

Package ‘bumphunter’

October 9, 2013

Version 1.0.2

Title Bump Hunter

Description Tools for finding bumps in genomic data

Author Rafael A. Irizarry, Martin Aryee, Hector Corrada Bravo, Kasper D. Hansen, Harris A. Jaffee

Maintainer Rafael A. Irizarry <rafa@jhu.edu>

Depends R (>= 2.10), IRanges, GenomicRanges, foreach, iterators, methods, parallel

Suggests

limma, RUnit, BiocGenerics, doParallel, org.Hs.eg.db, TxDb.Hsapiens.UCSC.hg19.knownGene

Imports matrixStats, limma, itertools, doRNG

Collate AllGenerics.R annotateNearest.R dummyData.R

limmaForPermutations.R regionFinder.R smooth.R bumphunter.R

nearestgene.R fuzzy.match2.R matchGenes.R pointMatch.R

known_transcripts.R dataframe_to.R utils.R

License Artistic-2.0

LazyData yes

biocViews DNAMethylation, Epigenetics, Infrastructure, MultipleComparisons

R topics documented:

annotateNearest	2
bumphunter	4
clusterMaker	6
dummyData	7
getSegments	9
loessByCluster	10
pointMatch	12
regionFinder	13
runmedByCluster	14
smoother	15
TT	17

annotateNearest	<i>Annotate Nearest</i>
-----------------	-------------------------

Description

Finds the nearest region in a table of target regions and annotates the relationship.

Usage

```
annotateNearest(x, subject, annotate=TRUE, ...)
```

```
matchGenes(x, build, ...)
```

Arguments

x	An IRanges or GenomicRanges object, or a data.frame with columns for start, end, and, optionally, chr or seqnames.
subject	For annotateNearest to operate in matchGenes mode, a single word indicating the desired species/build; this is currently limited to hg19. Otherwise, an object of the same class as x.
annotate	Whether to annotate the result.
build	For example, hg19.
...	For matchGenes or matchGenes mode, promoterDist (default 2500), verbose (TRUE), all (FALSE), and mc.cores (1). mc.cores is crucial for parallelizing the annotation process, assuming the hardware facilities. But note that the default is NOTgetOption("mc.cores", 2L), yet, which leaves the explicit decision to the user. up and down limits on the nearest gene location are planned but not implemented yet. Otherwise, options for nearest, except that select="all" is not implemented yet and is currently rejected.

Details

The name matchGenes is provided for backward compatibility. The preferred call is to nearest with subject set to the build name. annotateNearest and matchGenes are basically wrappers for nearest(select="arbitrary") with specialized annotation of the result.

Value

Normally, a data frame with columns c("dist", "matchIndex", "type", "amountOverlap", "insideDist", "size1", ...), unless annotate is FALSE, in which case the first two columns, only, are returned as an integer matrix.

For matchGenes or matchGenes mode, with all unset or FALSE, a data frame with one row for each query and with columns c("name", "annotation", "description", "region", "distance", "subregion", "insideDist", ...). The first column is the `_gene_` nearest the query, by virtue of owning the transcript nearest the query. Note that this nearest transcript (name), per se, does not exist in the result. If all is TRUE, then each

row of the previous result is expanded into perhaps multiple rows, reflecting multiple (equally) nearest transcripts for a given query. In this case, two new columns at the beginning hold the query index and nearest transcript name. The nearest gene to a given query (in column 3) may or may not be unique.

dist	Signed distance to the nearest target. Queries downstream from (i.e. past) their nearest target are given a negative distance.
matchIndex	The index of the nearest target.
type	one of c("inside", "cover", "disjoint", "overlap").
amountOverlap	The width of the overlap region, if any.
insideDist	When a query is contained in its nearest target, the signed minimum of the two distances target-start-to-query-start and query-end-to-target-end. The former is taken positive, and the latter, which wins in ties, negative. dist will be 0 in this case.
size1	equals width(x).
size2	equals width(subject).
name	nearest gene
annotation	RefSeq ID
description	a factor with levels c("upstream", "promoter", "overlaps 5'", "inside intron", "inside exon",
region	a factor with levels c("upstream", "promoter", "overlaps 5'", "inside", "overlaps 3'", "close
distance	distance before 5' end of gene
subregion	a factor with levels c("inside intron", "inside exon", "covers exon(s)", "overlaps exon upstr
insidedistance	distance past 5' end of gene
exonnumber	which exon
nexons	number of exons
UTR	a factor with levels c("inside transcription region", "5' UTR", "overlaps 5' UTR", "3'UTR", "
strand	"+" or "-"
geneL	the gene length
codingL	the coding length

