

Skeleton Key for RNAseq analysis

Written By: Ciera Martinez

libraries

```
library(edgeR)
```

Read in YAML guide

```
library(yaml)
yaml1 <- yaml.load_file("./de.yml")
```

```
sample1 <- yaml1$sample1
sample2 <- yaml1$sample2
```

```
sample1
```

```
## [1] "wtbmbr"
```

```
sample2
```

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

Read in Data

Read in raw count data per gene.

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

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

Subset per DE experiment

I am going to start by subsetting the particular treatments I am looking at.

```
colnames(counts)
```

```
## [1] "tf2ambr1"      "tf2ambr3"      "tf2ambr4"      "tf2ambr6"
## [5] "tf2aother1"    "tf2aother2"    "tf2aother4"    "tf2aother7"
## [9] "tf2bmr2"       "tf2bmr5"       "tf2bmr6"       "tf2bmr1"
## [13] "tf2bmr3"       "tf2bmr4"       "tf2bmr6"       "tf2cmbr1.4"
## [17] "tf2cmbr3"      "tf2cmbr6"      "tf2cmbr7"      "tf2cother2"
## [21] "tf2cother5"    "tf2cother6"    "tf2cother7"    "wtambr2"
## [25] "wtambr4"       "wtambr5"       "wtambr1"       "wtambr5"
## [29] "wtambr6"       "wtambr7"       "wtambr8"       "wtbmr2"
## [33] "wtbmr3"       "wtbmr6"       "wtbmr8"       "wtbmr1.4"
## [37] "wtbmr3"       "wtbmr5"       "wtbmr8"       "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)
```

```
##                wtambr2 wtambr3 wtambr6 wtambr8 wtambr1.4 wtambr3
## Solyc00g005040.2.1      2      4      3      0      0      8
## Solyc00g005050.2.1     20      5     18      0      0     25
## Solyc00g005060.1.1      1      2      1      1      0      0
## Solyc00g005070.1.1     14      6     12     14      0      6
## Solyc00g005080.1.1     25     15     27      0      0     29
## Solyc00g005150.1.1      0      0      3      0      0      2
##                wtambr5 wtambr8
## Solyc00g005040.2.1      0      3
## Solyc00g005050.2.1      0     14
## Solyc00g005060.1.1      0      0
## Solyc00g005070.1.1      2      4
## Solyc00g005080.1.1      0     11
## Solyc00g005150.1.1      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)
```

```
d$samples
```

```
##                group lib.size norm.factors
## wtambr2          wtambr 1355352          1
## wtambr3          wtambr 1213142          1
## wtambr6          wtambr 1598917          1
```

```
## wtbmbr      wtbmbr      48352      1
## wtbother1.4 wtbother      1421      1
## wtbother3   wtbother  1076939      1
## wtbother5   wtbother   200587      1
## wtbother8   wtbother   499487      1
```

```
cpm.d <- cpm(d)
d <- d[rowSums(cpm.d>5)>=3,] #change to 5
d <- estimateCommonDisp(d,verbose=T)
```

```
## Disp = 0.4689 , BCV = 0.6848
```

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

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

```
##           logFC logCPM    PValue
## Solyc00g005050.2.1  1.0479  4.688 3.542e-01
## Solyc00g005070.1.1 -2.6995  5.950 5.907e-03
## Solyc00g005080.1.1  0.3842  5.334 8.386e-01
## Solyc00g005440.1.1 -0.4347  5.254 5.182e-01
## Solyc00g005840.2.1  4.2263  7.983 1.360e-06
## Solyc00g006470.1.1  1.7289  8.848 3.324e-02
```

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

```
## Comparison of groups: wtbother-wtbmbr
##           logFC logCPM    PValue    FDR
## Solyc04g074380.2.1 13.146 12.477 3.432e-22 5.009e-18
## Solyc04g074390.2.1 10.771 10.887 7.628e-18 5.566e-14
## Solyc02g076780.2.1  7.454 11.097 6.038e-14 2.937e-10
## Solyc01g014280.2.1  9.819  9.759 5.369e-13 1.959e-09
## Solyc01g065610.1.1  8.382  8.501 3.374e-12 9.849e-09
## Solyc10g078540.1.1  6.421 13.411 9.554e-12 2.067e-08
```

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

```
## [1] 14594      4
```

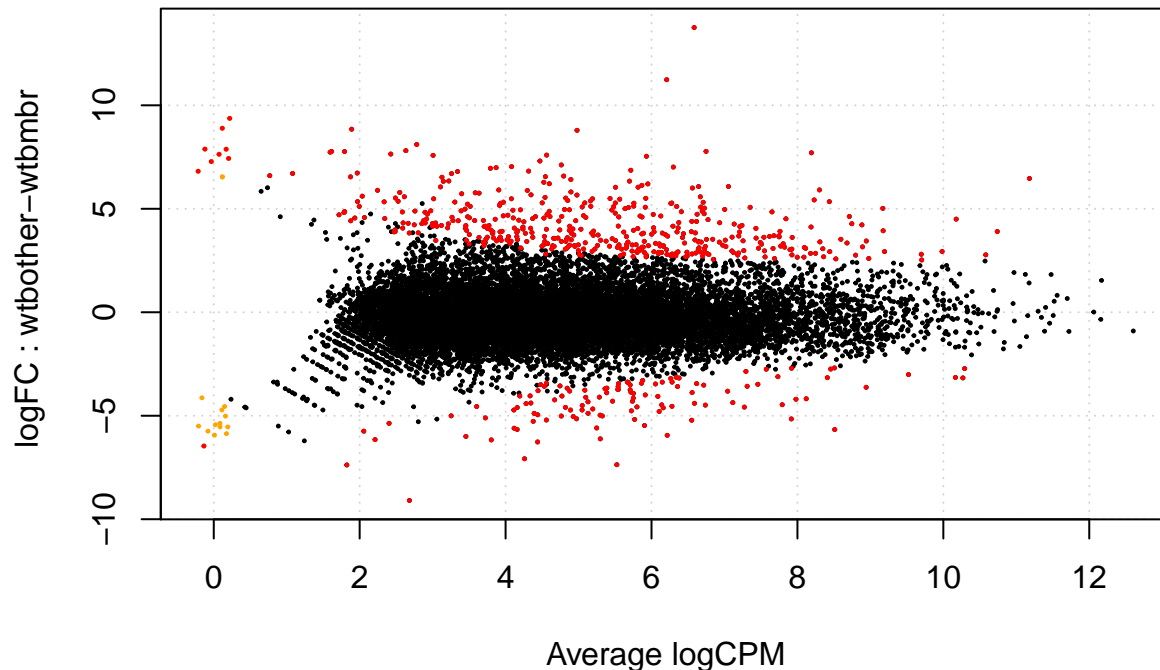
```
sum(results$table$FDR<.05) # How many are DE genes?
```

```
## [1] 556
```

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

```
##      [,1]
## -1     116
##  0    14038
##  1      440
```

```
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 ones with a significant 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("../ITAG2.3_all_Arabidopsis_ITAG_annotations.tsv", header=FALSE) #Changed to
colnames(annotation1) <- c("ITAG", "SGN_annotation")
annotation2<- read.delim("../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
```

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