

Class13

Jae Kim

Import countData and colData

Begin a new R script and use the read.csv() function to read these count data and metadata files.

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

head(counts)
```

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

Sanity check on correspondence of counts and metadata

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

```
[1] TRUE
```

Q1. How many genes are in this dataset?

```
nrow(counts)
```

```
[1] 38694
```

There are 38694 genes in this dataset.

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

```
table(metadata$dex)
```

```
control treated
        4      4
```

```
nControl<- sum(metadata$dex == "control")
```

There are 4 control cell lines in this dataset.

Toy Differential Gene Expression

First, we will take the mean of just the control variables. Before doing so, we must get all the counts. - Step 1. Calculate the mean of the control samples (i.e. columns in countData) Calculate the mean of the treated samples.

- (a) We need to find which columns are “control” samples.
 - We need to look in the metadata (a.k.a colData), specifically in the \$dex column

```
control <- metadata[metadata[, "dex"] == "control", ]
```

- (b) extract all the control columns from countData and call it control.counts

```
control.counts <- counts[, control$id]
```

- (c) calculate the mean value across the rows of control.counts i.e. calculate the mean count values for each gene in the control samples.

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

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

- step 2: calculate the mean of the treated side:

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)

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

	SRR1039509	SRR1039513	SRR1039517	SRR1039521
ENSG00000000003	486	445	1097	604
ENSG00000000005	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.counts <- counts[, treated inds]
treated.means <- rowMeans(treated.counts)
head(treated.means)
```

```

ENSG000000000003 ENSG000000000005 ENSG00000000419 ENSG00000000457 ENSG00000000460
       658.00          0.00        546.00        316.50        78.75
ENSG00000000938
       0.00

```

We now have control and treated mean count values. For ease of book-keeping I will combine these vectors into a new data.frame called `meancounts` ‘

```

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

```

	control.mean	treated.means
ENSG000000000003	900.75	658.00
ENSG000000000005	0.00	0.00
ENSG00000000419	520.50	546.00
ENSG00000000457	339.75	316.50
ENSG00000000460	97.25	78.75
ENSG00000000938	0.75	0.00

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

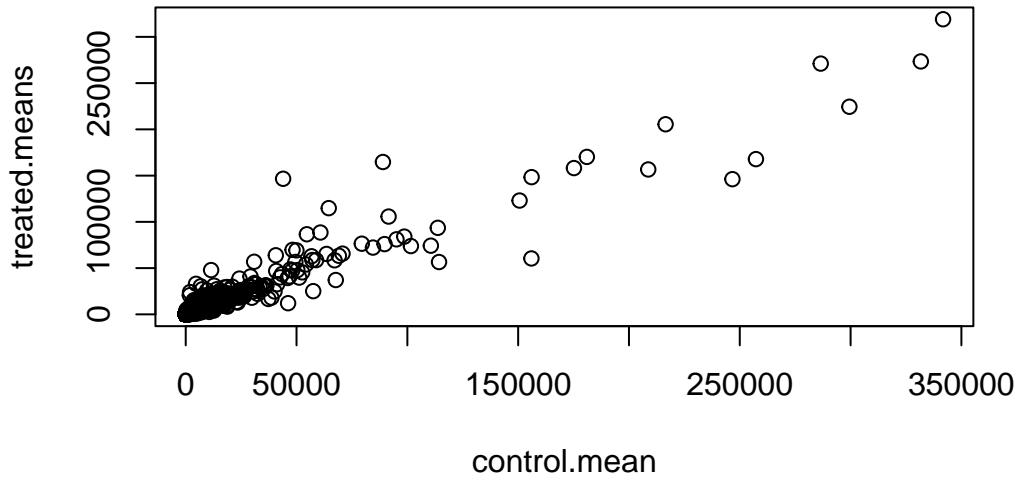
Based on the code from the lab worksheet, we have done it in a more robust and easier way as a class, all above !

