

Skeleton Key for RNAseq analysis

Written By: Ciera Martinez

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

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] "tf2aother"
```

```
sample2
```

```
## [1] "tf2cother"
```

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"     "tf2bmbbr8"
## [13] "tf2bmbbr3"     "tf2bmbbr4"     "tf2bmbbr6"     "tf2bmbbr7"
## [17] "tf2bmbbr5"     "tf2bmbbr6"     "tf2bmbbr7"     "tf2bmbbr8"
## [21] "tf2bmbbr6"     "tf2bmbbr7"     "tf2bmbbr8"     "wtambr2"
## [25] "wtambr4"       "wtambr5"       "wtambr6"       "wtambr7"
## [29] "wtambr6"       "wtambr7"       "wtambr8"       "wtambr9"
## [33] "wtbmbbr3"      "wtbmbbr6"      "wtbmbbr8"      "wtbmbbr9"
## [37] "wtbmbbr3"      "wtbmbbr5"      "wtbmbbr8"      "wtbmbbr9"
## [41] "wtbmbbr1.4.6"  "wtbmbbr2"      "wtbmbbr3"      "wtbmbbr7"
## [45] "wtbmbbr9"      "wtbmbbr1.3.4" "wtbmbbr2"      "wtbmbbr6"
```

```
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)
```

```
##          tf2aother1 tf2aother2 tf2aother4 tf2aother7 tf2cother2
## Solyc00g005040.2.1      0         1         0         2         3
## Solyc00g005050.2.1      0         2         3         0         4
## Solyc00g005060.1.1      0         0         0         0         1
## Solyc00g005070.1.1      3         6        33         2         4
## Solyc00g005080.1.1      0        12        10         3         9
## Solyc00g005150.1.1      0         0         2         1         1
##          tf2cother5 tf2cother6 tf2cother7
## Solyc00g005040.2.1      8         4         1
## Solyc00g005050.2.1     10        16        12
## Solyc00g005060.1.1      1         2         1
## Solyc00g005070.1.1     11         5         5
## Solyc00g005080.1.1     21        14         3
## Solyc00g005150.1.1      6         1         0
```

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
## tf2aother1 tf2aother  263117          1
## tf2aother2 tf2aother  698710          1
## tf2aother4 tf2aother  792325          1
## tf2aother7 tf2aother  142504          1
## tf2cother2 tf2cother  723602          1
## tf2cother5 tf2cother 1216379          1
## tf2cother6 tf2cother  838942          1
## tf2cother7 tf2cother  676969          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 = 0.5435 , BCV = 0.7372
```

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

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

```
##           logFC logCPM PValue
## Solyc00g005050.2.1  2.0848  3.302 0.02842
## Solyc00g005070.1.1 -1.7318  4.321 0.05302
## Solyc00g005080.1.1 -0.1499  4.116 0.90260
## Solyc00g005160.1.1 -1.2478  2.898 0.29472
## Solyc00g005440.1.1  0.1200  4.460 0.96820
## Solyc00g005840.2.1 -1.0769  6.351 0.17492
```

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

```
## Comparison of groups: tf2cother-tf2aother
##           logFC logCPM PValue FDR
## Solyc11g064800.1.1 -6.279  6.315 1.110e-09 1.295e-05
## Solyc12g009110.1.1  5.653  8.464 1.726e-09 1.295e-05
## Solyc03g112640.2.1 -5.896  6.636 3.623e-09 1.813e-05
## Solyc06g072480.1.1  6.447  5.897 9.779e-09 3.670e-05
## Solyc05g021410.1.1 -6.010  5.353 3.313e-08 9.947e-05
## Solyc03g118770.2.1 -5.473  5.443 6.848e-08 1.713e-04
```

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

```
## [1] 15012      4
```

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

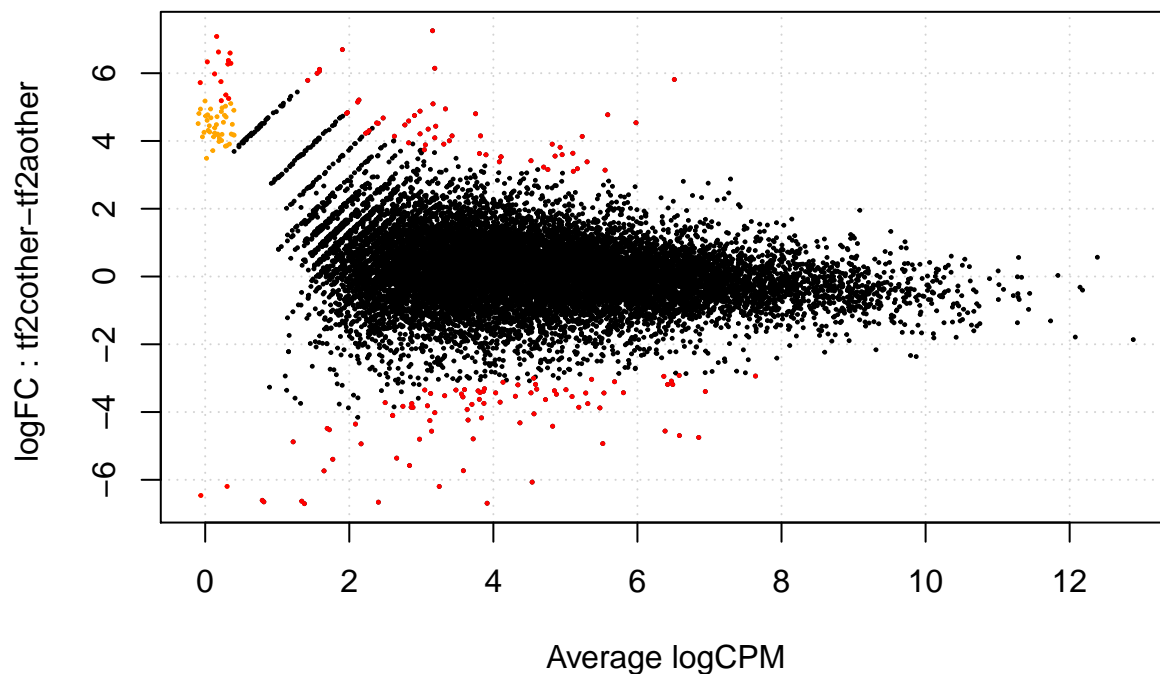
```
## [1] 156
```

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

```
##      [,1]  
## -1      86  
##  0    14856  
##  1       70
```

```
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)
```

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")

#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") #This is merging to only sig genes

#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.

```

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

```