

# exomeCopy

March 24, 2012

---

ExomeCopy-class      *Class "ExomeCopy"*

---

## Description

Object returned by exomeCopy

## Objects from the Class

Objects can be created by calls of the form `new ("ExomeCopy")`.

## Slots

`type`: Object of class "character": the type of model used, either "exomeCopy" or "exomeCopyVar"

`path`: Object of class "numeric": the index of the predicted state for each window

`ranges`: Object of class "IRangesList": the corresponding ranges for the observed counts and covariates

`O`: Object of class "numeric": the input vector of counts

`O.norm`: Object of class "numeric": the input vector of counts divided by  $X * \beta$

`mu`: Object of class "numeric":  $X * \beta$

`phi`: Object of class "numeric":  $Y * \gamma$

`fx.par`: Object of class "list": a list of the settings `S`, `d`, `normal.state` and `fit.var`

`init.par`: Object of class "list": a list of the initial parameters `goto.cnv`, `goto.normal`, `beta.hat` and `phi.hat`

`final.par`: Object of class "list": a list of the final parameters `goto.cnv`, `goto.normal`, `beta` (and `gamma` for `exomeCopyVar`)

`counts`: Object of class "numeric": the number of evaluations of the log likelihood performed by `optim`

`convergence`: Object of class "numeric": the integer for convergence of `optim`, 0 for convergence

`nll`: Object of class "numeric": the final value of the negative log likelihood

**Methods**

```
plot signature(x = "ExomeCopy", y = "missing"):...
```

```
show signature(object = "ExomeCopy"):...
```

**See Also**

[exomeCopy](#)

**Examples**

```
showClass("ExomeCopy")
```

---

copyCountSegments *Segments of identical copy count from exomeCopy*

---

**Description**

Unpacks an ExomeCopy object and returns a RangedData object with segments of identical predicted copy count in genomic coordinates.

**Usage**

```
copyCountSegments(object)
```

**Arguments**

object            ExomeCopy object

**Value**

Returns a RangedData object with the predicted copy count and the number of genomic ranges spanned by the segment.

**See Also**

[exomeCopy](#) [ExomeCopy-class](#) [RangedData](#)

**Examples**

```
example(exomeCopy)
copyCountSegments(fit)
```

---

countBamInGRanges *Count reads from BAM file in genomic ranges*

---

### Description

Counts the number of reads with a specified minimum mapping quality from BAM files in genomic ranges specified by a GRanges object. This is a convenience function for counting the reads in ranges covering the targeted regions, such as the exons in exome enrichment experiments, from each sample. These read counts are used by [exomeCopy](#) in predicting CNVs in samples.

With the default setting (`read.width=1`), only the read starts are used for counting purposes (the leftmost position regardless of the strandedness of the read).

The function [subdivideGRanges](#) can be used first to subdivide ranges of different size into ranges of nearly equal width.

The BAM file requires an associated index file (see the man page for [indexBam](#) in the Rsamtools package).

### Usage

```
countBamInGRanges(bam.file, granges, min.mapq=1, read.width=1)
```

### Arguments

<code>bam.file</code>	The path of the BAM file for the sample to be counted.
<code>granges</code>	An object of type GRanges with the ranges in which to count reads.
<code>min.mapq</code>	The minimum mapping quality to count a read. Defaults to 1. Set to 0 for counting all reads.
<code>read.width</code>	The width of a read, used in counting overlaps of mapped reads with the genomic ranges. The default is 1, resulting in the counting of only read starts in genomic ranges. If the length of fixed width reads is used, e.g. 100 for 100bp reads, then the function will return the count of all overlapping reads with the genomic ranges. However, counting all overlapping reads introduces dependency between the counts in adjacent windows.

