

# Package ‘uclust’

July 22, 2025

**Title** Clustering and Classification Inference with U-Statistics

**Version** 1.0.0

**Description** Clustering and classification inference for high dimension low sample size (HDLSS) data with U-statistics. The package contains implementations of nonparametric statistical tests for sample homogeneity, group separation, clustering, and classification of multivariate data. The methods have high statistical power and are tailored for data in which the dimension  $L$  is much larger than sample size  $n$ . See Gabriela B. Cybis, Marcio Valk and Sílvia RC Lopes (2018) <[doi:10.1080/00949655.2017.1374387](https://doi.org/10.1080/00949655.2017.1374387)>, Marcio Valk and Gabriela B. Cybis (2020) <[doi:10.1080/10618600.2020.1796398](https://doi.org/10.1080/10618600.2020.1796398)>, Debora Z. Bello, Marcio Valk and Gabriela B. Cybis (2021) <[doi:10.48550/arXiv.2106.09115](https://doi.org/10.48550/arXiv.2106.09115)>.

**Depends** R (>= 3.4.0),dendextend,robcor

**License** GPL-3

**Encoding** UTF-8

**RoxygenNote** 7.1.1

**Suggests** testthat

**NeedsCompilation** no

**Author** Gabriela Cybis [aut, cre],  
Marcio Valk [aut],  
Kazuki Yokoyama [ctb],  
Debora Zava Bello [ctb]

**Maintainer** Gabriela Cybis <[gcybis@gmail.com](mailto:gcybis@gmail.com)>

**Repository** CRAN

**Date/Publication** 2021-06-18 22:10:02 UTC

## Contents

bn	2
bn3	3
is_homo	4
is_homo3	5
plot_uhclust	7

print.utest_classify . . . . .	7
rep_optimBn . . . . .	8
uclust . . . . .	8
uclust3 . . . . .	10
uhclust . . . . .	11
utest . . . . .	13
utest3 . . . . .	14
utest_classify . . . . .	15
var_bn . . . . .	16
<b>Index</b>	<b>18</b>

---

bn	<i>Computes Bn Statistic.</i>
----	-------------------------------

---

## Description

Returns the value for the Bn statistic that measures the degree of separation between two groups. The statistic is computed through the difference of average within group distances to average between group distances. Large values of Bn indicate large group separation. Under overall sample homogeneity we have  $E(Bn)=0$ .

## Usage

```
bn(group_id, md = NULL, data = NULL)
```

## Arguments

group_id	A vector of 0s and 1s indicating to which group the samples belong. Must be in the same order as data or md.
md	Matrix of distances between all data points.
data	Data matrix. Each row represents an observation.

## Details

Either data OR md should be provided. If data are entered directly, Bn will be computed considering the squared Euclidean distance, which is compatible with [is\\_homo](#), [uclust](#) and [uhclust](#).

For more detail see Cybis, Gabriela B., Marcio Valk, and Sílvia RC Lopes. "Clustering and classification problems in genetics through U-statistics." *Journal of Statistical Computation and Simulation* 88.10 (2018) and Valk, Marcio, and Gabriela Bettella Cybis. "U-statistical inference for hierarchical clustering." *Journal of Computational and Graphical Statistics* 30(1) (2021).

## Value

Value of the Bn statistic.

**Examples**

```
n=5
x=matrix(rnorm(n*10),ncol=10)
bn(c(1,0,0,0,0),data=x)      # option (a) entering the data matrix directly
md=as.matrix(dist(x))^2
bn(c(0,1,1,1,1),md)        # option (b) entering the distance matrix
```

bn3

*Computes Bn Statistic for 3 Groups.***Description**

Returns the value for the Bn statistic that measures the degree of separation between three groups. The statistic is computed as a combination of differences of average within group and between group distances. Large values of Bn indicate large group separation. Under overall sample homogeneity we have  $E(Bn)=0$ .

