# Class 14: RNASeq Mini Project

Lourd Yakoob (PID: A18543409)

## Table of contents

## Background

Here we work through a complete RNASeq analysis project. The input data comes from a knock-down experiment of a HOX gene.

## Data Import

Reading the `counts` and `metadata` CSV files

```
metaFile <- "GSE37704_metadata.csv"
countFile <- "GSE37704_featurecounts.csv"
colData = read.csv(metaFile, row.names=1)
countData = read.csv(countFile, row.names=1)
```

Check on data structure

Some book-keeping is required as there looks to be a mis-match between metadata and counts columns.

Looks like we need to get rid of the first "length" column of our `counts` object.

```
countData <- read.csv(countFile, row.names = 1)
countData[] <- lapply(countData, function(x) as.numeric(as.character(x)))
cleancounts <- as.matrix(countData[, -1])
nonzero_counts <- cleancounts[rowSums(cleancounts) > 0, ]
head(nonzero_counts)
```

|                 | 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        |

## DESeq Analysis

Load the package

```
library(DESeq2)
```

Setup DESeq

```
dds = DESeqDataSetFromMatrix(countData=nonzero_counts,
                             colData=colData,
                             design=~condition)
```

```
Warning in DESeqDataSet(se, design = design, ignoreRank): some variables in
design formula are characters, converting to factors
```

Run DESeq

```
dds <- DESeq(dds)
```
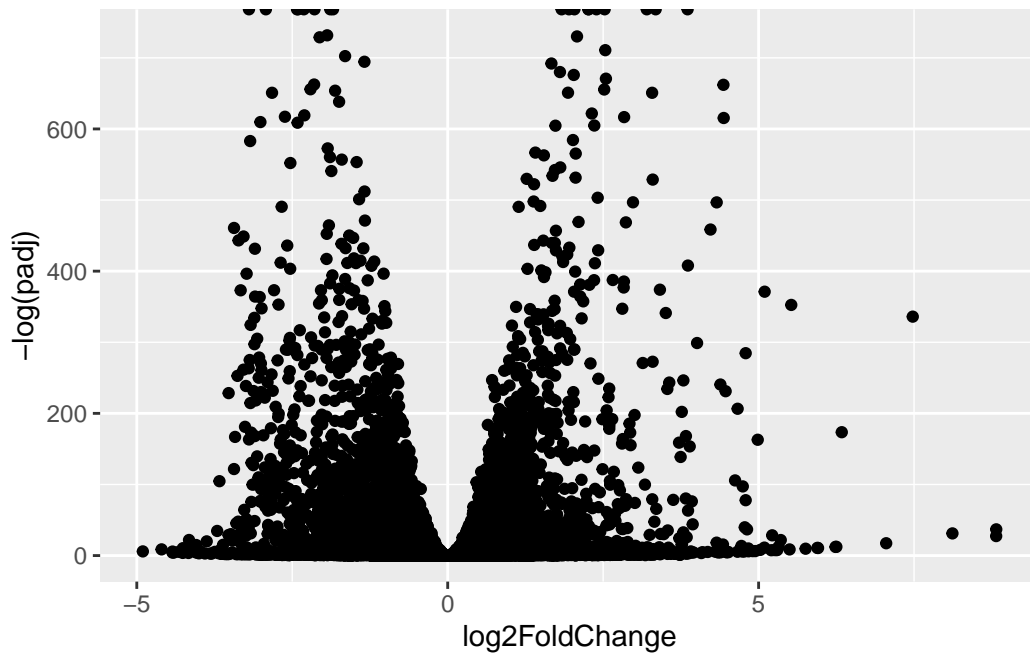
Get results

