

SCO_Analysis.R

srhilz

Mon Apr 30 20:34:51 2018

MAIN

```
# read in config file for analysis - change config to analyze a different
# subset of genes (choices for config - "Sertoli", "Leydig", "Union")
config <- "Sertoli"
source(paste0('config/',config,'Config.R'))

# set up file names
rawCountsFile <- 'data/merged_counts_noNC.txt'
CPMOutputFile <- paste0('output/',tag,'_CPMs.csv')
DAOutputFile <- paste0('output/',tag,'_NHvSCO_edgeR_results.csv')
PCAOutputFile <- paste0('output/',tag,'_PCA.pdf')
MAFile <- paste0('output/',tag,'_logCPM_v_logFC.pdf')
BoxplotRawOutputFile <- paste0('output/',tag,'_PreNorm_CPMs.pdf')
BoxplotNormOutputFile <- paste0('output/',tag,'_PostNorm_CPMs.pdf')
GOSpecificFile <- paste0('output/',tag,'_SpecificSubset_GeneOntology.csv')
GOUpFile <- paste0('output/',tag,'_Up_GeneOntology.csv')
GODownFile <- paste0('output/',tag,'_Down_GeneOntology.csv')
scatterFile <- paste0('output/',tag,'_Scatter.pdf')

# read in raw counts file
sampleTable_edgeR<-read.delim(rawCountsFile, row.names='gene')

# check dimensions
dim(sampleTable_edgeR)

## [1] 19136    11

# build logical vector of rownames that are not genes but summary outputs of HTSeq
noint = rownames(sampleTable_edgeR) %in% c("__ambiguous",
                                           "__too_low_aQual",
                                           "__not_aligned",
                                           "__no_feature",
                                           "__alignment_not_unique")

# set grouping - first four are normal, remaining are SCO
group<-factor(c(1,1,1,1,2,2,2,2,2,2))

# build DGEList object
d<-DGEList(counts=sampleTable_edgeR,group=group)

# subset original matrix by genes that are expressed over a CPM cutoff, and,
# if toFilter==1, that are in the provided gene list
if (toFilter==1){
  specific_list <- scan(file=specificListFile, what=character())
  specific = toupper(rownames(sampleTable_edgeR)) %in% toupper(specific_list)
  paste0('In specific list: ',
```

```

length(specific_list[toupper(specific_list) %in% toupper(rownames(sampleTable_edgeR))]))
paste('Not in specific list: ',
length(specific_list[!toupper(specific_list) %in% toupper(rownames(sampleTable_edgeR))]))
keep <- !noint & specific
}else{
keep <- !noint
}
d<- d[keep,]

# check dimensions after filtering
dim(d)

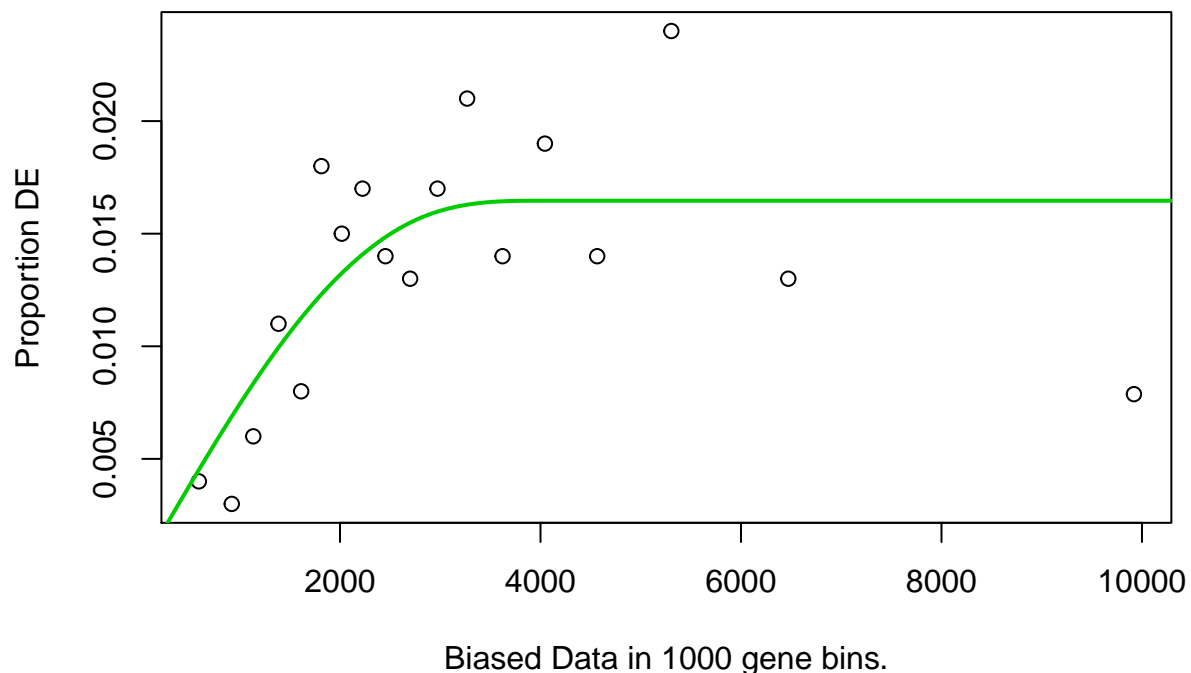
## [1] 247 11

# perform GO on specific gene list compared to all genes
if (toFilter == 1){
specificGenes=as.integer(rownames(sampleTable_edgeR) %in% specific_list)
names(specificGenes) <- rownames(sampleTable_edgeR)
performGO(specificGenes, GOSpecificFile)
}

## [1] "Table of input values"
## binaryList
##      0      1
## 18889  247

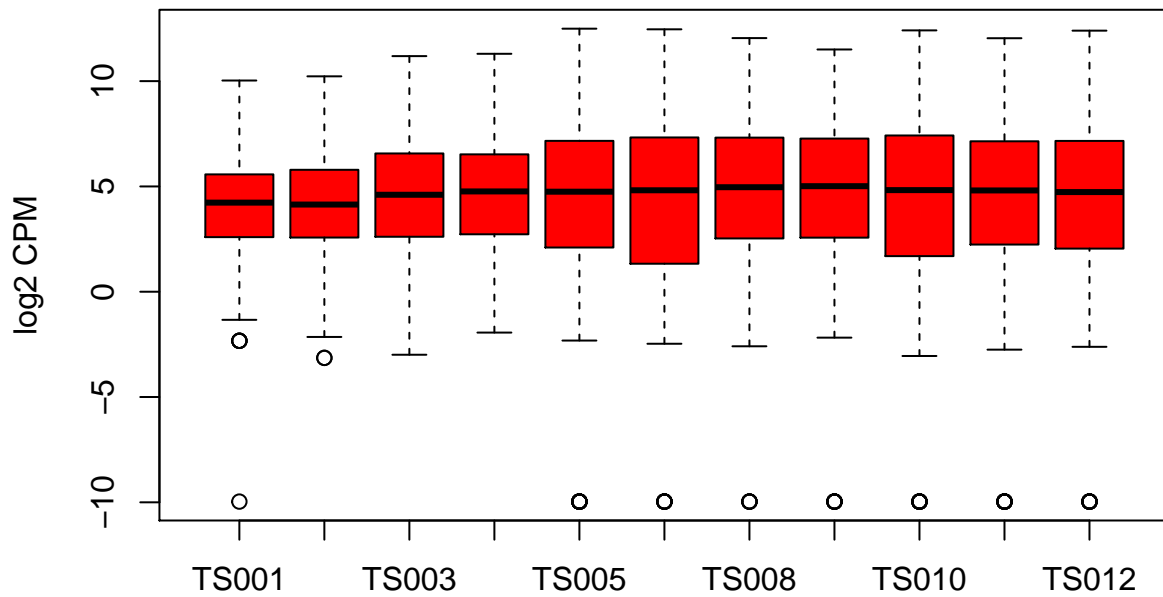
## Warning in pcls(G): initial point very close to some inequality constraints
## Warning: package 'AnnotationDbi' was built under R version 3.4.1
## Warning: package 'BiocGenerics' was built under R version 3.4.1
## Warning: package 'IRanges' was built under R version 3.4.1
## Warning: package 'S4Vectors' was built under R version 3.4.1

