

Class13_Transcriptomics and the analysis of RNA-Seq data

Changcheng Li (PID: A69027828)

The data for this hands-on session comes from a published RNA-seq experiment where airway smooth muscle cells were treated with dexamethasone, a synthetic glucocorticoid steroid with anti-inflammatory effects (Himes et al. 2014).

```
library(BiocManager)
library(DESeq2)
```

Data import

```
# Complete the missing code
counts <- read.csv("airway_scaledcounts.csv", row.names=1)
metadata <- read.csv("airway_metadata.csv")
```

Q1. How many genes are in this dataset?

```
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 | | |

| | | | |
|-----------------|-----|-----|-----|
| ENSG00000000419 | 781 | 417 | 509 |
| ENSG00000000457 | 447 | 330 | 324 |
| ENSG00000000460 | 94 | 102 | 74 |
| ENSG00000000938 | 0 | 0 | 0 |

```
dim(counts)
```

```
[1] 38694      8
```

There are 38694 genes

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

```
sum(metadata$dex == "control")
```

```
[1] 4
```

There are 4 “control” cell lines

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

```
control treated
      4      4
```

```
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 |
| 7 | SRR1039520 | control | N061011 | GSM1275874 |
| 8 | SRR1039521 | treated | N061011 | GSM1275875 |

I want to compare the control to the treated columns. To do this I will

-Step 1. Identify and extract the “control” columns. -Step 2. Calculate the mean value per gene for all these “control” columns -Step 3. Do the same for treated -Step 4. Compare the ‘control.mean’ and ‘treated.mean’ values

Step 1:

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

```
control.inds <- metadata$dex == "control"
```

```
metadata[control.inds, ]
```

| | 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.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 | | | | |

