

# Package ‘HMMcopy’

September 6, 2011

**Type** Package

**Title** Copy number prediction with correction for GC and mappability bias for HTS data

**Version** 0.1.0

**Date** 2011-09-06

**Author** Daniel Lai

**Maintainer** Daniel Lai <jujubix@cs.ubc.ca>

**Import** IRanges

**Depends** R (>= 2.10.0), IRanges (>= 1.4.16), geneplotter (>= 1.24.0)

**Description** Corrects GC and mappability biases for readcounts (i.e. coverage) in non-overlapping windows of fixed length for single whole genome samples, yielding a rough estimate of copy number for further analysis. Designed for rapid correction of high coverage whole genome tumour and normal samples.

**License** GPL-3

**biocViews** Sequencing, Preprocessing, Visualization, CopyNumberVariants, HighThroughputSequencing, Microarray, Classification08-2508-25

## R topics documented:

HMMcopy-package . . . . .	2
correctReadcount . . . . .	3
HMMcopy example dataset . . . . .	4
HMMcopy Segmentation . . . . .	4
HMMcopy Visualization . . . . .	7
WIG Import Functions . . . . .	8
WIG Output Functions . . . . .	9
wigsToRangedData . . . . .	10
<b>Index</b>	<b>12</b>

---

HMMcopy-package      *Bias-free copy number estimation and robust CNA detection in tumour samples from WGS HTS data*

---

## Description

HMMcopy is a package for making bias-free copy number estimation by correcting for GC-content and mappability bias in HTS readcounts. It also contains an implementation of the Hidden Markov Model to robustly segment a copy number profile into non-overlapping segments predicted to be of the same copy number state, and attributes a biological copy number aberration events to the segments.

## Details

HMMcopy takes as input WIG format files generated by fast C++ tools distributed as part of the *HMMcopy Suite*, namely readcount, GC-content and mappability values for non-overlapping fixed width “bins” across the reference genome of interest. It then uses a filtering and LOESS model to correct the GC-content and mappability biases observed in the readcounts (Benjamini and Speed, 2011), and uses the corrected readcounts as a proxy of copy number. The resultant copy number profile is then segmented with a six state Hidden Markov Model, with a handful of quick visualization functions for quick viewing.

Package: HMMcopy  
Type: Package  
Version: 0.1.0  
Date: 2011-09-06  
License: GPL-3

`example("HMMcopy-package")` for quick tour of functionality and visualization  
`vignette("HMMcopy")` for detailed example

## Author(s)

Daniel Lai  
Maintainer: Daniel Lai <jujubix@cs.ubc.ca>

## References

Y. Benjamini and T.P. Speed (2011) Estimation and correction for GC-content bias in high throughput sequencing. UC Berkley Department of Statistics Technical Reports, Report ID 804.

## Examples

```
# Read WIG file input
rfile <- system.file("extdata", "tumour.wig", package = "HMMcopy")
gfile <- system.file("extdata", "gc.wig", package = "HMMcopy")
mfile <- system.file("extdata", "map.wig", package = "HMMcopy")
uncorrected_reads <- wigsToRangedData(rfile, gfile, mfile)
```

```
# Correct reads into copy number
corrected_copy <- correctReadcount(uncorrected_reads)

# Segment copy number profile
segmented_copy <- HMMsegment(corrected_copy)

# Visualize one at a time
par(ask = TRUE)
plotBias(corrected_copy)
plotCorrection(corrected_copy)
plotSegments(corrected_copy, segmented_copy)
```

---

correctReadcount     *Readcount correction for GC and mappability bias*

---

### Description

Corrects readcounts for GC and mappability bias using the binning/loess method optimized for speed.

### Usage

```
correctReadcount(x, mappability = 0.9, samplesize = 50000, verbose = TRUE)
```

### Arguments

x	<a href="#">RangedData</a> object returned by <a href="#">wigsToRangedData</a>
mappability	Mappability threshold [0, 1] below which points are ignored during creating the correction curve.
samplesize	The number of points sampled during LOESS fitting, decreasing reduces runtime and memory usage, at the expense of robustness to data randomness.
verbose	Set to FALSE if messages are not desired.

