

# Package ‘NBBttest’

May 18, 2022

**Type** Package

**Title** Negative Binomial Beta t-Test

**Version** 1.0.1

**Date** 2022-05-08

**Suggests** BiocGenerics,DEXSeq

**Author** Yuan-De Tan

**Maintainer** Yuan-De Tan <tanyuande@gmail.com>

**Description** We constructed 'NBBttest' for identifying genes or RNA isoforms differentially expressed between two conditions on RNA-seq count data. Package 'NBBttest' can perform data quality check, data normalization, differential analysis, annotation and graphic analysis. In differential analysis, 'NBBttest' can identify differentially expressed genes and differential RNA isoforms in alternative splicing sites and alternative polyadenylation sites, differential sgRNA, and differential CRISPR (clustered regularly interspaced short palindromic repeats) screening genes. In graphic analysis, 'NBBttest' provides two types of heatmaps to visualize differential expression at gene or isoform level using z-score and n-score and creates pathway heatmap. 'NBBttest' can plot differentially expressed exons within a specified gene. In addition, 'NBBttest' provides a tool for annotation of genes and exons. The methods used in 'NBBttest' were new statistical methods developed from Tan and others (2015) <[doi:10.1371/journal.pone.0123658](https://doi.org/10.1371/journal.pone.0123658)>. The statistical methods are robust and very powerful in identifying genes or RNA isoforms differentially expressed between two conditions in small samples.

**License** GPL-3

**Depends** R (>= 3.5.0),gplots,gtools,stats,utils,graphics,grDevices

**LazyLoad** yes

**biocViews** Sequencing, DifferentialExpression, MultipleComparison,  
SAGE, GeneExpression, Transcription,  
AlternativeSplicing,Coverage, DifferentialSplicing

**NeedsCompilation** no

**RoxygenNote** 6.1.1

## R topics documented:

NBBttest-package . . . . .	2
annotat . . . . .	3
betaparametab . . . . .	4

betaparametVP . . . . .	5
betaparametw . . . . .	7
betatetest . . . . .	8
DDX39_100 . . . . .	9
exondata . . . . .	10
gbetatetest . . . . .	11
gtfa . . . . .	12
jktcell . . . . .	13
mbetatetest . . . . .	14
mtpvadjust . . . . .	16
myheatmap . . . . .	17
myheatmap2 . . . . .	19
NBBplot . . . . .	22
normalized . . . . .	23
oddratio . . . . .	24
omega . . . . .	25
pathwayHeatmap . . . . .	26
pathwy.A.up . . . . .	27
pratio . . . . .	28
prime3_PRP8_50 . . . . .	29
QC . . . . .	30
result . . . . .	32
sgRNA . . . . .	33
simSplicing . . . . .	34
simulat . . . . .	34
skjt . . . . .	36
smbetatetest . . . . .	37
subdata . . . . .	38
upGAm . . . . .	39
<b>Index</b>	<b>40</b>

NBBttest-package

*Negative Beta Binormal t-Test Package*

## Description

This package consists of 20 functions, of which betaparametab.R, betaparametVP.R, betaparametw.R, gbetatetest.R, betatetest.R, mbetatetest.R, normalize.R, oddratio.R, pratio.R, omega.R, simulat.R, and smbetatetest.R are used to estimate beta, alpha, weight, t-statistics, rho and omega, p-value, and multiple test at gene level or RNA isoform level. NBBplot.R is used to visualize count data of exons within a specified gene in given conditions. QC.R is used to check quality of data, and myheatmap.R, myheatmap2.R, and pathwayheatmap.R are used to show heatmap of differential expressions of DE genes, DE RNA isoform or pathways. Run of mbetatetest.R would output beta t-test results including geneid or isoformid, gene name, the other information, t-value, p-value, rho, and omega(W).

## Details

Package: NBBttest  
 Type: Package  
 Version: 1.0.0  
 Date: 2018-01-11  
 License: GPL-3

## Author(s)

Yuan-De Tan

Maintainer: Yuan-De Tan <tanyuande@gmail.com>

## References

Baggerly KA, Deng L, Morris JS, Aldaz CM (2003) Differential expression in SAGE: accounting for normal between-library variation. *Bioinformatics*, **19**: 1477-1483.

Yuan-De Tan Anita M. Chandler, Arindam Chaudhury, and Joel R. Neilson(2015) A Powerful Statistical Approach for Large-scale Differential Transcription Analysis. *Plos One*,10.1371/journal.pone.0123658.

## See Also

[betaparametab](#), [betaparametVP](#), [betaparametw](#), [gbetattest](#), [betattest](#), [mbetattest](#), [myheatmap](#), [pathwayHeatmap](#) [myheatmap2](#), [normalized](#), [oddratio](#), [NBBplot](#), [omega](#), [pratio](#), [QC](#), [simulat](#), [smbetattest](#)

## Examples

```
data(jkttcell)
res<-mbetattest(X=jkttcell[1:100, ], nci=7, na=3,
nb=3, alpha=0.05, norm="yes", side="both",
level="isoform",padjust_methods="fdr",C=0)
```

---

annotat

*Annotation of genes within which alternative splicing occurs*

---

## Description

Alternative splicing is detected in any element of 3'UTR, 5'UTR, exons and introns within a gene using RNA-seq data where RNA reads are mapped to a reference genome. As an example for annotation, the RNA-seq reads derived from human samples can be mapped onto human genome reference (GRCh38) using different methods, for example, HTSeq, spladder, rMAT, cufflinks, etc. These methods can detect alternative splicing sites within genes. However, none of these methods does gene annotation for users. Our NBBttest offers a R function for annotating genes with exons or isoforms.

## Usage

```
annotat(infile, mfile, type="gene", columnset = "NLL")
```

**Arguments**

<code>infile</code>	input data file with ENSGid column for annotation
<code>mfile</code>	reference genome file with ENSG id column.
<code>type</code>	has three options: "gene", "isoform" or ""isof" and "exon", see details. The default is "gene".
<code>columnset</code>	EGNSid column set. If the RNA count data are made by using HTSeq and DEXSeq annotation file, then some of genes have many different ENSGids. For example, ENSG00000285476+ENSG00000182230+ENSG00000251623 has three ENSG ids ENSG00000285476,ENSG00000182230, and ENSG00000251623 that share one gene, so it is splited into three columns (2,3,4)in excel. The default is 2(column 2).

**Details**

If type = "gene", then count data of RNA reads are obtained at gene level, annotation would be executed at gene level or if type = "isof" or "isoform", then RNA reads were mapped onto elements (for example, 3'UTR, 5'UTR, exon or cassette, intron) within genes and annotation would be executed at isoform level or if type = "exon", then RNA count data were obtained by mapping RNA reads onto exome by DEXseq and annotation would be done at exon level defined by DEXSeq. Note that GRCh38 is too big so it was removed from data. User may request to get it from yuande/github.

**Value**

return original data with an additional column for gene.

**Author(s)**

Yuan-De Tan <tanyuande@gmail.com>

**Examples**

```
data(DDX39_100)
data(gtfa)
DDX39_30<-annotat(infile=DDX39_100[1:30,],mfile=gtfa,type="gene")
```

---

betaparametab

---

*Estimation of parameters alpha ( $\alpha$ ) and beta ( $\beta$ ) of beta distribution*


---

**Description**

Parameters alpha(a) and beta (b) in betat distribution are estimated by using an iteration algorithm.

**Usage**

```
betaparametab(xn, w, P, V)
```

**Arguments**

xn	column vector, a set of library sizes.
w	column vector, a set of weights.
P	proportion of counts of a gene or an isoform.
V	variance of proportions of counts of a gene or an isoform over m replicate libraries in a condition.

**Value**

return parameters a and b.

**Author(s)**

Yuan-De Tan <tanyuande@gmail.com>

**References**

Baggerly KA, Deng L, Morris JS, Aldaz CM (2003) Differential expression in SAGE: accounting for normal between-library variation. *Bioinformatics* **19**: 1477-1483.

**See Also**

[betaparametVP](#), [betaparametw](#)

**Examples**

```
XX<-c(2000,2000,2000)
p<-0.15
V=0.004
w<-c(0.3,0.3,0.3)
betaparametab(xn=XX,w=w,P=p,V=V)
```

---

betaparametVP

*Estimation of parameters V and P in count data of RNA reads*


---

**Description**

This function is used to estimate parameters P and V by optimizing estimates of parameters: alpha and beta.