**Usage**

```
bn3(group_id, md = NULL, data = NULL)
```

**Arguments**

group_id	A vector of 1s, 2s and 3s indicating to which group the samples belong. Must be in the same order as data or md.
md	Matrix of distances between all data points.
data	Data matrix. Each row represents an observation.

**Details**

Either data OR md should be provided. If data are entered directly, Bn will be computed considering the squared Euclidean distance.

For more detail see Bello, Debora Zava, Marcio Valk and Gabriela Bettella Cybis. "Clustering inference in multiple groups." arXiv preprint arXiv:2106.09115 (2021).

**Value**

Value of the Bn3 statistic.

**Examples**

```
n=7
set.seed(1234)
x=matrix(rnorm(n*10),ncol=10)
bn3(c(1,2,2,2,3,3,3),data=x)      # option (a) entering the data matrix directly
md=as.matrix(dist(x))^2
bn3(c(1,2,2,2,3,3,3),md)        # option (b) entering the distance matrix
```

---

 is\_homo

*U-statistic based homogeneity test*


---

### Description

Homogeneity test based on the statistic  $b_n$ . The test assesses whether there exists a data partition for which group separation is statistically significant according to the U-test. The null hypothesis is overall sample homogeneity, and a sample is considered homogeneous if it cannot be divided into two statistically significant subgroups.

### Usage

```
is_homo(md = NULL, data = NULL, rep = 10)
```

### Arguments

md	Matrix of distances between all data points.
data	Data matrix. Each row represents an observation.
rep	Number of times to repeat optimization procedure. Important for problems with multiple optima.

### Details

This is the homogeneity test of Cybis et al. (2017) extended to account for groups of size 1. The test is performed through two steps: an optimization procedure that finds the data partition that maximizes the standardized  $B_n$  and a test for the resulting maximal partition. Should be used in high dimension small sample size settings.

Either data or md should be provided. If data are entered directly,  $B_n$  will be computed considering the squared Euclidean distance.

Variance of  $b_n$  is estimated through resampling, and thus, p-values may vary a bit in different runs.

For more detail see Cybis, Gabriela B., Marcio Valk, and Sílvia RC Lopes. "Clustering and classification problems in genetics through U-statistics." *Journal of Statistical Computation and Simulation* 88.10 (2018) and Valk, Marcio, and Gabriela Bettella Cybis. "U-statistical inference for hierarchical clustering." *Journal of Computational and Graphical Statistics* 30(1) (2021).

### Value

Returns a list with the following elements:

**minFobj** Test statistic. Minimum of the objective function for optimization ( $-\text{std}B_n$ ).

**group1** Elements in group 1 in the maximal partition. (obs: this is not the best partition for the data, see uclust)

**group2** Elements in group 2 in the maximal partition.

**p.MaxTest** P-value for the homogeneity test.

**Rep.Fobj** Values for the minimum objective function on all rep optimization runs.

**bootB** Resampling variance estimate for partitions with groups of size  $n/2$  (or  $(n-1)/2$  and  $(n+1)/2$  if  $n$  is odd).

**bootB1** Resampling variance estimate for partitions with one group of size 1.

### Examples

```
x = matrix(rnorm(500000),nrow=50) #creating homogeneous Gaussian dataset
res = is_homo(data=x)

x[1:30,] = x[1:30,]+0.15 #Heterogeneous dataset (first 30 samples have different mean)
res = is_homo(data=x)

md = as.matrix(dist(x)^2) #squared Euclidean distances for the same data
res = is_homo(md)

# Multidimensional scaling plot of distance matrix
fit <- cmdscale(md, eig = TRUE, k = 2)
x <- fit$points[, 1]
y <- fit$points[, 2]
plot(x,y, main=paste("Homogeneity test: p-value =",res$p.MaxTest))
```

---

is\_homo3

*U-statistic based homogeneity test for 3 groups*


---

### Description

Homogeneity test based on the statistic  $bn_3$ . The test assesses whether there exists a data partition for which three group separation is statistically significant according to  $utest_3$ . The null hypothesis is overall sample homogeneity, and a sample is considered homogeneous if it cannot be divided into three groups with at least one significantly different from the others.