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

## Data Visualization

Volcano Plot

```
library(ggplot2)

ggplot(res) +
  aes(log2FoldChange, -log(padj)) +
  geom_point()
```

```
Warning: Removed 1237 rows containing missing values or values outside the scale range
(`geom_point()`).
```
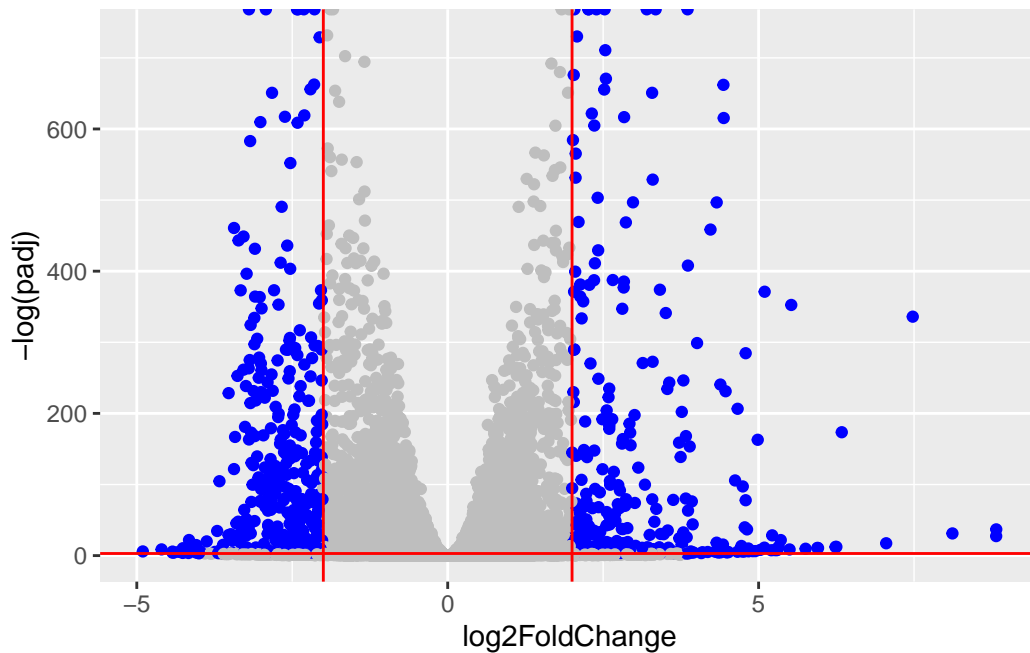


Add threshold lines for fold-change and P-valye and color our subset of genes that make these
threshold cut-offs in the plot.

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


ggplot(res) +
  aes(log2FoldChange, -log(padj)) +
  geom_point(col = mycols) +
  geom_vline(xintercept =c(-2,2), col="red") +
  geom_hline(yintercept = -log(0.05), col ="red")
```

```
Warning: Removed 1237 rows containing missing values or values outside the scale range
(`geom_point()`).
```



## Add Annotation

Add gene symbols and entrez ids

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

4

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

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

```
res$entrez = mapIds(org.Hs.eg.db,
                    keys = row.names(res),
                    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(res),
                  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      entrez                     name
                     <numeric> <character> <character>            <character>
ENSG00000279457 6.86555e-01          NA          NA                       NA
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 ..
```

Let's Reorder results by adjusted p-value and save as CSV file:

```
res <- res[order(res$padj), ]

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

## Pathway Analysis

### KEGG pathways

- run this in console: BiocManager::install( c("pathview", "gage", "gageData") ) *

Run gage analysis w/ KEGG

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

We need a named vector of fold-change value as input for gage.