**Usage**

```
betaparametVP(X, NX)
```

**Arguments**

X	count dataset derived from m replicate libraries in one condition.
NX	vector of m library sizes. Library size is sum of counts over the whole library.

## Details

Count data of *RNA* sequence reads are assumed to follow binomial distribution with parameters (P) and (n) or negative binomial distribution with parameters (P) and (r) , while P is frequency of a gene or an isoform in RNA sequence population and assumed to follow beta distribution with parameters alpha (a) and beta(b). Parameters P and V are estimated by using optimal estimation of parameters a and b. The optimal method is an iteration algorithm driven by weighting proportion of gene or isoform in each replicate library. This is a large-scale method for estimating these parameters. Estimation of parameters P and V is core of the multiple beta t-test method because P and V will be used to calculate t-value.

## Value

return a list:

P	N proportions estimated.
V	N variances estimated.

## Note

betaparametVP requires functions betaparametab and betaparametw.

## Author(s)

Yuan-DE Tan <tanyuande@gmail.com>

## References

Baggerly KA, Deng L, Morris JS, Aldaz CM (2003) Differential expression in SAGE: accounting for normal between-library variation. *Bioinformatics*, **19**: 1477-1483.  
 Yuan-De Tan, Anita M. Chandler, Arindam Chaudhury, and Joel R. Neilson(2015) A Powerful Statistical Approach for Large-scale Differential Transcription Analysis. *Plos One*,10.1371/journal.pone.0123658.

## See Also

[betaparametab](#), [betaparametw](#)

## Examples

```
data(jkttcell)
X<-jkttcell[1:100,]
na<-3
nb<-3
cn<-length(X[,])
rn<-length(X[,1])
XC<-X[,1:(cn-na-nb)]
XX<-X[, (cn-na-nb+1):cn]
n<-na+nb
XA<-XX[,1:na]
SA<-apply(XA,2,sum)
PA<-betaparametVP(XA,SA)
```

---

betaparametw	<i>Estimation of proportion weights</i>
--------------	---

---

## Description

Function betaparametw is used to calculate weights.

## Usage

```
betaparametw(xn, a, b)
```

## Arguments

xn	a vector of m library sizes. Library size is sum of counts over the whole library.
a	parameter alpha( $\alpha$ ) in beta distribution derived from output of function betaparametab.
b	parameter beta ( $\beta$ ) in beta distribution derived from output of function betaparametab.

## Details

alpha and beta ( $\alpha, \beta$ ) are used to calculate weight. Then weight is in turn used to correct bias of estimation of alpha and beta in betaparametab function.

## Value

return weight(W).

## Author(s)

Yuan-De Tan <tanyuande@gmail.com>

## References

Baggerly KA, Deng L, Morris JS, Aldaz CM (2003) Differential expression in SAGE: accounting for normal between-library variation. *Bioinformatics*, **19**: 1477-1483.  
Yuan-De Tan, Anita M. Chandler, Arindam Chaudhury, and Joel R. Neilson(2015) A Powerful Statistical Approach for Large-scale Differential Transcription Analysis. *Plos One*. 2015 DOI: 10.1371/journal.pone.0123658.

## See Also

[betaparametab](#) and [betaparametVP](#).

## Examples

```
XX<-c(2000,2000,2000)
a<-1.1458
b<-6.4932
betaparametw(xn=XX,a=a,b=b)
```

---

betatetest	<i>Beta t-test</i>
------------	--------------------

---

**Description**

Beta t-test and degree of freedom for each gene or isoform are calculated in this function.

**Usage**

```
betatetest(X, na, nb, NX=100, level)
```

**Arguments**

X	count data of RNA sequence reads containing N genes (or isoforms).
na	number of replicate libraries in condition A.
nb	number of replicate libraries in condition B.
NX	numeric value. It is used at level="isoform". NX=100 is default but does not used at any level.
level	string, has three options: "isoform" or "sgRNA"

**Details**

In beta t-test,

$$t = \frac{(P_A - P_B)}{\sqrt{(V_A + V_B)}}$$

where  $P_A$  and  $P_B$  are proportions of a gene or an isoform in conditions A and B,  $V_A$  and  $V_B$  are variances of this gene or isoform in conditions A and B, respectively. They are output of betaparam-etVP.

**Value**

return two lists:

t	t-value list.
df	df list. df is degree of freedom.

**Note**

In our method, pooled standard error > 0 in any case, so the t-statistics always has definition.

**Author(s)**

Yuan-De Tan <tanyuande@gmail.com>

**References**

Baggerly KA, Deng L, Morris JS, Aldaz CM (2003) Differential expression in SAGE: accounting for normal between-library variation. *Bioinformatics*, **19**: 1477-1483.  
Yuan-De Tan, Anita M. Chandler, Arindam Chaudhury, and Joel R. Neilson(2015) A Powerful Statistical Approach for Large-scale Differential Transcription Analysis. *Plos One*. 2015 DOI: 10.1371/journal.pone.0123658.



**See Also**

[pratio](#), [oddratio](#).

**Examples**

```
data(jkttcell)
X<-jkttcell[1:100,]
na<-3
nb<-3
cn<-ncol(X)
rn<-nrow(X)
XC<-X[,1:(cn-na-nb)]
XX<-X[(cn-na-nb+1):cn]
betatetest<-betatetest(X=XX, na=3,nb=3, level="isoform")
```

---

DDX39\_100

*DDX39 exon data with 100 exons*


---

**Description**

DDX39\_100 dataset is an example for implementing annotation of genes within there are some exons with gtf.

**Usage**

```
data("DDX39_100")
```

**Format**

A data frame with 100 observations on the following 14 variables.

X a numeric vector  
 gene\_id a factor with ENSG levels  
 CTL\_1\_S3 a numeric vector  
 CTL\_2\_S7 a numeric vector  
 CTL\_3\_S11 a numeric vector  
 DDX39\_1\_S1 a numeric vector  
 DDX39\_2\_S5 a numeric vector  
 DDX3\_S9 a numeric vector  
 tvalue a numeric vector  
 rho a numeric vector  
 df a numeric vector  
 pvalue a numeric vector  
 adjp a numeric vector  
 w a numeric vector

## Details

DDX39\_100 is dataset that was RNA-seq data from human cell breast cells with knockout DDX39 gene. DDX39 RNA-seq reads were mapped to human reference annotated with GRCh38 by HTseq with DEXseq annotation and generated count data of RNA reads of exons within genes. DDX39\_100 is output of NBBttest where X is row order column; gene\_id is ENSG ID; CTL is control cells with replicates 1, 2, and 3; DDX39 is treatment cells with knockout gene DDX39; tvalue is NBBttest t-statistic; rho is exon-wise variable of NBBttest; df is degree of freedom; pvalue is p-value for t-statistics and w is threshold of rho.

## Examples

```
data(DDX39_100)
```

---

exondata

*Exon data for NBBplot*

---

## Description

This dataset is an example for NBBplot to show differential exons within genes between two conditions.

## Usage

```
data("exondata")
```

## Format

A data frame with 16 observations on the following 22 variables.

ID a factor with levels ENSMUSG00000015597:E1 ENSMUSG00000015597:E10 ENSMUSG00000015597:E2 ENSMUSG00000015597:E3 ENSMUSG00000015597:E4 ENSMUSG00000015597:E5 ENSMUSG00000015597:E6 ENSMUSG00000015597:E7 ENSMUSG00000015597:E8 ENSMUSG00000015597:E9 ENSMUSG00000079547:E1 ENSMUSG00000079547:E2 ENSMUSG00000079547:E3 ENSMUSG00000079547:E4 ENSMUSG00000079547:E5 ENSMUSG00000079547:E6

chr a numeric vector

element a factor with levels exon

start a numeric vector

end a numeric vector

strand a factor with levels +

gene\_id a factor with levels ENSMUSG00000015597 ENSMUSG00000079547

exon a numeric vector

gene a factor with levels H2-DMb1 Zfp318

WT\_28 a numeric vector

WT\_43 a numeric vector

A25\_f0 a numeric vector

cKO\_30 a numeric vector

A26\_cKO a numeric vector

A6\_f0 a numeric vector  
 tv a numeric vector  
 rho a numeric vector  
 pvalue a numeric vector  
 w a numeric vector  
 order\_number a numeric vector  
 FDR0.05 a numeric vector  
 significance a numeric vector

## Details

This dataset is an object of NBBttest and annotated with a R function. It contains columns start and end, gene, isoform, tv and pvalue, data, significance that NBBplut uses. WT and A25 are count data from wild type cell lines 28,43 and fo and A and cKO are also count data from knockout a gene, tv is t-statistic, rho is exon-wise variable and w is threshold of rho. Column "significance" is 1 if  $pvalue < FDR0.05$ , 0, otherwise. Gene Zfp318 has 10 exons and H2-DMb1 has 6 exons.

