The bumphunter user's guide

Kasper Daniel Hansen khansen@jhsph.edu
Martin Aryee aryee.martin@mgh.harvard.edu
Rafael A. Irizarry rafa@jhu.edu

Modified: November 23, 2012. Compiled: October 7, 2025

Introduction

This package implements the statistical procedure described in [2] (with some small modifications). Notably, batch effect removal and the application of the bootstrap to linear models of Efron and Tibshirani [1] need additional code.

For any given type of data, it is usually necessary to make a number of choices and/or transformations before the bump hunting methodology is ready to be applied. Typically, these modifications resides in other packages. Examples are charm (for CHARM-like methylation microarrays), bsseq (for whole-genome bisulfite sequencing data) and minfi (for Illumina 450k methylation arrays). In some cases (specifically bsseq) only parts of the methodology as implemented in the bumphunter package is applied, although the conceptual approach is still build on bump hunting.

In other words, this package is mostly intended for developers wishing to adapt the general methodology to their specific applications.

The core of the package is encapsulated in the bumphunter method which uses the underlying bumphunterEngine to do the heavy lifting. However, bumphunterEngine consists of a number of useful functions that does much of the specific tasks involved in bump hunting. This document attempts to describe the overall workflow as well as the specific functions. The relevant functions are clusterMaker, getSegments, findRegions.

> library(bumphunter)

Note that this package is written with genomic data as an illustrative example but most of it is easily generalizable to other data types.

Other functions

Most of the bumphunter package is code for bump hunting. But we also include a number of convenience functions we have found useful, which are not necessarily part of the bump hunting exercise. At the time of writing, this include annotateNearest.

The Statistical Model

The bump hunter methodology is meant to work on data with several biological replicates, similar to the lmFit function in limma. While the package is written using genomic data as an illustrative example, most of it is generalizable to other data types (with some one-dimensional location information).

We assume we have data Y_{ij} where i represents (biological) replicate and l_j represents genomic location. The use of j and l_j is a convenience notation, allowing us to discuss the "next" observation j+1 which may be some distance $lj+1-l_j$ away. Note that we assume in this notation that all replicates have been observed at the same set of genomic locations.

The basic statistical model is the following:

$$Y_{ij} = \beta_0(l_j) + \beta_1(l_j)X_j + \varepsilon_{ij}$$

with i representing subject, l_j representing the jth location, X_j is the covariate of interest (for example $X_j = 1$ for cases and $X_j = 0$ otherwise), ε_{ij} is measurement error, $\beta_0(l)$ is a baseline function, and $\beta_1(l)$ is the parameter of interest, which is a function of location. We assume that $\beta_1(l)$ will be equal to zero over most of the genome, and we want to identify stretched where $\beta_1(l) \neq 0$, which we call *bumps*.

We want to share information between nearby locations, typically through some form of smoothing.

Creating clusters

For many genomic applications the locations are clustered. Each cluster is a distinct unit where the model fitting will be done separately, and each cluster does not depend on the data, only on the locations l_j . Typically there is some maximal distance, and we do not want to smooth between observations separated by more than this distance. The choice of distance is very application dependent.

"Clusters" are simply groups of locations such that two consecutive locations in the cluster are separated by less than some distance maxGap. For genomic applications, the biggest possible clusters are chromosomes.

The function clusterMaker defines such grouping locations.

Example: We first generate an example of typical genomic locations

Now we run the function to obtain the three clusters from the positions. We use the default gap of 300 base pairs (bps), i.e. any two points more than 300 bps away are put in a new cluster. Also, locations from different chromosomes are separated.

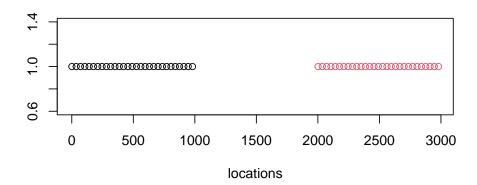
```
> cl <- clusterMaker(chr, pos, maxGap = 300)
> table(cl)

cl
    1    2    3
29    29    20
```

The output is an indexing variable telling us which cluster each location belongs to. Locations on different chromosomes are always on different clusters.

