

Class13

Nicholas Thiphakhinkeo A17686679

2024-11-12

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

```
head(counts)
```

	SRR1039508	SRR1039509	SRR1039512	SRR1039513	SRR1039516
ENSG00000000003	723	486	904	445	1170
ENSG00000000005	0	0	0	0	0
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
ENSG00000000003	1097	806	604
ENSG00000000005	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

Q1. How many genes are in this dataset ?

38694 genes

Q2. How many 'control' cell lines ?

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

```
control treated
      4      4
```

Finding Control Mean

```
control.inds <- metadata$dex == "control"
control.countd <- counts[, control.inds]
control.mean2 <- apply(control.countd, 1, mean)
```

Finding Treated Mean

```
treated.inds <- metadata$dex == "treated"
treated.countd <- counts[, treated.inds]
treated.mean2 <- apply(treated.countd, 1, mean)
```

Mean Counts

```
meancounts2 <- data.frame(control.mean2, treated.mean2)
head(meancounts2)
```

	control.mean2	treated.mean2
ENSG000000000003	900.75	658.00
ENSG000000000005	0.00	0.00
ENSG000000000419	520.50	546.00
ENSG000000000457	339.75	316.50
ENSG000000000460	97.25	78.75
ENSG000000000938	0.75	0.00

Calculating mean counts/gene across samples

```
control <- metadata[metadata[, "dex"]=="control",]  
control.counts <- counts[ ,control$id]  
control.mean <- rowSums( control.counts )/4  
head(control.mean)
```

```
ENSG000000000003 ENSG000000000005 ENSG000000000419 ENSG000000000457 ENSG000000000460  
          900.75          0.00          520.50          339.75          97.25  
ENSG000000000938  
          0.75
```

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

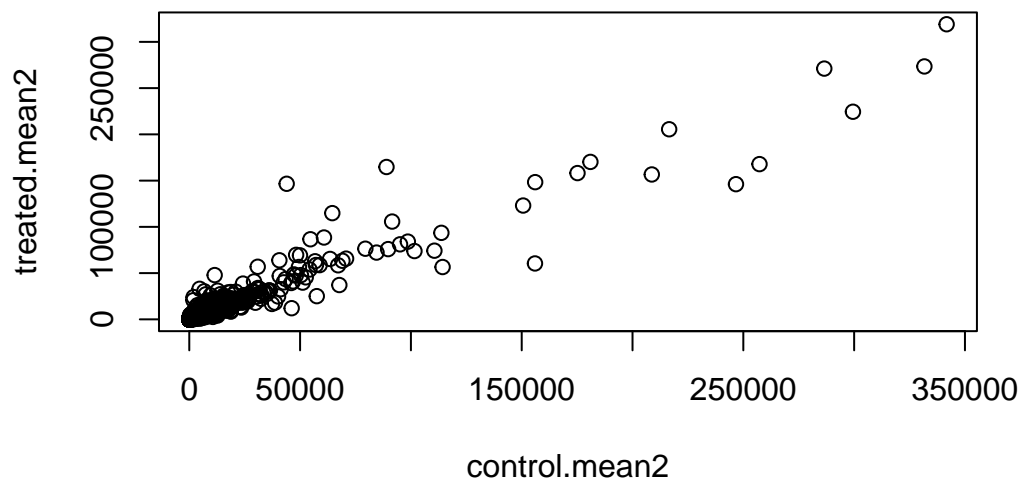
Q4. Follow the same procedure for the treated samples (i.e. calculate the mean per gene across drug treated samples and assign to a labeled vector called treated.mean)

```
treated <- metadata[metadata[, "dex"]=="treated",]  
treated.counts <- counts[ ,treated$id]  
treated.mean <- rowSums( treated.counts )/4  
head(treated.mean)
```

```
ENSG000000000003 ENSG000000000005 ENSG000000000419 ENSG000000000457 ENSG000000000460  
          658.00          0.00          546.00          316.50          78.75  
ENSG000000000938  
          0.00
```