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(counts[,metadata$dex=="treated"])
```

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

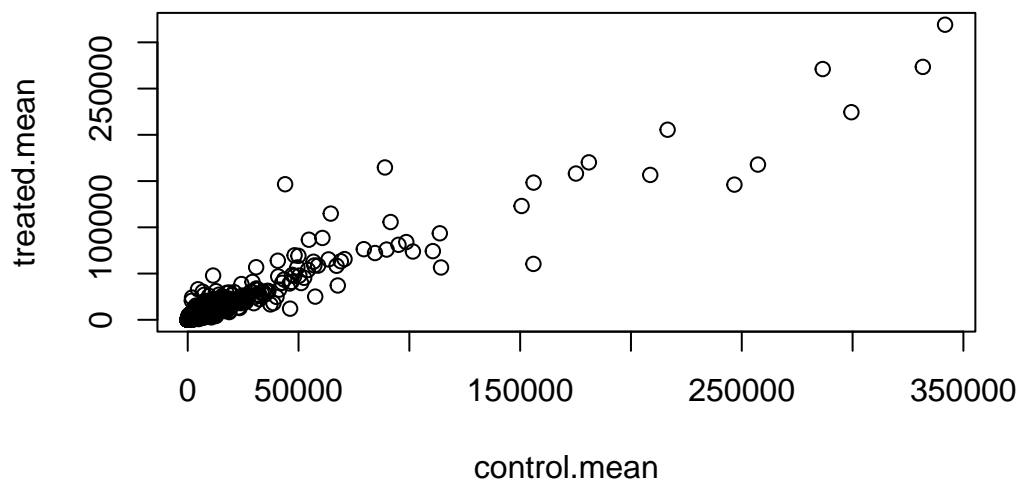
```
head(meancounts)
```

| | control.mean | treated.mean |
|------------------|--------------|--------------|
| ENSG000000000003 | 900.75 | 658.00 |

| | | |
|-------------------|--------|--------|
| ENSG000000000005 | 0.00 | 0.00 |
| ENSG0000000000419 | 520.50 | 546.00 |
| ENSG0000000000457 | 339.75 | 316.50 |
| ENSG0000000000460 | 97.25 | 78.75 |
| ENSG0000000000938 | 0.75 | 0.00 |

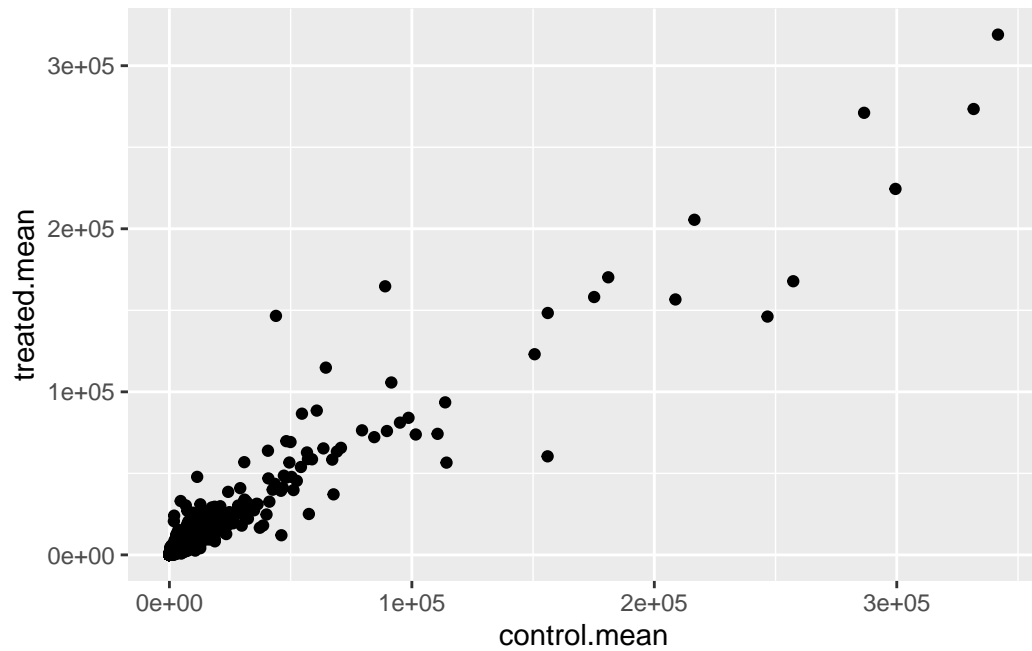
Q5 (a). Create a scatter plot showing the mean of the treated samples against the mean of the control samples. Your plot should look something like the following.

```
plot(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?

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

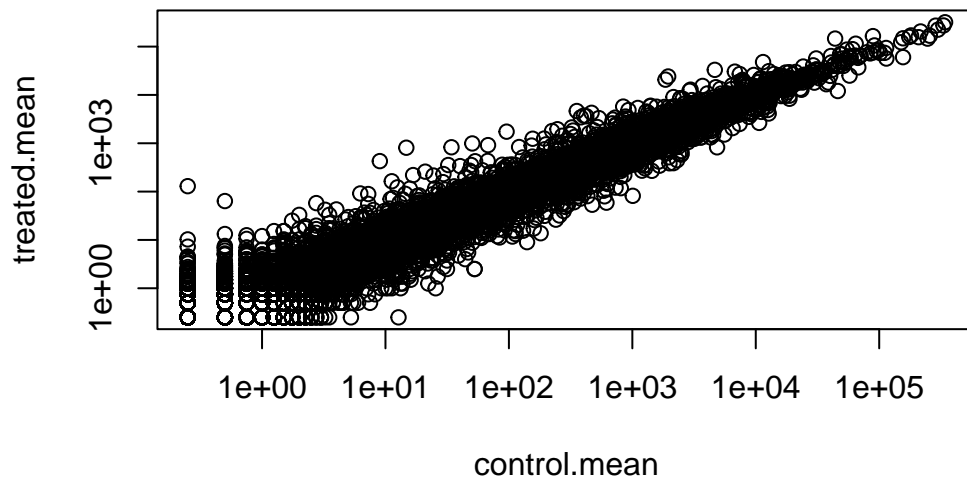


Q6. Try plotting both axes on a log scale. What is the argument to `plot()` that allows you to do this?

```
plot(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



Logs are super useful when we have such skewed data

```
# Treated / control
log2(20/10)
```

```
[1] 1
```

Add log2(Fold-change) values to our wee results table.