Note that data from the first chromosome has been split into two clusters:

```
> ind <- which(chr=="chr1")
> plot(pos[ind], rep(1,length(ind)), col=cl[ind],
+ xlab="locations", ylab="")
```



Breaking into segments

The function getSegments is used to find segments that are positive, near zero, and negative. Specifically we have a vector of numbers θ_j with each number associated with a genomic location l_j (thinks either test statistics or estimates of $\beta_i(l)$). A segment is a list of consecutive locations such that all θ_l in the segment are either "positive", "near zero" or "negative". In order to define "positive" etc we need a cutoff which is one number L (in which case "near zero" is [-L, L]) or two numbers L, U (in which case "near zero" is [L; U]).

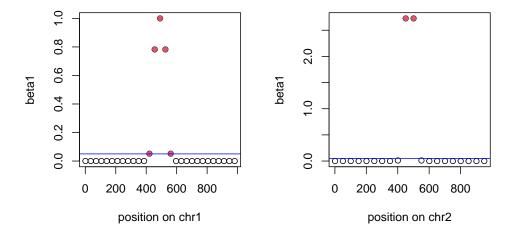
Example: we are going to create a simulated $\beta_1(l)$ with a couple of real bumps.

```
> Indexes <- split(seq_along(cl), cl)
> beta1 <- rep(0, length(pos))
> for(i in seq(along=Indexes)){
+        ind <- Indexes[[i]]
+        x <- pos[ind]
+        z <- scale(x, median(x), max(x)/12)
+        beta1[ind] <- i*(-1)^(i+1)*pmax(1-abs(z)^3,0)^3 ##multiply by i to value
+ }</pre>
```

We now find bumps of this functions by

```
> segs <- getSegments(beta1, cl, cutoff=0.05)</pre>
```

Now we can make, for example, a plot of all the positive bumps



This function is used by regionFinder which is described next.

regionFinder

This function packages up the results of getSegments into a table of regions with the location and characteristics of bumps.

```
> tab <- regionFinder(beta1, chr, pos, cl, cutoff=0.05)
> tab
   chr start
             end
                      value
                                 area cluster indexStart indexEnd
       2281 2701 -1.2636037 16.426848
3 chr1
                                            2
                                                      38
                                                               50 13
2 chr2
        451 501 2.7262463 5.452493
                                            3
                                                      68
                                                               69 2
1 chr1
        421 561 0.5336474 2.668237
                                            1
                                                      13
                                                               17 5
 clusterL
3
       29
2
       20
1
       29
```

In the plot in the preceding section we show two of these regions in red.

Note that regionFinder and getSegments do not really contain any statistical model. All it does is finding regions based on segmenting a vector θ_i associated with genomic locations l_i .

Bumphunting

Bumphunter is a more complicated function. In addition to regionFinder and clusterMaker it also implements a statistical model as well as permutation testing to assess uncertainty.

We start by creating a simulated data set of 10 cases and 10 controls (recall that beta1 was defined above).

```
> beta0 <- 3*sin(2*pi*pos/720)
> X <- cbind(rep(1,20), rep(c(0,1), each=10))
> error <- matrix(rnorm(20*length(beta1), 0, 1), ncol=20)
> y <- t(X[,1])%x%beta0 + t(X[,2])%x%beta1 + error</pre>
```

Now we can run bumphunter

```
> tab <- bumphunter(y, X, chr, pos, cl, cutoff=.5)
> tab
```

```
$table
                    value area cluster indexStart indexEnd
    chr start end
16 chr1
         2316 2631 -1.6615645 16.6156453
                                                  2
                                                            39
                                                                      48
                                                  3
11 chr2
         451
                501
                     3.0778683 6.1557366
                                                            68
                                                                      69
12 chr2
          601
                801
                     0.7788492
                                 3.8942462
                                                  3
                                                            71
                                                                      75
                     0.9675794
                                 2.9027382
                                                  1
3 chr1
          456
                526
                                                            14
                                                                      16
14 chr1
          141
               176 -0.7372597
                                1.4745194
                                                  1
                                                             5
                                                                      6
4 chr1
          631
                666
                     0.7062928
                                1.4125856
                                                  1
                                                            19
                                                                      20
6 chr1
         841
                841
                     1.1771912
                                 1.1771912
                                                  1
                                                            25
                                                                      25
8 chr1
         2876 2876
                     1.0069851
                                 1.0069851
                                                  2
                                                            55
                                                                      55
15 chr1
          876
               876 -1.0011715
                                1.0011715
                                                  1
                                                            26
                                                                      26
            1
                     0.9574119
                                 0.9574119
                                                  3
                                                                      59
10 chr2
                 1
                                                            59
                     0.9358925
                                 0.9358925
                                                  1
1 chr1
            1
                  1
                                                             1
                                                                       1
2 chr1
          281
                     0.9189571
                                 0.9189571
                                                  1
                                                             9
                                                                      9
                281
                     0.8133397
                                 0.8133397
                                                  1
                                                            22
                                                                      22
5 chr1
          736
               736
9 chr1
        2946 2946
                     0.7985793
                                 0.7985793
                                                  2
                                                            57
                                                                      57
13 chr2
          951
                951
                     0.5773292
                                 0.5773292
                                                  3
                                                            78
                                                                      78
7 chr1
          946
               946
                     0.5756542
                                0.5756542
                                                  1
                                                            28
                                                                      28
    L clusterL
16 10
            29
    2
            20
11
12
            20
    5
3
    3
            29
14
    2
            29
4
    2
            29
6
    1
            29
8
    1
            29
15
    1
            29
10
    1
            20
1
    1
            29
2
    1
            29
5
    1
            29
9
    1
            29
13
    1
            20
7
    1
            29
$coef
               [,1]
 [1,]
     0.935892451
 [2,] -0.195092438
 [3,] -0.040884079
```