```



```
## [1] "Top 20 most significant GO terms"
##
##                                     term
## 1                                cell periphery
## 2                                plasma membrane
## 3                                tissue development
## 4                                cell adhesion
## 5                                biological adhesion
## 6                                system development
## 7                                single-multicellular organism process
## 8                                multicellular organism development
## 9                                tissue morphogenesis
## 10                               cell junction
## 11                               single organism signaling
## 12                               signaling
## 13                               anatomical structure morphogenesis
## 14                               cell communication
## 15                               animal organ development
## 16                               receptor binding
## 17                               epithelium development
## 18 calcium-independent cell-cell adhesion via plasma membrane cell-adhesion molecules
## 19                               anatomical structure development
## 20                               single-organism developmental process
##
##          pvalue
## 1  1.705603e-21
## 2  1.021351e-20
## 3  2.086989e-19
## 4  2.028868e-18
## 5  2.663135e-18
## 6  5.051657e-18
## 7  3.593947e-17
## 8  1.297262e-16
## 9  1.326245e-16
## 10 2.394725e-16
## 11 3.282893e-16
## 12 3.723366e-16
## 13 5.888423e-16
## 14 1.260277e-15
## 15 1.774490e-15
## 16 1.786974e-15
## 17 1.797483e-15
## 18 1.838582e-15
## 19 2.198258e-15
## 20 5.751350e-15

# look at CPMs of selected genes before any normalization
boxplot(log(cpm(d)+.001,2), col='red', ylab="log2 CPM")
```



```
dev.copy2pdf(file=BoxplotRawOutputFile)
```

```
## pdf
## 2
```

```
# calculate the normalization factors (this will correct for overall differences in count
# means between samples)
```

```
d<-calcNormFactors(d, method="RLE")#normalizing by log median
```

