

```
1 Effects of temperature on EPEC expression
2 =====
3
4 Introduction
5 -----
6
7 This analysis describes differential expression analysis of an RNA-Seq data s
8 23 degrees C (C23EPEC) and 37 degrees C (T37EPEC)
9
10 The samples were from the laboratory of Christine White-Ziegler. You will be
11 instructions to create your own RStudio analysis with your own data.
12
13
14 ###Questions to answer in this data analysis
15 ###
16 * How similar were the samples within or across groups?
17 * How many genes were expressed in EPEC23 and in EPEC37?
18 * How many genes were differentially expressed due to the treatment, i.e., th
19 What were they?
20
21 Analysis
22 -----
23
24 To answer these questions, we'll import the edgeR package from the Bioconduct
25 perform differential expression analysis using the edgeR "classic" methods. A
26 manual and vignettes.
27
28 Note that this code assumes you've downloaded and installed edgeR. You also n
29 edgeR and limma, and then the library, instead of just running the chunk belo
30 enter, each line individually. For some reason, the library does not install
31 You have to try to install the library, wait for it to get stuck on installin
32 limma install line, then the library again.
33
34 ```{r message=FALSE}
35 source("http://bioconductor.org/biocLite.R")
36 biocLite("edgeR")
37 library(edgeR)
38 biocLite("limma")
39 library(edgeR)
40 ```
41
42 The edgeR user's guide is worth a look! Open it:
43
44 ```{r}
45 edgeRUsersGuide()
```

```
46 ```
47
48 Notice the guide contains many case studies illustrating different applicatio
49
50 See also
51
52 * http://www.bioconductor.org/packages/release/bioc/html/edgeR.html
53
54 ###About input data files
55 ###
56 We used Rockhopper to align the sample reads onto the EPEC genome, and count
57 overlapping annotated genes, requesting "Verbose Output" in the Rockhopper pa
58 the raw reads. The resulting data file for the sample data used for this wor
59 ~/RockhopperResults/EPEC_forRstudio/NC_011601_transcripts_forRstudio.txt.
60 We re-formatted this data from the 'transcripts.txt' file, produced by Rockho
61 Excel, sorting all data by 'Synonym', deleting all rows with no synonym, sort
62 'Translation Start' column, and copying and pasting the 'Synonym' column and
63 columns to a new workbook. Now close the original transcripts file and DO NO
64 was all to get rid of rows with the same name and keep in the same order as t
65 will use downstream, BUT we don't want to change the original file so we do n
66 Alternatively, you can duplicate the original so you don't risk messing it up
67 created file, we shortened the Raw Counts headings to <10 characters and inse
68 spaces (Easiest if you start all control names with C and all treatment names
69 this as a new tab delimited text .txt file with the same name but adding '_fo
70 this file in TextWrangler. Under 'View' go to 'Text Display' -> 'Show Invisib
71 and any tabs (depicted as triangles) that come before the first sample name.
72 this new file to perform differential expression analysis.
73
74
75 #####Read the data-put your new results file of interest into your folder yo w
76 #####otherwise, you will need to include the entire path name in the chunk bel
77 #####~/Rockhopper_Results/EPEC23vsEPEC37/NC_011601_transcripts_forRstudio.txt
78 #####filename is)
79 #####
80 ```{r}
81 d=read.delim('3newto2oldandnew_011601_transcripts_forRstudio.txt')
82 ```
83
84 View the first few rows of data:
85
86 ```{r}
87 head(d)
88 ```
89
90 Note that the data are nicely formatted, with gene names as row names and sam
```

```
91 names.  
92  
93 #####Make a DGEList object  
94 #####  
95 We'll use the DGEList function to create a DGEList object to contain the data  
96  
97 The DGEList function needs our table of counts (d) and a vector indicating wh  
98 belongs to. We also include an optional named argument (remove.zeros) that el  
99 zero counts. We eliminate genes with zero counts since it makes no sense to t  
100 differential expression if they were not expressed.  
101  
102 To learn more about the DGEList() function, type `help(DGEList)`  
103  
104 First, make a vector that indicates the group affiliation for each sample - C  
105 treatment.  
106  
107 ```{r}  
108 #define a new character vector using the "c" function  
109 group=c('C','C','C','C','T','T','T')  
110 ```  
111  
112 Then use the table of counts and the group affiliation vector to build the DG  
113  
114 ```{r}  
115 dge=DGEList(d,group=group,remove.zeros=TRUE)  
116 ```  
117  
118 Take a moment to look at the DGEList object:  
119  
120 ```{r}  
121 mode(dge) # find out what type it is  
122 dge # type it in the console to see what it is  
123 ```  
124  
125  
126 It's a list and it has two members: the original data frame that has been tur  
127 matrix and also an object called 'samples' that is also a list. Note that whe  
128 DGEList object, it used the groups object to assign group membership to each  
129 calculated the library sizes.  
130  
131 ###Normalization  
132 ###  
133 Some libraries had more reads than others and so counts per gene are not dire  
134 have to normalize. This is obvious and easy to understand. However, there is  
135 normalization that we need to take into account, which is that a treatment (l
```

