

Transcriptome Demo

Jair

2024-05-09

Load required packages

```
library(ballgown)
library(RColorBrewer)
library(genefilter)
library(dplyr)
library(devtools)
```

#Organizes data into a dataframe that has two columns. #The first column contains the ids while the second specifies the corresponding stage, either planktonic or biofilm.

```
pheno_data<-data.frame(ids = c("plank01", "plank02", "biofilm01", "biofilm02"),
                       stage = c("planktonic", "planktonic", "biofilm", "biofilm"))
```

#Creates a Ballgown object and check transcript number #Organizes the file path into a vector

```
samples.c <- paste('ballgown', pheno_data$ids, sep = '/')
bg <- ballgown(samples = samples.c, meas='all', pData = pheno_data)
bg
```

ballgown instance with 5737 transcripts and 4 samples

#Filters out transcripts with a small variance. #Only transcripts with a variance larger than 1 are kept.

```
bg_filt = subset(bg,"rowVars(expr(bg)) >1",genomesubset=TRUE)
bg_filt
```

ballgown instance with 5163 transcripts and 4 samples

#Creates a table of transcripts

```
results_transcripts<- stattest(bg_filt, feature = "transcript", covariate = "stage",
getFC = TRUE, meas = "FPKM")
results_transcripts<-data.frame(geneNames=geneNames(bg_filt),
transcriptNames=transcriptNames(bg_filt), results_transcripts)
```

#Shows the top of the table stored in results_transcripts

```
head(results_transcripts)
```

##	geneNames	transcriptNames	feature	id	fc	pval	qval
## 1	dnaA	gene-PA0001	transcript	1	5.247471e+01	0.3048003	0.9471885
## 2	dnaN	gene-PA0002	transcript	2	1.745401e+01	0.1001167	0.9471885
## 3	recF	gene-PA0003	transcript	3	5.229990e-01	0.8960742	0.9845954
## 4	gyrB	gene-PA0004	transcript	4	4.834298e+10	0.2743082	0.9471885
## 5	lptA	gene-PA0005	transcript	5	8.951420e+00	0.2455336	0.9471885

```
## 6          .      gene-PA0006 transcript  6 3.264697e+02 0.2949859 0.9471885
#Choose a transcript to examine more closely
results_transcripts[results_transcripts$transcriptNames == "gene-PA0004", ] #I chose to examine gene-PA

##   geneNames transcriptNames   feature id          fc      pval      qval
## 4      gyrB      gene-PA0004 transcript  4 48342981060 0.2743082 0.9471885

##This transcript is for the gene gyrB. ##The fold difference between planktonic and biofilm stages is
48342981060 ##The corrected q-value of 0.9471885 indicates that the difference is not significant.

#Filters out non-significant transcripts. #Only results with a significant p-value (less than 0.05) are kept
#The dim function gives the number of transcripts kept.
sigdiff <- results_transcripts %>% filter(pval<0.05)
dim(sigdiff)

## [1] 207   7

#Organizes the table by p-value and fold change (fc). #The p-values go from smallest to biggest (since
decreasing=FALSE) #The fold change is in increasing order since the negative sign in front and setting
decreasing=FALSE.
o = order(sigdiff[, "pval"], -abs(sigdiff[, "fc"]), decreasing=FALSE)
output = sigdiff[o, c("geneNames", "transcriptNames", "id", "fc", "pval", "qval")]
write.table(output, file="SigDiff.txt", sep="\t", row.names=FALSE, quote=FALSE)
head(output)

##      geneNames transcriptNames   id          fc      pval      qval
## 4091          .      gene-PA3992 4091 9.886091e+01 0.0003032315 0.9471885
## 4958          .      gene-PA4804 4958 3.563696e-04 0.0006661432 0.9471885
## 2745          .      gene-PA2690 2745 5.783390e-02 0.0014192618 0.9471885
## 2896      tpm      gene-PA2832 2896 1.786570e+03 0.0023414834 0.9471885
## 370          .      gene-PA0365 370 3.964652e-07 0.0023906201 0.9471885
## 3129      pelF      gene-PA3059 3129 1.687425e-03 0.0025838457 0.9471885

#Loads gene names
bg_table = texpr(bg_filt, 'all')
bg_gene_names = unique(bg_table[, 9:10])

#Pulls out gene expression data and visualize
gene_expression = as.data.frame(gexpr(bg_filt))
head(gene_expression)

##      FPKM.plank01 FPKM.plank02 FPKM.biofilm01 FPKM.biofilm02
## .      1.198359    0.9103059    2.526183    2.685373
## MSTRG.1      405.892761  400.8589780    232.324417  181.932617
## MSTRG.10      89.649139   78.5762100    35.010487   59.757320
## MSTRG.100     116.443428  106.2109530    92.206810   95.322479
## MSTRG.1000     7.833186   5.5019700    15.717344   42.342495
## MSTRG.1001     6.845010   4.7381980    38.199095   89.078876

#Renames the columns to the names specified within c(" ")
colnames(gene_expression) <- c("plank01", "plank02", "biofilm01", "biofilm02")
head(gene_expression)

##      plank01      plank02 biofilm01 biofilm02
```

```
## .          1.198359  0.9103059  2.526183  2.685373
## MSTRG.1    405.892761 400.8589780 232.324417 181.932617
## MSTRG.10   89.649139  78.5762100  35.010487  59.757320
## MSTRG.100 116.443428 106.2109530  92.206810  95.322479
## MSTRG.1000 7.833186  5.5019700  15.717344  42.342495
## MSTRG.1001 6.845010  4.7381980  38.199095  89.078876
```

```
dim(gene_expression)
```

```
## [1] 4592    4
```

```
#Loads the transcript to gene table and determine the number of transcripts and unique genes
```

```
transcript_gene_table = indexes(bg)$t2g
head(transcript_gene_table)
```

```
##   t_id  g_id
## 1    1 MSTRG.1
## 2    2 MSTRG.2
## 3    3 MSTRG.3
## 4    4 MSTRG.3
## 5    5 MSTRG.4
## 6    6 MSTRG.5
```

```
length(row.names(transcript_gene_table))
```

```
## [1] 5737
```

