

# baySeq: Empirical Bayesian analysis of patterns of differential expression in count data

Thomas J. Hardcastle

October 17, 2014

## 1 Introduction

---

This vignette is intended to give a rapid introduction to the commands used in implementing two methods of evaluating differential expression in Solexa-type, or *count* data by means of the baySeq R package. For fuller details on the methods being used, consult Hardcastle & Kelly [1]. The major improvement made in this release is the option to include region length in evaluating differential expression between genomic regions (e.g. genes). See Section 6.3.1 for more details.

We assume that we have discrete data from a set of sequencing or other high-throughput experiments, arranged in a matrix such that each column describes a sample and each row describes some entity for which counts exist. For example, the rows may correspond to the different sequences observed in a sequencing experiment. The data then consists of the number of times each sequence is observed in each sample. We wish to determine which, if any, rows of the data correspond to some patterns of differential expression across the samples. This problem has been addressed for pairwise differential expression by the edgeR [3] package.

However, baySeq takes an alternative approach to analysis that allows more complicated patterns of differential expression than simple pairwise comparison, and thus is able to cope with more complex experimental designs. We also observe that the methods implemented in baySeq perform at least as well, and in some circumstances considerably better than those implemented in edgeR [1].

baySeq uses empirical Bayesian methods to estimate the posterior likelihoods of each of a set of models that define patterns of differential expression for each row. This approach begins by considering a distribution for the row defined by a set of underlying parameters for which some prior distribution exists. By estimating this prior distribution from the data, we are able to assess, for a given model about the relatedness of our underlying parameters for multiple libraries, the posterior likelihood of the model.

In forming a set of models upon the data, we consider which patterns are biologically likely to occur in the data. For example, suppose we have count data from some organism in condition  $A$  and condition  $B$ . Suppose further that we have two biological replicates for each condition, and hence four libraries  $A_1, A_2, B_1, B_2$ , where  $A_1, A_2$  and  $B_1, B_2$  are the replicates. It is reasonable to suppose that at least some of the rows may be unaffected by our experimental conditions  $A$  and  $B$ , and the count data for each sample in these rows will be *equivalent*. These data need not in general be identical across each sample due to random effects and different library sizes, but they will share the same underlying parameters. However, some of the rows may be influenced by the different experimental conditions  $A$  and  $B$ . The count data for the samples  $A_1$  and  $A_2$  will then be equivalent, as will the count data for the samples  $B_1$  and  $B_2$ . However, the count data between samples  $A_1, A_2, B_1, B_2$  will not be equivalent. For such a row, the data from samples  $A_1$  and  $A_2$  will then share the same set of underlying parameters, the data from samples  $B_1$  and  $B_2$  will share the same set of underlying parameters, but, crucially, the two sets will not be identical.

Our task is thus to determine the posterior likelihood of each model for each row of the data.

## 2 Preparation

---

We begin by loading the baySeq package.

```
> library(baySeq)
```

Note that because the experiments that baySeq is designed to analyse are usually massive, we should use (if possible) parallel processing as implemented by the snow package. We use the parallel package (if it exists), and define a cluster. If parallel is not present, we can proceed anyway with a NULL cluster. Results may be slightly different depending on whether or not a cluster is used owing to the non-deterministic elements of the method.

```
> if(require("parallel")) cl <- makeCluster(8) else cl <- NULL
```

We load a simulated data set consisting of count data on one thousand counts.

```
> data(simData)
> simData[1:10,]
      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9] [,10]
[1,]    4    1    5    2    3    0    1    1    1    0
[2,]    1    0    9    6    5    0    1    0    0    1
[3,]    9    2    5    5   14    2    3    1    0    4
[4,]    7    3    8    2    2    0    1    0    1    0
[5,]    2    2    4    7    0    0    0    0    0    1
[6,]    2    1    0    1    0    4    3    5    5    3
[7,]    9    8    8    8    9    1    2    1    0    0
[8,]    9    5    7    8    7    1    2    0    1    2
[9,]    6    2    2    3    0    0    0    0    0    0
[10,]   1    0    2    0    1    3   17    2    2   10
```

The data are simulated such that the first hundred counts show differential expression between the first five libraries and the second five libraries. Our replicate structure, used to estimate the prior distributions on the data, can thus be defined as

```
> replicates <- c("simA", "simA", "simA", "simA", "simA",
+                "simB", "simB", "simB", "simB", "simB")
```

We can also establish two group structures for the data.

Each member (vector) contained within the 'groups' list corresponds to one model upon the data. In this setting, a model describes the patterns of data we expect to see at least some of the tags correspond to. In this simple example, we expect that some of the tags will be equivalently expressed between all ten libraries. This corresponds to the 'NDE' model, or vector  $c(1,1,1,1,1,1,1,1,1,1)$  - all libraries belong to the same group for these tags.

We also expect that some tags will show differential expression between the first five libraries and the second five libraries. For these tags, the two sets of libraries belong to different groups, and so we have the model 'DE', or vector  $c(1,1,1,1,1,2,2,2,2,2)$  - the first five libraries belong to group 1 and the second five libraries to group 2. We thus have the following group structure

```
> groups <- list(NDE = c(1,1,1,1,1,1,1,1,1,1),
+               DE = c(1,1,1,1,1,2,2,2,2,2))
```

In a more complex experimental design (Section ??) we might have several additional models. The key to constructing vectors corresponding to a model is to see for which groups of libraries we expect equivalent expression of tags.

We note that the group for DE corresponds to the replicate structure. This will often be the case, but need not be in more complex experimental designs.

The ultimate aim of the baySeq package is to evaluate posterior likelihoods of each model for each row of the data.

