Package 'monaLisa'

October 16, 2025

Type Package

Title Binned Motif Enrichment Analysis and Visualization

Version 1.14.1

Description Useful functions to work with sequence motifs in the analysis of genomics data. These include methods to annotate genomic regions or sequences with predicted motif hits and to identify motifs that drive observed changes in accessibility or expression. Functions to produce informative visualizations of the obtained results are also provided.

Depends R (>= 4.1)

Imports BiocGenerics, BiocParallel, Biostrings, BSgenome, circlize, ComplexHeatmap (>= 2.11.1), GenomeInfoDb, GenomicRanges, cli, ggplot2 (>= 4.0.0), glmnet, grDevices, grid, IRanges, methods, rlang, RSQLite, stabs, stats, SummarizedExperiment, S4Vectors, TFBSTools, tidyr, tools, utils, XVector

Suggests BiocManager, BiocStyle, BSgenome.Mmusculus.UCSC.mm10, ggrepel, gridExtra, JASPAR2020, JASPAR2024, knitr, rmarkdown, testthat, TxDb.Mmusculus.UCSC.mm10.knownGene

License GPL (>= 3) Encoding UTF-8 RoxygenNote 7.3.3

VignetteBuilder knitr

biocViews MotifAnnotation, Visualization, FeatureExtraction, Epigenetics

Config/testthat/edition 3

URL https://github.com/fmicompbio/monaLisa,
 https://bioconductor.org/packages/monaLisa/,
 https://fmicompbio.github.io/monaLisa/

BugReports https://github.com/fmicompbio/monaLisa/issues git_url https://git.bioconductor.org/packages/monaLisa git_branch RELEASE_3_21 2 Contents

| git_last_commit d208b44 |
|--|
| git_last_commit_date 2025-09-15 |
| Repository Bioconductor 3.21 |
| Date/Publication 2025-10-15 |
| Author Dania Machlab [aut] (ORCID: https://orcid.org/0000-0002-2578-6930) Lukas Burger [aut] (ORCID: https://orcid.org/0000-0003-3833-2169), Charlotte Soneson [aut] (ORCID: https://orcid.org/0000-0003-3833-2169), Dany Mukesha [ctb] (ORCID: https://orcid.org/0009-0001-9514-751X), Michael Stadler [aut, cre] (ORCID: https://orcid.org/0000-0002-2269-4934) |

Maintainer Michael Stadler <michael.stadler@fmi.ch>

Contents

| monaLisa-package | |
|----------------------------|--|
| .calcKmerEnrichment | |
| .calcMotifEnrichment | |
| .calculateGCweight | |
| .checkDfValidity | |
| .checkIfSeqsAreEqualLength | |
| .cons2matrix | |
| .defineBackground | |
| .filterSeqs | |
| .glmnetRandomizedLasso | |
| .iterativeNormForKmers | |
| .normForKmers | |
| annoSeqlogo | |
| bin | |
| calcBinnedKmerEnr | |
| calcBinnedMotifEnrHomer | |
| calcBinnedMotifEnrR | |
| dumpJaspar | |
| findHomer | |
| findMotifHits | |
| getColsByBin | |
| getKmerFreq | |
| getSetZeroBin | |
| homerToPFMatrixList | |
| motifKmerSimilarity | |
| motifSimilarity | |
| parseHomerOutput | |
| plotBinDensity | |
| plotBinDiagnostics | |
| plotBinHist | |
| plotBinScatter | |

| monaLisa-packa | nga | 2 |
|----------------|-----|---|
| шонашва-раска | age | J |

| | plotMotifHeatmaps | 42 |
|-------|--|------------|
| | plotSelectionProb | 45 |
| | plotStabilityPaths | 46 |
| | prepareHomer | 48 |
| | randLassoStabSel | 5 0 |
| | sampleRandomRegions | 5 3 |
| | seqLogoGrob | 54 |
| Index | | 55 |
| | | |
| monaL | isa-package monaLisa - MOtif aNAlysis with Lisa. | |

Description

monaLisa is a collection of tools that simplify motif enrichment analyses in genomic regions of interest.

Details

She makes use of her father Homer (http://homer.ucsd.edu/homer/index.html) and other algorithms to search for motif hits and look for enriched motifs in sets of genomic regions, compared to all other regions.

Known motifs can for example be obtained from a collection of transcription factor binding site specificities, such as **JASPAR2020**.

Author(s)

Dania Machlab, Lukas Burger, Charlotte Soneson and Michael Stadler

See Also

Useful links:

- https://github.com/fmicompbio/monaLisa
- https://bioconductor.org/packages/monaLisa/
- https://fmicompbio.github.io/monaLisa/
- Report bugs at https://github.com/fmicompbio/monaLisa/issues

4 .calcKmerEnrichment

.calcKmerEnrichment Calculate k-mer enrichment

Description

Given sequences, foreground/background labels and weights, calculate the enrichment of each kmer in foreground compared to background. This function is called by calcBinnedKmerEnr() for each bin if background != "model".

The default type of test is "fisher". Alternatively, a binomial test can be used by test = "binomial". Using Fisher's exact test has the advantage that special cases such as zero background counts are handled without ad-hoc adjustments to the k-mer frequencies.

For test = "fisher", fisher. test is used with alternative = "greater", making it a one-sided test for enrichment, as is the case with the binomial test.

Usage

```
.calcKmerEnrichment(k, df, test = c("fisher", "binomial"), verbose = FALSE)
```

Arguments

k Numeric scalar giving the length of k-mers to analyze.

df A DataFrame with sequence information as returned by .iterativeNormForKmers().

test Type of motif enrichment test to perform.

verbose A logical scalar. If TRUE, report on progress.

Details

The function works in ZOOPS mode, which means only one or zero occurrences of a k-mer are considered per sequence. This is helpful to reduce the impact of simple sequence repeats occurring in few sequences.

Value

A data. frame containing the motifs as rows and the columns:

motifName: the motif name

logP : the log p-value for enrichment (natural logarithm). If test="binomial" (default), this log p-value is identical to the one returned by Homer.

sumForegroundWgtWithHits: the weighted number of k-mer hits in foreground sequences.

sumBackgroundWgtWithHits: the weighted number of k-mer hits in background sequences.

totalWgtForeground: the total sum of weights of foreground sequences.

totalWgtBackground: the total sum of weights of background sequences.

.calcMotifEnrichment 5

.calcMotifEnrichment Calculate motif enrichment

Description

Given motif counts, foreground/background labels and weights for a set of sequences, calculate the enrichment of each motif in foreground compared to background. This function is called by calcBinnedMotifEnrR() for each bin.

The default type of test is "fisher", which is also what Homer uses if "-h" is specified for a hypergeometric test. Alternatively, a binomial test can be used by test = "binomial" (what Homer does by default). Using Fisher's exact test has the advantage that special cases such as zero background counts are handled without ad-hoc adjustments to the frequencies.

For test = "fisher", fisher.test is used with alternative = "greater", making it a one-sided test for enrichment, as is the case with the binomial test.

Usage

```
.calcMotifEnrichment(
  motifHitMatrix,
  df,
  test = c("fisher", "binomial"),
  verbose = FALSE
)
```

Arguments

motifHitMatrix Matrix with 0 and 1 entries for absence or presence of motif hits in each sequence.

df A DataFrame with sequence information as returned by .iterativeNormForKmers().

test Type of motif enrichment test to perform.
verbose A logical scalar. If TRUE, report on progress.

Value

A data. frame containing the motifs as rows and the columns:

motifName: the motif name

logP: the log p-value for enrichment (natural logarithm). If test="binomial" (default), this log p-value is identical to the one returned by Homer.

sumForegroundWgtWithHits: the sum of the weights of the foreground sequences that have at least one instance of a specific motif (ZOOPS mode).

sumBackgroundWgtWithHits: the sum of the weights of the background sequences that have at least one instance of a specific motif (ZOOPS mode).

totalWgtForeground: the total sum of weights of foreground sequences.

totalWgtBackground: the total sum of weights of background sequences.

6 .checkDfValidity

.calculateGCweight

Get background sequence weights for GC bins

Description

The logic is based on Homer (version 4.11). All sequences binned depending on GC content (GCbreaks). The numbers of foreground and background sequences in each bin are counted, and weights for background sequences in bin i are defined as: weight_i = (number_fg_seqs_i / number_bg_seqs_i) * (number_bg_seqs_total / number_fg_seqs_total)

Usage

```
.calculateGCweight(
   df,
   GCbreaks = c(0.2, 0.25, 0.3, 0.35, 0.4, 0.45, 0.5, 0.6, 0.7, 0.8),
   verbose = FALSE
)
```

Arguments

df A DataFrame with sequence information.

GCbreaks The breaks between GC bins. The default value is based on the hard-coded bins

used in Homer.

verbose A logical scalar. If TRUE, report on GC weight calculation.

Value

A DataFrame of the same dimensions as the input df, with the columns GCfrac, GCbin and GCwgt filled in with the sequence GC content, assigned GC bins and weights to correct differences in GC distributions between foreground and background sequences.

.checkDfValidity

Check if seginfo DataFrame is valid

Description

Check if the DataFrame with sequence information is valid, i.e. is of the correct object type (DataFrame) and has all expected columns and attributes.

Usage

```
.checkDfValidity(df)
```

Arguments

df Input object to be checked. It should have an attribute err and columns:

seqs: a DNAStringSet object.

isForeground that indicates if a sequence is in the foreground group.

GCfrac : the fraction of G+C bases per sequence.

GCbin: the GC bin for each sequence.

GCwgt: the sequence weight to adjust for GC differences between foreground

and background sequences.

seqWgt: the sequence weight to adjust for k-mer differences between foreground and background sequences.

Value

TRUE (invisibly) if df is valid, otherwise it raises an exception using cli::cli_abort()

.checkIfSeqsAreEqualLength

Check if elements of 'x' are have equal lengths

Description

Check if the elements of 'x' are all equally long. If not, generate a warning.

Usage

.checkIfSeqsAreEqualLength(x)

Arguments

x An object that implements a width method, typically a GRanges or DNAStringSet object.

Value

NULL (invisibly). The function is called for its side-effect of generating a warning if elements of the input are not of equal lengths.

8 .defineBackground

.cons2matrix

Create matrix from consensus sequence

Description

Given a nucleotide sequence of A,C,G,T letter corresponding to a motif's consensus string, construct a positional frequency matrix. This matrix can for example be used as the profileMatrix argument in the constructor for a TFBSTools::PFMatrix object.