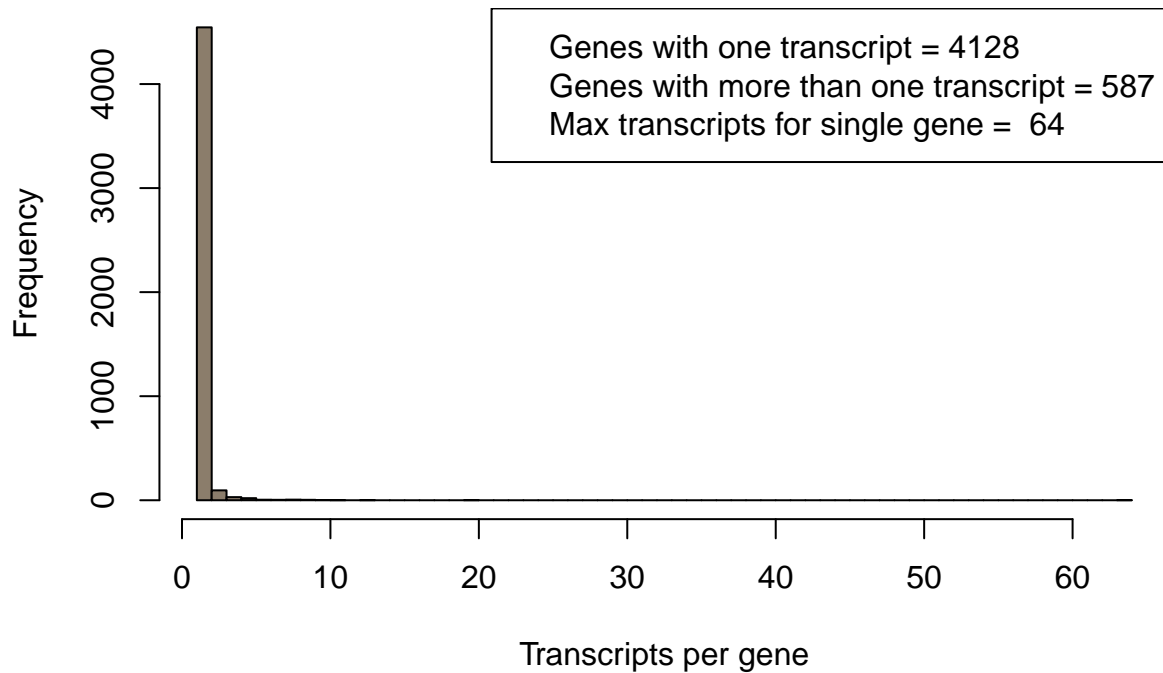
```
length(unique(transcript_gene_table[, "g_id"]))
```

```
## [1] 4715
```

```
#Plots the number of transcripts per gene
```

```
counts=table(transcript_gene_table[, "g_id"])
c_one = length(which(counts == 1))
c_more_than_one = length(which(counts > 1))
c_max = max(counts)
hist(counts, breaks=50, col="bisque4", xlab="Transcripts per gene",
main="Distribution of transcript count per gene")
legend_text = c(paste("Genes with one transcript =", c_one),
paste("Genes with more than one transcript =", c_more_than_one),
paste("Max transcripts for single gene = ", c_max))
legend("topright", legend_text, lty=NULL)
```

