

Class 14: DESeq2 Analysis Mini-Proj

Nathaniel Lightle (A16669288)

Finishing up Class 13

OLD CODE

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

	SRR1039508	SRR1039509	SRR1039512	SRR1039513	SRR1039516
ENSG000000000003	723	486	904	445	1170
ENSG000000000005	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		
ENSG000000000003	1097	806	604		
ENSG000000000005	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

nrow(count)

[1] 38694

table(metadata$dex)

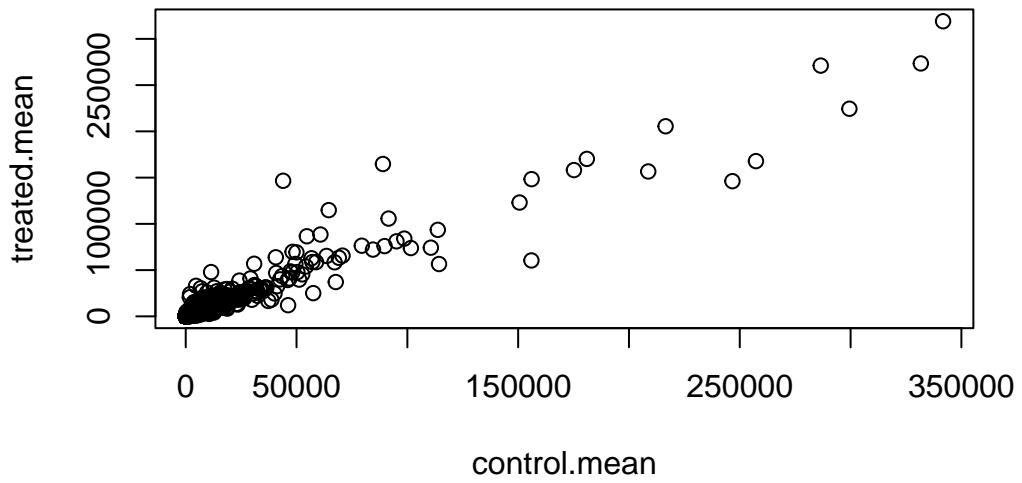
control treated
4      4

control.mean <- rowMeans(count[, metadata$dex == "control"])
treated.mean <- rowMeans(count[, metadata$dex == "treated"])
meancount <- data.frame(control.mean,treated.mean)
head(meancount)

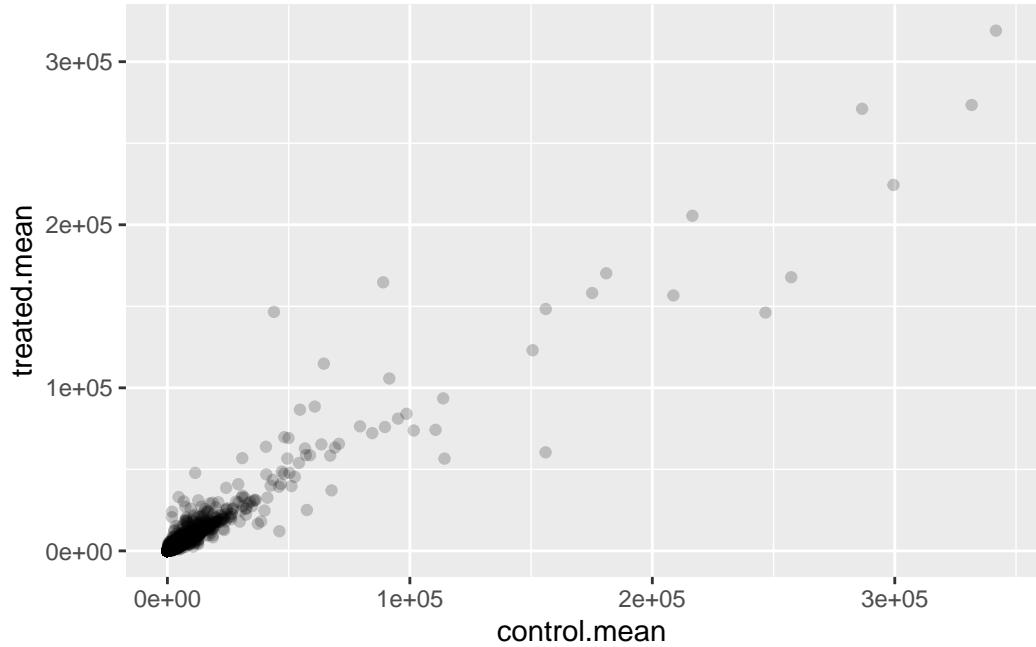
control.mean treated.mean
ENSG00000000003     900.75    658.00
ENSG00000000005      0.00      0.00
ENSG00000000419     520.50    546.00
ENSG00000000457     339.75    316.50
ENSG00000000460      97.25     78.75
ENSG00000000938      0.75      0.00

