

segmentSeq: methods for identifying small RNA loci from high-throughput sequencing data

Thomas J. Hardcastle

April 24, 2017

1 Introduction

High-throughput sequencing technologies allow the production of large volumes of short sequences, which can be aligned to the genome to create a set of *matches* to the genome. By looking for regions of the genome which to which there are high densities of matches, we can infer a segmentation of the genome into regions of biological significance. The methods we propose allows the simultaneous segmentation of data from multiple samples, taking into account replicate data, in order to create a consensus segmentation. This has obvious applications in a number of classes of sequencing experiments, particularly in the discovery of small RNA loci and novel mRNA transcriptome discovery.

We approach the problem by considering a large set of potential *segments* upon the genome and counting the number of tags that match to that segment in multiple sequencing experiments (that may or may not contain replication). We then adapt the empirical Bayesian methods implemented in the `baySeq` package [1] to establish, for a given segment, the likelihood that the count data in that segment is similar to background levels, or that it is similar to the regions to the left or right of that segment. We then rank all the potential segments in order of increasing likelihood of similarity and reject those segments for which there is a high likelihood of similarity with the background or the regions to the left or right of the segment. This gives us a large list of overlapping segments. We reduce this list to identify non-overlapping loci by choosing, for a set of overlapping segments, the segment which has the lowest likelihood of similarity with either background or the regions to the left or right of that segment and rejecting all other segments that overlap with this segment. For fuller details of the method, see Hardcastle *et al.* [2].

2 Preparation

We begin by loading the `segmentSeq` package.

```
> library(segmentSeq)
```

Note that because the experiments that `segmentSeq` is designed to analyse are usually massive, we should use (if possible) parallel processing as implemented by the `parallel` package. If using this approach, we need to begin by define a *cluster*. The following command will use eight processors on a single machine; see the help page for 'makeCluster' for more information. If we don't want to parallelise, we can proceed anyway with a `NULL` cluster.

```
> if(require("parallel"))
+ {
+   numCores <- min(8, detectCores())
+   cl <- makeCluster(numCores)
+ } else {
+   cl <- NULL
+ }
```

The `readGeneric` function is able to read in tab-delimited files which have appropriate column names, and create an `alignmentData` object. Alternatively, if the appropriate column names are not present, we can specify which columns to use for the data. In either case, to use this function we pass a character vector of files, together with information on which data are to be treated as replicates to the function. We also need to define the lengths of

the chromosome and specify the chromosome names as a character. The data here, drawn from text files in the 'data' directory of the segmentSeq package are taken from the first million bases of an alignment to chromosome 1 and the first five hundred thousand bases of an alignment to chromosome 2 of *Arabidopsis thaliana* in a sequencing experiment where libraries 'SL9' and 'SL10' are replicates, as are 'SL26' and 'SL32'. Libraries 'SL9' and 'SL10' are sequenced from an Argonaute 6 IP, while 'SL26' and 'SL32' are an Argonaute 4 IP.

A similar function, readBAM performs the same operation on files in the BAM format. Please consult the help page for further details.

```
> datadir <- system.file("extdata", package = "segmentSeq")
> libfiles <- c("SL9.txt", "SL10.txt", "SL26.txt", "SL32.txt")
> libnames <- c("SL9", "SL10", "SL26", "SL32")
> replicates <- c("AG06", "AG06", "AG04", "AG04")
> aD <- readGeneric(files = libfiles, dir = datadir,
+                   replicates = replicates, libnames = libnames,
+                   polyLength = 10, header = TRUE, gap = 200)
> aD
```

An object of class "alignmentData"
3149 rows and 4 columns

Slot "libnames":
[1] "SL9" "SL10" "SL26" "SL32"

