

# Package ‘BRGenomics’

December 29, 2019

**Type** Package

**Title** Tools for the Efficient Analysis of High-Resolution Genomics Data

**Version** 0.2.0

**Description** This package provides useful and efficient utilites for the analysis of high-resolution genomic data using standard Bioconductor methods and classes.

**License** Artistic-2.0

**Encoding** UTF-8

**LazyData** FALSE

**RoxygenNote** 7.0.2

**Depends** R (>= 4.0),  
GenomicRanges,  
GenomeInfoDb

**Imports** rtracklayer,  
parallel,  
IRanges,  
stats,  
DESeq2,  
SummarizedExperiment

**Suggests** BiocStyle,  
knitr,  
rmarkdown,  
TxDb.Hsapiens.UCSC.hg38.knownGene,  
TxDb.Hsapiens.UCSC.hg19.knownGene,  
TxDb.Mmusculus.UCSC.mm10.knownGene,  
TxDb.Mmusculus.UCSC.mm9.knownGene,  
TxDb.Dmelanogaster.UCSC.dm6.ensGene,  
TxDb.Dmelanogaster.UCSC.dm3.ensGene,  
testthat

**biocViews** Software,  
DataImport,  
RNASeq,  
ATACSeq,  
ChIPSeq,  
Transcription,  
GeneRegulation,  
GeneExpression,  
Normalization

VignetteBuilder knitr

## R topics documented:

binNdimensions . . . . .	2
genebodies . . . . .	3
getCountsByPositions . . . . .	5
getCountsByRegions . . . . .	7
getDESeqDataSet . . . . .	8
getDESeqResults . . . . .	10
getMaxPositionsBySignal . . . . .	12
getPausingIndices . . . . .	14
getStrandedCoverage . . . . .	15
import-functions . . . . .	16
makeGRangesBPRES . . . . .	18
mergeGRangesData . . . . .	19
metaSubsample . . . . .	20
metaSubsampleMatrix . . . . .	22
PROseq . . . . .	23
PROseq_paired . . . . .	24
subsampleGRanges . . . . .	25
subsetRegionsBySignal . . . . .	25
tidyChromosomes . . . . .	27
txs_dm6_chr4 . . . . .	28
<b>Index</b>	<b>29</b>

---

binNdimensions	<i>N-dimensional binning</i>
----------------	------------------------------

---

### Description

This function takes in data along 1 or more dimensions, and for each dimension the data is divided into evenly-sized bins from the minimum value to the maximum value, and bin numbers are returned. For instance, if each index of the input data were a gene, the input dimensions would be various quantitative measures of that gene, e.g. expression level, number of exons, length, etc. If plotted in cartesian coordinates, each gene would be a single datapoint, and each measurement would be a separate dimension. The bin numbers for each datapoint in each dimension are returned in a dataframe, with a column for each dimension and a row for each index.

### Usage

```
binNdimensions(..., nbins = 10)
```

### Arguments

... A single dataframe, or any number of lists or vectors containing different measurements across the same datapoints. If a dataframe is given, columns should correspond to measurements (dimensions). If lists or vectors are given, they must all have the same lengths. Other input classes will be coerced into a single dataframe.

**nbins** Either a number giving the number of bins to use for all dimensions (default = 10), or a vector containing the number of bins to use for each dimension of input data given.

### Value

A dataframe containing indices in 1:nbins for each datapoint in each dimension.

### Author(s)

Mike DeBerardine

### Examples

```
data("PROseq") # import included PROseq data
data("txs_dm6_chr4") # import included transcripts

#-----#
# find counts in promoter, early genebody, and near CPS
#-----#

pr <- promoters(txs_dm6_chr4, 0, 100)
early_gb <- genebodies(txs_dm6_chr4, 500, 1000, fix.end = "start")
cps <- genebodies(txs_dm6_chr4, -500, 500, fix.start = "end")

counts_pr <- getCountsByRegions(PROseq, pr)
counts_gb <- getCountsByRegions(PROseq, early_gb)
counts_cps <- getCountsByRegions(PROseq, cps)

#-----#
# divide genes into 20 bins for each measurement
#-----#

count_bins <- binNdimensions(counts_pr, counts_gb, counts_cps, nbins = 20)

length(txs_dm6_chr4)
nrow(count_bins)
count_bins[1:10, ]
```

---

genebodies

*Extract Genebodies*

---

### Description

This function returns ranges that are defined relative to the strand-specific start and end sites of regions of interest (usually genes). Unlike [GenomicRanges::promoters](#), distances can be upstream or downstream based on the sign, and both the start and end of the returned regions can be defined in terms of either the start or end site of the input ranges. For example, `genebodies(txs, -50, 150, fix.end = "start")` is equivalent to `promoters(txs, 50, 150)`. The default arguments return ranges that begin 300 bases downstream of the original start positions, and end 300 bases upstream of the original end positions.

**Usage**

```
genebodies(
  genelist,
  start = 300,
  end = -300,
  fix.start = "start",
  fix.end = "end",
  min.window = 0
)
```

**Arguments**

<code>genelist</code>	A GRanges object containing genes of interest.
<code>start</code>	Depending on <code>fix.start</code> , the distance from either the strand-specific start or end site to begin the returned ranges. If positive, the returned range will begin downstream of the reference position; negative numbers are used to return sites upstream of the reference. Set <code>start = 0</code> to return the reference position.
<code>end</code>	Identical to the <code>start</code> argument, but defines the strand-specific end position of returned ranges. <code>end</code> must be downstream of <code>start</code> .
<code>fix.start</code>	The reference point to use for defining the strand-specific start positions of returned ranges, either "start" or "end".
<code>fix.end</code>	The reference point to use for defining the strand-specific end positions of returned ranges, either "start" or "end". Cannot be set to "start" if <code>fix.start = "end"</code> .
<code>min.window</code>	When <code>fix.start = "start"</code> and <code>fix.end = "end"</code> , <code>min.window</code> defines the minimum size (width) of a returned range. However, when <code>fix.end = fix.start</code> , all returned ranges have the same width, and <code>min.window</code> simply size-filters the input ranges.