[4,] 0.035317753

- [5,] -0.512870863
- [6,] -0.961648537
- [7,] 0.018013482
- [8,] -0.146145859
- [9,] 0.918957051
- [10,] -0.330707316
- [11,] -0.047394007
- [12,] -0.267630605
- [13,] -0.070622163
- [15,] 0.070022105
- [14,] 1.021032087
- [15,] 0.936978394
- [16,] 0.944727696
- [17,] -0.394285782
- [18,] -0.006352622
- [19,] 0.906467185
- [20,] 0.506118419
- [21,] 0.194702263
- [22,] 0.813339746
- [23,] -0.250281073
- [24,] -0.143506222
- [25,] 1.177191216
- [26,] -1.001171473
- [27,] -0.003404589
- [28,] 0.575654195
- [29,] 0.123169050
- [30,] 0.414326874
- [31,] -0.281526197
- [32,] 0.001598158
- [33,] -0.420383148
- [34,] -0.476372363
- [35,] -0.390940701
- [36,] -0.208374551
- [37,] -0.029013902
- [38,] -0.054997386
- [39,] -0.518410293
- [40,] -1.286283494
- [41,] -2.185304337
- [42,] -2.386144655
- [43,] -1.474930742
- [44,] -2.618726560
- [45,] -1.537705177
- [46,] -1.475340202

- [47,] -2.366774470
- [48,] -0.766025321
- [49,] 0.235014876
- [50,] 0.416629002
- [51,] 0.387486402
- [52,] 0.238246172
- [53,] 0.426453884
- [54,] -0.388524280
- [55,] 1.006985091
- [56,] -0.199011910
- [57,] 0.798579279
- [58,] -0.023429311
- [59,] 0.957411915
- [60,] 0.439719522
- [61,] -0.303432897
- [01,] 0.303432037
- [62,] 0.184989278
- [63,] -0.315693063
- [64,] -0.374865595
- [65,] -0.024454534
- [66,] 0.263271634
- [67,] -0.364258216
- [68,] 3.276401420
- [69,] 2.879335208
- [70,] -0.387245544
- [71,] 0.731542360
- [72,] 0.564200907
- [73,] 0.840736895
- [74,] 1.010437839
- [75,] 0.747328175
- [76,] 0.235442771
- [77,] -0.265476686
- [78,] 0.577329197

\$fitted

- [,1]
- [1,] 0.935892451
- [2,] -0.195092438
- [3,] -0.040884079
- [4,] 0.035317753
- [5,] -0.512870863
- [6,] -0.961648537
- [7,] 0.018013482