### Value

A vector giving the number of reads over the input GRanges

### See Also

[Rsamtools GRanges subdivideGRanges](#)

### Examples

```
## get subdivided genomic ranges covering targeted region
## using subdivideGRanges()
example(subdivideGRanges)

## BAM file included in Rsamtools package
bam.file <- system.file("extdata", "mapping.bam", package="exomeCopy")
```

```
## create RangedData object to store read counts
rdata <- RangedData(space=seqnames(target.sub), ranges=ranges(target.sub))

## extract read counts from the BAM file in these genomic ranges
rdata[["sample"]] <- countBamInGRanges(bam.file, target.sub)
```

---

exomeCopy-package *Detection of CNV in exome/targeted sequencing data*

---

## Description

A hidden Markov model for the detection of copy number variants (CNV) in exome/targeted sequencing read depth data. The package uses positional covariates, such as background read depth and GC-content, to simultaneously normalize and segment the samples into regions of constant copy count.

## Details

Package:	exomeCopy
Type:	Package
Version:	1.0.3
Date:	2011-10-27
License:	GPL (>= 2)
LazyLoad:	yes
Depends:	methods, graphics, IRanges, GenomicRanges, Rsamtools (>= 1.4.3)
Suggests:	Biostrings

exomeCopy fits a hidden Markov model to observed read counts using covariates. It returns the Viterbi path, the most likely path of hidden states, which is the predicted copy count at each window.

## Author(s)

Michael Love <love@molgen.mpg.de>

## References

Love, Michael I.; Mysickova, Alena; Sun, Ruping; Kalscheuer, Vera; Vingron, Martin; and Haas, Stefan A. (2011) "Modeling Read Counts for CNV Detection in Exome Sequencing Data," *Statistical Applications in Genetics and Molecular Biology*: Vol. 10 : Iss. 1, Article 52. DOI: 10.2202/1544-6115.1732 [http://cmb.molgen.mpg.de/publications/Love\\_2011\\_exomeCopy.pdf](http://cmb.molgen.mpg.de/publications/Love_2011_exomeCopy.pdf).

## See Also

[exomeCopy](#)

exomeCopy

*Fit the exomeCopy or exomeCopyVar model to the observed counts.***Description**

Fits a hidden Markov model to observed read counts using positional covariates. It returns an object containing the fitted parameters and the Viterbi path, the most likely path of hidden states, which is the predicted copy count at each window. `exomeCopy` is designed to run on read counts from a single chromosome. Please see the vignette for an example of how to prepare input data for `exomeCopy` and how to loop the function over multiple chromosomes and samples.

`exomeCopy` requires as input a `RangedData` object containing read counts in genomic ranges along with the covariates. Two convenience functions are provided for preparing input for `exomeCopy`:

1. `subdivideGRanges`, to subdivide a `GRanges` object containing the genomic ranges of the targeted region into genomic ranges of nearly equal width, and
2. `countBamInGRanges`, to count the number of read starts from a BAM read mapping file in a `GRanges` object.

The GC-content (ratio of G and C bases to total number of bases) for the input ranges can be obtained using `scanFa` in the `Rsamtools` package to obtain a `DNAStringSet` object and `letterFrequency` in the `Biostrings` package. See the vignette for an example.

**Usage**

```
exomeCopy(rdata, sample.name, X.names, Y.names, fit.var=FALSE, reltol
= 0.0001, S = 0:4, d = 2, goto.cnv = 1e-4, goto.normal = 1/20,
init.phi="norm")
```

