

# Package ‘genefilter’

April 16, 2019

**Title** genefilter: methods for filtering genes from high-throughput experiments

**Version** 1.64.0

**Author** R. Gentleman, V. Carey, W. Huber, F. Hahne

**Description** Some basic functions for filtering genes.

**Maintainer** Bioconductor Package Maintainer

<maintainer@bioconductor.org>

**Suggests** class, hgu95av2.db, tkWidgets, ALL, ROC, DESeq, pasilla, RColorBrewer, BiocStyle, knitr

**Imports** S4Vectors (>= 0.9.42), AnnotationDbi, annotate, Biobase, graphics, methods, stats, survival

**License** Artistic-2.0

**LazyLoad** yes

**LazyData** yes

**Collate** AllClasses.R AllGenerics.R all.R dist2.R eSetFilter.R fastT.R filter\_volcano.R filtered\_p.R genefinder.R half.range.mode.R kappa\_p.R nsFilter.R rejection\_plot.R rowROC-accessors.R rowSds.R rowpAUCs-methods.R rowttests-methods.R shorth.R zzz.R

**biocViews** Microarray

**VignetteBuilder** knitr

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

**git\_branch** RELEASE\_3\_8

**git\_last\_commit** 82e91b7

**git\_last\_commit\_date** 2018-10-30

**Date/Publication** 2019-04-15

## R topics documented:

Anova	2
coxfilter	3
cv	4
dist2	5
eSetFilter	6
filtered_p	7

filterfun	8
filter_volcano	9
findLargest	10
gapFilter	11
genefilter	12
genefinder	13
genescale	14
half.range.mode	15
kappa_p	17
kOverA	18
maxA	18
nsFilter	19
pOverA	22
rejection_plot	23
rowFtests	24
rowpAUCs-methods	27
rowROC-class	29
rowSds	30
shorth	31
tdata	33
ttest	33

<b>Index</b>	<b>35</b>
--------------	-----------

---

Anova

*A filter function for Analysis of Variance*


---

## Description

Anova returns a function of one argument with bindings for cov and p. The function, when evaluated, performs an ANOVA using cov as the covariate. It returns TRUE if the p value for a difference in means is less than p.

## Usage

```
Anova(cov, p=0.05, na.rm=TRUE)
```

## Arguments

cov	The covariate. It must have length equal to the number of columns of the array that Anova will be applied to.
p	The p-value for the test.
na.rm	If set to TRUE any NA's will be removed.

## Details

The function returned by Anova uses lm to fit a linear model of the form  $lm(x \sim cov)$ , where x is the set of gene expressions. The F statistic for an overall effect is computed and if it has a p-value less than p the function returns TRUE, otherwise it returns FALSE for that gene.

**Value**

Anova returns a function with bindings for cov and p that will perform a one-way ANOVA.  
The covariate can be continuous, in which case the test is for a linear effect for the covariate.

**Author(s)**

R. Gentleman

**See Also**

[kOverA](#), [lm](#)

**Examples**

```
set.seed(123)
af <- Anova(c(rep(1,5),rep(2,5)), .01)
af(rnorm(10))
```

---

coxfilter

*A filter function for univariate Cox regression.*

---

**Description**

A function that performs Cox regression with bindings for surt, cens, and p is returned. This function filters genes according to the attained p-value from a Cox regression using surt as the survival times, and cens as the censoring indicator. It requires survival.

**Usage**

```
coxfilter(surt, cens, p)
```

**Arguments**

surt	Survival times.
cens	Censoring indicator.
p	The p-value to use in filtering.

**Value**

Calls to the [coxph](#) function in the survival library are used to fit a Cox model. The filter function returns TRUE if the p-value in the fit is less than p.

**Author(s)**

R. Gentleman

**See Also**

[Anova](#)

## Examples

```
set.seed(-5)
sfun <- coxfilter(rexp(10), ifelse(runif(10) < .7, 1, 0), .05)
ffun <- filterfun(sfun)
dat <- matrix(rnorm(1000), ncol=10)
out <- genefilter(dat, ffun)
```

---

cv

*A filter function for the coefficient of variation.*

---

## Description

cv returns a function with values for a and b bound. This function takes a single argument. It computes the coefficient of variation for the input vector and returns TRUE if the coefficient of variation is between a and b. Otherwise it returns FALSE

## Usage

```
cv(a=1, b=Inf, na.rm=TRUE)
```

## Arguments

a	The lower bound for the cv.
b	The upper bound for the cv.
na.rm	If set to TRUE any NA's will be removed.

## Details

The coefficient of variation is the standard deviation divided by the absolute value of the mean.

## Value

It returns a function of one argument. The function has an environment with bindings for a and b.

## Author(s)

R. Gentleman

## See Also

[pOverA](#), [kOverA](#)

## Examples

```
set.seed(-3)
cvfun <- cv(1,10)
cvfun(rnorm(10,10))
cvfun(rnorm(10))
```

---

dist2	<i>Calculate an n-by-n matrix by applying a function to all pairs of columns of an m-by-n matrix.</i>
-------	---

---

### Description

Calculate an n-by-n matrix by applying a function to all pairs of columns of an m-by-n matrix.

### Usage

```
dist2(x, fun, diagonal=0)
```

### Arguments

x	A matrix.
fun	A symmetric function of two arguments that may be columns of x.
diagonal	The value to be used for the diagonal elements of the resulting matrix.

### Details

With the default value of fun, this function calculates for each pair of columns of x the mean of the absolute values of their differences (which is proportional to the L1-norm of their difference). This is a distance metric.

The implementation assumes that fun(x[, i], x[, j]) can be evaluated for all pairs of i and j (see examples), and that fun is symmetric, i.e. fun(a, b) = fun(b, a). fun(a, a) is not actually evaluated, instead the value of diagonal is used to fill the diagonal elements of the returned matrix.

Note that `dist` computes distances between rows of x, while this function computes relations between columns of x (see examples).

### Value

A symmetric matrix of size n x n.

### Author(s)

Wolfgang Huber, James Reid

### Examples

```
# example matrix
z = matrix(1:15693, ncol=3)
matL1 = dist2(z)
matL2 = dist2(z, fun=function(a,b) sqrt(sum((a-b)^2, na.rm=TRUE)))

euc = as.matrix(dist(t(z)))

stopifnot(identical(dim(matL2), dim(euc)),
           all(euc==matL2))
```

---

`eSetFilter`*A function to filter an eSet object*

---

**Description**

Given a Bioconductor's ExpressionSet object, this function filters genes using a set of selected filters.

**Usage**

```
eSetFilter(eSet)
getFilterNames()
getFuncDesc(lib = "genefilter", funcs = getFilterNames())
getRdAsText(lib)
parseDesc(text)
parseArgs(text)
showESet(eSet)
setESetArgs(filter)
isESet(eSet)
```

**Arguments**

<code>eSet</code>	<code>eSet</code> an ExpressionSet object
<code>lib</code>	<code>lib</code> a character string for the name of an R library where functions of interests reside
<code>funcs</code>	<code>funcs</code> a vector of character strings for names of functions of interest
<code>text</code>	<code>text</code> a character of string from a filed (e. g. description, argument, ..) filed of an Rd file for a fuction
<code>filter</code>	<code>filter</code> a character string for the name of a filter function

**Details**

A set of filters may be selected to filter genes in through each of the filters in the order the filters have been selected

**Value**

A logical vector of length equal to the number of rows of 'expr'. The values in that vector indicate whether the corresponding row of 'expr' passed the set of filter functions.

