

# Package ‘highfrequency’

December 26, 2012

**Version** 0.1

**Date** 2011-04-04

**Title** highfrequency

**Author** Jonathan Cornelissen, Kris Boudt, Scott Payseur

**Maintainer** Jonathan Cornelissen <Jonathan.cornelissen@kuleuven.be>

**Description** The highfrequency package contains an extensive toolkit for the use of highfrequency financial data in R. It contains functionality to manage, clean and match highfrequency trades and quotes data. Furthermore, it enables users to: calculate easily various liquidity measures, estimate and forecast volatility, and investigate microstructure noise and intraday periodicity.

**License** GPL ( $\geq 2$ )

**Depends** R ( $\geq 2.12.0$ ), xts, zoo

**Suggests** realized, robustbase, cubature, mvtnorm, chron, timeDate, quantmod

**LazyLoad** yes

**Archs** i386, x86\_64

## R topics documented:

highfrequency-package . . . . .	3
aggregatePrice . . . . .	3
aggregateQuotes . . . . .	4
aggregateTrades . . . . .	5
aggregatets . . . . .	6
autoSelectExchangeQuotes . . . . .	8
autoSelectExchangeTrades . . . . .	9
convert . . . . .	10
exchangeHoursOnly . . . . .	12
getTradeDirection . . . . .	12
harModel . . . . .	13
heavyModel . . . . .	15
lltc.xts . . . . .	17
makePsd . . . . .	17
makeReturns . . . . .	18

matchTradesQuotes . . . . .	19
medRV . . . . .	20
mergeQuotesSameTimestamp . . . . .	21
mergeTradesSameTimestamp . . . . .	22
minRV . . . . .	22
noZeroPrices . . . . .	23
noZeroQuotes . . . . .	24
previoustick . . . . .	24
quotesCleanup . . . . .	24
rAccumulation . . . . .	26
rAVGCov . . . . .	27
rBPCov . . . . .	29
rCov . . . . .	30
rCumSum . . . . .	31
realized_library . . . . .	32
refreshTime . . . . .	32
rHYCov . . . . .	33
rKernel.available . . . . .	34
rKernelCov . . . . .	35
rMarginal . . . . .	36
rmLargeSpread . . . . .	38
rmNegativeSpread . . . . .	38
rmOutliers . . . . .	39
rmTradeOutliers . . . . .	40
rOWCov . . . . .	41
rRTSCov . . . . .	42
rScatterReturns . . . . .	45
rThresholdCov . . . . .	46
rTSCov . . . . .	47
rZero . . . . .	49
salesCondition . . . . .	50
sample_5minprices . . . . .	51
sample_5minprices_jumps . . . . .	51
sample_qdata . . . . .	52
sample_qdataraw . . . . .	52
sample_real5minprices . . . . .	53
sample_tdata . . . . .	53
sample_tdataraw . . . . .	53
sbux.xts . . . . .	54
selectExchange . . . . .	54
spotVol . . . . .	55
TAQLoad . . . . .	57
tqLiquidity . . . . .	58
tradesCleanup . . . . .	62
tradesCleanupFinal . . . . .	64

---

highfrequency-package    *highfrequency: Toolkit for the analysis of highfrequency financial data in R.*

---

### Description

The highfrequency package contains an extensive toolkit for the use of highfrequency financial data in R. It contains functionality to manage, clean and match highfrequency trades and quotes data. Furthermore, it enables users to: calculate easily various liquidity measures, estimate and forecast volatility, and investigate microstructure noise and intraday periodicity.

### Details

Package: highfrequency  
Type: Package  
Version: 0.1  
Date: 2012-06-01  
License:

### Author(s)

Jonathan Cornelissen, Kris Boudt, Scott Payseur

Maintainer: Jonathan Cornelissen <Jonathan.cornelissen@econ.kuleuven.be>

### Examples

```
# see users manual
```

---

aggregatePrice    *Aggregate a time series but keep first and last observation*

---

### Description

Function returns new time series as xts object where first observation is always the opening price and subsequent observations are the closing prices over the interval with as endpoint the timestamp of the result.

### Usage

```
aggregatePrice(ts,FUN = previoustick,on="minutes",k=1,  
marketopen,marketclose)
```

**Arguments**

ts	xts object to be aggregated, containing the intraday price series of a stock for one day.
FUN	function to apply over each interval. By default, previous tick aggregation is done.
on	character, indicating the time scale in which "k" is expressed. Possible values are: "secs", "seconds", "mins", "minutes", "hours".
k	positive integer, indicating the number of periods to aggregate over. E.g. to aggregate a xts object to the 5 minute frequency set k=5 and on="minutes".
marketopen	the market opening time, by default: marketopen = "09:30:00".
marketclose	the market closing time, by default: marketclose = "16:00:00".

**Value**

An xts object containing the aggregated time series.

**Details**

The timestamps of the new time series are the closing times and/or days of the intervals.

In case of previous tick aggregation, for on="seconds"/"minutes"/"hours", the element of the returned series with e.g. timestamp 09:35:00 contains the last observation up to that point, excluding the value at 09:35:00 itself. An exception is marketclose (i.e. 16:00:00 ET by default), where the observation at 16:00:00 is included in the interval, since this is the end of a trading day at the NYSE.

Please input an object containing ONE day of data.

**Author(s)**

Jonathan Cornelissen and Kris Boudt

**Examples**

```
#load data
data("sample_tdata");
#aggregate price data to the 30 second frequency
head(sample_tdata$PRICE)
head(aggregatePrice(sample_tdata$PRICE,on="secs",k=30));
```

---

aggregateQuotes

*Aggregate an xts object containing quote data*

---

**Description**

Function returns an xts object containing the aggregated quote data with columns "SYMBOL", "EX", "BID", "BIDSIZ", "OFR", "OFRSIZ". See [sample\\_qdata](#) for an example of the argument qdata.

**Usage**

```
aggregateQuotes(qdata,on="minutes",k=5,marketopen,marketclose)
```

**Arguments**

qdata	xts object to be aggregated, containing the intraday quote data of a stock for one day.
on	character, indicating the time scale in which "k" is expressed. Possible values are: "secs", "seconds", "mins", "minutes", "hours".
k	positive integer, indicating the number of periods to aggregate over. E.g. to aggregate a xts object to the 5 minute frequency, set k=5 and on="minutes".
marketopen	the market opening time, by default: marketopen= "09:30:00".
marketclose	the market closing time, by default: marketclose = "16:00:00".

**Value**

An xts object containing the aggregated quote data.

**Details**

The output "BID" and "OFR" columns are constructed using previous tick aggregation.

The variables "BIDSIZ" and "OFRSIZ" are aggregated by taking the sum of the respective inputs over each interval.

For the variables "SYMBOL" and "EX" aggregation doesn't really make sense, thus the first value of the input is taken as the value for the complete output.

Column "MODE" is dropped because aggregating makes absolutely no sense.

The timestamps of the new time series are the closing times of the intervals.

Please Note: Returned objects always contain the first observation (i.e. opening quotes,...).

**Author(s)**

Jonathan Cornelissen and Kris Boudt

**Examples**

```
#load data
data("sample_qdata");
#aggregate quote data to the 30 second frequency
xx = aggregateQuotes(sample_qdata,on="seconds",k=30);
head(xx);
```

---

aggregateTrades

*Aggregate an xts object containing trade data*

---

**Description**

Function returns an xts object containing the aggregated trade data with columns "SYMBOL", "EX", "PRICE", "SIZE". See [sample\\_tdata](#) for an example of the argument tdata.

**Usage**

```
aggregateTrades(tdata,on="minutes",k=5,marketopen,marketclose)
```

**Arguments**

tdata	xts object to be aggregated, containing the intraday trade data of a stock for one day.
on	character, indicating the time scale in which "k" is expressed. Possible values are: "secs", "seconds", "mins", "minutes", "hours".
k	positive integer, indicating the number of periods to aggregate over. E.g. to aggregate a xts object to the 5 minute frequency set k=5 and on="minutes".
marketopen	the market opening time, by default: marketopen= "09:30:00".
marketclose	the market closing time, by default: marketclose = "16:00:00".

**Value**

An xts object containing the aggregated trade data.

**Details**

The output "PRICE" column is constructed using previous tick aggregation.

The variable "SIZE" is aggregated by taking the sum over each interval.

For the variables "SYMBOL" and "EX" aggregation doesn't really make sense, thus the first value of the input is taken as the value for the complete output.

Columns "COND", "CORR", "G127" are dropped because aggregating them makes absolutely no sense.

The timestamps of the new time series are the closing times of the intervals.

Please Note:

Returned objects always contain the first observation (i.e. opening price,...).

Please input an object containing ONE day of data.

**Author(s)**

Jonathan Cornelissen and Kris Boudt

**Examples**

```
data("sample_tdata");
#aggregate trade data to 5 minute frequency
x = aggregateTrades(sample_tdata,on="minutes",k=5)
head(x);
```

---

aggregatets

---

*Aggregate a time series*


---

**Description**

Function returns aggregated time series as xts object. It can handle irregularly spaced timeseries and returns a regularly spaced one. Use univariate timeseries as input for this function, and check out [aggregateTrades](#) and [aggregateQuotes](#) to aggregate Trade or Quote data objects.

**Usage**

```
aggregatets(ts, FUN="previousstick", on="minutes", k=1, weights=NULL,
            dropna=FALSE)
```

**Arguments**

ts	xts object to aggregate.
FUN	function to apply over each interval. By default, previous tick aggregation is done. Alternatively one can set e.g. FUN="mean". In case weights are supplied, this argument is ignored and a weighted average is taken.
on	character, indicating the time scale in which "k" is expressed. Possible values are: "secs", "seconds", "mins", "minutes", "hours", "days", "weeks".
k	positive integer, indicating the number of periods to aggregate over. E.g. to aggregate a xts object to the 5 minute frequency set k=5 and on="minutes".
weights	By default, no weighting scheme is used. When you assign an xts object with weights to this argument, a weighted mean is taken over each interval. Of course, the weights should have the same timestamps as the supplied time series.
dropna	boolean, which determines whether empty intervals should be dropped. By default, an NA is returned in case an interval is empty, except when the user opts for previous tick aggregation, by setting FUN="previousstick" (default).

**Value**

An xts object containing the aggregated time series.

**Details**

The timestamps of the new time series are the closing times and/or days of the intervals. E.g. for a weekly aggregation the new timestamp is the last day in that particular week (namely sunday).

In case of previous tick aggregation, for on="seconds"/"minutes"/"hours", the element of the returned series with e.g. timestamp 09:35:00 contains the last observation up to that point, excluding the value at 09:35:00 itself.

**Author(s)**

Jonathan Cornelissen and Kris Boudt

**Examples**

```
#load sample price data
data("sample_tdata");
ts = sample_tdata$PRICE;

#Previous tick aggregation to the 5-minute sampling frequency:
tsagg5min = aggregatets(ts,on="minutes",k=5);
head(tsagg5min);

#Previous tick aggregation to the 30-second sampling frequency:
tsagg30sec = aggregatets(ts,on="seconds",k=30);
tail(tsagg30sec);
```

---

`autoSelectExchangeQuotes`*Retain only data from the stock exchange with the highest volume*

---

## Description

Function returns an xts object containing only observations of the exchange with highest value for the sum of "BIDSIZ" and "OFRSIZ", i.e. the highest quote volume.

## Usage

```
autoSelectExchangeQuotes(qdata)
```

## Arguments

`qdata` an xts object with at least a column "EX", indicating the exchange symbol and columns "BIDSIZ" and "OFRSIZ", indicating the volume available at the bid and ask respectively. The chosen exchange is printed on the console. The possible exchanges are:

- A: AMEX
- N: NYSE
- B: Boston
- P: Arca
- C: NSX
- T/Q: NASDAQ
- D: NASD ADF and TRF
- X: Philadelphia
- I: ISE
- M: Chicago
- W: CBOE
- Z: BATS

## Value

xts object

## Author(s)

Jonathan Cornelissen and Kris Boudt



---

`autoSelectExchangeTrades`*Retain only data from the stock exchange with the highest trading volume*

---

### Description

Function returns an xts object containing only observations of the exchange with the highest value for the variable "SIZE", i.e. the highest trade volume.

