

# Package ‘BEclear’

October 15, 2018

**Type** Package

**Title** Correct for batch effects in DNA methylation data

**Version** 1.12.1

**Date** 2014-11-03

**Author** Markus Merl, Ruslan Akulenko

**Maintainer** David Rasp <David.j.Rasp@gmail.com>

**Description** Provides some functions to detect and correct for batch effects in DNA methylation data. The core function “BEclear” is based on latent factor models and can also be used to predict missing values in any other matrix containing real numbers.

**License** GPL-2

**Depends** snowfall, Matrix

**biocViews** BatchEffect, DNAMethylation, Software

**git\_url** <https://git.bioconductor.org/packages/BEclear>

**git\_branch** RELEASE\_3\_7

**git\_last\_commit** 4ec4f42

**git\_last\_commit\_date** 2018-09-26

**Date/Publication** 2018-10-15

## R topics documented:

BEclear-package . . . . .	2
BEclear . . . . .	4
BEclear example data . . . . .	7
calcMedians . . . . .	8
calcPvalues . . . . .	9
calcScore . . . . .	10
calcSummary . . . . .	11
clearBEgenes . . . . .	12
correctBatchEffect . . . . .	13
countValuesToPredict . . . . .	16
ex.corrected.data . . . . .	17
findWrongValues . . . . .	17
makeBoxplot . . . . .	18
replaceWrongValues . . . . .	19

<b>Index</b>	<b>21</b>
--------------	-----------

---

BEclear-package

*Correct for batch effects in DNA methylation data*

---

## Description

Provides some functions to detect and correct for batch effects in DNA methylation data. The core function [BEclear](#) is based on Latent Factor Models and can also be used to predict missing values in any other matrix containing real numbers.

## Details

Package: BEclear  
Type: Package  
Version: 1.0  
Date: 2014-11-03  
License: GPL-2

[correctBatchEffect](#): The function combines most functions of the [BEclear-package](#) to one. This function performs the whole process of searching for batch effects and automatically correct them for a matrix of beta values stemming from DNA methylation data.

[BEclear](#): This function predicts the missing entries of an input matrix (NA values) through the use of a Latent Factor Model.

[calcMedians](#): Compares the median value of all beta values belonging to one batch with the median value of all beta values belonging to all other batches. Returns a matrix containing this median difference value for every gene in every batch, columns define the batch numbers, rows the gene names.

[calcPvalues](#): Compares the distribution of all beta values corresponding to one batch with the distribution of all beta values corresponding to all other batches and returns a p-value which defines if the distributions are the same or not.

[calcSummary](#): Summarizes the results of the median comparison function [calcMedians](#) and the p-value calculation function [calcPvalues](#). Should be used with the matrices originating from these two functions.

[calcScore](#): Returns a table with the number of found genes with found p-values less or equal to 0.01 and median values greater or equal to 0.05. A score is calculated depending on the number of found genes as well as the magnitude of the median difference values, this score is divided by the overall number of genes in the data and returned as "BEScore". See the methods details for further information and details about the score calculation.

[makeBoxplot](#): A simple [boxplot](#) is done with boxes either separated by batches or by samples and describe the five number summary of all beta values corresponding to a batch or a sample, respectively. The batch\_ids are shown on the x-axis with a coloring corresponding to the BEScore.

[clearBEGenes](#): A function that simply sets all values to NA which were previously found by median value comparison and p-value calculation and are stored in a summary. The summary defines which values in the data matrix are set to NA.

[countValuesToPredict](#): Simple function that counts all values in a matrix which are NA

[findWrongValues](#): A method which lists values below 0 or beyond 1 contained in the input matrix. The wrong entries are stored in a data.frame together with the corresponding row and column position of the matrix.

`replaceWrongValues`: A method which replaces values below 0 or beyond 1 contained in the input matrix. These wrong entries are replaced by 0 or 1, respectively.

### Author(s)

Ruslan Akulenko, Markus Merl

Maintainer: Markus Merl <merl.markus@googlemail.com>

### References

Y. Koren, R. Bell, C. Volinsky, Matrix factorization techniques for recommender systems, IEEE Computer, 42(8), p. 30-37, 2009, <http://research.yahoo.com/pub/2859>

E. Candes, B. Recht, Exact matrix completion via convex optimization, Communications of the ACM, 55(6), p. 111-119, 2012, <http://doi.acm.org/10.114/2184319.2184343>

### Examples

```
data(BEclearData)

## Calculate median comparison values in non-parallel mode
med <- calcMedians(data=ex.data, samples=ex.samples, parallel=FALSE)

## Calculate fdr-adjusted p-values in non-parallel mode
pvals <- calcPvalues(data=ex.data, samples=ex.samples, parallel=FALSE,
  adjusted=TRUE, method="fdr")

## Summarize p-values and median differences for batch affected genes
sum <- calcSummary(medians=med, pvalues=pvals)

## Calculates the score table
score.table <- calcScore(data=ex.data, samples=ex.samples, summary=sum)

## Simple boxplot for the example data separated by batch
makeBoxplot(data=ex.data, samples=ex.samples, score=score.table,
  bySamples=FALSE, main="Some box plot")

## Simple boxplot for the example data separated by samples
makeBoxplot(data=ex.data, samples=ex.samples, score=score.table,
  bySamples=TRUE, main="Some box plot")

## Sets assumed batch affected entries to NA
cleared <- clearBEgenes(data=ex.data, samples=ex.samples, summary=sum)

## Counts and stores number of entries to predict
numberOfEntries <- countValuesToPredict(data=cleared)

## Not run:
## Predicts the missing entries
predicted <- BEclear(data=cleared)

## Find wrongly predicted entries
wrongEntries <- findWrongValues(data=predicted)

## Replace wrongly predicted entries
corrected <- replaceWrongValues(data=predicted)
```

```
## End(Not run)
```

