

1 Introduction

1.1 Disclaimer

Use at your own risk. Some errors may exist and so it is imperative that you understand each line of code and what it's doing with respect to your experimental needs before copying and pasting line for line.

1.2 What is edgeR?

edgeR is a package written for [R] that performs differential tag/gene expression using count data under a negative binomial model. The methods used in edgeR do NOT support FPKM, RPKM or other types of data that are not counts.

You can obtain more information about edgeR from the bioconductor website <http://www.bioconductor.org/packages/2.8/bioc/html/edgeR.html> . Here you'll find a user's guide/vignette with detailed examples, a reference manual for all edgeR functions, and the [R] code for installing the package. A useful discussion on a forum about edgeR can be found at <http://seqanswers.com/forums/showthread.php?t=5591> .

1.3 Getting Help with edgeR

The bioconductor website provides access to a mailing list that all can subscribe to. Any questions concerning any bioconductor package can be sent to this mailing list. They are very responsive and very helpful; I have used it several times. To subscribe go to <https://stat.ethz.ch/mailman/listinfo/bioconductor> . For guides to posting questions go to <http://www.bioconductor.org/help/mailling-list/posting-guide/> .

If you just need help understanding an edgeR function or want to see an example of how it's used simply type `?functionname` into the [R] command line. This will bring up a text page with detailed information about the function.

1.4 Getting Help with [R]

Google your question! Write a short description and add on at the end "R-help". This will most likely bring up a thread in the R-help forum that deals with the same or similar problem. Almost every problem you could have can be found answered on this forum. You can also subscribe to the [R] mailing list by going to <https://stat.ethz.ch/mailman/listinfo/r-help>. A guide for posting questions can be found at <http://www.r-project.org/posting-guide.html>

Elizabeth Purdom has several [R] tutorials that can be found at <http://www.stat.berkeley.edu/~epurdom/>. In particular, see <http://www.stat.berkeley.edu/~epurdom/Saving/Saving.html> for using directories and saving objects and <http://www.stat.berkeley.edu/~epurdom/BrownLab/> for creating and manipulating data sets in [R].

1.5 Installing edgeR

With an internet connection, edgeR can be easily installed by running the following commands in [R]

```
source("http://www.bioconductor.org/biocLite.R")
biocLite("edgeR")
```

2 Example 1: A Straightforward RNA-Seq Differential Expression Analysis using edgeR

2.1 Description

This is an RNA-Seq data set from an article by Li et al. [2008] titled, "Determination of tag density required for digital transcriptome analysis: Application to an androgen-sensitive prostate cancer model." The study compared hormone treated cancer cells to non-treated cancer cells. For the non-treated, there are 4 biological replicates. For the treated, there are 3 biological replicates. A full analysis of this data set can also be found in section 10 of the edgeR vignette.

2.2 Setting Up the Count Matrix

First, we load the edgeR package into [R]:

```
library(edgeR)
```

Set the directory to the folder where the data is contained and read in the count data. If using Windows “\\” is needed between directories; if using Unix “/” is sufficient.

```
setwd( "C:\\Documents and Settings\\All Users\\Documents\\Shared Files\\School\\Seminars\\  
RNA-Seq Workshop\\Data" )  
raw.data <- read.table( file = "pnas_expression.txt" , header = TRUE )  
head( raw.data )
```

edgeR requires the the dataset to contain only the counts with the row names as the gene ids and the column names as the sample ids:

```
lane.cols <- grep( "lane" , colnames( raw.data ) , fixed = TRUE )  
lane.cols  
  
counts <- raw.data[ , lane.cols ]  
rownames( counts ) <- raw.data[ , 1 ] # gene names  
colnames( counts ) <- paste(c(rep("C_R",4),rep("T_R",3)),c(1:4,1:3),sep="") # sample names  
head( counts )
```

Some quick summaries:

```
dim( counts )  
colSums( counts ) # Library Sizes  
colSums( counts ) / 1e06 # Library Sizes in millions of reads  
table( rowSums( counts ) )[ 1:30 ] # Number of genes with low counts
```

2.3 Building the edgeR Object

DGEList is the function that coverts the count matrix into an edgeR object. First, we create a group variable that tells edgeR with samples belong to which group and supply that to DGEList in addition to the count matrix. We can then see the elements that the object contains by using the names() function. These elements can be accessed using the \$ symbol:

```
group <- c(rep("C", 4) , rep("T", 3))  
cds <- DGEList( counts , group = group )  
names( cds )  
  
head(cds$counts) # original count matrix  
cds$samples # contains a summary of your samples  
sum( cds$all.zeros ) # How many genes have 0 counts across all samples  
  
cds # or type the name of the object
```