### Usage

```
autoSelectExchangeTrades(tdata)
```

### Arguments

`tdata` an xts object with at least a column "EX", indicating the exchange symbol and "SIZE", indicating the trade volume. The chosen exchange is printed on the console. The possible exchanges are:

- A: AMEX
- N: NYSE
- B: Boston
- P: Arca
- C: NSX
- T/Q: NASDAQ
- D: NASD ADF and TRF
- X: Philadelphia
- I: ISE
- M: Chicago
- W: CBOE
- Z: BATS

### Value

xts object

### Author(s)

Jonathan Cornelissen and Kris Boudt

---

 convert

---

*Convert trade or quote data into xts object saved in the RData format*


---

## Description

Function converts both trade and quote data stored as "txt" or "csv" and structured as illustrated in the pdf documentation into xts-objects and stores them in the "RData" format. Subsequently, the data can be loaded into the R console by [TAQLoad](#).

## Usage

```
convert(from,to,datasource,datadestination,trades=TRUE,quotes=TRUE,
        ticker,dir=FALSE,extension="txt",header=FALSE,
        tradecolnames=NULL,quotecolnames=NULL,
        format="%Y%m%d %H:%M:%S",onefile=FALSE)
```

## Arguments

from	first date to convert e.g. "2008-01-30".
to	last date to convert e.g. "2008-01-31".
datasource	folder in which the original data is stored.
datadestination	folder in which the converted data should be stored.
trades	boolean, determines whether trades are converted.
quotes	boolean, determines whether quotes are converted.
ticker	vector with tickers to be converted.
dir	boolean, if TRUE the datadestination folder and subfolders will be created automatically.
extension	character, indicating the data format of the original data. Can be either "txt" or "csv". In case your data comes from "www.tickdata.com", set extension="tickdatacom".
header	boolean, indicating whether each data file contains a header.
tradecolnames	vector containing column names of your trade data. By default, the standard NYSE data format is taken, see pdf documentation for more details.
quotecolnames	vector containing column names of your quote data. By default, the standard NYSE data format is taken, see pdf documentation for more details.
format	character, indicates in what format TIME and DATE are recorded in the original data. By default, "%Y%m%d %H:%M:%S" is taken, which means the date is denoted by e.g. "20080130" and the time by e.g. "09:30:00".
onefile	indicates the way the source data is stored. By default FALSE, which means the function expects that the source data is saved in a folder structure with a folder for each date (see vignette for more info). In case the data for multiple days is stored in one file, set onefile=TRUE. The naming convention for files is, e.g. for ticker="AAPL", "AAPL_trades.extension" and "AAPL_quotes.extension" for trades and quotes respectively.

**Value**

For each day an xts object is saved into the folder of that date, containing the converted data. This procedure is performed for each stock in "ticker". More information can be found in the pdf documentation.

**Author(s)**

Jonathan Cornelissen and Kris Boudt

**Examples**

```
#In order for these examples to work, the folder datasource
#should contain two folders named 2008-01-02 and 2008-01-03.
#These folder contain the files with the trade data,
#which are named "AAPL_trades.txt" or "AA_trades.txt".

from="2008-01-02";
to = "2008-01-03";
## Not run: datasource=datadestination="C:\data"

### txt files from NYSE:

convert(from,to,datasource,datadestination,trades=TRUE,
        quotes=FALSE,ticker=c("AA","AAPL"),dir=FALSE,extension="txt",
        header=FALSE,tradecolnames=NULL,quotecolnames=NULL,
        format="%Y%m%d %H:%M:%S");

#Now, the folder datadestination will contain two folders
#named 2008-01-02 and 2008-01-03 containing
#the files AAPL_trades.RData and AAPL_trades.RData containing the trades.
#The data can now be loaded with the TAQLoad function.

##### Csv file from WRDS
#Suppose the datasource folder contains one csv file from WRDS
#with data on IBM for multiple days.
#The file should be named "IBM_trades.csv" and can be easily converted into xts
#and then saved in RData format by:

convert(from=from, to=to, datasourcesource=datasource, datadestination=datadestination, trades = T,
        quotes = T, ticker="IBM", dir = FALSE, extension = "csv", header = TRUE,
        tradecolnames = NULL, quotecolnames = NULL, format = format, onefile = TRUE )

##### ASC file from www.tickdata.com
#Suppose the datasource folder contains asc files for trades and quotes
#from "www.tickdata.com" for GLP.
#The files "GLP_quotes.asc" and "GLP_trades.asc" should be saved in datasourcesource folder.

convert(from=from, to=to, datasourcesource=datasourcesource, datadestination=datadestination, trades = T,
        quotes = T, ticker="GLP", dir = TRUE, extension = "tickdatacom", header = TRUE,
        tradecolnames = NULL, quotecolnames = NULL, format = "

## End(Not run)
```

---

exchangeHoursOnly	<i>Extract data from an xts object for the Exchange Hours Only</i>
-------------------	--

---

### Description

The function returns data within exchange trading hours "daybegin" and "dayend". By default, daybegin and dayend are set to "09:30:00" and "16:00:00" respectively (see Brownlees and Gallo (2006) for more information on good choices for these arguments).

### Usage

```
exchangeHoursOnly(data, daybegin = "09:30:00", dayend = "16:00:00")
```

### Arguments

data	an xts object containing the time series data.
daybegin	character in the format of "HH:MM:SS", specifying the starting hour, minute and second of an exchange trading day.
dayend	character in the format of "HH:MM:SS", specifying the closing hour, minute and second of an exchange trading day.

### Value

xts object

### Author(s)

Jonathan Cornelissen and Kris Boudt

### References

Brownlees, C.T. and Gallo, G.M. (2006). Financial econometric analysis at ultra-high frequency: Data handling concerns. *Computational Statistics & Data Analysis*, 51, pages 2232-2245.

---

getTradeDirection	<i>Get trade direction</i>
-------------------	----------------------------

---

### Description

Function returns a vector with the inferred trade direction which is determined using the Lee and Ready algorithm (Lee and Ready, 1991).

### Usage

```
getTradeDirection(tqdata,...);
```

### Arguments

tqdata	xts object, containing joined trades and quotes (e.g. using <a href="#">matchTradesQuotes</a> )
...	additional arguments.

**Value**

A vector which has values 1 or (-1) if the inferred trade direction is buy or sell respectively.

**Details**

NOTE: The value of the first (and second) observation of the output should be ignored if price==midpoint for the first (second) observation.

**Author(s)**

Jonathan Cornelissen and Kris Boudt

**References**

Lee, C. M. C. and M. J. Ready (1991). Inferring trade direction from intraday data. *Journal of Finance* 46, 733-746.

---

harModel	<i>HAR model estimation (Heterogeneous Autoregressive model for Realized volatility)</i>
----------	--

---

**Description**

Function returns the estimates for the Heterogeneous Autoregressive model for Realized volatility discussed in Andersen et al. (2007) and Corsi (2009). This model is mainly used to forecast the next days' volatility based on the high-frequency returns of the past. Consult the vignette for more information.

**Usage**

```
harModel(data, periods = c(1, 5, 22), periodsJ = c(1,5,22), leverage=NULL, RVest = c("rCov", "r
transform = NULL, ...)
```

**Arguments**

data	an xts-object containing the intraday (log-)returns.
periods	a vector of integers indicating over how days the realized measures in the model should be aggregated. By default periods = c(1,5,22), which corresponds to one day, one week and one month respectively. This default is in line with Andersen et al. (2007).
periodsJ	a vector of integers indicating over what time periods the jump components in the model should be aggregated. By default periodsJ = c(1,5,22), which corresponds to one day, one week and one month respectively.
leverage	a vector of integers indicating over what periods the negative returns should be aggregated. See Corsi and Reno (2012) for more information. By default leverage=NULL and the model assumes the absence of a leverage effect. Set leverage=c(1,5,22) to mimic the analysis in Corsi and Reno (2012).

<code>RVest</code>	a character vector with one or two elements. The first element refers to the name of the function to estimate the daily integrated variance (non-jump-robust), while the second element refers to the name of the function to estimate the continuous component of daily volatility (jump-robust). By default <code>RVest = c("rCov","rBPCov")</code> , i.e. using the Realized Volatility and Realized Bi-Power Variance.
<code>type</code>	a string referring to the type of HAR model you would like to estimate. By default <code>type = "HARRV"</code> , the most basic model. Other valid options are <code>type = "HARRVJ"</code> or <code>type = "HARRVCJ"</code> .
<code>jumpstest</code>	the function name of a function used to test whether the test statistic which determines whether the jump variability is significant that day. By default <code>jumpstest = "ABDJumpstest"</code> , hence using the test statistic in Equation or Equation (18) of Andersen et al. (2007).
<code>alpha</code>	a real indicating the confidence level used in testing for jumps. By default <code>alpha = 0.05</code> .
<code>h</code>	an integer indicating the number over how many days the dependent variable should be aggregated. By default, <code>h=1</code> , i.e. no aggregation takes place, you just model the daily realized volatility.
<code>transform</code>	optionally a string referring to a function that transforms both the dependent and explanatory variables in the model. By default <code>transform=NULL</code> , so no transformation is done. Typical other choices in this context would be <code>"log"</code> or <code>"sqrt"</code> .
<code>...</code>	extra arguments

### Value

The function outputs an object of class `harModel` and `lm` (so `harModel` is a subclass of `lm`). So far I only added a print method as you can see in the examples. Input here is welcome, what should a plot of an "harModel" object look like? What other methods are useful?

### Details

See vignette.

### Author(s)

Jonathan Cornelissen and Kris Boudt

### References

- Andersen, T. G., T. Bollerslev, and F. Diebold (2007). Roughing it up: including jump components in the measurement, modelling and forecasting of return volatility. *The Review of Economics and Statistics* 89, 701-720.
- Corsi, F. (2009). A simple approximate long memory model of realized volatility. *Journal of Financial Econometrics* 7, 174-196.
- Corsi, F. and Reno R. (2012). Discrete-time volatility forecasting with persistent leverage effect and the link with continuous-time volatility modeling. *Journal of Business and Economic Statistics*, forthcoming.

## Examples

```
##### Example 1: HARRVCJ #####
data("sample_5minprices_jumps");
data = sample_5minprices_jumps[,1];
data = makeReturns(data); #Get the high-frequency return data

x = harModel(data, periods = c(1,5,10), periodsJ=c(1,5,10), RVest = c("rCov","rBPCov"),
             type="HARRVCJ",transform="sqrt"); # Estimate the HAR model of type HARRVCJ
class(x);
x

##### Example 2: #####
# Forecasting daily Realized volatility for DJI 2008 using the basic harModel: HARRV
data(realized_library); #Get sample daily Realized Volatility data
DJI_RV = realized_library$Dow.Jones.Industrials.Realized.Variance; #Select DJI
DJI_RV = DJI_RV[!is.na(DJI_RV)]; #Remove NA's
DJI_RV = DJI_RV['2008'];

x = harModel(data=DJI_RV , periods = c(1,5,22), RVest = c("rCov"), type="HARRV",h=1,transform=NULL);
class(x);
x;
summary(x);
plot(x);
```

heavyModel

*HEAVY Model estimation*

## Description

This function calculates the High frEQUENCY bAsed VolatilitY (HEAVY) model proposed in Shephard and Sheppard (2010).

## Usage

```
heavyModel(data, p=matrix( c(0,0,1,1),ncol=2 ), q=matrix( c(1,0,0,1),ncol=2 ),
           startingvalues = NULL, LB = NULL, UB = NULL,
           backcast = NULL, compconst = FALSE);
```

## Arguments

data	a (T x K) matrix containing the data, with T the number of days. For the traditional HEAVY model: K = 2, the first column contains the squared daily de-meaned returns, the second column contains the realized measures.
p	a (K x K) matrix containing the lag length for the model innovations. Position (i, j) in the matrix indicates the number of lags in equation i of the model for the innovations in data column j. For the traditional heavy model p is given by matrix( c(0,0,1,1),ncol=2 ) (default).
q	a (K x K) matrix containing the lag length for the conditional variances. Position (i, j) in the matrix indicates the number of lags in equation i of the model for conditional variances corresponding to series j. For the traditional heavy model introduced above q is given by matrix( c(1,0,0,1),ncol=2 ) (default).

startingvalues	a vector containing the starting values to be used in the optimization to find the optimal parameters estimates.
LB	a vector of length K indicating the lower bounds to be used in the estimation. If NULL it is set to a vector of zeros by default.
UB	a vector of length K indicating the upper bounds to be used in the estimation. If NULL it is set to a vector of Inf by default.
backcast	a vector of length K used to initialize the estimation. If NULL the unconditional estimates are taken.
compconst	a boolean variable. In case TRUE, the omega values are estimated in the optimization. In case FALSE, volatility targeting is done and omega is just 1 minus the sum of all relevant alpha's and beta's multiplied by the unconditional variance.