plot(meancount)
library(ggplot2)

```



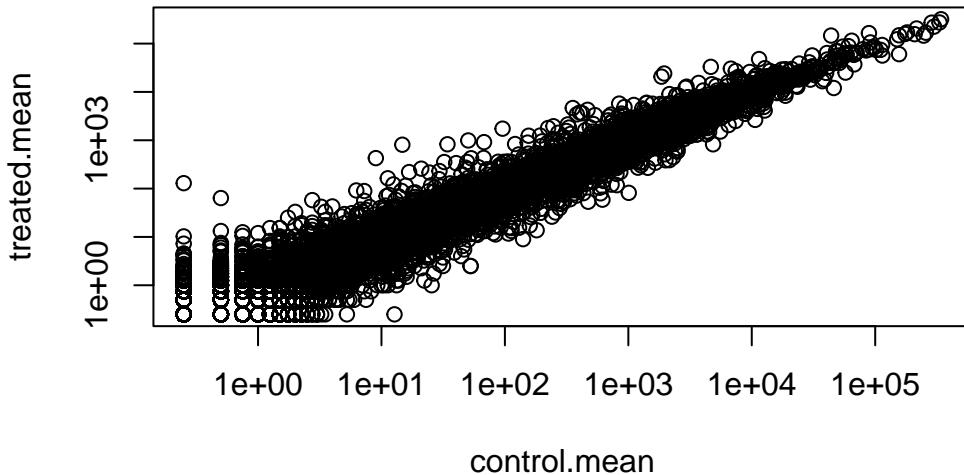
```
ggplot(meancount) +  
  aes(control.mean, treated.mean) +  
  geom_point(alpha=.2)
```



```
plot(meancount, 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



```
meancount$log2fc <- log2(meancount[,"treated.mean"]/meancount[,"control.mean"])
head(meancount)
```

	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

```
to.rm inds <- rowSums(meancount[,1:2] == 0) > 0
#meancount[!to.rm inds,]
mycount <- meancount[-to.rm inds,]
head(mycount)
```

	control.mean	treated.mean	log2fc
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
ENSG00000000971     5219.00    6687.50  0.35769358
```

```
dim(mycount)
```

```
[1] 38693      3
```

```
up.ind <- mycount$log2fc > 2
down.ind <- mycount$log2fc < (-2)
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, 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
```

```
dds <- DESeqDataSetFromMatrix(countData = count,  
                                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 <- DESeq(dds)
```

```
estimating size factors
```

```
estimating dispersions
```

```
gene-wise dispersion estimates
```

```
mean-dispersion relationship
```

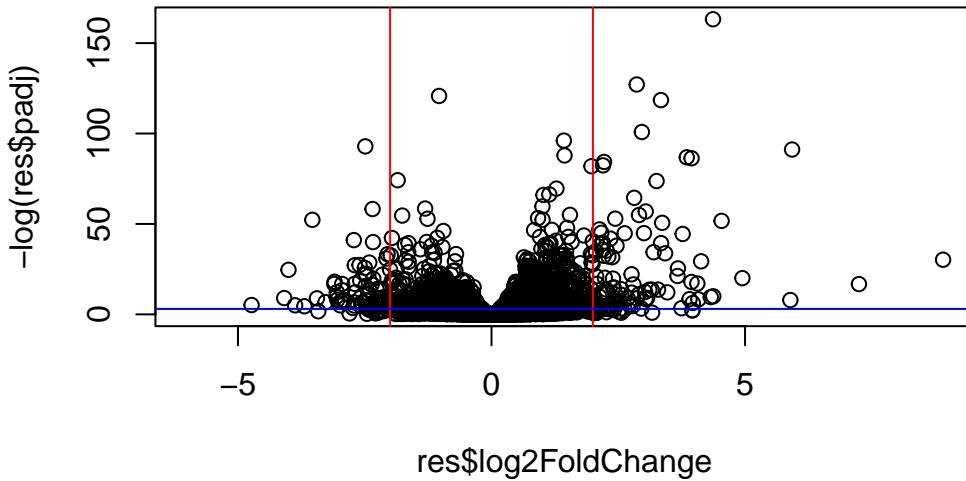
```
final dispersion estimates
```

```
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
```

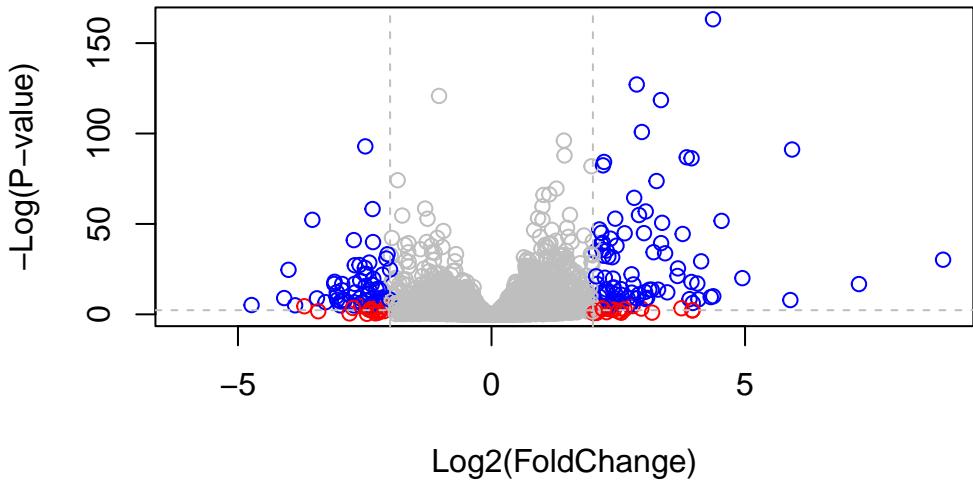
```
plot(res$log2FoldChange, -log(res$padj))
abline(v=2, col="red") #everything to the right up regulated
abline(v=-2,col="red") #everything to the left down regulated
abline(h=-log(.05), col="blue")
```



```

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)

```



NEW CODE

Section 8: Adding Annotation Data

Installing packages I need

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

Warning: package 'AnnotationDbi' was built under R version 4.3.2

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

Gettting a list of all available key types to map between

```
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

Q11.