Distribution of transcript count per gene

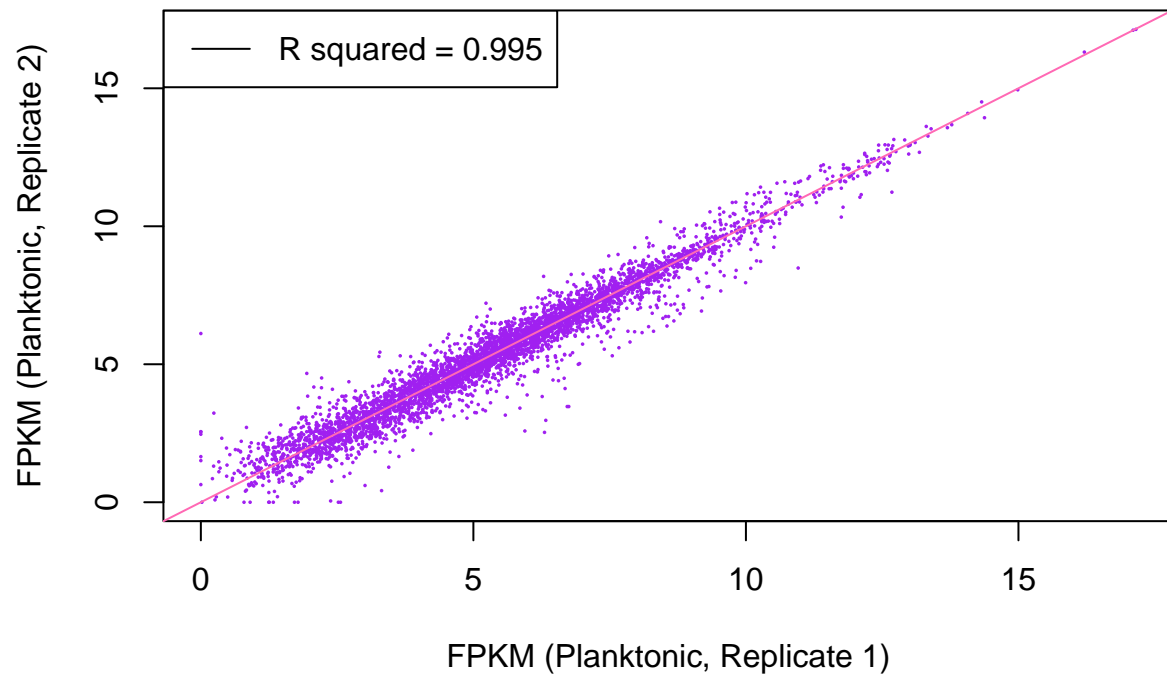


##A vast majority of the genes have a single transcript. ##Few genes have multiple transcripts ##The most transcripts for a single gene were 64.