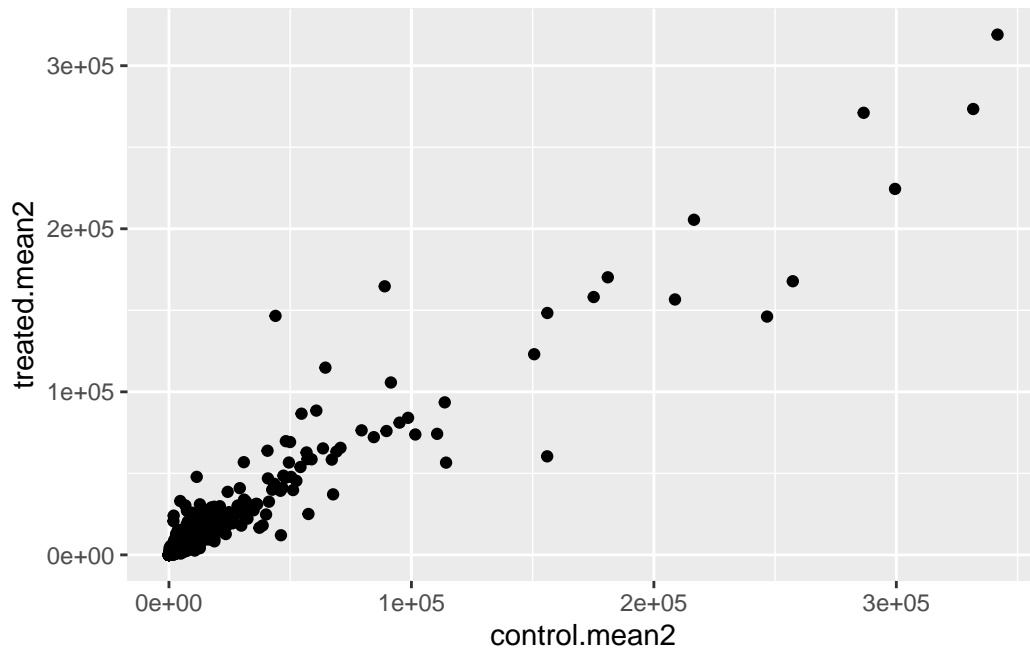
Q5a. Scatterplot of Mean Treated v Control Samples

```
plot(meancounts2)
```



Q5b. Using ggplot2

```
library(ggplot2)
ggplot() + aes(control.mean2,treated.mean2) +
  geom_point()
```



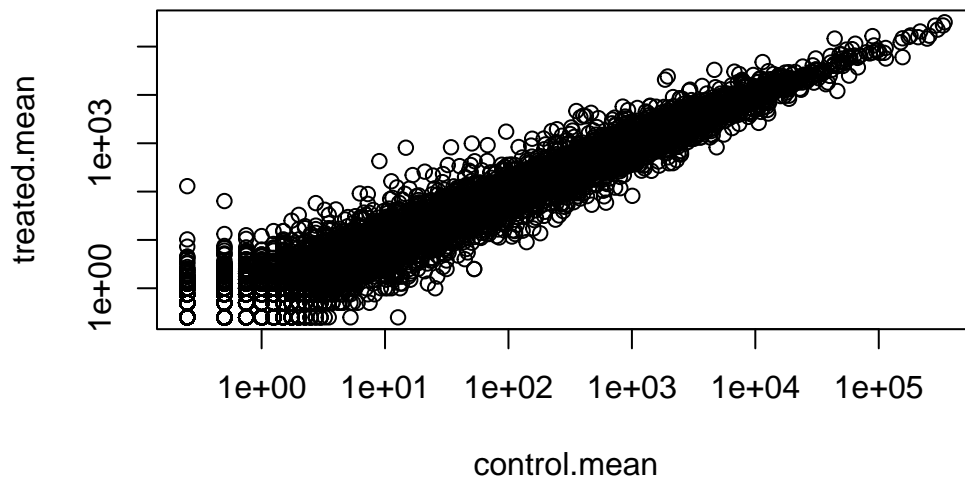
Q6. Plotting axes on log scale

Meancounts

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



rule of thumb: “fold-change” of either ± 2 is where we start to pay attention

```
log2(40/10)
```

```
[1] 2
```

Lets calculate the $\log_2(\text{fold-change})$ and add it to our “mean.counts” df

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

	control.mean	treated.mean	log2fc
ENSG000000000003	900.75	658.00	-0.45303916
ENSG000000000005	0.00	0.00	NaN
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

Filtering Unusable Data

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

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

	control.mean	treated.mean	log2fc
ENSG000000000003	900.75	658.00	-0.45303916
ENSG0000000000419	520.50	546.00	0.06900279
ENSG0000000000457	339.75	316.50	-0.10226805
ENSG0000000000460	97.25	78.75	-0.30441833
ENSG0000000000971	5219.00	6687.50	0.35769358
ENSG0000000001036	2327.00	1785.75	-0.38194109

Q7. What is the purpose of the `arr.ind` argument in `which()` function 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. In this case this will tell us which genes (rows) and samples (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.

