

Lab_14

Kendrick Nguyen

Table of contents

```
library(DESeq2)
```

Load In the Data

```
metadata<- read.csv("GSE37704_metadata.csv", row.names = 1)
head(metadata)
```

```
              condition
SRR493366 control_sirna
SRR493367 control_sirna
SRR493368 control_sirna
SRR493369      hoxa1_kd
SRR493370      hoxa1_kd
SRR493371      hoxa1_kd
```

```
#import countdata
```

```
countData <- read.csv("GSE37704_featurecounts.csv", 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 | | | | | |

```
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 |

```
# Filter count data where you have 0 read count across all samples.
countData <- countData[(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 = metadata,
                              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)
```

```
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%

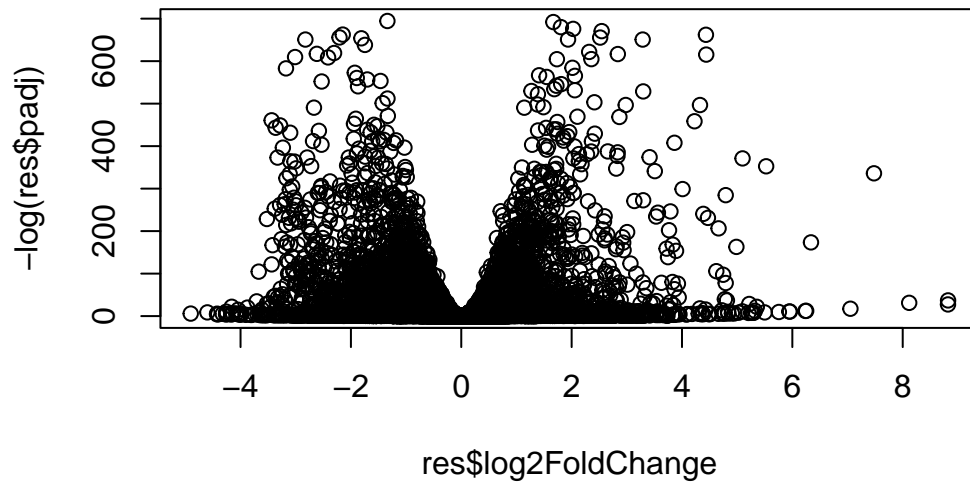
low counts [2] : 1237, 7.7%

(mean count < 0)

[1] see 'cooksCutoff' argument of ?results

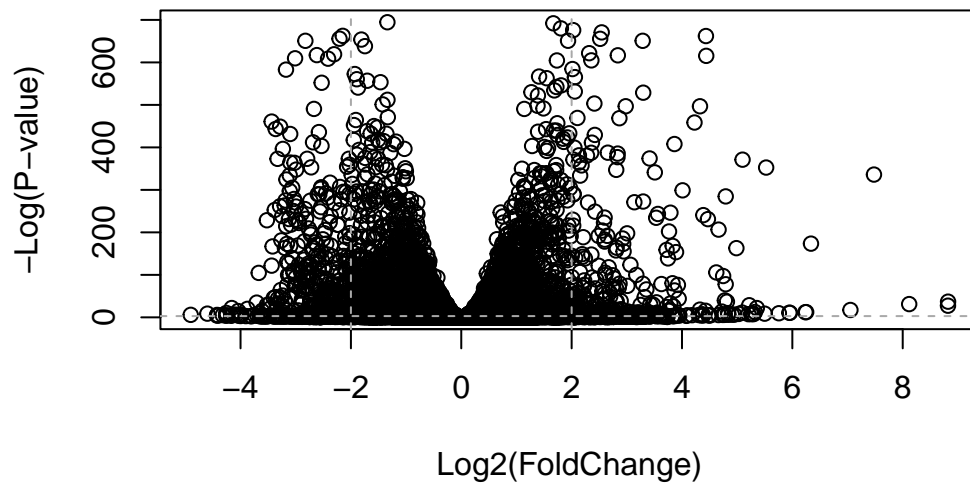
[2] see 'independentFiltering' argument of ?results

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



