

Package ‘SuperExactTest’

June 15, 2017

Type Package

Title Exact Test and Visualization of Multi-Set Intersections

Version 0.99.4

Date 2017-06-14

Author Minghui Wang, Yongzhong Zhao and Bin Zhang

Maintainer Minghui Wang <minghui.wang@mssm.edu>

Contact Minghui Wang <minghui.wang@mssm.edu>, Bin Zhang
<bin.zhang@mssm.edu>

Description

Identification of sets of objects with shared features is a common operation in all disciplines. Analysis of intersections among multiple sets is fundamental for in-depth understanding of their complex relationships. This package implements a theoretical framework for efficient computation of statistical distributions of multi-set intersections based upon combinatorial theory, and provides multiple scalable techniques for visualizing the intersection statistics. The statistical algorithm behind this package was published in Wang et al. (2015) <doi:10.1038/srep16923>.

License GPL-3

Depends grid (>= 3.1.0), methods, R (>= 3.1.0)

Suggests knitr, rmarkdown

VignetteBuilder knitr

NeedsCompilation yes

Repository CRAN

Date/Publication 2017-06-15 14:49:51 UTC

R topics documented:

Cancer	2
cis.eqtls	3
cpsets	3
deBarcode	4
dpdiff	5
GWAS	6

intersect	7
intersectElements	8
jaccard	9
MSET	9
msets	11
plot.msets	12
summary.msets	14
SuperExactTest	16
supertest	16

Index	18
--------------	-----------

Cancer	<i>Cancer Census Dataset</i>
--------	------------------------------

Description

This example dataset contains a list of seven cancer predisposition gene sets.

Usage

Cancer

Details

The seven cancer predisposition gene sets are:

- 1) NRG (Rahman, N. Realizing the promise of cancer predisposition genes. *Nature* 2014, 505:302-308);
- 2) NBG (Tamborero, D. et al. Comprehensive identification of mutational cancer driver genes across 12 tumor types. *Scientific reports* 2013, 3:2650);
- 3) LDG (Kandoth, C. et al. Mutational landscape and significance across 12 major cancer types. *Nature* 2013, 502:333-339);
- 4) GGG (Lawrence, M. S. et al. Discovery and saturation analysis of cancer genes across 21 tumour types. *Nature* 2014, 505:495-501);
- 5) ELG (Garraway, L. A. & Lander, E. S. Lessons from the cancer genome. *Cell* 2013, 153:17-37);
- 6) CCG (Futreal, P. A. et al. A census of human cancer genes. *Nature reviews. Cancer* 2004, 4:177-183);
- 7) BVG (Vogelstein, B. et al. Cancer genome landscapes. *Science* 2013, 339:1546-1558).

References

Minghui Wang, Yongzhong Zhao, and Bin Zhang (2015). Efficient Test and Visualization of Multi-Set Intersections. *Scientific Reports* 5: 16923.

See Also

[supertest](#)

 cis.eqtls

cis-eQTLs

Description

This example dataset contains a list of cis-eQTL genes.

Details

A list is included in this dataset: `cis.eqtls`, which contains four sets of cis-eQTL genes published by Gibbs et al (PLOS Genetics 2010, 6:e1000952) as deposited in the eQTL Browser (<http://www.ncbi.nlm.nih.gov/projects/g>). The four sets of cis-eQTL genes were detected in four different brain regions from Gibbs: brain cerebellum (CB), brain frontal cortex region (FC), brain temporal cortex region (TC), and brain pons region (PONS) respectively.

See Also

[supertest](#)

 cpsets

Multi-Set Intersection Probability

Description

Density and distribution function of multi-set intersection test.