**Value**

A GRanges object that may be shorter than `genelist` due to loss of short ranges.

**Author(s)**

Mike DeBerardine

**See Also**

[intra-range-methods](#)

**Examples**

```
data("txs_dm6_chr4") # load included transcript data
len <- length(txs_dm6_chr4)
txs <- txs_dm6_chr4[c(1:2, len - 1, len)]
txs

#-----#
# genebodies from +300 (300 bp after TSS) to 300 bp before the poly-A site
#-----#
```

```

genebodies(txs, 300, -300)

#-----#
# promoters from -50 to +100
#-----#

promoters(txs, 50, 100)

genebodies(txs, -50, 100, fix.end = "start")

#-----#
# region from 500 to 1000 bases after the poly-A site
#-----#

genebodies(txs, 500, 1000, fix.start = "end")

```

---

getCountsByPositions    *Get signal counts at each position within regions of interest*

---

## Description

Generate a matrix containing a row for each region of interest, and columns for each position (each base if binsize = 1) within each region.

## Usage

```

getCountsByPositions(
  dataset.gr,
  regions.gr,
  binsize = 1,
  FUN = sum,
  simplify.multi.widths = c("list", "pad 0", "pad NA"),
  field = "score",
  ncores = detectCores()
)

```

## Arguments

dataset.gr	A GRanges object in which signal is contained in metadata (typically in the "score" field).
regions.gr	A GRanges object containing all the regions of interest.
binsize	Size of bins (in bp) to use for counting within each range of regions.gr. Note that counts will <i>not</i> be length-normalized.
FUN	If binsize > 1, the function used to aggregate the signal within each bin. By default, the signal is summed, but any function operating on a numeric vector can be used.
simplify.multi.widths	A string indicating the output format if the ranges in regions.gr have variable widths. Default = "list". See details below.
field	The metadata field of dataset.gr to be counted. If length(field) > 1, the output is a list whose elements contain the output for generated each field.
ncores	Multiple cores can only be used if length(field) > 1.

## Details

If the widths of all ranges in `regions.gr` are equal, a matrix is returned containing a row for each range in `regions.gr`, and a column for each bin. For input `regions.gr` with varying widths, setting `simplify.multi.widths = "list"` will output a list of variable-length vectors, with each vector corresponding to an input region. If `simplify.multi.widths = "pad 0"` or `"pad NA"`, the output is a matrix containing a row for each range in `regions.gr`, and a column for each position in each range. The number of columns is determined by the largest range in `regions.gr`, and columns corresponding to positions outside of each range are either set to 0 or NA, depending on the argument.

## Author(s)

Mike DeBerardine

## See Also

[getCountsByRegions](#)

## Examples

```
data("PROseq") # load included PROseq data
data("txs_dm6_chr4") # load included transcripts

#-----#
# counts from 0 to 50 bp after the TSS
#-----#

txs_pr <- promoters(txs_dm6_chr4, 0, 50) # first 50 bases
countsmat <- getCountsByPositions(PROseq, txs_pr)
countsmat[10:15, 40:50] # show only 40-50 bp after TSS

#-----#
# redo with 10 bp bins from 0 to 100
#-----#

txs_pr <- promoters(txs_dm6_chr4, 0, 100)
countsmat <- getCountsByPositions(PROseq, txs_pr, binsize = 10)
countsmat[10:15, ]

#-----#
# same as the above, but with the average signal in each bin
#-----#

countsmat <- getCountsByPositions(PROseq, txs_pr, binsize = 10, FUN = mean)
countsmat[10:15, ]

#-----#
# standard deviation of signal in each bin
#-----#

countsmat <- getCountsByPositions(PROseq, txs_pr, binsize = 10, FUN = sd)
round(countsmat[10:15, ], 2)
```

---

getCountsByRegions	<i>Get signal counts in regions of interest</i>
--------------------	---

---

## Description

Returns a vector the same length as `regions.gr` containing signal found in each range.

## Usage

```
getCountsByRegions(  
  dataset.gr,  
  regions.gr,  
  field = "score",  
  ncores = detectCores()  
)
```

## Arguments

<code>dataset.gr</code>	A GRanges object in which signal is contained in metadata (typically in the "score" field).
<code>regions.gr</code>	A GRanges object containing all the regions of interest.
<code>field</code>	The metadata field of <code>dataset.gr</code> to be counted. If <code>length(field) &gt; 1</code> , a dataframe is returned containing the counts for each region in each field.
<code>ncores</code>	Multiple cores can only be used if <code>length(field) &gt; 1</code> .

## Author(s)

Mike DeBerardine

## See Also

[getCountsByPositions](#)

## Examples

```
data("PROseq") # load included PROseq data  
data("txs_dm6_chr4") # load included transcripts  
  
counts <- getCountsByRegions(PROseq, txs_dm6_chr4)  
  
length(txs_dm6_chr4)  
length(counts)  
head(counts)  
  
# Assign as metadata to the transcript GRanges  
txs_dm6_chr4$PROseq <- counts  
txs_dm6_chr4[1:6]
```

---

getDESeqDataSet

*Get DESeqDataSet objects for downstream analysis*


---

## Description

This is a convenience function for generating DESeqDataSet objects, but this function also adds support for counting reads across non-contiguous regions.