**Author(s)**

Jianhua Zhang

**See Also**

[genefilter](#)

**Examples**

```

if( interactive() ) {
  data(sample.ExpressionSet)
  res <- eSetFilter(sample.ExpressionSet)
}

```

filtered\_p

*Compute and adjust p-values, with filtering***Description**

Given filter and test statistics in the form of unadjusted p-values, or functions able to compute these statistics from the data, filter and then correct the p-values across a range of filtering stringencies.

**Usage**

```

filtered_p(filter, test, theta, data, method = "none")
filtered_R(alpha, filter, test, theta, data, method = "none")

```

**Arguments**

alpha	A cutoff to which p-values, possibly adjusted for multiple testing, will be compared.
filter	A vector of stage-one filter statistics, or a function which is able to compute this vector from data, if data is supplied.
test	A vector of unadjusted p-values, or a function which is able to compute this vector from the filtered portion of data, if data is supplied. The option to supply a function is useful when the value of the test statistic depends on which hypotheses are filtered out at stage one. (The <b>limma</b> t-statistic is an example.)
theta	A vector with one or more filtering fractions to consider. Actual cutoffs are then computed internally by applying <a href="#">quantile</a> to the filter statistics contained in (or produced by) the filter argument.
data	If filter and/or test are functions rather than vectors of statistics, they will be applied to data. The functions will be passed the whole data object, and must work over rows, etc. themselves as appropriate.
method	The unadjusted p-values contained in (or produced by) test will be adjusted for multiple testing after filtering, using the <a href="#">p.adjust</a> function in the <b>stats</b> package. See the method argument there for options. p

**Value**

For `filtered_p`, a matrix of p-values, possible adjusted for multiple testing, with one row per null hypothesis and one column per filtering fraction given in `theta`. For a given column, entries which have been filtered out are NA.

For `filtered_R`, a count of the entries in the `filtered_p` result which are less than alpha.

**Author(s)**

Richard Bourgon <bourgon@ebi.ac.uk>

**See Also**

See [rejection\\_plot](#) for visualization of filtered\_p results.

**Examples**

```
# See the vignette: Diagnostic plots for independent filtering
```

---

filterfun	<i>Creates a first FALSE exiting function from the list of filter functions it is given.</i>
-----------	--

---

**Description**

This function creates a function that takes a single argument. The filtering functions are bound in the environment of the returned function and are applied sequentially to the argument of the returned function. When the first filter function evaluates to FALSE the function returns FALSE otherwise it returns TRUE.

**Usage**

```
filterfun(...)
```

**Arguments**

... Filtering functions.

**Value**

filterfun returns a function that takes a single argument. It binds the filter functions given to it in the environment of the returned function. These functions are applied sequentially (in the order they were given to filterfun). The function returns FALSE (and exits) when the first filter function returns FALSE otherwise it returns TRUE.

**Author(s)**

R. Gentleman

**See Also**

[genefilter](#)

**Examples**

```
set.seed(333)
x <- matrix(rnorm(100,2,1),nc=10)
cvfun <- cv(.5,2.5)
ffun <- filterfun(cvfun)
which <- genefilter(x, ffun)
```



---

filter_volcano	<i>Volcano plot for overall variance filtering</i>
----------------	--

---

### Description

Generate a volcano plot contrasting p-value with fold change (on the log scale), in order to visualize the effect of filtering on overall variance and also assign significance via p-value.

### Usage

```
filter_volcano(
  d, p, S,
  n1, n2,
  alpha, S_cutoff,
  cex = 0.5, pch = 19,
  xlab = expression(paste(log[2], " fold change")),
  ylab = expression(paste("-", log[10], " p")),
  cols = c("grey80", "grey50", "black"),
  ltys = c(1, 3),
  use_legend = TRUE,
  ...
)
```

### Arguments

d	Fold changes, typically on the log scale, base 2.
p	The p-values
S	The overall standard deviation filter statistics, i.e., the square roots of the overall variance filter statistics.
n1	Sample size for group 1.
n2	Sample size for group 2.
alpha	Significance cutoff used for p-values.
S_cutoff	Filter cutoff used for the overall standard deviation in S.
cex	Point size for plotting.
pch	Point character for plotting.
xlab	Label for x-axis.
ylab	Label for y-axis.
cols	A vector of three colors used for plotting. These correspond to filtered data, data which pass the filter but are insignificant, and data pass the filter and are also statistically significant.
ltys	The induced bound on log-scale fold change is plotted, as is the significance cutoff for data passing the filter. The ltys argument gives line styles for these drawing these two thresholds on the plot.
use_legend	Should a legend for point color be produced?
...	Other arguments for plot.

**Author(s)**

Richard Bourgon <bourgon@ebi.ac.uk>

**Examples**

```
# See the vignette: Diagnostic plots for independent filtering
```

---

findLargest

*Find the Entrez Gene ID corresponding to the largest statistic*

---

**Description**

Most microarrays have multiple probes per gene (Entrez). This function finds all replicates, and then selects the one with the largest value of the test statistic.

**Usage**

```
findLargest(gN, testStat, data = "hgu133plus2")
```

**Arguments**

gN	A vector of probe identifiers for the chip.
testStat	A vector of test statistics, of the same length as gN with the per probe test statistics.
data	The character string identifying the chip.

**Details**

All the probe identifiers, gN, are mapped to Entrez Gene IDs and the duplicates determined. For any set of probes that map to the same Gene ID, the one with the largest test statistic is found. The return vector is the named vector of selected probe identifiers. The names are the Entrez Gene IDs.

This could be extended in different ways, such as allowing the user to use a different selection criterion. Also, matching on different identifiers seems like another alternative.

**Value**

A named vector of probe IDs. The names are Entrez Gene IDs.

**Author(s)**

R. Gentleman

**See Also**

[sapply](#)

**Examples**

```
library("hgu95av2.db")
set.seed(124)
gN <- sample(1s(hgu95av2ENTREZID), 200)
stats <- rnorm(200)
findLargest(gN, stats, "hgu95av2")
```

---

**gapFilter***A filter to select genes based on there being a gap.*

---

**Description**

The gapFilter looks for genes that might usefully discriminate between two groups (possibly unknown at the time of filtering). To do this we look for a gap in the ordered expression values. The gap must come in the central portion (we exclude jumps in the initial Prop values or the final Prop values). Alternatively, if the IQR for the gene is large that will also pass our test and the gene will be selected.

**Usage**

```
gapFilter(Gap, IQR, Prop, na.rm=TRUE, neg.rm=TRUE)
```

**Arguments**

Gap	The size of the gap required to pass the test.
IQR	The size of the IQR required to pass the test.
Prop	The proportion (or number) of samples to exclude at either end.
na.rm	If TRUE then NA's will be removed before processing.
neg.rm	If TRUE then negative values in x will be removed before processing.

**Details**

As stated above we are interested in

**Value**

A function that returns either TRUE or FALSE depending on whether the vector supplied has a gap larger than Gap or an IQR (inter quartile range) larger than IQR. For computing the gap we want to exclude a proportion, Prop from either end of the sorted values. The reason for this requirement is that genes which differ in expression levels only for a few samples are not likely to be interesting.

**Author(s)**

R. Gentleman

**See Also**

[ttest](#), [genefilter](#)

## Examples

```
set.seed(256)
x <- c(rnorm(10,100,3), rnorm(10, 100, 10))
y <- x + c(rep(0,10), rep(100,10))
tmp <- rbind(x,y)
Gfilter <- gapFilter(200, 100, 5)
ffun <- filterfun(Gfilter)
genefilter(tmp, ffun)
```

---

genefilter

*A function to filter genes.*

---

## Description

genefilter filters genes in the array `expr` using the filter functions in `flist`. It returns an array of logical values (suitable for subscripting) of the same length as there are rows in `expr`. For each row of `expr` the returned value is TRUE if the row passed all the filter functions. Otherwise it is set to FALSE.