Usage

```
.cons2matrix(x, n = 100L)
```

Arguments

x Character scalar with the motif the consensus sequence.

n Integer scalar giving the columns sums in the constructed matrix (number of observed bases at each position).

Value

A positional frequency matrix.

. define Background

Define background sequence set for a single motif enrichment calculation

Description

Define the background set for the motif enrichment calculation in a single bin, depending on the background mode and given foreground sequences.

Usage

```
.defineBackground(
    sqs,
    bns,
    bg,
    currbn,
    gnm,
    gnm.regions,
    gnm.oversample,
    maxFracN = 0.7,
    GCbreaks = c(0.2, 0.25, 0.3, 0.35, 0.4, 0.45, 0.5, 0.6, 0.7, 0.8)
)
```

.filterSeqs 9

Arguments

sqs, bns, bg The seqs, bins and background arguments from calcBinnedMotifEnrR. currbn An integer scalar with the current bin defining the foreground sequences.

gnm, gnm.regions, gnm.oversample

The genome, genome.regions and genome.oversample arguments from calcBinnedMotifEnrR.

maxFracN The maxFracN argument from calcBinnedMotifEnrR.

GCbreaks The breaks between GC bins. The default value is based on the hard-coded bins

used in Homer.

Value

A DataFrame with sequences represented by rows and columns seqs, isForeground, GCfrac, GCbin, GCwgt and seqWgt. Only the first three are already filled in.

.filterSeqs Filter Sequences

Description

Filter sequences that are unlikely to be useful for motif enrichment analysis. The current defaults are based on HOMER (version 4.11).

Usage

```
.filterSeqs(
   seqs,
   maxFracN = 0.7,
   minLength = 5L,
   maxLength = 100000L,
   verbose = FALSE
)
```

Arguments

seqs A DNAStringSet object.

maxFracN A numeric scalar with the maximal fraction of N bases allowed in a sequence

(defaults to 0.7).

minLength The minimum sequence length (default from Homer). Sequences shorter than

this will be filtered out.

maxLength The maximum sequence length (default from Homer). Sequences bigger than

this will be filtered out.

verbose A logical scalar. If TRUE, report on filtering.

Details

The filtering logic is based on removePoorSeq.pl from Homer.

Value

A logical vector of the same length as seqs with TRUE indicated to keep the sequence and FALSE to filter it out.

```
.glmnetRandomizedLasso
```

Randomized Lasso

Description

This function performs randomized lasso using the glmnet package. The function present in the stabs package that runs the lasso version was adapted for the randomized lasso here. Randomized lasso stability selection uses this function repeatedly to select predictors.

Usage

```
.glmnetRandomizedLasso(
    x,
    y,
    q,
    weakness = 1,
    type = c("conservative", "anticonservative"),
    ...
)
```

Arguments

| X | The predictor matrix. Passed to x of glmnet. lasso from stabs package. |
|----------|---|
| у | The response vector. Passed to y of glmnet.lasso from stabs package. |
| q | The number of variables that are selected on each subsample. Passed to ${\tt q}$ of ${\tt glmnet.lasso}$ from stabs package. |
| weakness | Weakness parameter used in randomized lasso (see details). |
| type | Parameter passed to type of glmnet.lasso from stabs package. It is a character vector specifying how much the PFER should be controlled. If type is "conservative" (default), then the number of selected variables per subsample is $<=q$. If type is "anticonservative" then the number of selected variables per subsample is $>=q$. |

Additional parameters for glmnet.

.iterativeNormForKmers 11

Details

This function is identical to glmnet.lasso from the stabs package. The only addition/modification is the weakness parameter which has been added when calling the glmnet function by setting penalty.factor = 1/runif(ncol(x)), weakness, 1), where ncol(x) is the number of predictors.

Value

The regression output which consists of a list of length 2. The list contains the following:

selected - a logical vector of length equal to the total number of predictors. The predictors that were chosen have a value of TRUE.

path - a logical matrix containing the regularization steps as columns and the predictors as rows.An entry of TRUE indicates selection.

See Also

```
glmnet.lasso and glmnet
```

.iterativeNormForKmers

Adjust for k-mer composition (multiple iterations)

Description

Here we run '.normForKmers' multiple times to converge to the final weights that will be used to correct the background sequences for k-mer composition differences compared to the foreground. We closely follow HOMER's normalizeSequence() function found in Motif2.cpp. Note that HOMER runs the normalizeSequence() one last time after going through all iterations or reaching a low error, which we do not do here.

Usage

```
.iterativeNormForKmers(
   df,
   maxKmerSize = 3L,
   minSeqWgt = 0.001,
   maxIter = 160L,
   verbose = FALSE
)
```

Arguments

df
maxKmerSize

A DataFrame with sequence information as returned by .calculateGCweight.

Integer scalar giving the maximum k-mer size to consider. The default is set to 3 (like in HOMER), meaning that k-mers of size 1, 2 and 3 are considered.

12 .normForKmers

minSeqWgt Numeric scalar greater than zero giving the minimal weight of a sequence. The

default value (0.001) was also used by HOMER (HOMER_MINIMUM_SEQ_WEIGHT

constant in Motif2.h).

maxIter An integer scalar giving the maximum number if times to run .normForKmers.

the default is set to 160 (as in HOMER).

verbose A logical scalar. If TRUE, report on k-mer composition adjustment.

Value

A DataFrame containing:

sequenceWeights: a dataframe containing the sequence GC content, GC bins they were assigned to, the weight to correct for GC differences between foreGround and background sequences, the weight to adjust for kmer composition, and the the error term

sequenceNucleotides: a DNAStringSet object containing the raw sequences

.normForKmers

Adjust for k-mer composition (single iteration)

Description

Adjust background sequence weights for differences in k-mer composition compared to the foreground sequences. This function implements a single iteration, and is called iteratively by .iterativeNormForKmers to get to the final set of adjusted weights, which will be the result of adjusting for GC and k-mer composition. The logic is based on Homer's normalizeSequenceIteration() function found in Motif2.cpp.

Usage

```
.normForKmers(
  kmerFreq,
  goodKmers,
  kmerRC,
  seqWgt,
  isForeground,
  minSeqWgt = 0.001,
  maxSeqWgt = 1000
)
```

Arguments

kmerFreq A list with of matrices. The matrix at index i in the list contains the probability

of k-mers of length i, for each k-mer (columns) and sequence (rows).

goodKmers A list of numeric vectors; the element at index i contains the number of good

(non-N-containing) k-mers of length i for each sequence.

kmerRC A list of character vectors; the element at index i contains the reverse comple-

ment sequences of all k-mers of length i.

annoSeqlogo 13

| seqWgt | A numeric vector with starting sequence weights at the beginning of the iteration. |
|--------------|---|
| isForeground | Logical vector of the same length as seqs. TRUE indicates that the sequence is from the foreground, FALSE that it is a background sequence. |
| minSeqWgt | Numeric scalar greater than zero giving the minimal weight of a sequence. The default value (0.001) is based on Homer (HOMER_MINIMUM_SEQ_WEIGHT constant in Motif2.h). |
| maxSeqWgt | Numeric scalar greater than zero giving the maximal weight of a sequence. The default value (1000) is based on HOMER (1/HOMER_MINIMUM_SEQ_WEIGHT constant in Motif2.h). |

Value

A named list with elements seqWgt (updated weights) and err (error measuring difference of foreground and weighted background sequence compositions).

annoSeqlogo Sequence logo annotation

Description

Create an annotation for a Heatmap containing sequence logos.

Usage

```
annoSeqlogo(
  grobL,
  which = c("column", "row"),
  space = unit(0.5, "mm"),
  width = NULL,
  height = NULL,
  gp = gpar(fill = NA, col = NA)
)
```

Arguments

| grobL | A list of sequence logo grobs, typically created using seqLogoGrob. |
|--------|--|
| which | Whether it is a column annotation or a row annotation? |
| space | The space around the image to the annotation grid borders. The value should be a unit object. |
| width | Width of the annotation. The value should be an absolute unit. Width is not allowed to be set for column annotation. |
| height | Height of the annotation. The value should be an absolute unit. Height is not allowed to be set for row annotation. |
| gp | Graphic parameters for annotation grids. Can be used to control the background color in the annotation grids. |

bin bin

Value

An annotation function which can be used in HeatmapAnnotation.

Examples

```
if (require(JASPAR2020) && require(TFBSTools) && require(gridExtra)) {
    pfm1 <- getMatrixByID(JASPAR2020, "MA0139")

    g1 <- seqLogoGrob(pfm1)

    anno <- annoSeqlogo(list(g1))
}</pre>
```

bin

Bin elements of x.

Description

bin groups elements of x into bins with either a constant number of elements per bin, a constant bin width or according to user-provided bin boundaries.

Usage

```
bin(
    x,
    binmode = c("equalN", "equalWidth", "breaks"),
    nElements = round(length(x)/5),
    nBins = NULL,
    minAbsX = NULL,
    breaks = NULL,
    ...
)
```

Arguments

| х | A numerical vector with the values used for binning. |
|-----------|--|
| binmode | The algorithm to be used for binning. Possible values are: "equalN" (default), "equalWidth" or "breaks" (see Details). |
| nElements | The number of elements per bin (only for binmode="equalN"). The width of bins is adjusted accordingly. |
| nBins | The number of bins (only for binmode="equalWidth"). The number of elements per bin will be variable. |
| minAbsX | The minimal absolute value in x for elements to be binned using the binmode="equalN" or binmode="equalWidth" (ignored for other values of binmode). Elements with x values in [-minAbsX, minAbsX] will be collected in a single bin. |

| breaks | Numerical vector with bin boundaries (only for binmode="breaks"). breaks |
|--------|--|
| | has to be ordered and strictly increasing, and has to be of length (number of |
| | bins) + 1. |
| | Further arguments to be passed to cut(x, breaks,include.lowest = TRUE,), such as labels=FALSE. |

Details

Elements are binned according to the values in x depending on binmode:

equalN Items are grouped into a variable number of bins with nElements elements each. If minAbsX is not NULL, elements with x-values in [-minAbsX, minAbsX] will first be collected in a single bin before binning the remaining elements. The boundaries of this single bin may be slightly adjusted in order to respect the nElements elements in the other bins.

equalWidth Items are group into nBins bins with a variable number of elements each.

breaks Items are grouped into bins using cut(x, breaks, include.lowest = TRUE)

Value

The return value from cut(x, ...), typically a factor of the same length as x. Binning mode, bin boundaries and the "neutral" bin are available from attr(..., "binmode"), attr(..., "breaks") and attr(..., "bin0"). For binmode = "breaks", the latter will be NA.