Author(s)

Harris Jaffee, Peter Murakami and Rafael A. Irizarry

See Also

nearest and distanceToNearest

Examples

```

query <- IRanges(c(1, 4, 9), c(5, 7, 10))
subject <- IRanges(c(2, 2, 10), c(2, 3, 12))
nearest(query, subject)
distanceToNearest(query, subject)
#
## showing 'cover' and 'disjoint', and 'amountOverlap'
annotateNearest(query, subject)
#
## showing 'inside' and 'insideDist', and 'amountOverlap'
annotateNearest(subject, query)
annotateNearest(IRanges(3,3), IRanges(2,5))
annotateNearest(IRanges(3,4), IRanges(2,5))
annotateNearest(IRanges(4,4), IRanges(2,5))

```

bumphunter

Bumphunter

Description

Estimate regions for which a genomic profile deviates from its baseline value. Originally implemented to detect differentially methylated genomic regions between two populations.

Usage

```

## S4 method for signature 'matrix'
bumphunter(object, design, chr=NULL, pos, cluster=NULL,
            coef=2, cutoff=NULL, cutoffQ=0.975, maxGap=500, smooth=TRUE,
            smoothFunction=loessByCluster, useWeights=TRUE, B=100,
            verbose=TRUE, ...)

bumphunterEngine(mat, design, chr = NULL, pos, cluster = NULL, coef = 2,
                 cutoff = NULL, cutoffQ = 0.975, maxGap = 500,
                 smooth = TRUE, smoothFunction = loessByCluster,
                 useWeights = TRUE, B = 100, verbose = TRUE, ...)

```

Arguments

object	An object of class matrix.
mat	A matrix with rows representing genomic locations and columns representing samples.
design	Design matrix with rows representing samples and columns representing covariates. Regression is applied to each row of mat.
chr	A character vector with the chromosomes of each location.
pos	A numeric vector representing the chromosomal position.

cluster	The clusters of locations that are to be analyzed together. In the case of microarrays, the clusters are many times supplied by the manufacturer. If not available the function <code>clusterMaker</code> can be used to cluster nearby locations.
coef	An integer denoting the column of the design matrix containing the covariate of interest. The hunt for bumps will be only be done for the estimate of this coefficient.
cutoff	A numeric value. Values of the estimate of the genomic profile above the cutoff or below the negative of the cutoff will be used as candidate regions. It is possible to give two separate values (upper and lower bounds). If one value is given, the lower bound is minus the value.
cutoffQ	If cutoff is not provided then this quantile of the estimated genomic profile will be used as a cutoff.
maxGap	If cluster is not provided this maximum location gap will be used to define cluster via the <code>clusterMaker</code> function.
smooth	A logical value. If TRUE the estimated profile will be smoothed with the smoother defined by <code>smoothFunction</code>
smoothFunction	A function to be used for smoothing the estimate of the genomic profile. Two functions are provided by the package: <code>loessByCluster</code> and <code>runmedByCluster</code> .
useWeights	A logical value. If TRUE then the standard errors of the point-wise estimates of the profile function will be used as weights in the loess smoother <code>loessByCluster</code> . If the <code>runmedByCluster</code> smoother is used this argument is ignored.
B	An integer denoting the number of resamples to use when computing null distributions: 1000 is recommended but can be quite slow if this process is not parallelized.
verbose	logical value. If TRUE, it writes out some messages indicating progress. If FALSE nothing should be printed.
...	further arguments to be passed to the smoother functions.

Details

This function performs the bump hunting approach described by Jaffe et al. *International Journal of Epidemiology* (2012). The main output is a table of candidate regions with permutation-based family-wide error rates (FWER) and p-values assigned.

