

Package ‘EXPERT’

December 6, 2010

Title A package for Extremely small P-value Evaluation for Resampling-based Tests

Version 0.0.1

Date 2010-12-06

Author Kai Yu and Faming Liang

Description A package for efficient p-value evaluation for resampling-based tests

Maintainer Kai Yu <yuka@mail.nih.gov>

Depends R (>= 2.10.1)

License GPL-2

R topics documented:

EXPERT	1
MAX.test.statistic	2
proposal.permute.matrix	3
proposal.permute.vector	4
SAMC.adapt	5
SAMC.fun	7
t.test.statistic	8
Index	9

EXPERT	<i>A package for Extremely small P-value Evaluation for Resampling-based Tests</i>
--------	--

Description

An R package for efficient p-value estimation for resampling-based tests.

Details

The main function in this package is [SAMC.adapt](#) which uses the SAMC (Stochastic Approximation Monte Carlo) method to estimate the p-value for an observed test statistic. In order to use this function, the user must provide two functions: one function for calculating the test statistic, and another function for updating the input data.

Author(s)

Kai Yu <yuka@mail.nih.gov> and Faming Liang

References

Yu K, Liang F, Ciampa J, Chatterjee N. Efficient p-value evaluation for resampling-based tests (Submitted manuscript).

MAX.test.statistic *Calculate the maximum value of multiple Cochran-Armitage trend test statistics*

Description

A function to calculate the MAX test statistic, i.e., the maximum value of multiple Cochran-Armitage trend test statistics performed on a set genetic markers (e.g. SNPs).

Usage

```
MAX.test.statistic(data.input)
```

Arguments

`data.input` A list of two components, `x` and `y`, which are matrices containing genotypes on a set of genetic markers from cases and controls, where the rows are for subjects and columns for the markers. Each element takes the value 0, 1, or 2, representing the number of copies of a certain allele in the genotype.

Value

The value of the test statistic

See Also

[proposal.permute.matrix](#)

Examples

```
MAX.test.statistic
```

`proposal.permute.matrix`

Example function for the fun.proposal argument of the SAMC.adapt() function

Description

This function updates the input data by randomly switching any chosen pair of rows from two matrices.

Usage

```
proposal.permute.matrix(data.input, prop.change)
```

Arguments

`data.input` A list of two components, `x` and `y`, which are matrices containing genotypes on a set of genetic markers from cases and controls, where the rows are for subjects and columns for the markers. Each element takes the value 0, 1, or 2, representing the number of copies of a certain allele in the genotype.

`prop.change` The proportion of the rows in the smaller matrix to be switched with chosen rows from the other matrix.

Details

This function can be used together with [MAX.test.statistic](#). See the reference on exactly how the permutations are conducted.

Value

The updated version of `data.input`.

References

Yu K, Liang F, Ciampa J, Chatterjee N. Efficient p-value evaluation for resampling-based tests (Submitted manuscript).

See Also

[MAX.test.statistic](#), [SAMC.adapt](#)

Examples

```
proposal.permute.matrix
```

```
proposal.permute.vector
```

Example function for the fun.proposal argument of the SAMC.adapt() function

Description

This function updates the input data by randomly switching any chosen pair of elements from two vectors.

Usage

```
proposal.permute.vector(data.input, prop.change)
```

Arguments

`data.input` The example input data object
`prop.change` The proportion of data to be updated

Details

This function can be used together with `t.test.statistic`. See the reference on exactly how the permutations are conducted.

Value

The updated version of `data.input`.

References

Yu K, Liang F, Ciampa J, Chatterjee N. Efficient p-value evaluation for resampling-based tests (Submitted manuscript).

See Also

`t.test.statistic`, `SAMC.adapt`

Examples

```
proposal.permute.vector
```

SAMC.adapt	<i>SAMC (Stochastic Approximation Monte Carlo) method for small p-value estimation</i>
------------	--

Description

To estimate the p-value for an observed test statistic based on the SAMC (Stochastic Approximation Monte Carlo) method.