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)

```



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?

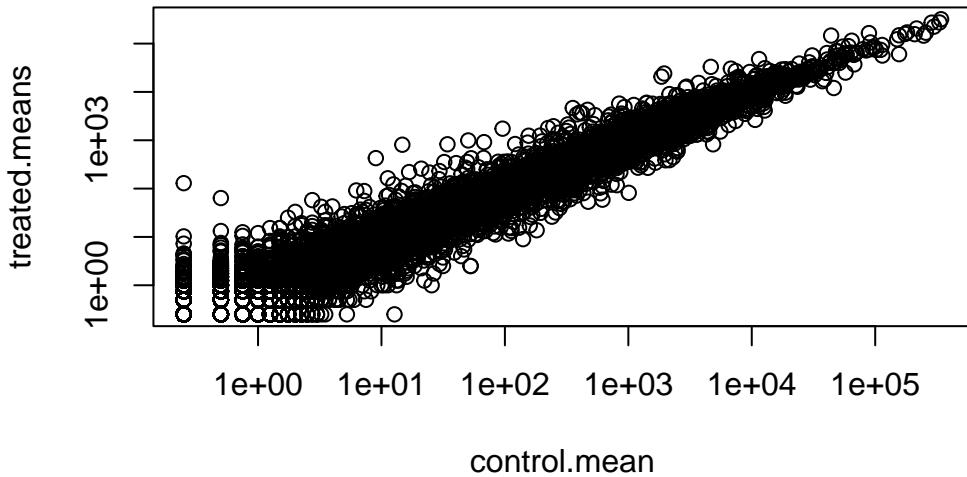
`geom_point()`

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

```
plot(meancounts, log="xy")
```

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

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



We used log transforms for skewed data such as this and because we really care most about relative changes in magnitude. We must often use \log_2 as our transform as the math is easier to interpret than \log_{10} or others. If we have no change - i.e. same values in control and treated will have a \log_2 value of zero.

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

```
[1] 0
```

If I have double the amount i.e. 20 compared to 10 for example, I will have a \log_2 fold-change of +1

```
log2(20/10)
```

```
[1] 1
```

If I have half the amount I will have a \log_2 fold-change of -1.

```
log2(10/20)
```

```
[1] -1

log2(40/10)

[1] 2

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

control.mean treated.means      log2fc
ENSG000000000003     900.75      658.00 -0.45303916
ENSG000000000005      0.00       0.00        NaN
ENSG000000000419     520.50      546.00  0.06900279
ENSG000000000457     339.75      316.50 -0.10226805
ENSG000000000460      97.25       78.75 -0.30441833
ENSG000000000938      0.75       0.00       -Inf
```

The NaN is returned when you divide by zero and try to take the log. The -Inf is returned when you try to take the log of zero. It turns out that there are a lot of genes with zero expression. Let's filter our data to remove these genes. Again inspect your result (and the intermediate steps) to see if things make sense to you

```
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.means      log2fc
ENSG000000000003     900.75      658.00 -0.45303916
ENSG000000000419     520.50      546.00  0.06900279
ENSG000000000457     339.75      316.50 -0.10226805
ENSG000000000460      97.25       78.75 -0.30441833
ENSG000000000971    5219.00     6687.50  0.35769358
ENSG00000001036     2327.00     1785.75 -0.38194109
```

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?

The `arr.ind=TRUE` argument will cause `which()` to return both the row and column indices (i.e. positions) where there are TRUE values. This will tell us which rows and columns have zero counts. We are going to ignore any genes that have zero counts in any sample so we just focus on the row answer

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

Q8. How many genes are upregulated at the common threshold of +2. log2FC values?

```
sum(meancounts$log2fc >= 2, na.rm=TRUE)
```

```
[1] 1910
```

```
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?

```
sum(meancounts$log2fc <= 2, na.rm=TRUE)
```

```
[1] 23412
```

```
down.ind <- mycounts$log2fc < (-2)  
sum(down.ind)
```

```
[1] 367
```

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

No, we do not trust these results, since we do not know and have not gone over the statistics to find out whether or not the changes are statistically significant.

Hold on, what about the statistics? Yes, these are big changes but are these changes significant?

To do this properly, we will turn to DESeq2 package.

DESeq2 analysis

```
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 objects are masked from 'package:base':

expand.grid, I, unname

Loading required package: IRanges

Attaching package: 'IRanges'

```
The following object is masked from 'package:grDevices':
```

```
windows
```

```
Loading required package: GenomicRanges
```

```
Loading required package: GenomeInfoDb
```

```
Loading required package: SummarizedExperiment
```

```
Loading required package: MatrixGenerics
```

```
Loading required package: matrixStats
```

```
Attaching package: 'MatrixGenerics'
```

```
The following objects are masked from 'package:matrixStats':
```

```
colAlls, colAnyNAs, colAnys, colAvgsPerRowSet, colCollapse,  
colCounts, colCummmaxs, 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, rowCummmaxs, 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
```

To use DESeq2 we need our input countData and colData in a specific format that DESeq2 wants:

```
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
```

To run the analysis, I can now use the main DESeq2 function called `DESeq()` with `dds` as input

```
dds <- DESeq(dds)

estimating size factors

estimating dispersions

gene-wise dispersion estimates

mean-dispersion relationship

final dispersion estimates

fitting model and testing
```

To get the results out of the `dds` object we can use the `results()` function from the package.