### Value

The original [RangedData](#) object appended with 4 new columns:

**valid** Valid bins, which have valid read, gc, and mappability values

**ideal** Ideal bins of high mappability and no outliers

**cor.gc** GC-corrected readcounts

**cor.map** Mappability corrected GC-corrected readcounts

**copy** cor.map values after log base 2

### Author(s)

Daniel Lai

### References

Y. Benjamini and T.P. Speed (2011) Estimation and correction for GC-content bias in high throughput sequencing. UC Berkeley Department of Statistics Technical Reports, Report ID 804.

**See Also**

[wigsToRangedData](#) to easily generate the proper input.

**Examples**

```
data(tumour) # Load tumour_reads
tumour_copy <- correctReadcount(tumour_reads)
```

---

HMMcopy example dataset

*HMMcopy example dataset*

---

**Description**

A set of data of chromosome 6 of matched tumour normal pair.

**tumour\_reads** The number of short reads in fixed width windows across the chromosome, generated with [wigsToRangedData](#) from WIG files

**tumour\_copy** Tumour copy number profile generated by correcting ‘tumour\_reads’ with [correctReadcount](#)

**normal\_copy** Normal copy number profile generated via [correctReadcount](#)

**tumour\_param** Parameters for segmenting ‘tumour\_copy’ in [HMMsegment](#)

**tumour\_segments** Segmented output of ‘tumour\_copy’ [HMMsegment](#) using ‘tumour\_param’

**Usage**

```
data(tumour)
```

**Format**

‘tumour\_reads’, ‘tumour\_copy’, and ‘normal\_copy’ are [RangedData](#) objects.

‘tumour\_param’ is a numeric matrix.

‘tumour\_segments’ is a list.

---

HMMcopy Segmentation

*Segmentation and classification of copy number profiles*

---

**Description**

Takes in a copy number profile and segments it into predicted regions of equal copy number, and assigns a biologically motivated copy number state to each region using a Hidden Markov Model.

**Usage**

```
HMMsegment(correctOut, param = NULL, autosomes = NULL,
            maxiter = 50, getparam = FALSE, verbose = TRUE)
```

**Arguments**

correctOut	Output value from <code>correctReadcount</code>
param	<p>If none is provided, will generate a reasonable set of parameters based on the input data, which can optionally be returned for inspection and manual adjustment by setting 'getparam' to TRUE.</p> <p>See Details for more information on parameters.</p> <p>A matrix with parameters values in columns for each state in rows:</p> <p><b>e</b> Probability of extending a segment, increase to lengthen segments, decrease to shorten segments. Range: (0, 1)</p> <p><b>strength</b> Strength of initial e suggestion, reducing allows e to change, increasing makes e undefiable. Range: [0, Inf)</p> <p><b>mu</b> Suggested median for copy numbers in state, change to readjust classification of states. Range: (-Inf, Inf)</p> <p><b>lambda</b> Suggested precision (inversed variance) for copy numbers in state, increase to reduce overlap between states. Range: [0, Inf)</p> <p><b>nu</b> Suggested degree of freedom between states, increase to reduce overlap between states. Range: [0, Inf)</p> <p><b>kappa</b> Suggested distribution of states. Should sum to 1.</p> <p><b>m</b> Optimal value for mu, difference from corresponding mu value determines elasticity of the mu value. <i>i.e.</i> Set to identical value as mu if you don't want mu to move much.</p> <p><b>eta</b> Mobility of mu, increase to allow more movement. Range: [0, Inf)</p> <p><b>gamma</b> Prior shape on lambda, gamma distribution. Effects flexibility of lambda.</p> <p><b>S</b> Prior scale on lambda, gamma distribution. Effects flexibility of lambda.</p>
autosomes	Array of LOGICAL values corresponding to the 'chr' argument where an element is TRUE if the chromosome is an autosome, otherwise FALSE. If not provided, will automatically set the following chromosomes to false: "X", "Y", "23", "24", "chrX", chrY", "M", "MT", "chrM".
maxiter	The maximum number of iterations allows for the Maximum-Expectation algorithm, reduce to decrease running time at the expense of robustness.
getparam	If TRUE, generates and returns parameters without running segmentation.
verbose	Set to FALSE if messages are not desired

**Details**

`HMMsegment` is a two stage algorithm that first runs an Expectation-Maximization algorithm to find the optimal set of parameters based on suggested parameter inputs, and allowed flexibilities. After iteratively finding the optimal parameters, the actual segmentation of the data is conducted with the Viterbi algorithm, finally output segmented states.