Usage

```
SAMC.adapt(data.input, t.obs, t.start=0, n.iter.1=200000, n.iter.2=1000000,
           n.region.1=101, n.region.2=301, prop.change=0.05, gain.factor.t0=1000,
           fun.test.statistic=NULL, fun.proposal=NULL, ...)
```

Arguments

<code>data.input</code>	The input data, which can be any data object, e.g. data frame, list, etc., as long as it is compatible with the two user provided functions <code>fun.test.statistic</code> and <code>fun.proposal</code> .
<code>t.obs</code>	The threshold at which the tail probability is to be calculated. For example, this value can be set to the observed test statistic. The estimated p-value for this value will be computed.
<code>t.start</code>	The minimum value for the considered test statistic. The default is 0.
<code>n.iter.1</code>	The number of iterations used for the initial run with subregions defined uniformly in the interval $[t.start, t.obs]$, with final region $[t.obs, \infty)$. The default is 200000.
<code>n.iter.2</code>	The number of iterations used for the final run with subregions defined according to the results from the initial run. If <code>n.iter.2 = 0</code> , then only the initial run is performed. The default is 1000000.
<code>n.region.1</code>	The number of subregions used by SAMC in the initial run. The default is 101.
<code>n.region.2</code>	The number of subregions used by SAMC in the final run. The default is 301.
<code>prop.change</code>	The proportion of data to be updated by <code>fun.proposal</code> in each iteration. The default is 0.05.
<code>gain.factor.t0</code>	This is used for defining the gain factor sequence. The default is 1000.
<code>fun.test.statistic</code>	A function written by the user to calculate and return the test statistic based on <code>data.input</code> . It has the following format: <code>fun.test.statistic(data.input, ...)</code> . See the example function <code>t.test.statistic</code> .
<code>fun.proposal</code>	A function written by the user to update the input data and return the updated data. It should be in the following format: <code>fun.proposal(data.input, prop.change, ...)</code> . See the example function <code>proposal.permute.vector</code> .
<code>...</code>	Other arguments for <code>fun.test.statistic</code> and/or <code>fun.proposal</code>

Value

A list with two components `hist.mat` and `p.value`. The object `hist.mat` will be a data frame of dimension `n.region.2` by 4 (or `n.region.1` by 4 if `n.iter.2 = 0`), which contains information used for checking convergence. The other component of each sublist is `p.value`, which is the estimated p-value. The `samplefreq` column of `hist.mat` provides the number of times the SAMC iterations generate a sample falling into each individual sub-region. A barplot of `samplefreq` (`barplot(hist.mat[, "samplefreq"])`) should show a flat pattern.

Author(s)

Kai Yu

References

Yu K, Liang F, Ciampa J, Chatterjee N. Efficient p-value evaluation for resampling-based tests (Submitted manuscript).

See Also

[t.test.statistic](#), [proposal.permute.vector](#), [SAMC.fun](#)

Examples

```
# Example 1: An example of two-sample t test.
# Suppose we do not know the theoretical distribution of t test statistic,
# but instead want to rely on a permutation procedure to evaluate the significance
# level for a given observed statistic value. In this permutation procedure, we
# randomly shuffle the group IDs among the two samples and generate new versions of
# the data. We can apply the SAMC algorithm to estimate the p-value (tail probability)
# at a much reduced number of iterations than the standard permutation procedure.
# First, we need to define the two functions for running SAMC, fun.test.statistic,
# and fun.proposal. In the SPERT package, we provide such two functions. They are
# t.test.statistic, which calculates the two sided t-test statistic, and
# proposal.permute.vector, which update the data through permutation.

set.seed(1)
x<-rnorm(200, mean=0, sd=1)
y<-rnorm(200, mean=0.5, sd=1)
data.input<-list(x=x, y=y)
t.obs<-t.test.statistic(data.input)

# True p-value
p.true <- 2.0*(1-pt(t.obs, df=398))
p.true

# Run SAMC to estimate the p-value
res<-SAMC.adapt(data.input, t.obs, t.start=0, n.iter.1=5000,
               n.iter.2=10000, n.region.1=101, n.region.2=201,
               prop.change=0.05, gain.factor.t0=1000,
               fun.test.statistic=t.test.statistic, fun.proposal=proposal.permute.vector)

# The estimated p-value
res$p.value

# To check for convergence:
barplot(res$hist.mat[, "samplefreq"])
```

```

# Note that the barplot is not flat, so the algorithm has not converged yet.
# To get a more precise estimate run the code below.
# WARNING: the following call to SAMC.adapt could take around 20 minutes to run.
ptm <- proc.time()
#res<-SAMC.adapt(data.input, t.obs, t.start=0, n.iter.1=200000,
#               n.iter.2=1000000, n.region.1=101, n.region.2=301,
#               prop.change=0.05, gain.factor.t0=1000,
#               fun.test.statistic=t.test.statistic, fun.proposal=proposal.permute.vector)
proc.time()-ptm