### Usage

```
is_homo3(md = NULL, data = NULL, rep = 20, test_max = TRUE, alpha = 0.05)
```

### Arguments

md	Matrix of distances between all data points.
data	Data matrix. Each row represents an observation.
rep	Number of times to repeat optimization procedure. Important for problems with multiple optima.
test_max	Logical indicating whether to employ the max test
alpha	Significance level

## Details

This is the homogeneity test of Bello et al. (2021). The test is performed through two steps: an optimization procedure that finds the data partition that maximizes the standardized Bn and a test for the resulting maximal partition. Should be used in high dimension small sample size settings.

Either data or md should be provided. If data are entered directly, Bn will be computed considering the squared Euclidean distance.

Variance of bn is estimated through resampling, and thus, p-values may vary a bit in different runs.

For more detail see Bello, Debora Zava, Marcio Valk and Gabriela Bettella Cybis. "Clustering inference in multiple groups." arXiv preprint arXiv:2106.09115 (2021).

## Value

Returns a list with the following elements:

**stdBn** Test statistic. Maximum standardized Bn.

**group1** Elements in group 1 in the maximal partition. (obs: this is not the best partition for the data, see uclust3)

**group2** Elements in group 2 in the maximal partition.

**group3** Elements in group 3 in the maximal partition.

**pvalue.Bonferroni** P-value for the homogeneity test.

**alpha\_Bonferroni** Alpha after Bonferroni correction

**bootB** Resampling variance estimate for partitions with central group sizes.

**bootB1** Resampling variance estimate for partitions with one group of size 1.

**varBn** Estimated variance of Bn for maximal standardized Bn configuration.

## Examples

```
set.seed(123)
x = matrix(rnorm(70000),nrow=7) #creating homogeneous Gaussian dataset
res = is_homo3(data=x)
res

#uncomment to run
# x = matrix(rnorm(18000),nrow=18)
# x[1:5,] = x[1:5,]+0.5 #Heterogeneous dataset (first 5 samples have different mean)
# x[6:9,] = x[6:9,]+1.5
# res = is_homo3(data=x)
# res
# md = as.matrix(dist(x)^2) #squared Euclidean distances for the same data
# res = is_homo3(md) # uncomment to run

# Multidimensional scaling plot of distance matrix
#fit <- cmdscale(md, eig = TRUE, k = 2)
#x <- fit$points[, 1]
#y <- fit$points[, 2]
#plot(x,y, main=paste("Homogeneity test: p-value =",res$p.MaxTest))
```

---

plot_uhclust	<i>Plot function for the result of uhclust</i>
--------------	--

---

**Description**

This function plots the p-value annotated dendrogram resulting from uhclust

**Usage**

```
plot_uhclust(  
  uhclust,  
  pvalues_cex = 0.8,  
  pvalues_dx = 2,  
  pvalues_dy = 0.08,  
  print_pvalues = TRUE  
)
```

**Arguments**

uhclust	Result from uhclust
pvalues_cex	Graphical parameter for p-value font size.
pvalues_dx	Graphical parameter for p-value position shift on x axis.
pvalues_dy	Graphical parameter for p-value position shift on y axis.
print_pvalues	Logical. Should the p-values be printed?

**Examples**

```
x = matrix(rnorm(100000),nrow=50)  
x[1:35,] = x[1:35,]+0.7  
x[1:15,] = x[1:15,]+0.4  
res = uhclust(data=x, plot=FALSE)  
plot_uhclust(res)
```

---

print.utest_classify	<i>Simple print method for utest_classify objects.</i>
----------------------	--

---

**Description**

Simple print method for utest\_classify objects.

**Usage**

```
## S3 method for class 'utest_classify'  
print(x, ...)
```

**Arguments**

x	utest_classify object
...	additional parameters passed to the function

---

rep_optimBn	<i>Optimization function with multiple starting points (for local optima)</i>
-------------	---

---

**Description**

Finds the configuration with max Bn among all configurations.