## Usage

```
getDESeqDataSet(
  dataset.list,
  regions.gr,
  sample_names = names(dataset.list),
  gene_names = NULL,
  sizeFactors = NULL,
  field = "score",
  ncores = detectCores(),
  quiet = FALSE
)
```

## Arguments

dataset.list	A list of GRanges datasets that can be individually passed to <a href="#">getCountsByRegions</a> .
regions.gr	A GRanges object containing regions of interest.
sample_names	Names for each dataset in dataset.list are required, and by default the names of the list elements are used. The names must each contain the string "_rep#", where "#" is a single character (usually a number) indicating the replicate. Sample names across different replicates must be otherwise identical.
gene_names	An optional character vector giving gene names, or any other identifier over which reads should be counted. Gene names are required if counting is to be performed over non-contiguous ranges, i.e. if any genes have multiple ranges. If supplied, gene names are added to the resulting DESeqDataSet object.
sizeFactors	DESeq2 sizeFactors can be optionally applied in to the DESeqDataSet object in this function, or they can be applied later on, either by the user or in a call to getDESeqResults. Applying the sizeFactors later is useful if multiple sets of factors will be explored, although sizeFactors can be overwritten at any time.
field	Argument passed to getCountsByRegions.
ncores	Number of cores to use for read counting across all samples. Default is the total number of cores available.
quiet	If TRUE, all output messages from call to <a href="#">DESeqDataSet</a> will be suppressed.

## Value

A DESeqData object in which rowData are given as rowRanges, which are equivalent to regions.gr, unless there are non-contiguous gene regions (see note below). Samples (as seen in colData) are factored so that samples are grouped by replicate and condition, i.e. all non-replicate samples are treated as distinct, and the DESeq2 design = ~condition.



### Use of non-contiguous gene regions

In DESeq2, genes must be defined by single, contiguous chromosomal locations. This function allows individual genes to be encompassed by multiple distinct ranges in `regions.gr`. To use non-contiguous gene regions, provide `gene_names` in which some names are duplicated. For each unique gene in `gene_names`, this function will generate counts across all ranges for that gene, but be aware that it will only keep the largest range for each gene in the resulting DESeqDataSet object's `rowRanges`.

### A note on DESeq2 sizeFactors

DESeq2 `sizeFactors` are sample-specific normalization factors that are applied by division, i.e.  $counts_{norm,i} = counts_i / sizeFactor_i$ . This is in contrast to normalization factors as defined in this package (and commonly elsewhere), which are applied by multiplication. Also note that DESeq2's "normalizationFactors" are not sample specific, but rather gene specific factors used to correct for ascertainment bias across different genes (e.g. as might be relevant for GSEA or Go analysis).

### Author(s)

Mike DeBerardine

### See Also

[DESeq2::DESeqDataSet](#), [getDESeqResults](#)

### Examples

```
suppressPackageStartupMessages(require(DESeq2))
data("PROseq") # import included PROseq data
data("txs_dm6_chr4") # import included transcripts

# divide PROseq data into 6 toy datasets
ps_a_rep1 <- PROseq[seq(1, length(PROseq), 6)]
ps_b_rep1 <- PROseq[seq(2, length(PROseq), 6)]
ps_c_rep1 <- PROseq[seq(3, length(PROseq), 6)]

ps_a_rep2 <- PROseq[seq(4, length(PROseq), 6)]
ps_b_rep2 <- PROseq[seq(5, length(PROseq), 6)]
ps_c_rep2 <- PROseq[seq(6, length(PROseq), 6)]

ps_list <- list(A_rep1 = ps_a_rep1,
               A_rep2 = ps_a_rep2,
               B_rep1 = ps_b_rep1,
               B_rep2 = ps_b_rep2,
               C_rep1 = ps_c_rep1,
               C_rep2 = ps_c_rep2)

# make flawed dataset (ranges in txs_dm6_chr4 not disjoint)
# this means there is double-counting
# also using discontinuous gene regions, as gene_ids are repeated
dds <- getDESeqDataSet(ps_list,
                      txs_dm6_chr4,
                      gene_names = txs_dm6_chr4$gene_id,
                      quiet = TRUE,
                      ncores = 2)
```

dds

---

getDESeqResults

*Get DESeq2 results using reduced dispersion matrices*

---

## Description

This function calls `DESeq2::DESeq` and `DESeq2::results` on a pre-existing `DESeqDataSet` object and returns a `DESeqResults` table for one or more pairwise comparisons. However, unlike a standard call to `DESeq2::results` using the `contrast` argument, this function subsets the dataset so that `DESeq2` only estimates dispersion for the samples being compared, and not for all samples present.

## Usage

```
getDESeqResults(
  dds,
  contrast.numer,
  contrast.denom,
  comparisons.list = NULL,
  sizeFactors = NULL,
  alpha = 0.1,
  args.DESeq = NULL,
  args.results = NULL,
  ncores = detectCores(),
  quiet = FALSE
)
```

## Arguments

- |                  |  |
|------------------|--|
| dds              | A <code>DESeqDataSet</code> object, produced using either <code>getDESeqDataSet</code> from this package or <code>DESeqDataSet</code> from <code>DESeq2</code> . If dds was not created using <code>getDESeqDataSet</code> , dds must be made with <code>design = ~condition</code> such that a unique condition level exists for each sample/treatment condition. |
| contrast.numer   | A string naming the condition to use as the numerator in the <code>DESeq2</code> comparison, typically the perturbative condition.   |
| contrast.denom   | A string naming the condition to use as the denominator in the <code>DESeq2</code> comparison, typically the control condition.  |
| comparisons.list | As an optional alternative to supplying a single <code>contrast.numer</code> and <code>contrast.denom</code> , users can supply a list of character vectors containing numerator-denominator pairs, e.g. <code>list(c("B", "A"), c("C", "A"), c("C", "B"))</code> .  |
| sizeFactors      | A vector containing <code>DESeq2</code> sizeFactors to apply to each sample. Each sample's readcounts are <i>divided</i> by its respective <code>DESeq2</code> sizeFactor. A warning will be generated if the <code>DESeqDataSet</code> already contains sizeFactors, and the previous sizeFactors will be over-written.   |
| alpha            | The significance threshold passed to <code>DESeqResults</code> . This won't affect the output results, but is used as a performance optimization by <code>DESeq2</code> .  |