#Creates a plot of how similar the two planktonic replicates are to one another.

```
x = gene_expression[, "plank01"]
y = gene_expression[, "plank02"]
min_nonzero=1
plot(x=log2(x+min_nonzero), y=log2(y+min_nonzero), pch=16, col="purple", cex=0.25,
xlab="FPKM (Planktonic, Replicate 1)", ylab="FPKM (Planktonic, Replicate 2)",
main="Comparison of expression values for a pair of replicates")
abline(a=0,b=1, col = "hotpink")
rs=cor(x,y)^2
legend("topleft", paste("R squared = ", round(rs, digits=3), sep=""), lwd=1, col="black")
```

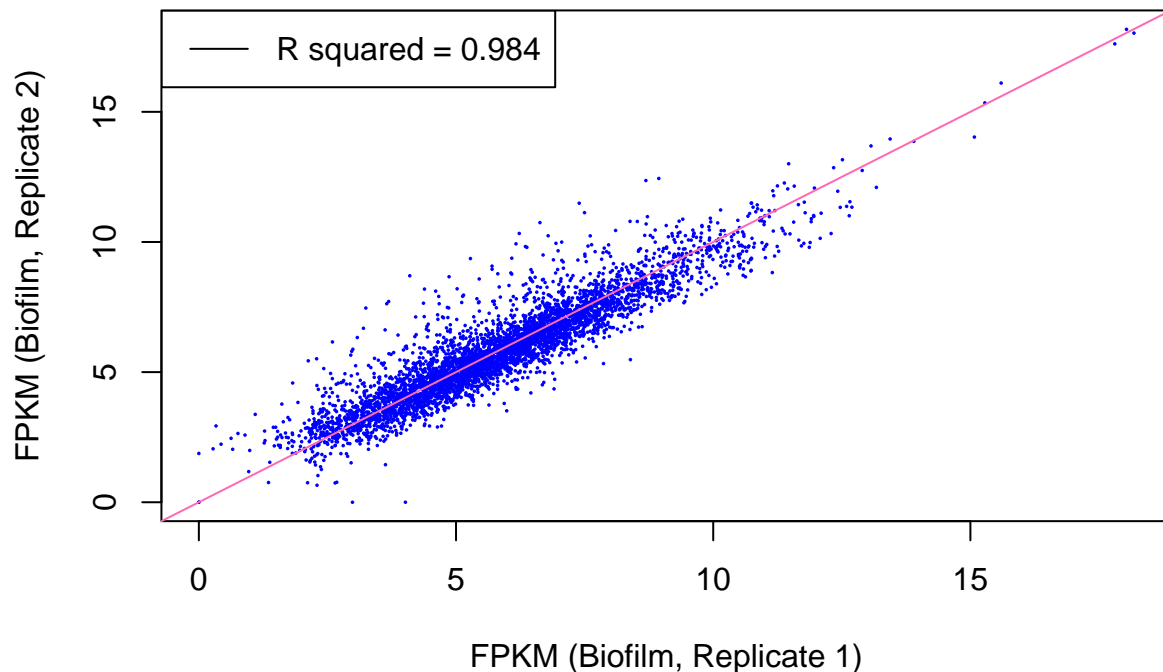
Comparison of expression values for a pair of replicates



#Creates a plot of how similar the two biofilm replicates are to one another.

```
x = gene_expression["biofilm01"]
y = gene_expression["biofilm02"]
min_nonzero=1
plot(x=log2(x+min_nonzero), y=log2(y+min_nonzero), pch=16, col="blue", cex=0.25,
xlab="FPKM (Biofilm, Replicate 1)", ylab="FPKM (Biofilm, Replicate 2)",
main="Comparison of expression values for a pair of replicates")
abline(a=0,b=1, col = "hotpink")
rs=cor(x,y)^2
legend("topleft", paste("R squared = ", round(rs, digits=3), sep=""), lwd=1, col="black")
```