136 increase in the expression of a subset of genes, thus "consuming" counts that
137 come from less highly expressed genes and making those genes appear down-regu
138 they were not.

139
140 For a deeper discussion, see [Robinson and Oshlack (2010) A scaling normaliza
141 differential expression analysis of RNA-Seq data](http://genomebiology.com/20
142

143 Note that when we normalize in edgeR, the software does not change the counts
144 object. Intead, it calculates *normalization factors* that will be used later

```
145  
146 ```{r}  
147 dge=calcNormFactors(dge)  
148 ```
```

149
150 Note that when we apply the calcNormFactors function to the dge object, it bo
151 the dge object while also changing it. Specifically, it added new information
152 case, it updated the normalization factors that were stored in the samples co
153

154 Look at the normalization factors - notice that the norm.factors now are bigg
155 depending on the read depth of the corresponding sample libraries:

```
156  
157 ```{r}  
158 dge$samples  
159 ```
```

160
161 **###Visualizing differences between samples**
162 **###**

163 We expect that 23 degree samples will more closely resemble other 23 degree s
164 samples, and vice versa. If they don't, this can indicate flaws in the experi
165 protocol. In general, the greater the separation between experimental groups,
166 of differentially expressed genes.

167
168 Two types of plots are useful in this step - MDS plots and hierarchical clust
169