## Examples

```
data(exondata)
```

---

gbetattest	<i>Beta t-tests within groups</i>
------------	-----------------------------------

---

## Description

Beta t-tests are conducted within groups, genes, or libraries.

## Usage

```
gbetattest(xx, W, nci, na, nb, level, padjust_methods, C=1.222, side)
```

## Arguments

xx	a datasheet consisting of n columns and m rows. Columns contain information and count data columns n must be 1 or more and m must be over 100.
W	numeric value. It is omega estimated from null simulation.
nci	int numeric value indicating number of information columns.
na	int numeric value indicating number of replicates in condition a.
nb	int numeric value indicating number of replicates in condition b.
level	string value. It has 6 options: "isoform", "sgRNA", "RNA", "polyA.gene", "CRISPR.gene" and "splicing.gene".
padjust_methods	padjust.methods can choose one of "holm", "hochberg", "hommel", "bonferroni", "BH", "BY", "fdr", "TX", and "none" where "fdr" = "BH", "TX" is Tan and Xu's method (2015) with C=1.222 for adjusting p-value.
C	float numeric value for specifying a multiple procedure. C=0 tells mbetattest to perform single tests, C=1.222 tells mbetattest to perform BH correction of pvalues, C>1000 tells mbetattest to perform Bonferroni correction of pvalues.
side	string value for specifying one-tail test or two-tail test: side="up" for left-tail test, side="down" for right-tail test and side="both" for two-tail tests.

## Details

Beta t-test will be conducted within a specified group or at a specified level. If level="RNA", then beta t-tests will be conducted within a whole library or the whole data. If level= "isoform", then data will be separated in two parts: single-isoform and multi-isoform datasets. Single-isoform RNA indicates that there is only one RNA isoform within a gene, while multi-isoform RNAs indicate that there are at least two RNA isoforms within a gene. For single-isoforms, data are as a group and beta t-tests will be performed in the group. For the multi-isoforms, t-test will be performed within genes. If level="polyA.gene" or "CRISPR.gene", then t-test will be performed at gene level. If level="splicing.gene", then t-values and p-values will be selected from gene groups with the least p-values.

## Value

return a list containing dataset, t-values, corrected p-values, rhos and w.

## Author(s)

Yuan-De Tan <tanyuande@gmail.com>

## References

Baggerly KA, Deng L, Morris JS, Aldaz CM (2003) Differential expression in SAGE: accounting for normal between-library variation.  
 Benjamini, Y., and Hochberg, Y. (1995). Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society Series, B*, **57**, 289-300. doi:10.1111/j.2517-6161.1995.tb02031.x, <https://www.jstor.org/stable/2346101>.  
 Benjamini, Y., and Yekutieli, D. (2001). The control of the false discovery rate in multiple testing under dependency. *Annals of Statistics*, **29**, 1165 -1188. doi:10.1214/aos/1013699998.  
 Tan YD, Xu H. A general method for accurate estimation of false discovery rates in identification of differentially expressed genes. *Bioinformatics*. 2014 Jul 15;30(14):2018-25. doi:10.1093/bioinformatics/btu124. Epub 2014 Mar 14. PMID: 24632499.

## See Also

[betatetest](#)

## Examples

```
data(jkttcell)
colnames(jkttcell)[3]<-"Gene"
res.isfo<-gbetatest(xx=jkttcell[1:100,], W=1, nci=7,
na=3, nb=3, level="isoform", padjust_methods="fdr", C=0, side="both")
```

---

gtfa

gtf

---

## Description

gtf is annotation file for genes and exons

**Usage**

```
data("gtfa")
```

**Format**

A data frame with 1461 observations on the following 8 variables.

V1 a factor with levels  
 V2 a factor with levels  
 V3 a numeric vector  
 V4 a numeric vector  
 V5 a factor with levels  
 V6 a factor with levels  
 V7 a numeric vector  
 V8 a factor with levels

**Details**

V1:chr,V2:exon,V3:start position,V4:end position,V5:strand,V6:gene id,V7: exon,V8:gene

**Examples**

```
data(gtfa)
```

---

jkttcell

*Jurkat T-cell transcriptomic data*


---

**Description**

The data are transcriptomic count data of *RNA* reads generated by next generation sequencing from Jurkat T-cells.

**Usage**

```
data("jkttcell")
```

**Format**

A data frame with 13409 observations on the following 13 variables.

tagid : a numeric vector  
 geneid : a numeric vector  
 name : a string vector  
 chr : a string vector  
 strand : a character vector  
 pos : a numeric vector  
 anno : a string vector  
 Jurk.NS.A : a numeric vector

Jurk.NS.B : a numeric vector  
 Jurk.NS.C : a numeric vector  
 Jurk.48h.A : a numeric vector  
 Jurk.48h.B : a numeric vector  
 Jurk.48h.C : a numeric vector

## Details

The data are count data generated by next generation sequencing from Jurkat T-cells. The T-cells were treated by resting and stimulating with *CD3/CD28* for 48 hours. The data have 7 columns for the information of *poly(A)* sites: tagid, geneid, gene name, chromosome, strand, *poly(A)* site position, *poly(A)* site annotation and 6 columns for count data: Jurk.NS.A, Jurk.NS.B, Jurk.NS.C, Jurk.48h.A, Jurk.48h.B, Jurk.48h.C where NS means normal state or no stimulation and 48h means 48 hours after *CD3/CD28* stimulation of T-cells. 13409 RNA isoforms were detected to have alternative *poly(A)* sites.

## Value

ID, information, count data of RNA reads

## Source

Real transcriptomic count data

## References

Yuan-De Tan Anita M. Chandler, Arindam Chaudhury, and Joel R. Neilson(2015) A Powerful Statistical Approach for Large-scale Differential Transcription Analysis. *Plos One*. DOI: 10.1371/journal.pone.0123658.

## Examples

```
data(jkttcell)
```

---

mbetattest

---

*Performance of multiple beta t-test on count data*


---

## Description

This function is used to perform multiple beta t-test method on real count data. The result lists "geneid" or "isoformid", gene name, the other information, t-value, p-value, rho, and w.

## Usage

```
mbetattest(X, nci, na, nb, alpha=0.05, norm="no",
  side="both", level="sgRNA", padjust_methods, C=1.222)
```

**Arguments**

X	count data of RNA sequence reads with na replicates in condition A and nb replicates in condition B.
nci	nonnegative int value: number of columns for data information, such as geneID, isoformID, gene name etc.
na	nonnegative int value: number of replicate libraries in condition A.
nb	int numeric value: number of replicate libraries in condition B.
alpha	float numeric value, a probabilistic threshold. The value must be in [0,1]. User can set alpha=0.05 or 0.01 or the other values. Defalt value is 0.05
norm	logistic value:"yes" or "no". If norm="yes", the count data will be normalized and mbetattest will work on the normalized data, if norm="no", then mbetattest will work on the unnormalized data.
side	string for specifying tail(s) of t-distribution. If side="up", then p-value is given with t-test in the left tail. If side="down", p-value is given with t-test in right tail. If side ="both", p-value is given with t-test in both sides.
level	string for specifying which level mbetattest work on. In the current version, level has 6 options: "isoform", "sgRNA", "RNA", "splicing.gene", "polyA.gene", and "CRISPR.gene".
padjust_methods	string for specifying a method for a multiple procedure. padjust_methods can choose one of "holm", "hochberg", "hommel", "bonferroni", "BH", "BY", "fdr", "TX", and "none" where "fdr" = "BH", "TX" is Tan and Xu's method (2015) with C=1.222 for adjusting p-value.
C	float numeric value for specifying a multiple procedure. C=0 tells mbetattest to perform single tests, C=1.222 tells mbetattest to perform BH correction of pvalues, C>1000 tells mbetattest to perform Bonferroni correction of pvalues.

**Details**

see MBttest2-manual.

**Value**

return a data and result list: data columns, t-values, rho.