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

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

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
```

	<code><numeric></code>	<code><numeric></code>	<code><numeric></code>	<code><numeric></code>	<code><numeric></code>
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
	<code>padj</code>	<code>symbol</code>	<code>entrez</code>	<code>uniprot</code>	
	<code><numeric></code>	<code><character></code>	<code><character></code>	<code><character></code>	
ENSG000000000003	0.163035	TSPAN6	7105	AOA024RC10	
ENSG000000000005	NA	TNMD	64102	Q9H2S6	
ENSG000000000419	0.176032	DPM1	8813	060762	
ENSG000000000457	0.961694	SCYL3	57147	Q8IZE3	
ENSG000000000460	0.815849	FIRRM	55732	AOA024R922	
ENSG000000000938	NA	FGR	2268	P09769	
		<code>genename</code>			
		<code><character></code>			
ENSG000000000003		tetraspanin 6			
ENSG000000000005		tenomodulin			
ENSG000000000419		dolichyl-phosphate m..			
ENSG000000000457		SCY1 like pseudokina..			
ENSG000000000460		FIGNL1 interacting r..			
ENSG000000000938		FGR proto-oncogene, ..			

Ordering results by adjusted p-value

```
ord <- order( res$padj )
#View(res[ord,])
head(res[ord,])
```

```
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>
ENSG00000152583   954.771     4.36836  0.2371268   18.4220 8.74490e-76
ENSG00000179094   743.253     2.86389  0.1755693   16.3120 8.10784e-60
ENSG00000116584  2277.913    -1.03470  0.0650984  -15.8944 6.92855e-57
ENSG00000189221  2383.754     3.34154  0.2124058   15.7319 9.14433e-56
ENSG00000120129  3440.704     2.96521  0.2036951   14.5571 5.26424e-48
ENSG00000148175 13493.920     1.42717  0.1003890   14.2164 7.25128e-46
  padj      symbol      entrez      uniprot
```

```

<numeric> <character> <character> <character>
ENSG00000152583 1.32441e-71 SPARCL1 8404 AOA024RDE1
ENSG00000179094 6.13966e-56 PER1 5187 015534
ENSG00000116584 3.49776e-53 ARHGEF2 9181 Q92974
ENSG00000189221 3.46227e-52 MAOA 4128 P21397
ENSG00000120129 1.59454e-44 DUSP1 1843 B4DU40
ENSG00000148175 1.83034e-42 STOM 2040 F8VSL7
genename
<character>
ENSG00000152583 SPARC like 1
ENSG00000179094 period circadian reg..
ENSG00000116584 Rho/Rac guanine nucl..
ENSG00000189221 monoamine oxidase A
ENSG00000120129 dual specificity pho..
ENSG00000148175 stomatin

```

Writing results to new CSV

```
write.csv(res[ord,], "deseq_results2.csv")
```

10: Pathway Analysis

```
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)
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

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

attributes(keggres)

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

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/Nate/Desktop/Fall Quarter 2023 /BIMM 143/Class 14: RNA seq
```

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

Q12. Can you do the same procedure as above to plot the pathview figures for the top 2 down-regulated pathways?

Yes, yes i can

```
pathview(gene.dat=foldchanges, pathway.id="hsa05332")
```

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

```
Info: Working in directory /Users/Nate/Desktop/Fall Quarter 2023 /BIMM 143/Class 14: RNA seq
```

```
Info: Writing image file hsa05332.pathview.png
```

```
pathview(gene.dat=foldchanges, pathway.id="hsa04940")
```

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

```
Info: Working in directory /Users/Nate/Desktop/Fall Quarter 2023 /BIMM 143/Class 14: RNA seq
```

```
Info: Writing image file hsa04940.pathview.png
```

OPTIONAL: Plotting counts for genes of interest

Finding GeneID for CRISPLD2

```
i <- grep("CRISPLD2", res$symbol)
res[i,]

log2 fold change (MLE): dex treated vs control
Wald test p-value: dex treated vs control
DataFrame with 1 row and 10 columns
  baseMean log2FoldChange      lfcSE      stat     pvalue
  <numeric>      <numeric> <numeric> <numeric>   <numeric>
ENSG00000103196    3096.16      2.62603  0.267444  9.81899 9.32747e-23
  padj      symbol      entrez      uniprot
  <numeric> <character> <character> <character>
ENSG00000103196 3.36344e-20    CRISPLD2      83716 AOA140VK80
  genename
  <character>
ENSG00000103196 cysteine rich secret..
```

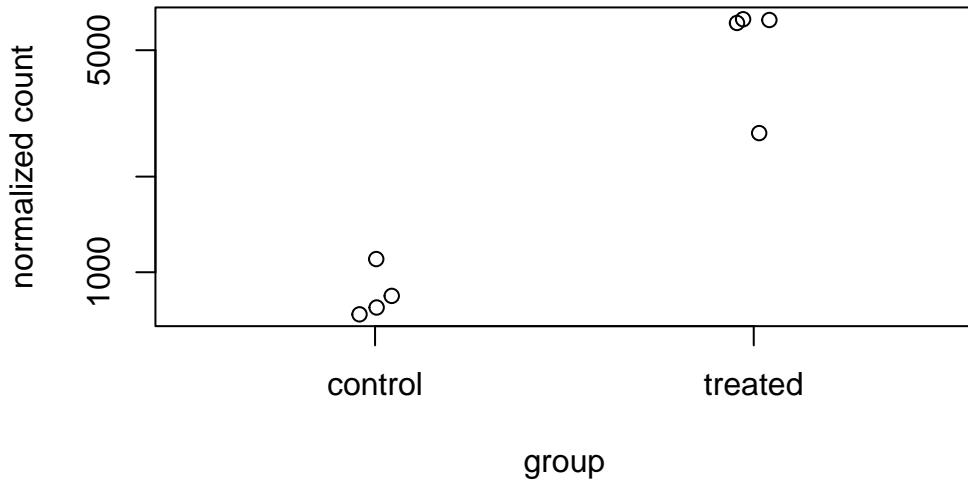
```
rownames(res[i,])
```