170 **####MDS Plot**
171 **####**

172 A multi-dimensional scaling (MDS) plot can show similarity between samples.

173
174 From the edgeR guide:

```
175  
176 >>>The function plotMDS draws a multi-dimensional scaling plot of the RNA sam  
177 >>>correspond to leading log-fold-changes between each pair of RNA samples. T  
178 >>>log-fold-change is the average (root-mean-square) of the largest absolute  
179 >>>between each pair of samples. This plot can be viewed as a type of unsuper  
180 >>>
```

```
181 Create the plot:
182
183 ```{r}
184 #color for controls
185 cn.color='blue'
186 #color for treatments
187 tr.color='brown'
188 #define a title for the plot
189 main='MDS Plot for Count Data'
190 #par(las=1) # makes y axis labels horizontal not vertical
191 colors=c(rep(cn.color,4),rep(tr.color,3))
192 plotMDS(dge,main=main,labels=colnames(dge$counts),
193         col=colors,las=1)
194 ```
195
196 The second dimension does a good job of separating samples but often is not
197 their similarities or differences if there is a low N. Let's see what Hierar
198 do.
199
200 ###Hierarchical clustering
201 ###
202 You can use hierarchical clustering to view the relationships between samples
203 calculates a distance value between each sample using measurements from all t
204 then clustered according to how similar they are with respect to that distanc
205
206 Note we need to use normalized "counts per million" instead of the raw count
207 that we have to transpose the matrix so that columns become rows and rows bec
208 words, to do the clustering, the entities that should be clustered need to be
209
210 ```{r}
211 normalized.counts=cpm(dge)
212 transposed=t(normalized.counts) # transposes the counts matrix
213 distance=dist(transposed) # calculates distance
214 clusters=hclust(distance) # does hierarchical clustering
215 plot(clusters) # plots the clusters as a dendrogram
216 ```
217
218 Looks like the samples grown at 37 degrees clustered together and away from s
219 degrees, which formed their own cluster as expected. This is a useful plot f
220 or something inherently wrong with the experiment overall. We will continue a
221 studies, if you see different groups clustering together or a sample off on i
222 to eliminate samples (or verify their labelling throughout the experiment) an
223
224 ###Differential expression analysis
225 ###
```

```
226 Now let's use functions in edgeR to identify genes whose expression changed d
227 change. For our analysis, we'll treat all the 23 and 37 degree samples as rep
228 variables as well we would ignore this for the first pass and treat all sampl
229 condition as a replicate. This corresponds to what the edgeR manual calls th
230
231 An alternative approach would be to use edgeR's generalized linear models met
232 differential expression analysis that takes batch effects (batches collected
233 into account. The Arabidopsis study in the User's Guide gives a nice example
234 this Markdown, however, we'll stick with the classic approach for demonstrati
235
236 #####Estimating dispersion
237 #####
238 To start, we need to estimate dispersion, which reflects the degree to which
239 depends on expression level. We need this to model gene expression and test w
240 has changed due to the treatment.
241
242 The details of how this works are explained in the edgeR user's guide. In a n
243 is that we use the negative binomial model to estimate a dispersion parameter
244 edgeR calls "the degree of inter-library variability" for that particular gen
245
246 Using edgeR to calculate common and tagwise dispersions:
247
248 ```{r}
249 dge=estimateCommonDisp(dge)
250 dge=estimateTagwiseDisp(dge)
251 ```
252
253 Now you can look again at the dge object – notice there are now new component
254 common and tagwise dispersions.
255
256 #####Differential expression
257 #####
258 Now that we have estimates for dispersion for every gene, we can use the exac
259 each gene individually using the exact negative binomial test.
260 ```{r}
261 dex=exactTest(dge,pair=c("C","T"),dispersion="tagwise")
262 #design = model.matrix(~group)
263 #dge=estimateGLMCommonDisp(dge,design)
264 #dge=estimateGLMTrendedDisp(dge,design)
265 #dge=estimateGLMTagwiseDisp(dge,design)
266 #fit=glmFit(dge,design)
267 #dex=glmLRT(fit,coef=2)
268 ```
269
270 *Note* This is a "classic" edgeR analysis according to the User's Guide. To t
```

```
271 linear modeling) edgeR analysis, comment the first line and uncomment the rem
272
273 The exactTest function returns a new DGEXact test object that we can pass to
274 Now we have a new object (dex) that captures the results from testing for dif
275 genes across the 23 degree and 37 degree samples. Type it into the console to
276 Notice that it has logFC, logCPM (counts per million), and PValue for each ge
277
278 #####Multiple Hypothesis testing correction
279 #####
280 Even if the data were totally random, around 5% of the p values would be 0.05
281 would be 0.10 or less, and so on. This is because p values reflect the probab
282 given result by chance. So even if our data were random, we would expect to o
283 positives if we use a p value threshold of 0.05 to determine significance. Th
284 multiple hypothesis testing problem which afflicts testing highly parallel, h
285 correct for this, we can calculate a false discovery rate (FDR) value for eac
286 values as inputs.
287
288 ```{r}
289 #use Benjamini-Hochberg method of calculated false discovery
290 #rate for each gene
291 fdrvalues=p.adjust(dex$table$PValue, method='BH')