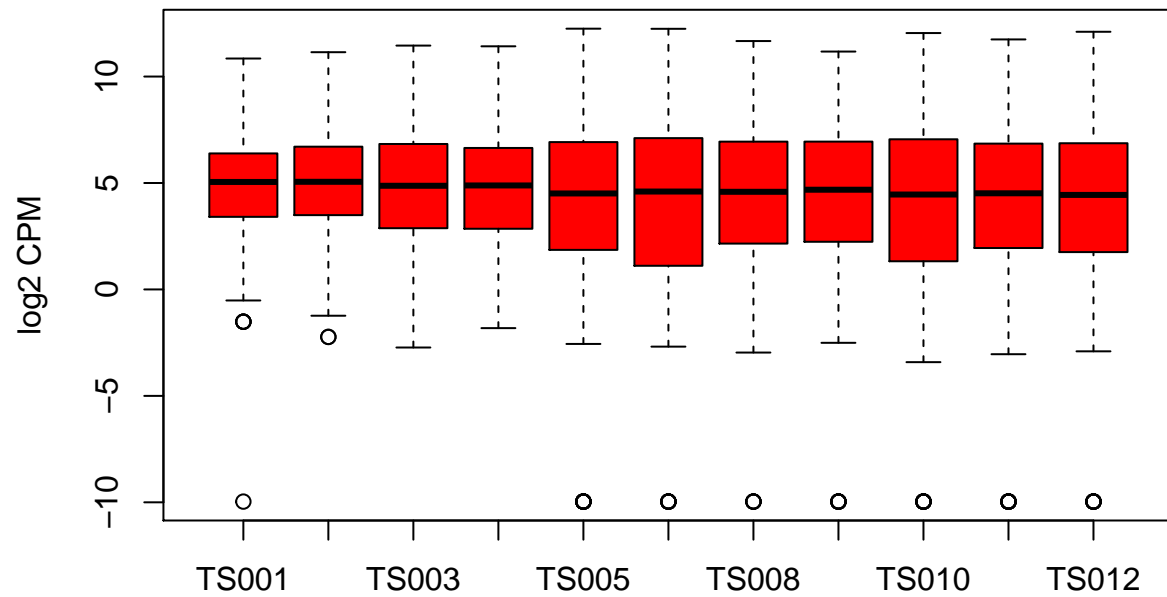
```
# show the normalization factor calculated for each library
```

```
d$samples
```

```
##      group lib.size norm.factors
## TS001     1  5057955   0.5674241
## TS002     1  8903776   0.5299568
## TS003     1  8018457   0.8318373
## TS004     1  7702072   0.9185077
## TS005     2  5015308   1.1831560
## TS007     2  5571521   1.1635257
## TS008     2  6060155   1.2984905
## TS009     2  4555248   1.2560846
## TS010     2  8361587   1.2898527
## TS011     2  6773761   1.2263488
## TS012     2  6169394   1.2254647
```

```
# look at CPMs of selected genes after normalization to mean counts
```

```
boxplot(log(cpm(d)+.001,2), col='red', ylab="log2 CPM")
```



```
dev.copy2pdf(file=BoxplotNormOutputFile)
```

```
## pdf
```

```
## 2
```

```
# estimate common dispersion across all samples
```

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

```
# view common dispersion
```

```
sqrt(d$common.disp)
```

```
## [1] 0.3467653
```

```
# estimate individual dispersion for each gene
```