See Also

cut which is used internally.

Examples

```
set.seed(1)
x <- rnorm(100)
summary(bin(x, "equalN", nElements=10))
summary(bin(x, "equalN", nElements=10, minAbsX=0.5))
summary(bin(x, "equalWidth", nBins=5))
summary(bin(x, "breaks", breaks=c(-10,-1,0,1,10)))</pre>
```

calcBinnedKmerEnr

Calculate k-mer enrichment in bins of sequences.

Description

Given a set of sequences and corresponding bins, identify enriched k-mers (n-grams) in each bin. The sequences can be given either directly or as genomic coordinates.

Usage

```
calcBinnedKmerEnr(
  seqs,
  bins = NULL,
  kmerLen = 5,
  background = c("otherBins", "allBins", "zeroBin", "genome", "model"),
 MMorder = 1,
  test = c("fisher", "binomial"),
  includeRevComp = TRUE,
 maxFracN = 0.7,
 maxKmerSize = 3L,
 GCbreaks = c(0.2, 0.25, 0.3, 0.35, 0.4, 0.45, 0.5, 0.6, 0.7, 0.8)
  pseudocount.kmers = 1,
  pseudocount.log2enr = 8,
  p.adjust.method = "BH",
  genome = NULL,
  genome.regions = NULL,
  genome.oversample = 2,
 BPPARAM = SerialParam(),
  verbose = FALSE
)
```

Arguments

seqs DNAStringSet object with sequences to test

bins Factor of the same length and order as seqs, indicating the bin for each se-

quence. Typically the return value of bin. For background = "genome" or

background = "model", bins can be omitted.

kmerLen A numeric scalar giving the k-mer length.

background A character scalar specifying the background sequences to use. One of "otherBins"

(default), "allBins", "zeroBin", "genome" or "model" (see "Details").

MMorder A numeric scalar giving the order of the Markov model used to calculate the

expected frequencies for background = "model".

test A character scalar specifying the type of enrichment test to perform. One of

"fisher" (default) or "binomial". The enrichment test is one-sided (enriched

in foreground).

includeRevComp A logical scalar. If TRUE (default), count k-mer occurrences in both seqs and

their reverse-complement, by concatenating seqs and their reverse-complemented versions before the counting. This is useful if motifs can be expected to occur on any strand (e.g. DNA sequences of ChIP-seq peaks). If motifs are only expected on the forward strand (e.g. RNA sequences of CLIP-seq peaks), includeRevComp = FALSE should be used. Note that bins will be recycled for the reverse complemented sequences, which means that each reverse-complemented sequence will be assigned to the same bib as the corresponding forward se-

quence.

maxFracN A numeric scalar with the maximal fraction of N bases allowed in a sequence

(defaults to 0.7). Sequences with higher fractions are excluded from the analy-

sis.

maxKmerSize The maximum k-mer size to consider, when adjusting background sequence

weights for k-mer composition compared to the foreground sequences. The

default value (3) will correct for mono-, di- and tri-mer composition.

GCbreaks The breaks between GC bins. The default value is based on the hard-coded bins

used in Homer.

pseudocount.kmers

A numeric scalar - will be added to the observed and expected counts for each

k-mer to avoid zero values.

pseudocount.log2enr

A numerical scalar with the pseudocount to add to foreground and background

counts when calculating log2 motif enrichments

p.adjust.method

A character scalar selecting the p value adjustment method (used in p. adjust).

genome A BSgenome or DNAStringSet object with the genome sequence. Only used for

background = "genome" for extracting background sequences.

genome regions An optional GRanges object defining the intervals in genome from which back-

ground sequences are sampled for background = "genome". If NULL, back-

ground sequences are sampled randomly from genome.

genome.oversample

A numeric scalar of at least 1.0 defining how many background sequences will

be sampled per foreground sequence for background = "genome". Larger values will take longer but improve the sequence composition similarity between

foreground and background (see "Details").

BPPARAM An optional BiocParallelParam instance determining the parallel back-end to

be used during evaluation.

verbose A logical scalar. If TRUE, report on progress.

Details

This function implements a binned k-mer enrichment analysis. In each enrichment analysis, the sequences in a specific bin are used as foreground sequences to test for k-mer enrichments comparing to background sequences (defined by background, see below), similarly as in done for motifs in calcBinnedMotifEnrR. Sequences are weighted to correct for GC and shorter k-mer composition differences between fore- and background sets.

The background sequences are defined according to the value of the background argument:

otherBins: sequences from all other bins (excluding the current bin)

allBins: sequences from all bins (including the current bin)

zeroBin: sequences from the "zero bin", defined by the maxAbsX argument of bin. If bins does not define a "zero bin", for example because it was created by bin(..., maxAbsX = NULL), selecting this background definition will abort with an error.

genome: sequences randomly sampled from the genome (or the intervals defined in genome.regions if given). For each foreground sequence, genome.oversample background sequences of the same size are sampled (on average). From these, one per foreground sequence is selected trying to match the G+C composition. In order to make the sampling deterministic, a seed number needs to be provided to the RNGseed parameter in SerialParam or MulticoreParam when creating the BiocParallelParam instance in BPPARAM.

model: a Markov model of the order MMorder is estimated from the foreground sequences and used to estimate expected k-mer frequencies. K-mer enrichments are then calculated comparing observed to these expected frequencies. In order to make the process deterministic, a seed number needs to be provided to the RNGseed parameter in SerialParam or MulticoreParam when creating the BiocParallelParam instance in BPPARAM.

For each k-mer, the weights of sequences is multiplied with the number of k-mer occurrences in each sequence and summed, separately for foreground (sumForegroundWgtWithHits) and background (sumBackgroundWgtWithHits) sequences. The function works in ZOOPS (Zero-Or-One-Per-Sequence) mode, so at most one occurrence per sequence is counted, which helps reduce the impact of sequence repeats. The total foreground (totalWgtForeground) and background (totalWgtBackground) sum of sequence weights is also calculated. If a k-mer has zero sumForegroundWgtWithHits and sumBackgroundWgtWithHits, then any values (p-values and enrichment) that are calculated using these two numbers are set to NA.

Two statistical tests for the calculation of enrichment log p-value are available: test = "fisher" (default) to perform Fisher's exact tests, or test = "binomial" to perform binomial tests, using:

fisher: fisher.test(x = tab, alternative = "greater"), where tab is the contingency table with the summed weights of sequences in foreground or background sets (rows), and with or without a occurrences of a particular k-mer (columns).

binomial : pbinom(q = sumForegroundWgtWithHits - 1, size = totalWgtForeground, prob =
 sumBackgroundWgtWithHits / totalWgtBackground, lower.tail = FALSE, log.p = TRUE)

Value

A SummarizedExperiment object with motifs in rows and bins in columns, containing seven assays:

negLog10P : -log10 P values

negLog10Padj: -log10 adjusted P values

pearsonResid: k-mer enrichments as Pearson residuals

expForegroundWgtWithHits: expected number of foreground sequences with motif hits

log2enr: k-mer enrichments as log2 ratios

sumForegroundWgtWithHits: Sum of foreground sequence weights in a bin that have k-mer occurrences

sumBackgroundWgtWithHits: Sum of background sequence weights in a bin that have k-mer occurrences

#' The rowData of the object contains annotations (name, PFMs, PWMs and GC fraction) for the k-mers, while the colData slot contains summary information about the bins.

calcBinnedMotifEnrHomer

See Also

getKmerFreq used to calculate k-mer enrichments; getSeq,BSgenome-method which is used to extract sequences from genomepkg if x is a GRanges object; bplapply that is used for parallelization; bin for binning of regions

Examples

```
seqs <- Biostrings::DNAStringSet(c("GCATGCATGC", "CATGCGCATG"))
bins <- factor(1:2)
calcBinnedKmerEnr(seqs = seqs, bins = bins, kmerLen = 3)</pre>
```

calcBinnedMotifEnrHomer

Prepare and run HOMER motif enrichment analysis.

Description

Run complete HOMER motif enrichment analysis, consisting of calls to prepareHomer, system2 and parseHomerOutput. This function requires HOMER to be installed (see http://homer.ucsd.edu/homer/index.html) and the path to the tool to be provided (homerfile argument).

Usage

```
calcBinnedMotifEnrHomer(
   gr,
   b,
   genomedir,
   outdir,
   motifFile,
   homerfile = findHomer(),
   regionsize = "given",
   pseudocount.log2enr = 8,
   p.adjust.method = "BH",
   Ncpu = 2L,
   verbose = FALSE,
   verbose.Homer = FALSE
)
```

Arguments

| gr | A GRanges object (or an object that can be coerced to one) with the genomic regions to analyze. |
|-----------|---|
| b | A vector of the same length as gr that groups its elements into bins (typically a factor, such as the one returned by bin). |
| genomedir | Directory containing sequence files in Fasta format (one per chromosome). |

calcBinnedMotifEnrHomer

outdir A path specifying the folder into which the output files will be written.

motiffile A file with HOMER formatted PWMs to be used in the enrichment analysis.

homerfile Path and file name of the findMotifsGenome.pl HOMER script.

regionsize The peak size to use in HOMER ("given" keeps the coordinate region, an inte-

ger value will keep only that many bases in the region center).

pseudocount.log2enr

A numerical scalar with the pseudocount to add to foreground and background

counts when calculating log2 motif enrichments

p.adjust.method

A character scalar selecting the p value adjustment method (used in p.adjust).

Ncpu Number of parallel threads that HOMER can use. verbose A logical scalar. If TRUE, print progress messages.

verbose. Homer A logical scalar. If TRUE, print the console output when running Homer.

Value

A SummarizedExperiment object with motifs in rows and bins in columns, containing seven assays:

negLog10P: -log10 P values

negLog10Padj: -log10 adjusted P values

pearsonResid: motif enrichments as Pearson residuals

expForegroundWgtWithHits: expected number of foreground sequences with motif hits

log2enr: motif enrichments as log2 ratios

sumForegroundWgtWithHits: Sum of foreground sequence weights in a bin that have motif hitssumBackgroundWgtWithHits: Sum of background sequence weights in a bin that have motif hits

The rowData of the object contains annotations (name, PFMs, PWMs and GC fraction) for the motifs, while the colData slot contains summary information about the bins.