---

 BEclear

---

*Matrix prediction method using a Latent Factor Model*


---

## Description

This function predicts the missing entries of an input matrix (NA values) through the use of a Latent Factor Model. You can run the function also in parallel mode and split up the matrix to a certain number of smaller matrices to speed up the prediction process. If you set the `rowBlockSize` and `colBlockSize` both to 0, the function is running on the whole matrix. Take a look at the details section for some deeper information about this. The default parameters are chosen with the intention to make an accurate prediction within an affordable time interval.

## Usage

```
BEclear(data, parallel=TRUE, cores=4, rowBlockSize=60, colBlockSize=60,
        epochs=50, outputFormat="RData", dir=getwd())
```

## Arguments

<code>data</code>	any matrix filled e.g. with beta values. The missing entries you want to predict have to be set to NA.
<code>parallel</code>	should the calculation be done in parallel mode? <a href="#">snowfall</a> package is needed to run the function in parallel mode.
<code>cores</code>	if running in parallel mode, define the number of cores used for the calculation. <a href="#">snowfall</a> package is needed to run the function in parallel mode
<code>rowBlockSize</code>	the number of rows that is used in a block if the function is run in parallel mode and/or not on the whole matrix. Set this, and the " <code>colBlockSize</code> " parameter to 0 if you want to run the function on the whole input matrix. We suggest to use a block size of 60 but you can also use any other block size, but the size has to be bigger than the number of samples in the biggest batch. Look at the details section for more information about this feature.
<code>colBlockSize</code>	the number of columns that is used in a block if the function is run in parallel mode and/or not on the whole matrix. Set this, and the " <code>rowBlockSize</code> " parameter to 0 if you want to run the function on the whole input matrix. We suggest to use a block size of 60 but you can also use any other block size, but the size has to be bigger than the number of samples in the biggest batch. Look at the details section for more information about this feature.
<code>epochs</code>	the number of iterations used in the gradient descent algorithm to predict the missing entries in the data matrix.
<code>outputFormat</code>	you can choose if the finally returned data matrix should be saved as an <code>.RData</code> file or as a tab-delimited <code>.txt</code> file in the specified directory. Allowed values are " <code>RData</code> " and " <code>txt</code> ".
<code>dir</code>	set the path to a directory the predicted matrix should be stored. The current working directory is defined as default parameter.

## Details

The method used to predict the missing entries in the matrix is called "latent factor model". In the following sections, the method itself is described as well as the correct usage of the parameters. The parameters are described in the same order as they appear in the usage section.

The method originally stems from recommender systems where the goal is to predict user ratings of products. It is based on matrix factorization and uses a discrete gradient descent (GDE) algorithm that stepwise predicts two matrices L and R with matching dimensions to the input matrix. These two matrices are initialized with random numbers and stepwise adjusted towards the values of the input matrix through the GDE algorithm. After every adjustment step, the global loss is calculated and the parameters used for the adjustment are possibly also adjusted so that the global loss is getting minimized and the prediction is getting accurate. After a predefined number of steps (called epochs) are executed by the GDE algorithm, the predicted matrix is calculated by matrix multiplication of L and R. Finally, all missing values in the input matrix are replaced with the values from the predicted matrix and the already known values from the input matrix are maintained. The completed input matrix is then returned at the end.

Description of the parameters: - data: simply the input matrix with missing values set to NA

- parallel: if TRUE, the function is running in parallel mode. The `snowfall` package has to be installed. You can also define the number of cores to be used. If running in parallel mode, the blocks are distributed over the cores and after all blocks are processed, the finally returned matrix is assembled from these blocks. Note that every block is intermediately saved in the specified directory after its processing has been done. Look at the section "About the blocks" for more details. Note that if you run the function in non-parallel mode (parallel = FALSE), the matrix is still divided into blocks of specified size which are processed sequentially.

- cores: Specify the number of cores if you run the function in parallel mode.

- rowBlockSize and colBlockSize: Here you can define the dimensions of the smaller matrices, the input matrix is divided into if the function is working in parallel mode. For details about these so called blocks, see the section "About the blocks" below.

- epochs: Defines the number of steps the gradient descent algorithm performs until the prediction ends. Note that the higher this number is, the more precisely is the prediction and the more time is needed to perform the prediction. If the step size is too small, the prediction would not be very good. We suggest to use a step size of 50 since we did not get better predictions if we took higher step sizes during our testing process.