We begin by combining the count data and user-defined groups into a countData object.

```
> CD <- new("countData", data = simData, replicates = replicates, groups = groups)
```

Library sizes can be inferred from the data if the user is not able to supply them.

```
> libsizes(CD) <- getLibsizes(CD)
```

We can then plot the data in the form of an MA-plot, suitable modified to plot those data where the data are uniformly zero (and hence the log-ratio is infinite) (Figure 1). Truly differentially expressed data can be identified in the plot by coloring these data red, while non-differentially expressed data are colored black.

```
> plotMA.CD(CD, samplesA = "simA", samplesB = "simB",
+           col = c(rep("red", 100), rep("black", 900)))
```

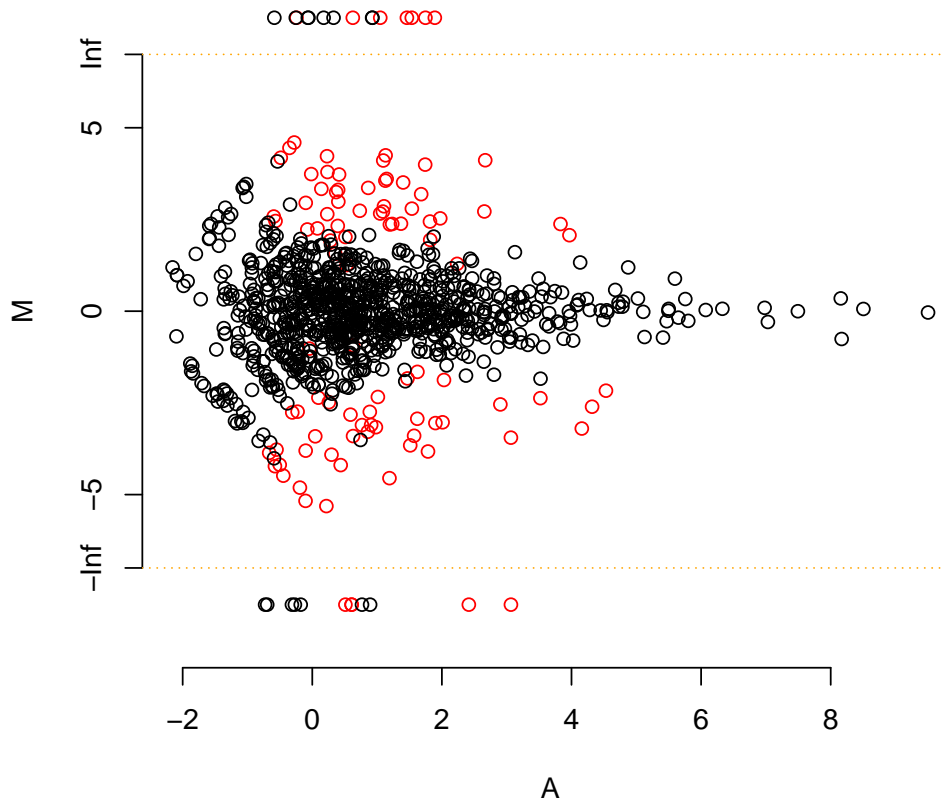


Figure 1: 'MA'-plot for count data. Where the log-ratio would be infinite (because the data in one of the sample groups consists entirely of zeros, we plot instead the log-values of the other group). Truly differentially expressed data are colored red, and non-differentially expressed data black.

We can also optionally add annotation details into the `@annotation` slot of the `countData` object.

```
> CD@annotation <- data.frame(name = paste("count", 1:1000, sep = "_"))
```

### 3 Negative-Binomial Approach

We first estimate an empirical distribution on the parameters of the Negative Binomial distribution by bootstrapping from the data, taking individual counts and finding the quasi-likelihood parameters for a Negative Binomial distribution. By taking a sufficiently large sample, an empirical distribution on the parameters is estimated. A sample size of around 10000 iterations is suggested, depending on the data being used), but 1000 is used here to rapidly generate the plots and tables.

```
> CD <- getPriors.NB(CD, samplesize = 1000, estimation = "QL", cl = cl)
```

The calculated priors are stored in the `@priors` slot of the `countData` object produced as before. For the negative-binomial method, we are unable to form a conjugate prior distribution. Instead, we build an empirical prior distribution which we record in the list object `$priors` of the slot `@priors`. Each member of this list object corresponds to one of the models defined by the `group` slot of the `countData` object and contains the estimated parameters for each of the individual counts selected under the models. The vector `$sampled` contained in the slot `@priors` describes which rows were sampled to create these sets of parameters.

We then acquire posterior likelihoods, estimating the proportions of differentially expressed counts.

```
> CD <- getLikelihoods(CD, pET = 'BIC', c1 = c1)
```

```
.
```

```
> CD@estProps
```

```
numeric(0)
```

```
> CD@posteriors[1:10,]
```

	NDE	DE
[1,]	-0.6379573	-0.751561829
[2,]	-0.9002858	-0.521639691
[3,]	-0.7711396	-0.620799948
[4,]	-2.2055580	-0.116746224
[5,]	-0.6307673	-0.759678681
[6,]	-1.1445501	-0.383264074
[7,]	-5.5615869	-0.003850076
[8,]	-3.8695783	-0.021087963
[9,]	-0.8983994	-0.522933530
[10,]	-1.6797002	-0.206323141

```
> CD@posteriors[101:110,]
```

	NDE	DE
[1,]	-6.066791e-02	-2.832521
[2,]	-7.991582e-05	-9.434577
[3,]	-4.640167e-02	-3.093531
[4,]	-1.701222e-02	-4.082318
[5,]	-4.290052e-03	-5.453601
[6,]	-5.109052e-02	-2.999593
[7,]	-7.850153e-02	-2.583631
[8,]	-4.885146e-02	-3.043297
[9,]	-6.376703e-02	-2.784233
[10,]	-8.884361e-03	-4.727902

Here the assumption of a Negative Binomial distribution with priors estimated by maximum likelihood gives an estimate of

```
[1] NA
```

as the proportion of differential expressed counts in the simulated data, where in fact the proportion is known to be 0.1.