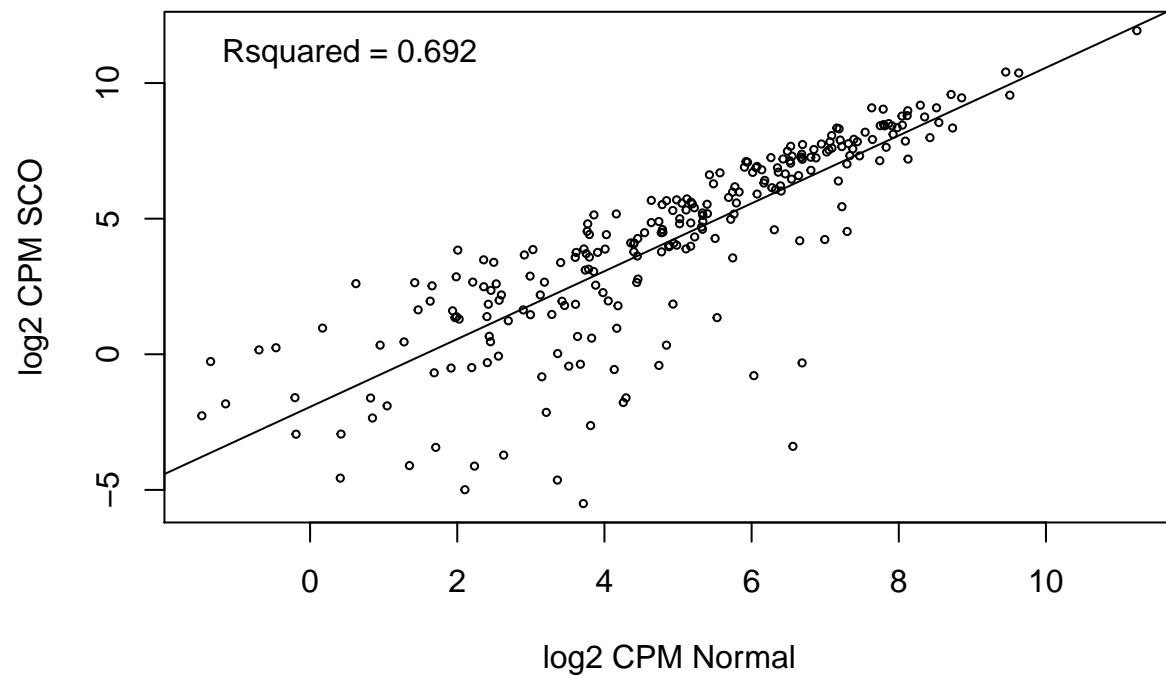
```
d<-estimateTagwiseDisp(d)
```

```
# ouptut CPMs to file
```

```
write.csv(cpm(d),CPMOutputFile)
```

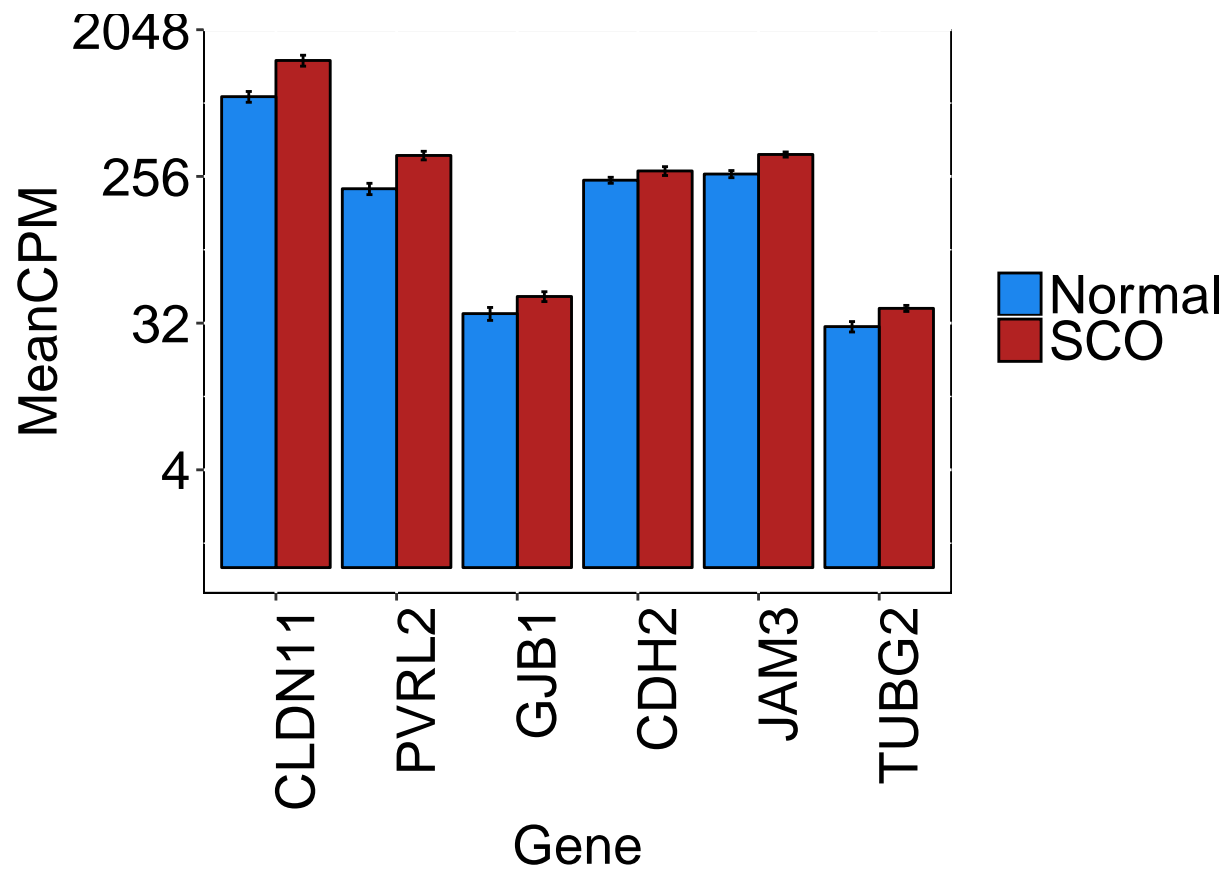
```
# examine CPMs as normal vs SCD scatterplot
```

```
plotScatter(cpm(d), group, scatterFile)
```