**Author(s)**

Yuan-De Tan <tanyuande@gmail.com>

**References**

- Baggerly KA, Deng L, Morris JS, Aldaz CM (2003) Differential expression in SAGE: accounting for normal between-library variation. *Bioinformatics* **19**: 1477-1483.
- Benjamini, Y., and Hochberg, Y. (1995). Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society Series, B*, **57**, 289-300. doi:10.1111/j.2517-6161.1995.tb02031.x, <https://www.jstor.org/stable/2346101>.
- Benjamini, Y., and Yekutieli, D. (2001). The control of the false discovery rate in multiple testing under dependency. *Annals of Statistics*, **29**, 1165 -1188. doi:10.1214/aos/1013699998.
- Tan YD, Xu H. A general method for accurate estimation of false discovery rates in identification of differentially expressed genes. *Bioinformatics*. 2014 Jul 15;30(14):2018-25. doi:10.1093/bioinformatics/btu124. Epub 2014 Mar 14. PMID: 24632499.

**See Also**

[smbetatest](#), [mtpvadjust](#), [normalized](#), [omega](#).

**Examples**

```
data(jkttcell)
res<-mbetatest(X=jkttcell[1:70, ], nci=7, na=3,
nb=3, alpha=0.05, norm="yes", side="both",
level="isoform", padjust_methods="fdr", C=0)
```

---

mtpvadjust	<i>Adjust p-values for multiple comparisons</i>
------------	---

---

**Description**

Given a set of p-values and chosen a C-value, returns a set of adjusted p-values

**Usage**

```
mtpvadjust(pv, C)
```

**Arguments**

pv	numeric vector of p-values (possibly with NAs). Any other R object is coerced by as.numeric.
C	real numeric value for specifying a multiple procedure.

**Details**

$C=0$  indicates that p-values are not adjusted,  $C=1.22$  indicates that p-values are adjusted with Benjamini and Hochberg (1995) ("BH"). The adjusted p-values are called "fdr". When  $C \geq 1000$ , p-values are adjusted with the Bonferroni method.  $C < 1.22$  indicates that p-values are adjusted by a relaxed BH method while  $C > 1.22$ , p-values are adjusted by a more strict BH method.

**Value**

A numeric vector of corrected p-values (of the same length as p, with names copied from p)

**Author(s)**

Yuan-De Tan  
<tanyuande@gmail.com>

**References**

Benjamini, Y., and Hochberg, Y. (1995). Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society Series B*, 57, 289 - 300.  
Yuan-De Tan, Hongyan Xu; A general method for accurate estimation of false discovery rates in identification of differentially expressed genes, *Bioinformatics*, Volume 30, Issue 14, 15 July 2014, Pages 2018 - 2025.



**See Also**[p.adjust](#)**Examples**

```
set.seed(123)
x <- rnorm(50, mean = c(rep(0, 25), rep(3, 25)))
p <- 2*pnorm(sort(-abs(x)))
p.1.22<-mtpvadjust(pv=p, C=1.22)
```

---

myheatmap	<i>Heatmap</i>
-----------	----------------

---

**Description**

This function is used to display heatmap of differential expressions of genes or isoforms detected by NBBttest in the real count data.

**Usage**

```
myheatmap(dat, IDcol, nci, r, r1, r2,
  colrs="greenred", rwcex=2.0, clcex=2.0, x=10, tree="both",
  method="euclidean", ky=1.5, rowBarColor=NULL,
  colBarColor=NULL, labrow="yes", labcol="yes", rsort="yes",
  adjrow=c(0.3, 0.0), adjcol = c(1, 1), rwangle=30,
  clangle=30, maptitle)
```

**Arguments**

dat	is NBBttest object that is outputted by mbetattest, includes information columns, data columns, t-value column, pvalue column, selection column.
IDcol	nonnegative int value, indicating which column is ID column where IDs will be shown in heatmap rows, is required.
nci	nonnegative int value, indicating column number of information (gene or isoform annotation) columns in the data, is required.
r	nonnegative int value, indicating column number of count data including empty columns if there are multiple datasets with the same row names.
r1	nonnegative int value, indicating replicate number in condition 1.
r2	nonnegative int value, indicating replicate number in condition 2.
colrs	heatmap colors. User has 8 options: "redgreen", "greenred", "redblue", "bluered", "cm.colors", "terrain.colors", "topo.colors", and "heat.colors".
rwcex	numeric argument: nonnegative value, used as cex.axis in for the row axis labeling. Default=1.8.
clcex	numeric argument: nonnegative value, used as cex.axis for the column axis labeling. Default=1.8.
x	numeric argument: nonnegative value, used as argument of cm.colors(x), terrain.colors(x) and topo.colors(x), the default value is 10.

tree	tree(s) drawn on row or column or both. User has four options: "both" for drawing trees on both row and column, "row" for drawing tree only on row, "column" for drawing tree only on column, and "none" for no tree specified on rows and columns. Default is "both".
method	method to be chosen to calculate distance between columns or rows. It has four options: "euclidean", "pearson", "spearman" and "kendall". The latter three are $d=1-cc$ where $cc$ is correlation coefficients. Default is "euclidean".
ky	numeric argument: nonnegative value is used to determine key size. The default is 1.5
rowBarColor	(optional) character vector for RowSideColors and colRow. Length of rowBarColor equals to gene or isoform number or row number. rowBarColor contains the color names of classes or types or groups of row names and show row color side bars and color row names. If row names, for example, genes in row are not classified, then we suggest user use its default value: rowBarColor is NULL.
colBarColor	(optional) character vector for colSideColors and colCol. Length of colBarColor equals to sample number. colBarColor contains the color names of classes or types or groups of column names and show column color side bars and color column names. If column names, for example, samples in column are just two types, then we suggest user use its default value: colBarColor=NULL.
labrow	string, logistical value. Rows on heatmap are labeled with genes or targets if labrow="yes", otherwise, the rows are not labeled.
labcol	string, logistical value. columns on heatmap are labeled with samples or treatments if labcol="yes", otherwise, the columns are not labeled.
rsort	logistical value. If choose tree="both" or "row", then rsort does not work. However, if tree="none" or "column", then rsort="yes" will force rows to be sorted in descent way. The default is "yes"
adjrow	two numeric values. The first value is used to adjust left or right position (x-axis) of labels, and the second value is used to adjust up or down position (y-axis) of labels. The default values are c(0.4, 1).
adjcol	two numeric values. The first value is used to adjust left or right position (x-axis) of labels, and the second value is used to adjust up or down position (y-axis) of labels. The default values are c(1, 1).
rwangle	angle of xlab under heatmap. The Default value is 30.
clangle	angle of ylab. The default value is 30
maptitle	string argument for giving heatmap title, default is set to be empty string ''.

### Details

This function uses significance to choose data and then to normalize the selected data by using z-scale. This function has multiple options to select map color, distance, cluster and x- and y-lab angles. If  $r > (r1 + r2)$ , then data are multiple datasets.

### Value

no return value but create a heatmap.

### Note

requires gplots and grDevices.

**Author(s)**

Yuan-De Tan &lt;tanyuande@gmail.com&gt;

**See Also**[heatmap.2](#), [myheatmap2](#)**Examples**

```
data(result)
colclass=c("1","1","1","2","2","2")
oldpar <- par(no.readonly =TRUE)
on.exit(par(oldpar))
par(mar=c(7.5,5.5,3.5,1.2))
par(oma=c(3,1,1,3))
myheatmap (dat=result, IDcol=1, nci=7, r=6, r1=3,
r2=3, colrs="terrain.colors", rowBarColor=NULL,
colBarColor=colclass, labrow="no", labcol="yes",
rsort="yes", adjrow=c(0.3, 0.0 ),
adjcol = c(1, 1) , maptitle="My heatmap")
```

myheatmap2

*Heatmap2***Description**

This function is used to display heatmap of differential expressions of genes or isoforms detected by NBBttest in the real count data.