```

SAMC.fun

SAMC algorithm

Description

This function is to run the SAMC algorithm with user defined subregion definition.

Usage

```

SAMC.fun(data.input, n.iter, partition.vec, prop.change, theta.start,
         gain.factor.t0, fun.test.statistic, fun.proposal, ...)

```

Arguments

<code>data.input</code>	The input data, which can be any data object, e.g. data frame, list, etc., as long as it is compatible with the two user provided functions <code>fun.test.statistic</code> and <code>fun.proposal</code> .
<code>n.iter</code>	The number of iterations
<code>partition.vec</code>	A vector of length <code>m</code> of cut points at the domain of the test statistic. The first subregion is defined as <code>[partition.vec[1], partition.vec[2]]</code> , the second subregion is <code>[partition.vec[2], partition.vec[3]]</code> , ..., with final subregion being <code>[partition.vec[m], ∞)</code> . Usually <code>partition.vec[1]</code> is the minimum value of the test statistic (e.g. 0), and <code>partition.vec[m]</code> can be set as the observed test statistic whose p-value is to be estimated.
<code>prop.change</code>	The percentage of data to be updated by <code>fun.proposal</code> in each iteration.
<code>theta.start</code>	<code>theta.start</code>
<code>gain.factor.t0</code>	This is used for defining the gain factor sequence.
<code>fun.test.statistic</code>	A function written by the user to calculate and return the test statistic based on <code>data.input</code> . It has the following format: <code>fun.test.statistic(data.input, ...)</code> . See the example function <code>t.test.statistic</code> .
<code>fun.proposal</code>	A function written by the user to update the input data and return the updated data. It should be in the following format: <code>fun.proposal(data.input, prop.change, ...)</code> . See the example function <code>proposal.permute.vector</code> .
<code>...</code>	Other arguments for <code>fun.test.statistic</code> and/or <code>fun.proposal</code>

Details

This function estimates the tail probability

$$P(\lambda(x) \geq \text{partition.vec}[m])$$

Value

A list with two components `hist.mat` and `p.value` (see [SAMC.adapt](#)).

Author(s)

Kai Yu

References

Yu K, Liang F, Ciampa J, Chatterjee N. Efficient p-value evaluation for resampling-based tests (Submitted manuscript).

See Also

[SAMC.adapt](#), [t.test.statistic](#), [proposal.permute.vector](#)

`t.test.statistic` *Calculate a 2 sample t-test statistic*

Description

A function to calculate a 2 sample t-test statistic which can be used as an example function for the `fun.test.statistic` argument of the `SAMC.adapt` function.

Usage

```
t.test.statistic(data.input)
```

Arguments

`data.input` A list of two components, `x` and `y`, which are vectors containing measurements from one of the two samples.

Value

The value of the test statistic

See Also

[proposal.permute.vector](#)

Examples

```
t.test.statistic
```


Index

*Topic **htest**

SAMC.adapt, 4

SAMC.fun, 6

*Topic **package**

EXPERT, 1

EXPERT, 1

MAX.test.statistic, 2, 3

proposal.permute.matrix, 2, 2

proposal.permute.vector, 3, 5, 7, 8

SAMC.adapt, 1, 3, 4, 4, 7

SAMC.fun, 5, 6

t.test.statistic, 4, 5, 7, 7