About the blocks: You have the possibility to change the size of the blocks in which the input matrix can be divided. if you choose e.g. the rowBlockSize = 50 and the colBlockSize = 60 your matrix will be cut into smaller matrices of the size approximately 50x60. Note that this splitting algorithm works with every possible matrix size! If both size parameters do not fit to the dimensions of the input matrix, the remaining rows and columns of the input matrix are distributed over some blocks, so that the block sizes are roughly of the same size. All blocks are saved at the specified directory after the processing of a block has been done within an RData file. These RData files are continuously numbered and contain the row and column start and stop positions in their name. Next, these blocks are assembled into the returned matrix and this matrix is saved in the specified directory. Finally, single blocks are deleted. To see how this is done, simply run the example at the end of this documentation. We suggest to use the block size of 60 (default) but you can also use any other block size, as far as it is bigger than the number of samples in the biggest batch. This avoids having an entire row of NA values in a block which leads to a crash of the BEclear method. In order to process the complete matrix without dividing into blocks, specify rowBlockSize = 0 and colBlockSize = 0. But if the input matrix is large (more than 200x200), it is not recommended due to exponential increase of computation time required.

Note that the size of the blocks affect the prediction accuracy. In case of very small blocks, the information obtained from neighbor entries is not sufficient. Thus, the larger the size of the block

is, the more accurately those entries are predicted. Default size 60 is enough to have accurate prediction in a reasonable amount of time.

### Value

Returns a data matrix with the same dimensions as well as same row and column names as the input matrix. According to the "outputFormat" parameter, either a .RData file containing only the returned matrix or a tab-delimited .txt file containing the content of the returned matrix is saved in the specified directory.

### References

Y. Koren, R. Bell, C. Volinsky, Matrix factorization techniques for recommender systems, *IEEE Computer*, 42(8), p. 30-37, 2009, <http://research.yahoo.com/pub/2859>

E. Candes, B. Recht, Exact matrix completion via convex optimization, *Communications of the ACM*, 55(6), p. 111-119, 2012, <http://doi.acm.org/10.114/2184319.2184343>

### See Also

[snowfall](#) [findWrongValues](#) [replaceWrongValues](#) [correctBatchEffect](#)

### Examples

```
## Shortly running example. For a more realistic example that takes
## some more time, run the same procedure with the full BEclearData
## dataset.

## Whole procedure that has to be done to use this function.
data(BEclearData)
ex.data <- ex.data[31:90,7:26]
ex.samples <- ex.samples[7:26,]

# Calculate median difference values and p-values
meds <- calcMedians(data=ex.data, samples=ex.samples)
pvals <- calcPvalues(data=ex.data, samples=ex.samples)

# Summarize p-values and median differences for batch affected genes
sum <- calcSummary(medians=meds, pvalues=pvals)

# Set entries defined by the summary to NA
clearedMatrix <- clearBEgenes(data=ex.data, samples=ex.samples, summary=sum)

# Predict the missing entries with standard values, row- and block sizes are
# just set to 10 to get a short runtime. To use these parameters, either use
# the default values or please note the description in the details section
# above
predicted <- BEclear(data=clearedMatrix, rowBlockSize=10, colBlockSize=10)
```

---

BEclear example data *Example data set for the BEclear-package*

---

### Description

Example data needed to run all examples contained in the documentation files of the BEclear-package

### Usage

```
data(BEclearData)
```

### Format

Contains a data matrix and a samples data frame which can be used to run the examples of the BEclear package.

`ex.data` An example data matrix that is filled with beta values originally stemming from breast cancer data from the TCGA portal [1], `colnames` are sample ids, `rownames` are gene names. Generally, beta values are calculated by dividing the methylated signal by the sum of the unmethylated and methylated signals from a DNA methylation microarray. The sample data used here contains averaged beta values of probes that belong to promoter regions of single genes. Another possibility would be to use beta values of single probes, whereby the probe names should then be used instead of the gene names as `rownames` of the matrix.

`ex.samples` An example data frame containing a column for the sample id and a column for the corresponding batch id.

### Value

A data frame (`ex.samples`) and a matrix (`ex.data`)

### References

[1] <http://cancergenome.nih.gov/>

### Examples

```
## Not run:
data(BEclearData)

## Whole procedure that has to be done to make a prediction using the BEclear
## function

# Calculate median difference values and p-values
meds <- calcMedians(data=ex.data, samples=ex.samples)
pvals <- calcPvalues(data=ex.data, samples=ex.samples)

# Summarize p-values and median differences for batch affected genes
sum <- calcSummary(medians=meds, pvalues=pvals)

# Set entries defined by the summary to NA
clearedMatrix <- clearBEgenes(data=ex.data, samples=ex.samples, summary=sum)
```

```
# Predict the missing entries with standard parameters
predicted <- BEclear(data=clearedMatrix)

## End(Not run)
```

---

calcMedians                      *Calculate median difference values*

---

### Description

Compares the median value of all beta values belonging to one batch with the median value of all beta values belonging to all other batches. Returns a matrix containing this median difference value for every gene in every batch, columns define the batch numbers, rows the gene names.