**Usage**

```
myheatmap2(dat, IDcol=1, nci=NULL, r, colrs=
"greenred", rwcex=2.8, clcex=2.8, x=10, tree=
"both", method="euclidean", ky=1.5, rowBarColor
=NULL, colBarColor=NULL, labrow="yes",
labcol="yes", rsort="yes", adjrow=c(0.2, 0.0 ),
adjcol = c(1, 1) , rwangle=30, clangle=30,
maptitle="",keyvalue)
```

**Arguments**

dat	object of NBBttest, including information columns, data columns, tvalue column, pvalue and selection column.
nci	int value, number of columns for information of genes or isoforms, like "geneid", ""tagetid"", "strainid", annotation etc.
IDcol	int value, indicating which column is ID column where IDs will be shown in heatmap rows.
r	int value, indicating column number of count data including empty columns if there are multiple datasets with the same row names.

colrs	heatmap colors. User has 8 options: "redgreen", "greenred", "redblue", "bluered", "cm.colors", "terrain.colors", "topo.colors", and "heat.colors". Default colrs is "redgreen".
rwcecx	numeric argument: nonnegative number, used as cex.axis for the row axis labeling. Default value is 1.8.
clcecx	numeric argument: nonnegative number, used as cex.axis for the column axis labeling. Default value is 1.8.
x	numeric argument: nonnegative number, used as argument of cm.colors(x), terrain.colors(x) and topo.colors(x), Default value is 10.
tree	tree(s) drawn on row or column or both. User has four options: "both" for drawing trees in both row and column, "row" for drawing tree only in row, "column" for drawing tree in only column, and "none" for no tree specified on rows and columns. If tree = "both", then columns and rows are sorted by trees. If tree = "row", the columns are not sorted. If tree = "column", then rows are not sorted. Default is "both".
method	method to be chosen to calculate distance between columns or rows. It has four options: "euclidean", "pearson", "spearman" and "kendall". The latter three are d=1-cc where cc is correlation coefficients. Default is "euclidean".
ky	numeric argument: nonnegative number, is used to determine key size. The default = 1.5.
rowBarColor	(optional) character vector for RowSideColors and colRow. Length of rowBarColor equals to gene or isoform number or row number. rowBarColor contains the color names of classes or types or groups of row names and show row color side bars and color row names. If row names, for example, genes in row are not classified, then we suggest user use its default value: rowBarColor=NULL.
colBarColor	(optional) character vector for colSideColors and colCol. Length of colBarColor equals to sample number. colBarColor contains the color names of classes or types or groups of column names and show column color side bars and color column names. If column names, for example, samples in column are just two types, then we suggest user use its default value: colBarColor=NULL.
labrow	a string, logistical value. Rows on heatmap are labeled with genes or targets if labrow="yes", otherwise, the rows are not labeled.
labcol	a string, logistical value. Columns on heatmap are labeled with samples or treatments if labcol="yes", otherwise, the columns are not labeled.
adjrow	two numeric values. The first value used to adjust left or right position (x-axis) of labels and the second value is used to adjust up or down position (y-axis) of labels. The default values are c(0.3, 0.0).
adjcol	two numeric values. The first value used to adjust left or right position (x-axis) of labels, and the second value is used to adjust up or down position (y-axis) of labels. The default values are c(1, 1).
rsort	logistical value. If choose tree="both" or "row", then rsort does not work. However, if tree="none" or "column", then rsort="yes" will force rows to be sorted in descent way. The default is "yes"
rwangle	heatmap object: angle of xlab. The default value is 30.
clangle	heatmap object: angle of ylab. The default value is 30.
maptitle	string argument for giving heatmap title. The default value is set to be empty string ''.
keyvalue	string argument for name of key x-axis.

## Details

This function uses selection to choose genes or isoforms in the data and then to normalize the selected data by using n-scale. Different from z-score, n-score does not follow standard normal distribution with mean = 0 and variance =1 for all rows but it has the same largest count in all rows and shows multiple colors for numeric difference between two conditions. This function has multiple options to select map color, distance, cluster and x-lab and y-lab angles. This function can be able to display multiple datasets in two ways: if multiple datasets have the same row names or features, the these datasets are put onto the different columns separated with empty column named with dataset names. If multiple datasets have the same column names of the datasets, then put them on different rows separated with empty rows named with dataset names or whatever names user specifies.

## Value

not return value but create a heatmap.

## Note

requires gplots and grDevices. If the data for heatmap are multiple datasets, then tree="none" and sort="no", otherwise, myheatmap2 will get error. So before performing myheatmap2, user should sort the data in excel.

## Author(s)

<tanyuande@gmail.com>

## See Also

[heatmap.2](#), [grDevices](#), and [myheatmap](#).

## Examples

```
data(result)
colclass=c("1","1","1","2","2","2")
oldpar <- par(no.readonly =TRUE)
on.exit(par(oldpar))
par(mar=c(7.5,5.5,3.5,1.2))
par(oma=c(3,1,1,3))
myheatmap2(dat=result, IDcol=1, nci=7, r=6,
  colrs="greenred", rwcex=1.8, clcex=1.8, x=10,
  tree="both", method="euclidean", ky=1.5,
  rowBarColor=NULL, colBarColor=colclass,
  labrow="no", labcol="yes", adjrow=c(0.2, 0.0 ),
  adjcol = c(1, 1) , rwangle=0, clangle=30,
  maptitle="My heatmap2",keyvalue="count")
```

---

NBBplot

---

*Plot differential expression of exons within a specified gene using result outputed by NBBttest.*


---

### Description

After performing NBBttest, NBBplot can be used to show differential expression of exons within a specified gene in na and nb replicates between conditions A and B.

### Usage

```
NBBplot(res, gene, nci, na, nb, C1, C2)
```

### Arguments

res	object of NBBttest containing information of genes including gene name, strand, chromosome, exons, introns, data, and t-value, p-value, significance/selection etc.
gene	gene name or symbol specified by user.
nci	number of columns for gene information.
na	replicate number in condition A
nb	replicate number in condition B
C1	name for condition 1(A)
C2	name for condition 2(B)

### Details

NBBplot consists of two parts: top is expression value of each exon in each replicate in two conditions marked in red and blue and bottom is boxes for exon and solid lines for introns. Differential expression of an exon is marked in red.

### Value

output NBBplot figure for given gene.

### Author(s)

Yuan-De Tan <tanyuande@gmail.com>

### See Also

[plotDEXSeq](#)

### Examples

```
data(exondata)
```

```
NBBplot(res=exondata, gene="H2-DMb1", nci=9, na=3, nb=3, C1="WT", C2="KO")
```

---

normalized	<i>Normalization of data</i>
------------	------------------------------

---

### Description

Function normalize makes all libraries in dataset have the same library size.

### Usage

```
normalized(dat, nci, m=0, lg2="no")
```

### Arguments

dat	count data of RNA reads.
nci	number of columns for the information of genes or isoforms in dataset.
m	numeric value for choosing genes or isoforms. If user wants to discard genes or isoforms with mean < 5, then m = 5. The default value is 0.
lg2	logistic value. lg2="yes" indicates that data are transformed in logarithm of 2.

### Details

Due to difference in RNA abstraction between libraries or cell samples or tissues, PCR amounts of RNA libraries would have difference that is not due to biological effects. To correctly compare differential expressions of genes between conditions or samples, one must should give the same RNA abstraction in all given samples. This is impossible. To address this problem, only one way is to normalize these count data across all given samples so that all experimental samples (libraries) have the same total counts.

### Value

output a standard datasheet.

### Author(s)

Yuan-De Tan <tanyuande@gmail.com>

### References

Yuan-De Tan Anita M. Chandler, Arindam Chaudhury, and Joel R. Neilson(2015) A Powerful Statistical Approach for Large-scale Differential Transcription Analysis. *Plos One*, 10.1371/journal.pone.0123658.  
Anders S, Huber W (2010) Differential expression analysis for sequence count data. *Genome Biol*, 11: R106.

### Examples

```
data(jkttcell)
njkttcell<-normalized(dat=jkttcell[1:50,],nci=7)
```

---

oddratio	<i>Calculation of zeta (<math>\zeta</math>)</i>
----------	---

---

### Description

$\zeta$  is used to measure homogeneity intensity of two subdatasets. If  $\zeta$  is larger than 1, these two subdatasets have good homogeneity; otherwise,  $\zeta < 1$  indicates that two subdatasets have poor homogeneity (big noise).

### Usage

```
oddratio(XX, na, nb)
```

### Arguments

XX	nonnegative count data of RNA reads generated by next generation sequencing.
na	number of replicate libraries in condition A.
nb	number of replicate libraries in condition B.

### Details

Zeta ( $\zeta$ ) is defined as

$$\zeta = \ln \left( 1 + \frac{\bar{X} \times \sigma^2 + 1}{\bar{X}_A \times \sigma_A^2 + \bar{X}_B \times \sigma_B^2 + 1} \right)$$

where  $\zeta$  is different from  $\psi$ . If two subdatasets have big a gap and good homogeneity, then set a value has much larger than 1.

### Value

zeta vector.

### Author(s)

Yuan-De Tan <tanyuande@gmail.com>

### References

Yuan-De Tan Anita M. Chandler, Arindam Chaudhury, and Joel R. Neilson(2015) A Powerful Statistical Approach for Large-scale Differential Transcription Analysis. *Plos One*. 2015 DOI: 10.1371/journal.pone.0123658.

### See Also

[pratio](#), [mbetatest](#).