**Arguments**

<code>rdata</code>	A <code>RangedData</code> object with the sample counts and positional covariates over the genomic ranges.
<code>sample.name</code>	The name of the value column of <code>rdata</code> with the sample read counts.
<code>X.names</code>	The names of the value columns of <code>rdata</code> with covariates for estimating $\mu$ .
<code>Y.names</code>	(optional) the names of the value columns of <code>rdata</code> with covariates for estimating $\phi$ , only required if <code>fit.var = TRUE</code> .
<code>fit.var</code>	A logical, whether the model should fit the overdispersion parameter $\phi$ with a linear combination of covariates ( <code>exomeCopyVar</code> ) or with a scalar ( <code>exomeCopy</code> ). Defaults to <code>FALSE</code> ( <code>exomeCopy</code> ).
<code>reltol</code>	The relative tolerance for convergence used in the <code>optim</code> function for optimizing the parameter settings. From testing, the default value was sufficient for fitting parameters, but lower relative tolerances can be used.
<code>S</code>	A vector of possible copy numbers for the different states.
<code>d</code>	The expected copy number for the normal state. This should be set to 2 for autosomes and 1 for haploid data.
<code>goto.cnv</code>	The initial setting for probability to transfer to a CNV state.
<code>goto.normal</code>	The initial setting for probability to transfer to a normal state.
<code>init.phi</code>	Either "norm" or "counts": initialize $\phi$ with the moment estimate using residuals from a linear model of read counts on covariates or with the raw counts.

## Details

exomeCopy fits transitional and emission parameters of an HMM to best explain the observed counts of a sample from exome or targeted sequencing. The set of underlying copy number states,  $S$ , in the sample must be provided before running the algorithm.

The emission probabilities are given as a negative binomial distribution using positional covariates, such as background read depth, quadratic terms for GC-content, and range width, which are stored in a matrix  $X$ . Optionally, for fitting the variance of the distribution, the standard deviation and/or variance of the background set can be included in a matrix  $Y$ . All covariates are normalized within exomeCopy for improved optimization.

For the observed count at range  $t$ ,  $O_t$ , the emission probability is given by:

$$f \sim \text{NB}(O_t, \mu_{ti}, \phi)$$

The mean parameter  $\mu_{ti}$  is given by:

$$\mu_{ti} = \frac{S_i}{d} (x_{t*} \vec{\beta})$$

Here  $S_i$  is the  $i$ -th possible copy number state,  $d$  is the expected background copy number ( $d = 2$  for diploid sequence), and  $\vec{\beta}$  is a vector of coefficients fitted by the model.  $x_{t*}$  is the  $t$ -th row of the matrix  $X$ .

$\mu$  must be positive, so it is replaced with a small positive number if the value is less than zero.

For exomeCopyVar, which also fits the variance, the emission probability is given by:

$$f \sim \text{NB}(O_t, \mu_{ti}, \phi_t)$$

where

$$\phi_t = y_{t*} \vec{\gamma}$$

or a small positive number if this is less than zero.

Two transition probabilities are fitted in the model: the probabilities of transitioning to a normal state and to a CNV state.

exomeCopy calls `negLogLike` to evaluate the likelihood of the HMM. The parameters are fit using Nelder-Mead optimization with the `optim` function on the negative likelihood. The viterbi path is calculated by calling `viterbiPath`.

## Value

Returns an ExomeCopy object with the following slots:

`type`: the type of model used, either "exomeCopy" or "exomeCopyVar"

`path`: the index of the predicted state for each genomic range

`ranges`: the IRangesList for ranges

`O`: the input vector of counts

`O.norm`: the input vector of counts divided by the estimated mean vector,  $\mu$

`mu`: the estimated mean vector, matrix multiplication of  $X$  and  $\beta$

`phi`: a scalar estimate of  $\phi$  ( or matrix multiplication of  $Y$  times  $\gamma$  for exomeCopyVar)

`fx.par`: a list of the settings `S`, `d`, `cnv.states`, and the logical variable `fit.var`  
`init.par`: a list of the initial parameters `goto.cnv`, `goto.normal`, `beta.hat` and `phi.hat`  
`final.par`: a list of the final parameters `goto.cnv`, `goto.normal`, `beta` (and `gamma` for `exomeCopyVar`)  
`counts`: the number of evaluations of the log likelihood performed by `optim`  
`convergence`: the integer for convergence of `optim`, 0 for convergence  
`nll`: the final value of the negative log likelihood

