ENSG00000283115 0.000000      NA        NA        NA        NA
## 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>
## ENSG0000000003 0.163035
## ENSG0000000005 NA
## ENSG00000000419 0.176032
## ENSG00000000457 0.961694
## ENSG00000000460 0.815849
## ...
## ...
## ENSG00000283115 NA
## ENSG00000283116 NA
## ENSG00000283119 NA
## ENSG00000283120 NA
## ENSG00000283123 NA

```

Main results figure as a volcano plot

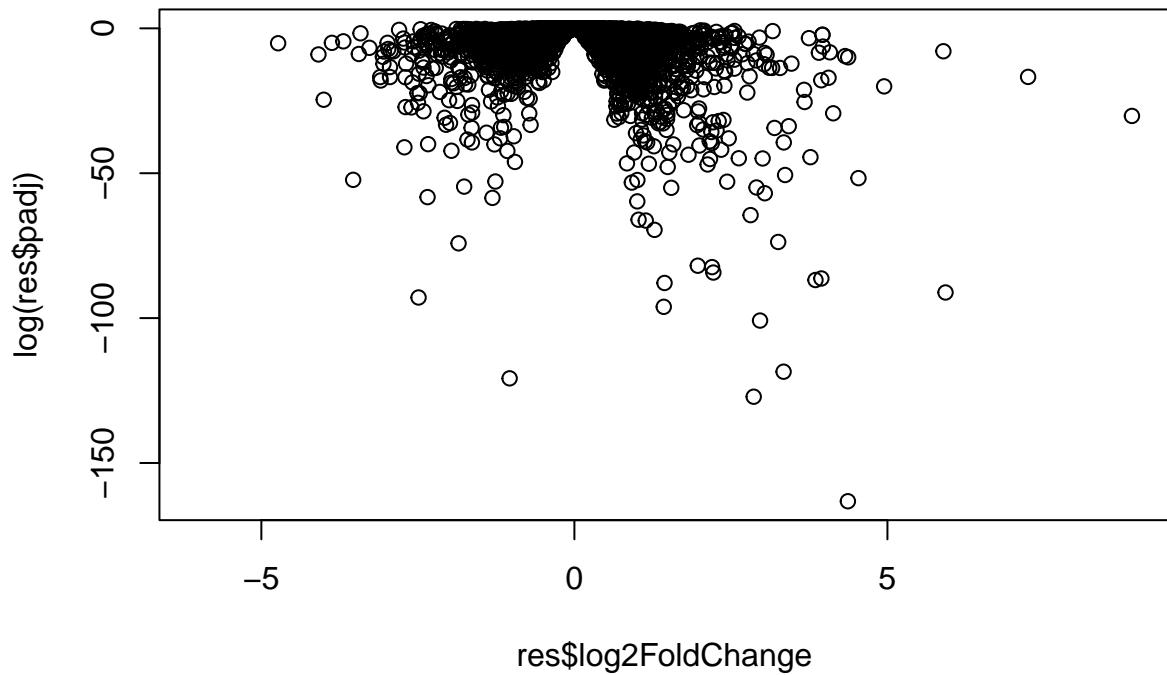
Log2 fold change vs p- value

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



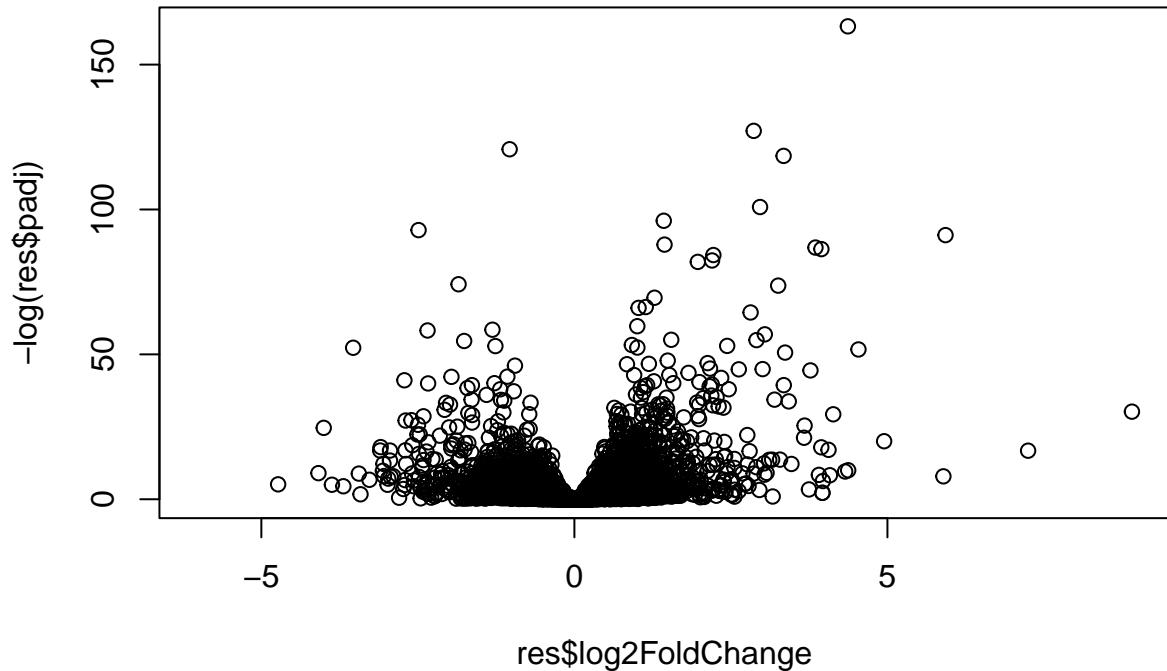
Log transformation of this plot

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



Flip the plot by putting a - in front of our log padj

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

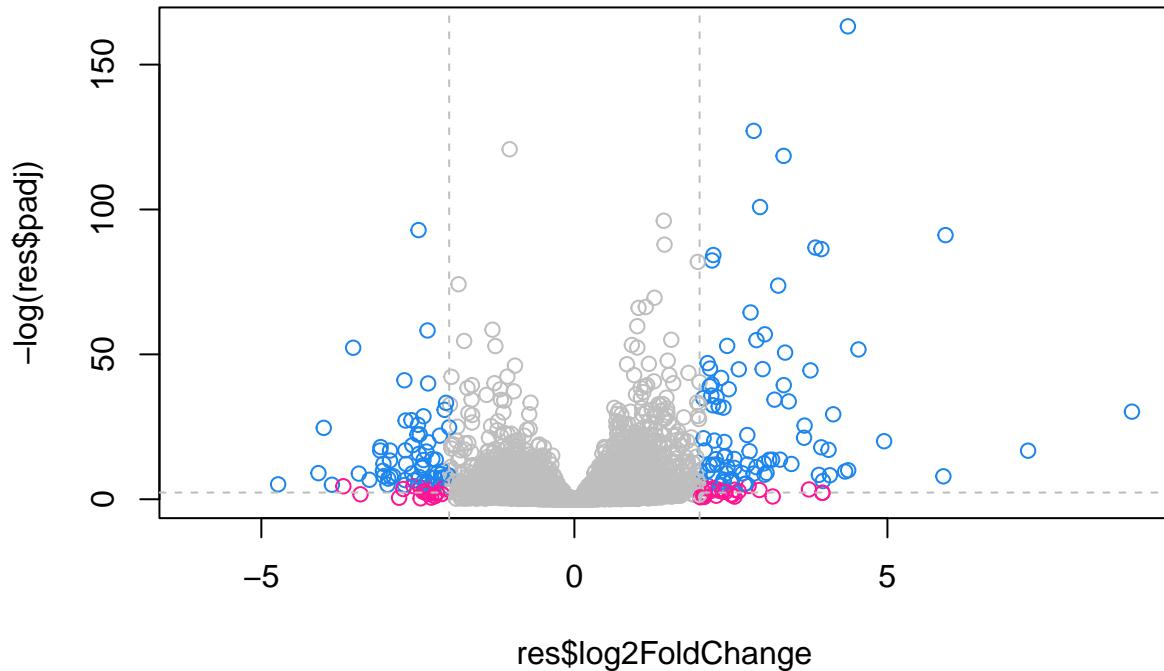