See Also

The functions that are wrapped: prepareHomer, system2 and parseHomerOutput, bin for binning of regions

Examples

calcBinnedMotifEnrR 21

```
# GRanges of regions used in binned motif enrichment analysis
gr <- GenomicRanges::tileGenome(
    seqlengths = c(chr1 = 10000L, chr2 = 10000L, chr3 = 10000L),
    tilewidth = 200, cut.last.tile.in.chrom = TRUE)

# create bins (motif enrichment analysis will be per bin)
bins <- factor(GenomicRanges::seqnames(gr))
table(bins)

# run calcBinnedMotifEnrHomer
outdir <- tempfile()
se <- calcBinnedMotifEnrHomer(gr = gr, b = bins, genomedir = genome,
    outdir = outdir, motifFile = motiffile)
list.files(outdir)
}</pre>
```

calcBinnedMotifEnrR

Binned Motif Enrichment Analysis with monaLisa

Description

This function performs a motif enrichment analysis on bins of sequences. For each bin, the sequences in all other bins are used as background.

Usage

```
calcBinnedMotifEnrR(
  segs,
  bins = NULL,
  pwmL = NULL.
  background = c("otherBins", "allBins", "zeroBin", "genome"),
  test = c("fisher", "binomial"),
  maxFracN = 0.7,
 maxKmerSize = 3L,
 min.score = 10,
 matchMethod = "matchPWM",
  GCbreaks = c(0.2, 0.25, 0.3, 0.35, 0.4, 0.45, 0.5, 0.6, 0.7, 0.8),
  pseudocount.log2enr = 8,
  p.adjust.method = "BH",
  genome = NULL,
  genome.regions = NULL,
  genome.oversample = 2,
  BPPARAM = SerialParam(),
  verbose = FALSE,
)
```

22 calcBinnedMotifEnrR

Arguments

seqs DNAStringSet object with sequences to test

bins Factor of the same length and order as seqs, indicating the bin for each se-

quence. Typically the return value of bin. For background = "genome", bins

can be omitted.

pwmL PWMatrixList with motifs for which to calculate enrichments.

background A character scalar specifying the background sequences to use. One of "otherBins"

(default), "allBins", "zeroBin" or "genome" (see "Details").

test A character scalar specifying the type of enrichment test to perform. One of

"fisher" (default) or "binomial". The enrichment test is one-sided (enriched

in foreground).

maxFracN A numeric scalar with the maximal fraction of N bases allowed in a sequence

(defaults to 0.7). Sequences with higher fractions are excluded from the analy-

sis.

maxKmerSize The maximum k-mer size to consider, when adjusting background sequence

weights for k-mer composition compared to the foreground sequences. The

default value (3) will correct for mono-, di- and tri-mer composition.

min.score The minimal score for motif hits, used in findMotifHits.

matchMethod The method used to scan for motif hits, passed to the method parameter in

findMotifHits.

GCbreaks The breaks between GC bins. The default value is based on the hard-coded bins

used in Homer.

pseudocount.log2enr

A numerical scalar with the pseudocount to add to foreground and background

counts when calculating log2 motif enrichments

p.adjust.method

A character scalar selecting the p value adjustment method (used in p.adjust).

genome A BSgenome or DNAStringSet object with the genome sequence. Only used for

background = "genome" for extracting background sequences.

genome.regions An optional GRanges object defining the intervals in genome from which back-

ground sequences are sampled for background = "genome". If NULL, back-

ground sequences are sampled randomly from genome.

genome.oversample

A numeric scalar of at least 1.0 defining how many background sequences will be sampled per foreground sequence for background = "genome". Larger val-

ues will take longer but improve the sequence composition similarity between

foreground and background (see "Details").

BPPARAM An optional BiocParallelParam instance determining the parallel back-end to

be used during evaluation.

verbose A logical scalar. If TRUE, print progress messages.

... Additional arguments for findMotifHits.

calcBinnedMotifEnrR 23

Details

This function implements a binned motif enrichment analysis. In each enrichment analysis, the sequences in a specific bin are used as foreground sequences to test for motif enrichments comparing to background sequences (defined by background, see below). The logic follows the findMotifsGenome.pl tool from Homer version 4.11, with -size given -nomotif -mknown and additionally -h if using test = "fisher", and gives very similar results. As in the Homer tool, sequences are weighted to correct for GC and k-mer composition differences between fore- and background sets.

The background sequences are defined according to the value of the background argument:

otherBins: sequences from all other bins (excluding the current bin)

allBins: sequences from all bins (including the current bin)

zeroBin: sequences from the "zero bin", defined by the maxAbsX argument of bin. If bins does not define a "zero bin", for example because it was created by bin(..., maxAbsX = NULL), selecting this background definition will abort with an error.

genome: sequences randomly sampled from the genome (or the intervals defined in genome.regions if given). For each foreground sequence, genome.oversample background sequences of the same size are sampled (on average). From these, one per foreground sequence is selected trying to match the G+C composition. In order to make the sampling deterministic, a seed number needs to be provided to the RNGseed parameter in SerialParam or MulticoreParam when creating the BiocParallelParam instance in BPPARAM.

Motif hits are predicted using findMotifHits and multiple hits per sequence are counted as just one hit (ZOOPS mode). For each motif, the weights of sequences that have a hit are summed separately for foreground (sumForegroundWgtWithHits) and background (sumBackgroundWgtWithHits). The total foreground (totalWgtForeground) and background (totalWgtBackground) sum of sequence weights is also calculated. If a motif has zero sumForegroundWgtWithHits and sumBackgroundWgtWithHits, then any values (p-values and enrichment) that are calculated using these two numbers are set to

Two statistical tests for the calculation of enrichment log p-value are available: test = "fisher" (default) to perform Fisher's exact tests, or test = "binomial" to perform binomial tests (default in Homer), using:

fisher: fisher.test(x = tab, alternative = "greater"), where tab is the contingency table with the summed weights of sequences in foreground or background sets (rows), and with or without a hit for a particular motif (columns).

binomial : pbinom(q = sumForegroundWgtWithHits - 1, size = totalWgtForeground, prob =
 sumBackgroundWgtWithHits / totalWgtBackground, lower.tail = FALSE, log.p = TRUE)

Value

A SummarizedExperiment object with motifs in rows and bins in columns, containing seven assays:

negLog10P: -log10P values

negLog10Padj: -log10 adjusted P values

pearsonResid: motif enrichments as Pearson residuals

expForegroundWgtWithHits: expected number of foreground sequences with motif hits

log2enr: motif enrichments as log2 ratios

24 dumpJaspar

sumForegroundWgtWithHits: Sum of foreground sequence weights in a bin that have motif hitssumBackgroundWgtWithHits: Sum of background sequence weights in a bin that have motif hits

The rowData of the object contains annotations (name, PFMs, PWMs and GC fraction) for the motifs, while the colData slot contains summary information about the bins.

Examples

dumpJaspar

Dump Jaspar motifs into a HOMER motif file.

Description

Get motifs from a Jaspar database package (e.g. JASPAR2020) and write them into a HOMER-compatible motif file as positional probability matrices.

Usage

```
dumpJaspar(
  filename,
  pkg = "JASPAR2020",
  opts = list(tax_group = "vertebrates"),
  pseudocount = 1,
  relScoreCutoff = 0.8,
  verbose = FALSE
)
```

Arguments

filename Name of the output file to be created.

pkg Name of the Jaspar package to use (default: JASPAR2020).

findHomer 25

opts A list with search options used in getMatrixSet. By default, only vertebrate

motifs are included in the output using opts = list(tax_group = "vertebrates").

pseudocount A numerical scalar with the pseudocount to be added to each element of the po-

sition frequency matrix extracted from Jaspar, before its conversion to a position

probability matrix (default: 1.0).

relScoreCutoff Currently ignored. numeric(1) in [0,1] that sets the default motif log-odds score

cutof to relScoreCutoff * maximal score for each PWM (default: 0.8).

verbose A logical scalar. If TRUE, print progress messages.

Value

TRUE if successful.

See Also

getMatrixSet for details on the argument opts. homerToPFMatrixList to read a file with HOMER-formatted motifs into a PFMatrixList.

Examples

findHomer

Find HOMER script file.

Description

Find absolute path to HOMER script file.

Usage

```
findHomer(homerfile = "findMotifsGenome.pl", dirs = NULL)
```

Arguments

homerfile Name of the script file to search.

dirs Directory names to look for homerfile. If dirs=NULL, all directories listed in

the PATH environment variable will be searched.

Details

In addition to dirs, findHomer will also look in the directory provided in the environment variable MONALISA_HOMER.

Value

Absolute path to homerfile, or NA if none or several were found.

Examples

```
homer_path <- findHomer()</pre>
```

findMotifHits

Find motif matches in sequences.

Description

findMotifHits scans sequences (either provided as a file, an R object or genomic coordinates) for matches to positional weight matrices (provided as a file or as R objects)

Usage

```
findMotifHits(
  query,
  subject,
 min.score,
 method = c("matchPWM", "homer2"),
  homerfile = findHomer("homer2"),
 BPPARAM = SerialParam(),
  genome = NULL
## S4 method for signature 'character, character'
findMotifHits(
  query,
  subject,
 min.score,
 method = c("matchPWM", "homer2"),
 homerfile = findHomer("homer2"),
 BPPARAM = SerialParam(),
  genome = NULL
## S4 method for signature 'character, DNAString'
findMotifHits(
  query,
  subject,
 min.score,
 method = c("matchPWM", "homer2"),
  homerfile = findHomer("homer2"),
 BPPARAM = SerialParam(),
```

```
genome = NULL
## S4 method for signature 'character, DNAStringSet'
findMotifHits(
  query,
  subject,
 min.score,
 method = c("matchPWM", "homer2"),
 homerfile = findHomer("homer2"),
 BPPARAM = SerialParam(),
  genome = NULL
)
## S4 method for signature 'PWMatrix, character'
findMotifHits(
  query,
  subject,
 min.score,
 method = c("matchPWM", "homer2"),
 homerfile = findHomer("homer2"),
 BPPARAM = SerialParam(),
  genome = NULL
)
## S4 method for signature 'PWMatrix, DNAString'
findMotifHits(
 query,
  subject,
 min.score,
 method = c("matchPWM", "homer2"),
  homerfile = findHomer("homer2"),
 BPPARAM = SerialParam(),
  genome = NULL
)
## S4 method for signature 'PWMatrix, DNAStringSet'
findMotifHits(
  query,
  subject,
 min.score,
 method = c("matchPWM", "homer2"),
  homerfile = findHomer("homer2"),
 BPPARAM = SerialParam(),
  genome = NULL
)
## S4 method for signature 'PWMatrixList, character'
```

```
findMotifHits(
  query,
  subject,
 min.score,
 method = c("matchPWM", "homer2"),
 homerfile = findHomer("homer2"),
 BPPARAM = SerialParam(),
  genome = NULL
)
## S4 method for signature 'PWMatrixList,DNAString'
findMotifHits(
  query,
  subject,
 min.score,
 method = c("matchPWM", "homer2"),
  homerfile = findHomer("homer2"),
 BPPARAM = SerialParam(),
  genome = NULL
)
## S4 method for signature 'PWMatrixList,DNAStringSet'
findMotifHits(
  query,
  subject,
 min.score,
 method = c("matchPWM", "homer2"),
  homerfile = findHomer("homer2"),
 BPPARAM = SerialParam(),
  genome = NULL
)
## S4 method for signature 'PWMatrix, GRanges'
findMotifHits(
  query,
  subject,
 min.score,
 method = c("matchPWM", "homer2"),
  homerfile = findHomer("homer2"),
 BPPARAM = SerialParam(),
  genome = NULL
)
## S4 method for signature 'PWMatrixList, GRanges'
findMotifHits(
  query,
  subject,
  min.score,
```

```
method = c("matchPWM", "homer2"),
homerfile = findHomer("homer2"),
BPPARAM = SerialParam(),
genome = NULL
)
```