### Examples

```
XX<-matrix(NA,2,8)
XX[1,]<-c(112,122, 108,127,302, 314, 322, 328)
XX[2,]<-c(511, 230, 754, 335,771, 842, 1014,798)
oddratio(XX=XX,na=4,nb=4)
```



---

omega	<i>Omega calculation</i>
-------	--------------------------

---

### Description

Omega function is a function that is used to estimate omega using simulate null data from negative binomial distribution. is a null rho that is used as a threshold for real rho. Simulation is dependent on the original data.

### Usage

```
omega(XX, nci, r1, r2, sn, alpha = 0.05)
```

### Arguments

XX	the real dataset.
nci	number of columns for information of data, like gene id, isoform id, gene name, etc.
r1	size of sample 1 or number of replicates in condition 1.
r2	size of sample 2 or number of replicates in condition 2.
sn	number of simulations specified.
alpha	significance level of test. Default value is 0.05.

### Details

This function is to use null data to calculate omega value with rho = 1.

### Value

return a numeric value.

### Author(s)

Yuan-De Tan <tanyuande@gmail.com>

### References

Yuan-De Tan Anita M. Chandler, Arindam Chaudhury, and Joel R. Neilson(2015) A Powerful Statistical Approach for Large-scale Differential Transcription Analysis. *Plos One*. 2015 DOI: 10.1371/journal.pone.0123658.

### See Also

[pratio](#) and [oddratio](#)

### Examples

```
data(jkttcell)

w<-omega(XX=jkttcell[1:100,],nci=7,r1=3,r2=3,sn=2,alpha=0.05)
```

---

pathwayHeatmap

*Heatmap for pathways found by gene ontology analysis*


---

### Description

This function is used to show differential expressions of pathways or functions between conditions. These pathways or functions were detected by function annotation or gene ontology methods such as David function analysis tools or Ingenuity pathway analysis. Pathway score or pathway value is a weighted expression value across genes in a pathway or a function. The weights of genes are given by p-values of enrichment or hit in function analysis.

### Usage

```
pathwayHeatmap(dat, pathway, nci, r1, r2, colclass, rowclass,
               colrs, maptitle)
```

### Arguments

dat	count dataset that contains a column for gene name and count data columns in two conditions.
pathway	a list that lists a function column and at least two gene columns. Row is function name.
nci	number of columns for gene information such as gene name, strand, chromosome, etc in dataset dat.
r1	number of replicates in condition 1 in dataset dat.
r2	number of replicates in condition 2 in dataset dat.
colclass	column class, equivalent to replicates in two conditions, such as c(1,1,1,1,2,2,2,2), meaning that condition A has 4 replicates and condition B has 4 replicates. Default = NULL.
rowclass	row class, pathway class, or pathway module. Default = NULL.
colrs	heatmap colors. User has 8 options: "redgreen", "greenred", "redblue", "bluered", "cm.colors", "terrain.colors", "topo.colors", and "heat.colors".
maptitle	heatmap title, default = "".

### Details

This function uses significance to choose gene function/ pathway data and then to normalize the selected data by using z-scale. This function has multiple options to select map color, distance, cluster and x- and y-axis angles.

### Value

return a graph for heatmap.

### Note

requires gplots and grDevices.

**Author(s)**

Yuan-De Tan <tanyuande@gmail.com>

**See Also**

[heatmap.2](#), [myheatmap](#), [myheatmap2](#)

**Examples**

```
data(upGAm)
data(pathwy.A.up)
pathwayup<-pathwy.A.up
colclass=c("1","1","1","1","2","2","2","2","2","2")
oldpar <- par(no.readonly =TRUE)
on.exit(par(oldpar))
par(mar=c(7.5,5.5,3.5,3))
par(oma=c(3,1,1,10))
pathwayHeatmap(dat=upGAm,pathway=pathwayup,nci=1,
r1=4,r2=6,colclass=colclass,rowclass=NULL,
colrs="greenred",maptitle="pathway up-expression
in Group A")
```

---

pathwy.A.up

*Pathway or function data*

---

**Description**

pathway data are functions or pathways of up-regulated genes given by gene annotation tool such as DAVID gene function annotation.

**Usage**

```
data("pathwy.A.up")
```

**Format**

A data frame with 39 observations on the following 159 variables.

39 observations, 159 variables

**Details**

The 39 observations are 39 functions/pathways listed each being a set up-regulated genes found by NBBttest in knockdowned DDX39 cell line. Log10 of pvalue of gene in differential expression detection is used to weight compoment of this gene playing role in this function. Then heatmap of pathways or functions are made using myheatmap. V1-V159 are gene sets of function or pathways

**Examples**

```
data(pathwy.A.up)
```

---

pratio	<i>Calculation of psi (<math>\psi</math>)</i>
--------	---

---

### Description

$\psi$  is also called polar ratio. It is used to measure overlap of two subdatasets. If  $\psi > 1$ , these two subdatasets have a gap, while  $\psi < 1$  indicates that two subdatasets overlap.

### Usage

```
pratio(xx, na, nb)
```

### Arguments

xx	count data of RNA reads generated by next generation sequencing.
na	number of replicate libraries in condition A.
nb	number of replicate libraries in condition B.

### Details

Psi ( $\psi$ ) is defined as

$$\psi = \max\left(\frac{\min(X_A)}{\max(X_B)}, \frac{\min(X_B)}{\max(X_A)}\right)$$

### Value

pratio	pratio list
--------	-------------

### Author(s)

Yuan-De Tan <tanyuande@gmail.com>

### References

Yuan-De Tan Anita M. Chandler, Arindam Chaudhury, and Joel R. Neilson(2015) A Powerful Statistical Approach for Large-scale Differential Transcription Analysis. *Plos One*. 2015 DOI: 10.1371/journal.pone.0123658.

### See Also

[omega](#), and [oddratio](#)

### Examples

```
XX<-matrix(NA,2,8)
XX[1,]<-c(112,122, 108,127,302, 314, 322, 328)
XX[2,]<-c(511, 230, 754, 335,771, 842, 1014,798)
pratio(xx=XX,na=4,nb=4)
```

---

prime3_PRP8_50	<i>3'UTR splicing data of 50 genes detected in the knockdowned PRP8 cell line</i>
----------------	---

---

### Description

Splicing events occurring in 3'UTR of 50 genes were detected in the knockdowned PRP8 cell line.

### Usage

```
data("prime3_PRP8_50")
```

### Format

A data frame with 50 observations on the following 24 variables.

chr a factor with levels  
strand a factor with levels  
isoform a factor with levels  
Gene a factor with levels  
exon\_const\_start a numeric vector  
exon\_const\_end a numeric vector  
exon\_alt1\_start a numeric vector  
exon\_alt1\_end a numeric vector  
exon\_alt2\_start a numeric vector  
exon\_alt2\_end a numeric vector  
Siha\_Ctl\_I\_S3 a numeric vector  
Siha\_Ctl\_II\_S7 a numeric vector  
Siha\_Ctl\_III\_S11 a numeric vector  
PRP8\_I\_S2 a numeric vector  
PRP8\_II\_S6 a numeric vector  
PRP8\_III\_S10 a numeric vector  
tvalue a numeric vector  
rho a numeric vector  
pvalue a numeric vector  
adjp a numeric vector  
w a numeric vector  
X a numeric vector  
X.1 a numeric vector  
X.2 a numeric vector

Details

This dataset is an example for annotation of genes within which splicing events occurring in 3'UTR were detected by spladder-NBBttest in knockdowned PRP8 cell line. Gene information columns are assigned to "chr", "Gene", "isoform", "exon\_const\_start" and "exon\_const\_end" that are start and end positions of constituted exons and "exon\_alt1\_start" and "exon\_alt1\_end" that are start and end positions of alternative exons. Count data of RNA reads are assigned to three replicate control columns (Siha\_Ctl\_I\_S3,II-S7, III\_S11) and three replicate knockdowned PRP8 columns (PRP8\_I\_S2,II\_S6 and III\_S10). The tvalue is t-statistic, rho is gene-wise variable, pvalue is p-value for t-statistic, w is  $\omega$ , a threshold for rho, FDR (X.1) is false discovery rate with  $\alpha=0.05$ , selection (X.2) = 1 if pvalue < FDR, 0, otherwise.

Examples

```
data(prime3_PRP8_50)
```

QC	<i>Count data quality check</i>
----	---------------------------------

Description

QC is used to check quality of count data in two ways: scatter plot of two replicated samples and correlation heatmap of all samples.