```

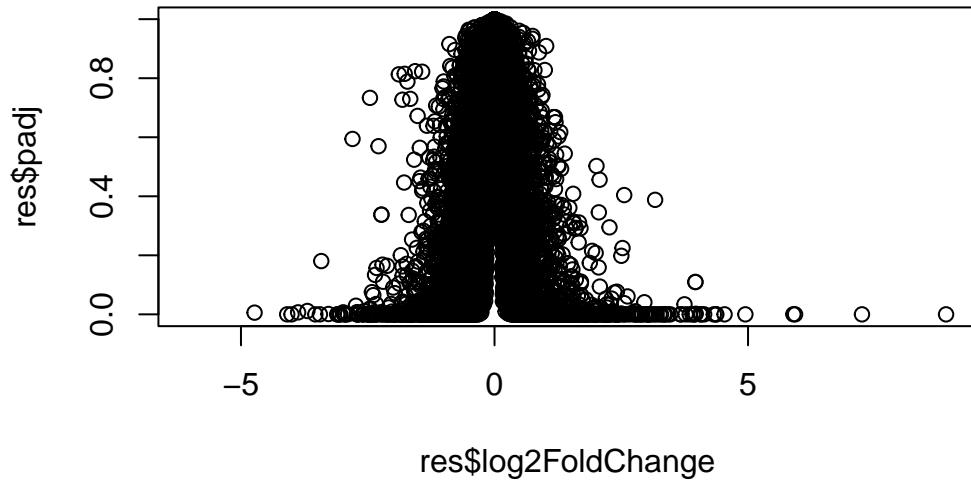
res <- results(dds)
head(res)

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

```

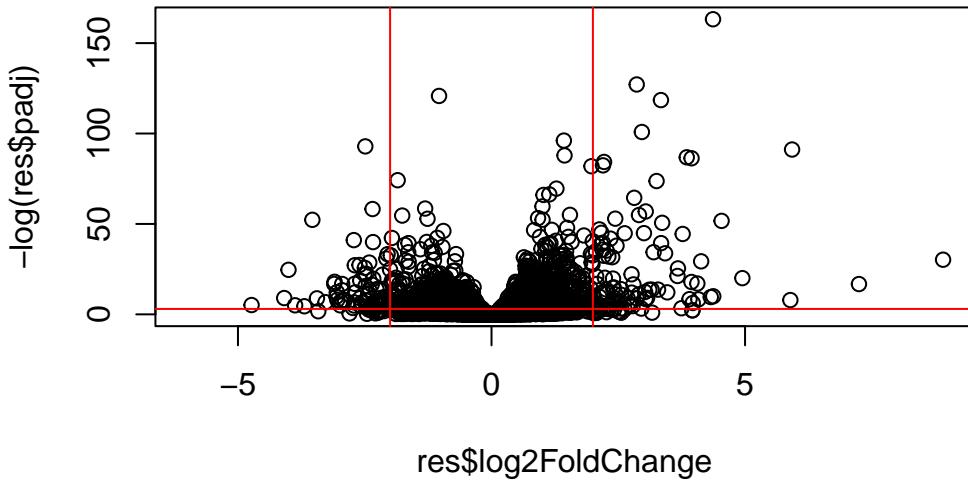
Let's make a final plot of the log2 fold change versus the adjusted p value.

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



It is the low p-values that we care about and these are lost in the skewed plot above. Let's take the log of the \$padj values for our plot.

```
plot(res$log2FoldChange, -log(res$padj))
abline(v=c(+2,-2), col="red")
abline(h=-log(0.05), col= "red")
```