## Usage

```
genefilter(expr, flist)
```

## Arguments

`expr` A matrix or ExpressionSet that the filter functions will be applied to.  
`flist` A list of filter functions to apply to the array.

## Details

This package uses a very simple but powerful protocol for *filtering* genes. The user simply constructs any number of tests that they want to apply. A test is simply a function (as constructed using one of the many helper functions in this package) that returns TRUE if the gene of interest passes the test (or filter) and FALSE if the gene of interest fails.

The benefit of this approach is that each test is constructed individually (and can be tested individually). The tests are then applied sequentially to each gene. The function returns a logical vector indicating whether the gene passed all tests functions or failed at least one of them.

Users can construct their own filters. These filters should accept a vector of values, corresponding to a row of the `expr` object. The user defined function should return a length 1 logical vector, with value TRUE or FALSE. User-defined functions can be combined with `filterfun`, just as built-in filters.

## Value

A logical vector of length equal to the number of rows of `expr`. The values in that vector indicate whether the corresponding row of `expr` passed the set of filter functions.

## Author(s)

R. Gentleman

**See Also**

[genefilter](#), [kOverA](#)

**Examples**

```
set.seed(-1)
f1 <- kOverA(5, 10)
f1ist <- filterfun(f1)
exprA <- matrix(rnorm(10000, 10), ncol = 10)
ans <- genefilter(exprA, f1ist)
```

---

genefinder

*Finds genes that have similar patterns of expression.*

---

**Description**

Given an ExpressionSet or a matrix of gene expressions, and the indices of the genes of interest, `genefinder` returns a list of the `numResults` closest genes. The user can specify one of the standard distance measures listed below. The number of values to return can be specified. The return value is a list with two components: `genes` (measured through the desired distance method) to the genes of interest (where `X` is the number of desired results returned) and their distances.

**Usage**

```
genefinder(X, ilist, numResults=25, scale="none", weights, method="euclidean")
```

**Arguments**

<code>X</code>	A numeric matrix where columns represent patients and rows represent genes.
<code>ilist</code>	A vector of genes of interest. Contains indices of genes in matrix <code>X</code> .
<code>numResults</code>	Number of results to display, starting from the least distance to the greatest.
<code>scale</code>	One of "none", "range", or "zscore". Scaling is carried out separately on each row.
<code>weights</code>	A vector of weights applied across the columns of <code>X</code> . If no weights are supplied, no weights are applied.
<code>method</code>	One of "euclidean", "maximum", "manhattan", "canberra", "correlation", "binary".

**Details**

If the `scale` option is "range", then the input matrix is scaled using `genescale()`. If it is "zscore", then the input matrix is scaled using the `scale` builtin with no arguments.

The `method` option specifies the metric used for gene comparisons. The metric is applied, row by row, for each gene specified in `ilist`.

The "correlation" option for the distance method will return a value equal to `1-correlation(x)`.

See [dist](#) for a more detailed description of the distances.

**Value**

The returned value is a list containing an entry for each gene specified in `ilist`. Each list entry contains an array of distances for that gene of interest.

**Author(s)**

J. Gentry and M. Kajen

**See Also**

[genescale](#)

**Examples**

```
set.seed(12345)

#create some fake expression profiles
m1 <- matrix (1:12, 4, 3)
v1 <- 1
nr <- 2

#find the 2 rows of m1 that are closest to row 1
genefinder (m1, v1, nr, method="euc")

v2 <- c(1,3)
genefinder (m1, v2, nr)

genefinder (m1, v2, nr, scale="range")

genefinder (m1, v2, nr, method="manhattan")

m2 <- matrix (rnorm(100), 10, 10)
v3 <- c(2, 5, 6, 8)
nr2 <- 6
genefinder (m2, v3, nr2, scale="zscore")
```

---

genescale

*Scales a matrix or vector.*

---

**Description**

`genescale` returns a scaled version of the input matrix `m` by applying the following formula to each column of the matrix:

$$y[i] = (x[i] - \min(x)) / (\max(x) - \min(x))$$

**Usage**

```
genescale(m, axis=2, method=c("Z", "R"), na.rm=TRUE)
```

**Arguments**

<code>m</code>	Input a matrix or a vector with numeric elements.
<code>axis</code>	An integer indicating which axis of <code>m</code> to scale.
<code>method</code>	Either "Z" or "R", indicating whether a Z scaling or a range scaling should be performed.
<code>na.rm</code>	A boolean indicating whether NA's should be removed.

**Details**

Either the rows or columns of `m` are scaled. This is done either by subtracting the mean and dividing by the standard deviation ("Z") or by subtracting the minimum and dividing by the range.

**Value**

A scaled version of the input. If `m` is a matrix or a dataframe then the dimensions of the returned value agree with that of `m`, in both cases the returned value is a matrix.

**Author(s)**

R. Gentleman

**See Also**

[genefinder](#), [scale](#)

**Examples**

```
m <- matrix(1:12, 4, 3)
genescale(m)
```

---

half.range.mode

*Mode estimation for continuous data*

---

**Description**

For data assumed to be drawn from a unimodal, continuous distribution, the mode is estimated by the "half-range" method. Bootstrap resampling for variance reduction may optionally be used.

**Usage**

```
half.range.mode(data, B, B.sample, beta = 0.5, diag = FALSE)
```

**Arguments**

<code>data</code>	A numeric vector of data from which to estimate the mode.
<code>B</code>	Optionally, the number of bootstrap resampling rounds to use. Note that <code>B = 1</code> resamples 1 time, whereas omitting <code>B</code> uses <code>data</code> as is, without resampling.
<code>B.sample</code>	If bootstrap resampling is requested, the size of the bootstrap samples drawn from <code>data</code> . Default is to use a sample which is the same size as <code>data</code> . For large data sets, this may be slow and unnecessary.
<code>beta</code>	The fraction of the remaining range to use at each iteration.
<code>diag</code>	Print extensive diagnostics. For internal testing only... best left FALSE.

**Details**

Briefly, the mode estimator is computed by iteratively identifying densest half ranges. (Other fractions of the current range can be requested by setting beta to something other than 0.5.) A densest half range is an interval whose width equals half the current range, and which contains the maximal number of observations. The subset of observations falling in the selected densest half range is then used to compute a new range, and the procedure is iterated. See the references for details.

If bootstrapping is requested, B half-range mode estimates are computed for B bootstrap samples, and their average is returned as the final estimate.

**Value**

The mode estimate.

**Author(s)**

Richard Bourgon <bourgon@stat.berkeley.edu>

**References**

- DR Bickel, “Robust estimators of the mode and skewness of continuous data.” *Computational Statistics & Data Analysis* 39:153-163 (2002).
- SB Hedges and P Shah, “Comparison of mode estimation methods and application in molecular clock analysis.” *BMC Bioinformatics* 4:31-41 (2003).