### Author(s)

Michael Love

### References

Love, Michael I.; Mysickova, Alena; Sun, Ruping; Kalscheuer, Vera; Vingron, Martin; and Haas, Stefan A. (2011) "Modeling Read Counts for CNV Detection in Exome Sequencing Data," *Statistical Applications in Genetics and Molecular Biology*: Vol. 10 : Iss. 1, Article 52. DOI: 10.2202/1544-6115.1732 [http://cmb.molgen.mpg.de/publications/Love\\_2011\\_exomeCopy.pdf](http://cmb.molgen.mpg.de/publications/Love_2011_exomeCopy.pdf).

References for HMM algorithms and use of HMM for segmentation of genomic data by copy number:

Rabiner, L. R. (1989): "A tutorial on hidden Markov models and selected applications in speech recognition," *Proceedings of the IEEE*, 77, 257, 286, <http://dx.doi.org/10.1109/5.18626>.

Fridlyand, J., A. M. Snijders, D. Pinkel, D. G. Albertson, and Jain (2004): "Hidden Markov models approach to the analysis of array CGH data," *Journal of Multivariate Analysis*, 90, 132, 153, <http://dx.doi.org/10.1016/j.jmva.2004.02.008>.

Marioni, J. C., N. P. Thorne, and S. Tavare (2006): "BioHMM: a heterogeneous hidden Markov model for segmenting array CGH data." *Bioinformatics*, 22, 1144, 1146, <http://view.ncbi.nlm.nih.gov/pubmed/16533818>.

### See Also

[subdivideGRanges](#) [countBamInGRanges](#) [copyCountSegments](#) [plot.ExomeCopy](#)  
[negLogLike](#) [IRanges](#) [RangedData](#)

### Examples

```
## The following is an example of running exomeCopy on simulated
## read counts using the model parameters defined above. For an example
## using real exome sequencing read counts (with simulated CNV) please
## see the vignette.

## create RangedData for storing genomic ranges and covariate data
## (background, background stdev, GC-content)
m <- 5000
rdata <- RangedData(IRanges(start=0:(m-1)*100+1,width=100),space=rep("chr1",m),universe="

## create read depth distributional parameters mu and phi
rdata$gc.sq <- rdata$gc^2
X <- cbind(bg=rdata$bg,gc=rdata$gc,gc.sq=rdata$gc.sq)
```

```

Y <- cbind(bg.sd=rdata$bg.sd)
beta <- c(20,10,2,-.01)
gamma <- c(.1,.05)
rdata$mu <- beta[1] + scale(X) %*% beta[2:4]
rdata$mu[rdata$mu<1e-8] <- 1e-8
rdata$phi <- gamma[1] + scale(Y) %*% gamma[2]
rdata$phi[rdata$phi<1e-8] <- 1e-8

## create observed counts with simulated heterozygous duplication
cnv.nranges <- 200
bounds <- (round(m/2)+1):(round(m/2)+cnv.nranges)
O <- rbinom(nrow(rdata),mu=rdata$mu,size=1/rdata$phi)
O[bounds] <- O[bounds] + rbinom(cnv.nranges,prob=0.5,size=O[bounds])
rdata[["sample1"]] <- O

## run exomeCopy() and list segments
fit <- exomeCopy(rdata,"sample1",X.names=c("bg","gc","gc.sq"))

## see man page for copyCountSegments() for summary of
## the predicted segments of constant copy count, and
## for plot.ExomeCopy() for plotting fitted objects

```

---

exomecounts

*Sample counts from 16 exome sequencing samples from 1000 Genomes Project*

---

## Description

This data set gives sample read counts in 1000 genomic ranges for 16 exome sequencing samples from the PUR population of the 1000 Genomes Project, along with the GC-content in the ranges. For instructions on how to prepare read count and covariate data, please see the example code in the man pages for [subdivideGRanges](#) and [countBamInGRanges](#).

