

lab15

Ziyuan_Han

11/17/2021

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

```
## Bioconductor version 3.11 (BiocManager 1.30.10), ?BiocManager::install for help
```

```
## Bioconductor version '3.11' is out-of-date; the current release version '3.14'
##   is available with R version '4.1'; see https://bioconductor.org/install
```

```
library(DESeq2)
```

```
## Loading required package: S4Vectors
```

```
## Loading required package: stats4
```

```
## Loading required package: BiocGenerics
```

```
## Loading required package: parallel
```

```
##
```

```
## Attaching package: 'BiocGenerics'
```

```
## The following objects are masked from 'package:parallel':
```

```
##
```

```
##   clusterApply, clusterApplyLB, clusterCall, clusterEvalQ,
##   clusterExport, clusterMap, parApply, parCapply, parLapply,
##   parLapplyLB, parRapply, parSapply, parSapplyLB
```

```
## 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, which.max, which.min
```

```

##
## Attaching package: 'S4Vectors'

## The following object is masked from 'package:base':
##
##     expand.grid

## Loading required package: IRanges

## Loading required package: GenomicRanges

## Loading required package: GenomeInfoDb

## Loading required package: SummarizedExperiment

## 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)".

## Loading required package: DelayedArray

## Loading required package: matrixStats

##
## Attaching package: 'matrixStats'

## The following objects are masked from 'package:Biobase':
##
##     anyMissing, rowMedians

##
## Attaching package: 'DelayedArray'

## The following objects are masked from 'package:matrixStats':
##
##     colMaxs, colMins, colRanges, rowMaxs, rowMins, rowRanges

## The following objects are masked from 'package:base':
##
##     aperm, apply, rowsum

#input the data for RNAseq

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
## 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
```

#Check the correspondence of the metadata and count data

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

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

```
all(colnames(counts)==metadata$id) #check whether they are matched
```

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

#Q1. How many genes are in this dataset? 38694

```
nrow(counts)
```

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

#Q2. How many 'control' cell lines do we have? 4

```
length(which(metadata$dex=="control"))
```

```
## [1] 4
```

#Compare control to treated first we need to access all the control columns in our counts data.

```
control.inds <- metadata$dex=="control"
control.inds<-metadata[control.inds,]$id
metadata$dex=="treated"
```

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

```
#Use the ids to access just the control fcolumns of our 'counts' data
```

```
head(counts[,control.inds])
```

```
##                SRR1039508 SRR1039512 SRR1039516 SRR1039520
## ENSG000000000003         723         904         1170         806
## ENSG000000000005          0          0          0          0
## ENSG000000000419         467         616         582         417
## ENSG000000000457         347         364         318         330
## ENSG000000000460          96          73         118         102
## ENSG000000000938          0           1           2           0
```

```
control.mean<-rowMeans(counts[,control.inds])
head(control.mean)
```

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

```
#Do the same for the drug treated
```

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

```
##                SRR1039509 SRR1039513 SRR1039517 SRR1039521
## ENSG000000000003         486         445         1097         604
## ENSG000000000005          0          0          0          0
## ENSG000000000419         523         371         781         509
## ENSG000000000457         258         237         447         324
## ENSG000000000460          81          66          94          74
## ENSG000000000938          0           0           0           0
```

```
treated.mean<-rowMeans(counts[,treated.inds])
head(treated.mean)
```

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

```
#combine our meancount data for bookkeeping
```

