

Package ‘univOutl’

January 9, 2018

Type Package

Title Detection of Univariate Outliers

Version 0.1-4

Date 2017-12-27

Author Marcello D'Orazio

Maintainer Marcello D'Orazio <mdo.statmatch@gmail.com>

Depends robustbase, Hmisc

Description Well known outlier detection techniques in the univariate case. Methods to deal with skewed distribution are included too. The Hidioglou-Berthelot (1986) method to search for outliers in ratios of historical data is implemented as well. When available, survey weights can be used in outliers detection.

License GPL (>= 2)

URL <https://github.com/marcellodo/univOutl>

NeedsCompilation no

Repository CRAN

Date/Publication 2018-01-09 10:45:01 UTC

R topics documented:

univOutl-package	2
boxB	2
HBmethod	5
LocScaleB	8
ratioSize	11

Index	14
--------------	-----------

univOutl-package

Detection of Univariate Outliers.

Description

Well known outlier detection techniques in the univariate case. Ratios of two variables are covered too. When available, survey weights can be considered.

Details

The package provides few simple functions implementing well known outlier detection techniques in the univariate case. Methods to deal with skewed distribution are included. The Hidiroglou-Berthelot (1986) method to search for outliers in ratios of historical data is implemented as well. When available, survey weights can be used in outliers detection.

Author(s)

Author and Maintainer: Marcello D'Orazio <mdo.statmatch@gmail.com>

References

Hidiroglou, M.A. and Berthelot, J.-M. (1986) 'Statistical editing and Imputation for Periodic Business Surveys'. *Survey Methodology*, Vol 12, pp. 73-83.

McGill, R., Tukey, J. W. and Larsen, W. A. (1978) 'Variations of box plots'. *The American Statistician*, 32, pp. 12-16.

Rousseeuw, P.J. and Croux, C. (1993) 'Alternatives to the Median Absolute Deviation', *Journal of the American Statistical Association* 88, pp. 1273-1283.

Vanderviere, E. and Huber, M. (2008) 'An Adjusted Boxplot for Skewed Distributions', *Computational Statistics & Data Analysis*, 52, pp. 5186-5201

boxB

BoxPlot based outlier detection

Description

Identifies univariate outliers by using methods based on BoxPlots

Usage

```
boxB(x, k=1.5, method='asymmetric', weights=NULL, id=NULL,
     exclude=NA, logt=FALSE)
```

Arguments

x	Numeric vector that will be searched for outliers.
k	Nonnegative constant that determines the extension of the 'whiskers'. Commonly used values are 1.5 (default), 2, or 3. Note that when method="adjbox" then k is set automatically equal to 1.5
method	Character, identifies the method to be used: method="resistant" provides the 'standard' boxplot fences; method="asymmetric" is a modification of standard method to deal with (moderately) skewed data; method="adjbox" uses Hubert and Vandervieren (2008) adjusted boxplot for skewed distributions.
weights	Optional numeric vector with units' weights associated to observations of x. Only nonnegative weights are allowed. Weights are used in estimating the quartiles (see Details).
id	Optional vector with identifiers of units in x. If missing (id=NULL, default) the identifiers will be set equal to the positions in the vector (i.e. id=1:length(x)).
exclude	Values of x that will be excluded by the analysis. By default missing values are excluded (exclude=NA).
logt	Logical, if TRUE, before searching outliers the x variable is log-transformed (log(x+1) is considered). In this case the summary outputs (bounds, etc.) will refer to the log-transformed x

Details

When method="resistant" the outlying observations are those outside the interval:

$$[Q_1 - k \times IQR; \quad Q_3 + k \times IQR]$$

where Q_1 and Q_3 are respectively the 1st and the 3rd quartile of x, while $IQR = (Q_3 - Q_1)$ is the Inter-Quartile Range. The value $k = 1.5$ (said 'inner fences') is commonly used when drawing a boxplot. Values $k = 2$ and $k = 3$ provide middle and outer fences, respectively.

When method="asymmetric" the outlying observations are those outside the interval:

$$[Q_1 - 2k \times (Q_2 - Q_1); \quad Q_3 + 2k \times (Q_3 - Q_2)]$$

being Q_2 the median; such a modification allows to account for slight skewness of the distribution.