```
[1] "ENSG00000103196"
```

Plotting the counts

```
plotCounts(dds, gene="ENSG00000103196", intgroup="dex")
```

ENSG00000103196



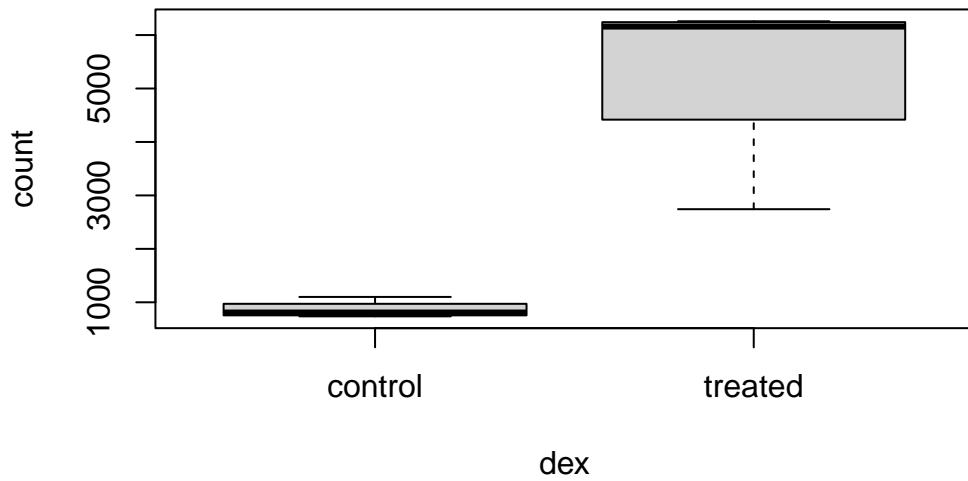
Getting the data instead

```
d <- plotCounts(dds, gene="ENSG00000103196", intgroup="dex", returnData=TRUE)
head(d)
```

	count	dex
SRR1039508	774.5002	control
SRR1039509	6258.7915	treated
SRR1039512	1100.2741	control
SRR1039513	6093.0324	treated
SRR1039516	736.9483	control
SRR1039517	2742.1908	treated

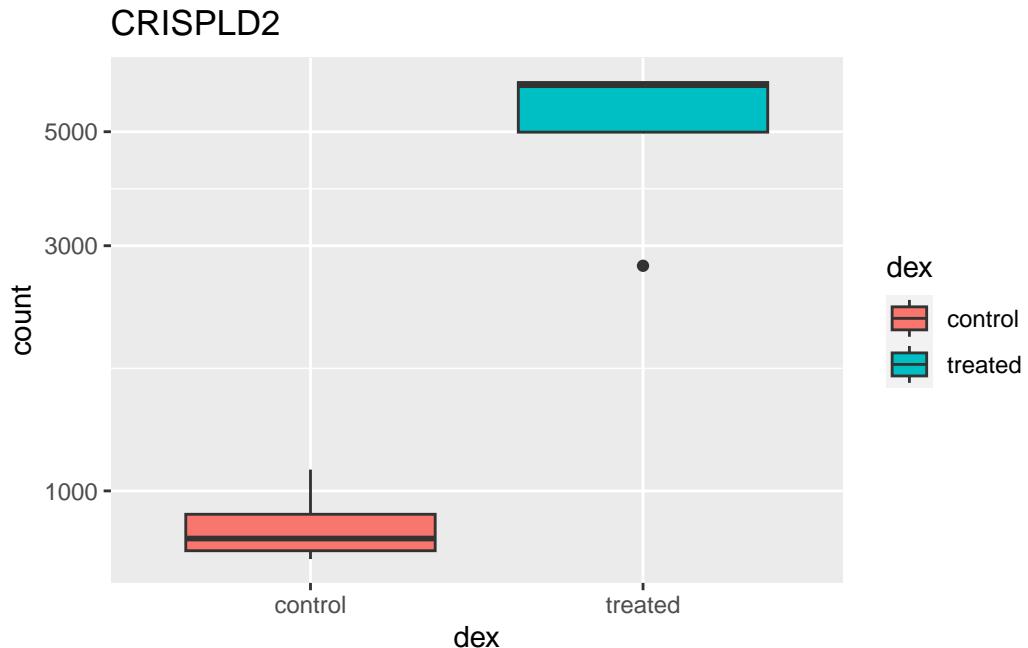
Making a boxplot

```
boxplot(count ~ dex , data=d)
```



New and improved boxplot w/ ggplot2

```
library(ggplot2)
ggplot(d, aes(dex, count, fill=dex)) +
  geom_boxplot() +
  scale_y_log10() +
  ggtitle("CRISPLD2")
```



Class 14:RNA-Seq Analysis Mini-Project

Section 1: Differential Expression Analysis

```
library(DESeq2)
```

Downloading data...