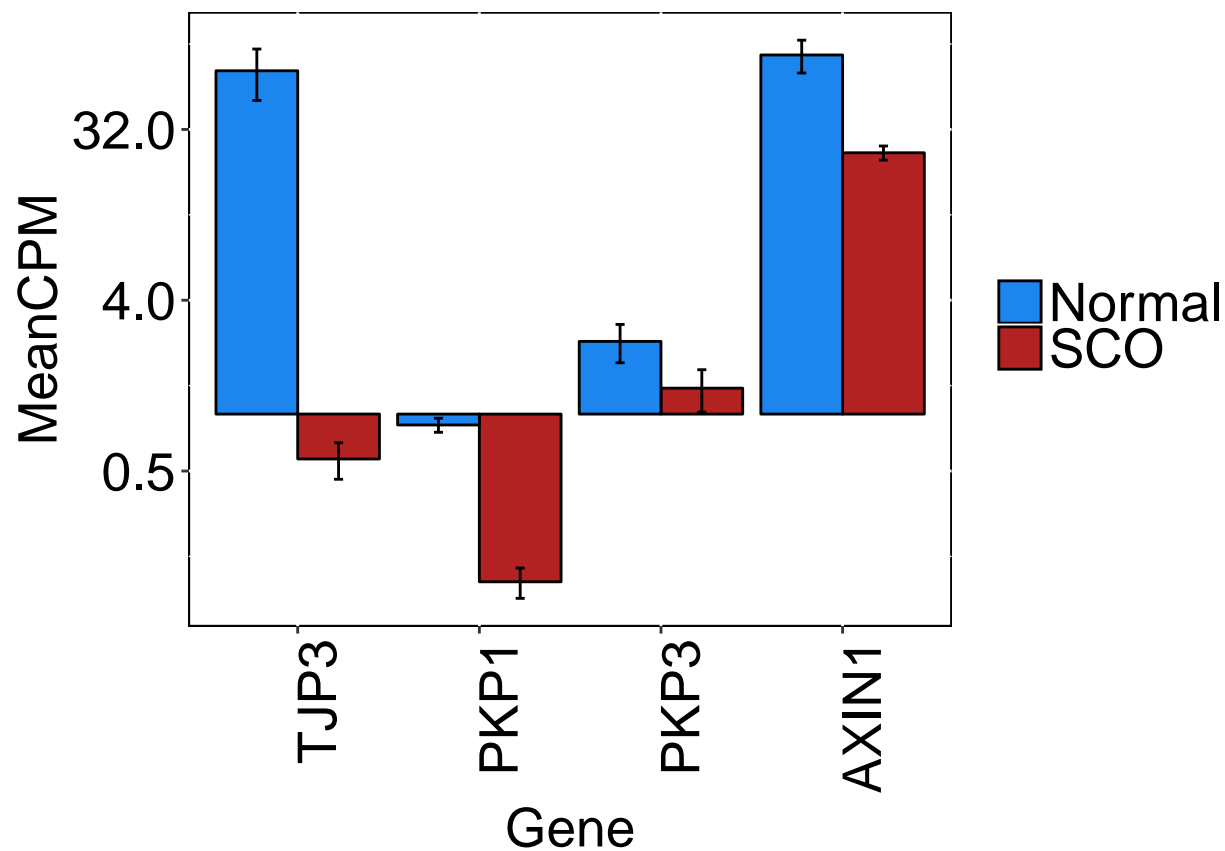


```
## pdf
## 2
# examine CPMs for proteins of interest
for (list in names(ofInterest)){
  print(paste0("Generating barplot for ",list," proteins"))
  barFile <- paste0("output/",tag,"_",list,"_Bar.pdf")
  plotBarchart(ofInterest[list], group, cpm(d), barFile)
}

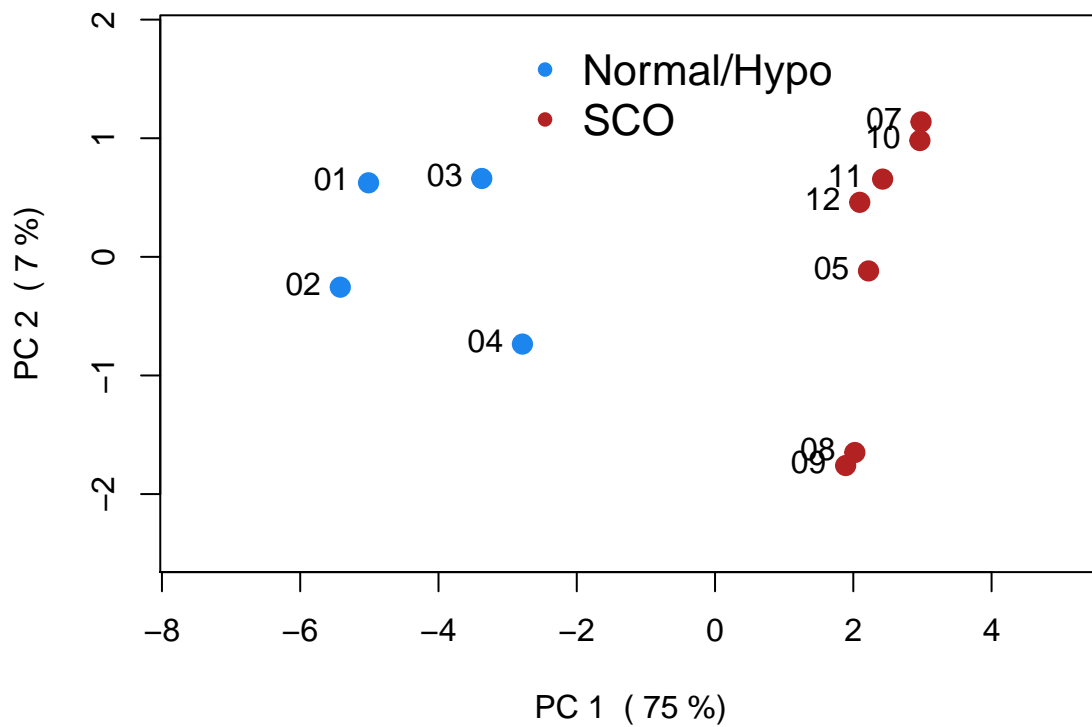
## [1] "Generating barplot for Transmembrane proteins"
```



```
## [1] "Generating barplot for Adapter proteins"
```



```
# PubQuality PCA plot with % explained
plotPCA(CPMOutputFile, PCAOutputFile)
```