Usage

```
dpsets(x,L,n,log.p =FALSE)
cpsets(x,L,n,lower.tail=TRUE,log.p=FALSE,
       simulation.p.value=FALSE,number.simulations=1000000)
```

Arguments

<code>x</code>	integer, number of elements overlap among all sets.
<code>L</code>	vector, set sizes.
<code>n</code>	integer, background population size.
<code>lower.tail</code>	logical; if TRUE, probability is $P[\text{overlap} \leq x]$, otherwise, $P[\text{overlap} > x]$.
<code>log.p</code>	logical; if TRUE, probability p is given as $\log(p)$.
<code>simulation.p.value</code>	logical; if TRUE, probability p is computed from simulation.
<code>number.simulations</code>	integer; number of simulations.

Value

dpsets gives the density and cpsets gives the distribution function.

Author(s)

Minghui Wang <minghui.wang@mssm.edu>, Bin Zhang <bin.zhang@mssm.edu>

References

Minghui Wang, Yongzhong Zhao, and Bin Zhang (2015). Efficient Test and Visualization of Multi-Set Intersections. *Scientific Reports* 5: 16923.

See Also

[supertest](#), [MSET](#), [cpdiff](#), [dpdiff](#)

Examples

```
## Not run:
#set up fake data
n=500; A=260; B=320; C=430; D=300; x=170
(d=dpsets(x,c(A,B,C,D),n))
(p=cpsets(x,c(A,B,C,D),n,lower.tail=FALSE))

## End(Not run)
```

deBarcode

Decrypt Barcode

Description

Decrypt barcode information.

Usage

```
deBarcode(barcode, setnames, collapse=' & ')
```

Arguments

barcode a vector of character strings, encoding the intersection combination.
setnames set names.
collapse an optional character string to separate the results. See [paste](#).

Details

barcode are character strings of '0' and '1', indicating absence or presence of each set in a intersection combination.

Value

A vector.

Author(s)

Minghui Wang <minghui.wang@mssm.edu>

Examples

```
deBarcode(c('01011', '10100'), c('S1', 'S2', 'S3', 'S4', 'S5'))
```

 dpdiff

Set Difference Probability

Description

Density and distribution functions of set difference tests.

Usage

```
dpdiff(x, s1, s2, s3, n, log.p=FALSE)
dpone(x, s1, s2, s3, n, log.p=FALSE)
cpdiff(x, s1, s2, s3, n, lower.tail=TRUE, log.p=FALSE)
cpone(x, s1, s2, s3, n, lower.tail=TRUE, log.p=FALSE)
```

Arguments

<code>x</code>	integer, number of elements that are different.
<code>s1, s2, s3</code>	set sizes.
<code>n</code>	integer, background population size.
<code>lower.tail</code>	logical; if TRUE, probability is $P[X \leq x]$, otherwise, $P[X > x]$.
<code>log.p</code>	logical; if TRUE, probability p is given as $\log(p)$.

Value

`dpdiff` gives the density and `cpdiff` gives the distribution function for set difference of $\text{intersect}(A, B) \setminus C$, while `dpone` gives the density and `cpone` gives the distribution function of set difference of $A \setminus \text{union}(B, C)$, where $|A|=s1$, $|B|=s2$, and $|C|=s3$.

Author(s)

Minghui Wang <minghui.wang@mssm.edu>, Bin Zhang <bin.zhang@mssm.edu>

References

Minghui Wang, Yongzhong Zhao, and Bin Zhang (2015). Efficient Test and Visualization of Multi-Set Intersections. *Scientific Reports* 5: 16923.

See Also

[cpsets](#), [dpsets](#)

Examples

```
## Not run:
#set up fake data
n=500; s1=260; s2=320; s3=430; x=80
(d1=sapply(0:80,function(x) SuperExactTest:::dpdiff(x,s1,s2,s3,n)))
(p1=sapply(0:80,function(x) SuperExactTest:::cpdiff(x,s1,s2,s3,n,lower.tail=TRUE)))
par(mfrow=c(1,2))
plot(0:80,d1,main='Density');
plot(0:80,p1,main='Distribution');

## End(Not run)
```

GWAS

GAWS Catalog Dataset

Description

This example dataset contains a list of gene sets associated with six types of clinical traits curated in the GWAS Catalog.