Slot "replicates":
[1] AG06 AG06 AG04 AG04
Levels: AG04 AG06

Slot "alignments":
GRanges object with 3149 ranges and 2 metadata columns:

	seqnames	ranges	strand	tag	multireads
	<Rle>	<IRanges>	<Rle>	<character>	<numeric>
[1]	>Chr1	[265, 284]	-	AAATGAAGATAAACCATCCA	1
[2]	>Chr1	[405, 427]	-	AAGGAGTAAGAATGACAATAAAT	1
[3]	>Chr1	[406, 420]	-	AAGAATGACAATAAAA	1
[4]	>Chr1	[600, 623]	+	AAGGATTGGTGGTTGAAGACACA	1
[5]	>Chr1	[665, 688]	+	ATCCTTGTAGCACACATTTGGCA	1
...
[3145]	>Chr1	[991569, 991589]	-	CCGATAAACGCATACTTCCT	1
[3146]	>Chr1	[992039, 992054]	-	AAGGAAATTAGAAAAT	1
[3147]	>Chr1	[995357, 995372]	+	AGAGACATGGGCGACA	1
[3148]	>Chr1	[995493, 995507]	+	AAACTCGTGAAGAAG	1
[3149]	>Chr1	[995817, 995840]	-	AGAGATCAAGTATATAGAATTAAG	1

seqinfo: 1 sequence from an unspecified genome; no seqlengths

Slot "data":
Matrix with 3149 rows.

	SL9	SL10	SL26	SL32
1	1	0	0	0
2	0	0	0	2
3	0	1	0	0
4	0	1	0	0
5	7	1	0	0
...
3145	1	0	0	0
3146	0	1	0	0
3147	0	1	0	0
3148	0	1	0	0
3149	1	0	0	0

```
Slot "libsizes":
[1] 1193 1598 1818 1417
```

Next, we process this alignmentData object to produce a segData object. This segData object contains a set of potential segments on the genome defined by the start and end points of regions of overlapping alignments in the alignmentData object. It then evaluates the number of tags that hit in each of these segments.

```
> sD <- processAD(aD, c1 = c1)
> sD
```

```
GRanges object with 1452 ranges and 0 metadata columns:
```

	seqnames	ranges	strand
	<Rle>	<IRanges>	<Rle>
[1]	>Chr1	[265, 284]	*
[2]	>Chr1	[265, 420]	*
[3]	>Chr1	[265, 623]	*
[4]	>Chr1	[265, 688]	*
[5]	>Chr1	[265, 830]	*
...
[1448]	>Chr1	[992039, 992054]	*
[1449]	>Chr1	[995357, 995372]	*
[1450]	>Chr1	[995357, 995507]	*
[1451]	>Chr1	[995493, 995507]	*
[1452]	>Chr1	[995817, 995840]	*

```
-----
```

```
seqinfo: 1 sequence from an unspecified genome; no seqlengths
```

```
An object of class "lociData"
```

```
1452 rows and 4 columns
```

```
Slot "replicates"
```

```
AG06 AG06 AG04 AG04
```

```
Slot "groups":
```

```
list()
```

```
Slot "data":
```

```
SL9 SL10 SL26 SL32
```

```
Slot "annotation":
```

```
data frame with 0 columns and 0 rows
```

```
Slot "locLikelihoods" (stored on log scale):
```

```
Matrix with 0 rows.
```

```
<0 x 0 matrix>
```

We can now construct a segment map from these potential segments.

Segmentation by heuristic methods

A fast method of segmentation can be achieved by exploiting the bimodality of the densities of small RNAs in the potential segments. In this approach, we assign each potential segment to one of two clusters for each replicate group, either as a segment or a null based on the density of sequence tags within that segment. We then combine these clusterings for each replicate group to gain a consensus segmentation map.

```
> hS <- heuristicSeg(sD = sD, aD = aD, RKPM = 1000, largeness = 1e8, getLikes = TRUE, c1 = c1)
```

```
.....
```