**See Also**

[shorth](#)

**Examples**

```
## A single normal-mixture data set

x <- c( rnorm(10000), rnorm(2000, mean = 3) )
M <- half.range.mode( x )
M.bs <- half.range.mode( x, B = 100 )

if(interactive()){
  hist( x, breaks = 40 )
  abline( v = c( M, M.bs ), col = "red", lty = 1:2 )
  legend(
    1.5, par("usr")[4],
    c( "Half-range mode", "With bootstrapping (B = 100)" ),
    lwd = 1, lty = 1:2, cex = .8, col = "red"
  )
}

# Sampling distribution, with and without bootstrapping

X <- rbind(
  matrix( rnorm(1000 * 100), ncol = 100 ),
  matrix( rnorm(200 * 100, mean = 3), ncol = 100 )
)
M.list <- list(
  Simple = apply( X, 2, half.range.mode ),
```



```
BS = apply( X, 2, half.range.mode, B = 100 )
)

if(interactive()){
  boxplot( M.list, main = "Effect of bootstrapping" )
  abline( h = 0, col = "red" )
}
```

---

kappa\_p

*Compute proportionality constant for fold change bound.*

---

## Description

Filtering on overall variance induces a lower bound on fold change. This bound depends on the significance of the evidence against the null hypothesis, and is a multiple of the cutoff used for an overall variance filter. It also depends on sample size in both of the groups being compared. These functions compute the multiplier for the supplied p-values or t-statistics.

## Usage

```
kappa_p(p, n1, n2 = n1)
kappa_t(t, n1, n2 = n1)
```

## Arguments

p	The p-values at which to compute the multiplier.
t	The t-statistics at which to compute the multiplier.
n1	Sample size for class 1.
n2	Sample size for class 2.

## Value

A vector of multipliers: one per p-value or t-static in p or t.

## Author(s)

Richard Bourgon <bourgon@ebi.ac.uk>

## Examples

```
# See the vignette: Diagnostic plots for independent filtering
```

---

kOverA *A filter function for k elements larger than A.*

---

**Description**

kOverA returns a filter function with bindings for k and A. This function evaluates to TRUE if at least k of the arguments elements are larger than A.

**Usage**

```
kOverA(k, A=100, na.rm=TRUE)
```

**Arguments**

A	The value you want to exceed.
k	The number of elements that have to exceed A.
na.rm	If set to TRUE any NA's will be removed.

**Value**

A function with bindings for A and k.

**Author(s)**

R. Gentleman

**See Also**

[pOverA](#)

**Examples**

```
fg <- kOverA(5, 100)
fg(90:100)
fg(98:110)
```

---

maxA *A filter function to filter according to the maximum.*

---

**Description**

maxA returns a function with the parameter A bound. The returned function evaluates to TRUE if any element of its argument is larger than A.

**Usage**

```
maxA(A=75, na.rm=TRUE)
```

**Arguments**

A                    The value that at least one element must exceed.  
na.rm                If TRUE then NA's are removed.

**Value**

maxA returns a function with an environment containing a binding for A.

**Author(s)**

R. Gentleman

**See Also**

[pOverA](#)

**Examples**

```
ff <- maxA(30)
ff(1:10)
ff(28:31)
```

---

nsFilter

*Filtering of Features in an ExpressionSet*

---

**Description**

The function nsFilter tries to provide a one-stop shop for different options of filtering (removing) features from an ExpressionSet. Filtering features exhibiting little variation, or a consistently low signal, across samples can be advantageous for the subsequent data analysis (Bourgon et al.). Furthermore, one may decide that there is little value in considering features with insufficient annotation.

**Usage**

```
nsFilter(eset, require.entrez=TRUE,
         require.GOBP=FALSE, require.GOCC=FALSE,
         require.GOMF=FALSE, require.CytoBand=FALSE,
         remove.dupEntrez=TRUE, var.func=IQR,
         var.cutoff=0.5, var.filter=TRUE,
         filterByQuantile=TRUE, feature.exclude="^AFFX", ...)

varFilter(eset, var.func=IQR, var.cutoff=0.5, filterByQuantile=TRUE)

featureFilter(eset, require.entrez=TRUE,
              require.GOBP=FALSE, require.GOCC=FALSE,
              require.GOMF=FALSE, require.CytoBand=FALSE,
              remove.dupEntrez=TRUE, feature.exclude="^AFFX")
```

**Arguments**

<code>eset</code>	an ExpressionSet object
<code>var.func</code>	The function used as the per-feature filtering statistic. This function should return a numeric vector of length one when given a numeric vector as input.
<code>var.filter</code>	A logical indicating whether to perform filtering based on <code>var.func</code> .
<code>filterByQuantile</code>	A logical indicating whether <code>var.cutoff</code> is to be interpreted as a quantile of all <code>var.func</code> values (the default), or as an absolute value.
<code>var.cutoff</code>	A numeric value. If <code>var.filter</code> is TRUE, features whose value of <code>var.func</code> is less than either: the <code>var.cutoff</code> -quantile of all <code>var.func</code> values (if <code>filterByQuantile</code> is TRUE), or <code>var.cutoff</code> (if <code>filterByQuantile</code> is FALSE) will be removed.
<code>require.entrez</code>	If TRUE, filter out features without an Entrez Gene ID annotation. If using an annotation package where an identifier system other than Entrez Gene IDs is used as the central ID, then that ID will be required instead.
<code>require.GOBP</code> , <code>require.GOCC</code> , <code>require.GOMF</code>	If TRUE, filter out features whose target genes are not annotated to at least one GO term in the BP, CC or MF ontology, respectively.
<code>require.CytoBand</code>	If TRUE, filter out features whose target genes have no mapping to cytoBand locations.
<code>remove.dupEntrez</code>	If TRUE and there are features mapping to the same Entrez Gene ID (or equivalent), then the feature with the largest value of <code>var.func</code> will be retained and the other(s) removed.
<code>feature.exclude</code>	A character vector of regular expressions. Feature identifiers (i.e. value of <code>featureNames(eset)</code> ) that match one of the specified patterns will be filtered out. The default value is intended to filter out Affymetrix quality control probe sets.
<code>...</code>	Unused, but available for specializing methods.

**Details**

In this Section, the effect of filtering on the type I error rate estimation / control of subsequent hypothesis testing is explained. See also the paper by Bourgon et al.

*Marginal type I errors:* Filtering on the basis of a statistic which is independent of the test statistic used for detecting differential gene expression can increase the detection rate at the same marginal type I error. This is clearly the case for filter criteria that do not depend on the data, such as the annotation based criteria provided by the `nsFilter` and `featureFilter` functions. However, marginal type I error can also be controlled for certain types of data-dependent criteria. Call  $U^I$  the stage 1 filter statistic, which is a function that is applied feature by feature, based on whose value the feature is or is not accepted to pass to stage 2, and which depends only on the data for that feature and not any other feature, and call  $U^{II}$  the stage 2 test statistic for differential expression. Sufficient conditions for marginal type-I error control are:

- $U^I$  the overall (across all samples) variance or mean,  $U^{II}$  the t-statistic (or any other scale and location invariant statistic), data normal distributed and exchangeable across samples.
- $U^I$  the overall mean,  $U^{II}$  the moderated t-statistic (as in limma's `eBayes` function), data normal distributed and exchangeable.

- $U^I$  a sample-class label independent function (e.g. overall mean, median, variance, IQR),  $U^{II}$  the Wilcoxon rank sum statistic, data exchangeable.

*Experiment-wide type I error:* Marginal type-I error control provided by the conditions above is sufficient for control of the family wise error rate (FWER). Note, however, that common false discovery rate (FDR) methods depend not only on the marginal behaviour of the test statistics under the null hypothesis, but also on their joint distribution. The joint distribution can be affected by filtering, even when this filtering leaves the marginal distributions of true-null test statistics unchanged. Filtering might, for example, change correlation structure. The effect of this is negligible in many cases in practice, but this depends on the dataset and the filter used, and the assessment is in the responsibility of the data analyst.

*Annotation Based Filtering* Arguments `require.entrez`, `require.GOBP`, `require.GOCC`, `require.GOMF` and `require.CytoBand` filter based on available annotation data. The annotation package is determined by calling `annotation(eset)`.