Comparison of expression values for a pair of replicates

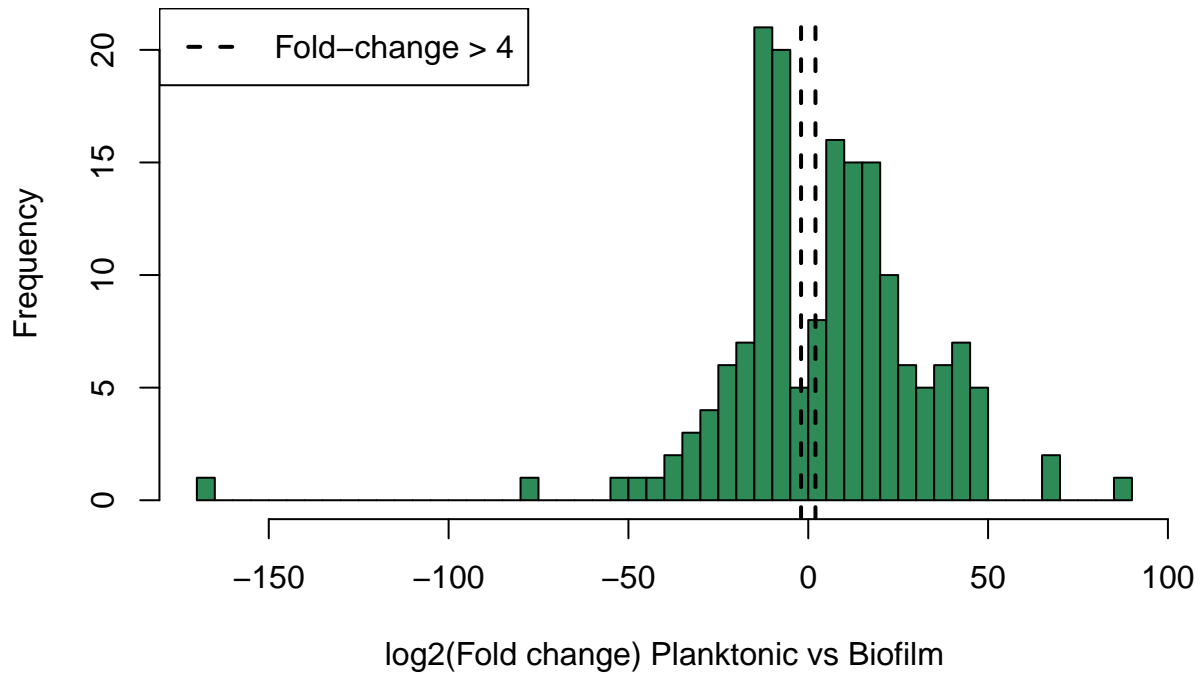


##Both replicates within each condition were similar. ##If both replicates within a single condition are similar, this means we can reliably compare between conditions (planktonic and biofilm)

#Creates plot of differential gene expression between the conditions

```
results_genes = statstest(bg_filt, feature="gene", covariate="stage", getFC=TRUE, meas="FPKM")
results_genes = merge(results_genes, bg_gene_names, by.x=c("id"), by.y=c("gene_id"))
sig=which(results_genes$pval<0.05)
results_genes[, "de"] = log2(results_genes[, "fc"])
hist(results_genes[sig, "de"], breaks=50, col="seagreen",
xlab="log2(Fold change) Planktonic vs Biofilm",
main="Distribution of differential expression values")
abline(v=-2, col="black", lwd=2, lty=2)
abline(v=2, col="black", lwd=2, lty=2)
legend("topleft", "Fold-change > 4", lwd=2, lty=2)
```