```
metaFile <- "Data/GSE37704_metadata.csv"  
countFile <- "Data/GSE37704_featurecounts.csv"
```

Takin' a peek, so to speak

```
colData = read.csv(metaFile, row.names=1)
head(colData)
```

```
                condition  
SRR493366 control_sirna  
SRR493367 control_sirna  
SRR493368 control_sirna
```

```
SRR493369      hoxa1_kd  
SRR493370      hoxa1_kd  
SRR493371      hoxa1_kd
```

```
countData = read.csv(countFile, row.names=1)  
head(countData)
```

	length	SRR493366	SRR493367	SRR493368	SRR493369	SRR493370
ENSG00000186092	918	0	0	0	0	0
ENSG00000279928	718	0	0	0	0	0
ENSG00000279457	1982	23	28	29	29	28
ENSG00000278566	939	0	0	0	0	0
ENSG00000273547	939	0	0	0	0	0
ENSG00000187634	3214	124	123	205	207	212
		SRR493371				
ENSG00000186092		0				
ENSG00000279928		0				
ENSG00000279457		46				
ENSG00000278566		0				
ENSG00000273547		0				
ENSG00000187634		258				

Q1. Align the columns

```
countData <- as.matrix(countData[,-1])  
head(countData)
```

	SRR493366	SRR493367	SRR493368	SRR493369	SRR493370	SRR493371
ENSG00000186092	0	0	0	0	0	0
ENSG00000279928	0	0	0	0	0	0
ENSG00000279457	23	28	29	29	28	46
ENSG00000278566	0	0	0	0	0	0
ENSG00000273547	0	0	0	0	0	0
ENSG00000187634	124	123	205	207	212	258

Q2. Exclude gene samples w/ no data

```
countData = countData[-which(rowSums(countData) == 0), ]  
head(countData)
```

	SRR493366	SRR493367	SRR493368	SRR493369	SRR493370	SRR493371
ENSG00000279457	23	28	29	29	28	46
ENSG00000187634	124	123	205	207	212	258
ENSG00000188976	1637	1831	2383	1226	1326	1504
ENSG00000187961	120	153	180	236	255	357
ENSG00000187583	24	48	65	44	48	64
ENSG00000187642	4	9	16	14	16	16

```
dds = DESeqDataSetFromMatrix(countData=countData,
                             colData=colData,
                             design=~condition)
```

Warning in DESeqDataSet(se, design = design, ignoreRank): some variables in
design formula are characters, converting to factors

```
dds = DESeq(dds)
```

estimating size factors

estimating dispersions

gene-wise dispersion estimates

mean-dispersion relationship

final dispersion estimates

fitting model and testing

```
dds
```

```
class: DESeqDataSet
dim: 15975 6
metadata(1): version
assays(4): counts mu H cooks
rownames(15975): ENSG00000279457 ENSG00000187634 ... ENSG00000276345
ENSG00000271254
rowData names(22): baseMean baseVar ... deviance maxCooks
colnames(6): SRR493366 SRR493367 ... SRR493370 SRR493371
colData names(2): condition sizeFactor
```

```
res = results(dds, contrast=c("condition", "hoxa1_kd", "control_sirna"))
```

Q3. Use summary function on results

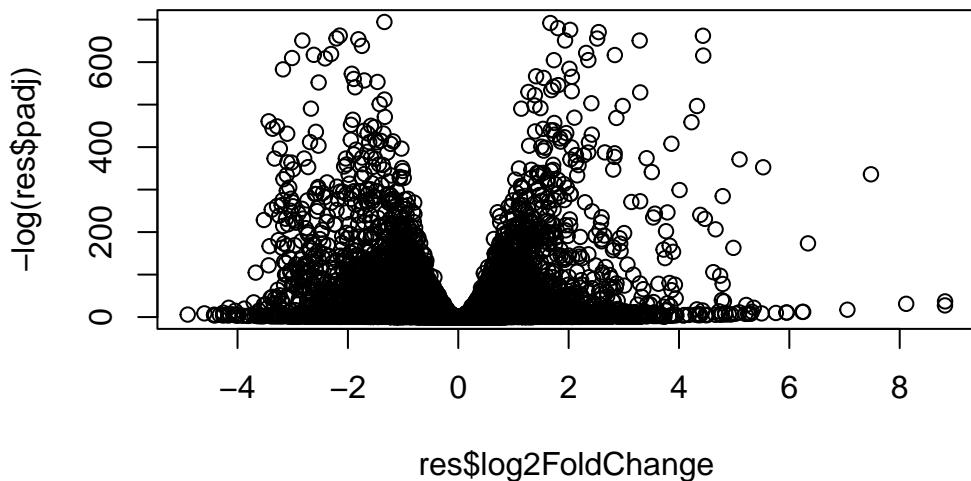
Summarizing results