Usage

GWAS

Details

The six clinical traits are:

- 1) NEU (Bipolar disorder and schizophrenia, Schizophrenia, Major depressive disorder, Alzheimer's disease, Parkinson's disease, Cognitive performance, Bipolar disorder);
- 2) INF (Crohn's disease, Ulcerative colitis, Inflammatory bowel disease, Rheumatoid arthritis, Multiple sclerosis, Systemic lupus erythematosus);
- 3) CVD (Type 2 diabetes, Coronary heart disease, Blood pressure, total Cholesterol, HDL cholesterol, Triglycerides);
- 4) HT (height);
- 5) IgG (IgG glycosylation);
- 6) OB (obesity, obesity related traits).

References

Minghui Wang, Yongzhong Zhao, and Bin Zhang (2015). Efficient Test and Visualization of Multi-Set Intersections. *Scientific Reports* 5: 16923.

See Also[supertest](#)

`intersect`*Set Operations*

Description

Performs set union and intersection on multiple input vectors.

Usage

```
union(x, y, ...)  
intersect(x, y, ...)
```

Arguments

`x, y, ...` vectors (of the same mode) containing a sequence of items (conceptually) with no duplicated values.

Details

These functions extend the the same functions in the base package to handle more than two input vectors.

Value

A vector of the same mode as `x` or `y` for `intersect`, and of a common mode for `union`.

Author(s)

Minghui Wang <minghui.wang@mssm.edu>, Bin Zhang <bin.zhang@mssm.edu>

References

Minghui Wang, Yongzhong Zhao, and Bin Zhang (2015). Efficient Test and Visualization of Multi-Set Intersections. *Scientific Reports* 5: 16923.

Examples

```
##not run##
```

intersectElements *Find Intersection Membership*

Description

Find intersections and assign element to intersection combinations.

Usage

```
intersectElements(x, mutual.exclusive=TRUE)
```

Arguments

x list; a collection of sets.
mutual.exclusive logical; see Details.

Details

See example below for the use of mutual.exclusive.

Value

A data.frame with two columns:

Entry set elements.
barcode intersection combination that each entry belongs to.

Author(s)

Minghui Wang <minghui.wang@mssm.edu>

Examples

```
set.seed(123)  
sets=list(S1=sample(letters,10), S2=sample(letters,5), S3=sample(letters,7))  
intersectElements(sets,mutual.exclusive=TRUE)  
intersectElements(sets,mutual.exclusive=FALSE)
```

`jaccard`*Calculate Jaccard Index*

Description

This function calculates Jaccard indices between pairs of sets.

Usage

```
jaccard(x)
```

Arguments

`x` list, a collect of sets.

Value

A matrix of pairwise Jaccard indices.

Author(s)

Minghui Wang <minghui.wang@mssm.edu>

Examples

```
## Not run:  
#set up fake data  
x=list(S1=letters[1:20], S2=letters[10:26], S3=sample(letters,10), S4=sample(letters,10))  
jaccard(x)  
  
## End(Not run)
```

`MSET`*Exact Test of Multi-Set Intersection*

Description

Calculate FE and significance of intersection among multiple sets.

Usage

```
MSET(x,n,lower.tail=TRUE,log.p=FALSE)
```

Arguments

<code>x</code>	list; a collection of sets.
<code>n</code>	integer; background population size.
<code>lower.tail</code>	logical; if TRUE, probability is $P[\text{overlap} < m]$, otherwise, $P[\text{overlap} \geq m]$, where m is the number of elements overlap between all sets.
<code>log.p</code>	logical; if TRUE, probability p is given as $\log(p)$.

Details

This function implements an efficient statistical test for multi-set intersections. The algorithm behind this function was described in Wang et al 2015.

Value

A list with the following elements:

<code>intersects</code>	a vector of intersect items.
<code>FE</code>	fold enrichment of the intersection.
<code>p.value</code>	one-tail probability of observing equal to or larger than the number of intersect items.