### Value

A list with the following values: (i) loglikelihood: The log likelihood evaluated at the parameter estimates. (ii) likelihoods: an xts object of length T containing the log likelihoods per day. (iii) condvar: a (T x K) xts object containing the conditional variances (iv) estparams: a vector with the parameter estimates. The order in which the parameters are reported is as follows: First the estimates for omega then the estimates for the non-zero alpha's with the most recent lags first in case  $\max(p) > 1$ , then the estimates for the non-zero beta's with the most recent lag first in case  $\max(q) > 1$ . (v) convergence: an integer code indicating the successfulness of the optimization. See `optim` for more information.

### Details

See vignette. NOTE: The implementation of the heavyModel is not completely finished. For the moment only bound constraints on the parameters are imposed in the optimization. Future developments also include outputting standard errors, and a c-implementation of the likelihood function to speed up the QML estimation.

### Author(s)

Jonathan Cornelissen and Kris Boudt

### References

Shephard, N. and K. Sheppard (2010). Realising the future: forecasting with high frequency based volatility (heavy) models. *Journal of Applied Econometrics* 25, 197-231.

### Examples

```
# Implementation of the heavy model on DJI:
data("realized_library");
returns = realized_library$Dow.Jones.Industrials>Returns;
rk      = realized_library$Dow.Jones.Industrials.Realized.Kernel;
returns = returns[!is.na(rk)]; rk = rk[!is.na(rk)]; # Remove NA's
data = cbind( returns^2, rk ); # Make data matrix with returns and realized measures
backcast = matrix( c(var(returns),mean(rk)) ,ncol=1);

startvalues = c(0.004,0.02,0.44,0.41,0.74,0.56); # Initial values
# output = heavyModel( data = as.matrix(data,ncol=2), compconst=FALSE,
#                      startingvalues = startvalues, backcast=backcast);
```



lltc.xts

*LLTC Data***Description**

Tick data for LLTC 2011/07/01, cleaned with [tradesCleanup](#).

**Usage**

```
data(lltc.xts)
```

**Examples**

```
data(lltc.xts)
plot(lltc.xts)
```

makePsd

*Returns the positive semidefinite projection of a symmetric matrix using the eigenvalue method*

**Description**

Function returns the positive semidefinite projection of a symmetric matrix using the eigenvalue method.

**Usage**

```
makePsd(S,method="covariance")
```

**Arguments**

S	matrix.
method	character, indicating whether the negative eigenvalues of the correlation or covariance should be replaced by zero. Possible values are "covariance" and "correlation".

**Value**

An xts object containing the aggregated trade data.

**Details**

We use the eigenvalue method to transform  $S$  into a positive semidefinite covariance matrix (see e.g. Barndorff-Nielsen and Shephard, 2004, and Rousseeuw and Molenberghs, 1993). Let  $\Gamma$  be the orthogonal matrix consisting of the  $p$  eigenvectors of  $S$ . Denote  $\lambda_1^+, \dots, \lambda_p^+$  its  $p$  eigenvalues, whereby the negative eigenvalues have been replaced by zeroes. Under this approach, the positive semi-definite projection of  $S$  is  $S^+ = \Gamma' \text{diag}(\lambda_1^+, \dots, \lambda_p^+) \Gamma$ .

If method="correlation", the eigenvalues of the correlation matrix corresponding to the matrix  $S$  are transformed. See Fan et al (2010).

**Author(s)**

Jonathan Cornelissen and Kris Boudt

**References**

Barndorff-Nielsen, O. and N. Shephard (2004). Measuring the impact of jumps in multivariate price processes using bipower covariation. Discussion paper, Nuffield College, Oxford University.

Fan, J., Y. Li, and K. Yu (2010). Vast volatility matrix estimation using high frequency data for portfolio selection. Working paper.

Rousseeuw, P. and G. Molenberghs (1993). Transformation of non positive semidefinite correlation matrices. Communications in Statistics - Theory and Methods 22, 965-984.

---

makeReturns

*Compute log returns*

---

**Description**

Function returns an xts object with the log returns as xts object.

$$\log \text{return}_t = (\log(\text{PRICE}_t) - \log(\text{PRICE}_{t-1})).$$

**Usage**

```
makeReturns(ts);
```

**Arguments**

ts                      xts object

**Value**

an xts object containing the log returns.

**Details**

Note: the first (row of) observation(s) is set to zero.

**Author(s)**

Jonathan Cornelissen and Kris Boudt

---

matchTradesQuotes	<i>Match trade and quote data</i>
-------------------	-----------------------------------

---

## Description

Function matches the trades and quotes and returns an xts-object containing both.

## Usage

```
matchTradesQuotes(tdata,qdata,adjustment=2)
```

## Arguments

tdata	xts-object containing the trade data.
qdata	xts-object containing the quote data.
adjustment	numeric, number of seconds the quotes are registered faster than the trades (should be round and positive). Based on the research of Vergote (2005), we set 2 seconds as the default.

## Value

xts-object containing the matched trade and quote data

## Author(s)

Jonathan Cornelissen and Kris Boudt

## References

Vergote, O. (2005). How to match trades and quotes for NYSE stocks? K.U.Leuven working paper.

## Examples

```
#load data samples
data("sample_tdata");
data("sample_qdata");
#match the trade and quote data
tqdata = matchTradesQuotes(sample_tdata,sample_qdata);
#have a look
head(tqdata);
```

---

medRV	<i>medRV</i>
-------	--------------

---

### Description

Function returns the medRV, defined in Andersen et al. (2009).

Let  $r_{t,i}$  be a return (with  $i = 1, \dots, M$ ) in period  $t$ .

Then, the medRV is given by

$$\text{medRV}_t = \frac{\pi}{6 - 4\sqrt{3} + \pi} \left( \frac{M}{M-2} \right) \sum_{i=2}^{M-1} \text{med}(|r_{t,i-1}|, |r_{t,i}|, |r_{t,i+1}|)^2$$

### Usage

```
medRV(rdata, align.by=NULL, align.period=NULL, makeReturns=FALSE, ...)
```

### Arguments

<code>rdata</code>	a zoo/xts object containing all returns in period $t$ for one asset.
<code>align.by</code>	a string, align the tick data to "seconds" "minutes" "hours".
<code>align.period</code>	an integer, align the tick data to this many [seconds minutes hours].
<code>makeReturns</code>	boolean, should be TRUE when <code>rdata</code> contains prices instead of returns. FALSE by default.
<code>...</code>	additional arguments.

### Value

numeric

### Details

The medRV belongs to the class of realized volatility measures in this package that use the series of high-frequency returns  $r_{t,i}$  of a day  $t$  to produce an ex post estimate of the realized volatility of that day  $t$ . medRV is designed to be robust to price jumps. The difference between RV and medRV is an estimate of the realized jump variability. Disentangling the continuous and jump components in RV can lead to more precise volatility forecasts, as shown in Andersen et al. (2007) and Corsi et al. (2010).

### Author(s)

Jonathan Cornelissen and Kris Boudt

### References

Andersen, T. G., D. Dobrev, and E. Schaumburg (2012). Jump-robust volatility estimation using nearest neighbor truncation. *Journal of Econometrics*, 169 (1), 75-93.

Andersen, T.G., T. Bollerslev, and F. Diebold (2007). Roughing it up: including jump components in the measurement, modelling and forecasting of return volatility. *The Review of Economics and Statistics* 89 (4), 701-720.

Corsi, F., D. Pirino, and R. Reno (2010). Threshold Bipower Variation and the Impact of Jumps on Volatility Forecasting. *Journal of Econometrics* 159 (2), 276-288.

**Examples**

```
data(sample_tdata);
medrv = medRV( rdata = sample_tdata$PRICE, align.by = "minutes",
               align.period = 5, makeReturns=TRUE);
medrv
```

---

mergeQuotesSameTimestamp

*Merge multiple quote entries with the same time stamp*

---

**Description**

Function replaces multiple quote entries that have the same time stamp by a single one and returns an xts object with unique time stamps only.

**Usage**

```
mergeQuotesSameTimestamp(qdata,selection="median")
```

**Arguments**

- |           |   |
|-----------|---|
| qdata     | an xts object containing the time series data, with at least two columns named "BID" and "OFR" indicating the bid and ask price and two columns "BIDSIZ", "OFRSIZ" indicating the number of round lots available at these prices.   |
| selection | <p>indicates how the bid and ask price for a certain time stamp should be calculated in case of multiple observation for a certain time stamp. By default, selection="median", and the median price is taken. Alternatively:</p> <ul style="list-style-type: none"> <li>• selection = "maxvolume": use the (bid/ask) price of the entry with largest (bid/ask) volume.</li> <li>• selection = "weightedaverage": take the weighted average of all bid (ask) prices, weighted by "BIDSIZ" ("OFRSIZ").</li> </ul> |

**Value**

xts object

**Author(s)**

Jonathan Cornelissen and Kris Boudt

---

mergeTradesSameTimestamp

*Merge multiple transactions with the same time stamp*


---

### Description

Function replaces multiple transactions that have the same time stamp by a single one and returns an xts object with unique time stamps only.

### Usage

```
mergeTradesSameTimestamp(tdata,selection="median")
```

### Arguments

- |           |   |
|-----------|---|
| tdata     | an xts object containing the time series data, with one column named "PRICE" indicating the transaction price and one column "SIZE" indicating the number of shares traded.   |
| selection | indicates how the price for a certain time stamp should be calculated in case of multiple observation for a certain time stamp. By default, selection="median", and the median price is taken. Alternatively: <ul style="list-style-type: none"> <li>• selection = "maxvolume": use the price of the transaction with largest volume.</li> <li>• selection = "weightedaverage": take the weighted average of all prices.</li> </ul> |

### Value

xts object

### Author(s)

Jonathan Cornelissen and Kris Boudt

---

minRV

*minRV*


---

### Description

Function returns the minRV, defined in Andersen et al. (2009).

Let  $r_{t,i}$  be a return (with  $i = 1, \dots, M$ ) in period  $t$ .

Then, the minRV is given by

$$\text{minRV}_t = \frac{\pi}{\pi - 2} \left( \frac{M}{M - 1} \right) \sum_{i=1}^{M-1} \min(|r_{t,i}|, |r_{t,i+1}|)^2$$

### Usage

```
minRV(rdata,align.by=NULL,align.period=NULL,makeReturns=FALSE,...)
```

**Arguments**

<code>rdata</code>	a zoo/xts object containing all returns in period <code>t</code> for one asset.
<code>align.by</code>	a string, align the tick data to "seconds" "minutes" "hours".
<code>align.period</code>	an integer, align the tick data to this many [seconds minutes hours].
<code>makeReturns</code>	boolean, should be TRUE when <code>rdata</code> contains prices instead of returns. FALSE by default.
<code>...</code>	additional arguments.

**Value**

numeric

**Author(s)**

Jonathan Cornelissen and Kris Boudt

**References**

Andersen, T. G., D. Dobrev, and E. Schaumburg (2012). Jump-robust volatility estimation using nearest neighbor truncation. *Journal of Econometrics*, 169 (1), 75-93.

**Examples**

```
data(sample_tdata);

minrv = minRV( rdata = sample_tdata$PRICE, align.by = "minutes",
               align.period = 5, makeReturns=TRUE);
minrv
```

---

noZeroPrices

---

*Delete the observations where the price is zero*


---

**Description**

Function deletes the observations where the price is zero.

**Usage**

```
noZeroPrices(tdata)
```

**Arguments**

<code>tdata</code>	an xts object at least containing a column "PRICE".
--------------------	---

**Value**

xts object

**Author(s)**

Jonathan Cornelissen and Kris Boudt

---

noZeroQuotes	<i>Delete the observations where the bid or ask is zero</i>
--------------	---

---

**Description**

Function deletes the observations where the bid or ask is zero.