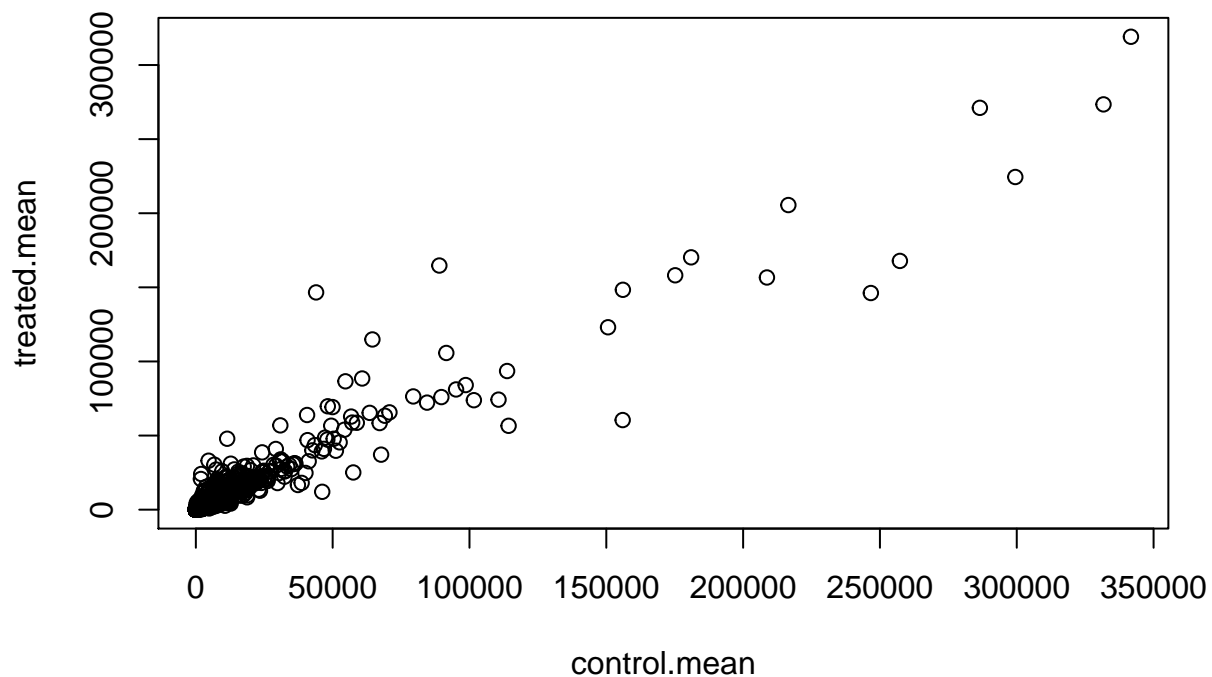
```
meancounts <- data.frame(control.mean,treated.mean)
nrow(counts)
```

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

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

a quick plot of our progress so far

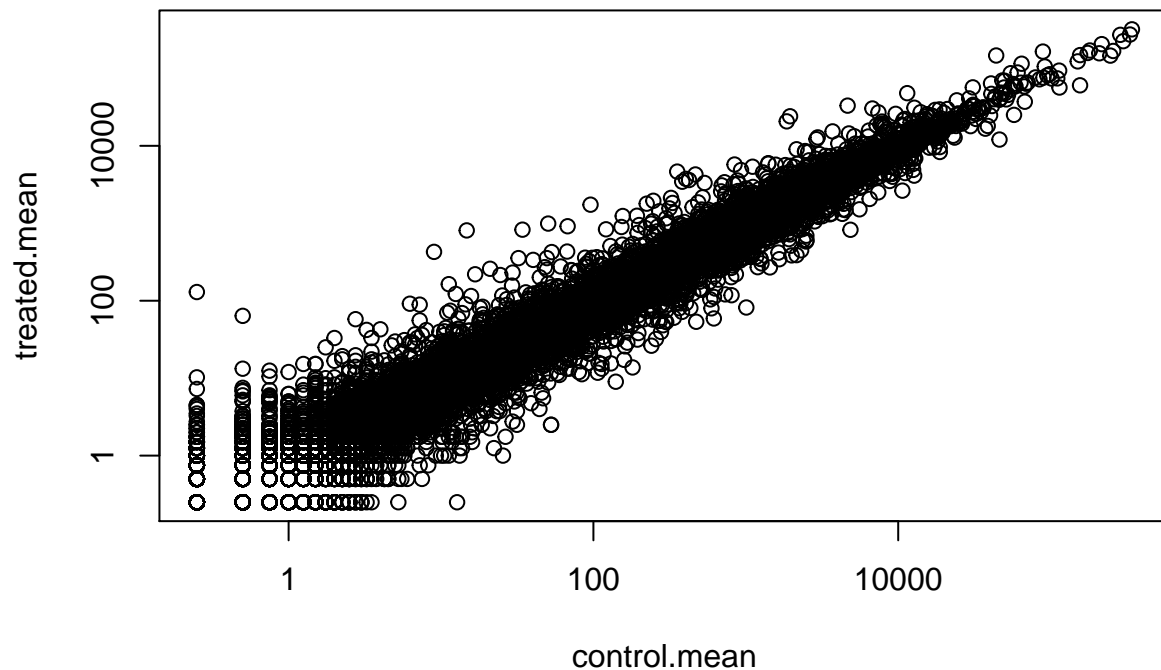
```
plot(meancounts)
```



```
#this would benefit from a single log transform from a plot
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
```



```
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
```