Author(s)

Minghui Wang <minghui.wang@mssm.edu>, Bin Zhang <bin.zhang@mssm.edu>

References

Minghui Wang, Yongzhong Zhao, and Bin Zhang (2015). Efficient Test and Visualization of Multi-Set Intersections. *Scientific Reports* 5: 16923.

See Also

[supertest](#), [cpsets](#), [dpsets](#)

Examples

```
## Not run:
#set up fake data
x=list(S1=letters[1:20], S2=letters[10:26], S3=sample(letters,10), S4=sample(letters,10))
MSET(x, 26, FALSE)

## End(Not run)
```

`msets`*Class to Contain Multi-Set Intersections*

Description

This object contains data regarding the intersections between multiple sets. This object is usually created by the `superTest` function.

Details

Intersection combination is denoted by a barcode string of '0' and '1', where a value of '1' in the i th position of the string indicates that the intersection is involved with the i th set, 0 otherwise. E.g., string '000101' indicates that the intersection is an overlap between the 4th and 6th sets. Function [deBarcode](#) can be used to decrypt the barcode. Generic `summary` and `plot` functions can be applied to extract and visualize the results.

Value

<code>x</code>	a list of sets from input.
<code>set.names</code>	names of the sets. If the input sets do not have names, they will be automatically named as SetX where X is an integer from 1 to the total number of sets.
<code>set.sizes</code>	a vector of set sizes.
<code>n</code>	background population size.
<code>overlap.sizes</code>	a named vector of intersection sizes. Each intersection component is named by a barcoded character string of '0' and '1'. See <code>Details</code> for barcode.
<code>P.value</code>	a vector of p values for the intersections.

Author(s)

Minghui Wang <minghui.wang@mssm.edu>, Bin Zhang <bin.zhang@mssm.edu>

References

Minghui Wang, Yongzhong Zhao, and Bin Zhang (2015). Efficient Test and Visualization of Multi-Set Intersections. *Scientific Reports* 5: 16923.

See Also

[superTest](#), [summary.msets](#), [plot.msets](#), [deBarcode](#)

plot.msets

*Draw Multi-Set Intersections***Description**

This function draws intersections among multiple sets.

Usage

```
## S3 method for class 'msets'
plot(x, Layout=c('circular','landscape'), degree=NULL,
keep.empty.intersections=TRUE,
sort.by=c('set','size','degree','p-value'),
min.intersection.size=0, max.intersection.size=Inf,
ylim=NULL, log.scale=FALSE, x.pos=c(0.05,0.95),
y.pos=c(0.025,0.975), yfrac=0.8, color.scale.pos=c(0.85, 0.9),
legend.pos=c(0.85,0.25), legend.col=2, legend.text.cex=1, color.scale.cex=1,
color.scale.title=expression(paste(-Log[10], '(', italic(P), ')')),
color.on='#2EFE64', color.off='#EEEEEE', show.overlap.size=TRUE,
show.set.size=TRUE, track.area.range=0.3, bar.area.range=0.2,
origin=if(sort.by[1]=='size'){c(0.45,0.5)}else{c(0.5,0.5)},
pos.size=0.005, new.gridPage=TRUE, minMinusLog10PValue=0,
maxMinusLog10PValue=NULL, ...)
```

Arguments

x	a msets object.
Layout	layout for plotting.
degree	a vector of intersection degrees for plotting. E.g., when degree=c(2:3), only those intersections involving two or three sets will be plotted. By default, degree=NULL, all possible intersections are plotted.
keep.empty.intersections	logical; if FALSE, empty intersection(s) will be discarded to save plotting space.
min.intersection.size	Minimum size of an intersection to be plotted.
max.intersection.size	Maximum size of an intersection to be plotted.
sort.by	how to sort intersections. It is one of "set", "size", "degree", and "p-value".
ylim	the limits c(y1, y2) of plotting overlap size.
log.scale	logical; whether to plot with log transformed intersection sizes.
x.pos	numeric; x coordinate (0 to 1) of the graph canvas for landscape Layout.
y.pos	numeric; y coordinate (0 to 1) of the graph canvas for landscape Layout.
yfrac	numeric; the fraction (0 to 1) of canvas used for plotting bars. Only used for landscape Layout.