292 dex$table$fdr=fdrvalues
293 ```
294
295 #####Picking an FDR cutoff
296 #####
297 For subsequent steps, we'll need to select an FDR cutoff for determination of
298 expression. This is somewhat arbitrary, but unavoidable for some analyses. Fo
299 Ontology enrichment analysis using G0Seq, we have to designate some genes as
300 DE.
301
302 Use the summary and decideTestsDGE commands to find out how many genes are up
303 using different Benjamini-Hockberg adjusted FDR cutoffs:
304
305 ```{r}
306 summary(decideTestsDGE(dex,p=0.05))
307 summary(decideTestsDGE(dex,p=0.01))
308 summary(decideTestsDGE(dex,p=0.005))
309 ```
310
311 Based on this, pick a cutoff:
312
313 ```{r}
314 cutoff=0.01
315 ```
```

```
316
317 ###Get an overview of the DE genes in the data set
318 ###
319 The plotSmear function illustrates the relationship between logFC and average
320 the differentially expressed genes, which are highlighted in red. The blue lines
321 are up or downregulated two fold.
322
323 ```{r}
324 de = decideTestsDGE(dex, p = cutoff, adjust = "BH")
325 detags = rownames(dex)[as.logical(de)]
326 plotSmear(dex, de.tags = detags)
327 abline(h = c(-1, 1), col = "blue")
328 ```
329
330 Not surprisingly, application of 37 degrees to EPEC results in differential
331 genes, hopefully many involved in virility and increased metabolism, and for
332 fold-change is greater than two.
333
334 Writing results files
335 =====
336
337 For the next steps, we'll write files containing lists of differentially expressed
338 then load into other programs like EcoCyc. Just run the following code so that
339 file. We are just using our .ptt file from our replicon folder. But ptt is a
340 it works for this exercise, you may find that many sequence feature files come
341 and many programs want to use that. A great little converter is BEDOPS. It
342 bed, that you can use here instead of the ptt file or use for further analysis
343 need to install BEDOPS on your Mac (follow instructions on the BEDOPS site).
344 a script called 'convert2bed'. Once installed, on your terminal you cd into
345 say a gff file. Then type:
346
347 convert2bed -i gff < /Users/lbierwer/Desktop/RtestOnRockhopperResults/NC_011601
348 /Users/lbierwer/Desktop/RtestOnRockhopperResults/NC_011601.bed
349
350 Here I told it the format I had, then my input file name with full path (on
351 drop the file), and what I wanted my output file name to be once converted.
352
353 ```{r}
354 #get normalized counts per million
355 cpms=cpm(dge$counts)
356 #find out which columns have controls
357 cn=grep('C',colnames(cpms))
358 #find out which columns have treatments
359 tr=grep('T',colnames(cpms))
360 #calculate the average expression for control samples
```



```
361 ave.cn=apply(cpms[,cn],1,mean)
362 #calculate average expression for treatment samples
363 ave.tr=apply(cpms[,tr],1,mean)
364 #make a data frame
365 res=data.frame(synonym=row.names(dex$table),
366               fdr=dex$table$fdr,
367               logFC=dex$table$logFC,
368               Cn=ave.cn,
369               Tr=ave.tr)
370 #read gene information (originally from the replicon info we gave Rockhopper
371 #this is a "sequence feature table" file
372 annots_file='NC_011601.ptt'
373 #keep gene id and gene description column
374 annots=read.delim(annots_file,sep='\t',header=F)[,5:6]
375 #name the column
376 names(annots)=c('gene','synonym')
377 #combine gene expression and annotations
378 res=merge(res,annots,by.x='synonym',by.y='synonym')
379 #put the results in order of significance (fdr)
380 res=res[order(res$fdr),]
381 #put columns in an order easy to browse
382 res=res[,c('fdr','logFC','Cn','Tr','gene','synonym')]
383 #write DE genes to a file:
384 out_file='results/3newto2oldandnew_011601_EPEC_DE.txt'
385 #select just the rows that made the cutoff
386 de=res$fdr<=cutoff
387 #write DE genes only
388 write.table(res[de,],file=out_file,row.names=F,sep='\t',quote=F)
389 #Make a file we can load into EcoCyc for pathways visualization
390 out_file='results/3newto2oldandnew_011601_forEcoCyc.txt'
391 write.table(res[de,c('synonym','logFC')],file=out_file,quote=FALSE,
392           sep='\t',col.names=FALSE,row.names=FALSE)
393 #write all DE genes
394 out_file='results/3newto2oldandnew_011601_EPEC_AllDE.txt'
395 write.table(res,file=out_file,row.names=F,sep='\t',quote=F)
396 ```
397
398 Summarizing the results
399 -----
400 Check out your tables. Cool tip for viewing big spreadsheets in excel: open
401 furthest right header go to 'Window' -> 'Freeze Pane' Now as you scroll thr
402 put. Go to 'Window' -> 'Unfreeze Pane' to stop.
403
404 Conclusion
405 -----
```

```
406 Of `r nrow(d)` annotated genes, there were `r nrow(dge)` genes with at least
407 Around `r round(nrow(dge)/nrow(d)*100,1)` % of EPEC37 was expressed in EPEC2
408 * At fdr of `r cutoff`, we observed `r sum(de)` differentially expressed genes
409
```

410 Limitations of the analysis

411 -----

- 412
- 413 * The FDR cutoff we picked might not be the most appropriate for your data.
- 414 * If the number of fragment counts per library is low. Power to detect differentially expressed genes also be low, leading to false negatives.
- 415
- 416 * The assumption about that collection time had no effect could be wrong. To account for this, we used generalized linear modeling features (glm functions in edgeR) to test for differentially expressed genes while conditioning on collection time.

417

419 Session information

420 -----

```
421
422
423 ```{r}
424 sessionInfo()
425 ```
```