Parameters are divided into two main categories:

- Initial parameters: e, mu, lambda, nu, kappa
- Flexibility parameters: strength, m, eta, gamma, S

Where *initial parameters* are treated as starting suggestions for the parameter optimization algorithm, and flexibility parameters (hyperparameters) define how much the initial parameters are allowed to deviate during the search for the optimal parameters.



## Examples

```
data(tumour) # Load tumour_copy
tumour_segments <- HMMsegment(tumour_copy)
```

---

HMMcopy Visualization  
*Visualization functions for correctReadcount results*

---

## Description

Convenience functions for creating plots for the analysis of the readcount correction process by [correctReadcount](#)

## Usage

```
plotBias(correctOutput, points = 10000, ...)
plotCorrection(correctOutput, chr = space(correctOutput)[1], ...)
plotSegments(correctOutput, segmentOutput, chr = space(correctOutput)[1], ...)
plotParam(segmentOutput, param, ...)
stateCols()
```

## Arguments

correctOutput	Output value from <a href="#">correctReadcount</a>
segmentOutput	Output value from <a href="#">HMMsegment</a>
points	Number of random sampled points to plot, decreasing reduces runtime
chr	Chromosome name to plot. A single value for <a href="#">plotCorrection</a> and a vector for <a href="#">plotSegments</a> .
param	Input parameters to call that produced segmentOutput
...	Further arguments are passed to <a href="#">plot</a> .

## Details

[plotBias](#) shows the effects of the correction process on GC bias and mappability bias in HTS readcounts.

[plotCorrection](#) shows the effects of the correction on the copy number profile. Defaultly plotting the entire first chromosome found in the list.

[plotSegments](#) shows the resultant segments and states assigned to each segment.

[plotParam](#) shows the initial suggested distribution of copy number in each state (dashed), and the optimal distribution of copy number in each state (solid)

[stateCols](#) returns a vector of six colours used in [plotSegments](#) and [plotParam](#)

## Author(s)

Daniel Lai

**See Also**

[correctReadcount](#) and [HMMsegment](#) for generating intended output.

**Examples**

```
data(tumour)

# Visualize one at a time
par(ask = TRUE)
plotBias(normal_copy)
plotCorrection(tumour_copy)
par(mfrow = c(1, 1))
plotSegments(tumour_copy, tumour_segments)
plotParam(tumour_segments, tumour_param)
```

---

WIG Import Functions

*WIG Import Functions*

---

**Description**

Fast fixedStep WIG file reading and parsing

**Usage**

```
wigToRangedData(wigfile, verbose = TRUE)
wigToArray(wigfile, verbose = TRUE)
```

**Arguments**

wigfile	Filepath to fixedStep WIG format file
verbose	Set to FALSE to suppress messages

**Details**

Reads the entire file into memory, then processes the file in rapid fashion, thus performance will be limited by memory capacity.

The WIG file is expected to conform to the minimal fixedStep WIG format (see References), where each chromosome is started by a “fixedStep” declaration line. The function assumes only a single track in the WIG file, and will ignore any lines before the first line starting with “fixedStep”.

**Value**

[RangedData](#) for [wigToRangedData](#) with chromosome and position information, sorted in decreasing chromosomal size and increasing position.

Numeric [array](#) for [wigToArray](#) sorted in decreasing chromosomal size and increasing position.

**Author(s)**

Daniel Lai



## References

**WIG** <http://genome.ucsc.edu/goldenPath/help/wiggle.html>

## See Also

[wigsToRangedData](#) is a wrapper around these functions for easy WIG file loading and structure formatting.

## Examples

```
wigfile <- system.file("extdata", "tumour.wig", package = "HMMcopy")
posAndValues <- wigToRangedData(wigfile)
justValues <- wigToArray(wigfile)
```

---

WIG Output Functions

*WIG Output Functions*

---

## Description

Fast fixedStep WIG file formatting and output

## Usage

```
rangedDataToWig(correctOutput, file, column = "copy", sample = "R",
  verbose = TRUE)
rangedDataToSeg(correctOutput, file, column = "copy", sample = "R",
  verbose = TRUE)
```

## Arguments

<code>correctOutput</code>	<code>RangedData</code> object for output, default options expect output from <code>correctReadcount</code> .
<code>file</code>	Filepath to write output to.
<code>column</code>	Column in input object to export. Defaults to corrected copy number.
<code>sample</code>	Sample name of the exported dataset, defaults to "R"
<code>verbose</code>	Set to FALSE to suppress messages.