Segmentation by empirical Bayesian methods

A more refined approach to the problem uses an existing segment map (or, if not provided, a segment map defined by the `hS` function) to acquire empirical distributions on the density of sequence tags within a segment. We can then estimate posterior likelihoods for each potential segment as being either a true segment or a null. We then identify all potential segments in the with a posterior likelihood of being a segment greater than some value 'lociCutoff' and containing no subregion with a posterior likelihood of being a null greater than 'nullCutoff'. We then greedily select the longest segments satisfying these criteria that do not overlap with any other such segments in defining our segmentation map.

```
> cS <- classifySeg(sD = sD, aD = aD, cD = hS, cI = cI)
.....
> cS
GRanges object with 64 ranges and 0 metadata columns:
      seqnames      ranges strand
      <Rle>        <IRanges> <Rle>
 [1] >Chr1 [ 1, 264]      *
 [2] >Chr1 [ 265, 967]     *
 [3] >Chr1 [ 968, 17054]   *
 [4] >Chr1 [17055, 18728]  *
 [5] >Chr1 [18729, 27656]  *
 ...      ...           ...
 [60] >Chr1 [889498, 889525] *
 [61] >Chr1 [889526, 944194] *
 [62] >Chr1 [944195, 944222] *
 [63] >Chr1 [944223, 958610] *
 [64] >Chr1 [958611, 959152] *
-----
seqinfo: 1 sequence from an unspecified genome; no seqlengths
An object of class "lociData"
64 rows and 4 columns

Slot "replicates"
AG06 AG06 AG04 AG04
Slot "groups":
list()

Slot "data":
      AG06.1 AG06.2 AG04.1 AG04.2
 [1,] 0      0      0      0
 [2,] 55     47     65     85
 [3,] 2      3      0      0
 [4,] 682    621    1405   1103
 [5,] 0      3      0      0
59 more rows...

Slot "annotation":
data frame with 0 columns and 64 rows

Slot "locLikelihoods" (stored on log scale):
Matrix with 64 rows.
      AG04      AG06
 1  0.040948 0.038297
 2  0.90347  0.96323
 3  0.016844 0.031157
 4  0.98952  0.99691
 5  0.018425 0.047827
```

```

...      ...      ...
60  0.08755  0.95692
61  0.11928  0.061966
62  0.98248  0.89929
63  0.16222  0.027597
64  0.29943  0.96233

```

Expected number of loci in each replicate group

```

AG04      AG06
28.59675  31.16924

```

By one of these methods, we finally acquire an annotated lociData object, with the annotations describing the co-ordinates of each segment.

We can use this lociData object, in combination with the alignmentData object, to plot the segmented genome.

```

> par(mfrow = c(2,1), mar = c(2,6,2,2))
> plotGenome(aD, hS, chr = ">Chr1", limits = c(1, 1e5),
+           showNumber = FALSE, cap = 50)
> plotGenome(aD, cS, chr = ">Chr1", limits = c(1, 1e5),
+           showNumber = FALSE, cap = 50)