```

mycols <- rep("gray", nrow(res))
mycols[ abs(res$log2FoldChange) > 2 ] <- "deeppink"
inds <- (res$padj < 0.01) & (abs(res$log2FoldChange) > 2 )
mycols[ inds ] <- "dodgerblue2"
plot(res$log2FoldChange, -log(res$padj), col= mycols)

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

```



Add annotation data

Install all the proper packages

```
#BiocManager::install("EnhancedVolcano")
library(EnhancedVolcano)

## Loading required package: ggrepel

## Registered S3 methods overwritten by 'ggalt':
##   method           from
##   grid.draw.absoluteGrob  ggplot2
##   grobHeight.absoluteGrob ggplot2
##   grobWidth.absoluteGrob ggplot2
##   grobX.absoluteGrob     ggplot2
##   grobY.absoluteGrob     ggplot2

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

##
```

Now use 'mapIDs()' function from the AnnotationDbi package to find

```
#First what is available
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
```

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.

Add in all of these new columns based on the mapID function

```

res$genename <- mapIds(org.Hs.eg.db,
                      keys=row.names(res),           # Our genenames
                      keytype="ENSEMBL",            # The format of our genenames
                      column="GENENAME",           # 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 8 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      genename
##          <numeric> <character> <character>
## ENSG00000000003  0.163035    TSPAN6      tetraspanin 6
## ENSG00000000005   NA         TNMD       tenomodulin
## ENSG00000000419  0.176032    DPM1       dolichyl-phosphate m..
## ENSG00000000457  0.961694    SCYL3      SCY1 like pseudokina..
## ENSG00000000460  0.815849    C1orf112    chromosome 1 open re..
## ENSG00000000938   NA         FGR        FGR proto-oncogene, ..

```

```

res$entrez <- mapIds(org.Hs.eg.db,
                      keys=row.names(res),           # Our genenames
                      keytype="ENSEMBL",            # The format of our genenames
                      column="ENTREZID",           # 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 9 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      genename      entrez
##          <numeric> <character>      <character> <character>
## ENSG000000000003 0.163035    TSPAN6      tetraspanin 6      7105
## ENSG000000000005 NA        TNMD      tenomodulin      64102
## ENSG000000000419 0.176032    DPM1      dolichyl-phosphate m..      8813
## ENSG000000000457 0.961694    SCYL3      SCY1 like pseudokina..      57147
## ENSG000000000460 0.815849    C1orf112    chromosome 1 open re..      55732
## ENSG000000000938 NA        FGR       FGR proto-oncogene, ..      2268

res$uniprot <- mapIds(org.Hs.eg.db,
                      keys=row.names(res),           # Our genenames
                      keytype="ENSEMBL",             # The format of our genenames
                      column="UNIPROT",              # 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 10 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      symbol      genename      entrez
##          <numeric> <character>      <character> <character>
## ENSG000000000003 0.163035    TSPAN6      tetraspanin 6      7105
## ENSG000000000005 NA        TNMD      tenomodulin      64102
## ENSG000000000419 0.176032    DPM1      dolichyl-phosphate m..      8813
## ENSG000000000457 0.961694    SCYL3      SCY1 like pseudokina..      57147
## ENSG000000000460 0.815849    C1orf112    chromosome 1 open re..      55732
## ENSG000000000938 NA        FGR       FGR proto-oncogene, ..      2268
##          uniprot
##          <character>
## ENSG000000000003 AOA024RCIO
## ENSG000000000005 Q9H2S6
## ENSG000000000419 O60762
## ENSG000000000457 Q8IZE3
## ENSG000000000460 AOA024R922
## ENSG000000000938 P09769

