

Skeleton Key for RNAseq analysis

Written By: Ciera Martinez

See README.md for more detailed instructions of how to use script

Run the `render()` function below and everything will be run with report at end.

```
library(rmarkdown)
render("skeletonDE.Rmd", "pdf_document", output_file = paste(sample1,"_",sample2,"_", "DE.pdf",sep=""))
```

Analysis

libraries

```
library(edgeR)
library(yaml)
```

Read in YAML guide

```
yamls <- yaml.load_file("de.yaml")
```

This part assigns your YMAL to a object in R. This will be used throughout the script to specify which sample types you are comparing.

```
sample1 <- yamls$sample1
sample2 <- yamls$sample2
```

```
sample1
```

```
## [1] "wtaother"
```

```
sample2
```

```
## [1] "wtbother"
```

Read in Data

Read in raw count data per gene.

```
counts <- read.delim("../requisiteData/sam2countsResults.tsv",row.names=1)

#check the file
head(counts)
colnames(counts)
#need to convert NA to 0 counts
counts[is.na(counts)] <- 0
```

Subset DE experiment

Start by subsetting the particular treatments which are being compared.

```
colnames(counts)
```

```
## [1] "tf2ambr1"      "tf2ambr3"      "tf2ambr4"      "tf2ambr6"
## [5] "tf2aother1"    "tf2aother2"    "tf2aother4"    "tf2aother7"
## [9] "tf2bmbbr2"     "tf2bmbbr5"     "tf2bmbbr6"     "tf2bother1"
## [13] "tf2bother3"    "tf2bother4"    "tf2bother6"    "tf2cmbr1.4"
## [17] "tf2cmbr3"      "tf2cmbr6"      "tf2cmbr7"      "tf2cother2"
## [21] "tf2cother5"    "tf2cother6"    "tf2cother7"    "wtambr2"
## [25] "wtambr4"       "wtambr5"       "wtaother1"     "wtaother5"
## [29] "wtaother6"     "wtaother7"     "wtaother8"     "wtbmbbr2"
## [33] "wtbmbbr3"      "wtbmbbr6"      "wtbmbbr8"      "wtbother1.4"
## [37] "wtbother3"     "wtbother5"     "wtbother8"     "wtcmbr10"
## [41] "wtcmbr1.4.6"   "wtcmbr2"       "wtcmbr3"       "wtcmbr7"
## [45] "wtcmbr9"       "wtcother1.3.4" "wtcother2"     "wtcother6"
```

```
counts1 <- counts[,grep(sample1, colnames(counts), value = TRUE)]
count1Len <- length(colnames(counts1)) #used in to specify library group in next step.
```

```
counts2 <- counts[,grep(sample2, colnames(counts), value = TRUE)]
count2Len <- length(colnames(counts2)) #used to specify library group in next step.
```

```
counts <- cbind(counts1, counts2)
```

```
head(counts)
```

```
##           wtaother1 wtaother5 wtaother6 wtaother7 wtaother8
## Solyc00g005040.2.1         1         1         1         0         2
## Solyc00g005050.2.1        17        16         9         2         3
## Solyc00g005060.1.1         0         0         0         0         2
## Solyc00g005070.1.1         8         6         5         5         6
## Solyc00g005080.1.1        18        37         6        10         7
## Solyc00g005150.1.1         2         5         0         0         2
##           wtbother1.4 wtbother3 wtbother5 wtbother8
## Solyc00g005040.2.1         0         8         0         3
## Solyc00g005050.2.1         0        25         0        14
## Solyc00g005060.1.1         0         0         0         0
## Solyc00g005070.1.1         0         6         2         4
## Solyc00g005080.1.1         0        29         0        11
## Solyc00g005150.1.1         0         2         0         2
```

Add column specifying library Group

Make a vector called group that will be used to make a new column named group to identify library region type.

```
group <- c(rep(sample1, count1Len), rep(sample2, count2Len))
d <- DGEList(counts=counts,group=group)
```