Distribution of differential expression values

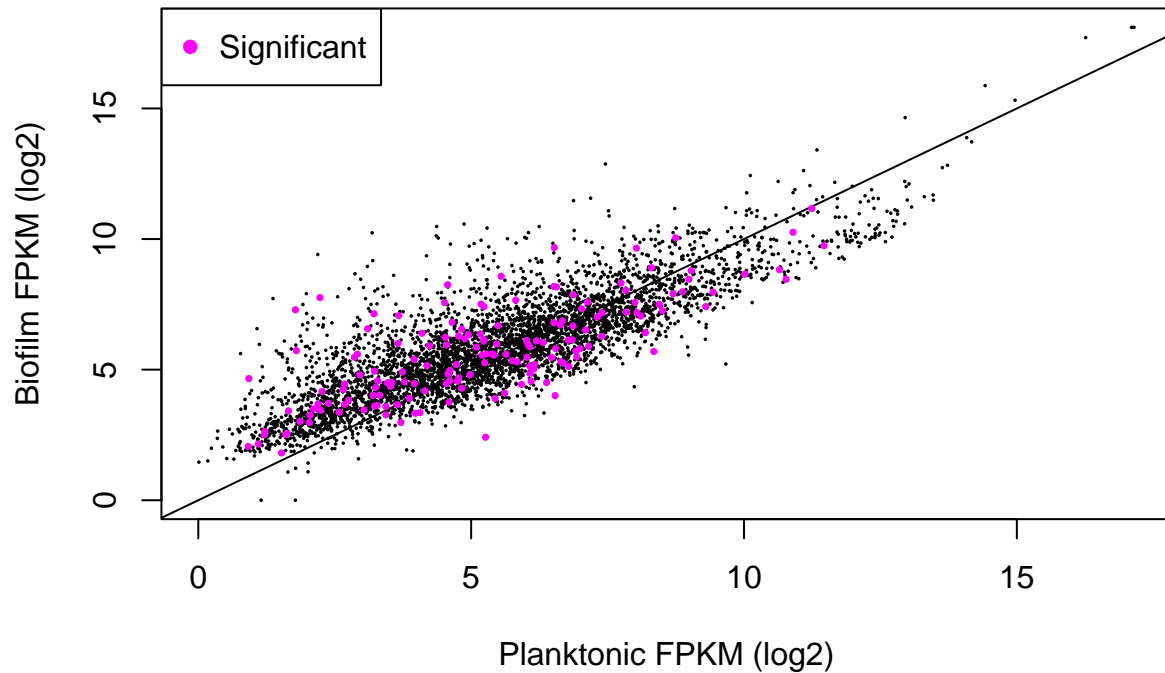


##The differential gene expression plot is bimodal, indicating a difference in gene expression between both conditions. #Many genes had a fold-change greater than 4, suggesting a large difference in expression base on the condition (planktonic or biofilm)

#Plots total gene expression highlighting differentially expressed genes

```
gene_expression[, "plank"] = apply(gene_expression[, c(1:2)], 1, mean)
gene_expression[, "biofilm"] = apply(gene_expression[, c(3:4)], 1, mean)
x = log2(gene_expression[, "plank"] + min_nonzero)
y = log2(gene_expression[, "biofilm"] + min_nonzero)
plot(x=x, y=y, pch=16, cex=0.25, xlab="Planktonic FPKM (log2)", ylab="Biofilm FPKM (log2)",
main="Planktonic vs Biofilm FPKMs")
abline(a=0, b=1)
xsig=x[sig]
ysig=y[sig]
points(x=xsig, y=ysig, col="magenta", pch=16, cex=0.5)
legend("topleft", "Significant", col="magenta", pch=16)
```