*Variance Based Filtering* The `var.filter`, `var.func`, `var.cutoff` and `varByQuantile` arguments control numerical cutoff-based filtering. Probes for which `var.func` returns NA are removed. The default `var.func` is IQR, which we here define as `rowQ(eset, ceiling(0.75 * ncol(eset))) - rowQ(eset, floor(0.25 * ncol(eset)))`. This choice is motivated by the observation that unexpressed genes are detected most reliably through low variability of their features across samples. Additionally, IQR is robust to outliers (see note below). The default `var.cutoff` is 0.5 and is motivated by a rule of thumb that in many tissues only 40% of genes are expressed. Please adapt this value to your data and question.

By default the numerical-filter cutoff is interpreted as a quantile, so with the default settings, 50% of the genes are filtered.

Variance filtering is performed last, so that (if `varByQuantile=TRUE` and `remove.dupEntrez=TRUE`) the final number of genes does indeed exclude precisely the `var.cutoff` fraction of unique genes remaining after all other filters were passed.

The stand-alone function `varFilter` does only `var.func`-based filtering (and no annotation based filtering). `featureFilter` does only annotation based filtering and duplicate removal; it always performs duplicate removal to retain the highest-IQR probe for each gene.

## Value

For `nsFilter` a list consisting of:

<code>eset</code>	the filtered <code>ExpressionSet</code>
<code>filter.log</code>	a list giving details of how many probe sets were removed for each filtering step performed.

For both `varFilter` and `featureFilter` the filtered `ExpressionSet`.

## Note

IQR is a reasonable variance-filter choice when the dataset is split into two roughly equal and relatively homogeneous phenotype groups. If your dataset has important groups smaller than 25% of the overall sample size, or if you are interested in unusual individual-level patterns, then IQR may not be sensitive enough for your needs. In such cases, you should consider using less robust and more sensitive measures of variance (the simplest of which would be `sd`).

## Author(s)

Seth Falcon (somewhat revised by Assaf Oron)

## References

R. Bourgon, R. Gentleman, W. Huber, Independent filtering increases power for detecting differentially expressed genes, Technical Report.

## Examples

```
library("hgu95av2.db")
library("Biobase")
data(sample.ExpressionSet)
ans <- nsFilter(sample.ExpressionSet)
ans$eset
ans$filter.log

## skip variance-based filtering
ans <- nsFilter(sample.ExpressionSet, var.filter=FALSE)

a1 <- varFilter(sample.ExpressionSet)
a2 <- featureFilter(sample.ExpressionSet)
```

---

pOverA	<i>A filter function to filter according to the proportion of elements larger than A.</i>
--------	---

---

## Description

A function that returns a function with values for  $A$ ,  $p$  and `na.rm` bound to the specified values. The function takes a single vector,  $x$ , as an argument. When the returned function is evaluated it returns TRUE if the proportion of values in  $x$  that are larger than  $A$  is at least  $p$ .

## Usage

```
pOverA(p=0.05, A=100, na.rm=TRUE)
```

## Arguments

$A$	The value to be exceeded.
$p$	The proportion that need to exceed $A$ for TRUE to be returned.
<code>na.rm</code>	If TRUE then NA's are removed.

## Value

`pOverA` returns a function with bindings for  $A$ ,  $p$  and `na.rm`. This function evaluates to TRUE if the proportion of values in  $x$  that are larger than  $A$  exceeds  $p$ .

## Author(s)

R. Gentleman

## See Also

[cv](#)

**Examples**

```
ff<- pOverA(p=.1, 10)
ff(1:20)
ff(1:5)
```

---

rejection_plot	<i>Plot rejections vs. p-value cutoff</i>
----------------	---

---

**Description**

Plot the number, or fraction, of null hypotheses rejected as a function of the p-value cutoff. Multiple sets of p-values are accepted, in a list or in the columns of a matrix, in order to permit comparisons.

**Usage**

```
rejection_plot(p,
               col, lty = 1, lwd = 1,
               xlab = "p cutoff", ylab = "number of rejections",
               xlim = c(0, 1), ylim,
               legend = names(p),
               at = c("all", "sample"),
               n_at = 100,
               probability = FALSE,
               ...
               )
```

**Arguments**

p	The p-values to be used for plotting. These may be in the columns of a matrix, or in the elements of a list. One curve will be generated for each column/element, and all NA entries will be dropped. If column or element names are supplied, they are used by default for a plot legend.
col	Colors to be used for each curve plotted. Recycled if necessary. If col is omitted, <a href="#">rainbow</a> is used to generate a set of colors.
lty	Line styles to be used for each curve plotted. Recycled if necessary.
lwd	Line widths to be used for each curve plotted. Recycled if necessary.
xlab	X-axis text label.
ylab	Y-axis text label.
xlim	X-axis limits.
ylim	Y-axis limits.
legend	Text for legend. Matrix column names or list element names (see p above) are used by default. If NULL, no legend is plotted.
at	Should step functions be plotted with a step at every value in p, or should linear interpolation be used at a sample of points spanning xlim? The latter looks when there are many p-values.
n_at	When at = "sample" is given, how many sample points should be used for interpolation and plotting?

probability Should the fraction of null hypotheses rejected be reported instead of the count? See the probability argument to [hist](#).

... Other arguments to pass to the [plot](#) call which sets up the axes. Note that the ... argument will not be passed to the [lines](#) calls which actually generate the curves.

### Value

A list of the step functions used for plotting is returned invisibly.

### Author(s)

Richard Bourgon <bourgon@ebi.ac.uk>

### Examples

```
# See the vignette: Diagnostic plots for independent filtering
```

---

rowFtests *t-tests and F-tests for rows or columns of a matrix*

---

### Description

t-tests and F-tests for rows or columns of a matrix, intended to be speed efficient.

### Usage

```
rowtttests(x, fac, tstatOnly = FALSE, na.rm = FALSE)
coltttests(x, fac, tstatOnly = FALSE, na.rm = FALSE)
fastT(x, ig1, ig2, var.equal = TRUE)

rowFtests(x, fac, var.equal = TRUE)
colFtests(x, fac, var.equal = TRUE)
```

### Arguments

**x** Numeric matrix. The matrix must not contain NA values. For `rowtttests` and `coltttests`, `x` can also be an [ExpressionSet](#).

**fac** Factor which codes the grouping to be tested. There must be 1 or 2 groups for the t-tests (corresponding to one- and two-sample t-test), and 2 or more for the F-tests. If `fac` is missing, this is taken as a one-group test (i.e. is only allowed for the t-tests). The length of the factor needs to correspond to the sample size: for the `row*` functions, the length of the factor must be the same as the number of columns of `x`, for the `col*` functions, it must be the same as the number of rows of `x`.  
If `x` is an [ExpressionSet](#), then `fac` may also be a character vector of length 1 with the name of a covariate in `x`.

**tstatOnly** A logical variable indicating whether to calculate p-values from the t-distribution with appropriate degrees of freedom. If TRUE, just the t-statistics are returned. This can be considerably faster.



<code>na.rm</code>	A logical variable indicating whether to remove NA values prior to calculation test statistics.
<code>ig1</code>	The indices of the columns of <code>x</code> that correspond to group 1.
<code>ig2</code>	The indices of the columns of <code>x</code> that correspond to group 2.
<code>var.equal</code>	A logical variable indicating whether to treat the variances in the samples as equal. If 'TRUE', a simple F test for the equality of means in a one-way analysis of variance is performed. If 'FALSE', an approximate method of Welch (1951) is used, which generalizes the commonly known 2-sample Welch test to the case of arbitrarily many samples.