## Details

Assumes that all ranges in data set are non-overlapping windows of fixed width covering the entire genome. Note that positions in WIG files are 1-based while those in SEG files are 0-based.

## Author(s)

Daniel Lai

## References

**WIG** <http://genome.ucsc.edu/goldenPath/help/wiggle.html>

**SEG** <http://www.broadinstitute.org/igv/SEG>

**See Also**

[correctReadcount](#) output is the intended input

**Examples**

```
data(tumour) # Load tumour_copy
rangedDataToWig(tumour_copy, file = "test.wig")
rangedDataToSeg(tumour_copy, file = "test.seg")
```

---

wigsToRangedData *Parsing and sorting of uncorrected read and sequence information files*

---

**Description**

Loads WIG files for readcount, GC, and mappability data for non-overlapping windows of fixed length (i.e. bins), and returns a structure ready to be used for readcount correction. See Details for specifics about file assumptions.

**Usage**

```
wigsToRangedData(readfile, gcfile, mapfile, verbose = FALSE)
```

**Arguments**

readfile	Pathname to WIG file containing readcounts per bin.
gcfile	Pathname to WIG file containing GC content per bin.
mapfile	Pathname to WIG file containing average mappability per bin.
verbose	Set to TRUE if messages are desired

**Details**

The number of lines in the three input files are expected to be identical, although the order and names of chromosomes in the file need not be identical. Chromosome lengths are required to be identical and unique, and if the latter is not true, the order of the chromosomes must then be identical.

At present, these three WIG files are expected to be generated by external programs, namely those from the HMMcopy suite (see See Also), rather than by existing R packages out of space and memory considerations when working with high coverage full genome samples.

**Value**

A [RangedData](#) object, where each row entry represents a bin, with the three values from the input as columns named reads, gc, and map.

**Author(s)**

Daniel Lai

## References

**correctedReadcount Suite** [TBA](#)

**WIG** <http://genome.ucsc.edu/goldenPath/help/wiggle.html>

## See Also

[correctReadcount](#), to correct the readcounts in the resultant value.

## Examples

```
rfile <- system.file("extdata", "tumour.wig", package = "HMMcopy")
gfile <- system.file("extdata", "gc.wig", package = "HMMcopy")
mfile <- system.file("extdata", "map.wig", package = "HMMcopy")

uncorrected_reads <- wigsToRangedData(rfile, gfile, mfile)
```

# Index

## \*Topic **IO**

HMMcopy Segmentation, 4  
HMMcopy-package, 2  
WIG Import Functions, 8  
WIG Output Functions, 9  
wigsToRangedData, 10

## \*Topic **datasets**

HMMcopy example dataset, 4

## \*Topic **hplot**

HMMcopy Visualization, 7

## \*Topic **manip**

correctReadcount, 3  
HMMcopy-package, 2

array, 8

correctReadcount, 3, 4-11

HMMcopy example dataset, 4  
HMMcopy Segmentation, 4  
HMMcopy Visualization, 7  
HMMcopy-dataset (HMMcopy example dataset), 4  
HMMcopy-package, 2  
HMMsegment, 4, 5, 7, 8  
HMMsegment (HMMcopy Segmentation), 4

normal\_copy (HMMcopy example dataset), 4

plot, 7  
plotBias, 7  
plotBias (HMMcopy Visualization), 7  
plotCorrection, 7  
plotCorrection (HMMcopy Visualization), 7  
plotParam, 7  
plotParam (HMMcopy Visualization), 7  
plotSegments, 7  
plotSegments (HMMcopy Visualization), 7

RangedData, 3, 4, 8-10  
rangedDataToSeg (WIG Output Functions), 9  
rangedDataToWig (WIG Output Functions), 9

stateCols, 7  
stateCols (HMMcopy Visualization), 7

tumour (HMMcopy example dataset), 4

tumour\_copy (HMMcopy example dataset), 4

tumour\_param (HMMcopy example dataset), 4

tumour\_reads (HMMcopy example dataset), 4

tumour\_segments (HMMcopy example dataset), 4

WIG Import Functions, 8  
WIG Output Functions, 9  
wigsToRangedData, 3, 4, 9, 10  
wigToArray, 8  
wigToArray (WIG Import Functions), 8  
wigToRangedData, 8  
wigToRangedData (WIG Import Functions), 8