Check to see if the group column matches your sample name and they are appropriate.

```
d$samples
```

```
##           group lib.size norm.factors
## wtaother1  wtaother  929017          1
## wtaother5  wtaother 1555921          1
## wtaother6  wtaother  498294          1
## wtaother7  wtaother  479003          1
## wtaother8  wtaother  510148          1
## wtbother1.4 wtbother   1421          1
## wtbother3  wtbother 1076939          1
## wtbother5  wtbother  200587          1
## wtbother8  wtbother  499487          1
```

Differential expression using edgeR

Make sure there is full understanding on each edgeR command being used. The manual is amazing so read it *before* running the DE analysis below [edgeR manual](#).

```
cpm.d <- cpm(d) #counts per mutant
d <- d[rowSums(cpm.d>5)>=3,] #This might be a line to adjust. It is removing genes with low counts.
d <- estimateCommonDisp(d,verbose=T)
```

```
## Disp = 0.3386 , BCV = 0.5818
```

```
d <- calcNormFactors(d)
d <- estimateCommonDisp(d)

DEtest <- exactTest(d,pair=c(sample1,sample2))
head(DEtest$table)
```

```
##           logFC logCPM   PValue
## Solyc00g005050.2.1  0.62285  4.611 7.226e-01
## Solyc00g005070.1.1 -0.19406  3.968 7.861e-01
## Solyc00g005080.1.1 -0.07833  4.603 6.041e-01
## Solyc00g005440.1.1 -0.37015  5.247 4.311e-01
## Solyc00g005840.2.1  3.56750  7.444 1.872e-07
## Solyc00g006470.1.1 -1.01326  9.987 7.075e-02
```

```
results <- topTags(DEtest, n=Inf)
head(results)
```

```
## Comparison of groups: wtbother-wtaother
##           logFC logCPM   PValue   FDR
## Solyc04g074380.2.1 10.013 11.898 1.005e-31 1.559e-27
## Solyc01g014280.2.1  9.015  9.148 8.652e-25 6.714e-21
## Solyc04g074390.2.1  8.239 10.304 1.684e-24 8.713e-21
## Solyc10g078540.1.1  7.619 13.383 7.914e-24 3.071e-20
## Solyc02g076780.2.1  6.601 10.530 2.781e-19 8.632e-16
## Solyc09g059140.1.1  6.627  8.461 1.211e-17 3.132e-14
```

```
dim(results$table)
```

```
## [1] 15522      4
```

```
sum(results$table$FDR<.05) # How many are DE genes?
```

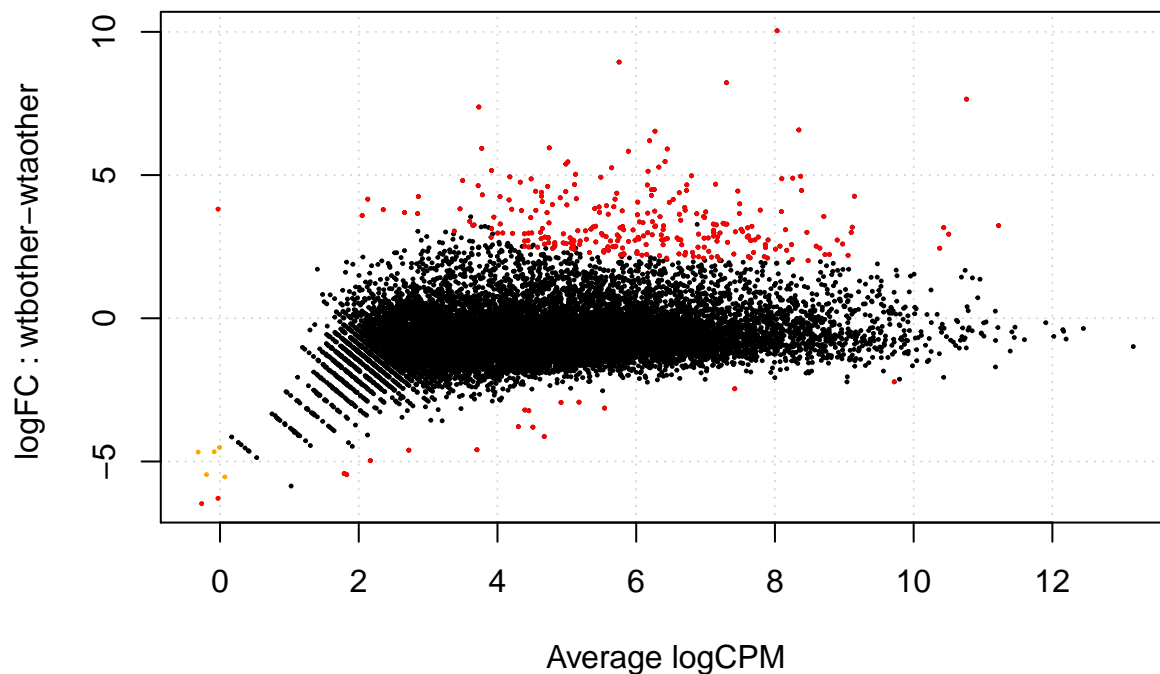
```
## [1] 291
```

```
summary(decideTestsDGE(DEtest,p.value=.05))
```

```
##      [,1]  
## -1      18  
##  0    15231  
##  1      273
```

```
sig.genes <- rownames(results$table[results$table$FDR<0.05,]) # outputs just significant gene names
```

```
plotSmea(d,de.tags=sig.genes)
```