## 4 Results

---

We can ask for the top candidates for differential expression using the `topCounts` function.

```
> topCounts(CD, group = "DE")
```

	annotation	simA.1	simA.2	simA.3	simA.4	simA.5	simB.1	simB.2	simB.3	simB.4	simB.5
1	count_80	1	1	0	1	1	13	21	8	6	20
2	count_78	1	1	0	1	1	8	13	7	9	10
3	count_66	0	0	0	0	0	15	10	4	4	10
4	count_21	2	0	1	1	0	15	15	6	5	11
5	count_7	9	8	8	8	9	1	2	1	0	0
6	count_26	13	4	11	5	7	1	1	1	0	0
7	count_72	0	0	1	0	0	7	6	4	3	8
8	count_64	6	6	8	11	9	1	1	0	0	1
9	count_83	14	6	9	2	9	1	0	0	1	1
10	count_27	5	3	6	4	7	0	0	0	1	0

	Likelihood	ordering	FDR.DE	FWER.DE
1	0.9988118	2>1	0.001188193	0.001188193
2	0.9985328	2>1	0.001327718	0.002653692
3	0.9977532	2>1	0.001634062	0.004894480
4	0.9968802	2>1	0.002005505	0.007999046
5	0.9961573	1>2	0.002372939	0.011810982
6	0.9949830	1>2	0.002813616	0.016768728
7	0.9938530	2>1	0.003289813	0.022812642
8	0.9918470	1>2	0.003897714	0.030779670
9	0.9891786	1>2	0.004667009	0.041267960
10	0.9861448	1>2	0.005585832	0.054551422

We can plot the posterior likelihoods against the log-ratios of the two sets of samples using the `plotPosteriors` function, coloring the truly differentially expressed data red and the non-differentially expressed data black (Figure 2).

```
> plotPosteriors(CD, group = "DE", col = c(rep("red", 100), rep("black", 900)))
```

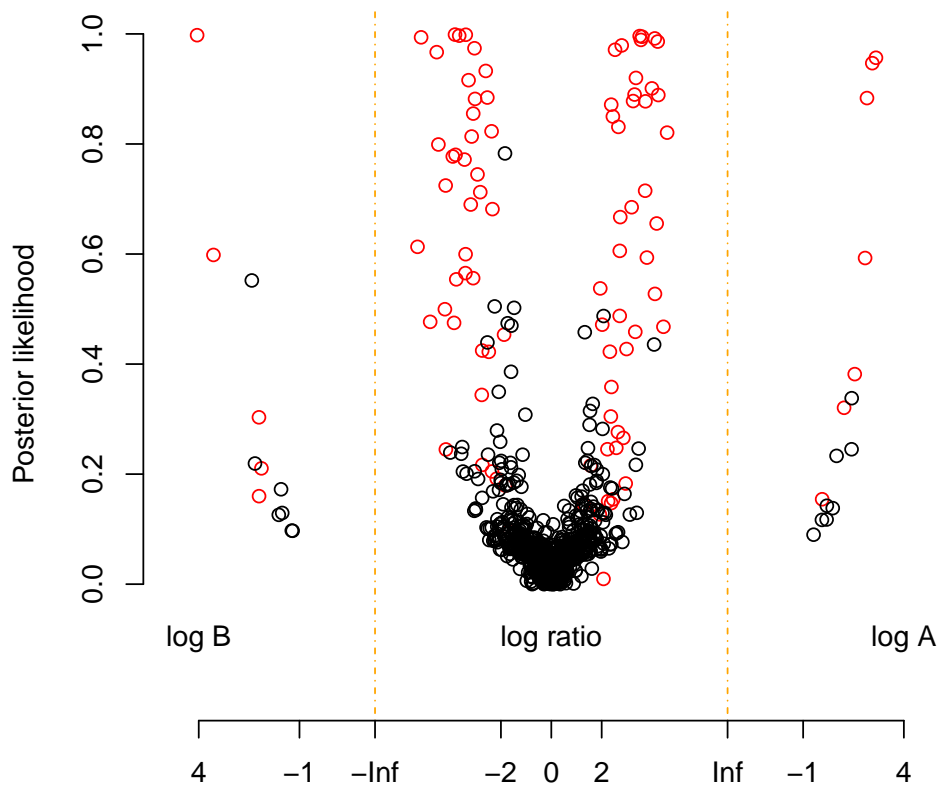


Figure 2: Posterior likelihoods of differential expression against log-ratio (where this would be non-infinite) or log values (where all data in the other sample group consists of zeros). Truly differentially expressed data are colored red, and non-differentially expressed data black.

## 5 Paired Data Analysis

There now exists functionality to analyse paired data through a similar process, using the beta-binomial distribution. The process for analysing paired data follows approximately the same steps as for analysing unpaired data. However, two different types of differential expression can exist within paired data. Firstly, we can find differential expression between replicate groups, as before. However, we can also find (consistent) differential expression between pairs; this would occur when for a single row of data, the first member of each pair differs from the second member of each pair. baySeq can identify both these types of differential expression simultaneously, and we implement this procedure below.