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

```
     1266        54855        1465        2034        2150        6659
-2.422719   3.201955  -2.313738  -1.888019   3.344508   2.392288
```

```r
data("kegg.sets.hs")
keggres = gage(foldchanges, gsets=kegg.sets.hs)
```
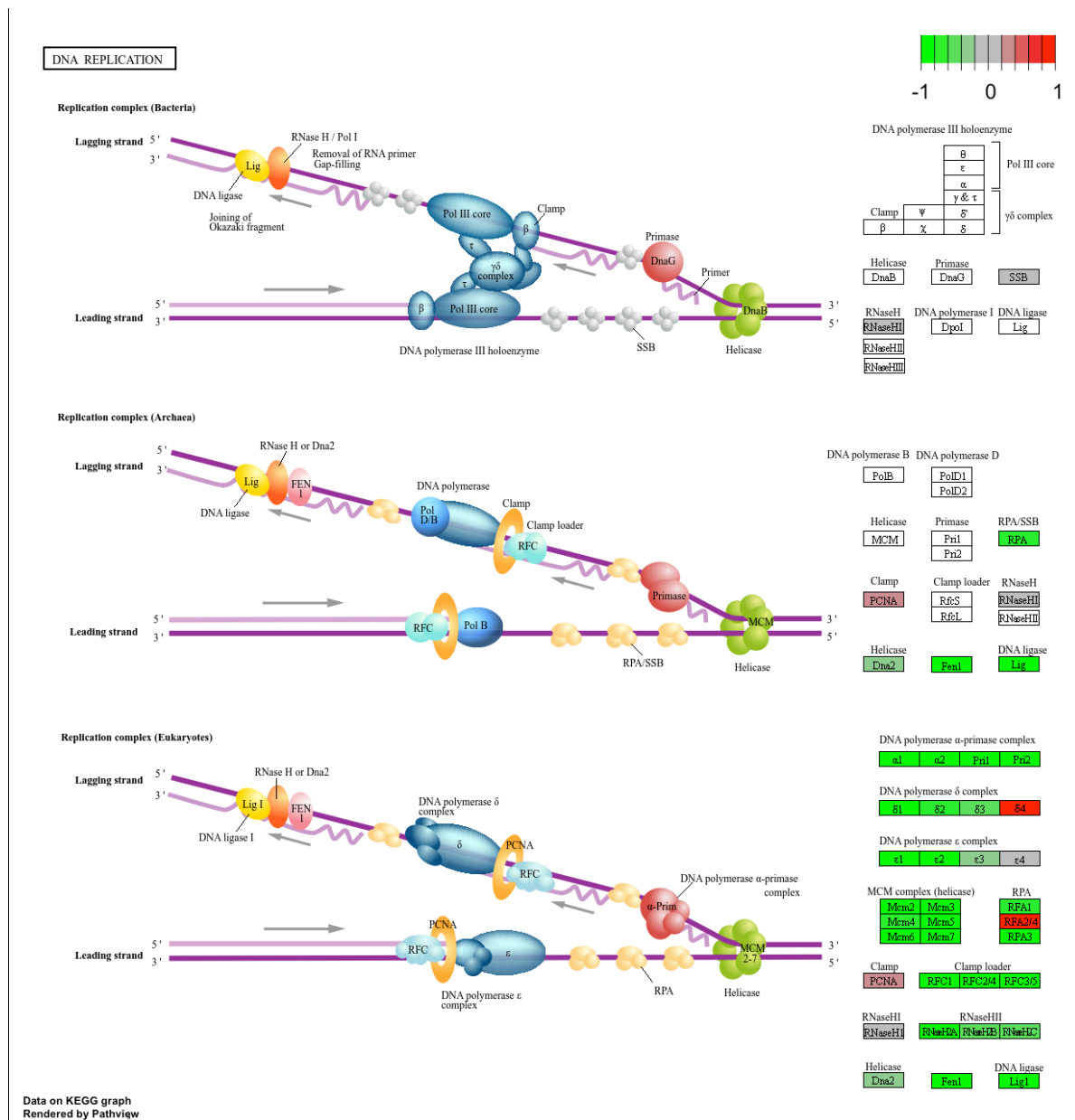
```r
head(keggres$less, 2)
```

```
                         p.geomean stat.mean       p.val       q.val
hsa04110 Cell cycle      8.995727e-06 -4.378644 8.995727e-06 0.001889103
hsa03030 DNA replication 9.424076e-05 -3.951803 9.424076e-05 0.009841047
                         set.size         exp1
hsa04110 Cell cycle           121 8.995727e-06
hsa03030 DNA replication       36 9.424076e-05
```

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

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

Info: Working in directory /Users/lourdyakoob/Desktop/Class 14

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

```
pathview(pathway.id = "hsa03030", gene.data = foldchanges)
```

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

Info: Working in directory /Users/lourdyakoob/Desktop/Class 14

Info: Writing image file hsa03030.pathview.png

8

## GO terms

Same analysis but using GO genesets rather than KEGG.