```
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 |
| ENSG0000000000419 | 520.50 | 546.00 | 0.06900279 |
| ENSG0000000000457 | 339.75 | 316.50 | -0.10226805 |
| ENSG0000000000460 | 97.25 | 78.75 | -0.30441833 |
| ENSG0000000000938 | 0.75 | 0.00 | -Inf |

I need to exclude any genes with zero counts as we cannot say anything about them anyway from this experiment and it causes me math pain.

```
# What values in the first two cols are zero
to.rm.inds <- rowSums(meancounts[,1:2] == 0) > 0
mycounts <- meancounts[!to.rm.inds, ]
```

```
which(c(TRUE, FALSE, TRUE))
```

```
[1] 1 3
```

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

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

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. In this case this will tell us which genes (rows) and samples (columns) have zero counts.

Q. How many genes do I have left

```
nrow(mycounts)
```

```
[1] 21817
```

There are 21817 genes left

Q8. How many genes are “up regulated” i.e. have a log2(fold-change) greater than +2?

```
sum(mycounts$log2fc > +2)
```

```
[1] 250
```

There are 250 up regulated genes.

Q9. How many genes are “down regulated” i.e. have a $\log_2(\text{fold-change})$ less than -2?

```
sum(mycounts$log2fc < -2)
```

```
[1] 367
```

There are 367 up regulated genes.

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

I do not fully trust these results. 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

Running DESeq

Like many bioconductor analysis packages DESeq wants its input in a very particular way

```
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 DESeq analysis we call the main function from the package called ‘DESeq(dds)’

```
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 this 'dds' object we can use the DESeq '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

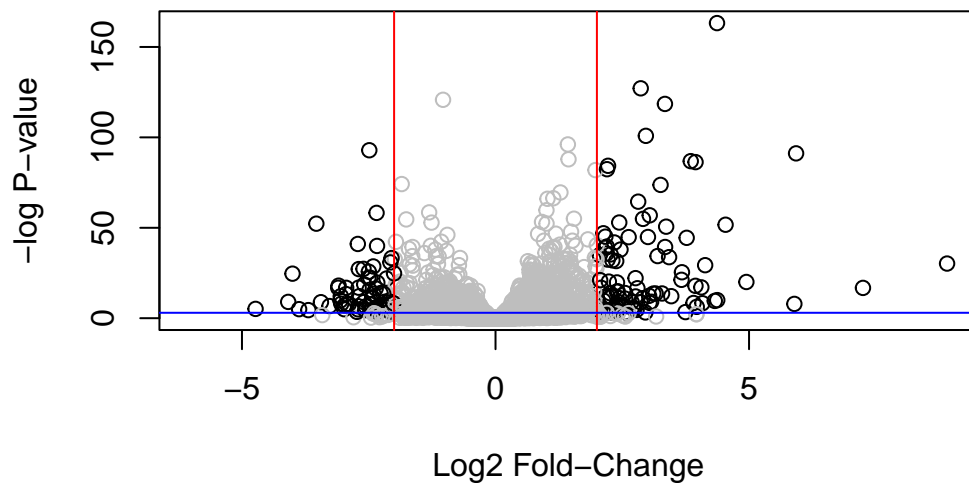
| | baseMean | log2FoldChange | lfcSE | stat | pvalue |
|-------------------|------------|----------------|-----------|-----------|-----------|
| | <numeric> | <numeric> | <numeric> | <numeric> | <numeric> |
| ENSG000000000003 | 747.194195 | -0.3507030 | 0.168246 | -2.084470 | 0.0371175 |
| ENSG000000000005 | 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 | | | | |