```
plot( res$log2FoldChange, -log(res$padj),  
      ylab="-Log(P-value)", xlab="Log2(FoldChange)")
```

```
# Add some cut-off lines  
abline(v=c(-2,2), col="darkgray", lty=2)  
abline(h=-log(0.05), col="darkgray", lty=2)
```

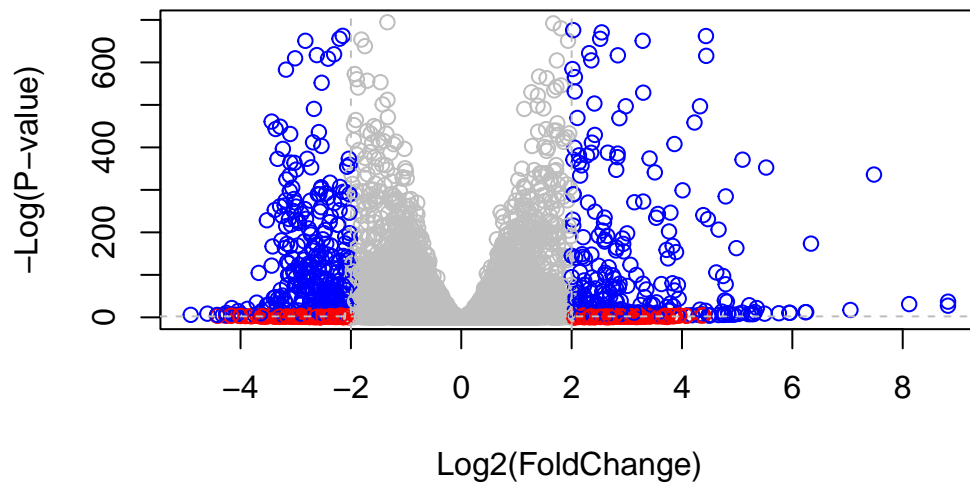


```
mycols <- rep("gray", nrow(res))
mycols[ abs(res$log2FoldChange) > 2 ] <- "red"

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)
```



```
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" | | | | |

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

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

```
res$entrezid = mapIds(org.Hs.eg.db,  
                      keys=row.names(countData),  
                      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(countData),  
                  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 | entrezid | name | |
| | <numeric> | <character> | <character> | <character> | |
| ENSG00000279457 | 6.86555e-01 | WASH9P | 102723897 | WAS protein family h.. | |

| | | | | |
|-----------------|-------------|---------|--------|------------------------|
| 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 .. |

```
library(gage)
library(gageData)
library(pathview)
```

The `gag()` function wants a vector of importance in our case here it will be the fold-change values with associated entrez gene names

```
foldchange<- res$log2FoldChange
names(foldchange)<- res$entrezid
```

```
data(kegg.sets.hs)
```

```
keggres= gage(foldchange, gsets =kegg.sets.hs )
```

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

| | p.geomean | stat.mean |
|--|--------------|-------------|
| hsa04110 Cell cycle | 8.995727e-06 | -4.378644 |
| hsa03030 DNA replication | 9.424076e-05 | -3.951803 |
| hsa05130 Pathogenic Escherichia coli infection | 1.405864e-04 | -3.765330 |
| hsa03013 RNA transport | 1.246882e-03 | -3.059466 |
| hsa03440 Homologous recombination | 3.066756e-03 | -2.852899 |
| hsa04114 Oocyte meiosis | 3.784520e-03 | -2.698128 |
| | p.val | q.val |
| hsa04110 Cell cycle | 8.995727e-06 | 0.001889103 |
| hsa03030 DNA replication | 9.424076e-05 | 0.009841047 |
| hsa05130 Pathogenic Escherichia coli infection | 1.405864e-04 | 0.009841047 |
| hsa03013 RNA transport | 1.246882e-03 | 0.065461279 |
| hsa03440 Homologous recombination | 3.066756e-03 | 0.128803765 |
| hsa04114 Oocyte meiosis | 3.784520e-03 | 0.132458191 |

| | set.size | exp1 |
|--|----------|--------------|
| hsa04110 Cell cycle | 121 | 8.995727e-06 |
| hsa03030 DNA replication | 36 | 9.424076e-05 |
| hsa05130 Pathogenic Escherichia coli infection | 53 | 1.405864e-04 |
| hsa03013 RNA transport | 144 | 1.246882e-03 |
| hsa03440 Homologous recombination | 28 | 3.066756e-03 |
| hsa04114 Oocyte meiosis | 102 | 3.784520e-03 |

hsa04110 cell cycle

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

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

Info: Working in directory E:/BGGN213 Intro to Bioinformatics/R projects/Lab 15/Lab_15_Bioinformatics

Info: Writing image file hsa04110.pathview.png

Picture wasn't showing up, Could not get Pathview to work and only was able to get this resolved at the end of class after an hour of troubleshooting