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

```
gobpsets = go.sets.hs[go.subs.hs$BP]

gobpres = gage(foldchanges, 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                      1.925222e-04  3.565432 1.925222e-04
GO:0060562 epithelial tube morphogenesis 5.932837e-04  3.261376 5.932837e-04
GO:0035295 tube development              5.953254e-04  3.253665 5.953254e-04
                                             q.val set.size        exp1
GO:0007156 homophilic cell adhesion      0.1951953      113 8.519724e-05
GO:0002009 morphogenesis of an epithelium 0.1951953      339 1.396681e-04
GO:0048729 tissue morphogenesis          0.1951953      424 1.432451e-04
GO:0007610 behavior                      0.1967577      426 1.925222e-04
GO:0060562 epithelial tube morphogenesis 0.3565320      257 5.932837e-04
GO:0035295 tube development              0.3565320      391 5.953254e-04


$less
                                            p.geomean stat.mean        p.val
GO:0048285 organelle fission            1.536227e-15 -8.063910 1.536227e-15
GO:0000280 nuclear division             4.286961e-15 -7.939217 4.286961e-15
GO:0007067 mitosis                      4.286961e-15 -7.939217 4.286961e-15
GO:0000087 M phase of mitotic cell cycle 1.169934e-14 -7.797496 1.169934e-14
GO:0007059 chromosome segregation       2.028624e-11 -6.878340 2.028624e-11
GO:0000236 mitotic prometaphase         1.729553e-10 -6.695966 1.729553e-10
                                             q.val set.size        exp1
GO:0048285 organelle fission            5.841698e-12      376 1.536227e-15
GO:0000280 nuclear division             5.841698e-12      352 4.286961e-15
GO:0007067 mitosis                      5.841698e-12      352 4.286961e-15
GO:0000087 M phase of mitotic cell cycle 1.195672e-11      362 1.169934e-14
GO:0007059 chromosome segregation       1.658603e-08      142 2.028624e-11
GO:0000236 mitotic prometaphase         1.178402e-07       84 1.729553e-10


$stats
                                       stat.mean     exp1
GO:0007156 homophilic cell adhesion     3.824205 3.824205
```

```
GO:0002009 morphogenesis of an epithelium  3.653886 3.653886
GO:0048729 tissue morphogenesis            3.643242 3.643242
GO:0007610 behavior                        3.565432 3.565432
GO:0060562 epithelial tube morphogenesis   3.261376 3.261376
GO:0035295 tube development                3.253665 3.253665
```

```
head(gobpres$less, 4)
```

```
                                           p.geomean stat.mean        p.val
GO:0048285 organelle fission            1.536227e-15 -8.063910 1.536227e-15
GO:0000280 nuclear division             4.286961e-15 -7.939217 4.286961e-15
GO:0007067 mitosis                      4.286961e-15 -7.939217 4.286961e-15
GO:0000087 M phase of mitotic cell cycle 1.169934e-14 -7.797496 1.169934e-14
                                             q.val set.size         exp1
GO:0048285 organelle fission            5.841698e-12      376 1.536227e-15
GO:0000280 nuclear division             5.841698e-12      352 4.286961e-15
GO:0007067 mitosis                      5.841698e-12      352 4.286961e-15
GO:0000087 M phase of mitotic cell cycle 1.195672e-11      362 1.169934e-14
```

**Reactome**

Lots of folkslike the reactome web interface. You can also run this as an R fucntion but lets look at the website first https://reactome.org/

The website wants a text file with one gene symbol per line of the genes ypu want to map to the pathways.

```
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"
```

and write out to a file:

```
write.table(sig_genes, file="significant_genes.txt", row.names=FALSE, col.names=FALSE, quote=
```

**Save Our Results**

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