**Usage**

```
noZeroQuotes(qdata)
```

**Arguments**

qdata                    an xts object at least containing the columns "BID" and "OFR".

**Value**

xts object

**Author(s)**

Jonathan Cornelissen and Kris Boudt

---

previoustick	<i>previoustick (internal function)</i>
--------------	---

---

**Description**

See [aggregatets](#).

---

quotesCleanup	<i>Cleans quote data</i>
---------------	--------------------------

---

**Description**

This is a wrapper function for cleaning the quote data of all stocks in "ticker" over the interval [from,to]. The result is saved in the folder datadestination.

In case you supply the argument "rawqdata", the on-disk functionality is ignored and the function returns a list with the cleaned quotes as xts object (see examples).

The following cleaning functions are performed sequentially: [noZeroQuotes](#), [selectExchange](#), [rmLargeSpread](#), [mergeQuotesSameTimestamp](#), [rmOutliers](#).

**Usage**

```
quotesCleanup(from,to,datasource,datadestination,ticker,exchanges,
               qdataraw,report,selection,maxi>window,type,
               rmoutliersmaxi,...)
```



**Arguments**

from	character indicating first date to clean, e.g. "2008-01-30".
to	character indicating last date to clean, e.g. "2008-01-31".
datasource	character indicating the folder in which the original data is stored.
datadestination	character indicating the folder in which the cleaned data is stored.
ticker	vector of tickers for which the data should be cleaned, e.g. <code>ticker = c("AAPL", "AIG")</code> .
exchanges	vector of stock exchange symbols for all tickers in vector "ticker". It thus should have the same length as the vector ticker. Only data from one exchanges will be retained for each stock respectively, e.g. <code>exchanges = c("T", "N")</code> . The possible exchange symbols are: <ul style="list-style-type: none"> <li>• A: AMEX</li> <li>• N: NYSE</li> <li>• B: Boston</li> <li>• P: Arca</li> <li>• C: NSX</li> <li>• T/Q: NASDAQ</li> <li>• D: NASD ADF and TRF</li> <li>• X: Philadelphia</li> <li>• I: ISE</li> <li>• M: Chicago</li> <li>• W: CBOE</li> <li>• Z: BATS</li> </ul>
qdataraw	xts object containing (ONE day and for ONE stock only) raw quote data. This argument is NULL by default. Enabling it means the arguments from, to, data-source and datadestination will be ignored. (only advisable for small chunks of data)
report	boolean and TRUE by default. In case it is true the function returns (also) a vector indicating how many quotes remained after each cleaning step.
selection	argument to be passed on to the cleaning routine <a href="#">mergeQuotesSameTimestamp</a> . The default is "median".
maxi	argument to be passed on to the cleaning routine <a href="#">rmLargeSpread</a> .
window	argument to be passed on to the cleaning routine <a href="#">rmOutliers</a> .
type	argument to be passed on to the cleaning routine <a href="#">rmOutliers</a> .
rmoutliersmaxi	argument to be passed on to the cleaning routine <a href="#">rmOutliers</a> .
...	additional arguments.

**Value**

For each day an xts object is saved into the folder of that date, containing the cleaned data. This procedure is performed for each stock in "ticker". The function returns a vector indicating how many quotes remained after each cleaning step.

In case you supply the argument "rawqdata", the on-disk functionality is ignored and the function returns a list with the cleaned quotes as xts object (see examples).

**Author(s)**

Jonathan Cornelissen and Kris Boudt

**References**

Barndorff-Nielsen, O. E., P. R. Hansen, A. Lunde, and N. Shephard (2009). Realized kernels in practice: Trades and quotes. *Econometrics Journal* 12, C1-C32.

Brownlees, C.T. and Gallo, G.M. (2006). Financial econometric analysis at ultra-high frequency: Data handling concerns. *Computational Statistics & Data Analysis*, 51, pages 2232-2245.

Falkenberry, T.N. (2002). High frequency data filtering. Unpublished technical report.

**Examples**

```
#Consider you have raw quote data for 1 stock for 1 day
#data("sample_qdataraw");
#head(sample_qdataraw);
#dim(sample_qdataraw);
#qdata_aftercleaning = quotesCleanup(qdataraw=sample_qdataraw,exchanges="N");
#qdata_aftercleaning$report;
#barplot(qdata_aftercleaning$report);
#dim(qdata_aftercleaning$qdata);

#In case you have more data it is advised to use the on-disk functionality
#via "from","to","datasource",etc. arguments
```

---

rAccumulation

*Realized Accumulation Plot*


---

**Description**

Plots the realized estimate as it accumulates over a time interval.

**Usage**

```
rAccumulation(x, period = 1, y = NULL, align.by="seconds", align.period = 1, plotit = FALSE, cts
```

**Arguments**

x	Tick data in xts object.
y	Tick data in xts object.
period	Sampling period
align.by	Align the tick data to seconds minutes hours
align.period	Align the returns to this period first
plotit	T for plot
cts	Create calendar time sampling if a non realizedObject is passed
makeReturns	Prices are passed make them into log returns

**Details**

Plots the realized estimate as it accumulates over a time interval. This is a good tool to determine what observations are adding (possibly subtracting for covariance) to the estimate.

**Value**

Realized accumulation vector if plotit = F

**Author(s)**

Scott Payseur <scott.payseur@gmail.com>

**References**

S. W. Payseur. A One Day Comparison of Realized Variance and Covariance Estimators. *Working Paper: University of Washington*, 2007

**See Also**

[rMarginal](#)

**Examples**

```
data(sbx.xts)

cumm <- list()
cumm[[1]] <- rCumSum(sbx.xts, period=1, align.by="seconds", align.period=60)
cumm[[2]] <- rCumSum(sbx.xts, period=10, align.by="seconds", align.period=60)
cumm[[3]] <- rCumSum(sbx.xts, period=20, align.by="seconds", align.period=60)
cumm[[4]] <- rCumSum(sbx.xts, period=30, align.by="seconds", align.period=60)
accum <- list()
accum[[1]] <- rAccumulation(sbx.xts, period=10, align.by="seconds", align.period=60)
accum[[2]] <- rAccumulation(sbx.xts, period=20, align.by="seconds", align.period=60)
accum[[3]] <- rAccumulation(sbx.xts, period=30, align.by="seconds", align.period=60)

par(mfrow=c(2,1))
plot(cumm[[1]], xlab="", ylab="Cumulative Returns", main="Starbucks (SBUX)", sub='20110701', type="p", col=1, lwd=2)
lines(cumm[[2]], col=2, lwd=2)
lines(cumm[[3]], col=3, lwd=2)
lines(cumm[[4]], col=4, lwd=2)
plot(accum[[1]], xlab="", ylab="Realized Accumulation", type="l", main="Starbucks (SBUX)", sub='20110701', col=1, lwd=2)
lines(accum[[2]], col=3, lwd=2)
lines(accum[[3]], col=4, lwd=2)
```

---

rAVGCov

*Realized Covariance: Average Subsample*


---

**Description**

Realized Covariance using average subsample.

**Usage**

```
rAVGCov(rdata, cor = FALSE, period = 1, align.by = "seconds",
        align.period = 1, cts = TRUE, makeReturns = FALSE, ...)
```

**Arguments**

rdata	In the multivariate case: a list. Each list-item <i>i</i> contains an xts object with the intraday data of stock <i>i</i> for day <i>t</i> . In the univariate case: an xts object containing the (tick) data for one day.
cor	boolean, in case it is TRUE, the correlation is returned. FALSE by default.
period	Sampling period
align.by	Align the tick data to seconds minutes hours
align.period	Align the tick data to this many [seconds minutes hours]
cts	Create calendar time sampling if a non realizedObject is passed
makeReturns	Prices are passed make them into log returns
...	...

**Value**

Realized covariance using average subsample.

**Author(s)**

Scott Payseur <scott.payseur@gmail.com>

**References**

L. Zhang, P.A Mykland, and Y. Ait-Sahalia. A tale of two time scales: Determining integrated volatility with noisy high-frequency data. *Journal of the American Statistical Association*, 2005.

Michiel de Pooter, Martin Martens, and Dick van Dijk. Predicting the daily covariance matrix for sp100 stocks using intraday data - but which frequency to use? *Econometric Reviews*, 2008.

**Examples**

```
# Average subsampled realized variance/covariance for CTS aligned at one minute returns at
# 5 subgrids (5 minutes).
data(sample_tdata);
data(1l1tc.xts);
data(sbox.xts);

# Univariate
rvSub = rAVGCov( rdata = sample_tdata$PRICE, period = 5, align.by = "minutes",
                align.period=5, makeReturns=TRUE);
rvSub

# Multivariate:
rcSub = rAVGCov( rdata = list(1l1tc.xts,sbox.xts), period = 5, align.by = "minutes",
                align.period=5, makeReturns=FALSE);
rcSub
```

rBPCov

*Realized BiPower Covariance***Description**

Function returns the Realized BiPower Covariance (rBPCov), defined in Barndorff-Nielsen and Shephard (2004).

Let  $r_{t,i}$  be an intraday  $N \times 1$  return vector and  $i = 1, \dots, M$  the number of intraday returns.

The rBPCov is defined as the process whose value at time  $t$  is the  $N$ -dimensional square matrix with  $k, q$ -th element equal to

$$\text{rBPCov}[k, q]_t = \frac{\pi}{8} \left( \sum_{i=2}^M |r_{(k)t,i} + r_{(q)t,i}| |r_{(k)t,i-1} + r_{(q)t,i-1}| - |r_{(k)t,i} - r_{(q)t,i}| |r_{(k)t,i-1} - r_{(q)t,i-1}| \right),$$

where  $r_{(k)t,i}$  is the  $k$ -th component of the return vector  $r_{i,t}$ .

**Usage**

```
rBPCov(rdata, cor = FALSE, align.by = NULL, align.period = NULL,
       makeReturns = FALSE, makePsd = FALSE, ...)
```

**Arguments**

rdata	a ( $M \times N$ ) matrix/zoo/xts object containing the $N$ return series over period $t$ , with $M$ observations during $t$ .
cor	boolean, in case it is TRUE, the correlation is returned. FALSE by default.
align.by	a string, align the tick data to "seconds" "minutes" "hours".
align.period	an integer, align the tick data to this many [seconds minutes hours].
makeReturns	boolean, should be TRUE when rdata contains prices instead of returns. FALSE by default.
makePsd	boolean, in case it is TRUE, the positive definite version of rBPCov is returned. FALSE by default.
...	additional arguments.

**Value**

an  $N \times N$  matrix

**Author(s)**

Jonathan Cornelissen and Kris Boudt

**References**

Barndorff-Nielsen, O. and N. Shephard (2004). Measuring the impact of jumps in multivariate price processes using bipower covariation. Discussion paper, Nuffield College, Oxford University.

**Examples**

```
# Realized Bipower Variance/Covariance for CTS aligned
# at 5 minutes.
data(sample_tdata);
data(sample_5minprices_jumps);

# Univariate:
rbpv = rBPCov( rdata = sample_tdata$PRICE, align.by = "minutes",
               align.period = 5, makeReturns=TRUE);
rbpv

# Multivariate:
rbpc = rBPCov( rdata = sample_5minprices_jumps['2010-01-04'], makeReturns=TRUE,makePsd=TRUE);
rbpc
```

---

rCov	<i>Realized Covariance</i>
------	----------------------------

---

**Description**

Function returns the Realized Covariation (rCov).

Let  $r_{t,i}$  be an intraday  $N \times 1$  return vector and  $i = 1, \dots, M$  the number of intraday returns.

Then, the rCov is given by

$$\text{rCov}_t = \sum_{i=1}^M r_{t,i} r'_{t,i}.$$

**Usage**

```
rCov(rdata, cor = FALSE, align.by = NULL, align.period = NULL, makeReturns = FALSE, ...)
```

**Arguments**

rdata	a ( $M \times N$ ) matrix/zoo/xts object containing the $N$ return series over period $t$ , with $M$ observations during $t$ .
cor	boolean, in case it is TRUE, the correlation is returned. FALSE by default.
align.by	a string, align the tick data to "seconds" "minutes" "hours".
align.period	an integer, align the tick data to this many [seconds minutes hours].
makeReturns	boolean, should be TRUE when rdata contains prices instead of returns. FALSE by default.
...	additional arguments.