<code>args.DESeq</code>	Additional arguments passed to <a href="#">DESeq</a> , given as a list of argument-value pairs, e.g. <code>list(test = "LRT", fitType = "local")</code> . All arguments given here will be passed to <code>DESeq</code> except for <code>object</code> and <code>parallel</code> . If no arguments are given, all defaults will be used.
<code>args.results</code>	Additional arguments passed to <a href="#">DESeq2::results</a> , given as a list of argument-value pairs, e.g. <code>list(altHypothesis = "greater", lfcThreshold = 1.5)</code> . All arguments given here will be passed to <code>results</code> except for <code>object</code> , <code>contrast</code> , <code>alpha</code> , and <code>parallel</code> . If no arguments are given, all defaults will be used.
<code>ncores</code>	The number of cores to use for parallel processing. Multicore processing is only used if more than one comparison is being made (i.e. argument <code>comparisons.list</code> is used), and the number of cores utilized will not be greater than the number of comparisons being performed.
<code>quiet</code>	If TRUE, all output messages from calls to <code>DESeq</code> and <code>results</code> will be suppressed, although passing option <code>quiet</code> in <code>args.DESeq</code> will supersede this option for the call to <code>DESeq</code> .

### Value

For a single comparison, the output is the `DESeqResults` result table. If a `comparisons.list` is used to make multiple comparisons, the output is a named list of `DESeqResults` objects, with elements named following the pattern "X\_vs\_Y", where X is the name of the numerator condition, and Y is the name of the denominator condition.

### Author(s)

Mike DeBerardine

### See Also

[getDESeqDataSet](#), [DESeq2::results](#)

### Examples

```
#-----#
# getDESeqDataSet
#-----#
suppressPackageStartupMessages(require(DESeq2))
data("PROseq") # import included PROseq data
data("txs_dm6_chr4") # import included transcripts

# divide PROseq data into 6 toy datasets
ps_a_rep1 <- PROseq[seq(1, length(PROseq), 6)]
ps_b_rep1 <- PROseq[seq(2, length(PROseq), 6)]
ps_c_rep1 <- PROseq[seq(3, length(PROseq), 6)]

ps_a_rep2 <- PROseq[seq(4, length(PROseq), 6)]
ps_b_rep2 <- PROseq[seq(5, length(PROseq), 6)]
ps_c_rep2 <- PROseq[seq(6, length(PROseq), 6)]

ps_list <- list(A_rep1 = ps_a_rep1,
               A_rep2 = ps_a_rep2,
               B_rep1 = ps_b_rep1,
               B_rep2 = ps_b_rep2,
               C_rep1 = ps_c_rep1,
```

```

C_rep2 = ps_c_rep2)

# make flawed dataset (ranges in txs_dm6_chr4 not disjoint)
#   this means there is double-counting
# also using discontinuous gene regions, as gene_ids are repeated
dds <- getDESeqDataSet(ps_list,
                      txs_dm6_chr4,
                      gene_names = txs_dm6_chr4$gene_id,
                      ncores = 2)

dds

#-----#
# getDESeqResults
#-----#

res <- getDESeqResults(dds, "B", "A")

res

reslist <- getDESeqResults(dds,
                          comparisons.list = list(c("B", "A"), c("C", "A")),
                          ncores = 1)

names(reslist)

reslist[[1]]

```

---

getMaxPositionsBySignal

*Find sites with max signal in regions of interest*


---

## Description

For each signal-containing region of interest, find the single site with the most signal. Sites can be found at base-pair resolution, or defined for larger bins.

## Usage

```

getMaxPositionsBySignal(
  regions.gr,
  dataset.gr,
  binsize = 1,
  bin.centers = FALSE,
  field = "score",
  keep.score = FALSE
)

```

## Arguments

regions.gr	A GRanges object containing regions of interest.
dataset.gr	A GRanges object in which signal is contained in metadata (typically in the "score" field).
binsize	The size of bin in which to calculate signal scores.

bin.centers	Logical indicating if the centers of bins are returned, as opposed to the entire bin. If TRUE,
field	The metadata field of dataset.gr to be counted.
keep.score	Logical indicating if the signal value at the max site should be reported. If set to TRUE, the values are kept as a new metadata column in regions.gr.

### Value

Output is a GRanges object with regions.gr metadata, but each range only contains the site within each regions.gr range that had the most signal. If binsize > 1, the entire bin is returned, unless bin.centers = TRUE, in which case a single-base site is returned. The site is set to the center of the bin, and if the binsize is even, the site is rounded to be closer to the beginning of the range.

If keep.score = TRUE, the output will also contain metadata for the signal at the max site. The output is *not* necessarily same length as regions.gr, as regions without signal are not returned. If *no regions* have signal (e.g. as could happen if running this function on a single region), the function will return an empty GRanges object with intact metadata columns.

### Author(s)

Mike DeBerardine

### See Also

[getCountsByPositions](#)

### Examples

```
data("PROseq") # load included PROseq data
data("txs_dm6_chr4") # load included transcripts

#-----#
# first 50 bases of transcripts
#-----#

pr <- promoters(txs_dm6_chr4, 0, 50)
pr[1:3]

#-----#
# max sites
#-----#

getMaxPositionsBySignal(pr[1:3], PROseq, keep.score = TRUE)

#-----#
# max sites in 5 bp bins
#-----#

getMaxPositionsBySignal(pr[1:3], PROseq, binsize = 5, keep.score = TRUE)
```

---

getPausingIndices	<i>Calculate pausing indices from user-supplied promoters &amp; genebodies</i>
-------------------	--

---

### Description

Pausing index (PI) is calculated for each gene (within matched promoters.gr and genebodies.gr) as promoter-proximal (or pause region) signal counts divided by genebody signal counts. If length.normalize = TRUE (recommended), the signal counts within each range in promoters.gr and genebodies.gr are divided by their respective range widths (region lengths) before pausing indices are calculated.