Finally, when method="adjbox" the outlying observations are identified using the method proposed by Hubert and Vandervieren (2008) and based on the Medcouple measure of skewness; in practice the bounds are:

$$[Q_1 - 1.5 \times e^{aM} \times IQR; \quad Q_3 + 1.5 \times e^{bM} \times IQR]$$

Where M is the medcouple; when $M > 0$ (positive skewness) then $a = -4$ and $b = 3$; on the contrary $a = -3$ and $b = 4$ for negative skewness ($M < 0$). This adjustment of the boxplot, according to Hubert and Vandervieren (2008), works with moderate skewness ($-0.6 \leq M \leq 0.6$). The bounds of the adjusted boxplot are derived by applying the function `adjboxStats` in the package **robustbase**.

When weights are available (passed via the argument `weights`) then they are used in the computation of the quartiles. In particular, the quartiles are derived using the function `wtd.quantile` in the package **Hmisc**.

Remember that when asking a log transformation (argument `logt=TRUE`) all the estimates (quartiles, etc.) will refer to $\log(x + 1)$.

Value

The output is a list containing the following components:

quartiles	The quartiles of x after discarding the values in the <code>exclude</code> argument. When weights are provided they are used in quartiles estimation through the function <code>wtd.quantile</code> in the package Hmisc .
fences	The bounds of the interval, values outside the interval are detected as outliers.
excluded	The identifiers or positions (when <code>id=NULL</code>) of units in x excluded by the computations, according to the argument <code>exclude</code> .
outliers	The identifiers or positions (when <code>id=NULL</code>) of units in x detected as outliers.

Author(s)

Marcello D'Orazio <mdo.statmatch@gmail.com>

References

McGill, R., Tukey, J. W. and Larsen, W. A. (1978) 'Variations of box plots'. *The American Statistician*, 32, pp. 12-16.

Vanderviere, E. and Huber, M. (2008) 'An Adjusted Boxplot for Skewed Distributions', *Computational Statistics and Data Analysis*, 52, pp. 5186-5201.

See Also

[adjboxStats](#), [wtd.quantile](#)

Examples

```
set.seed(321)
x <- rnorm(30, 50, 10)
x[10] <- 1
x[20] <- 100

out <- boxB(x = x, k = 1.5, method = 'asymmetric')
out$fences
out$outliers
x[out$outliers]

out <- boxB(x = x, k = 1.5, method = 'adjbox')
out$fences
out$outliers
```

```

x[out$outliers]

x[24] <- NA
x.ids <- paste0('obs',1:30)
out <- boxB(x = x, k = 1.5, method = 'adjbox', id = x.ids)
out$excluded
out$fences
out$outliers

set.seed(111)
w <- round(runif(n = 30, min=1, max=10))
out <- boxB(x = x, k = 1.5, method = 'adjbox', id = x.ids, weights = w)
out$excluded
out$fences
out$outliers

```

HBmethod

Hiroglou-Berthelot procedure for detecting outliers with periodic data

Description

This function implements the method proposed by Hidioglou and Berthelot (1986) to identify outliers in periodic data, i.e. when the same variable is measured at two time points.

Usage

```

HBmethod(yt1, yt2, U=0.5, A=0.05, C=4,
         id=NULL, std.score=FALSE, return.dataframe=FALSE)

```

Arguments

yt1	Numeric vector providing the values observed at time t1.
yt2	Numeric vector providing the values observed at time t2 ($t_2 > t_1$).
U	Numeric, parameter needed to determine the ‘importance’ of a ratio. The value should lie in $[0, 1]$ interval; commonly used values are 0.3, 0.4, or 0.5 (default) (see Details for further information).
A	Numeric, parameter needed when computing the scale measure used to derive the bounds. Hidioglou and Berthelot (1986) suggest setting $A = 0.05$ (default) (see Details for further information).
C	Numeric, parameter determining the extension of the interval (greater values will provide larger intervals, i.e. fewer expected outliers). Values commonly used are 4 (default) or 7, but also values close or greater than 40 can be used in some particular cases. Note that two C values can be provided instead of one, the first one will be used to determine the left tail bound, while the second determines the right tail bound; this setting can help in improving outlier detection in skewed distributions (see Details for further information).