Arguments

query The motifs to search for, either a

character(1) with the path and file name of a motif file with PWM in HOMER

format (currently only supported for method="homer2")

 ${\tt PWMatrix} \ \ with \ a \ single \ PWM$

PWMatrixList with several PWMs to search for.

subject The sequences to be searched, either a

character with the path and file name of a sequence file with DNA sequences

in FASTA format

DNAString with a single sequence
DNAStringSet with several sequences

GRanges object with the genomic coordinates of the sequences to be searched.

min.score The minimum score for counting a match. Can be given as a character string

containing a percentage (e.g. "85 highest possible score or as a single number.

method The internal method to use for motif searching. One of

"matchPWM" using Biostrings::matchPWM (optimized)

"homer2" call to the homer2 binary

Please note that the two methods might give slightly different results (see de-

tails).

homerfile Path and file name of the homer2 binary.

BPPARAM An optional BiocParallelParam instance determining the parallel back-end to

be used during evaluation.

genome BSgenome object that is the reference genome of the subject. This argument is

set to NULL by default and only used by the function when the subject is a GRanges object. It is then necessary to specify the genome so that the function

can internally convert the genomic regions into a DNAStringSet object.

Details

The implemented methods (matchPWM and homer2) are there for convenience (method="matchPWM" calls Biostrings::matchPWM internally in an optimized fashion, and method = "homer2" calls the command line tool from Homer and therefore requires an installation of Homer).

In general, running findMotifHits with the same parameters using any of the methods generates identical results. Some minor differences could occur that result from rounding errors during the necessary conversion of PWMs (log2-odd scores) to the probability matrices needed by Homer, and the conversion of scores from and to the natural log scale used by Homer. These conversions are implemented transparently for the user, so that the arguments of findMotifHits do not have to be

30 getColsByBin

adjusted (e.g. the PWMs should always contain log2-odd scores, and min. score is always on the log2 scale).

If there are bases with frequencies of less than 0.001 in a motif, Homer will set them to 0.001 and adjust the other frequencies at that motif position accordingly so that they sum to 1.0. This may differ from the adjustment used when scanning a PWM with matchPWM (e.g. the pseudocounts argument in the toPWM function), and thus can give rise to differences in reported motif hits and hit scores (typically only low-scoring hits).

Value

A GRanges object with the matches to query in subject.

Examples

getColsByBin

Get colors by bin.

Description

Get colors for elements according to their bin. Colors are assigned to bins forming a gradient from col1 to col2 in the order of levels{b}. col0 is assigned to the neutral bin (attribute "") if available.

Usage

```
getColsByBin(
   b,
   col1 = c("#003C30", "#01665E", "#35978F", "#80CDC1", "#C7EAE5"),
   col2 = c("#F6E8C3", "#DFC27D", "#BF812D", "#8C510A", "#543005"),
   col0 = "#F5F5F5"
)
```

getKmerFreq 31

Arguments

| b | A factor that groups elements into bins (typically the output of bin). |
|------|--|
| col1 | First color. |
| col2 | Second color. |
| col0 | Neutral color. |

Value

A character vector with colors for the elements in b.

See Also

bin.

Examples

```
set.seed(1)
x <- rnorm(100)
b <- bin(x, "equalN", nElements = 10)
cols <- getColsByBin(b)</pre>
```

getKmerFreq

Calculate observed and expected k-mer frequencies

Description

Given a set of sequences, calculate observed and expected k-mer frequencies. Expected frequencies are based on a Markov model of order MMorder.

Usage

```
getKmerFreq(
   seqs,
   kmerLen = 5,
   MMorder = 1,
   pseudocount = 1,
   zoops = TRUE,
   strata = rep(1L, length(seqs)),
   p.adjust.method = "BH",
   includeRevComp = TRUE
)
```

32 getSetZeroBin

Arguments

seqs Set of sequences, either a character vector or a DNAStringSet.

kmerLen A numeric scalar giving the k-mer length.

MMorder A numeric scalar giving the order of the Markov model used to calculate the

expected frequencies.

pseudocount A numeric scalar - will be added to the observed counts for each k-mer to avoid

zero values.

zoops A logical scalar. If TRUE (the default), only one or zero occurences of a k-mer

are considered per sequence.

strata A factor or a numeric scalar defining the strata of sequences. A separate

Markov model and expected k-mer frequencies are estimated for the set of sequences in each stratum (level in a strata factor). If strata is a scalar value, it will be interpreted as the number of strata to split the sequences into according to their CpG observed-over-expected counts using kmeans(CpGoe, centers =

strata).

p.adjust.method

A character scalar selecting the p value adjustment method (used in p.adjust).

includeRevComp A logical scalar. If TRUE (default), count k-mer occurrences in both seqs and

their reverse-complement, by concatenating seqs and their reverse-complemented versions before the counting. This is useful if motifs can be expected to occur on any strand (e.g. DNA sequences of ChIP-seq peaks). If motifs are only expected on the forward strand (e.g. RNA sequences of CLIP-seq peaks), includeRevComp = FALSE should be used. Note that if strata is a vector of the same length as seqs, each reverse-complemented sequence will be assigned to

the same stratum as the forward sequence.

Value

A list with observed and expected k-mer frequencies (freq.obs and freq.exp, respectively), and enrichment statistics for each k-mer.

Examples

```
res <- getKmerFreq(seqs = c("AAAAATT", "AAATTTT"), kmerLen = 3)
names(res)
head(res$freq.obs)
head(res$freq.exp)</pre>
```

getSetZeroBin

Get and set the zero bin manually

Description

Get and set the zero bin manually

homerToPFMatrixList 33

Usage

```
getZeroBin(bins)
setZeroBin(bins, zeroBin)
```

Arguments

bins Factor, typically the return value of bin.

zeroBin Numeric or character scalar indicating the level to use as the zero bin, or NA.

Value

For getZeroBin, the index of the level representing the zero bin. For setZeroBin, a modified factor with the zero bin set to the provided value.

Examples

```
set.seed(1)
x <- rnorm(100)
bins <- bin(x, "equalN", nElements = 10, minAbsX = 0.5)
getZeroBin(bins)
bins <- setZeroBin(bins, 2)</pre>
```

homerToPFMatrixList

Read a HOMER motif file and create a PFMatrixList

Description

Read motifs from a file in HOMER format and create a PFMatrixList from them.

Usage

```
homerToPFMatrixList(filename, n = 100L)
```

Arguments

filename Name of the input file with HOMER-formatted motifs.

n The number of observations (multiplied with base frequencies to create the num-

ber of observed bases at each position).

Value

A PFMatrixList with motifs from the file.

See Also

dumpJaspar for writing motifs from a Jaspar database package into a file in HOMER format.

motifKmerSimilarity

Examples

```
library(JASPAR2020)
optsL <- list(ID = c("MA0006.1"))
pfm1 <- TFBSTools::getMatrixSet(JASPAR2020, opts = optsL)
TFBSTools::Matrix(pfm1)

tmpfn <- tempfile()
dumpJaspar(filename = tmpfn, pkg = "JASPAR2020", opts = optsL)
pfm2 <- homerToPFMatrixList(tmpfn)
TFBSTools::Matrix(pfm2)
unlink(tmpfn)</pre>
```

motifKmerSimilarity

Calculate similarities between motifs and k-mers.

Description

For each motif, calculate it's similarity to all k-mers of length kmerLen, defined as the maximal probability of observing the k-mer given the base frequencies of the motif (the maximum is taken over for all possible ungapped alignments between motif and k-mer). If necessary matrices are padded on the sides with background base frequencies (assuming all bases to have a frequency of 0.25).

Usage

```
motifKmerSimilarity(
    x,
    kmerLen = 5,
    kmers = NULL,
    includeRevComp = FALSE,
    BPPARAM = SerialParam(),
    verbose = FALSE
)
```

Arguments

x Either a PFMatrixList, or a character scalar with a file containing motifs in

HOMER format (used directly method = "HOMER", loaded into a PFMatrixList

by homerToPFMatrixList for method = "R").

kmerLen A numeric scalar giving the k-mer length.

kmers Either a character vector of k-mers for which to calculate the similarity to each

motif, or NULL, in which case all k-mers of length kmerLen are used.

includeRevComp A logical scalar. If set to TRUE, each k-mer as well as its reverse complement

is compared to each motif, and the larger of the two similarities is returned.

motifSimilarity 35

BPPARAM An optional BiocParallelParam instance determining the parallel back-end to

be used during evaluation.

verbose A logical scalar. If TRUE, report on progress.

Value

A matrix of probabilties for each motif - k-mer pair.

See Also

bplapply used for parallelization.

Examples

motifSimilarity

Calculate similarities between pairs of motifs.

Description

For each pair of motifs, calculate the similarity defined as the maximal Pearson's correlation coefficient between base frequencies over all possible shifts (relative positions of the two matrices with at least one overlapping position). If necessary matrices are padded on the sides with background base frequencies (assuming all bases to have a frequency of 0.25) to enable comparison of all positions in both matrices.