The general idea is that for each genomic location we have a value for several individuals. We also have covariates for each individual and perform regression. This gives us one estimate of the coefficient of interest (a common example is case versus control). These estimates are then (optionally) smoothed. The smoothing occurs in clusters of locations that are ‘close enough’. This gives us an estimate of a genomic profile that is 0 when uninteresting. We then take values above (in absolute value) cutoff as candidate regions. Permutations are then performed to create null distributions for the candidate regions.

Uncertainty is assessed via permutations. Each of the B permutations will produce an estimated ‘null profile’ from which we can define ‘null candidate regions’. For each observed candidate region we determine how many null regions are ‘more extreme’ (longer and higher average value). The ‘p.value’ is the percent of candidate regions obtained from the permutations that as extreme as the observed region. These p-values should be interpreted with care as the theoretical properties

are not well understood. The ‘fwer’ is the proportion of permutations that had at least one region as extreme as the observed region. We compute p-values and FWER for the area of the regions (as opposed to length and value as a pair) as well.

Parallelization is implemented through the foreach package.

Value

A list with the following components:

tab	The table with candidate regions and annotation for these.
fitted	The estimated genomic profile used to determine the regions.
pvaluesMarginal	marginal p-value for each genomic location.

Author(s)

Rafael A. Irizarry and Martin J. Aryee

References

Jaffe AE, Murakami P, Lee H, Leek JT, Fallin MD, Feinberg AP, Irizarry RA (2012) Bump hunting to identify differentially methylated regions in epigenetic epidemiology studies. *International Journal of Epidemiology* 41(1):200-9.

Examples

```
dat <- dummyData()
# Enable parallelization
require(doParallel)
registerDoParallel()
# Find bumps
bumps <- bumhunter(dat$mat, design=dat$design, chr=dat$chr, pos=dat$pos,
                  cluster=dat$cluster, coef=2, cutoffQ=0.9,
                  smooth=TRUE, B=250, verbose=TRUE, smoothFunction=loessByCluster)
# cleanup on Windows
bumhunter:::closeSockets()
```

clusterMaker

Make clusters of genomic locations based on distance

Description

Genomic locations are grouped into clusters based on distance: locations that are close to each other are assigned to the same cluster. The operation is performed on each chromosome independently.

Usage

```
clusterMaker(chr, pos, assumeSorted = FALSE, maxGap = 300)
```

Arguments

chr	A vector representing chromosomes. This is usually a character vector, but may be a factor or an integer.
pos	A numeric vector with genomic locations.
assumeSorted	This is a statement that the function may assume that the vector pos is sorted (within each chr). Allowing the function to make this assumption may increase the speed of the function slightly.
maxGap	An integer. Genomic locations within maxGap from each other are placed into the same cluster.

Details

The main purpose of the function is to genomic location into clusters that are close enough to perform operations such as smoothing. A genomic location is a combination of a chromosome (chr) and an integer position (pos). Specifically, genomic intervals are not handled by this function.

Each chromosome is clustered independently from each other. Within each chromosome, clusters are formed in such a way that two positions belong to the same cluster if they are within maxGap of each other.

Value

A vector of integers to be interpreted as IDs for the clusters, such that two genomic positions with the same cluster ID is in the same cluster. Each genomic position receives one integer ID.

Author(s)

Rafael A. Irizarry

Examples

```
N <- 1000
chr <- sample(1:5, N, replace=TRUE)
pos <- round(runif(N, 1, 10^5))
o <- order(chr, pos)
chr <- chr[o]
pos <- pos[o]
regionID <- clusterMaker(chr, pos)
```

dummyData

Generate dummy data for use with bumhunter functions

Description

This function generates a small dummy dataset representing samples from two different groups (cases and controls) that is used in bumhunter examples.

Usage

```
dummyData(n1 = 5, n2 = 5, sd = 0.2, l = 100, spacing = 100,  
          clusterSpacing=1e5, numClusters=5)
```

Arguments

n1	Number of samples in group 1 (controls)
n2	Number of samples in group 2 (cases)
sd	Within group standard deviation to be used when simulating data
l	The number of genomic locations for which to simulate data
spacing	The average spacing between locations. The actual locations have a random component so the actual spacing will be non-uniform
clusterSpacing	The spacing between clusters. (Specifically, the spacing between the first location in each cluster.)
numClusters	Divide the genomic locations into this number of clusters, each of which will contain locations spaced spacing bp apart.