```

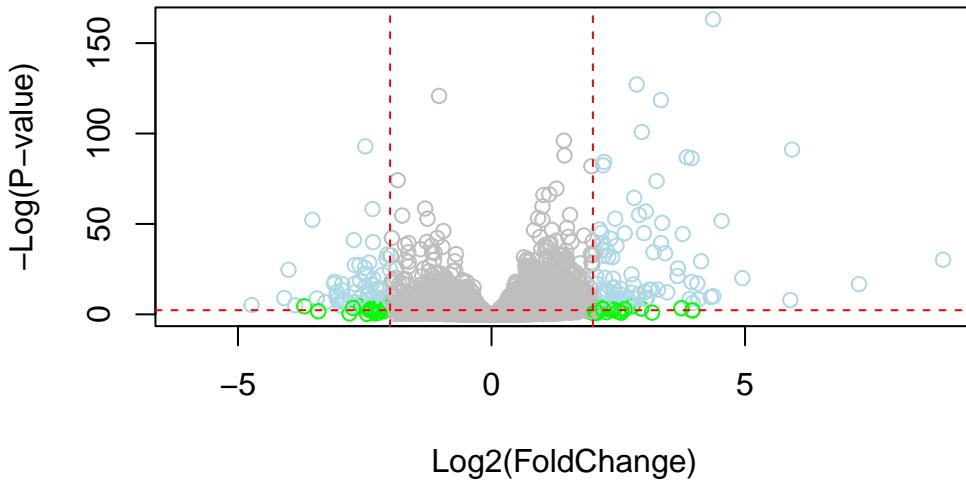
# Setup our custom point color vector
mycols <- rep("gray", nrow(res))
mycols[ abs(res$log2FoldChange) > 2 ] <- "green"

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

# 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="red", lty=2)
abline(h=-log(0.1), col="red", lty=2)

```



Adding Annotation Data

We can use the `AnnotationDbi` package to add annotation data such as gene identifiers from different sources to our results objects.

```
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  0.0000000    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

```

We can translate/map between all these database id formats.

```

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

```

Everything we want to translate (shortened with head())

```

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
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			
	<numeric>	<character>			
ENSG000000000003	0.163035	TSPAN6			
ENSG000000000005	NA	TNMD			
ENSG000000000419	0.176032	DPM1			
ENSG000000000457	0.961694	SCYL3			
ENSG000000000460	0.815849	C1orf112			
ENSG000000000938	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.

```
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  
  

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
```

Save our results to data

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

Pathway analysis

We can use the KEGG database of biological pathways to get some more insight into our differentially expressed genes and the kinds of biology they are involved in.

```
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)
```

Look at the first two KEGG pathways

```
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"  
  

head( res$entrez )  
  

ENSG00000000003 ENSG00000000005 ENSG00000000419 ENSG00000000457 ENSG00000000460  

"7105"          "64102"     "8813"     "57147"     "55732"  

ENSG00000000938  

"2268"
```

Make a new vector of fold-change values that I will use as input for `gage`. This will have the ENTREZ IDs as names

```
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

```
# Get the results  

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

attributes(keggres)
```

```
$names  

[1] "greater" "less"     "stats"
```

```
# 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

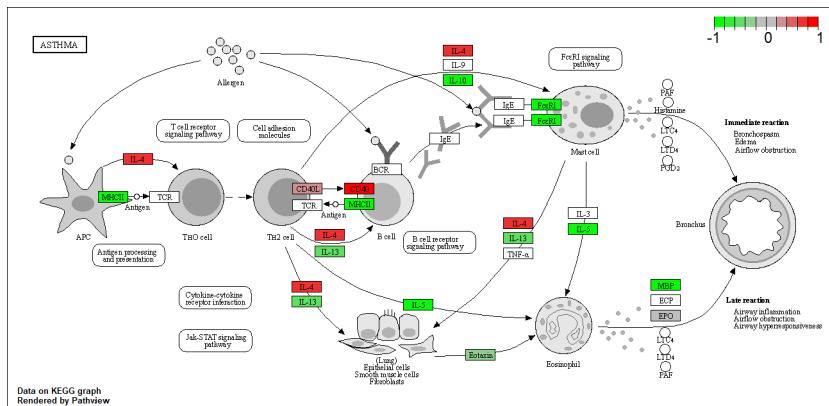
Now I can use the **KEGG IDs** ("hsa05310" etc.) of these pathways from gage to view our genes mapped to these pathways.

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

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

```
Info: Working in directory C:/Users/kjj01/OneDrive/Documents/UCSD 2022-2023/Spring2023/BIMM14
```

```
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?

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

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

```
Info: Working in directory C:/Users/kjj01/OneDrive/Documents/UCSD 2022-2023/Spring2023/BIMM14
```

Info: Writing image file hsa05332.pathview.png

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

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

Info: Working in directory C:/Users/kjj01/OneDrive/Documents/UCSD 2022-2023/Spring2023/BIMM14

Info: Writing image file hsa04940.pathview.png