Usage

```
motifSimilarity(
    x,
    y = NULL,
    method = c("R", "HOMER"),
    homerfile = findHomer("compareMotifs.pl"),
    homerOutfile = NULL,
    BPPARAM = SerialParam(),
    verbose = FALSE
)
```

36 motifSimilarity

Arguments

| x | Either a PFMatrixList, or a character scalar with a file containing motifs in HOMER format (used directly method = "HOMER", loaded into a PFMatrixList by homerToPFMatrixList for method = "R"). |
|--------------|--|
| у | Either a PFMatrixList or NULL (default). If $y = NULL$, then similarities will be calculated for all pairs of motifs within x. Otherwise, method must be "R" and similarities will be calculated between any motif from x to any motif from y. |
| method | A character scalar specifying the method for similarity calculations. Either "R" (pure R implementation) or "HOMER" (will call the compareMotifs.pl script from HOMER). Results are identical (apart from rounding errors), and the R implementation is usually faster and can be parallelized (BPPARAM argument). |
| homerfile | Path to the HOMER script compareMotifs.pl (only used for method = "HOMER". |
| homerOutfile | A character scalar giving the file to save the similarity scores (only for metho = "HOMER"). If NULL, scores will be stored into a temporary file. |
| BPPARAM | An optional BiocParallelParam instance determining the parallel back-end to be used during evaluation (only used for method = "R"). |
| verbose | A logical scalar. If TRUE, report on progress. |

Value

A matrix of Pearson's correlation coefficients for each pair of motifs.

See Also

bplapply used for parallelization for method = "R", documentation of HOMER's compareMotifs.pl
for details on method = "HOMER".

Examples

parseHomerOutput 37

parseHomerOutput

Load output from HOMER findMotifsGenome.pl into R

Description

Parse HOMER output files into R data structures.

Usage

```
parseHomerOutput(infiles, pseudocount.log2enr = 8, p.adjust.method = "BH")
```

Arguments

```
infiles HOMER output files to be parsed. pseudocount.log2enr
```

A numerical scalar with the pseudocount to add to foreground and background counts when calculating log2 motif enrichments

p.adjust.method

A character scalar selecting the p value adjustment method (used in p. adjust).

Value

A list of nine components (negLog10P, negLog10Padj, pearsonResid, expForegroundWgtWithHits, log2enr, sumForegroundWgtWithHits and sumBackgroundWgtWithHits), seven containing each a motif (rows) by bin (columns) matrix with raw -log10 P values, -log10 adjusted P values, the expected number of foreground sequences with hits, the observed number of foreground and background sequences with hits, and motif enrichments as Pearson residuals (pearsonResid) and as log2 ratios (log2enr), and two containing the total foreground and background weight (totalWgtForeground, totalWgtBackground).

38 plotBinDensity

plotBinDensity Density plot of binned elements.

Description

Plot the density of binned elements with binning information.

Usage

```
plotBinDensity(
    x,
    b,
    xlab = deparse(substitute(x, env = as.environment(-1))),
    ylab = "Density",
    main = "",
    legendPosition = "right",
    legend = NULL,
    legend.cex = NULL,
    ...
)
```

Arguments

Х A numerical vector with the values used for binning. b A factor that groups elements of x into bins (typically the output of bin). xlab, ylab, main character scalars that set the x-axis label, y-axis label and the main title. Use "" to suppress the label. legendPosition A character scalar. If not "none", draw a legend with binning information. The value is used to control the legend position and will be passed to theme(legend.position = legendPosition). legend Deprecated (ignored). Please use legendPosition to control the drawing and position of the legend. legend.cex Deprecated (ignored). You can use theme to set legend and other graphical parameters. Further arguments passed to getColsByBin.

Value

The generated density plot as a ggplot object.

See Also

```
getColsByBin, geom_density
```

plotBinDiagnostics 39

Examples

```
set.seed(1)
x <- rnorm(100)
b <- bin(x, "equalN", nElements = 10)
plotBinDensity(x, b)</pre>
```

plotBinDiagnostics

Plot diagnostics of binned sequences

Description

Plot various diagnostics of binned sequences. Three plot types are available:

length plots the distribution of sequence lengths within each bin.

GCfrac plots the distribution of GC fractions within each bin.

dinucfreq plots a heatmap of the relative frequency of each dinucleotide, averaged across the sequences within each bin. The values are centered for each dinucleotide to better highlight differences between the bins. The average relative frequency of each dinucleotide (across the bins) is indicated as well.

Usage

```
plotBinDiagnostics(
   seqs,
   bins,
   aspect = c("length", "GCfrac", "dinucfreq"),
   draw_quantiles = c(0.25, 0.5, 0.75),
   ...
)
```

Arguments

| seqs | DNAStringSet object with sequences. |
|----------------|--|
| bins | Factor of the same length and order as seqs, indicating the bin for each sequence. Typically the return value of bin. |
| aspect | The diagnostic to plot. Should be one of "length", "GCfrac" and "dinucfreq", to plot the distribution of sequence lengths, the distribution of GC fractions and the average relative dinucleotide frequencies across the bins. |
| draw_quantiles | For aspect="length" or "GCfrac", draw vertical lines at the given quantiles of the density estimate. If NULL, no quantile lines will be drawn. |
| | Additional argument passed to getColsByBin. |

Value

For aspect="length" or "GCfrac", returns a ggplot object. For aspect="dinucfreq", returns (invisibly) a Heatmap-class object.

40 plotBinHist

Examples

plotBinHist

Histogram of binned elements.

Description

Plot a histogram of binned elements with binning information.

Usage

```
plotBinHist(
    x,
    b,
    breaks = 6 * nlevels(b),
    xlab = deparse(substitute(x, env = as.environment(-1))),
    ylab = "Frequency",
    main = "",
    legendPosition = "right",
    legend = NULL,
    legend.cex = NULL,
    ...
)
```

Arguments

```
A numerical vector with the values used for binning.

A factor that groups elements of x into bins (typically the output of bin).

A numeric scalar controlling the histogram breaks (passed to geom_hist(..., bins = breaks)).

xlab, ylab, main character scalars that set the x-axis label, y-axis label and the main title. Use "" to suppress the label.

legendPosition A character scalar. If not "none", draw a legend with binning information.

The value is used to control the legend position and will be passed to theme(legend.position = legendPosition).
```

plotBinScatter 41

| legend | Deprecated (ignored). Please use legendPosition to control the drawing and position of the legend. | |
|------------|--|--|
| legend.cex | Deprecated (ignored). You can use theme to set legend and other graphical parameters. | |
| | Further arguments passed to getColsByBin. | |

Value

The generated histogram as a ggplot object.

See Also

```
getColsByBin, geom_histogram
```

Examples

```
set.seed(1)
x <- rnorm(100)
b <- bin(x, "equalN", nElements = 10)
plotBinHist(x, b)</pre>
```

plotBinScatter

Scatter plot (xy-plot) of binned elements.

Description

Plot a scatter (xy-plot) of binned elements with binning information.

Usage

```
plotBinScatter(
    x,
    y,
    b,
    cols = NULL,
    xlab = deparse(substitute(x, env = as.environment(-1))),
    ylab = deparse(substitute(y, env = as.environment(-1))),
    main = "",
    legendPosition = "right",
    legend = NULL,
    legend.cex = NULL,
    ...
)
```

42 plotMotifHeatmaps

Arguments

| x | A numerical vector with x values. |
|----------------|--|
| У | A numerical vector with y values (the values used for binning). |
| b | A factor that groups elements of x, y into bins (typically the output of $bin(y)$). |
| cols | NULL or a color vector defining the colors of points. If NULL, the colors will be computed based on b using getColsByBin(b)). |
| xlab | Label for x-axis. |
| ylab | Label for y-axis. |
| main | Main title. |
| legendPosition | A character scalar. If not "none", draw a legend with binning information. The value is used to control the legend position and will be passed to theme(legend.position = legendPosition). |
| legend | Deprecated (ignored). Please use legendPosition to control the drawing and position of the legend. |
| legend.cex | Deprecated (ignored). You can use theme to set legend and other graphical parameters. |
| • • • | Further arguments passed to getColsByBin (only used if cols is NULL). |

Value

The generated scatter plot as a ggplot object.

See Also

```
bin, getColsByBin, geom_point
```

Examples

```
set.seed(1)
x <- rnorm(100)
y <- rnorm(100)
b <- bin(y, "equalN", nElements = 10)
plotBinScatter(x, y, b)</pre>
```

plotMotifHeatmaps

Heatmap of motif enrichments.

Description

Plot motif enrichments (e.g. significance or magnitude) as a heatmap.

plotMotifHeatmaps 43

Usage

```
plotMotifHeatmaps(
  Х,
  which.plots = c("negLog10P", "pearsonResid", "negLog10Padj", "log2enr"),
  col.enr = c("#053061", "#2166AC", "#4393C3", "#92C5DE", "#D1E5F0", "#F7F7F7",
  "#FDDBC7", "#F4A582", "#D6604D", "#B2182B", "#67001F"),
col.sig = c("#F0F0F0", "#D9D9D9", "#BDBDBD", "#969696", "#737373", "#525252",
    "#252525", "#000000"),
 col.gc = c("#F7FCF5", "#E5F5E0", "#C7E9C0", "#A1D99B", "#74C476", "#41AB5D", "#238B45",
    "#006D2C", "#00441B"),
  maxEnr = NULL,
  maxSig = NULL,
  highlight = NULL,
  cluster = FALSE,
  show_dendrogram = FALSE,
  show_motif_GC = FALSE,
  show_seqlogo = FALSE,
  show_bin_legend = FALSE,
  width.seqlogo = 1.5,
  use_raster = FALSE,
  na_col = "white",
  doPlot = TRUE,
)
```

Arguments

| х | A SummarizedExperiment with numerical matrices (motifs-by-bins) in its assays(), typically the return value of calcBinnedMotifEnrR or calcBinnedMotifEnrHomer. |
|-------------|--|
| which.plots | Selects which heatmaps to plot (one or several from "negLog10P", "negLog10Padj", "pearsonResid" and "log2enr"). |
| width | The width (in inches) of each individual heatmap, without legend. |
| col.enr | Colors used for enrichment heatmap ("pearsonResid" and "log2enr"). |
| col.sig | Colors used for significance hetmaps ("negLog10P" and "negLog10Padj"). |
| col.gc | Colors used for motif GC content (for show_motif_GC = TRUE). |
| maxEnr | Cap color mapping at enrichment = maxEnr (default: 99.5th percentile). |
| maxSig | Cap color mapping at $-\log 10 \text{ P}$ value or $-\log 10 \text{ FDR} = \max \text{Sig}$ (default: 99.5th percentile). |
| highlight | A logical vector indicating motifs to be highlighted. |
| cluster | If TRUE, the order of transcription factors will be determined by hierarchical clustering of the "pearsonResid" component. Alternatively, an hclust-object can be supplied which will determine the motif ordering. No reordering is done for cluster = FALSE. |

44 plotMotifHeatmaps

show_dendrogram If cluster != FALSE, controls whether to show a row dendrogram for the clustering of motifs. Ignored for cluster = FALSE. show_motif_GC If TRUE, show a column with the percent G+C of the motif as part of the heatmap. If TRUE, show a sequence logo next to each motif label. This will likely only show_seqlogo make sense for a heatmap with a low number of motifs. show_bin_legend If TRUE, show a legend for the bin labels. If FALSE (default), the bin legend will be hidden. width.seqlogo The width (in inches) for the longest sequence logo (shorter logos are drawn to TRUE or FALSE (default). Passed to use_raster of Heatmap. use_raster na_col "white" (default). Passed to na_col of Heatmap. doPlot If TRUE (default), plot the generated heatmap(s) using Reduce(ComplexHeatmap::add_heatmap, heatmapList). If FALSE, just return the list of heatmap(s) (heatmapList) in example before), allowing to modify them further before plotting. Further arguments passed to Heatmap when creating the main heatmaps selected

by which.plots. For example, the following will set the font size of the motif names: plotMotifHeatmaps(..., row_names_gp = gpar(fontsize = 12))

Details

The heatmaps are created using the **ComplexHeatmap** package and plotted side-by-side.