NULL


```

## NULL
## $rect
## $rect$w
## [1] 4.969736
##
## $rect$h
## [1] 1.262631
##
## $rect$left
## [1] -3
##
## $rect$top
## [1] 2
##
##
## $text
## $text$x
## [1] -1.916906 -1.916906
##
## $text$y
## [1] 1.579123 1.158246

## pdf
## 2

# default of exactTest uses tag dispersion, does pairwise comp,
# comparing 2 to 1 Normal + Hypo vs SCO
NHvSCO_edgeR=exactTest(d, pair=c("1","2"))

# format results
results_NHvSCO<-topTags(NHvSCO_edgeR, n = nrow( NHvSCO_edgeR$table ) )$table

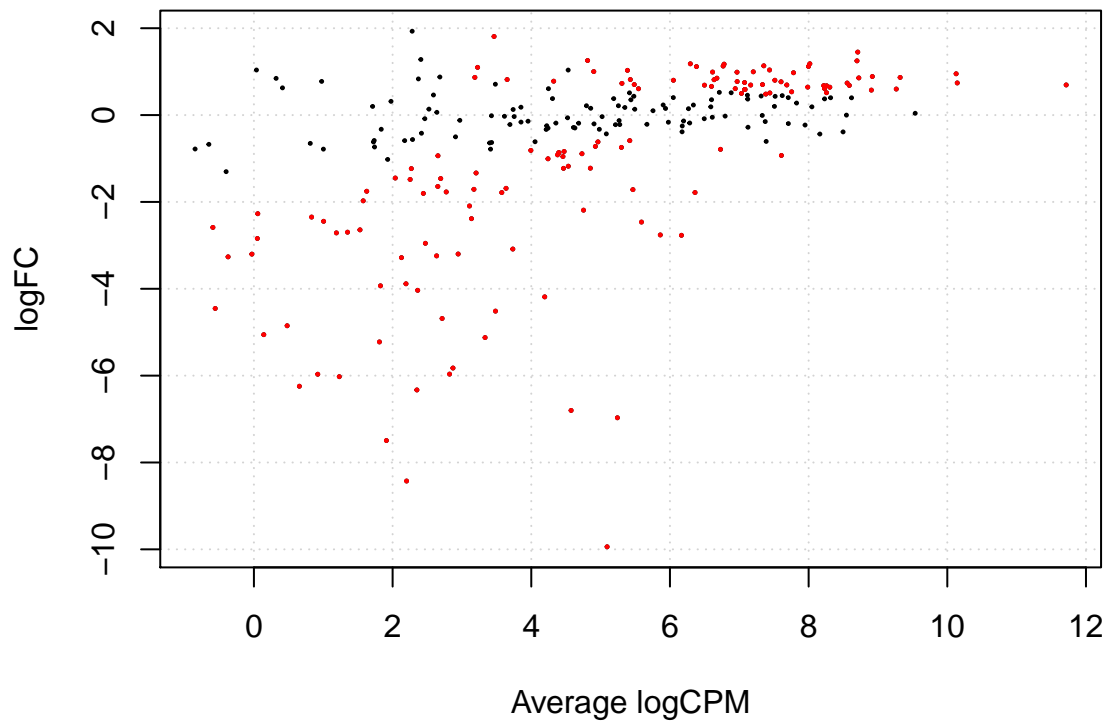
# make a vector of all differentially expressed genes
NHvSCO_detags <- rownames(results_NHvSCO)[results_NHvSCO$FDR < 0.05]

# summarize results
summary(decideTestsDGE(NHvSCO_edgeR, p=0.05, adjust="BH"))

##      1+2
## -1  76
## 0  110
## 1   61

# make a MA style plot
plotSmeaR(NHvSCO_edgeR, de.tags=NHvSCO_detags)