### Details

If `fac` is specified, `rowttests` performs for each row of `x` a two-sided, two-class t-test with equal variances. `fac` must be a factor of length `ncol(x)` with two levels, corresponding to the two groups. The sign of the resulting t-statistic corresponds to "group 1 minus group 2". If `fac` is missing, `rowttests` performs for each row of `x` a two-sided one-class t-test against the null hypothesis 'mean=0'.

`rowttests` and `colttests` are implemented in C and should be reasonably fast and memory-efficient. `fastT` is an alternative implementation, in Fortran, possibly useful for certain legacy code. `rowFtests` and `colFtests` are currently implemented using matrix algebra in R. Compared to the `rowttests` and `colttests` functions, they are slower and use more memory.

### Value

A data.frame with columns `statistic`, `p.value` (optional in the case of the t-test functions) and `dm`, the difference of the group means (only in the case of the t-test functions). The row.names of the data.frame are taken from the corresponding dimension names of `x`.

The degrees of freedom are provided in the attribute `df`. For the F-tests, if `var.equal` is 'FALSE', `nrow(x)+1` degree of freedoms are given, the first one is the first degree of freedom (it is the same for each row) and the other ones are the second degree of freedom (one for each row).

### Author(s)

Wolfgang Huber <whuber@embl.de>

### References

B. L. Welch (1951), On the comparison of several mean values: an alternative approach. *Biometrika*, \*38\*, 330-336

### See Also

[mt.teststat](#)

### Examples

```
##
## example data
##
x = matrix(runif(40), nrow=4, ncol=10)
f2 = factor(floor(runif(ncol(x))*2))
f4 = factor(floor(runif(ncol(x))*4))
```

```

##
## one- and two group row t-test; 4-group F-test
##
r1 = rowttests(x)
r2 = rowttests(x, f2)
r4 = rowFtests(x, f4)

## approximate equality
about.equal = function(x,y,tol=1e-10)
  stopifnot(is.numeric(x), is.numeric(y), length(x)==length(y), all(abs(x-y) < tol))

##
## compare with the implementation in t.test
##
for (j in 1:nrow(x)) {
  s1 = t.test(x[j,])
  about.equal(s1$statistic, r1$statistic[j])
  about.equal(s1$p.value, r1$p.value[j])

  s2 = t.test(x[j,] ~ f2, var.equal=TRUE)
  about.equal(s2$statistic, r2$statistic[j])
  about.equal(s2$p.value, r2$p.value[j])

  dm = -diff(tapply(x[j,], f2, mean))
  about.equal(dm, r2$dm[j])

  s4 = summary(lm(x[j,] ~ f4))
  about.equal(s4$fstatistic["value"], r4$statistic[j])
}

##
## colttests
##
c2 = colttests(t(x), f2)
stopifnot(identical(r2, c2))

##
## missing values
##
f2n = f2
f2n[sample(length(f2n), 3)] = NA
r2n = rowttests(x, f2n)
for(j in 1:nrow(x)) {
  s2n = t.test(x[j,] ~ f2n, var.equal=TRUE)
  about.equal(s2n$statistic, r2n$statistic[j])
  about.equal(s2n$p.value, r2n$p.value[j])
}

##
## larger sample size
##
x = matrix(runif(1000000), nrow=4, ncol=250000)
f2 = factor(floor(runif(ncol(x))*2))
r2 = rowttests(x, f2)
for (j in 1:nrow(x)) {
  s2 = t.test(x[j,] ~ f2, var.equal=TRUE)
  about.equal(s2$statistic, r2$statistic[j])
}

```

```

    about.equal(s2$p.value, r2$p.value[j])
  }

  ## single row matrix
  rowFtests(matrix(runif(10),1,10),as.factor(c(rep(1,5),rep(2,5))))
  rowTtests(matrix(runif(10),1,10),as.factor(c(rep(1,5),rep(2,5))))

```

rowpAUCs-methods

*Rowwise ROC and pAUC computation***Description**

Methods for fast rowwise computation of ROC curves and (partial) area under the curve (pAUC) using the simple classification rule  $x > \text{theta}$ , where  $\text{theta}$  is a value in the range of  $x$

**Usage**

```
rowpAUCs(x, fac, p=0.1, flip=TRUE, caseNames=c("1", "2"))
```

**Arguments**

<code>x</code>	ExpressionSet or numeric matrix. The matrix must not contain NA values.
<code>fac</code>	A factor or numeric or character that can be coerced to a factor. If <code>x</code> is an ExpressionSet, this may also be a character vector of length 1 with the name of a covariate variable in <code>x</code> . <code>fac</code> must have exactly 2 levels. For better control over the classification, use integer values in 0 and 1, where 1 indicates the "Disease" class in the sense of the Pepe et al paper (see below).
<code>p</code>	Numeric vector of length 1. Limit in (0,1) to integrate pAUC to.
<code>flip</code>	Logical. If TRUE, both classification rules $x > \text{theta}$ and $x < \text{theta}$ are tested and the (partial) area under the curve of the better one of the two is returned. This is appropriate for the cases in which the classification is not necessarily linked to higher expression values, but instead it is symmetric and one would assume both over- and under-expressed genes for both classes. You can set <code>flip</code> to FALSE if you only want to screen for genes which discriminate Disease from Control with the $x > \text{theta}$ rule.
<code>caseNames</code>	The class names that are used when plotting the data. If <code>fac</code> is the name of the covariate variable in the ExpressionSet the function will use its levels as <code>caseNames</code> .

**Details**

Rowwise calculation of Receiver Operating Characteristic (ROC) curves and the corresponding partial area under the curve (pAUC) for a given data matrix or ExpressionSet. The function is implemented in C and thus reasonably fast and memory efficient. Cutpoints ( $\text{theta}$ ) are calculated before the first, in between and after the last data value. By default, both classification rules  $x > \text{theta}$  and  $x < \text{theta}$  are tested and the (partial) area under the curve of the better one of the two is returned. This is only valid for symmetric cases, where the classification is independent of the magnitude of  $x$  (e.g., both over- and under-expression of different genes in the same class). For unsymmetric cases in which you expect  $x$  to be consistently higher/lower in one of the two classes (e.g. presence or absence of a single biomarker) set `flip=FALSE` or use the functionality provided in the ROC package. For better control over the classification (i.e., the choice of "Disease" and "Control" class in the sense of the Pepe et al paper), argument `fac` can be an integer in  $[0, 1]$  where 1 indicates "Disease" and 0 indicates "Control".

**Value**

An object of class `rowROC` with the calculated specificities and sensitivities for each row and the corresponding pAUCs and AUCs values. See `rowROC` for details.

**Methods**

Methods exist for rowPAUCs:

```
signature(x="matrix", fac="factor")
```

```
rowPAUCs signature(x="matrix", fac="numeric")
```

```
rowPAUCs signature(x="ExpressionSet")
```

```
rowPAUCs signature(x="ExpressionSet", fac="character")
```

**Author(s)**

Florian Hahne <fhahne@fhcrc.org>

**References**

Pepe MS, Longton G, Anderson GL, Schummer M.: Selecting differentially expressed genes from microarray experiments. *Biometrics*. 2003 Mar;59(1):133-42.