```
summary(res)
```

```
out of 15975 with nonzero total read count
adjusted p-value < 0.1
LFC > 0 (up)      : 4349, 27%
LFC < 0 (down)    : 4396, 28%
outliers [1]       : 0, 0%
low counts [2]     : 1237, 7.7%
(mean count < 0)
[1] see 'cooksCutoff' argument of ?results
[2] see 'independentFiltering' argument of ?results
```

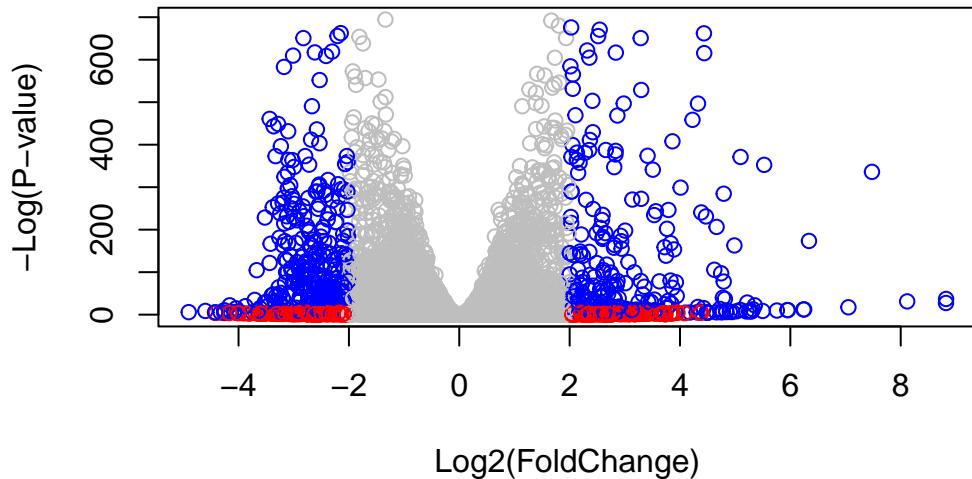
Making a volcano plot

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



Q4. Improve plot by adding color + axis labels

```
mycols <- rep("gray", nrow(res) )  
  
mycols[ abs(res$log2FoldChange) > 2 ] <- "red"  
  
inds <- (res$padj < .01) & (abs(res$log2FoldChange) > 2 )  
mycols[ inds ] <- "blue"  
  
plot( res$log2FoldChange, -log(res$padj), col= mycols, xlab="Log2(FoldChange)", ylab="-Log
```



Q5. Use the mapIDs() function multiple times to add SYMBOL, ENTREZID and GENENAME annotation to our results by completing the code below

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

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

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

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

```
res$name =   mapIds(org.Hs.eg.db,
                     keys=row.names(res),
                     keytype= "ENSEMBL",
                     column="GENENAME",
                     multiVals="first")
```

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

```
head(res, 10)
```

```
log2 fold change (MLE): condition hoxa1_kd vs control_sirna
Wald test p-value: condition hoxa1 kd vs control sirna
DataFrame with 10 rows and 9 columns
  baseMean log2FoldChange      lfcSE       stat      pvalue
  <numeric>      <numeric> <numeric>  <numeric>  <numeric>
ENSG00000279457  29.913579  0.1792571  0.3248216  0.551863 5.81042e-01
ENSG00000187634 183.229650  0.4264571  0.1402658  3.040350 2.36304e-03
```

ENSG00000188976	1651.188076	-0.6927205	0.0548465	-12.630158	1.43990e-36
ENSG00000187961	209.637938	0.7297556	0.1318599	5.534326	3.12428e-08
ENSG00000187583	47.255123	0.0405765	0.2718928	0.149237	8.81366e-01
ENSG00000187642	11.979750	0.5428105	0.5215598	1.040744	2.97994e-01
ENSG00000188290	108.922128	2.0570638	0.1969053	10.446970	1.51282e-25
ENSG00000187608	350.716868	0.2573837	0.1027266	2.505522	1.22271e-02
ENSG00000188157	9128.439422	0.3899088	0.0467163	8.346304	7.04321e-17
ENSG00000237330	0.158192	0.7859552	4.0804729	0.192614	8.47261e-01
	padj	symbol	entrez		name
	<numeric>	<character>	<character>		<character>
ENSG00000279457	6.86555e-01	NA	NA		NA
ENSG00000187634	5.15718e-03	SAMD11	148398	sterile alpha motif ..	
ENSG00000188976	1.76549e-35	NOC2L	26155	NOC2 like nucleolar ..	
ENSG00000187961	1.13413e-07	KLHL17	339451	kelch like family me..	
ENSG00000187583	9.19031e-01	PLEKHN1	84069	pleckstrin homology ..	
ENSG00000187642	4.03379e-01	PERM1	84808	PPARGC1 and ESRR ind..	
ENSG00000188290	1.30538e-24	HES4	57801	hes family bHLH tran..	
ENSG00000187608	2.37452e-02	ISG15	9636	ISG15 ubiquitin like..	
ENSG00000188157	4.21963e-16	AGRN	375790		agrin
ENSG00000237330	NA	RNF223	401934	ring finger protein ..	

Q6. Finally for this section let's reorder these results by adjusted p-value and save them to a CSV file in your current project directory.