- [8,] -0.146145859
- [9,] 0.918957051
- [10,] -0.330707316
- [11,] -0.047394007
- [12,] -0.267630605
- [13,] -0.070622163
- [14,] 1.021032087
- [15,] 0.936978394
- [16,] 0.944727696
- [17,] -0.394285782
- [18,] -0.006352622
- [19,] 0.906467185
- [20,] 0.506118419
- [21,] 0.194702263
- [22,] 0.813339746
- [23,] -0.250281073
- [24,] -0.143506222
- [25,] 1.177191216
- [26,] -1.001171473
- [27,] -0.003404589
- [28,] 0.575654195
- [29,] 0.123169050
- [30,] 0.414326874
- [31,] -0.281526197
- [32,] 0.001598158
- [33,] -0.420383148
- [34,] -0.476372363
- [35,] -0.390940701
- [36,] -0.208374551
- [37,] -0.029013902 [38,] -0.054997386
- [39,] -0.518410293
- [03,] 0.010110230
- [40,] -1.286283494 [41,] -2.185304337
- [42,] -2.386144655
- [43,] -1.474930742
- [44,] -2.618726560
- [45,] -1.537705177
- [46,] -1.475340202
- [47,] -2.366774470
- [40] 0 76600[201
- [48,] -0.766025321
- [49,] 0.235014876

```
[50,] 0.416629002
[51,] 0.387486402
[52,] 0.238246172
[53,] 0.426453884
[54,] -0.388524280
[55,] 1.006985091
[56,] -0.199011910
[57,] 0.798579279
[58,] -0.023429311
[59,] 0.957411915
[60,] 0.439719522
[61,] -0.303432897
[62,] 0.184989278
[63,] -0.315693063
[64,] -0.374865595
[65,] -0.024454534
[66,] 0.263271634
[67,] -0.364258216
[68,] 3.276401420
[69,] 2.879335208
[70,] -0.387245544
[71,] 0.731542360
[72,] 0.564200907
[73,] 0.840736895
[74,] 1.010437839
[75,] 0.747328175
```

\$pvaluesMarginal [1] NA

[76,] 0.235442771 [77,] -0.265476686 [78,] 0.577329197

> names(tab)

- [1] "table" "coef" "fitted"
- [4] "pvaluesMarginal"

> tab\$table

```
value
                                           area cluster indexStart indexEnd
    chr start
                 end
16 chr1
          2316 2631 -1.6615645 16.6156453
                                                        2
                                                                    39
                                                                               48
11 chr2
            451
                 501
                       3.0778683
                                     6.1557366
                                                        3
                                                                    68
                                                                               69
                                                        3
12 chr2
            601
                 801
                       0.7788492
                                     3.8942462
                                                                    71
                                                                              75
   chr1
           456
                 526
                                                        1
                                                                    14
                                                                              16
3
                       0.9675794
                                    2.9027382
                                                        1
                                                                     5
14 chr1
           141
                 176 - 0.7372597
                                                                                6
                                     1.4745194
                       0.7062928
            631
                 666
                                                        1
                                                                    19
                                                                              20
4
   chr1
                                     1.4125856
            841
                 841
                                                        1
                                                                    25
                                                                              25
6
   chr1
                       1.1771912
                                    1.1771912
                                                        2
8
   chr1
          2876 2876
                       1.0069851
                                     1.0069851
                                                                    55
                                                                              55
15 chr1
           876
                 876 -1.0011715
                                                        1
                                                                    26
                                                                              26
                                     1.0011715
10 chr2
              1
                    1
                       0.9574119
                                     0.9574119
                                                        3
                                                                    59
                                                                              59
              1
                    1
                                                        1
                                                                     1
1
   chr1
                       0.9358925
                                     0.9358925
                                                                                1
                                                                     9
                                                                                9
2
                 281
                                                        1
   chr1
           281
                       0.9189571
                                     0.9189571
5
                                                        1
                                                                    22
                                                                              22
   chr1
           736
                 736
                       0.8133397
                                     0.8133397
                                                        2
9
          2946 2946
                                                                    57
   chr1
                       0.7985793
                                     0.7985793
                                                                              57
13 chr2
            951
                 951
                       0.5773292
                                                        3
                                                                    78
                                                                              78
                                     0.5773292
7
   chr1
            946
                 946
                       0.5756542
                                     0.5756542
                                                        1
                                                                    28
                                                                              28
    L clusterL
16 10
              29
    2
11
              20
    5
12
              20
3
    3
              29
14
    2
              29
4
    2
              29
6
    1
              29
8
    1
              29
15
    1
              29
10
    1
              20
    1
              29
1
2
    1
              29
5
    1
              29
9
    1
              29
13
    1
              20
7
    1
              29
```

Briefly, the bumphunter function fits a linear model for each location (like lmFit from the limma package), focusing on one specific column (coefficient) of the design matrix. This coefficient of interest is optionally smoothed. Subsequently, a permutation can be used to test is formed for this specific coefficient.

The simplest way to use permutations to create a null distribution is to set B. If the number of samples is large this can be set to a large number, such as 1000. Note that this will be slow and

we have therefore provided parallelization capabilities. In cases were the user wants to define the permutations, for example cases in which all possible permutations can be enumerated, these can be supplied via the permutation argument.