### Usage

```
calcMedians(data, samples, parallel=TRUE, cores=4)
```

### Arguments

data	any matrix filled with beta values, column names have to be sample_ids corresponding to the ids listed in "samples", row names have to be gene names.
samples	data frame with two columns, the first column has to contain the sample numbers, the second column has to contain the corresponding batch number. Column names have to be named as "sample_id" and "batch_id".
parallel	should the calculation be done in parallel mode? <a href="#">snowfall</a> package is needed to run the function in parallel mode.
cores	if running in parallel mode, define the number of cores used for the calculation. <a href="#">snowfall</a> package is needed to run the function in parallel mode

### Value

a matrix containing median comparison values for all genes in all batches, the column names define the batch numbers, row names are the same gene names as contained in the input matrix.

### See Also

[snowfall](#) [correctBatchEffect](#)

### Examples

```
## Shortly running example. For a more realistic example that takes
## some more time, run the same procedure with the full BEclearData
## dataset.

## Calculate median comparison values in non-parallel mode
data(BEclearData)
ex.data <- ex.data[31:90,7:26]
ex.samples <- ex.samples[7:26,]
medians <- calcMedians(data=ex.data, samples=ex.samples, parallel=FALSE)
```



---

calcPvalues	<i>Calculate p-values</i>
-------------	---------------------------

---

### Description

Compares the distribution of all beta values corresponding to one batch with the distribution of all beta values corresponding to all other batches and returns a p-value which defines if the distributions are the same or not. Standard two sided Kolmogorov-Smirnov test is used to calculate the (adjusted) p-values.

### Usage

```
calcPvalues(data, samples, parallel=TRUE, cores=4, adjusted=TRUE, method="fdr")
```

### Arguments

data	any matrix filled with beta values, column names have to be sample_ids corresponding to the ids listed in "samples", row names have to be gene names.
samples	data frame with two columns, the first column has to contain the sample numbers, the second column has to contain the corresponding batch number. Column names have to be named as "sample_id" and "batch_id".
parallel	should the calculation be done in parallel mode? <a href="#">snowfall</a> package is needed to run the function in parallel mode.
cores	if running in parallel mode, define the number of cores used for the calculation. <a href="#">snowfall</a> package is needed to run the function in parallel mode.
adjusted	should the p-values be adjusted or not, see "method" for available adjustment methods.
method	adjustment method for p-value adjustment (if TRUE), default method is "false discovery rate adjustment", for other available methods see the description of the used standard R package <a href="#">p.adjust</a> .

### Value

a matrix containing p-values for all genes in all batches, the column names define the batch numbers, row names are the same gene names as contained in the input matrix.

### See Also

[snowfall](#) [ks.test](#) [p.adjust](#) [correctBatchEffect](#)

### Examples

```
## Shortly running example. For a more realistic example that takes
## some more time, run the same procedure with the full BEclearData
## dataset.

## Calculate fdr-adjusted p-values in non-parallel mode
data(BEclearData)
ex.data <- ex.data[31:90,7:26]
ex.samples <- ex.samples[7:26,]
pvals <- calcPvalues(data=ex.data, samples=ex.samples, parallel=FALSE,
  adjusted=TRUE, method="fdr")
```

---

calcScore	<i>calculate batch effect score</i>
-----------	-------------------------------------

---

### Description

Returns a table with the number of found genes with found p-values less or equal to 0.01 and median values greater or equal to 0.05. A score is calculated depending on the number of found genes as well as the magnitude of the median difference values, this score is divided by the overall number of genes in the data and returned as "BEScore". See details for further information and details about the score calculation. The returned data.frame is also stored in the specified directory as .RData file.

### Usage

```
calcScore(data, samples, summary, dir=getwd())
```

### Arguments

data	any matrix filled with beta values, column names have to be sample_ids corresponding to the ids listed in "samples", row names have to be gene names.
samples	data frame with two columns, the first column has to contain the sample numbers, the second column has to contain the corresponding batch number. Column names have to be named as "sample_id" and "batch_id".
summary	a summary data.frame containing the columns "gene", "batch", "median" and "p-value" and consists of all genes which were found in the median and p-value calculations, see <a href="#">calcSummary</a> function for more details.
dir	set the path to a directory the returned data.frame should be stored. The current working directory is defined as default parameter.

### Details

The returned data frame contains one column for the batch numbers, 11 columns containing the number of genes found in a certain range of the median difference value and a column with the calculated BEScore. These found genes are assumed to be batch affected due to their difference in median values and their different distribution of the beta values. The higher the found number of genes and the more extreme the median difference is, the more severe is the assumed batch effect supposed to be. We suggest that there is no need for a batch effect correction if the BEScore for a batch is less than 0.02. BEScores between 0.02 and 0.1 are lying in a "grey zone" for which we assume a not severe batch effect, and values beyond 0.1 certainly describe a batch effect and should definitely be corrected.

The 11 columns containing the numbers of found genes count the median difference values which are ranging from  $\geq 0.05$  to  $< 0.1$  ;  $\geq 0.1$  to  $< 0.2$  ;  $\geq 0.2$  to  $< 0.3$  and so on up to a limit of 1.

The BEScore is calculated by the sum of the weighted number of genes divided by the number of genes. Weightings are calculated by multiplying the number of found genes between 0.05 and 0.1 by 1, between 0.1 and 0.2 by 2, between 0.2 and 0.3 by 4, between 0.3 and 0.4 by 6 and so on.