**Usage**

```
rep_optimBn(mdm, rep = 15, bootB = -1)
```

**Arguments**

mdm	Matrix of squared Euclidean distances between all data points.
rep	Number of replications
bootB	Result of previous bootstrap (if available). If, -1, a new bootstrap is performed for the variance of Bn.

---

uclust	<i>U-statistic based significance clustering</i>
--------	--

---

**Description**

Partitions the sample into the two significant subgroups with the largest Bn statistic. If no significant partition exists, the test will return "homogeneous".

**Usage**

```
uclust(md = NULL, data = NULL, alpha = 0.05, rep = 15)
```

**Arguments**

md	Matrix of distances between all data points.
data	Data matrix. Each row represents an observation.
alpha	Significance level.
rep	Number of times to repeat optimization procedures. Important for problems with multiple optima.



## Details

This is the significance clustering procedure of Valk and Cybis (2018). The method first performs a homogeneity test to verify whether the data can be significantly partitioned. If the hypothesis of homogeneity is rejected, then the method will search, among all the significant partitions, for the partition that better separates the data, as measured by larger  $b_n$  statistic. This function should be used in high dimension small sample size settings.

Either data or md should be provided. If data are entered directly,  $B_n$  will be computed considering the squared Euclidean distance.

Variance of  $b_n$  is estimated through resampling, and thus, p-values may vary a bit in different runs.

For more detail see Cybis, Gabriela B., Marcio Valk, and Sílvia RC Lopes. "Clustering and classification problems in genetics through U-statistics." *Journal of Statistical Computation and Simulation* 88.10 (2018) and Valk, Marcio, and Gabriela Bettella Cybis. "U-statistical inference for hierarchical clustering." *Journal of Computational and Graphical Statistics* 30(1) (2021). See also `is_homo`, `uhclust`, `Utest_class`.

## Value

Returns a list with the following elements:

**cluster1** Elements in group 1 in the final partition. This is the significant partition with maximal  $B_n$ , if sample is heterogeneous.

**cluster2** Elements in group 2 in the final partition.

**p.value** P-value for the test that renders the final partition, if heterogeneous. Homogeneity test p-value, if homogeneous.

**alpha\_corrected** Bonferroni corrected significance level for the test that renders the final partition, if heterogeneous. Homogeneity test significance level, if homogeneous.

**n1** Size of the smallest cluster

**ishomo** Logical, returns TRUE when the sample is homogeneous.

**Bn** Value of  $B_n$  statistic for the final partition, if heterogeneous. Value of  $B_n$  statistic for the maximal homogeneity test partition, if homogeneous.

**varBn** Variance estimate for final partition, if heterogeneous. Variance estimate for the maximal homogeneity test partition, if homogeneous.

**ishomoResult** Result of homogeneity test (see `is_homo`).

## Examples

```
set.seed(17161)
x = matrix(rnorm(100000),nrow=50) #creating homogeneous Gaussian dataset
res = uclust(data=x)

x[1:30,] = x[1:30,]+0.25 #Heterogeneous dataset (first 30 samples have different mean)
res = uclust(data=x)

md = as.matrix(dist(x)^2) #squared Euclidean distances for the same data
res = uclust(md)
```

```
# Multidimensional scaling plot of distance matrix
fit <- cmdscale(md, eig = TRUE, k = 2)
x <- fit$points[, 1]
y <- fit$points[, 2]
col=rep(3,dim(md)[1])
col[res$cluster2]=2
plot(x,y, main=paste("Multidimensional scaling plot of data:
                    homogeneity p-value =",res$ishomoResult$p.MaxTest),col=col)
```

---

uclust3

*U-statistic based significance clustering for three way partitions*


---

### Description

Partitions data into three groups only when these partitions are statistically significant. If no significant partition exists, the test will return "homogeneous".

### Usage

```
uclust3(md = NULL, data = NULL, alpha = 0.05, rep = 15)
```

### Arguments

md	Matrix of distances between all data points.
data	Data matrix. Each row represents an observation.
alpha	Significance level.
rep	Number of times to repeat optimization procedures. Important for problems with multiple optima.