```



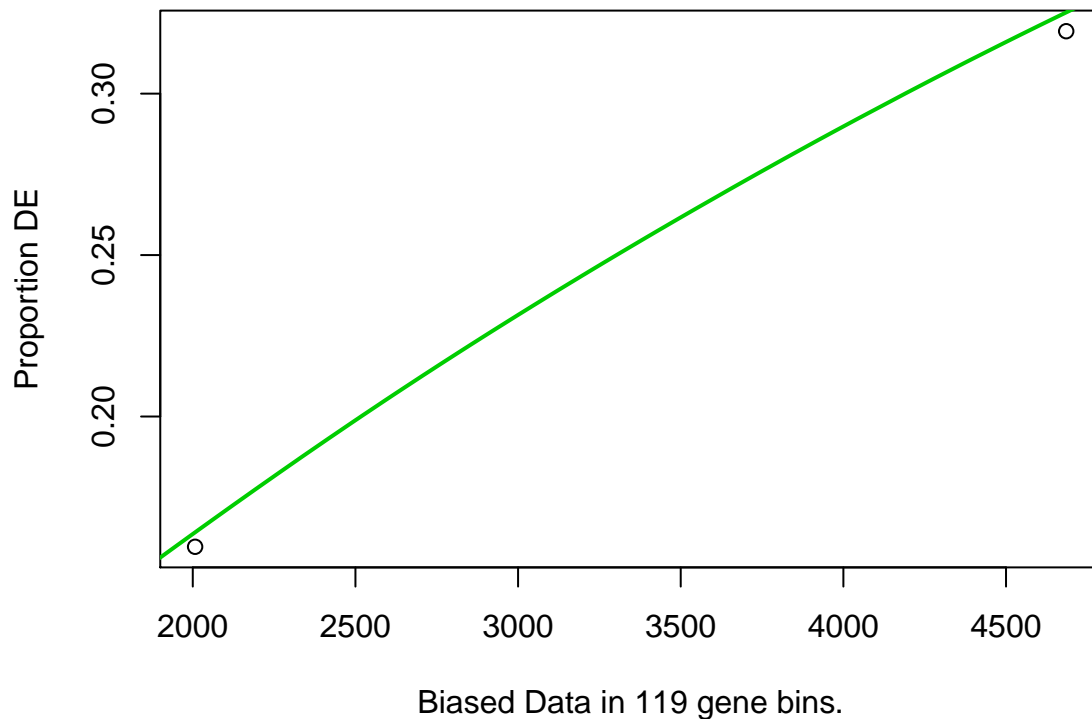
```
dev.copy2pdf(file=MAFile)
```

```
## pdf
## 2
# output to a file
write.csv(results_NHvSCO, DAOutputFile)

# perform GO on significantly upregulated genes
genes_up_NHvSCO = as.integer(results_NHvSCO$logFC > 0 &
                             rownames(results_NHvSCO) %in% NHvSCO_detags)
names(genes_up_NHvSCO) <- rownames(results_NHvSCO)
performGO(genes_up_NHvSCO, GOUpFile)

## [1] "Table of input values"
## binaryList
## 0 1
## 186 61

## Warning in pcls(G): initial point very close to some inequality constraints
```