Note that the function permits the matrix X to have more than two columns. This can be useful for those wanting to fit models that try to adjust for confounders, such as age and sex. However, when X has columns other than those representing an intercept term and the covariate of interest, the permutation test approach is not recommended. The function will run but give a warning. A method based on the bootstrap for linear models of Efron and Tibshirani [1] may be more appropriate but this is not currently implemented.

Faster bumphunting with multiple cores

bumphunter can be speeded up by using multiple cores. We use the foreach package which allows different parallel "back-ends" that will distribute the computation across multiple cores in a single machine, or across multiple machines in a cluster. The most straight-forward usage, illustrated below, involves multiple cores on a single machine. See the foreach documentation for more complex use cases, as well as the packages doParallel and doSNOW (among others). Finally, we use doRNG to ensure reproducibility of setting the seed within the parallel computations.

In order to use the foreach package we need to register a backend, in this case a multicore machine with 2 cores.

```
> library(doParallel)
> registerDoParallel(cores = 2)
```

bumphunter will now automatically use this backend

```
> tab <- bumphunter(y, X, chr, pos, cl, cutoff=.5, B=250, verbose = TRUE)
> tab
a 'bumps' object with 16 bumps
```

References

[1] B. Efron and R.J. Tibshirani.

[2] Andrew E Jaffe, Peter Murakami, Hwajin Lee, Jeffrey T Leek, M Daniele Fallin, Andrew P Feinberg, and Rafael A Irizarry. Bump hunting to identify differentially methylated regions in epigenetic epidemiology studies. *International Journal of Epidemiology*, 41(1):200–209, 2012.

Cleanup

This is a cleanup step for the vignette on Windows; typically not needed for users.

```
> bumphunter:::foreachCleanup()
```

SessionInfo

- R version 4.5.1 Patched (2025-08-23 r88802), x86_64-pc-linux-gnu
- Locale: LC_CTYPE=en_US.UTF-8, LC_NUMERIC=C, LC_TIME=en_GB, LC_COLLATE=C, LC_MONETARY=en_US.UTF-8, LC_MESSAGES=en_US.UTF-8, LC_PAPER=en_US.UTF-8, LC_NAME=C, LC_ADDRESS=C, LC_TELEPHONE=C, LC_MEASUREMENT=en_US.UTF-8, LC_IDENTIFICATION=C
- Time zone: America/New_York
- TZcode source: system (glibc)
- Running under: Ubuntu 24.04.3 LTS
- Matrix products: default
- BLAS: /home/biocbuild/bbs-3.22-bioc/R/lib/libRblas.so
- LAPACK: /usr/lib/x86_64-linux-gnu/lapack/liblapack.so.3.12.0
- Base packages: base, datasets, grDevices, graphics, methods, parallel, stats, stats4, utils
- Other packages: BiocGenerics 0.55.1, GenomicRanges 1.61.5, IRanges 2.43.5, S4Vectors 0.47.4, Seqinfo 0.99.2, bumphunter 1.51.1, doParallel 1.0.17, doRNG 1.8.6.2, foreach 1.5.2, generics 0.1.4, iterators 1.0.14, locfit 1.5-9.12, rngtools 1.5.2

• Loaded via a namespace (and not attached): AnnotationDbi 1.71.1, Biobase 2.69.1, BiocIO 1.19.0, BiocParallel 1.43.4, Biostrings 2.77.2, DBI 1.2.3, DelayedArray 0.35.3, GenomicAlignments 1.45.4, GenomicFeatures 1.61.6, KEGGREST 1.49.1, Matrix 1.7-4, MatrixGenerics 1.21.0, R6 2.6.1, RCurl 1.98-1.17, RSQLite 2.4.3, Rsamtools 2.25.3, S4Arrays 1.9.1, SparseArray 1.9.1, SummarizedExperiment 1.39.2, XML 3.99-0.19, XVector 0.49.1, abind 1.4-8, bit 4.6.0, bit64 4.6.0-1, bitops 1.0-9, blob 1.2.4, cachem 1.1.0, cli 3.6.5, codetools 0.2-20, compiler 4.5.1, crayon 1.5.3, curl 7.0.0, digest 0.6.37, fastmap 1.2.0, grid 4.5.1, httr 1.4.7, lattice 0.22-7, matrixStats 1.5.0, memoise 2.0.1, png 0.1-8, restfulr 0.0.16, rjson 0.2.23, rlang 1.1.6, rtracklayer 1.69.1, tools 4.5.1, vctrs 0.6.5, yaml 2.3.10