

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] "tf2ambr"
```

```
sample2
```

```
## [1] "tf2aother"
```

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] "tf2bmbbr3"     "tf2bmbbr6"     "tf2bmbbr7"     "tf2bmbbr8"
## [21] "tf2bmbbr5"     "tf2bmbbr6"     "tf2bmbbr7"     "tf2bmbbr8"
## [25] "tf2bmbbr6"     "tf2bmbbr7"     "tf2bmbbr8"     "wtambr2"
## [29] "wtambr4"       "wtambr5"       "wtambr6"       "wtambr7"
## [33] "wtambr6"       "wtambr7"       "wtambr8"       "wtambr9"
## [37] "wtambr8"       "wtambr9"       "wtambr10"      "wtambr11"
## [41] "wtambr9"       "wtambr10"      "wtambr11"      "wtambr12"
## [45] "wtambr11"      "wtambr12"      "wtambr13"      "wtambr14"
```

```
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)
```

```
##          tf2ambr1 tf2ambr3 tf2ambr4 tf2ambr6 tf2aother1
## Solyc00g005040.2.1      12         0         3         12         0
## Solyc00g005050.2.1      33         1        14         17         0
## Solyc00g005060.1.1         1         5         1         1         0
## Solyc00g005070.1.1      14        22        23         5         3
## Solyc00g005080.1.1      19         2        25        32         0
## Solyc00g005150.1.1         3         0         0         4         0
##          tf2aother2 tf2aother4 tf2aother7
## Solyc00g005040.2.1         1         0         2
## Solyc00g005050.2.1         2         3         0
## Solyc00g005060.1.1         0         0         0
## Solyc00g005070.1.1         6        33         2
## Solyc00g005080.1.1        12        10         3
## Solyc00g005150.1.1         0         2         1
```

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
## tf2ambr1      tf2ambr  1313540         1
## tf2ambr3      tf2ambr   91726         1
## tf2ambr4      tf2ambr  1438416         1
## tf2ambr6      tf2ambr  1088653         1
## tf2aother1 tf2aother   263117         1
## tf2aother2 tf2aother   698710         1
## tf2aother4 tf2aother   792325         1
## tf2aother7 tf2aother   142504         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 = 0.475 , BCV = 0.6892
```

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

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

```
##           logFC logCPM  PValue
## Solyc00g005040.2.1 -1.5139  3.083 0.497449
## Solyc00g005050.2.1 -2.4605  3.595 0.006189
## Solyc00g005070.1.1 -1.2337  5.449 0.092153
## Solyc00g005080.1.1 -0.5418  4.315 0.490062
## Solyc00g005160.1.1  1.1651  3.528 0.418041
## Solyc00g005440.1.1 -0.9575  4.881 0.248082
```

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

```
## Comparison of groups:  tf2aother-tf2ambr
##           logFC logCPM  PValue      FDR
## Solyc03g112640.2.1  5.366  6.648 3.965e-09 3.058e-05
## Solyc12g014030.1.1  5.600  7.996 4.177e-09 3.058e-05
## Solyc09g009620.1.1  7.982  4.538 5.688e-08 2.776e-04
## Solyc11g062200.1.1 -4.717  6.186 7.611e-08 2.786e-04
## Solyc08g081230.1.1  4.518  6.616 1.348e-07 3.573e-04
## Solyc10g051200.1.1  4.630  6.144 1.464e-07 3.573e-04
```

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

```
## [1] 14644      4
```

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

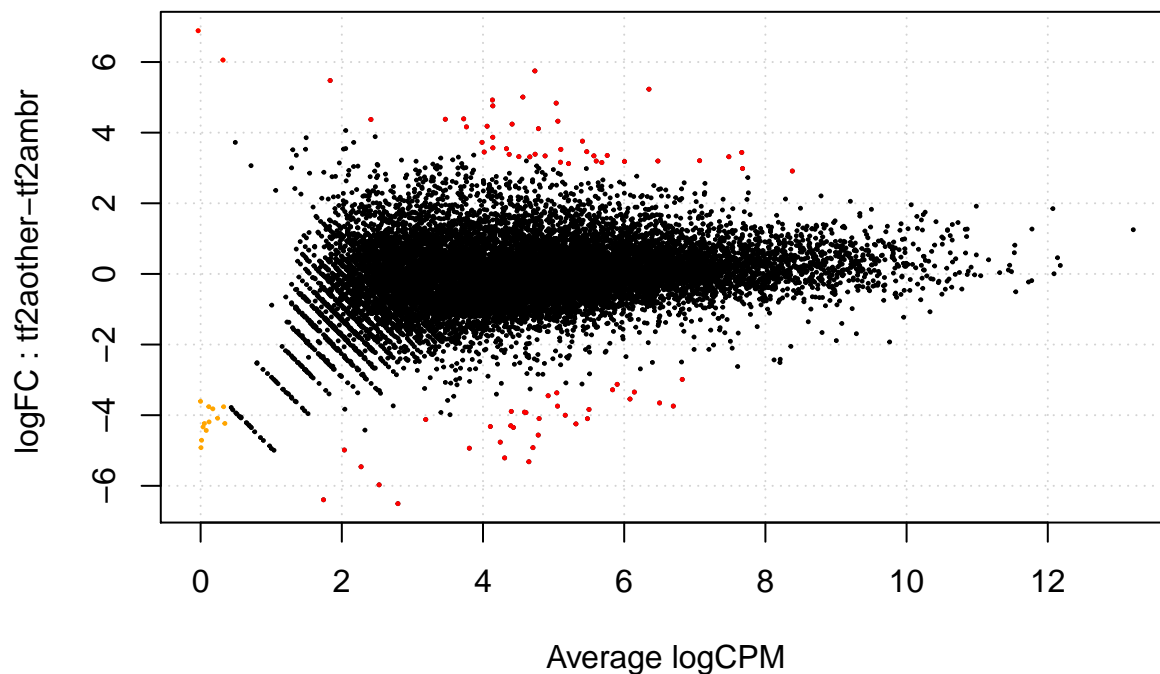
```
## [1] 76
```

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

```
##      [,1]  
## -1      33  
##  0    14568  
##  1      43
```

```
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=""))

```