#we need to drop the zero count genes/rows #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 a logical statement that judge whether should array indices be returned when x is an array. Calling unique() will ensure we dont count any row twice.

```
head(meancounts[, 1:2] == 0)
```

```
##           control.mean treated.mean
## ENSG000000000003      FALSE      FALSE
## ENSG000000000005       TRUE       TRUE
## ENSG000000000419      FALSE      FALSE
## ENSG000000000457      FALSE      FALSE
## ENSG000000000460      FALSE      FALSE
## ENSG000000000938      FALSE       TRUE
```

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

```
to.rm<-unique(sort(inds[, "row"]))  
mycounts<-meancounts[-to.rm,]  
head(mycounts)
```

```
##               control.mean treated.mean      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
```

```
#we now have genes remaining as 'r nrow(mycounts)'
```

```
nrow(mycounts)
```

```
## [1] 21817
```

```
sum(mycounts < -2)
```

```
## [1] 367
```

#how many of the genes are up regulated at the log2 fold-change threshold of +2 or greater and how to calculate the percentage

```
round((sum(mycounts$log2fc > +2) / nrow(mycounts))*100,2)
```

```
## [1] 1.15
```

#we first need to setup the DEseq object and run the pipeline

```
library(DESeq2)  
citation("DESeq2")
```

```
##  
## Love, M.I., Huber, W., Anders, S. Moderated estimation of fold change  
## and dispersion for RNA-seq data with DESeq2 Genome Biology 15(12):550  
## (2014)  
##  
## A BibTeX entry for LaTeX users is  
##  
## @Article{,  
## title = {Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2},  
## author = {Michael I. Love and Wolfgang Huber and Simon Anders},  
## year = {2014},  
## journal = {Genome Biology},  
## doi = {10.1186/s13059-014-0550-8},  
## volume = {15},  
## issue = {12},  
## pages = {550},  
## }
```

```
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
```

```
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)  
summary(res)
```

```
##  
## out of 25258 with nonzero total read count  
## adjusted p-value < 0.1  
## LFC > 0 (up)      : 1563, 6.2%  
## LFC < 0 (down)    : 1188, 4.7%  
## outliers [1]      : 142, 0.56%  
## low counts [2]    : 9971, 39%  
## (mean count < 10)  
## [1] see 'cooksCutoff' argument of ?results  
## [2] see 'independentFiltering' argument of ?results
```