How many genes left after filtering ?

```
nrow(mycounts)
```

```
[1] 21817
```

Q8. How many genes upregulated upon drug treatment at +2 log2-fold-change ?

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

```
[1] 250
```

Q9. How many genes downregulated upon drug treatment at -2 log2-fold-change ?

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

```
[1] 367
```

need to perform t-test and get p-value using DESeq Analysis

```
#|message: false
suppressPackageStartupMessages(library(DESeq2))
```

Warning: package 'matrixStats' was built under R version 4.4.2

```
library(DESeq2)
```

Formatting for DESeq

```
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

Running dds in DESeq()

```
dds<-DESeq(dds)
```

estimating size factors

estimating dispersions

gene-wise dispersion estimates

mean-dispersion relationship

final dispersion estimates

fitting model and testing

Results() Function:

```
res <- results(dds)
head(res)
```

log2 fold change (MLE): dex treated vs control

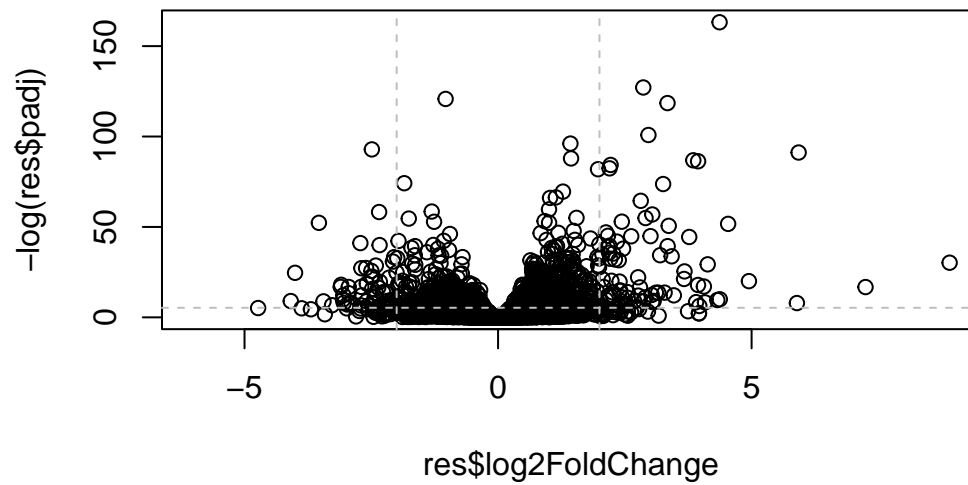
Wald test p-value: dex treated vs control

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
ENSG0000000000419	520.134160	0.2061078	0.101059	2.039475	0.0414026
ENSG0000000000457	322.664844	0.0245269	0.145145	0.168982	0.8658106
ENSG0000000000460	87.682625	-0.1471420	0.257007	-0.572521	0.5669691
ENSG0000000000938	0.319167	-1.7322890	3.493601	-0.495846	0.6200029
	padj				
	<numeric>				
ENSG000000000003	0.163035				
ENSG000000000005	NA				
ENSG0000000000419	0.176032				
ENSG0000000000457	0.961694				
ENSG0000000000460	0.815849				
ENSG0000000000938	NA				

Common Overall Results Figure. Plot Fold-Change vs p-Value