```
res = res[order(res$pvalue),]
write.csv(res, file ="deseq_results.csv")
```

Section 2: Pathway Analysis

```
library(pathview)

library(gage)
library(gageData)
data(kegg.sets.hs)
data(sigmet.idx.hs)
kegg.sets.hs = kegg.sets.hs[sigmet.idx.hs]
head(kegg.sets.hs, 3)

$`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"

$`hsa00230 Purine metabolism`
[1] "100"     "10201"   "10606"   "10621"   "10622"   "10623"   "107"     "10714"
[9] "108"     "10846"   "109"     "111"     "11128"   "11164"   "112"     "113"
[17] "114"     "115"     "122481"  "122622"  "124583"  "132"     "158"     "159"
[25] "1633"    "171568"  "1716"    "196883"  "203"     "204"     "205"     "221823"
[33] "2272"    "22978"   "23649"   "246721"  "25885"   "2618"    "26289"  "270"
[41] "271"     "27115"   "272"     "2766"    "2977"    "2982"    "2983"    "2984"
[49] "2986"    "2987"    "29922"   "3000"    "30833"   "30834"   "318"     "3251"
[57] "353"     "3614"    "3615"    "3704"    "377841"  "471"     "4830"    "4831"
[65] "4832"    "4833"    "4860"    "4881"    "4882"    "4907"    "50484"  "50940"
[73] "51082"   "51251"   "51292"   "5136"    "5137"    "5138"    "5139"    "5140"
[81] "5141"    "5142"    "5143"    "5144"    "5145"    "5146"    "5147"    "5148"
[89] "5149"    "5150"    "5151"    "5152"    "5153"    "5158"    "5167"    "5169"
[97] "51728"   "5198"    "5236"    "5313"    "5315"    "53343"  "54107"  "5422"
[105] "5424"    "5425"    "5426"    "5427"    "5430"    "5431"    "5432"    "5433"
[113] "5434"    "5435"    "5436"    "5437"    "5438"    "5439"    "5440"    "5441"
[121] "5471"    "548644"  "55276"   "5557"    "5558"    "55703"   "55811"  "55821"
[129] "5631"    "5634"    "56655"   "56953"   "56985"   "57804"   "58497"  "6240"
[137] "6241"    "64425"   "646625"  "654364"  "661"     "7498"    "8382"    "84172"
[145] "84265"   "84284"   "84618"   "8622"    "8654"    "87178"   "8833"    "9060"
[153] "9061"    "93034"   "953"     "9533"    "954"     "955"     "956"     "957"
[161] "9583"    "9615"

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

  1266      54855      1465      51232      2034      2317
-2.422719  3.201955 -2.313738 -2.059631 -1.888019 -1.649792

```

Runnin' gage pathway analysis

```

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

attributes(keggres)

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

head(keggres$less)

          p.geomean stat.mean      p.val
hsa04110 Cell cycle     8.995727e-06 -4.378644 8.995727e-06
hsa03030 DNA replication 9.424076e-05 -3.951803 9.424076e-05
hsa03013 RNA transport   1.375901e-03 -3.028500 1.375901e-03
hsa03440 Homologous recombination 3.066756e-03 -2.852899 3.066756e-03
hsa04114 Oocyte meiosis    3.784520e-03 -2.698128 3.784520e-03
hsa00010 Glycolysis / Gluconeogenesis 8.961413e-03 -2.405398 8.961413e-03

          q.val set.size      exp1
hsa04110 Cell cycle     0.001448312    121 8.995727e-06
hsa03030 DNA replication 0.007586381    36 9.424076e-05
hsa03013 RNA transport   0.073840037   144 1.375901e-03
hsa03440 Homologous recombination 0.121861535    28 3.066756e-03
hsa04114 Oocyte meiosis    0.121861535   102 3.784520e-03
hsa00010 Glycolysis / Gluconeogenesis 0.212222694    53 8.961413e-03

```

Using pathview()

```
pathview(gene.data=foldchanges, pathway.id="hsa04110")
```

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

```
Info: Working in directory /Users/Nate/Desktop/Fall Quarter 2023 /BIMM 143/Class 14: RNA seq
```

```
Info: Writing image file hsa04110.pathview.png
```

Focusing on top 5 upregulated pathways

```
keggrespathways <- rownames(keggres$greater) [1:5]
keggresids = substr(keggrespathways, start=1, stop=8)
keggresids
```

```
[1] "hsa04640" "hsa04630" "hsa00140" "hsa04142" "hsa04330"
```

Making the mega pathview()

```
pathview(gene.data=foldchanges, pathway.id=keggresids, species="hsa")
```

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

```
Info: Working in directory /Users/Nate/Desktop/Fall Quarter 2023 /BIMM 143/Class 14: RNA seq
```

```
Info: Writing image file hsa04640.pathview.png
```

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

```
Info: Working in directory /Users/Nate/Desktop/Fall Quarter 2023 /BIMM 143/Class 14: RNA seq
```

```
Info: Writing image file hsa04630.pathview.png
```

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

```
Info: Working in directory /Users/Nate/Desktop/Fall Quarter 2023 /BIMM 143/Class 14: RNA seq
```

```
Info: Writing image file hsa00140.pathview.png
```

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

```
Info: Working in directory /Users/Nate/Desktop/Fall Quarter 2023 /BIMM 143/Class 14: RNA seq
```

```
Info: Writing image file hsa04142.pathview.png
```

```
Info: some node width is different from others, and hence adjusted!
```

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

```
Info: Working in directory /Users/Nate/Desktop/Fall Quarter 2023 /BIMM 143/Class 14: RNA seq
```

```
Info: Writing image file hsa04330.pathview.png
```

Q7. Can you do the same procedure as above to plot the pathview figures for the top 5 down-regulated pathways?

Yes

```
keggrespathways2 <- rownames(keggres$less)[1:5]  
keggresids2 = substr(keggrespathways2, start=1, stop=8)  
keggresids2
```