### Details

This is the significance clustering procedure of Bello et al. (2021). The method first performs a homogeneity test to verify whether the data can be significantly partitioned. If the hypothesis of homogeneity is rejected, then the method will search, among all the significant partitions, for the partition that better separates the data, as measured by larger bn statistic. This function should be used in high dimension small sample size settings.

Either data or md should be provided. If data are entered directly, Bn will be computed considering the squared Euclidean distance.

Variance of bn is estimated through resampling, and thus, p-values may vary a bit in different runs.

For more detail see Bello, Debora Zava, Marcio Valk and Gabriela Bettella Cybis. "Clustering inference in multiple groups." arXiv preprint arXiv:2106.09115 (2021). See also is\_homo3, uclust.

**Value**

Returns a list with the following elements:

**groups** List with elements of final three groups

**p.value** P-value for the test that renders the final partition, if heterogeneous. Homogeneity test p-value, if homogeneous.

**alpha\_corrected** Bonferroni corrected significance level for the test that renders the final partition, if heterogeneous. Homogeneity test significance level, if homogeneous.

**ishomo** Logical, returns TRUE when the sample is homogeneous.

**Bn** Value of Bn statistic for the final partition, if heterogeneous. Value of Bn statistic for the maximal homogeneity test partition, if homogeneous.

**varBn** Variance estimate for final partition, if heterogeneous. Variance estimate for the maximal homogeneity test partition, if homogeneous.

**Examples**

```
set.seed(123)
x = matrix(rnorm(70000),nrow=7) #creating homogeneous Gaussian dataset
res = uclust3(data=x)
res

# uncomment to run
# x = matrix(rnorm(15000),nrow=15)
# x[1:6,] = x[1:6,]+1.5 #Heterogeneous dataset (first 5 samples have different mean)
# x[7:12,] = x[7:12,]+3
# res = uclust3(data=x)
# res$groups
```

---

uhclust

*U-statistic based significance hierarchical clustering*


---

**Description**

Hierarchical clustering method that partitions the data only when these partitions are statistically significant.

**Usage**

```
uhclust(md = NULL, data = NULL, alpha = 0.05, rep = 15, plot = TRUE)
```

### Arguments

<code>md</code>	Matrix of distances between all data points.
<code>data</code>	Data matrix. Each row represents an observation.
<code>alpha</code>	Significance level.
<code>rep</code>	Number of times to repeat optimization procedures. Important for problems with multiple optima.
<code>plot</code>	Logical, TRUE if p-value annotated dendrogram should be plotted.

### Details

This is the significance hierarchical clustering procedure of Valk and Cybis (2018). The data are repeatedly partitioned into two subgroups, through function `uclust`, according to a hierarchical scheme. The procedure stops when resulting subgroups are homogeneous or have fewer than 3 elements. This function should be used in high dimension small sample size settings.

Either `data` or `md` should be provided. If data are entered directly, `Bn` will be computed considering the squared Euclidean distance.

Variance of `bn` is estimated through resampling, and thus, p-values may vary a bit in different runs.

For more detail see Cybis, Gabriela B., Marcio Valk, and Sílvia RC Lopes. "Clustering and classification problems in genetics through U-statistics." *Journal of Statistical Computation and Simulation* 88.10 (2018) and Valk, Marcio, and Gabriela Bettella Cybis. "U-statistical inference for hierarchical clustering." *Journal of Computational and Graphical Statistics* 30(1) (2021).

See also `is_homo`, `uclust` and `Utest_class`.

### Value

Returns an object of class `hclust` with three additional attribute arrays:

**Pvalues** P-values from `uclust` for the final data partition at each node of the dendrogram. This array is in the same order of height, and only contains values for tests that were performed.

**alpha** Bonferroni corrected significance levels for `uclust` for the data partitions at each node of the dendrogram. This array is in the same order of height, and only contains values for tests that were performed.