### Usage

```
getPausingIndices(
  dataset.gr,
  promoters.gr,
  genebodies.gr,
  field = "score",
  length.normalize = TRUE,
  remove.empty = FALSE,
  ncores = detectCores()
)
```

### Arguments

dataset.gr	A GRanges object in which signal is contained in metadata (typically in the "score" field).
promoters.gr	A GRanges object containing promoter-proximal regions of interest.
genebodies.gr	A GRanges object containing genebody regions of interest.
field	The metadata field of dataset.gr to be counted. If length(field) > 1, a dataframe is returned containing the pausing indices for each region in each field.
length.normalize	A logical indicating if signal counts within regions of interest should be length normalized. The default is TRUE, which is recommended, especially if input regions don't all have the same width.
remove.empty	A logical indicating if genes without any signal in promoters.gr should be removed. No genes are filtered by default.
ncores	Multiple cores can only be used if length(field) > 1.

### Value

A vector of length given by the length of the genelist (or possibly shorter if remove.empty = TRUE). If length(field) > 1, a dataframe is returned, containing a column for each field.

### Author(s)

Mike DeBerardine

### See Also

[getCountsByRegions](#)

**Examples**

```

data("PROseq") # load included PROseq data
data("txs_dm6_chr4") # load included transcripts

#-----#
# Get promoter-proximal and genebody regions
#-----#

# genebodies from +300 to 300 bp before the poly-A site
gb <- genebodies(txs_dm6_chr4, 300, -300, min.window = 400)

# get the transcripts that are large enough (>1kb in size)
txs <- subset(txs_dm6_chr4, tx_name %in% gb$tx_name)

# for the same transcripts, promoter-proximal region from 0 to +100
pr <- promoters(txs, 0, 100)

#-----#
# Calculate pausing indices
#-----#

pidx <- getPausingIndices(PROseq, pr, gb)

length(txs_dm6_chr4)
length(pidx)
head(pidx)

#-----#
# Without length normalization
#-----#

head( getPausingIndices(PROseq, pr, gb, length.normalize = FALSE) )

#-----#
# Removing empty means the values no longer match the genelist
#-----#

pidx_signal <- getPausingIndices(PROseq, pr, gb, remove.empty = TRUE)

length(pidx_signal)

```

---

getStrandedCoverage	<i>Get strand-specific coverage</i>
---------------------	-------------------------------------

---

**Description**

Computes strand-specific coverage signal, and returns a GRanges object with signal in the "score" metadata column. Function also works for non-strand-specific data. Note that output is not automatically converted into a "basepair-resolution" GRanges object.

**Usage**

```
getStrandedCoverage(dataset.gr, field = "score")
```

**Arguments**

<code>dataset.gr</code>	A GRanges object either containing ranges for each read, or one in which readcounts for individual ranges are contained in metadata (typically in the "score" field).
<code>field</code>	The name of the metadata field that contains readcounts. If no metadata field contains readcounts, and each range represents a single read, set to NULL.

**Author(s)**

Mike DeBerardine

**See Also**

[makeGRangesBPres](#)

**Examples**

```
#-----#
# Using included full-read data
#-----#

data("PROseq_paired")

PROseq_paired[1:6]

getStrandedCoverage(PROseq_paired)[1:6]

#-----#
# Re-creating score for included single-base data
#-----#

data("PROseq")

PROseq[1:6]

# undo coverage for the first 100 positions
ps <- PROseq[1:100]

ps_reads <- rep(ps, times = ps$score)
mcols(ps_reads) <- NULL
ps_reads[1:6]

# re-create coverage
getStrandedCoverage(ps_reads, field = NULL)[1:6]
```

---

import-functions

---

*Import basepair-resolution files*


---

**Description**

Import basepair-resolution files



## Usage

```
import_bigWig(  
  plus_file,  
  minus_file,  
  genome = NULL,  
  keep.X = TRUE,  
  keep.Y = TRUE,  
  keep.M = FALSE,  
  keep.nonstandard = FALSE  
)  
  
import_bedGraph(  
  plus_file,  
  minus_file,  
  genome = NULL,  
  keep.X = TRUE,  
  keep.Y = TRUE,  
  keep.M = FALSE,  
  keep.nonstandard = FALSE  
)
```

## Arguments

plus_file, minus_file	Paths for strand-specific input files.
genome	Optional string for UCSC reference genome, e.g. "hg38". If given, non-standard chromosomes are trimmed, and options for sex and mitochondrial chromosomes are applied.
keep.X, keep.Y, keep.M, keep.nonstandard	Logicals indicating which non-autosomes should be kept. By default, sex chromosomes are kept, but mitochondrial and non-standard chromosomes are removed.

## Details

Imports a GRanges object containing base-pair resolution data, with the score metadata column indicating the number of reads represented by each range.

import\_bedGraph is useful for when both 5'- and 3'-end information is to be maintained for each sequenced molecule. It effectively imports the entire read.

For import\_bigWig, all ranges are of width = 1.

## Author(s)

Mike DeBerardine

## See Also

[tidyChromosomes](#), [rtracklayer::import](#)

**Examples**

```
# get local address for included bigWig files
p.bw <- system.file("extdata", "PROseq_dm6_chr4_plus.bw", package = "BRGenomics")
m.bw <- system.file("extdata", "PROseq_dm6_chr4_minus.bw", package = "BRGenomics")
# import bigWigs
PROseq <- import_bigWig(p.bw, m.bw, genome = "dm6")
```

---

makeGRangesBPres

*Make base-pair resolution GRanges object*


---

**Description**

Splits up all ranges in `gr` to be each 1 basepair wide. All information is preserved, including all metadata. To wit, `length(output.gr) = sum(width(dataset.gr))`.

**Usage**

```
makeGRangesBPres(dataset.gr)
```

**Arguments**

`dataset.gr`      A disjoint `GRanges` object

**Details**

Note that this function doesn't perform any transformation on the metadata in the input; for any ranges of `width > 1`, the metadata is simply copied to the daughters of that range (whose widths are all equal to 1).