Each heatmap will be width inches wide, so the total plot needs a graphics device with a width of at least length(which.plots) * width plus the space used for motif names and legend. The height will be auto-adjusted to the graphics device.

Value

A list of ComplexHeatmap::Heatmap objects.

References

Gu, Z. Complex heatmaps reveal patterns and correlations in multidimensional genomic data. Bioinformatics 2016.

See Also

bin, Heatmap

plotSelectionProb 45

plotSelectionProb

Plot selection probabilities of predictors

Description

This function plots the selection probabilities of predictors (for example the selected motifs), optionally multiplied with either +1 or -1 to give a sense of both the strength and the directionality of the associated effects. The directionality is estimated from the sign of the correlation coefficient between each predictor and the response vector.

Usage

```
plotSelectionProb(
    se,
    directional = TRUE,
    selProbMin = metadata(se)$stabsel.params.cutoff,
    selProbMinPlot = 0.4,
    showSelProbMin = TRUE,
    selColor = "cadetblue",
    notSelColor = "grey",
    selProbCutoffColor = "firebrick",
    method = c("pearson", "kendall", "spearman"),
    ylimext = 0.2
)
```

Arguments

| se | The SummarizedExperiment object with the results from stability selection |
|----|---|
| | |

(typically returned by randLassoStabSel).

directional A logical scalar. If TRUE, selection probabilities are plotted with the sign of the

marginal correlation between a predictor and the response.

selProbMin A numerical scalar in [0,1]. Predictors with a selection probability greater than

selProbMin are considered selected and colored by the input from selColor.

By default, selProbMin is extracted from the parameters stored in se.

selProbMinPlot A numerical scalar in [0,1] less than selProbMin. Predictors with a selection

probability greater than selProbMinPlot but less than selProbMin are shown as bars with color defined in notSelColor. selProbMinPlot is useful in order to include additional predictors in the barplot, that were not selected according to selProbMin but may be close to that cutoff or are simply nice to visualize alongside the selected predictors. Setting selProbMinPlot = 0 will include all

predictors.

showSelProbMin A logical scalar. If TRUE, the value of selProbMin is shown by a horizontal line

with the color defined by selProbCutoffColor.

selColor Color for the selected predictors which have a selection probability greater than

selProbMin.

46 plotStabilityPaths

notSelColor Color for the rest of the (unselected) predictors which will be show in the barplot. selProbCutoffColor

Color for the line depicting the selection probability cutoff.

method A character scalar with the correlation method to use in the calculation of predictor-

response marginal correlations. One of "pearson", "kendall" or "spearman" (see

cor).

ylimext A numeric scalar defining how much the y axis limits should be expanded be-

yond the plotted probabilities to allow for space for the bar labels. This value can be increased if the predictor names above the bars are too long and not showing

in the plot.

Details

This function creates a bar plot with ggplot. Each bar corresponds to a predictor (motif) and the colors correspond to whether or not it was selected. The y-axis shows the selection probabilities (directional=FALSE) or selection probabilities with the sign of the marginal correlation to the response (directional=TRUE).

Value

a ggplot2 object.

```
## create data set
set.seed(321)
Y <- rnorm(n = 500, mean = 2, sd = 1)
X <- matrix(data = NA, nrow = length(Y), ncol = 50)
for (i in seq_len(ncol(X))) {
    X[ ,i] <- runif(n = 500, min = 0, max = 3)
}
s_cols <- sample(x = seq_len(ncol(X)), size = 10,
    replace = FALSE)
for (s in s_cols) {
    X[ ,s] <- X[, s] + (Y + rnorm(500, 0, 4)) * ifelse(s %% 2, -1, 1)
}
## reproducible randLassoStabSel() with 1 core
set.seed(123)
ss <- randLassoStabSel(x = X, y = Y)
plotSelectionProb(ss)</pre>
```

plotStabilityPaths 47

Description

Plot the stability paths of each variable (predictor), showing the selection probability as a function of the regularization step.

Usage

```
plotStabilityPaths(
    se,
    selProbMin = metadata(se)$stabsel.params.cutoff,
    selColor = "cadetblue",
    notSelColor = "grey",
    selProbCutoffColor = "firebrick",
    linewidth = 0.5,
    alpha = 1,
    ylim = c(0, 1),
    labelPaths = FALSE,
    labels = NULL,
    labelNudgeX = 8,
    labelSize = 3
)
```

Arguments

se The SummarizedExperiment object resulting from stability selection, by run-

ning randLassoStabSel.

selProbMin A numerical scalar in [0,1]. Predictors with a selection probability greater than

selProbMin are shown as colored lines. The color is defined by the col argu-

ment.

selColor Color for the selected predictors which have a selection probability greater than

selProbMin.

notSelColor Color for the rest of the (un-selected) predictors.

selProbCutoffColor

Color for the line depicting the selection probability cutoff.

linewidth Line width.

alpha Line transparency of the stability paths.

ylim Limits for y-axis.

The predictor labels given in labels will be shown. If unspecified, the labels corresponding to the selected predictors will be added. If predictors have the same y-value in the last regularization step, the labels will be shown in a random

order. One needs to use set. seed to reproduce the plot in this case.

labels If labelPaths = TRUE, the predictors which should be labelled. If NULL, the

selected predictors greater than metadata(se)\$stabsel.params.cutoff will

be shown.

labelNudgeX If labelPaths = TRUE, how much to nudge the labels to the right of the x-axis.

labelSize If labelPaths = TRUE, the size of the labels.

48 prepareHomer

Value

a ggplot2 object.

See Also

stabsel

Examples

```
## create data set
Y <- rnorm(n = 500, mean = 2, sd = 1)
X <- matrix(data = NA, nrow = length(Y), ncol = 50)
for (i in seq_len(ncol(X))) {
    X[ ,i] <- runif(n = 500, min = 0, max = 3)
}
s_cols <- sample(x = seq_len(ncol(X)), size = 10,
    replace = FALSE)
for (i in seq_along(s_cols)) {
    X[ ,s_cols[i]] <- X[ ,s_cols[i]] + Y
}
## reproducible randLassoStabSel() with 1 core
set.seed(123)
ss <- randLassoStabSel(x = X, y = Y)
plotStabilityPaths(ss)</pre>
```

prepareHomer

Prepare input files for HOMER motif enrichment analysis.

Description

For each bin, write genomic coordinates for foreground and background regions into files for HOMER motif enrichment analysis.

Usage

```
prepareHomer(
   gr,
   b,
   genomedir,
   outdir,
   motifFile,
   homerfile = findHomer(),
   regionsize = "given",
   Ncpu = 2L,
   verbose = FALSE
)
```

prepareHomer 49

Arguments

| gr | A GRanges object (or an object that can be coerced to one) with the genomic regions to analyze. |
|------------|--|
| b | A vector of the same length as gr that groups its elements into bins (typically a factor). |
| genomedir | Directory containing sequence files in Fasta format (one per chromosome). |
| outdir | A path specifying the folder into which the output files (two files per unique value of b) will be written. |
| motifFile | A file with HOMER formatted PWMs to be used in the enrichment analysis. |
| homerfile | Path and file name of the findMotifsGenome.pl HOMER script. |
| regionsize | The peak size to use in HOMER ("given" keeps the coordinate region, an integer value will keep only that many bases in the region center). |
| Ncpu | Number of parallel threads that HOMER can use. |
| verbose | A logical scalar. If TRUE, print progress messages. |

Details

For each bin (unique value of b) this functions creates two files in outdir/bin_N_foreground. tab and outdir/bin_N_background. tab, where N is the number of the bin and foreground/background correspond to the ranges that are/are not within the current bin). The files are in the HOMER peak file format (see http://homer.ucsd.edu/homer/ngs/peakMotifs.html for details).

In addition, a shell script file is created containing the shell commands to run the HOMER motif enrichment analysis.

Value

The path and name of the script file to run the HOMER motif enrichment analysis.

50 randLassoStabSel

randLassoStabSel

Randomized Lasso Stability Selection

Description

This function runs randomized lasso stability selection as presented by Meinshausen and Bühlmann (2010) and with the improved error bounds introduced by Shah and Samworth (2013). The function uses the stabsel function from the stabs package, but implements the randomized lasso version.

Usage

```
randLassoStabSel(
    x,
    y,
    weakness = 0.8,
    cutoff = 0.8,
    PFER = 2,
    mc.cores = 1L,
    glmnet.args = list(),
    ...
)
```

Arguments

cutoff

| X | The predictor matrix. |
|----------|--|
| у | The response vector. |
| weakness | Value between 0 and 1 (default = 0.8). It affects how strict the method will be in selecting predictors. The closer it is to 0, the more stringent the selection. A weakness value of 1 is identical to performing lasso stability selection (not the randomized version). |

Value between 0 and 1 (default = 0.8) which is the cutoff for the selection probability. Any variable with a selection probability that is higher than the set cutoff

will be selected.