**groups** Final group assignments.

### Examples

```
x = matrix(rnorm(100000),nrow=50) #creating homogeneous Gaussian dataset
res = uhclust(data=x)
```

```
x[1:30,] = x[1:30,]+0.7 #Heterogeneous dataset
x[1:10,] = x[1:10,]+0.4
res = uhclust(data=x)
res$groups
```

---

utest	<i>U test</i>
-------	---------------

---

### Description

Test for the separation of two groups. The null hypothesis states that the groups are homogeneous and the alternative hypothesis states that they are separate.

### Usage

```
utest(group_id, md = NULL, data = NULL, numB = 1000)
```

### Arguments

group_id	A vector of 0s and 1s indicating to which group the samples belong. Must be in the same order as data or md.
md	Matrix of distances between all data points.
data	Data matrix. Each row represents an observation.
numB	Number of resampling iterations.

### Details

Either data or md should be provided. If data are entered directly, Bn will be computed considering the squared Euclidean distance, which is compatible with [is\\_homo](#), [uclust](#) and [uhclust](#).

For more details see Cybis, Gabriela B., Marcio Valk, and Sílvia RC Lopes. "Clustering and classification problems in genetics through U-statistics." *Journal of Statistical Computation and Simulation* 88.10 (2018)

### Value

Returns a list with the following elements:

**Bn** Test Statistic

**Pvalue** Replication based p-value

**Replication** Number of replications used to compute p-value

### See Also

[bn](#), [is\\_homo](#)

## Examples

```
# Simulate a dataset with two separate groups, the first 5 rows have mean 0 and
# the last 5 rows have mean 5.
data <- matrix(c(rnorm(75, 0), rnorm(75, 5)), nrow = 10, byrow=TRUE)

# U test for mixed up groups
utest(group_id=c(1,0,1,0,1,0,1,0,1,0), data=data, numB=3000)
# U test for correct group definitions
utest(group_id=c(1,1,1,1,1,0,0,0,0,0), data=data, numB=3000)
```

---

utest3

*U-test for three groups*


---

## Description

Test for the separation of three groups. The null hypothesis states that the groups are homogeneous and the alternative hypothesis states that at least one is separated from the others.

## Usage

```
utest3(group_id, md = NULL, data = NULL, alpha = 0.05, numB = 1000)
```

## Arguments

group_id	A vector of 1s, 2s and 3s indicating to which group the samples belong. Must be in the same order as data or md.
md	Matrix of distances between all data points.
data	Data matrix. Each row represents an observation.
alpha	Significance level
numB	Number of resampling iterations.

## Details

Either data or md should be provided. If data are entered directly, Bn will be computed considering the squared Euclidean distance.

For more detail see Bello, Debora Zava, Marcio Valk and Gabriela Bettella Cybis. "Clustering inference in multiple groups." arXiv preprint arXiv:2106.09115 (2021).

## Value

Returns a list with the following elements:

**is.homo** Logical of whether test indicates that data is homogeneous

**Pvalue** Replication based p-value

**Bn** Test Statistic

**sdBn** Standard error for Bn statistic computed through resampling

**See Also**

[bn3,utest.is\\_homo3](#)

**Examples**

```
# Simulate a dataset with two separate groups,
# the first row has mean -4, the next 5 rows have mean 0 and the last 5 rows have mean 4.
data <- matrix(c(rnorm(15, -4),rnorm(75, 0), rnorm(75, 4)), nrow = 11, byrow=TRUE)
# U test for mixed up groups
utest3(group_id=c(1,2,3,1,2,3,1,2,3,1,2), data=data, numB=3000)
# U test for correct group definitions
utest3(group_id=c(1,2,2,2,2,2,3,3,3,3,3), data=data, numB=3000)
```

---

utest_classify	<i>Test for classification of a sample in one of two groups.</i>
----------------	--

---

**Description**

The null hypothesis is that the new data is not well classified into the first group when compared to the second group. The alternative hypothesis is that the data is well classified into the first group.