Subset by all the genes with a significant FDR score.

```
results.sig <- subset(results$table, results$table$FDR < 0.05)
```

```
dim(results.sig)
```

What are the genes that are misexpressed? For this we need to add some annotation.

Essentially we are merging two annotations files to 1.) only sig genes 2.) all genes

```

annotation1<- read.delim("../requisiteData/ITAG2.3_all_Arabidopsis_ITAG_annotations.tsv", header=FALSE)
colnames(annotation1) <- c("ITAG", "SGN_annotation")
annotation2<- read.delim("../requisiteData/ITAG2.3_all_Arabidopsis_annotated.tsv")
annotation <- merge(annotation1,annotation2, by = "ITAG")
head(annotation)

```

```

##          ITAG
## 1 Solyc00g005000.2.1
## 2 Solyc00g005040.2.1
## 3 Solyc00g005050.2.1
## 4 Solyc00g005080.1.1
## 5 Solyc00g005900.1.1
## 6 Solyc00g006490.2.1
##
## 1          Aspartic proteinase nepenthesin I (AHRD V1 ***-
## 2          Potassium channel (AHRD V1 ***- DOEM91_9ROSI
## 3
## 4
## 5 Oxygen-evolving enhancer protein 1, chloroplastic (AHRD V1 ***- PSBO_SOLTU); contains Interpro dom
## 6 Serine/threonine-protein phosphatase 6 regulatory subunit 3 (AHRD V1 ***- SAPS3_HUMAN); contain
##      AGI symbol
## 1 AT3G20015   <NA>
## 2 AT5G46240   KAT1
## 3 AT5G11680   <NA>
## 4 ATCG01280   YCF2.2
## 5 AT5G66570   MSP-1
## 6 AT1G07990   <NA>
##
## 1
## 2
## 3
## 4
## 5
## 6 SIT4 phosphatase-associated family protein; similar to SIT4 phosphatase-associated family protein
## X..identity alignment.length e.value bit.score percent.query.align
## 1      63.76          447 7e-148      520          89.94
## 2      66.02          103 2e-37       150          85.71
## 3      76.96          204 1e-88       322          98.98
## 4      91.25           80 2e-38       153          79.80
## 5      69.62           79 4e-26       112          78.79
## 6      61.92          856 0e+00       979          99.77

```

```

#Making the only significant gene table
results.sig$ITAG <- rownames(results.sig) #change row.names to ITAG for merging
results.sig.annotated <- merge(results.sig,annotation,by = "ITAG", all.x=TRUE) #This is merging to only

#Making all table

results$table$ITAG <- rownames(results$table)
results.all.annotated <- merge(results$table, annotation,by = "ITAG")

```

Write table with results.

```
write.table(results.all.annotated, file=paste(sample1,"_",sample2,"_", "DE_all.txt", sep=""), sep="\t", row
write.table(results.sig.annotated, file=paste(sample1,"_",sample2,"_", "DE_sig.txt", sep=""), sep="\t", row
```

Now run the script below for a full `knitr` report of what was run and leave this report in the folder that the analysis was done with output files.