PFER Integer (default = 2) representing the absolute number of false positives that we

allow for in the final list of selected variables. For details see Meinshausen and

Bühlmann (2010).

randLassoStabSel 51

mc.cores Integer (default = 1) specifying the number of cores to use in mclapply, which is the default way stabsel does parallelization.

glmnet.args Named list with additional arguments to the internal .glmnetRandomizedLasso

function (beyond x, y and weakness, which are determined automatically, and q, which should not be specified (it will be determined from cutoff and PFER). The available arguments to .glmnetRandomizedLasso are the same as the ones for glmnet.lasso. A typical use case would be to define the family argument

to glmnet.

... Additional parameters that can be passed on to stabsel.

Details

Randomized lasso stability selection runs a randomized lasso regression several times on subsamples of the response variable and predictor matrix. N/2 elements from the response variable are randomly chosen in each regression, where N is the length of the vector. The corresponding section of the predictor matrix is also chosen, and the internal .glmnetRandomizedLasso function is applied. Stability selection results in selection probabilities for each predictor. The probability of a specific predictor is the number of times it was selected divided by the total number of subsamples that were done (total number of times the regression was performed).

We made use of the stabs package that implements lasso stability selection, and adapted it to run randomized lasso stability selection.

Value

A SummarizedExperiment object where the rows are the observations and the columns the predictors (same dimnames as the predictor matrix x). It contains:

```
assays :
```

x: the predictor matrix.

rowData: a DataFrame with columns:

y: the response vector.

colData: a DataFrame with columns:

selProb: the final selection probabilities for the predictors (from the last regularization step).

selected: logical indicating the predictors that made the selection with the specified cutoff.

selAUC: the normalized area under the seletion curve (mean of selection probabilities over regulatization steps).

reg'i': columns containing the selection probabilities for regularization step i.

metadata: a list of output returned from stabsel and randLassoStabSel:

stabsel.params.cutoff: probability cutoff set for selection of predictors (see stabsel).

stabsel.params.selected : elements with maximal selection probability greater cutoff (see stabsel).

stabsel.params.max: maximum of selection probabilities (see stabsel).

stabsel.params.q: average number of selected variables used (see stabsel).

stabsel.params.PFER: (realized) upper bound for the per-family error rate (see stabsel).

52 randLassoStabSel

```
stabsel.params.specifiedPFER : specified upper bound for the per-family error rate (see stabsel).
stabsel.params.p : the number of effects subject to selection (see stabsel).
stabsel.params.B : the number of subsamples (see stabsel).
stabsel.params.sampling.type : the sampling type used for stability selection (see stabsel).
stabsel.params.assumption : the assumptions made on the selection probabilities (see stabsel).
stabsel.params.call : stabsel the call.
randStabsel.params.weakness : the weakness parameter in the randomized lasso stability selection.
```

References

N. Meinshausen and P. Bühlmann (2010), Stability Selection, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **72**, 417–73.

R.D. Shah and R.J. Samworth (2013), Variable Selection with Error Control: Another Look at Stability Selection, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **75**, 55–80.

B. Hofner, L. Boccuto, and M. Göker (2015), Controlling False Discoveries in High-Dimensional Situations: Boosting with Stability Selection, *BMC Bioinformatics*, **16** 144.

See Also

stabsel

```
## create data set
Y <- rnorm(n = 500, mean = 2, sd = 1)
X <- matrix(data = NA, nrow = length(Y), ncol = 50)</pre>
for (i in seq_len(ncol(X))) {
  X[,i] \leftarrow runif(n = 500, min = 0, max = 3)
s_{cols} < - sample(x = seq_len(ncol(X)), size = 10,
  replace = FALSE)
for (i in seq_along(s_cols)) {
  X[,s\_cols[i]] \leftarrow X[,s\_cols[i]] + Y
## reproducible randLassoStabSel() with 1 core
set.seed(123)
ss <- randLassoStabSel(x = X, y = Y)</pre>
## reproducible randLassoStabSel() in parallel mode
## (only works on non-windows machines)
if (.Platform$OS.type == "unix") {
    RNGkind("L'Ecuyer-CMRG")
    set.seed(123)
    ss <- randLassoStabSel(x = X, y = Y, mc.preschedule = TRUE,</pre>
                            mc.set.seed = TRUE, mc.cores = 2L)
}
```

sampleRandomRegions

sampleRandomRegions

Sample random regions of fixed length.

Description

Sample random regions from the mappable parts of the genome with a given fraction from CpG islands.

Usage

```
sampleRandomRegions(allowedRegions = NULL, N = 100L, regWidth = 200L)
```

Arguments

allowedRegions An unstranded GRanges object of the "allowed" of the genome, usually the mappable regions.

N Number of regions to sample.

regWidth Region width.

Details

In order to make the results deterministic, set the random number seed before calling sampleRandomRegions using set.seed.

Value

A GRanges object with randomly sampled mappable regions of width regWidth with fractionCGI coming from CpG islands.

```
regs <- GenomicRanges::GRanges(
  seqnames = rep(c("chr1", "chr2"), each = 2),
  ranges = IRanges::IRanges(start = 1:4, end = 5:8))
set.seed(123)
sampleRandomRegions(regs, N = 2, regWidth = 3L)</pre>
```

54 seqLogoGrob

|--|

Description

Create a simple sequence logo grob (grid-graphics object) for a transcription factor from a position frequency matrix. The logo drawing code is a simplified version from seqLogo and for example can be used to embedd sequence logos within other plots.

Usage

```
seqLogoGrob(x, xmax = NULL, ymax = 2, xjust = c("left", "center", "right"))
```

Arguments

| Χ | A PFMatrix object |
|-------|--|
| xmax | A numeric scalar with the maximal width for the logo (in base-pairs). A value of NULL will scale the logo to the full width of the viewport. |
| ymax | A numeric scalar with the maximal height for the logo (in bits) A value of NULL will scale the logo to the full height of the viewport. |
| xjust | A character scalar specifying the horizontal adjustment of the sequence log withint the viewport; one of "left", "center" or "right". |

Value

A polygon grob.

See Also

seqLogo for the original, more flexible version of this function.

```
if (require(JASPAR2020) && require(TFBSTools) && require(gridExtra)) {
   pfm1 <- getMatrixByID(JASPAR2020, "MA0139")
   pfm2 <- getMatrixByID(JASPAR2020, "MA0531")

   g1 <- seqLogoGrob(pfm1)
   g2 <- seqLogoGrob(pfm2)

   gridExtra::grid.arrange(g1, g2)
}</pre>
```

Index

| * internal | findHomer, 25 |
|---|--|
| .calcKmerEnrichment,4 | findMotifHits, <i>22</i> , <i>23</i> , <i>26</i> |
| .calcMotifEnrichment,5 | findMotifHits,character,character-method |
| .calculateGCweight, 6 | (findMotifHits), 26 |
| .checkDfValidity, 6 | findMotifHits,character,DNAString-method |
| .checkIfSeqsAreEqualLength, 7 | (findMotifHits), 26 |
| .cons2matrix,8 | <pre>findMotifHits,character,DNAStringSet-method</pre> |
| .defineBackground,8 | (findMotifHits), 26 |
| .filterSeqs,9 | <pre>findMotifHits,PWMatrix,character-method</pre> |
| .glmnetRandomizedLasso, 10 | (findMotifHits), 26 |
| .iterativeNormForKmers, 11 | findMotifHits,PWMatrix,DNAString-method |
| .normForKmers, 12 | (findMotifHits), 26 |
| monaLisa-package, 3 | <pre>findMotifHits,PWMatrix,DNAStringSet-method</pre> |
| .calcKmerEnrichment, 4 | (findMotifHits), 26 |
| .calcMotifEnrichment, 5 | findMotifHits,PWMatrix,GRanges-method |
| .calculateGCweight, 6 | (findMotifHits), 26 |
| .checkDfValidity, 6 | <pre>findMotifHits,PWMatrixList,character-method</pre> |
| .checkIfSeqsAreEqualLength,7 | (findMotifHits), 26 |
| .cons2matrix, 8 | <pre>findMotifHits,PWMatrixList,DNAString-method</pre> |
| .defineBackground, 8 | (findMotifHits), 26 |
| .filterSeqs, 9 | findMotifHits,PWMatrixList,DNAStringSet-method |
| .glmnetRandomizedLasso, 10 | (findMotifHits), 26 |
| .iterativeNormForKmers, 11 | <pre>findMotifHits,PWMatrixList,GRanges-method</pre> |
| .normForKmers, 12 | (findMotifHits), 26 |
| annoSeqlogo, 13 | geom_density, 38 |
| | geom_histogram,41 |
| bin, 14, 16, 17, 19, 20, 22, 23, 31, 33, 38, 40, | geom_point, 42 |
| 42, 44 | getColsByBin, 30, 38, 41, 42 |
| BiocParallelParam, 17, 22, 29, 35, 36 | getKmerFreq, 19,31 |
| bplapply, 19, 35, 36 | getMatrixSet, <u>25</u> |
| | getSetZeroBin, 32 |
| calcBinnedKmerEnr, 15 | getZeroBin(getSetZeroBin), 32 |
| calcBinnedMotifEnrHomer, 19, 43 | ggplot, <i>39</i> |
| calcBinnedMotifEnrR, $17, 21, 43$ | glmnet, <i>11</i> , <i>51</i> |
| cor, <i>46</i> | glmnet.lasso, <i>11</i> , <i>51</i> |
| cut, 15 | GRanges, 17, 22 |
| DNAStringSet, <i>16</i> , <i>22</i> , <i>32</i> , <i>39</i> | Heatmap, <i>13</i> , <i>44</i> |
| dumpJaspar, 24, 33 | HeatmapAnnotation, 14 |

56 INDEX

```
homerToPFMatrixList, 25, 33, 34, 36
mclapply, 51
monaLisa(monaLisa-package), 3
monaLisa-package, 3
motifKmerSimilarity, 34
motifSimilarity, 35
MulticoreParam, 18, 23
p.adjust, 17, 20, 22, 32, 37
parseHomerOutput, 19, 20, 37
PFMatrix, 54
PFMatrixList, 25, 33, 34, 36
plotBinDensity, 38
plotBinDiagnostics, 39
plotBinHist, 40
plotBinScatter, 41
plotMotifHeatmaps, 42
plotSelectionProb, 45
plotStabilityPaths, 46
prepareHomer, 19, 20, 48
randLassoStabSel, 45, 47, 50
sampleRandomRegions, 53
seqLogo, 54
seqLogoGrob, 13, 54
SerialParam, 18, 23
setZeroBin (getSetZeroBin), 32
stabsel, 48, 50–52
SummarizedExperiment, 18, 23, 43
system2, 19, 20
theme, 38, 41, 42
toPWM, 30
```