### Value

A data.frame is returned containing the number of found genes assumed to be batch affected separated by batch and a BEScore for every batch. The data.frame is also stored in the specified directory as .RData file.

**See Also**

[calcMedians](#) [calcPvalues](#) [calcSummary](#) [correctBatchEffect](#)

**Examples**

```
## Shortly running example. For a more realistic example that takes
## some more time, run the same procedure with the full BEclearData
## dataset.

## Whole procedure that has to be done to use this function.
data(BEclearData)
ex.data <- ex.data[31:90,7:26]
ex.samples <- ex.samples[7:26,]

# Calculates median difference values and p-values from the example data
med <- calcMedians(data=ex.data, samples=ex.samples, parallel=FALSE)
pvals <- calcPvalues(data=ex.data, samples=ex.samples, parallel=FALSE,
  adjusted=TRUE, method="fdr")

# Summarize p-values and median differences for batch affected genes
sum <- calcSummary(medians=med, pvalues=pvals)

# Calculates the score table
score.table <- calcScore(data=ex.data, samples=ex.samples, summary=sum,
  dir=getwd())
```

---

calcSummary

*Summarize median comparison and p-value calculation results*

---

**Description**

Summarizes the results of the median comparison function [calcMedians](#) and the p-value calculation function [calcPvalues](#). Should be used with the matrices originating from these two functions.

**Usage**

```
calcSummary(medians, pvalues)
```

**Arguments**

medians	a matrix containing median difference values calculated by the <a href="#">calcMedians</a> function. For further details look at the documentation of this function.
pvalues	a matrix containing p-values calculated by the <a href="#">calcPvalues</a> function. For further details look at the documentation of this function.

**Details**

All genes with a median comparison value  $\geq 0.05$  and a p-value of  $\leq 0.01$  are summarized into a data.frame. These genes are assumed to contain a batch effect

**Value**

A data frame with the columns "gene" containing the gene name, "batch" containing the batch number from which the gene was found, "median" and "p-value" containing the calculated median difference values and the p-values, respectively.

**See Also**

[calcMedians](#) [calcPvalues](#) [correctBatchEffect](#)

**Examples**

```
## Shortly running example. For a more realistic example that takes
## some more time, run the same procedure with the full BEclearData
## dataset.

## Whole procedure that has to be done to use this function.
data(BEclearData)
ex.data <- ex.data[31:90,7:26]
ex.samples <- ex.samples[7:26,]

# Calculates median difference values and p-values from the example data
med <- calcMedians(data=ex.data, samples=ex.samples, parallel=FALSE)
pvals <- calcPvalues(data=ex.data, samples=ex.samples, parallel=FALSE,
  adjusted=TRUE, method="fdr")

# Summarize p-values and median differences for batch affected genes
sum <- calcSummary(medians=med, pvalues=pvals)
```

---

clearBEgenes

*Prepare a data matrix for the BEclear function*

---

**Description**

A function that simply sets all values to NA which were previously found by median value comparison and p-value calculation and are stored in a summary. The summary defines which values in the data matrix are set to NA.

**Usage**

```
clearBEgenes(data, samples, summary)
```

**Arguments**

data	any matrix filled with beta values, column names have to be sample_ids corresponding to the ids listed in "samples", row names have to be gene names.
samples	data frame with two columns, the first column has to contain the sample numbers, the second column has to contain the corresponding batch number. Column names have to be named as "sample_id" and "batch_id".
summary	a summary data.frame containing the columns "gene", "batch", "median" and "p-value" and consists of all genes which were found in the median and p-value calculations, see <a href="#">calcSummary</a> function for more details.

## Details

All entries belonging to genes stated in the summary are set to NA for the corresponding batches in the data matrix. Please look at the descriptions of [calcMedians](#) and [calcPvalues](#) for more detailed information about the data which should be contained in the summary data.frame.

## Value

A data matrix with the same dimensions as well as the same column and row names as the input data matrix is returned, all entries which are defined in the summary are now set to NA.

## See Also

[calcMedians](#) [calcPvalues](#) [calcSummary](#) [correctBatchEffect](#)

## Examples

```
## Shortly running example. For a more realistic example that takes
## some more time, run the same procedure with the full BEclearData
## dataset.

## Whole procedure that has to be done to use this function.
data(BEclearData)
ex.data <- ex.data[31:90,7:26]
ex.samples <- ex.samples[7:26,]

# Calculate median difference values and p-values
meds <- calcMedians(data=ex.data, samples=ex.samples)
pvals <- calcPvalues(data=ex.data, samples=ex.samples)

# Summarize p-values and median differences for batch affected genes
sum <- calcSummary(medians=meds, pvalues=pvals)

# Set values for summarized BEgenes to NA
clearedMatrix <- clearBEgenes(data=ex.data, samples=ex.samples, summary=sum)
```

---

correctBatchEffect      *Correct a batch effect in DNA methylation data*

---

## Description

This method combines most functions of the [BEclear-package](#) to one. The method performs the whole process of searching for batch effects and automatically correct them for a matrix of beta values stemming from DNA methylation data.