```
plot(res$log2FoldChange, -log(res$padj))
abline(v=c(-2,+2), col="grey", lty=2)
abline(h=-log(0.005), col="grey", lty=2)
```



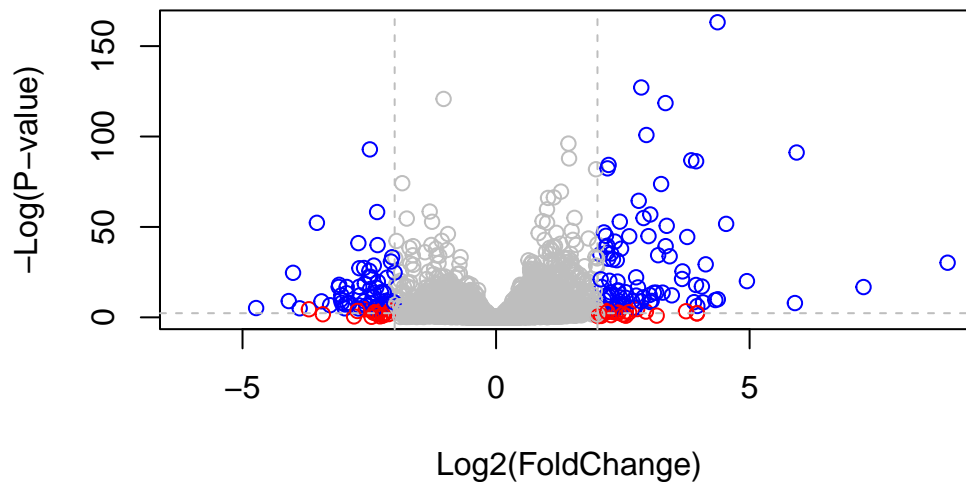
Adding Color

```
# 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)
```



Save Results

```
write.csv(res, file="class13results.csv")
```

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

All of our analysis has been done based on fold change. However, fold change can be large (e.g. »two-fold up- or down-regulation) without being statistically significant (e.g. based on p-values). We have not done anything yet to determine whether the differences we are seeing are significant. These results in their current form are likely to be very misleading.

Need to translate our gene identifiers "ENSG000..." into gene names that are more readily interpretable using *annotationDbi*

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

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

```
[1] "ACCNUM"      "ALIAS"      "ENSEMBL"    "ENSEMBLPROT" "ENSEMBLTRANS"
[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"
```

using *mapIds()* function to map my identifiers from different databases. I will go between “ENSEMBL” and “SYMBOL” (and then after “GENENAME”)

```
res$symbol <- mapIds(org.Hs.eg.db,
  keys = rownames(res),
  keytype = "ENSEMBL",
  column = "SYMBOL")
```

'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>
ENSG00000000003	747.194195	-0.3507030	0.168246	-2.084470	0.0371175
ENSG00000000005	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>			
ENSG00000000003	0.163035	TSPAN6			
ENSG00000000005	NA	TNMD			

ENSG000000000419	0.176032	DPM1
ENSG000000000457	0.961694	SCYL3
ENSG000000000460	0.815849	FIRRM
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`.

Add “GENENAME”

```
res$genename <- mapIds(org.Hs.eg.db,
  keys = rownames(res),
  keytype = "ENSEMBL",
  column = "GENENAME")
```

'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>
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
	<numeric>	<character>	<character>
ENSG000000000003	0.163035	TSPAN6	tetraspanin 6
ENSG000000000005	NA	TNMD	tenomodulin
ENSG000000000419	0.176032	DPM1	dolichyl-phosphate m..
ENSG000000000457	0.961694	SCYL3	SCY1 like pseudokina..
ENSG000000000460	0.815849	FIRRM	FIGNL1 interacting r..
ENSG000000000938	NA	FGR	FGR proto-oncogene, ..

Add “ENTREZID”

```
res$entrezid <- mapIds(org.Hs.eg.db,  
  keys = rownames(res),  
  keytype = "ENSEMBL",  
  column = "ENTREZID")
```

'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
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	entrezid
	<numeric>	<character>	<character>	<character>
ENSG00000000003	0.163035	TSPAN6	tetraspanin 6	7105
ENSG00000000005	NA	TNMD	tenomodulin	64102
ENSG000000000419	0.176032	DPM1	dolichyl-phosphate m..	8813
ENSG000000000457	0.961694	SCYL3	SCY1 like pseudokina..	57147
ENSG000000000460	0.815849	FIRRM	FIGNL1 interacting r..	55732
ENSG000000000938	NA	FGR	FGR proto-oncogene, ..	2268

Saving Annotated Results

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

Pathway Analysis using *gage* to look for KEGG pathways in our genes of interest

```
#|message: false
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)
```

gage wants “vector of importance.” For RNASeq daa like we have this is our log2FC values

```
foldchanges = res$log2FoldChange
names(foldchanges) = res$entrezid
head(foldchanges)
```

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

gage pathway analysis

```
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

Using Pathview() to look at highlighter KEGG Pathways. “hsa05310 Asthma”

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

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

Info: Working in directory C:/Users/Owner/Desktop/School/UCSD Q1 2024/BIMM143/Class13

Info: Writing image file hsa05310.pathview.png