The genomic ranges are generated from small portion of the CCDS regions of chromosome 1 (hg19). The CCDS regions are subdivided evenly into ranges around 100bp using the [subdivideGRanges](#) function with default settings. Only ranges with positive counts across samples are retained. These regions were downloaded as a BED file from the UCSC Genome Browser (<http://genome.ucsc.edu/cgi-bin/hgGateway>). The mapping files for the exome sequencing data and descriptions of the experiments are available at the 1000 Genomes Project website (<http://www.1000genomes.org/data>). The directories used are listed in the file `1000Genomes_files.txt` in the `extdata` directory.

The column names are the sample names from the 1000 Genomes Project. Library format is paired-end reads and sample counts reflect both sequenced reads counted in their respective genomic ranges.

## Usage

```
data(exomecounts)
```

## Format

A `RangedData` object.



**Source**

1000 Genomes Project and Consensus Coding Sequence Project

**References**

1000 Genomes Project Consortium. A map of human genome variation from population-scale sequencing. *Nature* 467, 1061-1073 (2010). <http://dx.doi.org/10.1038/nature09534>.

1000 Genomes Project: Release of phase 1 exome alignments <http://www.1000genomes.org/announcements/release-phase-1-exome-alignments-2011-07-19>

Pruitt, K. D. et al. The consensus coding sequence (CCDS) project: Identifying a common protein-coding gene set for the human and mouse genomes. *Genome research* 19, 1316-1323 (2009). <http://dx.doi.org/10.1101/gr.080531.108>.

**See Also**

[RangedData](#)

---

negLogLike

*Generalized negative log likelihood and Viterbi algorithms*

---

**Description**

negLogLike: Returns the negative log likelihood calculated with the forward equations.

viterbiPath: Calculates the most likely sequence of hidden states for the Markov model given the current parameters.

**Usage**

```
negLogLike(par, fx.par, data, nstates, stFn, trFn, emFn)
viterbiPath(par, fx.par, data, nstates, stFn, trFn, emFn)
```

**Arguments**

par	A list of parameters, over which the likelihood will be optimized.
fx.par	A list of fixed parameters.
data	A list of data objects, which must contain a vector O, which represents the observed sequence of the HMM.
nstates	The number of states of the HMM.
stFn	A function which takes arguments par, fx.par, data, and nstates, and returns a vector of length nstates of starting probabilities.
trFn	A function which takes arguments par, fx.par, data, and nstates, and returns a matrix of dimension (nstates,nstates) of the transition probabilities.
emFn	A function which takes arguments par, fx.par, data, and nstates, and returns a matrix of dimension (nstates,length(O)) of the emission probabilities.

**Value**

negLogLike: The negative log likelihood of the HMM. The likelihood is slightly modified to account for ranges with read counts which have zero probability of originating from any of the states. In this case the likelihood is lowered and the range is skipped.

viterbiPath: The Viterbi path through the states given the parameters.

**References**

On the forward equations and the Viterbi algorithm:

Rabiner, L. R. (1989): "A tutorial on hidden Markov models and selected applications in speech recognition," Proceedings of the IEEE, 77, 257, 286, <http://dx.doi.org/10.1109/5.18626>.

**Examples**

```
## functions for starting, transition, and emission probabilities
stFn <- function(par, fx.par, data, nstates) rep(1/nstates, nstates)
trFn <- function(par, fx.par, data, nstates) {
  A <- matrix(1/(nstates*10), ncol=nstates, nrow=nstates)
  diag(A) <- 1 - rowSums(A)
  A
}
emFn <- function(par, fx.par, data, nstates) {
  t(sapply(1:nstates, function(j) dnorm(data$O, par$means[j], fx.par$sdev)))
}

## simulate some observations from two states
Q <- c(rep(1, 100), rep(2, 100), rep(1, 100), rep(2, 100))
T <- length(Q)
means <- c(-0.5, 0.5)
sdev <- 1
O <- rnorm(T, means[Q], sdev)

## use viterbiPath() to recover the state chain using parameters
viterbi.path <- viterbiPath(par=list(means=means), fx.par=list(sdev=sdev), data=list(O=O), n=nstates)
plot(O, pch=Q, col=c("darkgreen", "orange")[viterbi.path])
```