## Usage

```
correctBatchEffect(data, samples, parallel=TRUE, cores=4, adjusted=TRUE,
  method="fdr", rowBlockSize=60, colBlockSize=60, epochs=50,
  outputFormat="RData", dir=getwd())
```

**Arguments**

data	any matrix filled with beta values, column names have to be sample_ids corresponding to the ids listed in "samples", row names have to be gene names.
samples	data frame with two columns, the first column has to contain the sample numbers, the second column has to contain the corresponding batch number. Column names have to be named as "sample_id" and "batch_id".
parallel	should the calculation be done in parallel mode? <a href="#">snowfall</a> package is needed to run the function in parallel mode.
cores	if running in parallel mode, define the number of cores used for the calculation. <a href="#">snowfall</a> package is needed to run the function in parallel mode
adjusted	should the p-values be adjusted or not, see "method" for available adjustment methods.
method	adjustment method for p-value adjustment (if TRUE), default method is "false discovery rate adjustment", for other available methods see the description of the used standard R package <a href="#">p.adjust</a> . See <a href="#">calcPvalues</a> for more information
rowBlockSize	the number of rows that is used in a block if the function is run in parallel mode and/or not on the whole matrix. Set this, and the "colBlockSize" parameter to 0 if you want to run the function on the whole input matrix. See <a href="#">BEclear</a> and especially the details section of See <a href="#">BEclear</a> for more information about this feature.
colBlockSize	the number of columns that is used in a block if the function is run in parallel mode and/or not on the whole matrix. Set this, and the "rowBlockSize" parameter to 0 if you want to run the function on the whole input matrix. See <a href="#">BEclear</a> and especially the details section of See <a href="#">BEclear</a> for more information about this feature.
epochs	the number of iterations used in the gradient descent algorithm to predict the missing entries in the data matrix. See <a href="#">BEclear</a> for more information.
outputFormat	you can choose if the finally returned data matrix should be saved as an .RData file or as a tab-delimited .txt file in the specified directory. Allowed values are "RData" and "txt". See <a href="#">BEclear</a> for more information.
dir	set the path to a directory the predicted matrix should be stored. The current working directory is defined as default parameter.

**Details**

The function performs the whole process of searching for batch effects and automatically correct them for a matrix of beta values stemming from DNA methylation data. Thereby, the function is running most of the functions from the [BEclear-package](#) in a logical order.

First, median comparison values are calculated by the [calcMedians](#) function, followed by the calculation of p-values by the [calcPvalues](#) function. With the results from the median comparison and p-value calculation, a summary data frame is build using the [calcSummary](#) function, and a scoring table is established by the [calcScore](#) function. Now, found entries from the summary are set to NA in the input matrix using the [clearBEGenes](#) function, then the [BEclear](#) function is used to predict the missing values and at the end, possibly existing wrongly predicted entries (values lower than 0 or greater than 1) are corrected using the [replaceWrongValues](#) function.

**Value**

A list containing the following fields (for detailed information look at the documentations of the corresponding functions):

medians	A data.frame containing all median comparison values corresponding to the input matrix.
pvalues	A data.frame containing all p-values corresponding to the input matrix.
summary	The summarized results of the median comparison and p-value calculation.
score.table	A data.frame containing the number of found genes and a BEscore for every batch.
cleared.data	the input matrix with all values defined in the summary set to NA.
predicted.data	the input matrix after all previously NA values have been predicted.
corrected.predicted.data	the predicted matrix after the correction for wrongly predicted values.

## References

Y. Koren, R. Bell, C. Volinsky, Matrix factorization techniques for recommender systems, IEEE Computer, 42(8), p. 30-37, 2009, <http://research.yahoo.com/pub/2859>

E. Candes, B. Recht, Exact matrix completion via convex optimization, Communications of the ACM, 55(6), p. 111-119, 2012, <http://doi.acm.org/10.114/2184319.2184343>

## See Also

[calcMedians](#) [calcPvalues](#) [calcSummary](#) [calcScore](#) [clearBEgenes](#) [BEclear](#) [replaceWrongValues](#)

## Examples

```
## Shortly running example. For a more realistic example that takes
## some more time, run the same procedure with the full BEclearData
## dataset.

## Whole procedure that has to be done to use this function.
## Correct the example data for a batch effect
data(BEclearData)
ex.data <- ex.data[31:90,7:26]
ex.samples <- ex.samples[7:26,]

# Note that row- and block sizes are just set to 10 to get a short runtime. To
# use these parameters, either use the default values or please note the
# description in the details section of \link{BEclear}
result <- correctBatchEffect(data=ex.data, samples=ex.samples, parallel=TRUE,
  cores=4, adjusted=TRUE, method="fdr", rowBlockSize=10, colBlockSize=10,
  epochs=50, outputFormat="RData", dir=getwd())

# Unlist variables
medians <- result$medians
pvals <- result$pvals
summary <- result$summary
score <- result$score.table
cleared <- result$clearedData
predicted <- result$predictedData
corrected <- result$correctedPredictedData
```

---

countValuesToPredict *Count NA entries in a matrix*

---

### Description

Simple function that counts all values in a matrix which are NA

### Usage

```
countValuesToPredict(data)
```

### Arguments

data            any kind of matrix

### Value

Returns a data frame with the number of NA entries per column. Since the function is mainly written for the usage in batch effect correction of DNA methylation data, the column names of the data frame are set to "sample" and "num\_pred\_values". Nevertheless, the function can be used with any other matrix containing anything but beta values.