We begin by loading a simulated dataset containing counts for four paired datasets.

```
> data(pairData)
```

The first four columns in these data are paired with the second four columns. We construct a pairedData in a similar fashion to the countData object.

```
> pairCD <- new("countData", data = list(pairData[,1:4], pairData[,5:8]),
+       replicates = c(1,1,2,2),
+       groups = list(NDE = c(1,1,1,1), DE = c(1,1,2,2)),
+       densityFunction = bbDensity)
```

We can find the library sizes for the data with the getLibsizes function.

```
> libsizes(pairCD) <- getLibsizes(pairCD)
```

We estimate an empirical distribution on the parameters of a beta-binomial distribution by bootstrapping from the data, taking individual counts and finding the maximum likelihood parameters for a beta-binomial distribution. By taking a sufficiently large sample, an empirical distribution on the parameters is estimated. A sample size of around 10000 iterations is suggested, depending on the data being used), but 1000 is used here to rapidly generate the plots and tables.

```
> pairCD <- getPriors(pairCD, samplesize = 1000, cl = cl)
```

We then acquire posterior likelihoods as before. The use of 'nullData = TRUE' in this context allows us to identify pairs which show no differential expression between replicate groups, but does show deviation from a one-to-one ratio of data between pairs.

```
> pairCD <- getLikelihoods(pairCD, pET = 'BIC', nullData = TRUE, cl = cl)
```

.

We can ask for the top candidates for differential expression between replicate groups using the topCounts function as before.

```
> topCounts(pairCD, group = 2)
```

rowID	X1.1	X1.2	X2.1	X2.2	Likelihood	ordering	FDR.DE	FWER.DE	
1	5	159:73	44:24	0:49	0:68	0.9968336	1>2	0.003166442	0.003166442
2	35	53:12	19:7	0:77	0:6	0.9931146	1>2	0.005025897	0.010029992
3	53	709:0	895:0	373:191	124:60	0.9929691	1>2	0.005694233	0.016990377
4	96	25:0	73:0	8:3	36:13	0.9882965	1>2	0.007196550	0.028495032
5	65	80:0	48:0	36:50	12:3	0.9840703	1>2	0.008943187	0.043970848
6	24	63:0	21:0	47:80	6:13	0.9829191	1>2	0.010299478	0.060300717
7	90	268:0	39:0	74:107	98:36	0.9753827	1>2	0.012344880	0.083433566
8	71	8:0	15:0	21:16	2:1	0.9434958	1>2	0.017864789	0.135223381
9	50	43:19	44:46	106:6	133:5	0.9315488	2>1	0.023485504	0.194418397
10	45	3:1	5:9	7:0	24:0	0.9314183	2>1	0.027995125	0.249666566

However, we can also look for consistent differential expression between the pairs.

```
> topCounts(pairCD, group = 1)
```

rowID	X1.1	X1.2	X2.1	X2.2	Likelihood	FDR.NDE	FWER.NDE	
1	116	17:70	1:40	9:117	3:45	0.9768200	0.02317999	0.02317999
2	146	1:38	0:68	0:28	0:26	0.9647937	0.02919312	0.05757016

3	123	1:4	1:11	0:5	1:14	0.9553069	0.03435978	0.09969028
4	193	69:1	10:1	119:17	53:5	0.9458452	0.03930854	0.14844637
5	101	0:30	0:5	0:60	0:24	0.9423483	0.04297717	0.19753987
6	138	0:12	0:4	0:4	0:13	0.9422982	0.04543128	0.24384330
7	144	0:4	0:21	0:2	0:12	0.9416347	0.04727899	0.28797660
8	127	0:3	0:12	0:15	0:4	0.9414991	0.04868174	0.32963063
9	180	1:2	1:16	2:41	0:2	0.9409184	0.04983727	0.36923710
10	118	0:6	0:31	0:5	0:6	0.9406590	0.05078764	0.40666717

## 6 Case Study: Analysis of sRNA-Seq Data

---

### 6.1 Introduction

We will look at data sequenced from small RNAs acquired from six samples of root stock from *Arabidopsis thaliana* in a grafting experiment [2]. Three different biological conditions exist within these data; one in which a Dicer 2,3,4 triple mutant shoot is grafted onto a Dicer 2,3,4 triple mutant root (**SL236** & **SL260**), one in which a wild-type shoot is grafted onto a wild-type root (**SL239** & **SL240**), and one in which a wild-type shoot is grafted onto a Dicer 2,3,4 triple mutant root (**SL237** & **SL238**). Dicer 2,3,4 is required for the production of 22nt and 24nt small RNAs, as well as some 21nt ones. Consequently, if we detect differentially expressed sRNA loci in the root stock of the grafts, we can make inferences about the mobility of small RNAs.

### 6.2 Reading in data

The data and annotation are stored in two text files. We can read them in using **R**'s standard functions.

```
> data(mobData)
> data(mobAnnotation)
```

### 6.3 Making a countData object

We can create a `countData` object containing all the information we need for a first attempt at a differential expression analysis.

#### 6.3.1 Including lengths

If two genes are expressed at the same level, but one is twice the length of the other, then (on average) we will sequence twice as many reads from the longer gene. The same is true for sRNA loci, and so in these analyses it is often useful to include the lengths of each feature. The lengths can be derived from the annotation of each feature, but we need to explicitly declare them within the 'countData' object.

```
> seglens <- mobAnnotation$end - mobAnnotation$start + 1
> cD <- new("countData", data = mobData, seglens = seglens, annotation = mobAnnotation)
```

Determining the best library scaling factor to use is a non-trivial task. The simplest approach would be to use the total number of sequenced reads aligning to the genome. However, this approach means that a few sequences that appear at very high levels can drastically skew the size of the scaling factor. Bullard *et al* suggest that good results can be obtained by taking the sum of the reads below the  $n$ th percentile of the data.

```
> libsizes(cD) <- getLibsizes(cD, estimationType = "quantile")
```

### 6.4 Pairwise Differential Expression

We start by looking at a pairwise differential expression analysis between two of the sample types. The analysis between samples 'SL236', 'SL260' and 'SL237', 'SL238' should be a first step in discovering sRNA loci associated with mobility.