---

plot.ExomeCopy

*Plot function for exomeCopy*

---

**Description**

Plots the predicted copy count segments of an ExomeCopy object

**Usage**

```
## S3 method for class 'ExomeCopy'
plot(x, points = TRUE, cols = NULL, show.legend = TRUE,
     main = "exomeCopy predicted segments", xlab = "genomic position",
     ylab = "normalized read count", xlim = NULL, ylim = NULL, cex = 1, lwd = 4, ...)
```

**Arguments**

<code>x</code>	The ExomeCopy object.
<code>points</code>	Logical, whether normalized read counts should be drawn.
<code>cols</code>	A vector of the same length as <code>b</code> , specifying a color for each of the states of the HMM.
<code>show.legend</code>	Logical, whether a default legend should be shown.
<code>main</code>	main title
<code>xlab</code>	x axis label
<code>ylab</code>	y axis label
<code>xlim</code>	x limits
<code>ylim</code>	y limits
<code>cex</code>	size of the points (if plotted)
<code>lwd</code>	line width
<code>...</code>	Other arguments passed to <code>plot()</code>

**See Also**

[exomeCopy](#) [ExomeCopy-class](#) [copyCountSegments](#)

**Examples**

```
example(exomeCopy)
plot(fit)
```

---

subdivideGRanges     *Subdivide ranges of a GRanges object into nearly equal width ranges*

---

**Description**

Takes an input `GRanges` object and, splits each range into multiple ranges of nearly equal width. For an input range of width `w` and subdividing size `s`, it will subdivide the range into  $\max(1, \text{floor}(w/s))$  nearly equal width ranges. The output is then a new `GRanges` object. This function can be used to split the targeted region (such as exons in exome enrichment experiments) into nearly equal width ranges.

**Usage**

```
subdivideGRanges(x, subsize=100)
```

**Arguments**

<code>x</code>	An object of type <code>GRanges</code> .
<code>subsize</code>	The desired width for the ranges in the output <code>GRanges</code> object.

**Value**

A GRanges object with ranges from the input GRanges object subdivided to nearly subsize.

**See Also**

[GRanges](#)

**Examples**

```
## read in target region BED file
target.file <- system.file("extdata", "targets.bed", package="exomeCopy")
target.df <- read.delim(target.file, header=FALSE,
  col.names=c("seqname", "start", "end"))

## create GRanges object with 5 ranges over 2 sequences
target <- GRanges(seqname=target.df$seqname,
  IRanges(start=target.df$start, end=target.df$end))

## subdivide into 7 smaller genomic ranges
target.sub <- subdivideGRanges(target)
```

# Index

## \*Topic classes

ExomeCopy-class, 1

## \*Topic datasets

exomecounts, 8

## \*Topic package

exomeCopy-package, 4

copyCountSegments, 2, 7, 11

countBamInGRanges, 3, 5, 7, 8

exomeCopy, 2-4, 5, 11

ExomeCopy-class, 2, 11

ExomeCopy-class, 1

exomeCopy-package, 4

exomecounts, 8

GRanges, 3, 5, 12

indexBam, 3

IRanges, 7

negLogLike, 6, 7, 9

optim, 1, 5-7

plot, ExomeCopy, missing-method  
(ExomeCopy-class), 1

plot.ExomeCopy, 7, 10

RangedData, 2, 5, 7, 9

Rsamtools, 3

show, ExomeCopy-method  
(ExomeCopy-class), 1

subdivideGRanges, 3, 5, 7, 8, 11

viterbiPath, 6

viterbiPath(negLogLike), 9