This function is intended to work on datasets at single-base resolution. Data of this type is often formatted as a bigWig file, and any data imported from a bigWig file by `rtracklayer` is suitable for processing. bigWig files will typically use run-length compression on the data signal (the 'score' column), such that when imported by `rtracklayer`, adjacent bases sharing the same signal will combined into a single range. The base-pair resolution `GRanges` objects produced by this function remove this compression, resulting in each index (each range) of the `GRanges` object addressing a single genomic position.

To properly use base-pair resolution information, the user should be selecting a single-base from each read, which can be accomplished using [GenomicRanges::resize\(\)](#). Then, single-base coverage can be calculated using [getStrandedCoverage](#).

**Author(s)**

Mike DeBerardine

**See Also**

[getStrandedCoverage](#), [GenomicRanges::resize\(\)](#)

**Examples**

```
data("PROseq") # load included PROseq data
range(width(PROseq))

# simulate the format of a bigWig file, using arbitrary scores
bw <- reduce(PROseq)
score(bw) <- score(PROseq)[seq_along(bw)]
range(width(bw))
length(bw)

gr <- makeGRangesBPRES(bw)
range(width(gr))
length(gr)
```

---

mergeGRangesData	<i>Merge base-pair resolution GRanges objects</i>
------------------	---

---

**Description**

Merges 2 or more GRanges objects. For each object, the range widths must all be 1, and the score metadata column contains coverage information at each site. This function returns a single GRanges object containing all sites of the input objects, and the sum of all scores at all sites.

**Usage**

```
mergeGRangesData(..., field = "score", ncores = detectCores())
```

**Arguments**

...	Any number of GRanges objects in which signal (e.g. readcounts) are contained within metadata.
field	One or more metadata fields to be combined, typically the "score" field. Fields typically contain coverage information.
ncores	More than one core can be used to coerce non-single-width GRanges objects using makeGRangesBPRES.

**Author(s)**

Mike DeBerardine

**See Also**

[makeGRangesBPRES](#)

**Examples**

```
data("PROseq") # load included PROseq data

#-----#
# divide PROseq data into thirds
#-----#
```

```

thirds <- floor( (1:3)/3 * length(PROseq) )
ps_1 <- PROseq[1:thirds[1]]
ps_2 <- PROseq[(thirds[1]+1):thirds[2]]
ps_3 <- PROseq[(thirds[2]+1):thirds[3]]

#-----#
# re-merge PROseq data
#-----#

length(PROseq)
length(mergeGRangesData(ps_1, ps_2))
length(mergeGRangesData(ps_1, ps_2, ps_3))

```

---

metaSubsample

*Iterative Subsampling for Metaplotting*


---

## Description

This function performs bootstrap subsampling of mean readcounts at different positions within regions of interest. Mean signal counts can be estimated at base-pair resolution, or smoothed over larger bins.

## Usage

```

metaSubsample(
  dataset.gr,
  regions.gr,
  binsize = 1,
  first.output.xval = 1,
  sample.name = deparse(substitute(dataset.gr)),
  n.iter = 1000,
  prop.sample = 0.1,
  lower = 0.125,
  upper = 0.875,
  NF = 1,
  field = "score",
  remove.empty = FALSE,
  ncores = 1
)

```

## Arguments

dataset.gr	A GRanges object in which signal is contained in metadata (typically in the "score" field).
regions.gr	A GRanges object containing intervals over which to metaplot. All ranges must have the same width.
binsize	The size of bin (number of columns, e.g. basepairs) to use for metaplotting. Especially important for metaplots over large/sparse regions.
first.output.xval	The relative start position of the first bin, e.g. if regions.gr begins at 50 bases upstream of the TSS, set first.output.xval = -50. This number only affects the x-values that are returned, which are provided as a convenience.

sample.name	Defaults to the name of dataset.gr. This is included in the output as a convenience for row-binding outputs from different samples.
n.iter	Number of random subsampling iterations to perform. Default is 1000.
prop.sample	The proportion of the ranges in regions.gr (e.g. the proportion of genes) to subsample in each iteration. The default is 0.1 (10 percent).
lower	The lower quantile of subsampled signal means to return. The default is 0.125 (12.5th percentile).
upper	The upper quantile of subsampled signal means to return. The default is 0.875 (87.5th percentile).
NF	Optional normalization factor by which to multiply the counts.
field	The metadata field of dataset.gr to be counted.
remove.empty	A logical indicating whether regions without signal should be removed from the analysis.
ncores	Number of cores to use for parallel computation. No parallel processing is used by default, as there's no performance benefit for typical usage with short computation times.

**Value**

Dataframe containing x-values, means, lower quantiles, upper quantiles, and the sample name (as a convenience for row-binding multiple of these dataframes).

**Author(s)**

Mike DeBerardine

**See Also**

[metaSubsampleMatrix](#), [getCountsByPositions](#)

**Examples**

```
data("PROseq") # import included PROseq data
data("txs_dm6_chr4") # import included transcripts

# for each transcript, use promoter-proximal region from TSS to +100
pr <- promoters(txs_dm6_chr4, 0, 100)

#-----#
# Bootstrap average signal in each 5 bp bin across all transcripts,
# and get confidence bands for middle 30% of bootstrapped means
#-----#

set.seed(11)
df <- metaSubsample(PROseq, pr, binsize = 5, lower = 0.35, upper = 0.65)
df[1:10, ]

#-----#
# Plot bootstrapped means with confidence intervals
#-----#

plot(mean ~ x, df, type = "l", main = "PROseq Signal",
      ylab = "Mean + 30% CI", xlab = "Distance from TSS")
```

```

polygon(c(df$x, rev(df$x)), c(df$lower, rev(df$upper)),
        col = adjustcolor("black", 0.1), border = FALSE)

```

---

metaSubsampleMatrix      *Iterative Subsampling for Metaplotting (On Count Matrices)*

---

## Description

In the most general sense, this function performs iterations of randomly subsampling rows of a matrix, and returns a summary of mean values calculated for each column. The typical application is for generating metaplots, with the typical input being a matrix in which each row is a gene or other region of interest, each column is a position within that gene (either a specific basepair or a bin), and element values are signal (e.g. read counts) within those positions.