We begin by selecting a subset of the available data:

```
> cDPair <- cD[,1:4]
```

We then need to define the replicate structure of the countData object. We do this by creating a vector that defines the replicate group that each sample belongs to.

```
> replicates(cDPair) <- as.factor(c("D3/D3", "D3/D3", "WT/D3", "WT/D3"))
```

We next need to define each of the models applicable to the data. In the first case, it is reasonable to suppose that at least some of the loci will be unaffected by the different experimental conditions prevailing in our replicate groups, and so we create one model of no differential expression.

We do this by defining a vector NDE.

```
> NDE <- c(1,1,1,1)
```

Each member of the NDE vector represents one sample in our experiment. By giving each item in the NDE vector the same number, we indicate that, under the hypothesis of no differential expression, all the samples belong to the same group.

We may also conjecture that some of the loci will be affected depending on whether the shoot is a Dicer mutant or a wild-type *Arabidopsis* sample.

```
> mobile <- c("non-mobile", "non-mobile", "mobile", "mobile")
```

This vector indicates that the third and fourth samples, which consist of the wild-type shoot samples, are in a separate expression group to the first and second samples, corresponding to the Dicer 2,3,4 mutant shoot.

We can now add these models to the locus data by modifying the @groups slot

```
> groups(cDPair) <- list(NDE = NDE, mobile = mobile)
```

Now that we have defined our models, we need to establish prior distributions for the data. We do this using the getPriors.NB function.

```
> cDPair <- getPriors.NB(cDPair, samplesize = 1e4, cl = cl)
```

The accuracy of the distribution is determined by the number of data points used to estimate the distribution; the 'samplesize'. Here we've used a small sample size to reduce the computational effort required, but higher values will give more accurate results (the default is 1e5).

Having found prior distributions for the data, we can identify posterior likelihoods for the data using the getLikelihoods function. Before we do this, however, it is worth considering the possibility that some loci will not be expressed at all in our data.

### 6.4.1 Null Data

We first examine the priors to see if any 'null' data, consisting of un-expressed sRNA loci, are present. If the distribution of priors for the non-differentially expressed group is bimodal, it is likely that some of the loci are expressed at substantially lower levels than others.

```
> plotNullPrior(cDPair)
```

There is some evidence for bimodality, with a small peak of lowly expressed data to the left of the distribution.

We can use the nullData = TRUE option in the getLikelihoods function to allow for the possibility that some of the loci are miscalled in our locus map, and should properly be identified as nulls.

```
> cDPair <- getLikelihoods(cDPair, nullData = TRUE, cl = cl)
```

If we now look at the cDPair object, we can see that we have acquired posterior likelihoods for the data

```
> cDPair
```

```
An object of class "countData"
3000 rows and 4 columns
```

```
Slot "replicates"
```



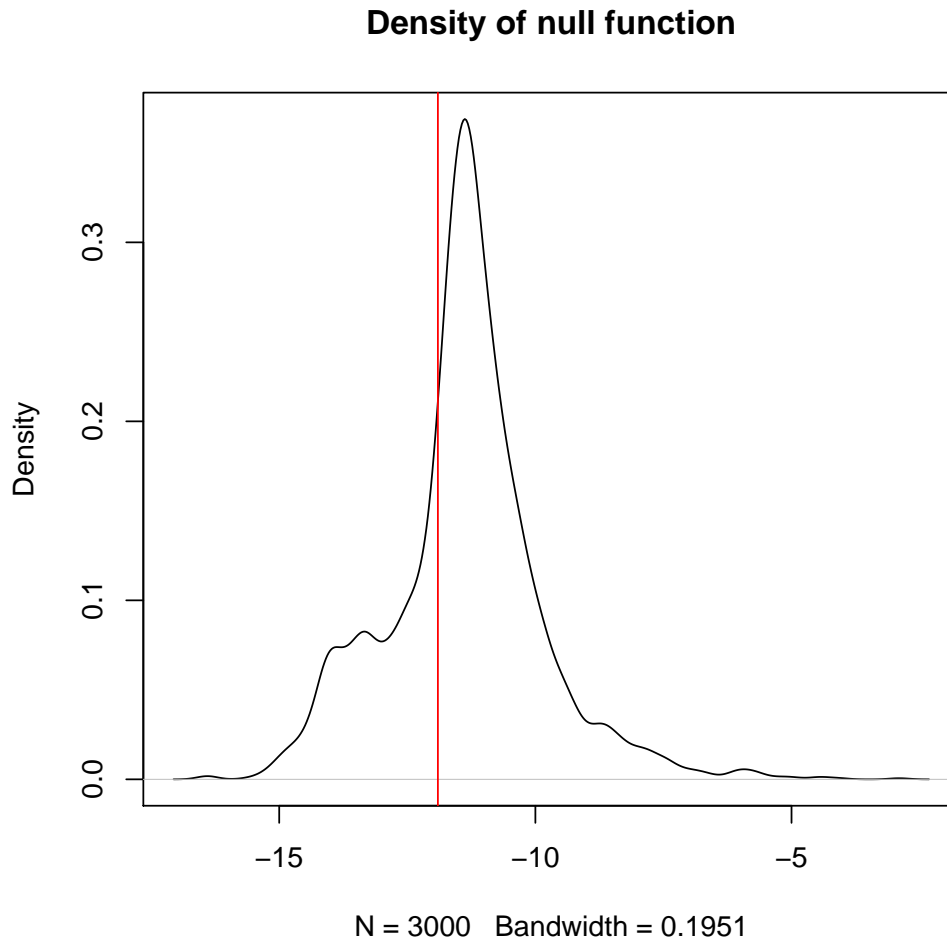


Figure 3: Distribution of  $\mu_{ij}$ . Bimodality suggests the presence of 'null', or un-expressed, data.

```
D3/D3 D3/D3 WT/D3 WT/D3
Slot "groups":
$NDE
[1] 1 1 1 1
Levels: 1
```

```
$mobile
[1] non-mobile non-mobile mobile mobile
Levels: mobile non-mobile
```

```
Slot "data":
  SL236 SL260 SL237 SL238
[1,]    0    0    0    0
[2,]   18   21    1    5
[3,]    1    2    2    3
[4,]   68   87  270  184
[5,]   68   87  270  183
2995 more rows...
```

```
Slot "annotation":
  chr start end
```

```

1 1 789 869
2 1 8641 8700
3 1 10578 10599
4 1 17041 17098
5 1 17275 17318
2995 more rows...
Slot "posteriors":
      NDE      mobile
[1,] 0.0007287778 0.03940907
[2,] 0.1424642330 0.85703432
[3,] 0.8033025998 0.12854417
[4,] 0.3244638539 0.67553615
[5,] 0.4276421287 0.57235787
2995 more rows...