Value

A list containing data that can be used with various bumphunter functions.

mat	A simulated data matrix with rows representing genomic locations and columns representing samples.
design	Design matrix with rows representing samples and columns representing covariates.
chr	A character vector with the chromosomes of each location.
pos	A numeric vector representing the chromosomal position.
cluster	A vector representing the cluster of each locations
n1	Number of samples in cluster 1
n2	Number of samples in cluster 2

Author(s)

Martin J. Aryee

Examples

```
dat <- dummyData()  
names(dat)  
head(dat$pos)
```

getSegments

Segment a vector into positive, zero, and negative regions

Description

Given two cutoffs, L and U, this function divides a numerical vector into contiguous parts that are above U, between L and U, and below L.

Usage

```
getSegments(x, f = NULL, cutoff = quantile(abs(x), 0.99),
            assumeSorted = FALSE, verbose = FALSE)
```

Arguments

x	A numeric vector.
f	A factor used to pre-divide x into pieces. Each piece is then segmented based on the cutoff. Setting this to NULL says that there is no pre-division. Often, clusterMaker is used to define this factor.
cutoff	a numeric vector of length either 1 or 2. If length is 1, U (see details) will be cutoff and L will be -cutoff. Otherwise it specifies L and U. The function will furthermore always use the minimum of cutoff for L and the maximum for U.
assumeSorted	This is a statement that the function may assume that the vector f is sorted. Allowing the function to make this assumption may increase the speed of the function slightly.
verbose	Should the function be verbose?

Details

This function is used to find the indexes of the ‘bumps’ in functions such as [bumphunter](#).

x is a numeric vector, which is converted into three levels depending on whether $x \geq U$ (‘up’), $L < x < U$ (‘zero’) or $x \leq L$ (‘down’), with L and U coming from cutoff. We assume that adjacent entries in x are next to each other in some sense. Segments, consisting of consecutive indices into x (ie. values between 1 and length(x)), are formed such that all indices in the same segment belong to the same level of f and have the same discretized value of x.

In other words, we can use getSegments to find runs of x belonging to the same level of f and with all of the values of x either above U, between L and U, or below L.

Value

A list with three components, each a list of indices. Each component of these lists represents a segment and this segment is represented by a vector of indices into the original vectors x and f.

upIndex: a list with each entry an index of contiguous values in the same segment. These segments have values of x above U.

- dnIndex: a list with each entry an index of contiguous values in the same segment. These segments have values of x below L.
- zeroIndex: a list with each entry an index of contiguous values in the same segment. These segments have values of x between L and U.

Author(s)

Rafael A Irizarry and Kasper Daniel Hansen

See Also

[clusterMaker](#)

Examples

```
x <- 1:100
y <- sin(8*pi*x/100)
chr <- rep(1, length(x))
indexes <- getSegments(y, chr, cutoff=0.8)
plot(x, y, type="n")
for(i in 1:3){
  ind <- indexes[[i]]
  for(j in seq(along=ind)) {
    k <- ind[[j]]
    text(x[k], y[k], j, col=i)
  }
}
abline(h=c(-0.8,0.8))
```

loessByCluster

Apply loess smoothing to values within each spatially-defined cluster.

Description

Loess smoothing is applied independently to each cluster of genomic locations. Locations within the same cluster are close together to warrant smoothing across neighbouring locations.

Usage

```
loessByCluster(y, x = NULL, cluster, weights = NULL, bpSpan = 1000,
               minNum = 7, minInSpan = 5, maxSpan = 1, verbose = TRUE)
```

Arguments

- y A vector or matrix of values to be smoothed. If a matrix, each column represents a sample.
- x The genomic location of the values in y

pointMatch

Fuzzy Matching of Genomic Locations

Description

Finds the nearest locations in a given table of genomic targets

Usage

```
pointMatch(object1, object2, col1=2, col2=2, verbose)
```

Arguments

object1	A data.frame holding the query positions in column col1, on the chromosomes given in the chr column.
object2	As object1, but holding the target positions in the col2 and chr columns.
col1	The column in object1 holding the query positions.
col2	The column in object2 holding the target positions.
verbose	For backward compatibility; ignored.

Details

Similar to `annotateNearest`, but for points instead of ranges.