A common summary visualization is called a Volcano plot.

```
mycols <- rep("gray", nrow(res))
mycols[res$log2FoldChange > 2] <- "black"
mycols[res$log2FoldChange < -2] <- "black"
mycols[res$padj > 0.05] <- "gray"
```

```
plot(res$log2FoldChange, -log(res$padj), col = mycols,
     xlab = "Log2 Fold-Change",
     ylab = "-log P-value")

abline(v = c(-2,2), col = "red")
abline(h = -log(0.05), col = "blue")
```



Save our results to date

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

Adding annotation data

We need to translate or “map” our ensemble IDs into more understandable gene names and the identifiers.

```
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"      "UCSCCKG"
[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 | 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`.

```
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$uniprot <- 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

```
res$genename <- 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 |

| | padj | symbol | entrez | uniprot |
|------------------|-------------|-------------|-------------|------------------------|
| | <numeric> | <character> | <character> | <character> |
| ENSG000000000457 | 322.664844 | | | |
| ENSG000000000460 | 87.682625 | | | |
| ENSG000000000938 | 0.319167 | | | |
| ENSG000000000003 | 0.163035 | TSPAN6 | 7105 | tetraspanin 6 |
| ENSG000000000005 | NA | TNMD | 64102 | tenomodulin |
| ENSG000000000419 | 0.176032 | DPM1 | 8813 | dolichyl-phosphate m.. |
| ENSG000000000457 | 0.961694 | SCYL3 | 57147 | SCY1 like pseudokina.. |
| ENSG000000000460 | 0.815849 | FIRRM | 55732 | FIGNL1 interacting r.. |
| ENSG000000000938 | NA | FGR | 2268 | FGR proto-oncogene, .. |
| | genename | | | |
| | <character> | | | |
| ENSG000000000003 | AOA024RCIO | | | |
| ENSG000000000005 | Q9H2S6 | | | |
| ENSG000000000419 | O60762 | | | |
| ENSG000000000457 | Q8IZE3 | | | |
| ENSG000000000460 | AOA024R922 | | | |
| ENSG000000000938 | P09769 | | | |

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)

# 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

# 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)

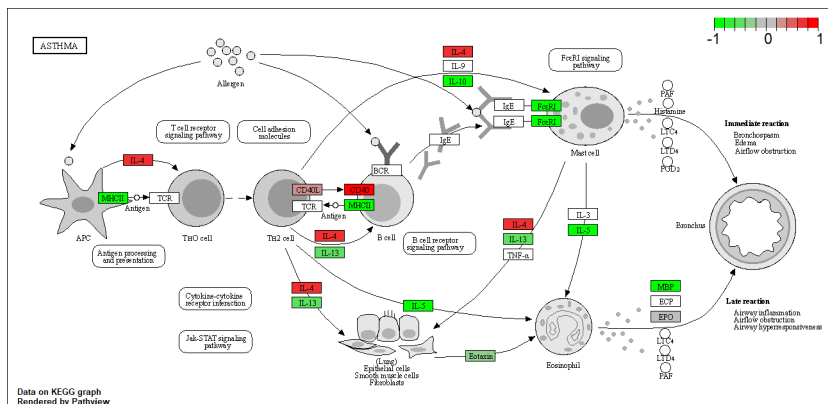
```

Lets have a look at one of these pathways

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

Info: Working in directory E:/ /UCSD/Course/BGGN213 Found Bioinfor/Week6/12 Genome informati

Info: Writing image file hsa05310.pathview.png



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

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

Info: Working in directory E:/ /UCSD/Course/BGGN213 Found Bioinfor/Week6/12 Genome informati

Info: Writing image file hsa05310.pathview.pdf

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

```
# Look at 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 E:/ /UCSD/Course/BGGN213 Found Bioinfor/Week6/12 Genome informati

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 E:/ /UCSD/Course/BGGN213 Found Bioinfor/Week6/12 Genome informati

Info: Writing image file hsa04940.pathview.png