Planktonic vs Biofilm FPKMs



```
#Makes a table of FPKM values
```

```
fpkm = texpr(bg_filt, meas="FPKM")
```

```
#Chooses a gene to determine individual expression
```

```
ballgown::transcriptNames(bg_filt)[10]
```

```
##      10
```

```
## "gene-PA0010"
```

```
ballgown::geneNames(bg_filt)[10]
```

```
##      10
```

```
## "tag"
```

```
#Transforms to log2
```

```
transformed_fpkm <- log2(fpkm[2, ] + 1)
```

```
#Makes sure values are properly coded as numbers
```

```
numeric_stages <- as.numeric(factor(pheno_data$stage))
```

```
jittered_stages <- jitter(numeric_stages)
```

```
#Plots expression of individual gene
```

```
boxplot(transformed_fpkm ~ pheno_data$stage,  
  main=paste(ballgown::geneNames(bg_filt)[10], ' : ', ballgown::transcriptNames(bg_filt)[10]),  
  xlab="Stage",
```

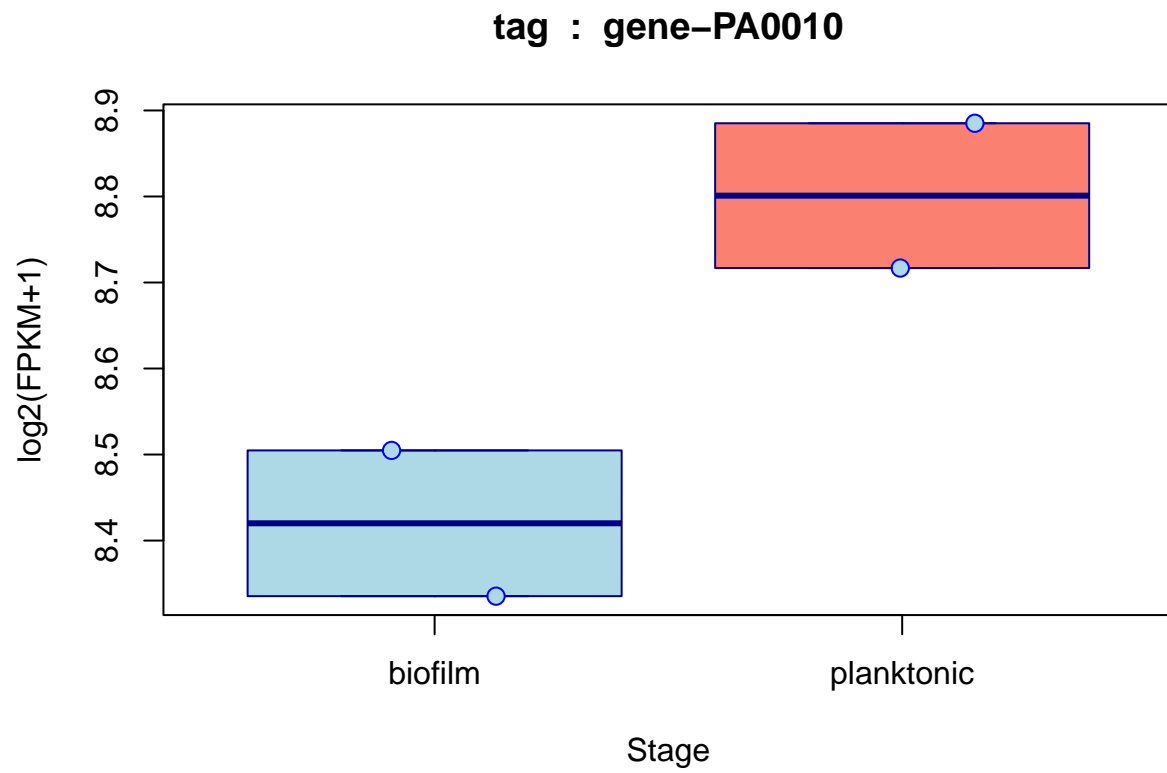


```

ylab="log2(FPKM+1)",
col=c("lightblue", "salmon"),
border="darkblue")

points(transformed_fpk ~ jittered_stages,
pch=21, col="blue", bg="lightblue", cex=1.2)

```



##gene-PA0010 had a higher expression in the planktonic than in the biofilm stage.