### See Also

[clearBEgenes](#) [BEclear](#) [correctBatchEffect](#)

### Examples

```
## Shortly running example. For a more realistic example that takes
## some more time, run the same procedure with the full BEclearData
## dataset.

## Whole procedure that has to be done to use this function.
data(BEclearData)
ex.data <- ex.data[31:90,7:26]
ex.samples <- ex.samples[7:26,]

# Calculate median difference values and p-values
meds <- calcMedians(data=ex.data, samples=ex.samples)
pvals <- calcPvalues(data=ex.data, samples=ex.samples)

# Summarize p-values and median differences for batch affected genes
sum <- calcSummary(medians=meds, pvalues=pvals)
clearedMatrix <- clearBEgenes(data=ex.data, samples=ex.samples, summary=sum)
numberOfEntries <- countValuesToPredict(data=clearedMatrix)
```



---

ex.corrected.data	<i>Example matrix of corrected data for the BEclear-package</i>
-------------------	---

---

**Description**

Example matrix containing a already batch effect corrected sample matrix of beta values from breast invasive carcinoma TCGA methylation data. The matrix contains a small amount of wrongly predicted beta values to show the operating principles of some of the methods from the BEclear package.

**Usage**

```
data(BEclearCorrected)
```

**Format**

A matrix containing already corrected beta values of some samples from the breast invasive carcinoma TCGA methylation data. The colnames denote samples, rownames denote gene names.

**Value**

A matrix (ex.corrected.data)

**Source**

<http://cancergenome.nih.gov/>

---

findWrongValues	<i>Find DNA methylation values out of the boundaries</i>
-----------------	--

---

**Description**

A method which lists values below 0 or beyond 1 contained in the input matrix. The wrong entries are stored in a data.frame together with the corresponding row and column position of the matrix. Note that this method is especially designed for DNA methylation data.

**Usage**

```
findWrongValues(data)
```

**Arguments**

data            any matrix filled with values that normally should be bounded between 0 and 1.

**Details**

Note that this method is especially designed to run after the batch effect correction of DNA methylation data, e.g. with the [BEclear](#) method. It can happen, that the predicted values are lying slightly below the lower bound of 0 or beyond the upper bound of 1. This method finds these inaccurately predicted entries. Another method called [replaceWrongValues](#) replaces these values either by 0 or 1, respectively.

**Value**

A data frame containing the columns "level", "row", "col" and "value" defining if the wrong value is below 0 or beyond 1 (level = 0 or level = 1), the row position and the column position in the input matrix and the wrong value itself, respectively.

**See Also**

[BEclear](#) [replaceWrongValues](#) [correctBatchEffect](#)

**Examples**

```
data(BEclearCorrected)
# Find wrongly predicted values
wrongEntries <- findWrongValues(data=ex.corrected.data)
```

---

makeBoxplot

*produce simple predefined boxplot for methylation data*

---

**Description**

A simple [boxplot](#) is done with boxes either separated by batches or by samples and describe the five number summary of all beta values corresponding to a batch or a sample, respectively. The `batch_ids` are shown on the x-axis with a coloring corresponding to the BEscore.

**Usage**

```
makeBoxplot(data, samples, score, bySamples=FALSE, col="standard", main="",
            xlab="Batch", ylab="Beta value", scoreCol=TRUE)
```

**Arguments**

<code>data</code>	any matrix filled with beta values, column names have to be <code>sample_ids</code> corresponding to the ids listed in "samples", row names have to be gene names.
<code>samples</code>	data frame with two columns, the first column has to contain the sample numbers, the second column has to contain the corresponding batch number. Column names have to be named as "sample_id" and "batch_id".
<code>score</code>	data frame produced by the <a href="#">calcScore</a> function. Contains the number of presumably batch affected genes and a BEscore which is needed for the coloring of the <code>batch_ids</code> .
<code>bySamples</code>	should the boxes be separated by samples or not. If not, boxes are separated by the <code>batch_ids</code> .
<code>col</code>	colors for the boxes, refers to the standard <a href="#">boxplot</a> R-function. If it is set to "standard", boxes are colored batch-wise (if separated by samples) or the standard color "yellow" is used (if separated by batches).
<code>main</code>	main title for the box plot. Default is an empty string.
<code>xlab</code>	label for the x-axis of the box plot. Default is "Batch".
<code>ylab</code>	label for the y-axis of the box plot. Default is "Beta value"
<code>scoreCol</code>	should the <code>batch_ids</code> on the x-axis be colored according to the BEscore or not? If not, black is used as color for all <code>batch_ids</code> .

**Details**

The color code for the batch\_ids on the x-axis provides a simple "traffic light" the user can use to decide if he wants to correct for an assumed batch effect or not. Green means no batch effect, yellow a possibly existing not severe batch effect and red stands for an obviously existing batch effect that should be corrected. The traffic light colors are set according to the BScore from the [calcScore](#) function, values from 0 to 0.02 are colored in green, from 0.02 to 0.1 in yellow and values over 0.1 are colored in red.

**Value**

Returns a boxplot on the graphic device with the features explained above.