```

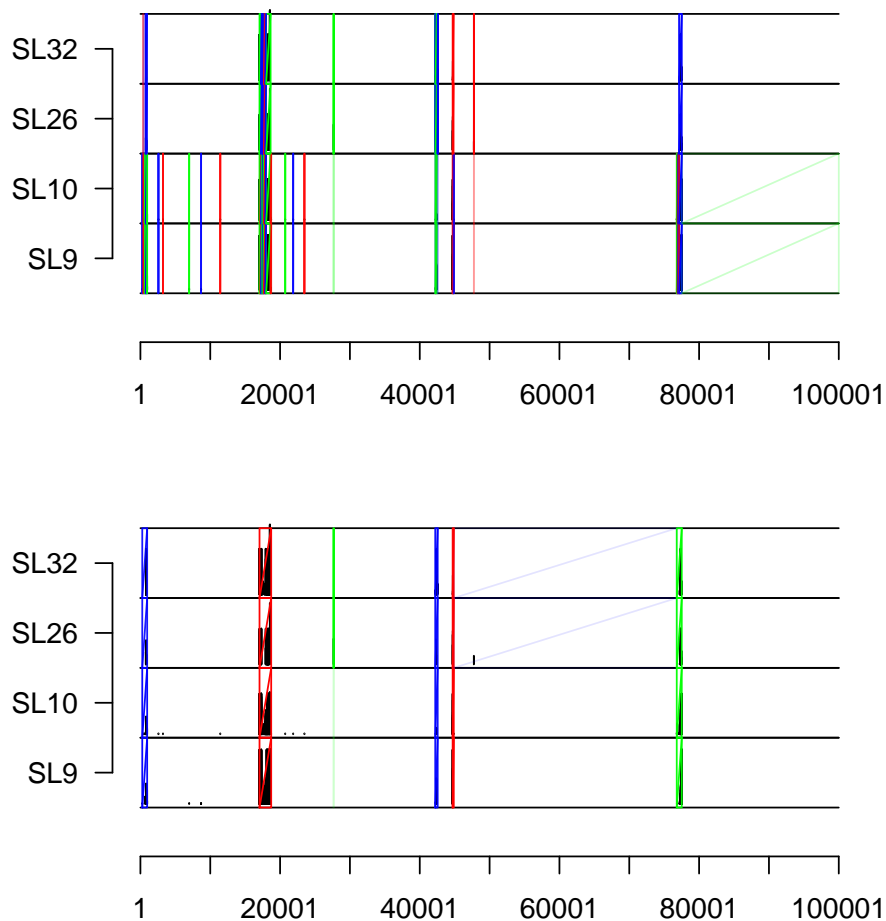


Figure 1: The segmented genome (first 10^5 bases of chromosome 1).

Given the calculated likelihoods, we can filter the segmented genome by controlling on likelihood, false discovery rate, or familywise error rate

```

> loci <- selectLoci(cS, FDR = 0.05)
> loci
GRanges object with 30 ranges and 0 metadata columns:
      seqnames          ranges strand
      <Rle>             <IRanges> <Rle>
 [1] >Chr1 [ 265,  967]      *
 [2] >Chr1 [17055, 18728]     *
 [3] >Chr1 [42217, 42570]     *
 [4] >Chr1 [44710, 44870]     *
 [5] >Chr1 [76799, 77519]     *
 ...
 [26] >Chr1 [758302, 760446]    *
 [27] >Chr1 [789508, 789548]    *
 [28] >Chr1 [889498, 889525]    *
 [29] >Chr1 [944195, 944222]    *
 [30] >Chr1 [958611, 959152]    *
-----
seqinfo: 1 sequence from an unspecified genome; no seqlengths
An object of class "lociData"
30 rows and 4 columns

Slot "replicates"
AG06 AG06 AG04 AG04
Slot "groups":
list()

Slot "data":
      AG06.1 AG06.2 AG04.1 AG04.2
 [1,]    55    47    65    85
 [2,]   682   621  1405  1103
 [3,]    31    11    48    68
 [4,]    73    58    47    21
 [5,]   182   168   275   131
25 more rows...

Slot "annotation":
data frame with 0 columns and 30 rows

Slot "locLikelihoods" (stored on log scale):
Matrix with 30 rows.
      AG04    AG06
 1  0.90347 0.96323
 2  0.98952 0.99691
 3  0.94658 0.93199
 4   0.9552 0.99631
 5  0.95362 0.95036
...
26 0.96388 0.94654
27 0.81632 0.98334
28 0.08755 0.95692
29 0.98248 0.89929
30 0.29943 0.96233

Expected number of loci in each replicate group
      AG04    AG06
24.96588 29.01355

```

The lociData objects can now be examined for differential expression with the baySeq package.

First we define the possible models of differential expression on the data. In this case, the models are of non-differential expression and pairwise differential expression.

```
> groups(cS) <- list(NDE = c(1,1,1,1), DE = c("AG06", "AG06", "AG04", "AG04"))
```

Then we get empirical distributions on the parameter space of the data.

```
> cS <- getPriors(cS, cl = cl)
```

Then we get the posterior likelihoods of the data conforming to each model. Since the 'cS' object contains null regions as well as true loci, we will use the 'nullData = TRUE' option to distinguish between non-differentially expressed loci and non-expressed regions. By default, the loci likelihoods calculated earlier will be used to weight the initial parameter fit in order to detect null data.

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

.

