

# Class 14: RNASeq mini Project

Ashley Mazon (PID: A17478903)

## 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

```
counts <- read.csv("GSE37704_featurecounts.csv",row.names = 1)
metadata <- read.csv("GSE37704_metadata.csv")
```

```
head(counts)
```

	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					

```
head(metadata)
```

```
      id      condition
1 SRR493366 control_sirna
2 SRR493367 control_sirna
3 SRR493368 control_sirna
4 SRR493369      hoxa1_kd
5 SRR493370      hoxa1_kd
6 SRR493371      hoxa1_kd
```

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

```
ncol(counts)
```

```
[1] 7
```

```
nrow(metadata)
```

```
[1] 6
```

We need to remove the first column for count

```
cleancounts <- counts[,-1]
head(cleancounts)
```

	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

```
all(colnames(cleancounts) == metadata$id)
```

```
[1] TRUE
```

## Remove zero count genes

There are lost of genes with zero counts. We can remove these from further analysis

```
head(cleancounts)
```

	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

```
to.keep.inds <- rowSums(cleancounts) >0  
nonzero_counts <- cleancounts[to.keep.inds,]
```

##DESeq Analysis

Load the package

```
library(DESeq2)
```

Warning: package 'matrixStats' was built under R version 4.5.2

Setup DESeq

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

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

run DESeq

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

estimating size factors

estimating dispersions

gene-wise dispersion estimates

mean-dispersion relationship

final dispersion estimates

fitting model and testing

get results

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

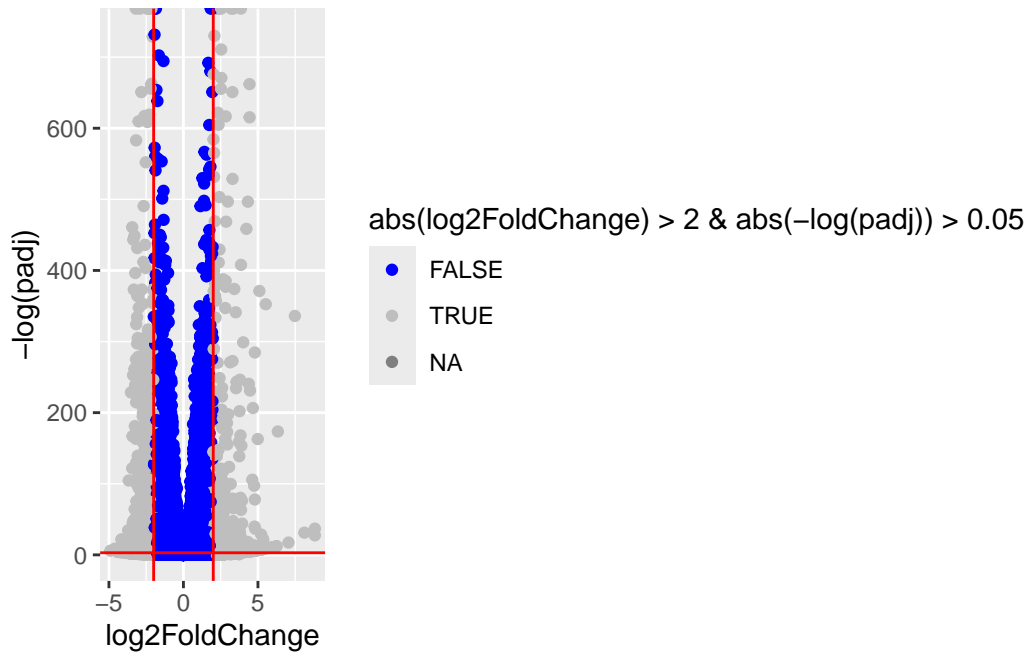
## Data Visualization

Volcano plot

```
library(ggplot2)

ggplot(res) +
  aes(log2FoldChange, -log(padj),
      color = abs(log2FoldChange) > 2 & abs(-log(padj)) > 0.05) +
  geom_point() +
  scale_color_manual(values = c("blue", "grey")) +
  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 threshold lines for fold-change and p-value and color our subset of genes that make these threshold cut-offs in the plot.

## Add annotation

Add gene symbols and entrez ids

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

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

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

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

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

## Pathway Analysis

### KEGG Pathways

Run gage analysis with KEGG

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

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

# Focus on signaling and metabolic pathways only
kegg.sets.hs = kegg.sets.hs[sigmet.idx.hs]

# Examine the first 3 pathways
head(kegg.sets.hs, 3)
```

```
$`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"
```

```
$`hsa00230 Purine metabolism`
[1] "100"    "10201" "10606" "10621" "10622" "10623" "107"    "10714"
```

[9]	"108"	"10846"	"109"	"111"	"11128"	"11164"	"112"	"113"
[17]	"114"	"115"	"122481"	"122622"	"124583"	"132"	"158"	"159"
[25]	"1633"	"171568"	"1716"	"196883"	"203"	"204"	"205"	"221823"
[33]	"2272"	"22978"	"23649"	"246721"	"25885"	"2618"	"26289"	"270"
[41]	"271"	"27115"	"272"	"2766"	"2977"	"2982"	"2983"	"2984"
[49]	"2986"	"2987"	"29922"	"3000"	"30833"	"30834"	"318"	"3251"
[57]	"353"	"3614"	"3615"	"3704"	"377841"	"471"	"4830"	"4831"
[65]	"4832"	"4833"	"4860"	"4881"	"4882"	"4907"	"50484"	"50940"
[73]	"51082"	"51251"	"51292"	"5136"	"5137"	"5138"	"5139"	"5140"
[81]	"5141"	"5142"	"5143"	"5144"	"5145"	"5146"	"5147"	"5148"
[89]	"5149"	"5150"	"5151"	"5152"	"5153"	"5158"	"5167"	"5169"
[97]	"51728"	"5198"	"5236"	"5313"	"5315"	"53343"	"54107"	"5422"
[105]	"5424"	"5425"	"5426"	"5427"	"5430"	"5431"	"5432"	"5433"
[113]	"5434"	"5435"	"5436"	"5437"	"5438"	"5439"	"5440"	"5441"
[121]	"5471"	"548644"	"55276"	"5557"	"5558"	"55703"	"55811"	"55821"
[129]	"5631"	"5634"	"56655"	"56953"	"56985"	"57804"	"58497"	"6240"
[137]	"6241"	"64425"	"646625"	"654364"	"661"	"7498"	"8382"	"84172"
[145]	"84265"	"84284"	"84618"	"8622"	"8654"	"87178"	"8833"	"9060"
[153]	"9061"	"93034"	"953"	"9533"	"954"	"955"	"956"	"957"
[161]	"9583"	"9615"						

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

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

<NA>	148398	26155	339451	84069	84808
0.17925708	0.42645712	-0.69272046	0.72975561	0.04057653	0.54281049

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

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

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

Info: Working in directory C:/Users/mazon/OneDrive/Documents/BIMM143 Data Sets/class14

Info: Writing image file hsa04110.pathview.png





## Reactome

Lots of folks like the reactome web interface. You can also run this as an R function but let's look at the website first. < <https://reactome.org/>>

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

```
sig_genes <- res[res$padj <= 0.05 & !is.na(res$padj),] $symbol  
head(sig_genes)
```

```
ENSG00000187634 ENSG00000188976 ENSG00000187961 ENSG00000188290 ENSG00000187608  
      "SAMD11"      "NOC2L"      "KLHL17"      "HES4"      "ISG15"  
ENSG00000188157  
      "AGRN"
```

and write out to a file

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

## Save our results

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