```
[1] "hsa04110" "hsa03030" "hsa03013" "hsa03440" "hsa04114"
```

```
pathview(gene.data=foldchanges, pathway.id=keggresids2, species="hsa")
```

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

```
Info: Working in directory /Users/Nate/Desktop/Fall Quarter 2023 /BIMM 143/Class 14: RNA seq
```

```
Info: Writing image file hsa04110.pathview.png
```

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

```
Info: Working in directory /Users/Nate/Desktop/Fall Quarter 2023 /BIMM 143/Class 14: RNA seq
```

```
Info: Writing image file hsa03030.pathview.png
```

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

```
Info: Working in directory /Users/Nate/Desktop/Fall Quarter 2023 /BIMM 143/Class 14: RNA seq
```

```
Info: Writing image file hsa03013.pathview.png
```

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

```
Info: Working in directory /Users/Nate/Desktop/Fall Quarter 2023 /BIMM 143/Class 14: RNA seq
```

```
Info: Writing image file hsa03440.pathview.png
```

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

```
Info: Working in directory /Users/Nate/Desktop/Fall Quarter 2023 /BIMM 143/Class 14: RNA seq
```

```
Info: Writing image file hsa04114.pathview.png
```

Section 3: Gene Ontology

```
data(go.sets.hs)
data(go.subs.hs)
gobpsets = go.sets.hs[go.subs.hs$BP]
gobpres = gage(foldchanges, gsets=gobpsets, same.dir=TRUE)
lapply(gobpres, head)

$greater
                               p.geomean stat.mean      p.val
GO:0007156 homophilic cell adhesion    8.519724e-05 3.824205 8.519724e-05
GO:0002009 morphogenesis of an epithelium 1.396681e-04 3.653886 1.396681e-04
GO:0048729 tissue morphogenesis        1.432451e-04 3.643242 1.432451e-04
GO:0007610 behavior                  1.925222e-04 3.565432 1.925222e-04
GO:0060562 epithelial tube morphogenesis 5.932837e-04 3.261376 5.932837e-04
GO:0035295 tube development          5.953254e-04 3.253665 5.953254e-04
                                         q.val set.size      exp1
GO:0007156 homophilic cell adhesion    0.1952430     113 8.519724e-05
GO:0002009 morphogenesis of an epithelium 0.1952430     339 1.396681e-04
GO:0048729 tissue morphogenesis        0.1952430     424 1.432451e-04
GO:0007610 behavior                  0.1968058     426 1.925222e-04
GO:0060562 epithelial tube morphogenesis 0.3566193     257 5.932837e-04
GO:0035295 tube development          0.3566193     391 5.953254e-04

$less
                               p.geomean stat.mean      p.val
```

```

GO:0048285 organelle fission          1.536227e-15 -8.063910 1.536227e-15
GO:0000280 nuclear division           4.286961e-15 -7.939217 4.286961e-15
GO:0007067 mitosis                   4.286961e-15 -7.939217 4.286961e-15
GO:0000087 M phase of mitotic cell cycle 1.169934e-14 -7.797496 1.169934e-14
GO:0007059 chromosome segregation      2.028624e-11 -6.878340 2.028624e-11
GO:0000236 mitotic prometaphase       1.729553e-10 -6.695966 1.729553e-10
                                         q.val set.size     exp1
GO:0048285 organelle fission          5.843127e-12    376 1.536227e-15
GO:0000280 nuclear division           5.843127e-12    352 4.286961e-15
GO:0007067 mitosis                   5.843127e-12    352 4.286961e-15
GO:0000087 M phase of mitotic cell cycle 1.195965e-11   362 1.169934e-14
GO:0007059 chromosome segregation      1.659009e-08    142 2.028624e-11
GO:0000236 mitotic prometaphase       1.178690e-07    84 1.729553e-10

$stats
                                         stat.mean     exp1
GO:0007156 homophilic cell adhesion   3.824205 3.824205
GO:0002009 morphogenesis of an epithelium 3.653886 3.653886
GO:0048729 tissue morphogenesis        3.643242 3.643242
GO:0007610 behavior                  3.565432 3.565432
GO:0060562 epithelial tube morphogenesis 3.261376 3.261376
GO:0035295 tube development           3.253665 3.253665

```

Section 4: Reactome Analysis

Listing the significant genes at the 0.05 level as plain text file

```

sig_genes <- res[res$padj <= 0.05 & !is.na(res$padj), "symbol"]
print(paste("Total number of significant genes:", length(sig_genes)))

[1] "Total number of significant genes: 8147"

write.table(sig_genes, file="significant_genes.txt", row.names=FALSE, col.names=FALSE, quo

```

Q8. What pathway has the most significant “Entities p-value”? Do the most significant pathways listed match your previous KEGG results? What factors could cause differences between the two methods?

Mitotic cell cycle has the most significant “Entities p-value”. The most significant pathways listed don’t match my previous KEGG results. I think the methodology of the KEGG system versus the Reactome system causes the difference in results.

Section 5: GO Online

Q9. What pathway has the most significant “Entities p-value”? Do the most significant pathways listed match your previous KEGG results? What factors could cause differences between the two methods?

The pathway with the most significant “Entities p-value” is “GO biological process complete”. These results definitely don’t match up to my KEGG results. I’m not exactly sure what cause the difference because I wasn’t in class, so I don’t completely understand whats happening. But, I’m going to take a stab in the dark and chalk it up to methodology differences between KEGG and Panther.