```

The estimated posterior likelihoods for each model are stored in the natural logarithmic scale in the `@posteriors` slot of the `countDataPosterior` object. The  $n$ th column of the posterior likelihoods matrix corresponds to the  $n$ th model as listed in the `group` slot of `CDPair`. In general, what we would like to do with this information is form a ranked list in which the loci most likely to be differentially expressed are at the top of the list.

Try looking at the proportions of data belonging to each group. Note that these no longer sum to 1, as some data are now classified as 'null'.

```

> summarisePosteriors(cD)
numeric(0)

```

The value contained in the `@estProps` slot is a best-guess figure for the proportion of data belonging to each model defined by the `@groups` slot. In this case, it is estimated that approximately 65% of the loci are not differentially expressed, while 35% are differentially expressed. These estimates should not be relied upon absolutely, but are a useful indicator of the global structure of the data.

We can ask for the rows most likely to be differentially expressed under our different models using the `topCounts` function. If we look at the second model, or grouping structure, we see the top candidates for differential expression. Because the library sizes of the different libraries differ, it can be unclear as to why some loci are identified as differentially expressed unless the data are normalised.

```

> topCounts(cDpair, group = 2, normaliseData = TRUE)
  chr  start      end SL236 SL260 SL237 SL238 Likelihood      ordering
1  1  447231  447298     0     0   174   174  0.9999683  mobile>non-mobile
2  1  8287590 8287674     0     0    85    79  0.9997268  mobile>non-mobile
3  1  9254068 9254167     0     0    69    66  0.9995797  mobile>non-mobile
4  1 13463357 13463459    10    10   109   110  0.9990213  mobile>non-mobile
5  1 11140107 11140158     0     0    78    59  0.9988097  mobile>non-mobile
6  1  6880517 6880553     0     0    75    57  0.9984281  mobile>non-mobile
7  1  6127755 6127808     0     0   102    63  0.9983938  mobile>non-mobile
8  1  5056092 5056161    80   132     1     0  0.9983935  non-mobile>mobile
9  1  2157113 2157287    13    11   291   187  0.9983082  mobile>non-mobile
10 1 14188044 14188079     2     0    89    74  0.9981891  mobile>non-mobile
      FDR.mobile  FWER.mobile
1  3.173037e-05  3.173037e-05
2  1.524783e-04  3.049479e-04
3  2.417443e-04  7.250962e-04
4  4.259944e-04  1.703131e-03
5  5.788587e-04  2.891420e-03
6  7.443651e-04  4.458772e-03
7  8.674774e-04  6.057762e-03
8  9.598502e-04  7.654490e-03
9  1.041180e-03  9.333361e-03
10 1.118155e-03  1.112739e-02

```

Observe how the data change in the normalised results; the effect is particularly noticeable in the SL236 and SL260 datasets, in which the normalised data is much less variable between these two samples.

We can also use `topCounts` to examine the data identified as 'null'.

```
> topCounts(cDPair, group = NULL, number = 500)
```

We can visualise the data in a number of ways. We can first examine the posterior likelihoods against log-ratio values.

```
> plotPosteriors(cDPair, group = 2, samplesA = 1:2, samplesB = 3:4)
```