id	Optional numeric or character vector, with identifiers of units. If id=NULL units identifiers will be set equal to their position.
std.score	Logical, if TRUE the output will include a standardized score variable (see Details, for further information)
return.dataframe	Logical, if TRUE the output will save all the relevant information for outlier detection in a dataframe with the following columns: 'id' (units' identifiers), 'yt1', 'yt2', 'ratio' (= yt1/yt2), 'Escore' (the E scores, see Details), 'std.Escore' (the standardized E scores when std.score=TRUE, see Details) and finally 'outliers', where value 1 indicates observations detected as an outlier, 0 otherwise.

Details

The method proposed by Hidioglou and Berthelot (1986) to identify outliers in periodic data consists in deriving a score variable based on the ratios $r_i = y_{i,t2}/y_{i,t1}$ ($yt2/yt1$) with $i = 1, 2, \dots, n$ being n the number of observations after discarding NAs and 0s in both yt1 and yt2.

At first the ratios are centered around their median r_M :

$$s_i = 1 - r_M/r_i \quad \text{if} \quad 0 < r_i < r_M$$

$$s_i = r_i/r_M - 1 \quad \text{if} \quad r_i \geq r_M$$

Then, in order to account for the magnitude of data, the following score is derived:

$$E_i = s_i \times [\max(y_{i,t1}, y_{i,t2})]^U$$

Finally, the interval is calculated as:

$$(E_M - C \times d_{Q1}, E_M + C \times d_{Q3})$$

where

$$d_{Q1} = \max(E_M - E_{Q1}, |A \times E_M|) \text{ and } d_{Q3} = \max(E_{Q3} - E_M, |A \times E_M|)$$

being d_{Q1} , E_M and d_{Q3} the quartiles of E scores.

In practice, all the units with an E score outside the interval are considered as outliers. Notice that when two C values are provided, then the first is used to derive the left bound while the second determines the right bound.

When std.score=TRUE a standardized score is derived in the following manner:

$$z_{E,i} = 0.6745 \times \frac{E_i - E_M}{d_{Q1}} \quad \text{if} \quad E_i < E_M$$

$$z_{E,i} = 0.6745 \times \frac{E_i - E_M}{d_{Q3}} \quad \text{if} \quad E_i \geq E_M$$

The constant 0.6745 makes d_{Q1} and d_{Q3} approximately unbiased estimators when the E scores follow the normal distribution.

Value

A list whose components depend on the `return.dataframe` argument. When `return.dataframe=FALSE` just the following components are provided:

<code>median.r</code>	the median of the ratios
<code>quartiles.E</code>	Quartiles of the E score
<code>bounds.E</code>	Bounds of the interval of the E score, values outside are considered outliers.
<code>excluded</code>	The identifiers or positions (when <code>id=NULL</code>) of units in both <code>yt1</code> and <code>yt2</code> that are excluded by the outliers detection, i.e. NAs and 0s.
<code>outliers</code>	The identifiers or positions (when <code>id=NULL</code>) of units in <code>yt1</code> or <code>yt2</code> identified as outliers.

When `return.dataframe=TRUE`, the first three components remain the same with, in addition, two dataframes:

<code>excluded</code>	A dataframe with the subset of observations excluded. The data frame has the following columns: 'id' (units' identifiers), 'yt1' columns 'yt2'
<code>data</code>	A dataframe with the the not excluded observations and the following columns: 'id' (units' identifiers), 'yt1', 'yt2', 'ratio' (= $yt1/yt2$), 'Escore' (the E scores, see Details), 'std.Escore' (the standardized E scores when <code>std.score=TRUE</code> , see Details) and finally 'outliers', where value 1 indicates observations detected as an outlier, 0 otherwise.

Author(s)

Marcello D'Orazio <mdo.statmatch@gmail.com>

References

Hidiroglou, M.A. and Berthelot, J.-M. (1986) 'Statistical editing and Imputation for Periodic Business Surveys'. *Survey Methodology*, Vol 12, pp. 73-83.