**Usage**

```
utest_classify(x, data, group_id, bootstrap_iter = 1000)
```

**Arguments**

- x                    A numeric vector to be classified.
- data                Data matrix. Each row represents an observation.
- group\_id            A vector of 0s (first group) and 1s indicating to which group the samples belong. Must be in the same order as data.
- bootstrap\_iter    Numeric scalar. The number of bootstraps. It's recommended  $1000 < bootstrap\_iter < 10000$ .

**Details**

The test is performed considering the squared Euclidean distance.

For more detail see Cybis, Gabriela B., Marcio Valk, and Sílvia RC Lopes. "Clustering and classification problems in genetics through U-statistics." *Journal of Statistical Computation and Simulation* 88.10 (2018) and Valk, Marcio, and Gabriela Bettella Cybis. "U-statistical inference for hierarchical clustering." arXiv preprint arXiv:1805.12179 (2018).

**Value**

A list with class "utest\_classify" containing the following components:

statistic        the value of the test statistic.  
 p\_value         The p-value for the test.  
 bootstrap\_iter  the number of bootstrap iterations.

**Examples**

```
# Example 1
# Five observations from each group, G1 and G2. Each observation has 60 dimensions.
data <- matrix(c(rnorm(300, 0), rnorm(300, 10)), ncol = 60, byrow=TRUE)
# Test data comes from G1.
x <- rnorm(60, 0)
# The test correctly indicates that the test data should be classified into G1 (p < 0.05).
utest_classify(x, data, group_id = c(rep(0,times=5),rep(1,times=5)))

# Example 2
# Five observations from each group, G1 and G2. Each observation has 60 dimensions.
data <- matrix(c(rnorm(300, 0), rnorm(300, 10)), ncol = 60, byrow=TRUE)
# Test data comes from G2.
x <- rnorm(60, 10)
# The test correctly indicates that the test data should be classified into G2 (p > 0.05).
utest_classify(x, data, group_id = c(rep(1,times=5),rep(0,times=5)))
```

var\_bn

*Variance of Bn***Description**

Estimates the variance of the Bn statistic using the resampling procedure described in Cybis, Gabriela B., Marcio Valk, and Sílvia RC Lopes. "Clustering and classification problems in genetics through U-statistics." *Journal of Statistical Computation and Simulation* 88.10 (2018) and Valk, Marcio, and Gabriela Bettella Cybis. "U-statistical inference for hierarchical clustering." *Journal of Computational and Graphical Statistics* 30(1) (2021).

**Usage**

```
var_bn(group_sizes, md = NULL, data = NULL, numB = 2000)
```

**Arguments**

group\_sizes     A vector with two entries: size of group 1 and size of group 2.  
 md              Matrix of distances between all data points.  
 data            Data matrix. Each row represents an observation.  
 numB            Number of resampling iterations. Only used if no groups are of size 1.



**Details**

Either data or md should be provided. If data are entered directly, Bn will be computed considering the squared Euclidean distance, which is compatible with [is\\_homo](#), [uclust](#) and [uhclust](#).

**Value**

Variance of Bn

**See Also**

[bn](#)

**Examples**

```
n=5
x=matrix(rnorm(n*20),ncol=20)
# option (a) entering the data matrix directly and considering a group of size 1
var_bn(c(1,4),data=x)

# option (b) entering the distance matrix and considering a groups of size 2 and 3
md=as.matrix(dist(x))^2
var_bn(c(2,3),md)
```

# Index

`bn`, [2](#), [13](#), [17](#)

`bn3`, [3](#), [15](#)

`is_homo`, [2](#), [4](#), [13](#), [17](#)

`is_homo3`, [5](#), [15](#)

`plot_uhclust`, [7](#)

`print.utest_classify`, [7](#)

`rep_optimBn`, [8](#)

`uclust`, [2](#), [8](#), [13](#), [17](#)

`uclust3`, [10](#)

`uhclust`, [2](#), [11](#), [13](#), [17](#)

`utest`, [13](#), [15](#)

`utest3`, [14](#)

`utest_classify`, [15](#)

`var_bn`, [16](#)