We need to filter out low count reads since it would be impossible to detect differential expression. The method used in the edgeR vignette is to keep only those genes that have at least 1 read per million in at least 3 samples. Once this filtering is done, we can calculate the normalization factors which correct for the different compositions of the samples. The effective library sizes are then the product of the actual library sizes and these factors.

```
cds <- cds[rowSums(1e+06 * cds$counts/expandAsMatrix(cds$samples$lib.size, dim(cds)) > 1) >= 3, ]
dim( cds )

cds <- calcNormFactors( cds )
cds$samples

# effective library sizes
cds$samples$lib.size * cds$samples$norm.factors
```

2.4 Multi-Dimensional Scaling Plot

An MDS plot measures the similarity of the samples and projects this measure into 2-dimensions. When creating plots in R, you can also save plots as pdf's, png's, etc. (see *?pdf*).

```
# To view the plot immediately
plotMDS.dge( cds , main = "MDS Plot for Count Data", labels = colnames( cds$counts ) )

# Output plot as a pdf
pdf( "MDS_plot_1_ex1.pdf" , width = 7 , height = 7 ) # in inches
plotMDS.dge( cds , main = "MDS Plot for Count Data", labels = colnames( cds$counts ) )
dev.off() # this tells [R] to close and stop writing to the pdf.
```

2.5 Estimating Dispersions

The first dispersion type to calculate is the common dispersion. In the common dispersion setting, each gene gets assigned the same dispersion estimate. The output of the estimation will include the estimate as well as some other elements added to the edgeR object, cds.

```
cds <- estimateCommonDisp( cds )
names( cds )

# The estimate
cds$common.dispersion
```

To understand what this value means, recall the parameterization for the variance of the negative binomial is $\nu(\mu) = \mu + \mu^2 \cdot \phi$. For poisson it's $\nu(\mu) = \mu$. The implied standard deviations are the square-roots of the variances. Now, suppose a gene had an average count of 200. Then the sd's under the two models would be,

```
sqrt( 200 ) # poisson sd
sqrt( 200 + 200^2 * cds$common.dispersion ) # negative binomial sd
sqrt( 200 + 200^2 * cds$common.dispersion ) / sqrt( 200 ) # NB sd is over 2 times larger
```

Once the common dispersion is estimated we can estimate the tagwise dispersions. In this scenario, each gene will get its own unique dispersion estimate, but the common dispersion is still used in the calculation. The tagwise dispersions are squeezed toward the common value. The amount of squeezing is governed by the parameter *prior.n*. The higher *prior.n*, the closer the estimates will be to the common dispersion. The recommended value is the nearest integer to $50/(\#samples - \#groups)$. For this data set that's $50/(7 - 2) = 10$.

```
# Default Setting
cds <- estimateTagwiseDisp( cds , prior.n = 10 )
names( cds )
summary( cds$tagwise.dispersion )

# More shrinkage/squeezing toward the common
cds <- estimateTagwiseDisp( cds , prior.n = 25 )
summary( cds$tagwise.dispersion ) # not much changed, but the ends got squeezed in quite a bit.

# The recommended setting for this data set is the default of 10. Let's stick with that.
cds <- estimateTagwiseDisp( cds , prior.n = 10 )
```

2.6 Mean-Variance Plot

Now that we've estimated the dispersion parameters we can see how well they fit the data by plotting the mean-variance relationship. Four things are shown in the plot: the raw variances of the counts (grey dots), the variances using the tagwise dispersions (light blue dots), the variances using the common dispersion (solid blue line), and the *variance = mean* a.k.a. poisson variance (solid black line). The plot function outputs the variances and that will be stored in the data set `meanVarPlot`.

```
meanVarPlot <- plotMeanVar( cds , show.raw.vars=TRUE ,
                             show.tagwise.vars=TRUE ,
                             show.binned.common.disp.vars=FALSE ,
                             show.ave.raw.vars=FALSE ,
                             dispersion.method = "qcm1" , NBline = TRUE ,
                             nbins = 100 ,
                             pch = 16 ,
                             xlab = "Mean Expression (Log10 Scale)" ,
                             ylab = "Variance (Log10 Scale)" ,
                             main = "Mean-Variance Plot" )
```