Value

An integer matrix with columns `c("dist", "matchIndex")`. `dist` is the signed distance from a query to its nearest target, which is given by its index `matchIndex` in `object2`. The distance is negative for queries past their nearest target. A row of NA is returned for queries on chromosomes not in `object2`.

Author(s)

Harris Jaffee, Peter Murakami and Rafael A. Irizarry

See Also

`annotateNearest`, `nearest`, and `distanceToNearest`

Examples

```
object1 <- data.frame(3, chr="chr1")
object2 <- data.frame(c(1,4), chr="chr1")
pointMatch(object1, object2, col1=1, col2=1)
object2 <- data.frame(c(2,4), chr="chr1")
pointMatch(object1, object2, col1=1, col2=1)
```

regionFinder	<i>Find non-zero regions in vector</i>
--------------	--

Description

Find regions for which a numeric vector is above (or below) predefined thresholds.

Usage

```
regionFinder(x, chr, pos, cluster = NULL, y = x, summary = mean,
             ind = seq(along = x), order = TRUE, oneTable = TRUE,
             maxGap = 300, cutoff=quantile(abs(x), 0.99),
             assumeSorted = FALSE, verbose = TRUE)
```

Arguments

x	A numeric vector.
chr	A character vector with the chromosomes of each location.
pos	A numeric vector representing the genomic location.
cluster	The clusters of locations that are to be analyzed together. In the case of microarrays, the cluster is many times supplied by the manufacturer. If not available the function clusterMaker can be used.
y	A numeric vector with same length as x containing values to be averaged for the region summary. See details for more.
summary	The function to be used to construct a summary of the y values for each region.
ind	an optional vector specifying a subset of observations to be used when finding regions.
order	if TRUE then the resulting tables are ordered based on area of each region. Area is defined as the absolute value of the summarized y times the number of features in the regions.
oneTable	if TRUE only one results table is returned. Otherwise, two tables are returned: one for the regions with positive values and one for the negative values.
maxGap	If cluster is not provided this number will be used to define clusters via the clusterMaker function.
cutoff	This argument is passed to getSegments . It represents the upper (and optionally the lower) cutoff for x.
assumeSorted	This argument is passed to getSegments and clusterMaker .
verbose	Should the function be verbose?

Details

This function is used in the final steps of [bumphunter](#). While [bumphunter](#) does many things, such as regression and permutation, [regionFinder](#) simply finds regions that are above a certain threshold (using [getSegments](#)) and summarizes them. The regions are found based on x and the summarized values are based on y (which by default equals x). The summary is used for the ranking so one might, for example, use t-tests to find regions but summarize using effect sizes.

Value

If oneTable is FALSE it returns two tables otherwise it returns one table. The rows of the table are regions. Information on the regions is included in the columns.

Author(s)

Rafael A Irizarry

See Also

[bumphunter](#) for the main usage of this function, [clusterMaker](#) for the typical input to the cluster argument and [getSegments](#) for a function used within regionFinder.

Examples

```
x <- seq(1:1000)
y <- sin(8*pi*x/1000) + rnorm(1000, 0, 0.2)
chr <- rep(c(1,2), each=length(x)/2)
tab <- regionFinder(y, chr, x, cutoff=0.8)
print(tab[tab$L>10,])
```

runmedByCluster	<i>Apply running median smoothing to values within each spatially-defined cluster</i>
-----------------	---

Description

Running median smoothing is applied independently to each cluster of genomic locations. Locations within the same cluster are close together to warrant smoothing across neighbouring locations.

Usage

```
runmedByCluster(y, x = NULL, cluster, weights = NULL, k = 5,
               endrule = "constant", verbose = TRUE)
```

Arguments

y	A vector or matrix of values to be smoothed. If a matrix, each column represents a sample.
x	The genomic location of the values in y.
cluster	A vector indicating clusters of locations. A cluster is typically defined as a region that is small enough that it makes sense to smooth across neighbouring locations. Smoothing will only be applied within a cluster, not across locations from different clusters.
weights	weights used by the smoother.
k	integer width of median window; must be odd. See runmed

endrule	character string indicating how the values at the beginning and the end (of the data) should be treated. See runmed .
verbose	Boolean. Should progress be reported?

Details

This function is typically called by [smoother](#), which is in turn called by [bumphunter](#).