## Usage

```

metaSubsampleMatrix(
  counts.mat,
  binsize = 1,
  first.output.xval = 1,
  sample.name = deparse(substitute(counts.mat)),
  n.iter = 1000,
  prop.sample = 0.1,
  lower = 0.125,
  upper = 0.875,
  NF = 1,
  ncores = 1
)

```

## Arguments

counts.mat	A matrix of signal counts in which rows are regions of interest and columns are sites/bins in each region.
binsize	The size of bin (number of columns, e.g. basepairs) to use for metaplotting. Especially important for metaplots over large/sparse regions.
first.output.xval	The relative start position of the first bin, e.g. if regions.gr begins at 50 bases upstream of the TSS, set first.output.xval = -50. This number only affects the x-values that are returned, which are provided as a convenience.
sample.name	Defaults to the name of dataset.gr.
n.iter	Number of random subsampling iterations to perform. Default is 1000.
prop.sample	The proportion of rows to subsample in each iteration. The default is 0.1.
lower	The lower quantile of subsampled signal means to return. The default is 0.125 (12.5th percentile).
upper	The upper quantile of subsampled signal means to return. The default is 0.875 (85.5th percentile).
NF	Optional normalization factor by which to multiply the counts.
ncores	Number of cores to use for parallel computation. As of writing, parallel processing doesn't show any benefit for short computation times (e.g. <1 minute for our typical experience on a laptop).

**Value**

Dataframe containing x-values, means, lower quantiles, upper quantiles, and the sample name (as a convenience for row-binding multiple of these dataframes).

**Author(s)**

Mike DeBerardine

**See Also**

[metaSubsample](#), [getCountsByPositions](#)

**Examples**

```
data("PROseq") # load included PROseq data
data("txs_dm6_chr4") # load included transcripts

# for each transcript, use promoter-proximal region from TSS to +100
pr <- promoters(txs_dm6_chr4, 0, 100)

# generate a matrix of counts in each region
countsmat <- getCountsByPositions(PROseq, pr)
dim(countsmat)

#-----#
# bootstrap average signal in 10 bp bins across all transcripts
#-----#

set.seed(11)
df <- metaSubsampleMatrix(countsmat, binsize = 10, sample.name = "PROseq")
df[1:10, ]

#-----#
# the same, using a normalization factor, and changing the x-values
#-----#

set.seed(11)
df <- metaSubsampleMatrix(countsmat, binsize = 10, first.output.xval = 0,
                           NF = 0.75, sample.name = "PROseq")
df[1:10, ]
```

---

PROseq

*PRO-seq data from Drosophila S2 cells*

---

**Description**

PRO-seq data of Drosophila S2 cells, chromosome 4.

**Usage**

PROseq

**Format**

A disjoint GRanges object with 47533 ranges with 1 metadata column:

**score** coverage of PRO-seq read 3'-ends ...

**Details**

Hojoong Kwak, Nicholas J. Fuda, Leighton J. Core, John T. Lis (2013). Precise Maps of RNA Polymerase Reveal How Promoters Direct Initiation and Pausing. *Science*, 339(6122), 950–953. <https://doi.org/10.1126/science.1229386>

**Source**

GEO Accession GSM1032758, run SRR611828.

---

PROseq_paired	<i>Paired PRO-seq data from Drosophila S2 cells</i>
---------------	---

---

**Description**

PRO-seq data of Drosophila S2 cells, chromosome 4. Entire mapped reads kept.

**Usage**

PROseq\_paired

**Format**

A GRanges object with 52464 ranges with 1 metadata column:

**score** number of reads sharing the same mapped 5' and 3' ends ...

**Details**

Hojoong Kwak, Nicholas J. Fuda, Leighton J. Core, John T. Lis (2013). Precise Maps of RNA Polymerase Reveal How Promoters Direct Initiation and Pausing. *Science*, 339(6122), 950–953. <https://doi.org/10.1126/science.1229386>

**Source**

GEO Accession GSM1032758, run SRR611828.



---

subsampleGRanges	<i>Randomly subsample reads from GRanges dataset</i>
------------------	--

---

### Description

Random subsampling is not performed on ranges, but on reads. Readcounts should be given as a metadata field (usually "score"), and should normally be integers. If normalized readcounts are given, an attempt will be made to infer the normalization factor based on the least-common-multiple of the signal found in the specified field. This function can also subsample ranges directly if `field = NULL`, but the `sample` function can be used in this scenario.

### Usage

```
subsampleGRanges(dataset.gr, n = NULL, prop = NULL, field = "score")
```

### Arguments

<code>dataset.gr</code>	A GRanges object in which signal (e.g. readcounts) are contained within meta-data.
<code>n</code>	Number of reads to subsample. Either <code>n</code> or <code>prop</code> can be given.
<code>prop</code>	Proportion of total signal to subsample.
<code>field</code>	The metadata field of <code>dataset.gr</code> that contains readcounts for reach position. If each range represents a single read, set <code>field = NULL</code>

### Author(s)

Mike DeBerardine

### Examples

```
data("PROseq") # load included PROseq data

length(PROseq)
sum(score(PROseq))

# sample 10% of the reads
ps_sample <- subsampleGRanges(PROseq, prop = 0.1)

length(ps_sample)
sum(score(ps_sample)) # 1/10th the score is sampled
```

---

subsetRegionsBySignal	<i>Subset regions of interest by quantiles of overlapping signal</i>
-----------------------	--

---

### Description

A convenience function to subset regions of interest by the amount of signal they contain, according to their quantile (i.e. their signal ranks).