Examples

```
set.seed(222)
x0 <- rnorm(30, 50, 5)
x0[1] <- NA
set.seed(333)
rr <- runif(30, 0.9, 1.2)
rr[10] <- 2
x1 <- x0 * rr
x1[20] <- 0

out <- HBmethod(yt1 = x0, yt2 = x1)
out$excluded
out$median.r
out$bounds.E
out$outliers
```

```

cbind(x0[out$outliers], x1[out$outliers])

out <- HBmethod(yt1 = x0, yt2 = x1,
               return.dataframe = TRUE)
out$excluded
head(out$data)

```

LocScaleB

Univariate outlier detection with bounds based on robust location and scale estimates

Description

This function identifies outliers in the tails of a distribution by detecting the observations outside the bounds built using a robust estimate of both location and scale parameters.

Usage

```

LocScaleB(x, k=3, method='MAD', weights=NULL, id=NULL,
          exclude=NA, logt=FALSE, return.dataframe=FALSE)

```

Arguments

x	Numeric vector that will be searched for outliers.
k	Nonnegative constant that determines the extension of bounds. Commonly used values are 2, 2.5 and 3 (default).
method	<p>character identifying how to estimate the scale of the distribution. Available choices are:</p> <p>method='IQR' for using the Inter-Quartile Range, i.e. Q75-Q25;</p> <p>method='IDR' for using the Inter-Decile Range; i.e. Q90-Q10</p> <p>method='MAD' for using the Median Absolute Deviation;</p> <p>method='Gini' robust scale estimate based on Gini's Mean Difference (see GiniMd);</p> <p>method='ScaleTau2' robust tau-estimate of univariate scale, as proposed by Maronna and Zamar (2002) (see also scaleTau2);</p> <p>method='Qn' for using the Qn estimator proposed by Rousseeuw and Croux (1993) (see also Qn);</p> <p>method='Sn' for using the Sn estimator proposed by Rousseeuw and Croux (1993) (see also Sn).</p> <p>When method='dQ' the estimated scale for the left tail is $(Q2-Q1)/0.6745$, while for the right tail it is considered $(Q3-Q2)/0.6745$; this double estimate should be able to account for slight skewness.</p> <p>Finally, when method='AdjOut', bounds are based on the adjusted outlyingness method as proposed by Hubert and Van der Veen (2008).</p>

weights	Optional numeric vector that provides weights associated to observations. Only nonnegative weights are allowed. Note that weights can only be used when <code>method='MAD'</code> , <code>method='IQR'</code> or <code>method='dQ'</code> .
id	Optional numeric or character vector, with identifiers of units in <code>x</code> . If <code>id=NULL</code> (default) units' identifiers will be set equal to their position in <code>x</code> .
exclude	Values of <code>x</code> that will be excluded by the analysis. By default missing values (<code>exclude = NA</code>)
logt	Logical, if TRUE, before searching outliers the <code>x</code> variable is log-transformed (<code>log(x+1)</code> is considered). Note that in this case that summary output (bounds, etc.) will refer to log-transformed variable.
return.dataframe	Logical, if TRUE the output will save all the relevant information for outlier detection in a dataframe with the following columns: 'id' (units' identifiers), 'x', 'log.x' (only if <code>logt=TRUE</code>), 'weight' (only when argument <code>weights</code> is provided), 'score' (the standardized scores, see Details) and, finally, 'outliers', where value 1 indicates observations detected as an outlier, 0 otherwise.

Details

The intervals are derived by considering the median as a robust location estimate while different robust scale estimators are considered:

$$[Q_2 - k \times \tilde{s}_L; \quad Q_2 + k \times \tilde{s}_R]$$

Q_2 is the median, while \tilde{s}_L and \tilde{s}_R are robust scale estimates. With most of the methods $\tilde{s}_L = \tilde{s}_R$ with exception of `method='dQ'` where:

$$\tilde{s}_L = (Q_2 - Q_1)/0.6745 \quad \text{and} \quad \tilde{s}_R = (Q_3 - Q_2)/0.6745$$

With `method='AdjOut'` the following estimates are considered:

$$\tilde{s}_L = (Q_2 - f_L) \quad \text{and} \quad \tilde{s}_R = (f_R - Q_2)$$

being f_R and f_L derived starting from the fences of the adjusted boxplot (Hubert and Vandervieren, 2008; see [adjboxStats](#)).