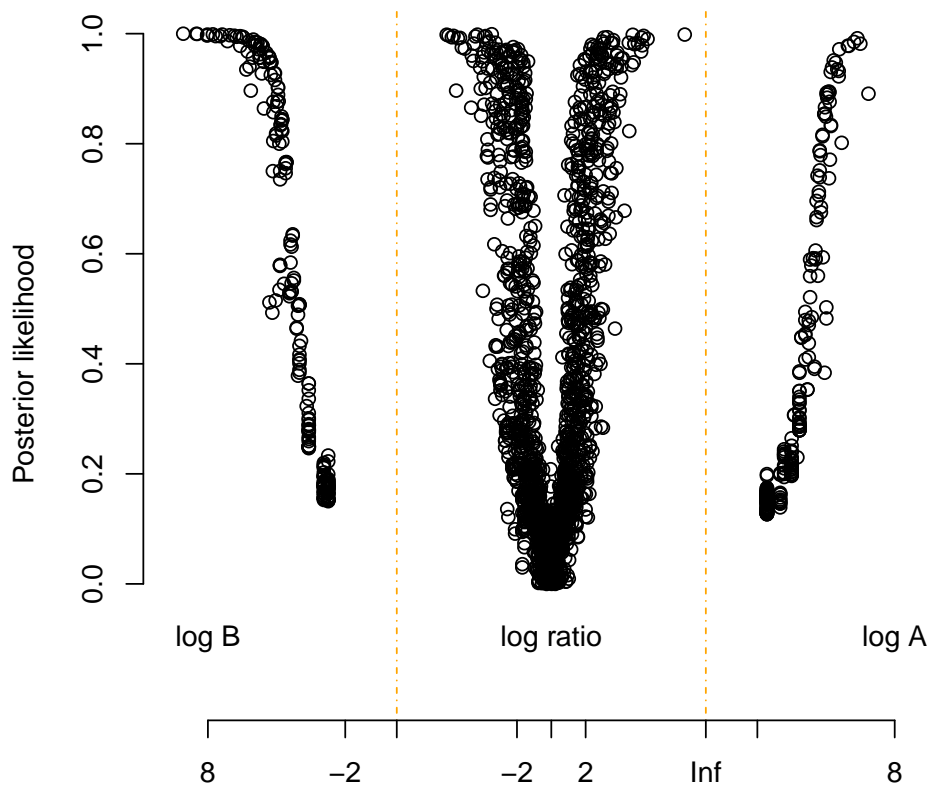


Figure 4: Posterior likelihoods of differential expression against log-ratios of the data. Where the data in one of the sample groups consists entirely of zeros, the log-ratio would be infinite. In this case, we plot instead the log-values of the non-zero group. Note the skew in the data; there are many more loci with a high-likelihood of differential expression over-expressed in the WT/D3 graft compared to the D3/D3 graft than vice versa.

Also informative is the MA-plot. We can color the data by the posterior likelihoods of differential expression.

```
> plotMA.CD(cDPair, samplesA = c(1,2), samplesB = c(3,4),
+           col = rgb(red = exp(cDPair@posteriors[,2]), green = 0, blue = 0))
```

## 6.5 Multiple Group Comparisons

We next examine all three experimental conditions simultaneously. We first need to define the replicate structure of the data.

```
> cD@replicates <- as.factor(c("D3/D3", "D3/D3", "WT/D3", "WT/D3", "WT/WT", "WT/WT"))
```

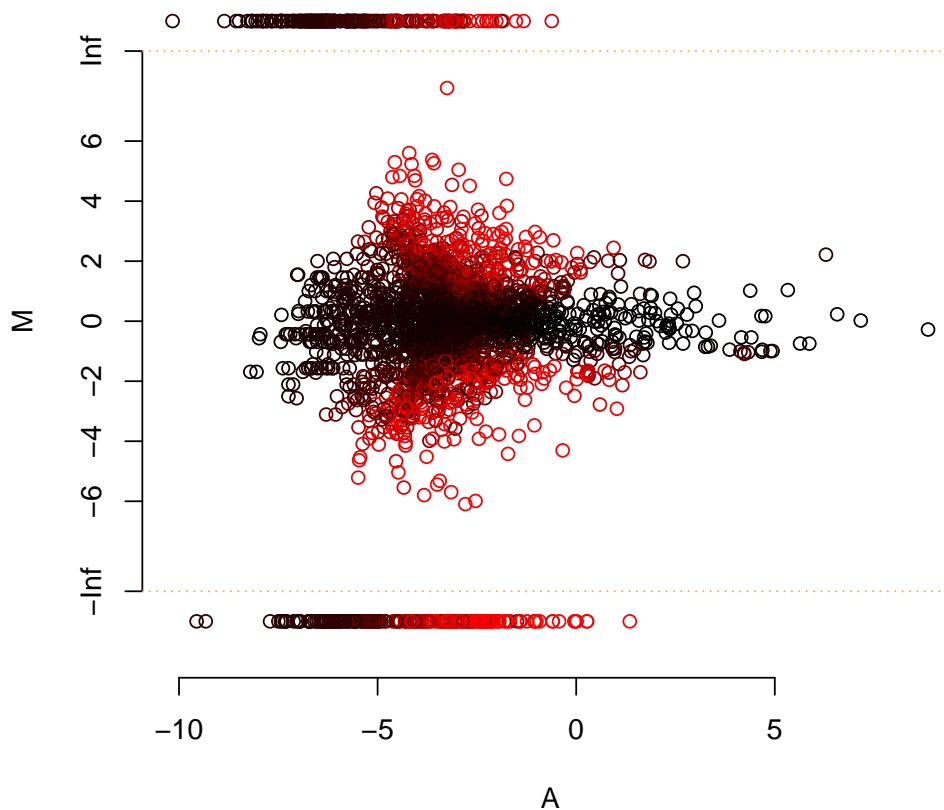


Figure 5: 'MA'-plot for count data. Where the data in one of the sample groups consists entirely of zeros, the log-ratio would be infinite. In this case, we plot instead the log-values of the non-zero group. Differentially expressed data are colored red, and non-differentially expressed data black.

As before, we begin by supposing that at least some of the loci will be unaffected by the different experimental conditions prevailing in our replicate groups, and so we create one model of no differential expression.

We do this by defining a vector NDE.

```
> NDE <- factor(c(1,1,1,1,1,1))
```

Each member of the NDE vector represents one sample in our experiment. By giving each item in the NDE vector the same number, we indicate that, under the hypothesis of no differential expression, all the samples belong to the same group.

We may also conjecture that some of the loci that are present in the wild-type root will not be present in the Dicer 2,3,4 mutant roots. We represent this conjecture with the vector

```
> d3dep <- c("wtRoot", "wtRoot", "wtRoot", "wtRoot", "dicerRoot", "dicerRoot")
```

This vector indicates that the fifth and sixth samples, which consist of the wild-type root samples, are in a separate expression group to the other samples, corresponding to the Dicer 2,3,4 mutant.

Finally, we hypothesise that some of the small RNAs generated in the wild-type shoot will move to the root. We represent this hypothesis with the vector

```
> mobile <- c("dicerShoot", "dicerShoot", "wtShoot", "wtShoot", "wtShoot", "wtShoot")
```

This vector shows that all samples with a wild-type shoot are distinct from those samples with a Dicer 2,3,4 shoot.

We can now add these models to the locus data by modifying the @groups slot

```
> groups(cD) <- list(NDE = NDE, d3dep = d3dep, mobile = mobile)
```

Note that in this case the replicate structure does not correspond to any biologically plausible model; we do not expect that any loci will be different between all three experimental groups.