**Usage**

```
subsetRegionsBySignal(
  regions.gr,
  dataset.gr,
  quantiles = c(0.5, 1),
  field = "score",
  order.by.rank = FALSE,
  density = FALSE
)
```

**Arguments**

<code>regions.gr</code>	A GRanges object containing regions of interest.
<code>dataset.gr</code>	A GRanges object in which signal is contained in metadata (typically in the "score" field).
<code>quantiles</code>	A value pair giving the lower quantile and upper quantile of regions to keep. Regions with signal quantiles below than the lower quantile are removed, while regions with signal quantiles above the upper quantile are removed. Quantiles must be in range (0,1). An empty GRanges object is returned if lower quantile = 1 or upper quantile = 0.
<code>field</code>	The metadata field of <code>dataset.gr</code> to be counted.
<code>order.by.rank</code>	If TRUE, the output regions are sorted based on the amount of signal contained (in decreasing order). If FALSE (the default), genes are sorted by their positions.
<code>density</code>	A logical indicating whether signal counts should be normalized to the width of ranges in <code>regions.gr</code> . By default, the function only considers the total signal in each range.

**Details**

Typical uses may include removing the 5 signal (`lower_quantile = 0.05`) and the 5 (`upper_quantile = 0.95`), or returning the middle 50 signal (`lower_quantile = 0.25`, `upper_quantile = 0.75`). If `lower_quantile = 0` and `upper_quantile = 1`, all regions are returned, but the returned regions will be sorted by position, or by score if `order.by.rank = TRUE`.

**Value**

A GRanges object of length `length(regions.gr) * (upper_quantile - lower_quantile)`.

**Author(s)**

Mike DeBerardine

**See Also**

[getCountsByRegions](#)

**Examples**

```
data("PROseq") # load included PROseq data
data("txs_dm6_chr4") # load included transcripts

txs_dm6_chr4
```

```

#-----#
# get the top 50% of transcripts by signal
#-----#

subsetRegionsBySignal(txs_dm6_chr4, PROseq)

#-----#
# get the middle 50% of transcripts by signal
#-----#

subsetRegionsBySignal(txs_dm6_chr4, PROseq, quantiles = c(0.25, 0.75))

#-----#
# get the top 10% of transcripts by signal, and sort them by highest signal
#-----#

subsetRegionsBySignal(txs_dm6_chr4, PROseq,
                      quantiles = c(0.9, 1),
                      order.by.rank = TRUE)

```

---

tidyChromosomes

---

*Remove odd chromosomes from GRanges objects*


---

## Description

This convenience function removes non-standard, mitochondrial, and/or sex chromosomes from any GRanges object. For the chromosomes being removed, any ranges found on those chromosomes are removed, and the chromosomes are also removed from seqinfo. Standard chromosomes are defined using the [standardChromosomes](#) function from the GenomeInfoDb package.

## Usage

```

tidyChromosomes(
  gr,
  keep.X = TRUE,
  keep.Y = TRUE,
  keep.M = FALSE,
  keep.nonstandard = FALSE
)

```

## Arguments

gr	Any GRanges object, however the object should have a standard genome set, e.g. <code>genome(gr) &lt;- "hg38"</code>
keep.X, keep.Y, keep.M, keep.nonstandard	Logicals indicating which non-autosomes should be kept. By default, sex chromosomes are kept, but mitochondrial and non-standard chromosomes are removed.

## Author(s)

Mike DeBerardine

**See Also**[GenomeInfoDb::standardChromosomes](#)**Examples**

```
data("PROseq") # load included PROseq data
# (only data on chr4, so nothing actually changes)
PROseq_tidy <- tidyChromosomes(PROseq)
```

---

txs_dm6_chr4	<i>Ensembl transcripts for Drosophila melanogaster, dm6, chromosome 4.</i>
--------------	--

---

**Description**

Transcripts obtained from annotation package TxDb.Dmelanogaster.UCSC.dm6.ensGene, which was in turn made by the Bioconductor Core Team from UCSC resources on 2019-04-25. Metadata columns were obtained from "TXNAME" and "GENEID" columns. Data exported from the TxDb package using GenomicFeatures version 1.35.11 on 2019-12-19.

**Usage**

```
txs_dm6_chr4
```

**Format**

A GRanges object with 339 ranges and 2 metadata columns:

**tx\_name** Flybase unique identifiers for transcripts

**gene\_id** Flybase unique identifiers for the associated genes

**Source**

TxDb.Dmelanogaster.UCSC.dm6.ensGene version 3.4.6

# Index

## \* datasets

PROseq, [23](#)  
PROseq\_paired, [24](#)  
txs\_dm6\_chr4, [28](#)

GenomeInfoDb::standardChromosomes, [28](#)

binNdimensions, [2](#)

DESeq, [11](#)  
DESeq2::DESeq, [10](#)  
DESeq2::DESeqDataSet, [9](#)  
DESeq2::results, [10](#), [11](#)  
DESeqDataSet, [8](#), [10](#)

genebodies, [3](#)  
GenomicRanges::promoters, [3](#)  
GenomicRanges::resize(), [18](#)  
getCountsByPositions, [5](#), [7](#), [13](#), [21](#), [23](#)  
getCountsByRegions, [6](#), [7](#), [8](#), [14](#), [26](#)  
getDESeqDataSet, [8](#), [10](#), [11](#)  
getDESeqResults, [9](#), [10](#)  
getMaxPositionsBySignal, [12](#)  
getPausingIndices, [14](#)  
getStrandedCoverage, [15](#), [18](#)

import-functions, [16](#)  
import\_bedGraph(import-functions), [16](#)  
import\_bigWig(import-functions), [16](#)

makeGRangesBPRES, [16](#), [18](#), [19](#)  
mergeGRangesData, [19](#)  
metaSubsample, [20](#), [23](#)  
metaSubsampleMatrix, [21](#), [22](#)

PROseq, [23](#)  
PROseq\_paired, [24](#)

rtracklayer::import, [17](#)

standardChromosomes, [27](#)  
subsampleGRanges, [25](#)  
subsetRegionsBySignal, [25](#)

tidyChromosomes, [17](#), [27](#)  
txs\_dm6\_chr4, [28](#)