Usage

```
QC(dat, nci, S1="NULL", S2="NULL", method="plot", colrs="greenred", rwcex=1.8, clcex=1.8, x=10, tree="none", log="none", col="blue", pch=19, labsize=1.5, axis=1.5)
```

Arguments

dat	count matric dataset.
nci	number of columns containing data information suach as gene id, library id, target id, gene name, strand etc.
S1	numeric int value, indicating which column in data matrix is specified in x-axis. S1 > nci.
S2	numeric int value, indicating which column in data matrix is specified in y-axis. S2 > nci.
method	string. Here two methods are given for choice: "plot" and "heatmap". Default is "plot".
colrs	string. 8 color sets are given for choice in this function: "redgreen", "heat.colors", "redblue", "greenred", "bluered", "cm.colors", "terrain.colors", "topo.colors". The default color set is "redgreen".
rwcex	positive numbers, used as cex.axis for the row axis labeling. The default value is 1.8.
clcex	positive numbers, used as cex.axis for the column axis labeling. The default value is 1.8.

x	numeric argument: positive number, used as argument of <code>cm.colors(x)</code> , <code>terrain.colors(x)</code> and <code>topo.colors(x)</code> , the default value is 10.
tree	tree(s) drawn on row or column or both. User has four options: "both" for drawing trees on both row and column, "row" for drawing tree only on row, "column" for drawing tree on only column, and "none" for no tree specified. Default is "none".
log	string. Two options are given for choice: "none" and "log". <code>log="log"</code> indicates that data value would be transformed with <code>log2</code> . Default is "none".
col	string, used to specify scatter plot dot color.
pch	numeric value, used to specify dot type.
labsize	numeric value for size of xlabel and ylabel
axis	numeric value for axis scale.

### Details

S1, S2 and nci must be given numeric int values for plot and heatmap. However, when method is chosen to be "heatmap", then S1 and S2 are not specified. Columns of information should be left of matrix and count data should be after columns of information.

### Value

not return values but create scatter plot or heatmap plot.

### Note

requires `gplots` and `grDevices`.

### Author(s)

Yuan-De Tan <tanyuande@gmail.com>

### See Also

[heatmap.2](#) and [grDevices](#)

### Examples

```
data(jkttcell)
QC(dat=jkttcell, nci=7, S1=8, S2=9, method = "plot",
  log = "log", col = "blue", pch = 19)
QC(dat=jkttcell, nci=7, S1=8, S2=9, method = "plot",
  log = "log", col = "blue", pch = 19)
QC(dat=jkttcell, nci=7, method = "heatmap", log = "log")
```

result

*Jurkat T-cell transcriptomic data with isoforms selected by NBBttest***Description**

A data consist of 7 information columns and 6 data (numeric) columns and 1953 RNA isoforms selected by NBBttest are used to make heatmaps.

**Usage**

```
data("result")
```

**Format**

A data frame with 1953 observations on the following 14 variables.

tagid : a numeric vector  
 geneid : a numeric vector  
 name : a string vector  
 chr : a string vector, a set of 24 chromosomes  
 strand : a factor with levels - +  
 pos : a numeric vector  
 anno : string vector, poly(A) site types  
 Jurk.NS.A : a numeric vector  
 Jurk.NS.B : a numeric vector  
 Jurk.NS.C : a numeric vector  
 Jurk.48h.A : a numeric vector  
 Jurk.48h.B : a numeric vector  
 Jurk.48h.C : a numeric vector  
 tvalue : a numeric vector

**Details**

The original data are count data generated by next generation sequencing from Jurkat T-cells. This dataset is a short dataset containing 1953 isoforms that show differential expression between rest status and stimulation status. The T-cells were treated by resting and stimulating with *CD3/CD28* for 48 hours. The data have 7 columns for the information of *poly(A)* site: "tagid", "geneid", "gene name", "chromosome", "strand", *poly(A)* site position, *poly(A)* site annotation and 6 columns for data: Jurk.NS.A, Jurk.NS.B, Jurk.NS.C, Jurk.48h.A, Jurk.48h.B, Jurk.48h.C. where NS means Normal state or no stimulation and 48h means 48 hours after *CD3/CD28* stimulation of T-cells. 13409 RNA isoforms were detected to have alternative *poly(A)* sites.

**Value**

datasheet contained ID, information, count data of RNA reads.



## References

Yuan-De Tan Anita M. Chandler, Arindam Chaudhury, and Joel R. Neilson(2015) A Powerful Statistical Approach for Large-scale Differential Transcription Analysis. *Plos One*. DOI: 10.1371/journal.pone.0123658.

## Examples

```
data(result)
```

---

sgRNA	<i>sgRNA dataset</i>
-------	----------------------

---

## Description

This dataset was created by simulating single guide RNAs to edit genes.

## Usage

```
data("sgRNA")
```

## Format

A data frame with 1000 observations on the following 11 variables.

gene a factor with levels  
sgRNA a factor with levels  
class a factor with levels  
L1 a numeric vector  
L2 a numeric vector  
L3 a numeric vector  
L4 a numeric vector  
H1 a numeric vector  
H2 a numeric vector  
H3 a numeric vector  
H4 a numeric vector

## Details

The dataset is CRISPR screening data. Each gene was edited by 10 sgRNAs that contains a targeting sequence (crRNA sequence) and a Cas9 nuclease-recruiting sequence (tracrRNA). L1-L4 samples were lowly targeted by sgRNAs and H1-H4 were highly targeted by sgRNAs.

## Examples

```
data(sgRNA)
```

---

simSplicing	<i>Simulated alternative splicing</i>
-------------	---------------------------------------

---

### Description

This alternative splicing count data were created by simulating dorsal and ventral RNA-sequence count data.

### Usage

```
data("simSplicing")
```

### Format

A data frame with 5000 observations on the following 9 variables.

Isoform a factor with levels

geneid a factor with levels

label a numeric vector

D1 a numeric vector

D2 a numeric vector

D3 a numeric vector

V1 a numeric vector

V2 a numeric vector

V3 a numeric vector

### Details

D1-D3 are samples from mouse dorsal tissue and V1-V3 are samples from mouse ventral tissue. This dataset was simulated with negative binomial distribution with 10 percent of isoforms of being differentially spliced and differential effect of 500U where U is uniform variable.

### Examples

```
data(simSplicing)
```

---

simulat	<i>Simulation</i>
---------	-------------------

---

### Description

This function uses negative binomial (NB) pseudorandom generator to create count datasets of RNA isoform reads based on real data.

### Usage

```
simulat(yy, nci, r1, r2, p, q, A)
```

**Arguments**

yy	real count data
nci	nonnegative int value: column number of information related to genes or isoforms.
r1	numeric argument: number of replicate libraries in condition 1.
r2	numeric argument: number of replicate libraries in condition 2.
p	numeric argument: proportion of genes or isoforms differentially expressed. The value is in range of 0~1. Default is 0.
q	numeric argument: proportion of genes or isoforms artificially noised. The value is in range of 0~1. Default is 0.
A	numeric argument: conditional effect value. The value is larger than or equal to 0. Default is 0.

**Details**

Null count data are created by using R negative binomial pseudorandom generator `rnbinom` with  $\mu$  and size. Parameters  $\mu$  and size are given by mean and variance drawn from real read counts of a gene set or an isoform set in a condition. Condition (or treatment) effect on differential transcription of isoforms is linearly and randomly assigned to genes or isoforms. The conditional effect = AU where U is uniform variable and A is input constant. P percent of genes or isoforms is set to be differentially expressed or differentially spliced. Q percent of genes or isoforms has technical noise. If P = 0, then simulation is null simulation, the data are null data or baseline data.

**Value**

Return count data.

**Author(s)**

Yuan-De Tan <tanyuande@gmail.com>

**References**

Yuan-De Tan Anita M. Chandler, Arindam Chaudhury, and Joel R. Neilson(2015) A Powerful Statistical Approach for Large-scale Differential Transcription Analysis. *Plos One*, 10.1371/journal.pone.0123658.

**See Also**

[NegBinomial](#)

**Examples**

```
data(jkttcell)
jknull<-simulat(yy=jkttcell[1:60,],nci=7,r1=3,
r2=3,p=0,q=0.2,A=0)
```

skjt

*Simulated Null Transcriptomic data***Description**

The dataset skjt generated by using R negative binomial pseudorandom generator rnbinom is used as an example for calculating omega.