**Value**

an  $N \times N$  matrix

**Author(s)**

Jonathan Cornelissen and Kris Boudt

**Examples**

```
# Realized Variance/Covariance for CTS aligned
# at 5 minutes.
data(sample_tdata);
data(sample_5minprices_jumps);

# Univariate:
rv = rCov( rdata = sample_tdata$PRICE, align.by = "minutes",
           align.period = 5, makeReturns=TRUE);
rv

# Multivariate:
rc = rCov( rdata = sample_5minprices_jumps['2010-01-04'], makeReturns=TRUE);
rc
```

rCumSum

*Plot cummulative returns***Description**

Plots cummulative returns at a certain alignment given a return series.

**Usage**

```
rCumSum(x, period = 1, align.by="seconds",align.period = 1, plotit = FALSE, type = "l", cts = TR
```

**Arguments**

x	Tick data in xts object.
period	Sampling period
align.by	Align the tick data to seconds minutes hours
align.period	Align the returns to this period first
plotit	T for plot
type	Line or points
cts	Create calendar time sampling if a non realizedObject is passed
makeReturns	Prices are passed make them into log returns

**Value**

Cummulative return vector if plotit = F

**Author(s)**

Scott Payseur <scott.payseur@gmail.com>

## Examples

```
data(sbox.xts)

cumm <- list()
cumm[[1]] <- rCumSum(sbox.xts, period=1, align.by="seconds", align.period=60)
cumm[[2]] <- rCumSum(sbox.xts, period=10, align.by="seconds", align.period=60)
cumm[[3]] <- rCumSum(sbox.xts, period=20, align.by="seconds", align.period=60)
cumm[[4]] <- rCumSum(sbox.xts, period=30, align.by="seconds", align.period=60)
plot(cumm[[1]], xlab="", ylab="Cumulative Returns", main="Starbucks (SBUX)", sub='20110701', type="p", col=
lines(cumm[[2]], col=2, lwd=2)
lines(cumm[[3]], col=3, lwd=2)
lines(cumm[[4]], col=4, lwd=2)
```

---

realized_library	<i>The realized library from the Oxford-Man Institute of Quantitative Finance</i>
------------------	---

---

## Description

An xts object containing the daily returns, daily Realized Variance and daily Realized Kernels ranging from 1996-01-03 up to 2009-03-01 for several indices and exchange rates. Use `colnames(realized_library)` to see which assets are included. The full library of the Oxford-Man Institute of Quantitative Finance can be found on their website: <http://realized.oxford-man.ox.ac.uk>.

## Usage

```
data("realized_library")
```

## Format

xts object

## References

Gerd Heber, Asger Lunde, Neil Shephard, and Kevin Sheppard (2009) "Oxford-Man Institute's realized library, version 0.1", Oxford-Man Institute, University of Oxford.

Shephard, N. and K. Sheppard (2010). Realising the future: forecasting with high frequency based volatility (heavy) models. *Journal of Applied Econometrics* 25, 197-231.

---

refreshTime	<i>Synchronize (multiple) irregular timeseries by refresh time</i>
-------------	--

---

## Description

This function implements the refresh time synchronization scheme proposed by Harris et al. (1995). It picks the so-called refresh times at which all assets have traded at least once since the last refresh time point. For example, the first refresh time corresponds to the first time at which all stocks have traded. The subsequent refresh time is defined as the first time when all stocks have again traded. This process is repeated until the end of one time series is reached.



**Usage**

```
refreshTime(pdata)
```

**Arguments**

`pdata` a list. Each list-item contains an xts object containing the original time series (one day only and typically a price series).

**Value**

An xts object containing the synchronized time series.

**Author(s)**

Jonathan Cornelissen and Kris Boudt

**References**

Harris, F., T. McNish, G. Shoesmith, and R. Wood (1995). Cointegration, error correction, and price discovery on informationally linked security markets. *Journal of Financial and Quantitative Analysis* 30, 563-581.

**Examples**

```
#suppose irregular timepoints:
start = as.POSIXct("2010-01-01 09:30:00")
ta = start + c(1,2,4,5,9);
tb = start + c(1,3,6,7,8,9,10,11);

#yielding the following timeseries:
a = as.xts(1:length(ta),order.by=ta);
b = as.xts(1:length(tb),order.by=tb);

#Calculate the synchronized timeseries:
refreshTime(list(a,b))
```

---

rHYCov

*Hayashi-Yoshida Covariance*


---

**Description**

Hayashi-Yoshida Covariance

**Usage**

```
rHYCov(rdata, cor = FALSE, period = 1, align.by = "seconds",
       align.period = 1, cts = TRUE, makeReturns = FALSE, makePsd = TRUE, ...)
```

**Arguments**

rdata	a list. Each list-item i contains an xts object with the intraday data of stock i for day t.
cor	boolean, in case it is TRUE, the correlation is returned. FALSE by default.
period	Sampling period
align.by	Align the tick data to seconds minutes hours
align.period	Align the tick data to this many [seconds minutes hours]
cts	Create calendar time sampling if a non realizedObject is passed
makeReturns	Prices are passed make them into log returns
makePsd	boolean, in case it is TRUE, the positive definite version of rTSCov is returned. FALSE by default.
...	...

**Author(s)**

Scott Payseur

**References**

T. Hayashi and N. Yoshida. On covariance estimation of non-synchronously observed diffusion processes. *Bernoulli*, 11:359-379, 2005.

**Examples**

```
# Average Realized Kernel Variance/Covariance for CTS aligned at one minute returns at
# 5 subgrids (5 minutes).
data(lltc.xts);
data(sbox.xts);
# Multivariate:
rHYCov = rHYCov( rdata = list(lltc.xts,sbox.xts), period = 5, align.by ="minutes",
                 align.period=5, makeReturns=FALSE);
rHYCov
#Note: for the diagonal elements the rCov is used.
```

---

rKernel.available	<i>Available Kernels</i>
-------------------	--------------------------

---

**Description**

Returns a vector of the available kernels.

**Usage**

```
rKernel.available()
```

**Author(s)**

Scott Payseur <scott.payseur@gmail.com>

## References

Ole E. Barndorff-Nielsen, Peter Reinhard Hansen, Asger Lunde, and Neil Shephard. Regular and modified kernel-based estimators of integrated variance: The case with independent noise. *Working Paper*, 2004.

## See Also

[rKernelCov](#)

## Examples

```
rKernel.available()
```

---

rKernelCov	<i>Realized Covariance: Kernel</i>
------------	------------------------------------

---

## Description

Realized covariance calculation using a kernel estimator.

## Usage

```
rKernelCov(rdata, cor = FALSE, kernel.type = "rectangular", kernel.param = 1,
            kernel.dofadj = TRUE, align.by = "seconds", align.period = 1,
            cts = TRUE, makeReturns = FALSE, type = NULL, adj = NULL,
            q = NULL, ...)
```

## Arguments

rdata	In the multivariate case: a list. Each list-item <i>i</i> contains an xts object with the intraday data of stock <i>i</i> for day <i>t</i> . In the univariate case: an xts object containing the (tick) data for one day.
cor	boolean, in case it is TRUE, the correlation is returned. FALSE by default.
kernel.type	Kernel name (or number)
kernel.param	Kernel parameter (usually lags)
kernel.dofadj	Kernel Degree of freedom adjustment
align.by	Align the tick data to seconds minutes hours
align.period	Align the tick data to this many [seconds minutes hours]
cts	Calendar Time Sampling is used
makeReturns	Convert to Returns
type	Deprecated, use kernel.type
adj	Deprecated, use kernel.dofadj
q	Deprecated, use kernel.param
...	...

**Details**

The different types of kernels can be found using [rKernel.available](#).

**Value**

Kernel estimate of realized covariance.

**Author(s)**

Scott Payseur <scott.payseur@gmail.com>

**References**

Ole E. Barndorff-Nielsen, Peter Reinhard Hansen, Asger Lunde, and Neil Shephard. Regular and modified kernel-based estimators of integrated variance: The case with independent noise. *Working Paper*, 2004.

B. Zhou. High-frequency data and volatility in foreign-exchange rates. *Journal of Business & Economic Statistics*, 14:45-52, 1996.

P. Hansen and A. Lunde. Realized variance and market microstructure noise. *Journal of Business and Economic Statistics*, 24:127-218, 2006.

**See Also**

[rKernel.available](#)

**Examples**

```
# Average Realized Kernel Variance/Covariance for CTS aligned at one minute returns at
# 5 subgrids (5 minutes).
data(sample_tdata);
data(1l1tc.xts);
data(sbox.xts);

# Univariate:
rvKernel = rKernelCov( rdata = sample_tdata$PRICE, period = 5, align.by = "minutes",
                      align.period=5, makeReturns=TRUE);
rvKernel

# Multivariate:
rcKernel = rKernelCov( rdata = list(1l1tc.xts,sbox.xts), period = 5, align.by = "minutes",
                      align.period=5, makeReturns=FALSE);
rcKernel
```

---

rMarginal

---

*Marginal Contribution to Realized Estimate*


---

**Description**

Plots the marginal contribution to the realized estimate.

**Usage**

```
rMarginal(x, y = NULL, period, align.by="seconds", align.period = 1, plotit = FALSE, cts = TRUE, makeReturns = FALSE)
```

**Arguments**

x	RealizedObject or TimeSeries for S+
y	RealizedObject or TimeSeries for S+
period	Sampling period
align.by	Align the tick data to seconds minutes hours
align.period	Align the returns to this period first
plotit	T for plot
cts	Create calendar time sampling if a non realizedObject is passed
makeReturns	Prices are passed make them into log returns

**Details**

Plots the marginal contribution to the realized estimate. This is a good tool to determine what observations are adding (possibly subtracting for covariance) to the estimate.

**Value**

Marginal contribution vector if plotit = F

**Author(s)**

Scott Payseur <scott.payseur@gmail.com>

**References**

S. W. Payseur. A One Day Comparison of Realized Variance and Covariance Estimators. *Working Paper: University of Washington*, 2007

**See Also**

[rAccumulation](#)

**Examples**

```
data(sbox.xts)
par(mfrow=c(2,1))
plot(rCumSum(sbox.xts, period=10, align.by="seconds", align.period=60), xlab="", ylab="Cumulative Returns")
barplot(rMarginal(sbox.xts, period=10, align.by="seconds", align.period=60)$y, main="Marginal Contribution")
```

---

rmLargeSpread	<i>Delete entries for which the spread is more than "maxi" times the median spread</i>
---------------	--

---

**Description**

Function deletes entries for which the spread is more than "maxi" times the median spread on that day.

**Usage**

```
rmLargeSpread(qdata,maxi)
```

**Arguments**

qdata	an xts object at least containing the columns "BID" and "OFR".
maxi	an integer. By default maxi="50", which means that entries are deleted if the spread is more than 50 times the median spread on that day.

**Value**

xts object

**Details**

NOTE: This function works only correct if supplied input data consists of 1 day!

**Author(s)**

Jonathan Cornelissen and Kris Boudt

---

rmNegativeSpread	<i>Delete entries for which the spread is negative</i>
------------------	--

---

**Description**

Function deletes entries for which the spread is negative.

**Usage**

```
rmNegativeSpread(qdata)
```

**Arguments**

qdata	an xts object at least containing the columns "BID" and "OFR".
-------	--

**Value**

xts object

**Author(s)**

Jonathan Cornelissen and Kris Boudt

---

rmOutliers	<i>Delete entries for which the mid-quote is outlying with respect to surrounding entries</i>
------------	---

---

## Description

If type="standard": Function deletes entries for which the mid-quote deviated by more than "maxi" median absolute deviations from a rolling centered median (excluding the observation under consideration) of "window" observations.

If type="advanced": Function deletes entries for which the mid-quote deviates by more than "maxi" median absolute deviations from the value closest to the mid-quote of these three options:

1. Rolling centered median (excluding the observation under consideration)
2. Rolling median of the following "window" observations
3. Rolling median of the previous "window" observations

The advantage of this procedure compared to the "standard" proposed by Barndorff-Nielsen et al. (2010) is that it will not incorrectly remove large price jumps. Therefore this procedure has been set as the default for removing outliers.