**See Also**

[calcScore](#) [boxplot](#) [correctBatchEffect](#)

**Examples**

```
## Shortly running example. For a more realistic example that takes
## some more time, run the same procedure with the full BEclearData
## dataset.

## Whole procedure that has to be done to use this function.
data(BEclearData)
ex.data <- ex.data[31:90,7:26]
ex.samples <- ex.samples[7:26,]

## Prepare the data for the box plots
# Calculate median difference values and p-values
meds <- calcMedians(data=ex.data, samples=ex.samples)
pvals <- calcPvalues(data=ex.data, samples=ex.samples)

# Summarize p-values and median differences for batch affected genes
sum <- calcSummary(medians=meds, pvalues=pvals)

# Calculate the BScore for the batch_id colorings of the x-axis
score <- calcScore(data=ex.data, samples=ex.samples, summary=sum)

## Simple boxplot for the example data separated by batch
makeBoxplot(data=ex.data, samples=ex.samples, score=score, bySamples=FALSE,
  main="Some box plot")

## Simple boxplot for the example data separated by samples
makeBoxplot(data=ex.data, samples=ex.samples, score=score, bySamples=TRUE,
  main="Some box plot")
```

---

replaceWrongValues

*Replace DNA methylation values out of the boundaries*

---

**Description**

A method which replaces values below 0 or beyond 1 contained in the input matrix. These wrong entries are replaced by 0 or 1, respectively. Note that this method is especially designed for DNA methylation data.

**Usage**

```
replaceWrongValues(data)
```

**Arguments**

`data` any matrix filled with values that normally should be bounded between 0 and 1.

**Details**

Note that this method is especially designed to run after the batch effect correction of DNA methylation data, e.g. with the [BEclear](#) method. It can happen, that the predicted values are lying slightly below the lower bound of 0 or beyond the upper bound of 1. This method finds these inaccurately predicted entries. Another method called [replaceWrongValues](#) replaces these values either by 0 or 1, respectively. Another method called [findWrongValues](#) returns a list of existing wrong values and can be run before the replacement.

**Value**

Returns the input matrix with every value previously below 0 changed to 0 and every value previously beyond 1 changed to 1.

**See Also**

[BEclear](#) [findWrongValues](#) [correctBatchEffect](#)

**Examples**

```
data(BEclearCorrected)
# Replace wrongly predicted values
corrected <- replaceWrongValues(data=ex.corrected.data)
```

# Index

- \*Topic **BEclear**
    - BEclear, 4
    - BEclear example data, 7
    - BEclear-package, 2
    - clearBEGenes, 12
    - correctBatchEffect, 13
    - countValuesToPredict, 16
    - ex.corrected.data, 17
    - findWrongValues, 17
    - replaceWrongValues, 19
  - \*Topic **BEscore**
    - calcScore, 10
    - makeBoxplot, 18
  - \*Topic **DNA methylation**
    - BEclear, 4
    - BEclear-package, 2
    - correctBatchEffect, 13
  - \*Topic **NA**
    - countValuesToPredict, 16
  - \*Topic **batch effect correction**
    - BEclear, 4
    - BEclear-package, 2
    - correctBatchEffect, 13
    - findWrongValues, 17
    - replaceWrongValues, 19
  - \*Topic **batch effect**
    - calcScore, 10
    - calcSummary, 11
  - \*Topic **batch**
    - makeBoxplot, 18
  - \*Topic **boxplot**
    - makeBoxplot, 18
  - \*Topic **cleared matrix**
    - clearBEGenes, 12
  - \*Topic **datasets**
    - BEclear example data, 7
    - ex.corrected.data, 17
  - \*Topic **data**
    - BEclear example data, 7
  - \*Topic **ex.data**
    - ex.corrected.data, 17
  - \*Topic **median comparison**
    - calcMedians, 8
    - correctBatchEffect, 13
  - \*Topic **methylation**
    - makeBoxplot, 18
  - \*Topic **p-value adjustment**
    - calcPvalues, 9
  - \*Topic **p-value**
    - calcPvalues, 9
    - correctBatchEffect, 13
  - \*Topic **summary**
    - calcSummary, 11
  - \*Topic **wrong values**
    - findWrongValues, 17
    - replaceWrongValues, 19
- BEclear, 2, 4, 14–18, 20
- BEclear example data, 7
- BEclear-package, 2
- BEclearCorrected (BEclear-package), 2
- BEclearData (BEclear example data), 7
- boxplot, 2, 18, 19
- calcMedians, 2, 8, 11–15
- calcPvalues, 2, 9, 11–15
- calcScore, 2, 10, 14, 15, 18, 19
- calcSummary, 2, 10, 11, 11, 12–15
- clearBEGenes, 2, 12, 14–16
- correctBatchEffect, 2, 6, 8, 9, 11–13, 13, 16, 18–20
- countValuesToPredict, 2, 16
- ex.corrected.data, 17
- ex.data (BEclear example data), 7
- ex.samples (BEclear example data), 7
- findWrongValues, 2, 6, 17, 20
- ks.test, 9
- makeBoxplot, 2, 18
- p.adjust, 9, 14
- replaceWrongValues, 3, 6, 14, 15, 17, 18, 19, 20
- snowfall, 4–6, 8, 9, 14