**Usage**

```
data("skjt")
```

**Format**

A data frame with 13409 observations on the following 14 variables.

```
geneid : a string vector
tagid : a numeric vector
geneid.1 : a numeric vector
name : a string vector
chr : a string vector
strand : a character vector
pos : a numeric vector
anno : a string vector
Jurk.NS.A : a numeric vector
Jurk.NS.B : a numeric vector
Jurk.NS.C : a numeric vector
Jurk.48h.A : a numeric vector
Jurk.48h.B : a numeric vector
Jurk.48h.C : a numeric vector
```

**Details**

The dataset skjt was generated by using R negative binomial pseudorandom generator rnbinom with  $\mu$  and size. Parameters  $\mu$  and size are given by mean and variance drawn from real Jurkat T cell transcriptomic count data . Condition (or treatment) effect on differential transcription of isoforms was set to zero. The data have 13409 genes and 7 information columns: geneid tagid name chr,strand, pos,anno, and 6 data columns.

**Value**

a datasheet contained ID, information, count data of RNA reads.

**References**

Yuan-De Tan Anita M. Chandler, Arindam Chaudhury, and Joel R. Neilson(2015) A Powerful Statistical Approach for Large-scale Differential Transcription Analysis. *Plos One*. DOI: 10.1371/journal.pone.0123658.

**Examples**

```
data(skjt)
```

---

smbetattest

*Performance of multiple beta t-test on simulated data*


---

**Description**

This function is to perform beta t-test with  $\rho = 1$  and  $\omega = 1$  on simulated data. The result lists differentially expressed genes or isoforms and their  $\rho$  values. The  $\rho$  values are used to calculate  $\omega$  value for performance of beta t-tests on the real data.

**Usage**

```
smbetattest(X, na, nb, alpha)
```

**Arguments**

X	simulated count data with N genes or isoforms.
na	number of replicate libraries in condition A.
nb	number of replicate libraries in condition B.
alpha	statistical probabilistic threshold, default is 0.05.

**Details**

Before performing NBBttest on real data, user needs  $\omega$  value for the threshold of  $\rho$ . To determine  $\omega$  value, user is required to generate a set of null data having the same gene or isoform number and the same numbers of replicate libraries in two conditions and then performs beta t-test on the null datasets by setting  $\rho = 1$  and  $\omega = 1$ . In current package, NBBttest can automatically perform the simulation of null data, multiple beta t-test to estimate  $\omega$ .

**Value**

Return a set of null  $\rho$  values.

**Author(s)**

Yuan-De Tan <tanyuande@gmail.com>

**References**

Yuan-De Tan Anita M. Chandler, Arindam Chaudhury, and Joel R. Neilson(2015) A Powerful Statistical Approach for Large-scale Differential Transcription Analysis. *Plos One*. DOI: 10.1371/journal.pone.0123658.

**See Also**

See Also as [mbetattest](#)

**Examples**

```
data(skjt)
nrho<-smbetattest(X=skjt[1:60,],na=3,nb=3,alpha=0.05)
```

---

subdata	<i>Split data into two subsets</i>
---------	------------------------------------

---

**Description**

Data are splited into two subsets: gene single-isoform data and gene multi-isoform data.

**Usage**

```
subdata(xx, sg)
```

**Arguments**

xx	real data containing single-isoforms and multi-isoforms of genes in rows.
sg	int value.

**Details**

For the RNA count data, some genes have only one isoform, some genes have multiple isoforms. so data are divided by subdata.R into two datasets: single-isoform data and multi-isoform data.

**Value**

return dataset by setting sg=1 or sg=2.

**Author(s)**

Yuan-De Tan <tanyuande@gmail.com>

**Examples**

```
data(jkttcell)
colnames(jkttcell)[3]<-"Gene"
jk.mg<-subdata(xx=jkttcell, sg=2)
```

upGAm

*Count data of group A treated breast cancer in mice***Description**

Data of RNA-seq reads were obtained from drug treated breast cancer in mice and mapped by using STAR onto mm10 to create count data. Dataset upGAm is output of NBBttest.

**Usage**

```
data("upGAm")
```

**Format**

A data frame with 263 observations on the following 19 variables.

gene a factor with 263 DE genes

A.2\_S35S a numeric vector

A.4\_S36S a numeric vector

A.42\_S45S a numeric vector

A.39\_S41S a numeric vector

A.9\_S44R a numeric vector

A.12\_S40R a numeric vector

A.18\_S37R a numeric vector

A.29\_S39R a numeric vector

A.31\_S38R a numeric vector

A.38\_S34R a numeric vector

tvalue a numeric vector

rho a numeric vector

pvalue a numeric vector

adjp a numeric vector

w a numeric vector

order a numeric vector

FDR a numeric vector

significance a numeric vector

**Details**

This dataset is a demo dataset having 263 DE genes identified by NBBttest. The dataset was created from types of two tumor cells: sensitive to drug (S) and resistant to drug (R). 4 sensitive cells and 6 resistant cells are available for differential analysis. Tvalue is t-statistic,  $\rho$  is gene-wise variable, pvalue is p-value for t-test,  $\omega$  is threshold for  $\rho$ , FDR is false discovery rate, significance = 1 if pvalue < FDR, 0, otherwise.

**Examples**

```
data(upGAm)
```

# Index

- \* **QC**
  - QC, 30
- \* **adjust pvalue**
  - mtpvadjust, 16
- \* **alpha**
  - betaparametab, 4
- \* **annotation**
  - annotat, 3
- \* **beta and negative binomial**
  - mbetattest, 14
- \* **beta distribution**
  - betattest, 8
- \* **beta t-test**
  - gbetattest, 11
- \* **beta**
  - betaparametab, 4
- \* **binomial**
  - simulat, 34
- \* **datasets**
  - DDX39\_100, 9
  - exondata, 10
  - gtfa, 12
  - jktcell, 13
  - pathwy.A.up, 27
  - prime3\_PRP8\_50, 29
  - result, 32
  - sgRNA, 33
  - simSplicing, 34
  - skjt, 36
  - upGAm, 39
- \* **data**
  - subdata, 38
- \* **differential expressions**
  - myheatmap2, 19
- \* **differential expression**
  - NBBplot, 22
- \* **exons**
  - NBBplot, 22
- \* **gap**
  - oddratio, 24
  - pratio, 28
- \* **gene**
  - annotat, 3
- \* **group**
  - gbetattest, 11
- \* **heatmap**
  - myheatmap, 17
  - myheatmap2, 19
  - pathwayHeatmap, 26
- \* **homogeneity**
  - oddratio, 24
- \* **library size**
  - normalized, 23
- \* **multiple test**
  - mtpvadjust, 16
- \* **negative**
  - simulat, 34
- \* **normalize**
  - normalized, 23
- \* **omega**
  - omega, 25
- \* **overlap**
  - pratio, 28
- \* **package**
  - NBBttest-package, 2
- \* **pathway**
  - pathwayHeatmap, 26
- \* **proportion**
  - betaparametVP, 5
- \* **rho**
  - omega, 25
- \* **scatter plot**
  - QC, 30
- \* **simulation**
  - simulat, 34
  - smbetattest, 37
- \* **split**
  - subdata, 38
- \* **t-tests**
  - mbetattest, 14
- \* **t-test**
  - smbetattest, 37
- \* **t-value**
  - betattest, 8
- \* **variance**
  - betaparametVP, 5



**\* weight**

betaparametw, 7

annotat, 3

betaparametab, 3, 4, 6, 7

betaparametVP, 3, 5, 5, 7

betaparametw, 3, 5, 6, 7

betattest, 3, 8, 12

DDX39\_100, 9

exondata, 10

gbetattest, 3, 11

grDevices, 21, 31

gtfa, 12

heatmap.2, 19, 21, 27, 31

jkttcell, 13

mbetattest, 3, 14, 24, 37

mtpvadjust, 16, 16

myheatmap, 3, 17, 21, 27

myheatmap2, 3, 19, 19, 27

NBBplot, 3, 22

NBBtatest (NBBttest-package), 2

NBBtatest-package (NBBttest-package), 2

NBBttest-package, 2

NegBinomial, 35

normalized, 3, 16, 23

oddratio, 3, 9, 24, 25, 28

omega, 3, 16, 25, 28

p.adjust, 17

pathwayHeatmap, 3, 26

pathwy.A.up, 27

plotDEXSeq, 22

pratio, 3, 9, 24, 25, 28

prime3\_PRP8\_50, 29

QC, 3, 30

result, 32

sgRNA, 33

simSplicing, 34

simulat, 3, 34

skjt, 36

smbetattest, 3, 16, 37

subdata, 38

upGAm, 39