Note that the median absolute deviation is taken over the entire sample. In case it is zero (which can happen if mid-quotes don't change much), the median absolute deviation is taken over a subsample without constant mid-quotes.

## Usage

```
rmOutliers(qdata,maxi=10,window=50,type="advanced")
```

## Arguments

qdata	an xts object at least containing the columns "BID" and "OFR".
maxi	an integer, indicating the maximum number of median absolute deviations allowed.
window	an integer, indicating the time window for which the "outlyingness" is considered.
type	should be "standard" or "advanced" (see description).

## Value

xts object

## Details

NOTE: This function works only correct if supplied input data consists of 1 day.

## Author(s)

Jonathan Cornelissen and Kris Boudt

## References

Barndorff-Nielsen, O. E., P. R. Hansen, A. Lunde, and N. Shephard (2009). Realized kernels in practice: Trades and quotes. *Econometrics Journal* 12, C1-C32.

Brownlees, C.T. and Gallo, G.M. (2006). Financial econometric analysis at ultra-high frequency: Data handling concerns. *Computational Statistics & Data Analysis*, 51, pages 2232-2245.

---

rmTradeOutliers	<i>Delete transactions with unlikely transaction prices</i>
-----------------	---

---

## Description

Function deletes entries with prices that are above the ask plus the bid-ask spread. Similar for entries with prices below the bid minus the bid-ask spread.

## Usage

```
rmTradeOutliers(tdata,qdata)
```

## Arguments

tdata	an xts object containing the time series data, with at least the column "PRICE", containing the transaction price.
qdata	an xts object containing the time series data, with at least the columns "BID" and "OFR", containing the bid and ask prices.

## Value

xts object

## Details

Note: in order to work correctly, the input data of this function should be cleaned trade (tdata) and quote (qdata) data respectively.

## Author(s)

Jonathan Cornelissen and Kris Boudt



rOWCov

*Realized Outlyingness Weighted Covariance***Description**

Function returns the Realized Outlyingness Weighted Covariance, defined in Boudt et al. (2008).

Let  $r_{t,i}$ , for  $i = 1, \dots, M$  be a sample of  $M$  high-frequency ( $N \times 1$ ) return vectors and  $d_{t,i}$  their outlyingness given by the squared Mahalanobis distance between the return vector and zero in terms of the reweighted MCD covariance estimate based on these returns.

Then, the rOWCov is given by

$$\text{rOWCov}_t = c_w \frac{\sum_{i=1}^M w(d_{t,i}) r_{t,i} r_{t,i}'}{\frac{1}{M} \sum_{i=1}^M w(d_{t,i})},$$

The weight  $w_{i,\Delta}$  is one if the multivariate jump test statistic for  $r_{i,\Delta}$  in Boudt et al. (2008) is less than the 99.9% percentile of the chi-square distribution with  $N$  degrees of freedom and zero otherwise. The scalar  $c_w$  is a correction factor ensuring consistency of the rOWCov for the Integrated Covariance, under the Brownian Semimartingale with Finite Activity Jumps model.

**Usage**

```
rOWCov(rdata, cor = FALSE, align.by = NULL, align.period = NULL,
       makeReturns = FALSE, seasadjR = NULL, wfunction = "HR", alphaMCD = 0.75,
       alpha = 0.001, ...)
```

**Arguments**

rdata	a ( $M \times N$ ) matrix/zoo/xts object containing the $N$ return series over period $t$ , with $M$ observations during $t$ .
cor	boolean, in case it is TRUE, the correlation is returned. FALSE by default.
align.by	a string, align the tick data to "seconds" "minutes" "hours".
align.period	an integer, align the tick data to this many [seconds minutes hours].
makeReturns	boolean, should be TRUE when rdata contains prices instead of returns. FALSE by default.
seasadjR	a ( $M \times N$ ) matrix/zoo/xts object containing the seasonally adjusted returns. This is an optional argument.
wfunction	determines whether a zero-one weight function (one if no jump is detected based on $d_{t,i}$ and 0 otherwise) or Soft Rejection ("SR") weight function is to be used. By default a zero-one weight function (wfunction = "HR") is used.
alphaMCD	a numeric parameter, controlling the size of the subsets over which the determinant is minimized. Allowed values are between 0.5 and 1 and the default is 0.75. See Boudt et al. (2008) or the covMcd function in the robustbase package.
alpha	is a parameter between 0 and 0.5, that determines the rejection threshold value (see Boudt et al. (2008) for details).
...	additional arguments.

## Details

Advantages of the rOWCov compared to the rBPCov include a higher statistical efficiency, positive semidefiniteness and affine equivariance. However, the rOWCov suffers from a curse of dimensionality. Indeed, the rOWCov gives a zero weight to a return vector if at least one of the components is affected by a jump. In the case of independent jump occurrences, the average proportion of observations with at least one component being affected by jumps increases fast with the dimension of the series. This means that a potentially large proportion of the returns receives a zero weight, due to which the rOWCov can have a low finite sample efficiency in higher dimensions

## Value

an  $N \times N$  matrix

## Author(s)

Jonathan Cornelissen and Kris Boudt

## References

Boudt, K., C. Croux, and S. Laurent (2008). Outlyingness weighted covariation. Mimeo.

## Examples

```
# Realized Outlyingness Weighted Variance/Covariance for CTS aligned
# at 5 minutes.
data(sample_tdata);
data(sample_5minprices_jumps);

# Univariate:
rvoutw = rOWCov( rdata = sample_tdata$PRICE, align.by = "minutes",
                 align.period = 5, makeReturns=TRUE);
rvoutw

# Multivariate:
rcoutw = rOWCov( rdata = sample_5minprices_jumps['2010-01-04'], makeReturns=TRUE);
rcoutw
```

---

rRTSCov

---

*Robust two time scale covariance estimation*


---

## Description

Function returns the robust two time scale covariance matrix proposed in Boudt and Zhang (2010). Unlike the rOWCov, but similarly to the rThresholdCov, the rRTSCov uses univariate jump detection rules to truncate the effect of jumps on the covariance estimate. By the use of two time scales, this covariance estimate is not only robust to price jumps, but also to microstructure noise and non-synchronic trading.

## Usage

```
rRTSCov(pdata, cor=FALSE, startIV=NULL, noisevar = NULL,
        K = 300 , J = 1, K_cov = NULL , J_cov = NULL,
        K_var = NULL , J_var = NULL, eta = 9, makePsd = FALSE)
```

### Arguments

pdata	a list. Each list-item $i$ contains an xts object with the intraday price data of stock $i$ for day $t$ .
cor	boolean, in case it is TRUE, the correlation is returned. FALSE by default.
startIV	vector containing the first step estimates of the integrated variance of the assets, needed in the truncation. Is NULL by default.
noisevar	vector containing the estimates of the noise variance of the assets, needed in the truncation. Is NULL by default.
K	positive integer, slow time scale returns are computed on prices that are $K$ steps apart.
J	positive integer, fast time scale returns are computed on prices that are $J$ steps apart.
K_cov	positive integer, for the extradiagonal covariance elements the slow time scale returns are computed on prices that are $K$ steps apart.
J_cov	positive integer, for the extradiagonal covariance elements the fast time scale returns are computed on prices that are $J$ steps apart.
K_var	vector of positive integers, for the diagonal variance elements the slow time scale returns are computed on prices that are $K$ steps apart.
J_var	vector of positive integers, for the diagonal variance elements the fast time scale returns are computed on prices that are $J$ steps apart.
makePsd	boolean, in case it is TRUE, the positive definite version of rRTSCov is returned. FALSE by default.
eta	positive real number, squared standardized high-frequency returns that exceed $\eta$ are detected as jumps.

### Details

The rRTSCov requires the tick-by-tick transaction prices. (Co)variances are then computed using log-returns calculated on a rolling basis on stock prices that are  $K$  (slow time scale) and  $J$  (fast time scale) steps apart.

The diagonal elements of the rRTSCov matrix are the variances, computed for log-price series  $X$  with  $n$  price observations at times  $\tau_1, \tau_2, \dots, \tau_n$  as follows:

$$(1 - \frac{\bar{n}_K}{\bar{n}_J})^{-1}(\{X, X\}_T^{(K)*} - \frac{\bar{n}_K}{\bar{n}_J}\{X, X\}_T^{(J)*}),$$

where  $\bar{n}_K = (n - K + 1)/K$ ,  $\bar{n}_J = (n - J + 1)/J$  and

$$\{X, X\}_T^{(K)*} = \frac{c_\eta^* \sum_{i=1}^{n-K+1} (X_{t_{i+K}} - X_{t_i})^2 I_X^K(i; \eta)}{\frac{1}{n-K+1} \sum_{i=1}^{n-K+1} I_X^K(i; \eta)}.$$

The constant  $c_\eta$  adjusts for the bias due to the thresholding and  $I_X^K(i; \eta)$  is a jump indicator function that is one if

$$\frac{(X_{t_{i+K}} - X_{t_i})^2}{(\int_{t_i}^{t_{i+K}} \sigma_s^2 ds + 2\sigma_{\varepsilon_X}^2)} \leq \eta$$

and zero otherwise. The elements in the denominator are the integrated variance (estimated recursively) and noise variance (estimated by the method in Zhang et al, 2005).

The extradiagonal elements of the rRTSCov are the covariances. For their calculation, the data is first synchronized by the refresh time method proposed by Harris et al (1995). It uses the function `refreshTime` to collect first the so-called refresh times at which all assets have traded at least once since the last refresh time point. Suppose we have two log-price series:  $X$  and  $Y$ . Let  $\Gamma = \{\tau_1, \tau_2, \dots, \tau_{N_T^X}\}$  and  $\Theta = \{\theta_1, \theta_2, \dots, \theta_{N_T^Y}\}$  be the set of transaction times of these assets. The first refresh time corresponds to the first time at which both stocks have traded, i.e.  $\phi_1 = \max(\tau_1, \theta_1)$ . The subsequent refresh time is defined as the first time when both stocks have again traded, i.e.  $\phi_{j+1} = \max(\tau_{N_{\phi_j}^X+1}, \theta_{N_{\phi_j}^Y+1})$ . The complete refresh time sample grid is  $\Phi = \{\phi_1, \phi_2, \dots, \phi_{M_N+1}\}$ , where  $M_N$  is the total number of paired returns. The sampling points of asset  $X$  and  $Y$  are defined to be  $t_i = \max\{\tau \in \Gamma : \tau \leq \phi_i\}$  and  $s_i = \max\{\theta \in \Theta : \theta \leq \phi_i\}$ .

Given these refresh times, the covariance is computed as follows:

$$c_N(\{X, Y\}_T^{(K)} - \frac{\bar{n}_K}{\bar{n}_J} \{X, Y\}_T^{(J)}),$$

where

$$\{X, Y\}_T^{(K)} = \frac{1}{K} \frac{\sum_{i=1}^{M_N-K+1} c_i(X_{t_{i+K}} - X_{t_i})(Y_{s_{i+K}} - Y_{s_i}) I_X^K(i; \eta) I_Y^K(i; \eta)}{\frac{1}{M_N-K+1} \sum_{i=1}^{M_N-K+1} I_X^K(i; \eta) I_Y^K(i; \eta)},$$

with  $I_X^K(i; \eta)$  the same jump indicator function as for the variance and  $c_N$  a constant to adjust for the bias due to the thresholding.

Unfortunately, the rRTSCov is not always positive semidefinite. By setting the argument `makePsd = TRUE`, the function `makePsd` is used to return a positive semidefinite matrix. This function replaces the negative eigenvalues with zeroes.

## Value

an  $N \times N$  matrix

## Author(s)

Jonathan Cornelissen and Kris Boudt

## References

- Boudt K. and Zhang, J. 2010. Jump robust two time scale covariance estimation and realized volatility budgets. Mimeo.
- Harris, F., T. McInish, G. Shoesmith, and R. Wood (1995). Cointegration, error correction, and price discovery on informationally linked security markets. *Journal of Financial and Quantitative Analysis* 30, 563-581.
- Zhang, L., P. A. Mykland, and Y. Ait-Sahalia (2005). A tale of two time scales: Determining integrated volatility with noisy high-frequency data. *Journal of the American Statistical Association* 100, 1394-1411.