**See Also**

[rocdemo.sca](#), [pAUC](#), [rowROC](#)

**Examples**

```
library(Biobase)
data(sample.ExpressionSet)

r1 = rowttests(sample.ExpressionSet, "sex")
r2 = rowpAUCs(sample.ExpressionSet, "sex", p=0.1)

plot(area(r2, total=TRUE), r1$statistic, pch=16)
sel <- which(area(r2, total=TRUE) > 0.7)
plot(r2[sel])

## this compares performance and output of rowpAUCs to function pAUC in
## package ROC
if(require(ROC)){
  ## performance
  myRule = function(x)
    pAUC(rocdemo.sca(truth = as.integer(sample.ExpressionSet$sex)-1 ,
      data = x, rule = dxrule.sca), t0 = 0.1)
  nGenes = 200
  cat("computation time for ", nGenes, "genes:\n")
  cat("function pAUC: ")
  print(system.time(r3 <- esApply(sample.ExpressionSet[1:nGenes, ], 1, myRule)))
  cat("function rowpAUCs: ")
  print(system.time(r2 <- rowpAUCs(sample.ExpressionSet[1:nGenes, ],
    "sex", p=1)))

  ## compare output
  myRule2 = function(x)
```

```

    pAUC(rocdemo.sca(truth = as.integer(sample.ExpressionSet$sex)-1 ,
                    data = x, rule = dxrule.sca), t0 = 1)
  r4 <- esApply(sample.ExpressionSet[1:nGenes, ], 1, myRule2)
  plot(r4, area(r2), xlab="function pAUC", ylab="function rowpAUCs",
       main="pAUCs")

  plot(r4, area(rowpAUCs(sample.ExpressionSet[1:nGenes, ],
                        "sex", p=1, flip=FALSE)), xlab="function pAUC", ylab="function rowpAUCs",
       main="pAUCs")

  r4[r4<0.5] <- 1-r4[r4<0.5]
  plot(r4, area(r2), xlab="function pAUC", ylab="function rowpAUCs",
       main="pAUCs")
}

```

rowROC-class

*Class "rowROC"***Description**

A class to model ROC curves and corresponding area under the curve as produced by rowpAUCs.

**Objects from the Class**

Objects can be created by calls of the form `new("rowROC", ...)`.

**Slots**

**data:** Object of class "matrix" The input data.

**ranks:** Object of class "matrix" The ranked input data.

**sens:** Object of class "matrix" Matrix of sensitivity values for each gene at each cutpoint.

**spec:** Object of class "matrix" Matrix of specificity values for each gene at each cutpoint.

**pAUC:** Object of class "numeric" The partial area under the curve (integrated from 0 to p).

**AUC:** Object of class "numeric" The total area under the curve.

**factor:** Object of class "factor" The factor used for classification.

**cutpoints:** Object of class "matrix" The values of the cutpoints at which specificity and sensitivity was calculated. (Note: the data is ranked prior to computation of ROC curves, the cutpoints map to the ranked data.)

**caseNames:** Object of class "character" The names of the two classification cases.

**p:** Object of class "numeric" The limit to which pAUC is integrated.

**Methods**

**show** signature(object="rowROC") Print nice info about the object.

[ signature(x="rowROC", j="missing") Subset the object according to rows/genes.

**plot** signature(x="rowROC", y="missing") Plot the ROC curve of the first row of the object along with the pAUC. To plot the curve for a specific row/gene subsetting should be done first (i.e. `plot(rowROC[1])`).

**pAUC** signature(object="rowROC", p="numeric", flip="logical") Integrate area under the curve from 0 to p. This method returns a new rowROC object.

**AUC** signature(object="rowROC") Integrate total area under the curve. This method returns a new rowROC object.

**sens** signature(object="rowROC") Accessor method for sensitivity slot.

**spec** signature(object="rowROC") Accessor method for specificity slot.

**area** signature(object="rowROC", total="logical") Accessor method for pAUC slot.

### Author(s)

Florian Hahne <fhahne@fhcrc.org>

### References

Pepe MS, Longton G, Anderson GL, Schummer M.: Selecting differentially expressed genes from microarray experiments. *Biometrics*. 2003 Mar;59(1):133-42.

### See Also

[rowpAUCs](#)

### Examples

```
library(Biobase)
require(genefilter)
data(sample.ExpressionSet)
roc <- rowpAUCs(sample.ExpressionSet, "sex", p=0.5)
roc
area(roc[1:3])

if(interactive()) {
  par(ask=TRUE)
  plot(roc)
  plot(1-spec(roc[1]), sens(roc[2]))
  par(ask=FALSE)
}

pAUC(roc, 0.1)
roc
```

---

rowSds

*Row variance and standard deviation of a numeric array*

---

### Description

Row variance and standard deviation of a numeric array

### Usage

```
rowVars(x, ...)
rowSds(x, ...)
```

**Arguments**

- `x` An array of two or more dimensions, containing numeric, complex, integer or logical values, or a numeric data frame.
- `...` Further arguments that get passed on to `rowMeans` and `rowSums`.

**Details**

These are very simple convenience functions, the main work is done in `rowMeans` and `rowSums`. See the function definition of `rowVars`, it is very simple.

**Value**

A numeric or complex array of suitable size, or a vector if the result is one-dimensional. The 'dimnames' (or 'names' for a vector result) are taken from the original array.

**Author(s)**

Wolfgang Huber <http://www.ebi.ac.uk/huber>

**See Also**

`rowMeans` and `rowSums`

**Examples**

```
a = matrix(rnorm(1e4), nrow=10)
rowSds(a)
```

---

shorth

*A location estimator based on the shorth*

---

**Description**

A location estimator based on the shorth

**Usage**

```
shorth(x, na.rm=FALSE, tie.action="mean", tie.limit=0.05)
```

**Arguments**

- `x` Numeric
- `na.rm` Logical. If TRUE, then non-finite (according to `is.finite`) values in `x` are ignored. Otherwise, presence of non-finite or NA values will lead to an error message.
- `tie.action` Character scalar. See details.
- `tie.limit` Numeric scalar. See details.

## Details

The shorth is the shortest interval that covers half of the values in  $x$ . This function calculates the mean of the  $x$  values that lie in the shorth. This was proposed by Andrews (1972) as a robust estimator of location.

Ties: if there are multiple shortest intervals, the action specified in `ties.action` is applied. Allowed values are `mean` (the default), `max` and `min`. For `mean`, the average value is considered; however, an error is generated if the start indices of the different shortest intervals differ by more than the fraction `tie.limit` of `length(x)`. For `min` and `max`, the left-most or right-most, respectively, of the multiple shortest intervals is considered.

Rate of convergence: as an estimator of location of a unimodal distribution, under regularity conditions, the quantity computed here has an asymptotic rate of only  $n^{-1/3}$  and a complicated limiting distribution.

See `half.range.mode` for an iterative version that refines the estimate iteratively and has a builtin bootstrapping option.

## Value

The mean of the  $x$  values that lie in the shorth.

## Author(s)

Wolfgang Huber <http://www.ebi.ac.uk/huber>, Ligia Pedroso Bras

## References

- G Sawitzki, "The Shorth Plot." Available at <http://lshorth.r-forge.r-project.org/TheShorthPlot.pdf>
- DF Andrews, "Robust Estimates of Location." Princeton University Press (1972).
- R Grueble, "The Length of the Shorth." *Annals of Statistics* 16, 2:619-628 (1988).
- DR Bickel and R Fruehwirth, "On a fast, robust estimator of the mode: Comparisons to other robust estimators with applications." *Computational Statistics & Data Analysis* 50, 3500-3530 (2006).

## See Also

[half.range.mode](#)

## Examples

```
x = c(rnorm(500), runif(500) * 10)
methods = c("mean", "median", "shorth", "half.range.mode")
ests = sapply(methods, function(m) get(m)(x))

if(interactive()) {
  colors = 1:4
  hist(x, 40, col="orange")
  abline(v=ests, col=colors, lwd=3, lty=1:2)
  legend(5, 100, names(ests), col=colors, lwd=3, lty=1:2)
}
```



---

tdata	<i>A small test dataset of Affymetrix Expression data.</i>
-------	--

---

**Description**

The tdata data frame has 500 rows and 26 columns. The columns correspond to samples while the rows correspond to genes. The row names are Affymetrix accession numbers.