```
data(go.sets.hs)
data(go.subs.hs)

# Focus on Biological Process subset of GO
gobpsets = go.sets.hs[go.subs.hs$BP]

gobpres = gage(foldchange, 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 | 2.195494e-04 | 3.530241 | 2.195494e-04 |
| GO:0060562 epithelial tube morphogenesis | 5.932837e-04 | 3.261376 | 5.932837e-04 |

| | | | | |
|------------|--------------------------------|--------------|----------|--------------|
| G0:0035295 | tube development | 5.953254e-04 | 3.253665 | 5.953254e-04 |
| | | q.val | set.size | exp1 |
| G0:0007156 | homophilic cell adhesion | 0.1951953 | 113 | 8.519724e-05 |
| G0:0002009 | morphogenesis of an epithelium | 0.1951953 | 339 | 1.396681e-04 |
| G0:0048729 | tissue morphogenesis | 0.1951953 | 424 | 1.432451e-04 |
| G0:0007610 | behavior | 0.2243795 | 427 | 2.195494e-04 |
| G0:0060562 | epithelial tube morphogenesis | 0.3711390 | 257 | 5.932837e-04 |
| G0:0035295 | tube development | 0.3711390 | 391 | 5.953254e-04 |

\$less

| | | p.geomean | stat.mean | p.val |
|------------|-------------------------------|--------------|-----------|--------------|
| G0:0048285 | organelle fission | 1.536227e-15 | -8.063910 | 1.536227e-15 |
| G0:0000280 | nuclear division | 4.286961e-15 | -7.939217 | 4.286961e-15 |
| G0:0007067 | mitosis | 4.286961e-15 | -7.939217 | 4.286961e-15 |
| G0:0000087 | M phase of mitotic cell cycle | 1.169934e-14 | -7.797496 | 1.169934e-14 |
| G0:0007059 | chromosome segregation | 2.028624e-11 | -6.878340 | 2.028624e-11 |
| G0:0000236 | mitotic prometaphase | 1.729553e-10 | -6.695966 | 1.729553e-10 |

| | | q.val | set.size | exp1 |
|------------|-------------------------------|--------------|----------|--------------|
| G0:0048285 | organelle fission | 5.841698e-12 | 376 | 1.536227e-15 |
| G0:0000280 | nuclear division | 5.841698e-12 | 352 | 4.286961e-15 |
| G0:0007067 | mitosis | 5.841698e-12 | 352 | 4.286961e-15 |
| G0:0000087 | M phase of mitotic cell cycle | 1.195672e-11 | 362 | 1.169934e-14 |
| G0:0007059 | chromosome segregation | 1.658603e-08 | 142 | 2.028624e-11 |
| G0:0000236 | mitotic prometaphase | 1.178402e-07 | 84 | 1.729553e-10 |

\$stats

| | | stat.mean | exp1 |
|------------|--------------------------------|-----------|----------|
| G0:0007156 | homophilic cell adhesion | 3.824205 | 3.824205 |
| G0:0002009 | morphogenesis of an epithelium | 3.653886 | 3.653886 |
| G0:0048729 | tissue morphogenesis | 3.643242 | 3.643242 |
| G0:0007610 | behavior | 3.530241 | 3.530241 |
| G0:0060562 | epithelial tube morphogenesis | 3.261376 | 3.261376 |
| G0:0035295 | tube development | 3.253665 | 3.253665 |

```
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
```