<code>color.scale.pos</code>	numeric; x and y coordinates (0 to 1) for packing the color scale guide. It could be a keyword "topright" or "topleft" in the landscape layout, and one of "topright", "topleft", "bottomright" and "bottomleft" in the circular layout.
<code>legend.pos</code>	numeric; x and y coordinates (0 to 1) for packing the legend in the circular layout. It could be one of the keywords "bottomright", "bottomleft", "topleft" and "topright".
<code>legend.col</code>	integer; number of columns of the legend in the circular layout.
<code>legend.text.cex</code>	numeric; specifying the amount by which legend text should be magnified relative to the default.
<code>color.scale.cex</code>	numeric; specifying the amount by which color scale text should be magnified relative to the default.
<code>color.scale.title</code>	character or expression; a title for the color scale guide.
<code>color.on</code>	color code; specifying the color for set(s) which are "present" for an intersection.
<code>color.off</code>	color code; specifying the color for set(s) which are "absent" for an intersection.
<code>show.overlap.size</code>	logical; whether to show overlap size in circular layout.
<code>show.set.size</code>	color code; whether to show set size in landscape layout.
<code>track.area.range</code>	the magnitude of track area from origin in the circular layout.
<code>bar.area.range</code>	the magnitude of bar area from edge of the track area in the circular layout. The sum of <code>track.area.range</code> and <code>bar.area.range</code> should not be larger than 0.5.
<code>origin</code>	the origin coordinates (0 to 1) in the circular layout.
<code>pos.size</code>	position offset when plotting the set size for circular Layout.
<code>new.gridPage</code>	logic; whether to start a new grid page. Set FALSE to allow for customized arrangement of the grid layout.
<code>minMinusLog10PValue</code>	numeric; minimum minus log ₁₀ P value for capping the scale of color map. Default 0.
<code>maxMinusLog10PValue</code>	numeric; maximum minus log ₁₀ P value for capping the scale of color map. Default maximum from the data.
<code>...</code>	additional arguments for plot function, including <code>heatmapColor</code> (a vector of customized heat colors), <code>cex</code> (scale of text font size), <code>phantom.traks</code> (number of phantom tracks in the middle in the circular layout, default 2), <code>gap.within.track</code> (ratio of gap width over block width on the same track, default 0.1), and <code>gap.between.track</code> (ratio of gap width over track width, default 0.1). Not fully implemented.

Details

The plot canvas has coordinates 0~1 for both x and y axes.

Value

No return.

Author(s)

Minghui Wang <minghui.wang@mssm.edu>, Bin Zhang <bin.zhang@mssm.edu>

References

Minghui Wang, Yongzhong Zhao, and Bin Zhang (2015). Efficient Test and Visualization of Multi-Set Intersections. *Scientific Reports* 5: 16923.

See Also

[msets](#)

Examples

```
## Not run:
#set up fake data
x=list(S1=letters[1:20], S2=letters[10:26], S3=sample(letters,10), S4=sample(letters,10))
obj=supertest(x,n=26)
plot(obj)

## End(Not run)
```

summary.msets

Summarize an msets Object

Description

This function outputs summary statistics of a msets object.

Usage

```
## S3 method for class 'msets'
summary(object, degree=NULL, ...)
```

Arguments

object	a msets object.
degree	a vector of intersection degrees to pull out.
...	additional arguments (not implemented).

Value

A list:

Barcode	a vector of 0/1 character strings, representing the set composition of each intersection.
otab	a vector of observed intersection size between any combination of sets.
etab	a vector of expected intersection size between any combination of sets if background population size is specified.
set.names	set names.
set.sizes	set sizes.
n	background population size.
P.value	upper tail p value for each intersection if background population size n is specified.
Table	a data.frame containing degree, otab, etab, fold change, p value and the overlap elements.