**Usage**

```
data(tdata)
```

**Format**

This data frame contains 26 columns.

**Source**

An unknown data set.

**Examples**

```
data(tdata)
```

---

ttest	<i>A filter function for a t.test</i>
-------	---------------------------------------

---

**Description**

ttest returns a function of one argument with bindings for cov and p. The function, when evaluated, performs a t-test using cov as the covariate. It returns TRUE if the p value for a difference in means is less than p.

**Usage**

```
ttest(m, p=0.05, na.rm=TRUE)
```

**Arguments**

m	If m is of length one then it is assumed that elements one through m of x will be one group. Otherwise m is presumed to be the same length as x and constitutes the groups.
p	The p-value for the test.
na.rm	If set to TRUE any NA's will be removed.

**Details**

When the data can be split into two groups (diseased and normal for example) then we often want to select genes on their ability to distinguish those two groups. The t-test is well suited to this and can be used as a filter function.

This helper function creates a t-test (function) for the specified covariate and considers a gene to have passed the filter if the p-value for the gene is less than the prespecified p.

**Value**

ttest returns a function with bindings for m and p that will perform a t-test.

**Author(s)**

R. Gentleman

**See Also**

[kOverA](#), [Anova](#), [t.test](#)

**Examples**

```
dat <- c(rep(1,5),rep(2,5))
set.seed(5)
y <- rnorm(10)
af <- ttest(dat, .01)
af(y)
af2 <- ttest(5, .01)
af2(y)
y[8] <- NA
af(y)
af2(y)
y[1:5] <- y[1:5]+10
af(y)
```

# Index

- \*Topic **arith**
  - shorth, [31](#)
- \*Topic **array**
  - rowSds, [30](#)
- \*Topic **classes**
  - rowROC-class, [29](#)
- \*Topic **datasets**
  - tdata, [33](#)
- \*Topic **manip**
  - Anova, [2](#)
  - coxfilter, [3](#)
  - cv, [4](#)
  - dist2, [5](#)
  - eSetFilter, [6](#)
  - filterfun, [8](#)
  - findLargest, [10](#)
  - gapFilter, [11](#)
  - genefilter, [12](#)
  - genefinder, [13](#)
  - genescale, [14](#)
  - kOverA, [18](#)
  - maxA, [18](#)
  - nsFilter, [19](#)
  - pOverA, [22](#)
  - rowSds, [30](#)
  - ttest, [33](#)
- \*Topic **math**
  - rowFtests, [24](#)
  - rowPAUCs-methods, [27](#)
- \*Topic **robust**
  - half.range.mode, [15](#)
- \*Topic **univar**
  - half.range.mode, [15](#)
- [, rowROC, ANY, ANY, ANY-method  
(rowROC-class), [29](#)
- Anova, [2](#), [3](#), [34](#)
- area (rowROC-class), [29](#)
- area, rowROC-method (rowROC-class), [29](#)
- AUC (rowROC-class), [29](#)
- AUC, rowROC-method (rowROC-class), [29](#)
- colFtests (rowFtests), [24](#)
- colFtests, ExpressionSet, character-method  
(rowFtests), [24](#)
- colFtests, ExpressionSet, factor-method  
(rowFtests), [24](#)
- colFtests, matrix, factor-method  
(rowFtests), [24](#)
- colttests (rowFtests), [24](#)
- colttests, ExpressionSet, character-method  
(rowFtests), [24](#)
- colttests, ExpressionSet, factor-method  
(rowFtests), [24](#)
- colttests, ExpressionSet, missing-method  
(rowFtests), [24](#)
- colttests, matrix, factor-method  
(rowFtests), [24](#)
- colttests, matrix, missing-method  
(rowFtests), [24](#)
- coxfilter, [3](#)
- coxph, [3](#)
- cv, [4](#), [22](#)
- dist, [5](#), [13](#)
- dist2, [5](#)
- eBayes, [20](#)
- eSetFilter, [6](#)
- ExpressionSet, [24](#)
- fastT (rowFtests), [24](#)
- featureFilter (nsFilter), [19](#)
- filter\_volcano, [9](#)
- filtered\_p, [7](#)
- filtered\_R (filtered\_p), [7](#)
- filterfun, [8](#), [12](#)
- findLargest, [10](#)
- gapFilter, [11](#)
- genefilter, [6](#), [8](#), [11](#), [12](#), [13](#)
- genefinder, [13](#), [15](#)
- genefinder, ExpressionSet, vector-method  
(genefinder), [13](#)
- genefinder, matrix, vector-method  
(genefinder), [13](#)
- genescale, [14](#), [14](#)

- getFilterNames (eSetFilter), 6
- getFuncDesc (eSetFilter), 6
- getRdAsText (eSetFilter), 6
- half.range.mode, 15, 32
- hist, 24
- is.finite, 31
- isESet (eSetFilter), 6
- kappa\_p, 17
- kappa\_t (kappa\_p), 17
- kOverA, 3, 4, 13, 18, 34
- lines, 24
- lm, 3
- maxA, 18
- mt.teststat, 25
- nsFilter, 19
- nsFilter, ExpressionSet-method  
(nsFilter), 19
- p.adjust, 7
- parseArgs (eSetFilter), 6
- parseDesc (eSetFilter), 6
- pAUC, 28
- pAUC (rowROC-class), 29
- pAUC, rowROC, numeric-method  
(rowROC-class), 29
- plot, 24
- plot, rowROC, missing-method  
(rowROC-class), 29
- pOverA, 4, 18, 19, 22
- quantile, 7
- rainbow, 23
- rejection\_plot, 8, 23
- rocdemo.sca, 28
- rowFtests, 24
- rowFtests, ExpressionSet, character-method  
(rowFtests), 24
- rowFtests, ExpressionSet, factor-method  
(rowFtests), 24
- rowFtests, matrix, factor-method  
(rowFtests), 24
- rowMeans, 31
- rowpAUCs, 30
- rowpAUCs (rowpAUCs-methods), 27
- rowpAUCs, ExpressionSet, ANY-method  
(rowpAUCs-methods), 27
- rowpAUCs, ExpressionSet, character-method  
(rowpAUCs-methods), 27
- rowpAUCs, matrix, factor-method  
(rowpAUCs-methods), 27
- rowpAUCs, matrix, numeric-method  
(rowpAUCs-methods), 27
- rowpAUCs-methods, 27
- rowROC, 28
- rowROC (rowROC-class), 29
- rowROC-class, 29
- rowSds, 30
- rowSums, 31
- rowttests (rowFtests), 24
- rowttests, ExpressionSet, character-method  
(rowFtests), 24
- rowttests, ExpressionSet, factor-method  
(rowFtests), 24
- rowttests, ExpressionSet, missing-method  
(rowFtests), 24
- rowttests, matrix, factor-method  
(rowFtests), 24
- rowttests, matrix, missing-method  
(rowFtests), 24
- rowVars (rowSds), 30
- sapply, 10
- scale, 15
- sens (rowROC-class), 29
- sens, rowROC-method (rowROC-class), 29
- setESetArgs (eSetFilter), 6
- shorth, 16, 31
- show, rowROC-method (rowROC-class), 29
- showESet (eSetFilter), 6
- spec (rowROC-class), 29
- spec, rowROC-method (rowROC-class), 29
- t.test, 34
- tdata, 33
- ttest, 11, 33
- varFilter (nsFilter), 19