We can examine the highest likelihood non-expressed ('null') regions

```
> topCounts(cS, NULL, number = 3)
```

	seqnames	start	end	width	strand	AG06.1	AG06.2	AG04.1	AG04.2	Likelihood	FDR.
1	>Chr1	372735	423138	50404	*	0	0	0	0	0.9723825	0.02761747
2	>Chr1	552693	587497	34805	*	0	0	0	0	0.9721573	0.02773009
3	>Chr1	641225	670379	29155	*	0	0	0	0	0.9720306	0.02780984

FWER.

1	0.02761747
2	0.05469123
3	0.08113090

The highest likelihood expressed but non-differentially expressed regions

```
> topCounts(cS, "NDE", number = 3)
```

	seqnames	start	end	width	strand	AG06.1	AG06.2	AG04.1	AG04.2	Likelihood	FDR.NDE
1	>Chr1	76799	77519	721	*	182	168	275	131	0.9187587	0.08124134
2	>Chr1	423139	423547	409	*	19	20	28	12	0.8645442	0.10834857
3	>Chr1	265	967	703	*	55	47	65	85	0.8324692	0.12807599

FWER.NDE

1	0.08124134
2	0.20569254
3	0.33876351

And the highest likelihood differentially expressed regions

```
> topCounts(cS, "DE", number = 3)
```

	seqnames	start	end	width	strand	AG06.1	AG06.2	AG04.1	AG04.2	Likelihood	ordering
1	>Chr1	634297	634350	54	*	65	90	12	17	0.9983699	AG06>AG04
2	>Chr1	238359	238417	59	*	9	9	0	0	0.9891103	AG06>AG04
3	>Chr1	587498	588015	518	*	132	333	897	764	0.9883517	AG04>AG06

	FDR.DE	FWER.DE
1	0.001630082	0.001630082
2	0.006259895	0.012502038
3	0.008056040	0.024004742

Finally, to be a good citizen, we stop the cluster we started earlier:

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

Session Info

```
> sessionInfo()
```

R version 3.4.0 (2017-04-21)
 Platform: x86_64-pc-linux-gnu (64-bit)
 Running under: Ubuntu 16.04.2 LTS

Matrix products: default
 BLAS: /home/biocbuild/bbs-3.5-bioc/R/lib/libRblas.so
 LAPACK: /home/biocbuild/bbs-3.5-bioc/R/lib/libRlapack.so

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

other attached packages:
 [1] segmentSeq_2.10.0 ShortRead_1.34.0 GenomicAlignments_1.12.0
 [4] SummarizedExperiment_1.6.0 DelayedArray_0.2.0 matrixStats_0.52.2
 [7] Biobase_2.36.0 Rsamtools_1.28.0 Biostrings_2.44.0
 [10] XVector_0.16.0 BiocParallel_1.10.0 baySeq_2.10.0
 [13] abind_1.4-5 GenomicRanges_1.28.0 GenomeInfoDb_1.12.0
 [16] IRanges_2.10.0 S4Vectors_0.14.0 BiocGenerics_0.22.0

loaded via a namespace (and not attached):
 [1] Rcpp_0.12.10 RColorBrewer_1.1-2 compiler_3.4.0
 [4] bitops_1.0-6 tools_3.4.0 zlibbioc_1.22.0
 [7] digest_0.6.12 evaluate_0.10 lattice_0.20-35
 [10] Matrix_1.2-9 yaml_2.1.14 GenomeInfoDbData_0.99.0
 [13] stringr_1.2.0 hwriter_1.3.2 knitr_1.15.1
 [16] locfit_1.5-9.1 rprojroot_1.2 grid_3.4.0
 [19] rmarkdown_1.4 limma_3.32.0 latticeExtra_0.6-28
 [22] edgeR_3.18.0 magrittr_1.5 backports_1.0.5
 [25] htmltools_0.3.5 BiocStyle_2.4.0 stringi_1.1.5
 [28] RCurl_1.95-4.8

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] Thomas J. Hardcastle and Krystyna A. Kelly and David C. Baulcombe. *Identifying small RNA loci from high-throughput sequencing data*. Bioinformatics (2012).