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)

##          7105      64102      8813      57147      55732      2268
## -0.35070302           NA  0.20610777  0.02452695 -0.14714205 -1.73228897

```

Now we can run gene set enrichment (aka pathway analysis) We can find out what keggres is using the ‘attributes()’ function

```

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

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

```

First look at the downregulated list (less)

```
# Look at the first three down (less) pathways
head(keggres$less, 3)
```

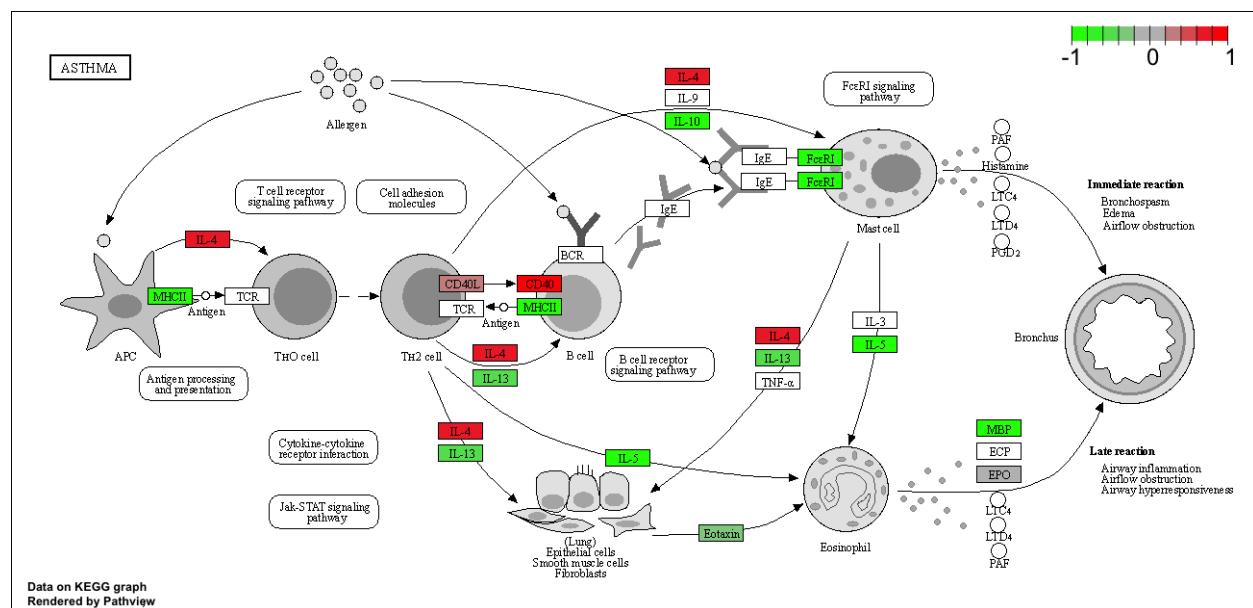
```
##                                     p.geomean stat.mean      p.val
## hsa05332 Graft-versus-host disease 0.0004250461 -3.473346 0.0004250461
## hsa04940 Type I diabetes mellitus 0.0017820293 -3.002352 0.0017820293
## hsa05310 Asthma                  0.0020045888 -3.009050 0.0020045888
##                                     q.val set.size      exp1
## hsa05332 Graft-versus-host disease 0.09053483      40 0.0004250461
## hsa04940 Type I diabetes mellitus 0.14232581      42 0.0017820293
## hsa05310 Asthma                  0.14232581      29 0.0020045888
```

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

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

```
## Info: Working in directory /Users/morganfarrell/Documents/PhD Research/UCSD Courses/BGGN213- Bioinfo
```

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



```
# A different PDF based output of the same data
```

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

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

```
## Info: Working in directory /Users/morganfarrell/Documents/PhD Research/UCSD Courses/BGGN213- Bioinfo
```

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

Save the results to a CSV

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