Author(s)

Minghui Wang <minghui.wang@mssm.edu>, Bin Zhang <bin.zhang@mssm.edu>

References

Minghui Wang, Yongzhong Zhao, and Bin Zhang (2015). Efficient Test and Visualization of Multi-Set Intersections. *Scientific Reports* 5: 16923.

See Also

[msets](#)

Examples

```
## Not run:
#set up fake data
x=list(S1=letters[1:20], S2=letters[10:26], S3=sample(letters,10), S4=sample(letters,10))
obj=supertest(x,n=26)
summary(obj)

## End(Not run)
```

SuperExactTest	<i>SuperExactTest Package</i>
----------------	-------------------------------

Description

Efficient Test and Visualization of Multi-set Intersections

Details

The main functions that most users may need from this package are [supertest](#) and [MSET](#). For a brief introduction of using this package, please see `vignette("set_html")`.

Author(s)

Minghui Wang <minghui.wang@mssm.edu>, Bin Zhang <bin.zhang@mssm.edu>

References

Minghui Wang, Yongzhong Zhao, and Bin Zhang (2015). Efficient Test and Visualization of Multi-Set Intersections. *Scientific Reports* 5: 16923.

See Also

[supertest](#), [MSET](#)

Examples

```
## Not run:
#See a brieft instroduction of using this package
vignette("set_html")

## End(Not run)
```

supertest	<i>Calculate Intersections Among Multiple Sets and Perform Statistical Tests</i>
-----------	--

Description

This function calculates intersection sizes among multiple sets and performs statistical tests of the intersections.

Usage

```
supertest(x, n=NULL, degree=NULL, ...)
```


Arguments

x	list; a collection of sets.
n	integer, background population size. Required for computing the statistical significance of intersections.
degree	a vector of intersection degrees for overlap analysis. E.g., when degree=c(2:3), only those intersections involving two or three sets will be computed. By default, degree=NULL, all possible intersections are computed.
...	additional arguments (not implemented).

Details

This function calculates intersection sizes between multiple sets and, if background population size n is specified, performs statistical tests of the intersections. For a brief introduction of using this package, please see `vignette("set_html")`.

Value

An object of class `msets`.

Author(s)

Minghui Wang <minghui.wang@mssm.edu>, Bin Zhang <bin.zhang@mssm.edu>

References

Minghui Wang, Yongzhong Zhao, and Bin Zhang (2015). Efficient Test and Visualization of Multi-Set Intersections. *Scientific Reports* 5: 16923.

See Also

[msets](#), [MSET](#), [Cancer](#), [cpsets](#), [dpsets](#)

Examples

```
## Not run:  
#Analyze the cancer gene sets  
data(Cancer)  
Result=supertest(Cancer, n=20687)  
summary(Result)  
plot(Result,degree=2:7,sort.by='size')  
  
## End(Not run)
```

Index

*Topic **classes**

msets, [11](#)

*Topic **datasets**

Cancer, [2](#)

cis.eqtls, [3](#)

GWAS, [6](#)

Cancer, [2](#), [17](#)

cis.eqtls, [3](#)

cpdiff, [4](#)

cpdiff (dpdiff), [5](#)

cpone (dpdiff), [5](#)

cpsets, [3](#), [6](#), [10](#), [17](#)

deBarcode, [4](#), [11](#)

dpdiff, [4](#), [5](#)

dpone (dpdiff), [5](#)

dpsets, [6](#), [10](#), [17](#)

dpsets (cpsets), [3](#)

GWAS, [6](#)

intersect, [7](#)

intersectElements, [8](#)

jaccard, [9](#)

MSET, [4](#), [9](#), [16](#), [17](#)

msets, [11](#), [14](#), [15](#), [17](#)

paste, [4](#)

plot.msets, [11](#), [12](#)

summary.msets, [11](#), [14](#)

SuperExactTest, [16](#)

supertest, [2-4](#), [7](#), [10](#), [11](#), [16](#), [16](#)

supertest, list-method (supertest), [16](#)

union (intersect), [7](#)