## Examples

```
# Robust Realized two timescales Variance/Covariance for CTS
data(sample_tdata);
data(1l1tc.xts);
data(sbox.xts);
```

```

# Univariate:
rvRTS = rRTSCov( pdata = sample_tdata$PRICE);
# Note: Prices as input
rvRTS

# Multivariate:
rcRTS = rRTSCov( pdata = list(cumsum(1l1tc.xts)+100,cumsum(sbox.xts)+100) );
# Note: List of prices as input
rcRTS

```

---

rScatterReturns	<i>Scatterplot of aligned returns</i>
-----------------	---------------------------------------

---

## Description

Creates a scatterplot of cross returns.

## Usage

```
rScatterReturns(x,y, period, align.by="seconds", align.period=1,numbers=FALSE,xlim= NULL, ylim=)
```

## Arguments

x	Tick data in xts object.
y	Tick data in xts object.
period	Sampling period
align.by	Align the tick data to seconds minutes hours
align.period	Align the returns to this period first
cts	Create calendar time sampling if a non realizedObject is passed
makeReturns	Prices are passed make them into log returns
plotit	T for plot
numbers	T for count
pch	type of point
ylim	ylimit
xlim	xlimit
scale.size	.
col.change	.
...	...

## Details

Scatterplot of returns.

## Author(s)

Scott Payseur <scott.payseur@gmail.com>

## References

S. W. Payseur. A One Day Comparison of Realized Variance and Covariance Estimators. *Working Paper: University of Washington, 2007*

## Examples

```
data(sbox.xts)
data(lltc.xts)
par(mfrow=c(2,1))
rScatterReturns(sbox.xts,y=lltc.xts, period=1, align.period=20,ylab="LLTC",xlab="SBUX",numbers=FALSE)
rScatterReturns(sbox.xts,y=lltc.xts, period=1, align.period=20,ylab="LLTC",xlab="SBUX",numbers=TRUE)
```

---

rThresholdCov	Threshold Covariance
---------------	----------------------

---

## Description

Function returns the threshold covariance matrix proposed in Gobbi and Mancini (2009). Unlike the [rOWCov](#), the THRESCov uses univariate jump detection rules to truncate the effect of jumps on the covariance estimate. As such, it remains feasible in high dimensions, but it is less robust to small cojumps.

Let  $r_{t,i}$  be an intraday  $N \times 1$  return vector and  $i = 1, \dots, M$  the number of intraday returns.

Then, the  $k, q$ -th element of the threshold covariance matrix is defined as

$$\text{thresholdcov}[k, q]_t = \sum_{i=1}^M r_{(k)t,i} 1_{\{r_{(k)t,i}^2 \leq TR_M\}} r_{(q)t,i} 1_{\{r_{(q)t,i}^2 \leq TR_M\}},$$

with the threshold value  $TR_M$  set to  $9\Delta^{-1}$  times the daily realized bi-power variation of asset  $k$ , as suggested in Jacod and Todorov (2009).

## Usage

```
rThresholdCov(rdata, cor = FALSE, align.by = NULL, align.period = NULL,
              makeReturns = FALSE, ...)
```

## Arguments

rdata	a $(M \times N)$ matrix/zoo/xts object containing the $N$ return series over period $t$ , with $M$ observations during $t$ .
cor	boolean, in case it is TRUE, the correlation is returned. FALSE by default.
align.by	a string, align the tick data to "seconds" "minutes" "hours".
align.period	an integer, align the tick data to this many [seconds minutes hours].
makeReturns	boolean, should be TRUE when rdata contains prices instead of returns. FALSE by default.
...	additional arguments.

## Value

an  $N \times N$  matrix

**Author(s)**

Jonathan Cornelissen and Kris Boudt

**References**

Barndorff-Nielsen, O. and N. Shephard (2004). Measuring the impact of jumps in multivariate price processes using bipower covariation. Discussion paper, Nuffield College, Oxford University.

Jacod, J. and V. Todorov (2009). Testing for common arrival of jumps in discretely-observed multidimensional processes. *Annals of Statistics* 37, 1792-1838.

Mancini, C. and F. Gobbi (2009). Identifying the covariation between the diffusion parts and the co-jumps given discrete observations. *Mimeo*.

**Examples**

```
# Realized threshold Variance/Covariance:
data(lltc.xts);
data(sbox.xts);

# Multivariate:
rcThreshold = rThresholdCov(cbind(lltc.xts,sbox.xts), align.by="minutes",align.period=1);
rcThreshold
```

---

rTSCov

---

*Two time scale covariance estimation*


---

**Description**

Function returns the two time scale covariance matrix proposed in Zhang et al (2005) and Zhang (2010). By the use of two time scales, this covariance estimate is robust to microstructure noise and non-synchronous trading.

**Usage**

```
rTSCov(pdata, cor=FALSE, K = 300 , J = 1, K_cov = NULL , J_cov = NULL,
       K_var = NULL , J_var = NULL, makePsd = FALSE);
```

**Arguments**

pdata	a list. Each list-item <i>i</i> contains an xts object with the intraday price data of stock <i>i</i> for day <i>t</i> .
cor	boolean, in case it is TRUE, the correlation is returned. FALSE by default.
K	positive integer, slow time scale returns are computed on prices that are <i>K</i> steps apart.
J	positive integer, fast time scale returns are computed on prices that are <i>J</i> steps apart.
K_cov	positive integer, for the extradiagonal covariance elements the slow time scale returns are computed on prices that are <i>K</i> steps apart.
J_cov	positive integer, for the extradiagonal covariance elements the fast time scale returns are computed on prices that are <i>J</i> steps apart.

K_var	vector of positive integers, for the diagonal variance elements the slow time scale returns are computed on prices that are K steps apart.
J_var	vector of positive integers, for the diagonal variance elements the fast time scale returns are computed on prices that are J steps apart.
makePsd	boolean, in case it is TRUE, the positive definite version of rTSCov is returned. FALSE by default.

### Details

The rTSCov requires the tick-by-tick transaction prices. (Co)variances are then computed using log-returns calculated on a rolling basis on stock prices that are  $K$  (slow time scale) and  $J$  (fast time scale) steps apart.

The diagonal elements of the rTSCov matrix are the variances, computed for log-price series  $X$  with  $n$  price observations at times  $\tau_1, \tau_2, \dots, \tau_n$  as follows:

$$(1 - \frac{\bar{n}_K}{\bar{n}_J})^{-1}([X, X]_T^{(K)} - \frac{\bar{n}_K}{\bar{n}_J}[X, X]_T^{(J)})$$

where  $\bar{n}_K = (n - K + 1)/K$ ,  $\bar{n}_J = (n - J + 1)/J$  and

$$[X, X]_T^{(K)} = \frac{1}{K} \sum_{i=1}^{n-K+1} (X_{t_{i+K}} - X_{t_i})^2.$$

The extradiagonal elements of the rTSCov are the covariances. For their calculation, the data is first synchronized by the refresh time method proposed by Harris et al (1995). It uses the function `refreshTime` to collect first the so-called refresh times at which all assets have traded at least once since the last refresh time point. Suppose we have two log-price series:  $X$  and  $Y$ . Let  $\Gamma = \{\tau_1, \tau_2, \dots, \tau_{N_T^X}\}$  and  $\Theta = \{\theta_1, \theta_2, \dots, \theta_{N_T^Y}\}$  be the set of transaction times of these assets. The first refresh time corresponds to the first time at which both stocks have traded, i.e.  $\phi_1 = \max(\tau_1, \theta_1)$ . The subsequent refresh time is defined as the first time when both stocks have again traded, i.e.  $\phi_{j+1} = \max(\tau_{N_{\phi_j}^X+1}, \theta_{N_{\phi_j}^Y+1})$ . The complete refresh time sample grid is  $\Phi = \{\phi_1, \phi_2, \dots, \phi_{M_N+1}\}$ , where  $M_N$  is the total number of paired returns. The sampling points of asset  $X$  and  $Y$  are defined to be  $t_i = \max\{\tau \in \Gamma : \tau \leq \phi_i\}$  and  $s_i = \max\{\theta \in \Theta : \theta \leq \phi_i\}$ .

Given these refresh times, the covariance is computed as follows:

$$c_N([X, Y]_T^{(K)} - \frac{\bar{n}_K}{\bar{n}_J}[X, Y]_T^{(J)}),$$

where

$$[X, Y]_T^{(K)} = \frac{1}{K} \sum_{i=1}^{M_N-K+1} (X_{t_{i+K}} - X_{t_i})(Y_{s_{i+K}} - Y_{s_i}).$$

Unfortunately, the rTSCov is not always positive semidefinite. By setting the argument `makePsd = TRUE`, the function `makePsd` is used to return a positive semidefinite matrix. This function replaces the negative eigenvalues with zeroes.

### Value

an  $N \times N$  matrix

### Author(s)

Jonathan Cornelissen and Kris Boudt



## References

Harris, F., T. McNish, G. Shoesmith, and R. Wood (1995). Cointegration, error correction, and price discovery on informationally linked security markets. *Journal of Financial and Quantitative Analysis* 30, 563-581.

Zhang, L., P. A. Mykland, and Y. Ait-Sahalia (2005). A tale of two time scales: Determining integrated volatility with noisy high-frequency data. *Journal of the American Statistical Association* 100, 1394-1411.

Zhang, L. (2011). Estimating covariation: Epps effect, microstructure noise. *Journal of Econometrics* 160, 33-47.

## Examples

```
# Robust Realized two timescales Variance/Covariance for CTS
data(sample_tdata);
data(lltc.xts);
data(sbox.xts);

# Univariate:
rvTS = rTSCov( pdata = sample_tdata$PRICE);
# Note: Prices as input
rvTS

# Multivariate:
rcTS = rTSCov( pdata = list(cumsum(lltc.xts)+100,cumsum(sbox.xts)+100) );
# Note: List of prices as input
rcTS
```

---

rZero	<i>Calculates the percentage of co-zero returns at a specified sampling period</i>
-------	--

---

## Description

Calculates the percentage of co-zero returns at a specified sampling period.

## Usage

```
rZero(rdata, period=1, align.by="seconds", align.period = 1, cts = TRUE, makeReturns = FALSE, ...)
```

## Arguments

rdata	an xts object containing the tick data or a list. In case of list: each list-item i contains an xts object with the intraday data of stock i for day t.
period	Sampling period
align.by	Align the tick data to seconds minutes hours
align.period	Align the returns to this period first
cts	Create calendar time sampling if a non realizedObject is passed
makeReturns	Prices are passed make them into log returns
...	...

**Value**

Percentage of co-zero returns.

**Author(s)**

Scott Payseur <scott.payseur@gmail.com>

**References**

S. W. Payseur. A One Day Comparison of Realized Variance and Covariance Estimators. *Working Paper: University of Washington*, 2007

**Examples**

```
data(sbox.xts)
data(lltc.xts)
rZero( rdata = list(sbox.xts, lltc.xts) , period = 60, align.by ="seconds", align.period=1)
```

---

salesCondition

*Delete entries with abnormal Sale Condition.*

---

**Description**

Function deletes entries with abnormal Sale Condition: trades where column "COND" has a letter code, except for "E" and "F".

**Usage**

```
salesCondition(tdata)
```

**Arguments**

tdata            an xts object containing the time series data, with one column named "COND" indicating the Sale Condition.

**Value**

xts object

**Author(s)**

Jonathan Cornelissen and Kris Boudt

---

sample_5minprices	<i>Ten artificial time series for the NYSE trading days during January 2010</i>
-------------------	---

---

### Description

Ten simulated price series for the 19 trading days in January 2010:

Ten hypothetical price series were simulated according to the factor diffusion process discussed in Barndorff-Nielsen et al. We assume that prices are only observed when a transaction takes place. The intensity of transactions follows a Poisson process and consequently, the inter transaction times are exponentially distributed. Therefore, we generated the inter transaction times of the price series by an independent exponential distributions with  $\lambda = 0.1$ , which we keep constant over time. This means we expect one transaction every ten seconds. In a final step, the time series were aggregated to the 5-minute frequency by previous tick aggregation.

### Usage

```
data("sample_5minprices")
```

### Format

xts object

### References

Barndorff-Nielsen, O. E., P. R. Hansen, A. Lunde and N. Shephard (2009). Multivariate realised kernels: consistent positive semi-definite estimators of the covariation of equity prices with noise and non-synchronous trading. *Journal of Econometrics*, forthcoming.