Value

fitted	The smoothed data values
smoothed	A boolean vector indicating whether a given position was smoothed
spans	The span used by the loess smoother. One per cluster.
clusterL	The number of locations in each cluster.
smoother	always set to 'runmed'.

Author(s)

Rafael A. Irizarry

See Also

[smoother](#), [loessByCluster](#). Also see [runmed](#).

Examples

```
dat <- dummyData()
smoothed <- runmedByCluster(y=dat$mat[,1], cluster=dat$cluster,
                           k=5, endrule="constant")
```

smoother	<i>Smooth genomic profiles</i>
----------	--------------------------------

Description

Apply smoothing to values typically representing the difference between two populations across genomic regions.

Usage

```
smoother(y, x = NULL, cluster, weights = NULL, smoothFunction,
         verbose = TRUE, ...)
```

Arguments

<code>y</code>	A vector or matrix of values to be smoothed. If a matrix, each column represents a sample.
<code>x</code>	The genomic location of the values in <code>y</code>
<code>cluster</code>	A vector indicating clusters of locations. A cluster is typically defined as a region that is small enough that it makes sense to smooth across neighbouring locations. Smoothing will only be applied within a cluster, not across locations from different clusters
<code>weights</code>	weights used by the smoother.
<code>smoothFunction</code>	A function to be used for smoothing the estimate of the genomic profile. Two functions are provided by the package: <code>loessByCluster</code> and <code>runmedByCluster</code> .
<code>verbose</code>	Boolean. Should progress be reported?
<code>...</code>	Further arguments to be passed to <code>smoothFunction</code>

Details

This function is typically called by `bumphunter` prior to identifying candidate bump regions. Smoothing is carried out within regions defined by the `cluster` argument.

Value

<code>fitted</code>	The smoothed data values
<code>smoothed</code>	A boolean vector indicating whether a given position was smoothed
<code>spans</code>	The span used by the loess smoother. One per cluster.
<code>clusterL</code>	The number of locations in each cluster.
<code>smoother</code>	The name of the smoother used

Author(s)

Rafael A. Irizarry and Martin J. Aryee

See Also

[loessByCluster](#), [runmedByCluster](#)

Examples

```
dat <- dummyData()

## loessByCluster
smoothed <- smoother(y=dat$mat[,1], cluster=dat$cluster, smoothFunction=loessByCluster,
                    bpSpan = 1000, minNum=7, minInSpan=5, maxSpan=1)

## runmedByCluster
smoothed <- smoother(y=dat$mat[,1], cluster=dat$cluster, smoothFunction=runmedByCluster,
                    k=5, endrule="constant")
```

TT	<i>HG19 Transcripts</i>
----	-------------------------

Description

Database of transcripts associated with "known" hg19 genes, namely those with Entrez ID gene identifiers associated.

Usage

```
data(TT)
```

Format

A [GRanges](#) object with 8 metadata columns, "CSS" (coding start), "CSE" (coding end), "Tx" (transcript name), "Entrez" (Entrez ID), "Gene" (gene name), "Refseq" (Refseq number), "Nexons" (number of exons), and "Exons" (the exons themselves, as [IRanges](#) objects). Note that CSS is always less than CSE, even on minus strands where their "start" and "end" meanings are reversed. Similarly with the exons.

Source

```
the result of bumphunter:::known_transcripts()
```

Index

*Topic **hg19, known transcripts**

TT, [17](#)

annotateNearest, [2](#)

bumphunter, [4](#), [9](#), [11](#), [13–15](#)

bumphunter, matrix-method (bumphunter), [4](#)

bumphunter-methods (bumphunter), [4](#)

bumphunterEngine (bumphunter), [4](#)

clusterMaker, [5](#), [6](#), [9](#), [10](#), [13](#), [14](#)

dummyData, [7](#)

getSegments, [9](#), [13](#), [14](#)

GRanges, [17](#)

loessByCluster, [10](#), [15](#), [16](#)

matchGenes (annotateNearest), [2](#)

pointMatch, [12](#)

regionFinder, [13](#)

regionMatch (annotateNearest), [2](#)

runmed, [14](#), [15](#)

runmedByCluster, [11](#), [14](#), [16](#)

smoother, [11](#), [15](#), [15](#)

TT, [17](#)