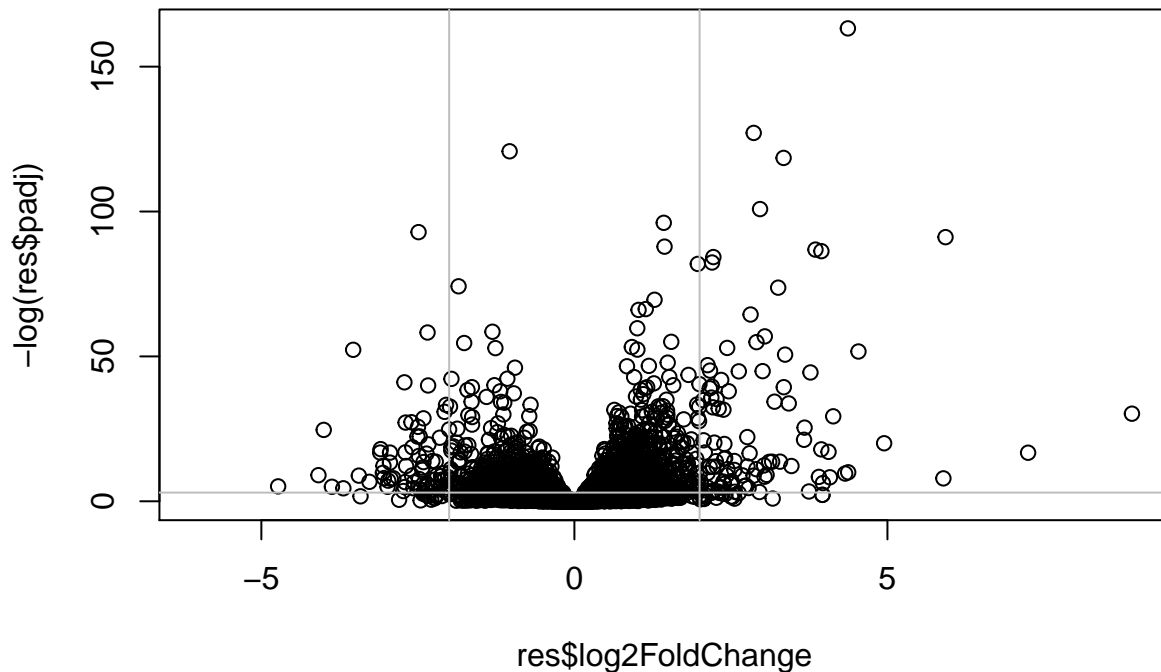
```
res05 <- results(dds, alpha=0.05)  
summary(res05)
```

```
##  
## out of 25258 with nonzero total read count  
## adjusted p-value < 0.05  
## LFC > 0 (up)      : 1236, 4.9%  
## LFC < 0 (down)    : 933, 3.7%  
## outliers [1]      : 142, 0.56%  
## low counts [2]    : 9033, 36%  
## (mean count < 6)  
## [1] see 'cooksCutoff' argument of ?results  
## [2] see 'independentFiltering' argument of ?results
```



```
#A Volcano plot
```

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



```
#add annotations #BiocManager::install("org.Hs.eg.db") #BiocManager::install("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] "GO"          "GOALL"       "IPI"          "MAP"          "OMIM"  
## [16] "ONTOLOGY"    "ONTOLOGYALL" "PATH"         "PFAM"         "PMID"  
## [21] "PROSITE"     "REFSEQ"      "SYMBOL"       "UCSCKG"       "UNIGENE"  
## [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
## 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

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

```
#Pathway analysis Let's try to bring some insights into this #BiocManager::install( c("pathview", "gage",
"gageData") )
```

```
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"
```

Before we can use KEGG we need to get our gene identifiers in the correct format for KEGG which is ENTREZ format in this case

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

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

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

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

```
foldchanges<-res$log2FoldChange
head(foldchanges)
```

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

```
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"
```

```
head(keggres$less)
```

```
##                                     p.geomean stat.mean
## hsa05332 Graft-versus-host disease 0.0004250461 -3.473346
## hsa04940 Type I diabetes mellitus 0.0017820293 -3.002352
## hsa05310 Asthma 0.0020045888 -3.009050
## hsa04672 Intestinal immune network for IgA production 0.0060434515 -2.560547
## hsa05330 Allograft rejection 0.0073678825 -2.501419
## hsa04340 Hedgehog signaling pathway 0.0133239547 -2.248547
##                                     p.val      q.val
## hsa05332 Graft-versus-host disease 0.0004250461 0.09053483
## hsa04940 Type I diabetes mellitus 0.0017820293 0.14232581
## hsa05310 Asthma 0.0020045888 0.14232581
## hsa04672 Intestinal immune network for IgA production 0.0060434515 0.31387180
## hsa05330 Allograft rejection 0.0073678825 0.31387180
## hsa04340 Hedgehog signaling pathway 0.0133239547 0.47300039
##                                     set.size      exp1
## hsa05332 Graft-versus-host disease 40 0.0004250461
## hsa04940 Type I diabetes mellitus 42 0.0017820293
## hsa05310 Asthma 29 0.0020045888
## hsa04672 Intestinal immune network for IgA production 47 0.0060434515
## hsa05330 Allograft rejection 36 0.0073678825
## hsa04340 Hedgehog signaling pathway 56 0.0133239547
```

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

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

```
## Info: Working in directory /Users/ziyuanhan/Desktop/UCSD PhD/Courses VS Assignment/2021Fall Courses/1
```

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