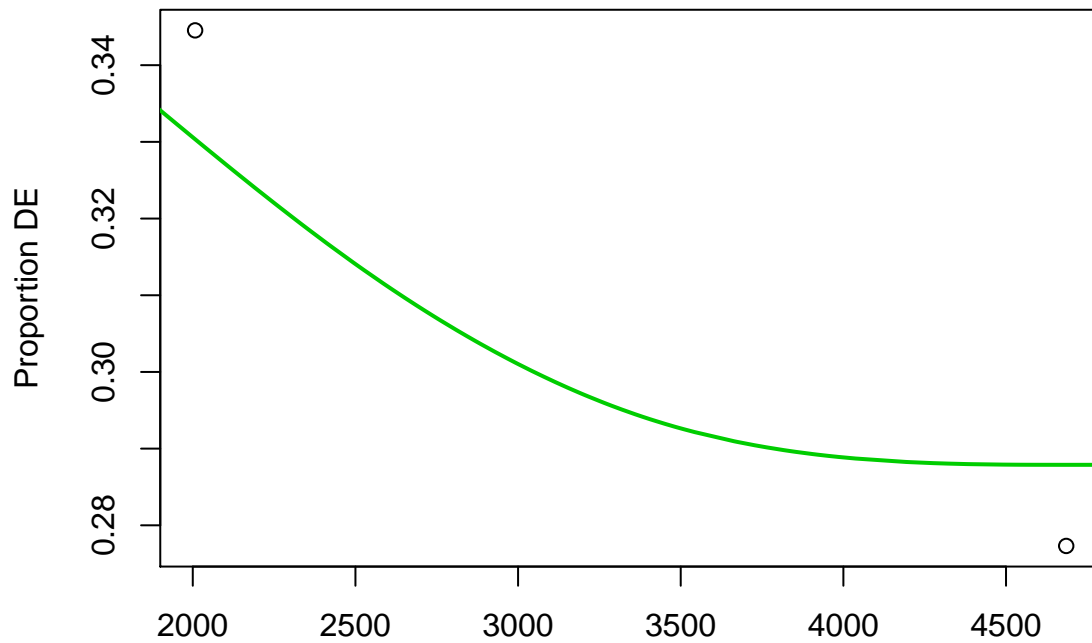


```
## [1] "Top 20 most significant GO terms"
##                                     term      pvalue
## 1                non-membrane-bounded organelle 0.002368962
## 2 intracellular non-membrane-bounded organelle 0.002368962
## 3                actin filament bundle assembly 0.005392929
## 4                actin filament bundle organization 0.005392929
## 5                      centrosome 0.011117972
## 6                Rho protein signal transduction 0.012380716
## 7                      RNA binding 0.012769383
## 8                      nucleolus 0.017380618
## 9                actin-based cell projection 0.017530118
## 10 contractile actin filament bundle assembly 0.017949996
## 11                stress fiber assembly 0.017949996
## 12      G2/M transition of mitotic cell cycle 0.018856008
## 13      cell cycle G2/M phase transition 0.018856008
## 14      Ras protein signal transduction 0.018940926
## 15      cytoskeletal protein binding 0.024618723
## 16                      cell cycle arrest 0.026956742
## 17      actomyosin structure organization 0.029013771
## 18                      cell recognition 0.030631501
## 19                      microvillus 0.031324746
## 20      regulation of cell cycle arrest 0.033707296

# perform GO on significantly downregulated genes
genes_down_NHvSCO=as.integer(results_NHvSCO$logFC < 0 &
                             rownames(results_NHvSCO) %in% NHvSCO_detags)
names(genes_down_NHvSCO) <- rownames(results_NHvSCO)
performGO(genes_down_NHvSCO, GODownFile)

## [1] "Table of input values"
## binaryList
##    0    1
```

```
## 171 76
```



Biased Data in 119 gene bins.

```
## [1] "Top 20 most significant GO terms"
```

```
##               term      pvalue
## 1             axonogenesis 0.0006439074
## 2             cell part morphogenesis 0.0010984568
## 3             neuron projection morphogenesis 0.0010984568
## 4             cell projection morphogenesis 0.0010984568
## 5             cellular response to lipid 0.0024508254
## 6 cell morphogenesis involved in neuron differentiation 0.0026668284
## 7             intracellular signal transduction 0.0026727068
## 8             embryonic hindlimb morphogenesis 0.0034519869
## 9             hindlimb morphogenesis 0.0034519869
## 10            axon guidance 0.0039540333
## 11            neuron projection guidance 0.0039540333
## 12 central nervous system neuron differentiation 0.0042953717
## 13            embryonic limb morphogenesis 0.0045688055
## 14            embryonic appendage morphogenesis 0.0045688055
## 15            response to lipid 0.0047924233
## 16            axon development 0.0048834683
## 17            MAPK cascade 0.0052817913
## 18 apoptotic process involved in morphogenesis 0.0065273454
## 19 apoptotic process involved in development 0.0065273454
## 20            actin filament capping 0.0068028875
```

```
""
```