When weights are available (passed via the argument `weights`) then they are used in the computation of the quartiles. In particular, the quartiles are derived using the function `wtd.quantile` in the package **Hmisc**. Note that their use is allowed just with `method='IQR'`, `method='dQ'` or `method='AdjOut'`.

The 'score' variable reported in the the data dataframe when `return.dataframe=TRUE` is the standardized score derived as $(x - \text{Median})/\text{scale}$.

Value

A list whose components depend on the `return.dataframe` argument. When `return.dataframe = FALSE` just the following components are provided:

<code>pars</code>	Vector with estimated median and scale parameters
<code>bounds</code>	The bounds of the interval, values outside the interval are considered outliers.
<code>excluded</code>	The position or identifiers of <code>x</code> values excluded by outlier detection, according to the argument <code>exclude</code>
<code>outliers</code>	The position or identifiers of <code>x</code> values detected as outliers (outside bounds).

When `return.dataframe=TRUE` the latter two components are substituted with two dataframes:

<code>excluded</code>	A dataframe with the subset of observations excluded.
<code>data</code>	A dataframe with the the not excluded observations and the following columns: 'id' (units' identifiers), 'x', 'log.x' (only if <code>logt=TRUE</code>), 'weight' (only when argument <code>weights</code> is provided), 'score' (the standardized scores, see Details) and, finally, 'outliers', where value 1 indicates observations detected as an outlier and 0 otherwise.

Author(s)

Marcello D'Orazio <mdo.statmatch@gmail.com>

References

Hubert, M. and Van der Veecken, S. (2008) 'Outlier Detection for Skewed Data'. *Journal of Chemometrics*, 22, pp. 235-246.

Maronna, R.A. and Zamar, R.H. (2002) 'Robust estimates of location and dispersion of high-dimensional datasets' *Technometrics*, 44, pp. 307-317.

Rousseeuw, P.J. and Croux, C. (1993) 'Alternatives to the Median Absolute Deviation', *Journal of the American Statistical Association* 88, pp. 1273-1283.

Vanderviere, E. and Huber, M. (2008) 'An Adjusted Boxplot for Skewed Distributions', *Computational Statistics & Data Analysis*, 52, pp. 5186-5201

See Also

[mad](#), [scaleTau2](#), [Qn](#), [Sn](#), [GiniMd](#)

Examples

```
set.seed(333)
x <- rnorm(30, 50, 1)
x[10] <- 1
x[20] <- 100

out <- LocScaleB(x = x, k = 3, method='MAD')
out$pars
```

```

out$bounds
out$outliers
x[out$outliers]

out <- LocScaleB(x = x, k = 3, method='MAD',
                 return.dataframe = TRUE)
head(out$data)

out <- LocScaleB(x = x, k = 3, method='AdjOut')
out$outliers

```

ratioSize	<i>Identifies outliers on ratios and filter them by a size measure</i>
-----------	--

Description

Identifies outliers on transformed ratios (centering with respect to their median) using the adjusted boxplot for skewed distributions. Outliers can be sorted/filtered according to a size measure.

Usage

```

ratioSize(numerator, denominator, id=NULL,
          size=NULL, U=1, size.th=NULL, return.dataframe=FALSE)

```

Arguments

numerator	Numeric vector with the values that go at numerator of the ratio
denominator	Numeric vector with the values that go at denominator of the ratio
id	Optional numeric or character vector, with identifiers of units. If id=NULL units identifiers will be set equal to their positions in x.
size	Optional numeric vector providing a measure of the importance of a ratio. If size = NULL the size measure is the maximum value between the numerator and the denominator of each ratio (makes sense if both the variables are observed using the same unit of measure). Observations' importance is also controlled by the argument U.
U	Numeric, constant with $0 < U \leq 1$ controlling importance of each unit, in practice the final size measure is derived as $(size^U)$. Commonly used values are 0.4, 0.5 or 1.
size.th	Numeric, size threshold. Can be specified when a size measure is used. In such a case just outliers with a size greater than the threshold will be returned. Note that when argument U is not set equal to 1, then the final threshold will be $size.th^U$.