We can now find the priors and likelihoods for this analysis as before.

```
> cD <- getPriors.NB(cD, cl = cl)
> cD <- getLikelihoods(cD, nullData = TRUE, cl = cl)
```

We can see if there are any potential candidates for mobile sRNA loci by using the 'topCounts' function.

```
> topCounts(cD, group = "mobile", normaliseData = TRUE)
  chr  start      end SL236 SL260 SL237 SL238 SL239 SL240 Likelihood
1   1  447231  447298     0     0   202   203   166   157  0.9999998
2   1  8287590 8287674     0     0   100    92    59    89  0.9999988
3   1 14188044 14188079     3     0   103    86    92    79  0.9999974
4   1  6127755  6127808     0     0   119    73   100    61  0.9999946
5   1  6880517  6880553     0     0    88    66    58    59  0.9999917
6   1  9254068  9254167     0     0    80    76    55    40  0.9999828
7   1 11140107 11140158     0     0    91    69    61    40  0.9999676
8   1 13042720 13042777     3     4    63    50    50    56  0.9999497
9   1  9373429  9373528     0     0    62    43    36    39  0.9999413
10  1  8766946  8767133    91   152     5     4     7     7  0.9999225
      ordering  FDR.mobile  FWER.mobile
1  wtShoot>dicerShoot 2.096725e-07 2.096725e-07
2  wtShoot>dicerShoot 7.068096e-07 1.413619e-06
3  wtShoot>dicerShoot 1.338094e-06 4.014277e-06
4  wtShoot>dicerShoot 2.357099e-06 9.428371e-06
5  wtShoot>dicerShoot 3.537660e-06 1.768820e-05
6  wtShoot>dicerShoot 5.812390e-06 3.487393e-05
7  wtShoot>dicerShoot 9.609260e-06 6.726329e-05
8  wtShoot>dicerShoot 1.469499e-05 1.175550e-04
9  wtShoot>dicerShoot 1.958490e-05 1.762523e-04
10 dicerShoot>wtShoot 2.537716e-05 2.537461e-04
```

We can also identify dicer-dependent root specific small RNA loci by examining our alternative model for differential expression.

```
> topCounts(cD, group = "d3dep", normaliseData = TRUE)
  chr  start      end SL236 SL260 SL237 SL238 SL239 SL240 Likelihood      ordering
1   1 12726934 12726976     5     5     6    10    42    41  0.9987233 dicerRoot>wtRoot
2   1  9013965  9014013     5     5     5     9    37    36  0.9987203 dicerRoot>wtRoot
3   1  8741412  8741466     5     4     1     0    37    46  0.9981447 dicerRoot>wtRoot
4   1 14154618 14154660    23    36    17    20   196   254  0.9979729 dicerRoot>wtRoot
5   1 13689324 13689396     6     5     5     7    30    25  0.9979293 dicerRoot>wtRoot
6   1 12824336 12824400     0     1     0     0     7     5  0.9926080 dicerRoot>wtRoot
7   1  8238064  8238106     7     5     8     5    30    23  0.9886580 dicerRoot>wtRoot
8   1  2105085  2105119     9     8     8     4     0     0  0.9865197 wtRoot>dicerRoot
9   1 14206419 14206455    32    34    46    32     9    13  0.9864306 wtRoot>dicerRoot
10  1  6263246  6263343     0     1     2     0     9     9  0.9862118 dicerRoot>wtRoot
      FDR.d3dep  FWER.d3dep
1  0.001276659  0.001276659
2  0.001278187  0.002554739
3  0.001470545  0.004405261
4  0.001609677  0.006423405
5  0.001701873  0.008480760
6  0.002650231  0.015810092
7  0.003891916  0.026972799
```

```
8 0.005090467 0.040089522
9 0.006032575 0.053114965
10 0.006808138 0.066170813
```

By including more experimental conditions in our analyses, increasingly complex patterns of expression can be detected from sequencing data.

Finally, we shut down the cluster (assuming it was started to begin with).

```
> if(!is.null(cl)) stopCluster(cl)
```

## Session Info

---

```
> sessionInfo()
```

```
R version 3.1.1 Patched (2014-09-25 r66681)
```

```
Platform: x86_64-unknown-linux-gnu (64-bit)
```

```
locale:
```

```
[1] LC_CTYPE=en_US.UTF-8      LC_NUMERIC=C               LC_TIME=en_US.UTF-8
[4] LC_COLLATE=C              LC_MONETARY=en_US.UTF-8   LC_MESSAGES=en_US.UTF-8
[7] LC_PAPER=en_US.UTF-8     LC_NAME=C                  LC_ADDRESS=C
[10] LC_TELEPHONE=C           LC_MEASUREMENT=en_US.UTF-8 LC_IDENTIFICATION=C
```

```
attached base packages:
```

```
[1] stats4      parallel    stats      graphics  grDevices  utils      datasets  methods
[9] base
```

```
other attached packages:
```

```
[1] baySeq_2.0.50      abind_1.4-0          GenomicRanges_1.18.1 GenomeInfoDb_1.2.0
[5] IRanges_2.0.0     S4Vectors_0.4.0     BiocGenerics_0.12.0
```

```
loaded via a namespace (and not attached):
```

```
[1] BiocStyle_1.4.1 XVector_0.6.0  tools_3.1.1
```

## References

---

- [1] Thomas J. Hardcastle and Krystyna A. Kelly. *baySeq: Empirical Bayesian Methods For Identifying Differential Expression In Sequence Count Data*. BMC Bioinformatics (2010).
- [2] Attila Molnar and Charles W. Bassett and Thomas J. Hardcastle and Ruth Dunn and David C. Bauclombe *Small silencing RNAs in plants are mobile and direct epigenetic modification in recipient cells*. Science (2010).
- [3] Mark Robinson edgeR: 'Methods for differential expression in digital gene expression datasets'. Bioconductor.