return.dataframe

Logical, if TRUE the output will save all the relevant information for outlier detection in a dataframe with the following columns: 'id' (units' identifiers), 'numerator', 'denominator', 'ratio' (= numerator/denominator), 'c.ratio' (centered ratios, see Details), 'size' (size^U values) and finally 'outliers', where value 1 indicates observations detected as an outlier and 0 otherwise.

Details

This function searches for outliers starting from ratios $r = \text{numerator}/\text{denominator}$. At first the ratios are centered around their median, as in Hidioglou Berthelot (1986) procedure (see [HBmethod](#)), then the outlier identification is based on the adjusted boxplot for skewed distribution (Hubert and Vandervieren 2008) (see [adjboxStats](#)). The subset of outliers is sorted in decreasing order according to the size measure. If a size threshold is provided then just outliers with $(\text{size}^U) > (\text{size.th}^U)$ will be returned.

Value

A list whose components depend on the return.dataframe argument. When return.dataframe = FALSE just the following components are returned:

median.r	the median of the ratios
bounds	The bounds of the interval for centered ratios
excluded	The position or the identifiers of the units with values excluded by the computations because of 0s or NAs.
outliers	The position or the identifiers of the units detected as outliers. Remember that when size.th is set, just outliers with $(\text{size}^U) > (\text{size.th}^U)$ will be returned.

When return.dataframe=TRUE the latter two components are substituted with two dataframes:

excluded	A dataframe with the subset of observations excluded
data	A dataframe with the not excluded observations with the following columns: 'id' (units' identifiers), 'numerator', 'denominator', 'ratio' (= numerator/denominator), 'c.ratio' (centered ratios, see Details), 'size' (size ^U values) and finally 'outliers', where value 1 indicates observations detected as an outlier and 0 otherwise. The data frame will be sorted in decreasing manner according to size ^U . Note that when a size threshold is provided then ONLY outliers with $(\text{size}^U) > (\text{size.th}^U)$ will be returned.

Author(s)

Marcello D'Orazio <mdo.statmatch@gmail.com>

References

- Hidioglou, M.A. and Berthelot, J.-M. (1986) 'Statistical editing and Imputation for Periodic Business Surveys'. *Survey Methodology*, Vol 12, pp. 73-83.
- Vanderviere, E. and Huber, M. (2008) 'An Adjusted Boxplot for Skewed Distributions', *Computational Statistics and Data Analysis*, 52, pp. 5186-5201.

See Also[HBmethod](#), [boxB](#), [adjboxStats](#)**Examples**

```
set.seed(444)
x1 <- rnorm(30, 50, 5)
set.seed(555)
rr <- runif(30, 0.9, 1.2)
rr[10] <- 2
x2 <- x1 * rr

out <- ratioSize(numerator = x2, denominator = x1)
out

out <- ratioSize(numerator = x2, denominator = x1,
                 return.dataframe = TRUE)
head(out$data)

out <- ratioSize(numerator = x2, denominator = x1,
                 size.th = 65, return.dataframe = TRUE)
head(out$data)
```

Index

- *Topic **package**
 - univOutl-package, 2
- *Topic **robust**
 - LocScaleB, 8
- *Topic **survey**
 - boxB, 2
 - HBmethod, 5
 - LocScaleB, 8
 - ratioSize, 11
- *Topic **univar**
 - boxB, 2
 - HBmethod, 5
 - LocScaleB, 8
 - ratioSize, 11

adjboxStats, 3, 4, 9, 12, 13

boxB, 2, 13

GiniMd, 8, 10

HBmethod, 5, 12, 13

LocScaleB, 8

mad, 10

Qn, 8, 10

ratioSize, 11

scaleTau2, 8, 10

Sn, 8, 10

univOutl (univOutl-package), 2

univOutl-package, 2

wtd.quantile, 4, 9