---

sample_5minprices_jumps	<i>Ten artificial time series (including jumps) for the NYSE trading days during January 2010</i>
-------------------------	---

---

### Description

Ten simulated price series for the 19 trading days in January 2010:

Ten hypothetical price series were simulated according to the factor diffusion process discussed in Barndorff-Nielsen et al. On top of this process we added a jump process, with jump occurrences governed by the Poisson process with 1 expected jump per day and jump magnitude modelled as in Boudt et al. (2008). We assume that prices are only observed when a transaction takes place. The intensity of transactions follows a Poisson process and consequently, the inter transaction times are exponentially distributed. Therefore, we generated the inter transaction times of the price series by an independent exponential distributions with  $\lambda = 0.1$ , which we keep constant over time. This means we expect one transaction every ten seconds. In a final step, the time series were aggregated to the 5-minute frequency by previous tick aggregation.

### Usage

```
data("sample_5minprices_jumps")
```

**Format**

xts object

**References**

Barndorff-Nielsen, O. E., P. R. Hansen, A. Lunde and N. Shephard (2009). Multivariate realised kernels: consistent positive semi-definite estimators of the covariation of equity prices with noise and non-synchronous trading. *Journal of Econometrics*, forthcoming.

Boudt, K., C. Croux, and S. Laurent (2008). Outlyingness weighted covariation. *Mimeo*.

---

sample_qdata	<i>Sample of cleaned quotes for stock XXX for 1 day</i>
--------------	---

---

**Description**

An xts object containing the raw quotes for the imaginary stock XXX for 1 day, in the typical NYSE TAQ database format. This is the cleaned version of the data sample [sample\\_qdataraw](#), using [quotesCleanup](#).

**Usage**

```
data("sample_qdata")
```

**Format**

xts object

---

sample_qdataraw	<i>Sample of raw quotes for stock XXX for 1 day</i>
-----------------	---

---

**Description**

An imaginary xts object containing the raw quotes for stock XXX for 1 day, in the typical NYSE TAQ database format.

**Usage**

```
data("sample_qdataraw")
```

**Format**

xts object

---

sample_real5minprices	<i>Sample of imaginary price data for 61 days</i>
-----------------------	---

---

**Description**

An xts object containing the 5-min aggregated imaginary price series for the trading days between 2005-03-04 and 2005-06-01.

**Usage**

```
data("sample_real5minprices")
```

**Format**

xts object

---

sample_tdata	<i>Sample of cleaned trades for stock XXX for 1 day</i>
--------------	---

---

**Description**

An xts object containing the trades for the imaginary stock XXX for 1 day, in the typical NYSE TAQ database format. This is the cleaned version of the data sample [sample\\_tdataraw](#), using [tradesCleanup](#).

**Usage**

```
data("sample_tdata")
```

**Format**

xts object

---

sample_tdataraw	<i>Sample of raw trades for stock XXX for 1 day</i>
-----------------	---

---

**Description**

An imaginary xts object containing the raw trades for stock XXX for 1 day, in the typical NYSE TAQ database format.

**Usage**

```
data("sample_tdataraw")
```

**Format**

xts object

sbux.xts

*Starbucks Data***Description**

Tick data for Starbucks 2011/07/01, cleaned with [tradesCleanup](#).

**Usage**

```
data(sbux.xts)
```

**Examples**

```
data(sbux.xts)
plot(sbux.xts)
```

selectExchange

*Retain only data from a single stock exchange***Description**

Function returns an xts object containing the data of only 1 stock exchange.

**Usage**

```
selectExchange(data,exch="N");
```

**Arguments**

- |      |  |
|------|--|
| data | an xts object containing the time series data. The object should have a column "EX", indicating the exchange by its symbol.  |
| exch | <p>The symbol of the stock exchange that should be selected. By default the NYSE is chosen (exch="N"). Other exchange symbols are:</p> <ul style="list-style-type: none"> <li>• A: AMEX</li> <li>• N: NYSE</li> <li>• B: Boston</li> <li>• P: Arca</li> <li>• C: NSX</li> <li>• T/Q: NASDAQ</li> <li>• D: NASD ADF and TRF</li> <li>• X: Philadelphia</li> <li>• I: ISE</li> <li>• M: Chicago</li> <li>• W: CBOE</li> <li>• Z: BATS</li> </ul> |

**Value**

xts object

**Author(s)**

Jonathan Cornelissen and Kris Boudt

---

spotVol	<i>Spot volatility estimation</i>
---------	-----------------------------------

---

**Description**

Function returns an estimate of the standard deviation  $\sigma_{t,i}$  of equispaced high-frequency returns  $r_{t,i}$  (read: the  $i$ th return on day  $t$ ). The underlying assumption is that, in the absence of price jumps, high-frequency returns are normally distributed with mean zero and standard deviation  $\sigma_{t,i}$ , where the standard deviation is the product between a deterministic periodic factor  $f_i$  (identical for every day in the sample) and a daily factor  $s_t$  (identical for all observations within a day).

For the estimation of  $s_t$  one can choose between the realized volatility, the bipower variation of Barndorff-Nielsen and Shephard (2004) or the MedRV of Andersen et al. (2009). The latter two have the advantage of being robust to price jumps.

The function takes as input the tick-by-tick price series. From these prices, equispaced returns are computed as the change in the log price of previous tick interpolated prices sampled every  $k$  minutes.

The estimation of  $f_i$  is either based on scale or regression estimators. The scale estimator can be the standard deviation or its jump robust version called the weighted standard deviation. For regression, choose OLS for the classical estimation and TML (truncated maximum likelihood) for jump robust regression. The regression specification consists either of one dummy for each intraday period (dummies=TRUE) or the flexible fourrier form with P1 cosinus and P2 sinus terms. For more details on the classical methods, see Taylor and Xu (1997) and Andersen et al. (1997). For the jump robust versions, see Boudt et al. (2010).

**Usage**

```
spotVol(pdata, dailyvol = "bipower", periodicvol = "TML",
        on = "minutes", k = 5, dummies = FALSE, P1 = 4, P2 = 2,
        marketopen = "09:30:00", marketclose = "16:00:00")
```

**Arguments**

pdata	xts object, containing the price series.
dailyvol	determines the estimation method for the component of intraday volatility that is constant over the day, but changes from day to day. Possible values are "bipower", "rv", "medrv".
periodicvol	determines the estimation method for the component of intraday volatility that depends in a deterministic way on the intraday time at which the return is observed. Possible values are "TML", "sd", "wsd", "OLS".
on	character, indicating the time scale in which "k" is expressed. Possible values are: "secs", "seconds", "mins", "minutes", "hours".

k	positive integer, indicating the number of periods to aggregate over. E.g. to aggregate a xts object to the 5 minute frequency set k=5 and on="minutes".
dummies	boolean, in case it is TRUE, the parametric estimator of periodic standard deviation specifies the periodicity function as the sum of dummy variables corresponding to each intraday period. If it false, the parametric estimator uses the Flexible Fourier specification. FALSE by default.
P1	is a positive integer valued parameter that corresponds to the number of cosinus terms used in the flexible fourrier specification for the periodicity function, see Andersen et al. (1997) for details.
P2	is a positive integer valued parameter that corresponds to the number of sinus terms used in the flexible fourrier specification for the periodicity function, see Andersen et al. (1997) for details.
marketopen	the market opening time, by default: marketopen = "09:30:00".
marketclose	the market closing time, by default: marketclose = "16:00:00".

### Details

Returns an xts object with first column equal to the high-frequency return series, second column is the estimated standard deviation, third column is the daily standard deviation factor and, finally, the fourth column is the periodic component.

### Author(s)

Jonathan Cornelissen and Kris Boudt

### References

- Andersen, T. G. and T. Bollerslev (1997). Intraday periodicity and volatility persistence in financial markets. *Journal of Empirical Finance* 4, 115-158.
- Andersen, T. G., D. Dobrev, and E. Schaumburg (2009). Jump-robust volatility estimation using nearest neighbor truncation. NBER Working Paper No. 15533.
- Barndorff-Nielsen, O. and N. Shephard (2004). Power and bipower variation with stochastic volatility and jumps. *Journal of Financial Econometrics* 2 (1), 1-37.
- Boudt K., Croux C. and Laurent S. (2011). Robust estimation of intraweek periodicity in volatility and jump detection. *Journal of Empirical Finance* 18, 353-367.
- Taylor, S. J. and X. Xu (1997). The incremental volatility information in one million foreign exchange quotations. *Journal of Empirical Finance* 4, 317-340.

### Examples

```
data("sample_real5minprices");

# Compute and plot intraday periodicity:
out = spotVol(sample_real5minprices,P1=6,P2=4,periodicvol="TML",k=5, dummies=FALSE);
head(out);
```



TAQLoad

*Load trade or quote data into R***Description**

Function to load the taq data into R. If only the trades (or quotes) should be loaded the function returns them directly as xts object. If both trades and quotes should be extracted, the function returns a list with two items named "trades" and "quotes" respectively.

The function assumes that the files are ordered in folders per day, and the folders contain: ticker\_trades.RData or ticker\_quotes.RData. In these files the xts object is stored. If you used the function [convert](#), this shall be the case. Please find more information in the pdf documentation.

**Usage**

```
TAQLoad(tickers,from,to,trades=TRUE,quotes=FALSE,
        datasource=NULL,variables=NULL)
```

**Arguments**

tickers	the ticker(s) to be loaded. It is recommended that you use only 1 ticker as input in case of non-synchronic observations. For synchronic data a vector of tickers can be used.
from	first day to load e.g. "2008-01-30".
to	last day to load e.g. "2008-01-30".
trades	boolean, determines whether trades are extracted.
quotes	boolean, determines whether quotes are extracted.
datasource	path to folder in which the files are contained.
variables	a character (or character vector) containing the name(s) of the variable(s) that should be loaded, e.g. c("SYMBOL","PRICE"). By default all data is loaded.

**Value**

see section description

**Author(s)**

Jonathan Cornelissen and Kris Boudt

**Examples**

```
#In order for these examples to work, the folder datasource
#should contain two folders named 2008-01-02 and 2008-01-03.
#These folder contain the files with the trade data,
#which are named "AAPL_trades.RData" or "AA_trades.RData".

from="2008-01-02";
to = "2008-01-03";
## Not run: datasource="C:\data";

#TAQLoad: load data for stock AAPL
```

```
## Not run: xx = TAQLoad(tickers="AAPL", from, to, trades=TRUE, quotes=FALSE,
datasource=datasource, variables=NULL)
## End(Not run)
## Not run: head(xx);

#Load only price data for stocks AA and AAPL
## Not run: xx = TAQLoad(tickers=c("AA", "AAPL"), from, to, trades=TRUE,
quotes=FALSE, datasource=datasource, variables="PRICE")
## End(Not run)
## Not run: head(xx);
## Not run: tail(xx);
```

---

tqLiquidity

---

*Calculate numerous (23) liquidity measures*


---

### Description

Function returns an xts object containing one of the following liquidity measures:

es, rs, value\_trade, signed\_value\_trade, di\_diff, di\_div, pes, prs, price\_impact, prop\_price\_impact, tsread, pts, p\_return\_sqr, p\_return\_abs, qs, pqs, logqs, logsize, qslope, logqslope, mq\_return\_sqr, mq\_return\_abs.

### Usage

```
tqLiquidity(tqdata, tdata, qdata, type, ...)
```

### Arguments

tqdata	xts object, containing joined trades and quotes (e.g. using <code>matchTradesQuotes</code> )
tdata	xts-object containing the trade data.
qdata	xts-object containing the quote data.
type	a character vector containing the name of one of the above mentioned liquidity measures, e.g. type="es" to get the effective spread.
...	additional arguments.

### Value

an xts object containing one of the above mentioned liquidity measures.

### Details

The respective liquidity measures are defined as follows:

- es: effective spread

$$\text{effective spread}_t = 2 * D_t * \left( \text{PRICE}_t - \frac{(\text{BID}_t + \text{OFR}_t)}{2} \right),$$

where  $D_t$  is 1 (-1) if  $\text{trade}_t$  was buy (sell) (see Boehmer (2005), Bessembinder (2003)). Note that the input of this function consists of the matched trades and quotes, so this is where the time indication refers to (and thus not to the registered quote timestamp).

- rs: realized spread

$$\text{realized spread}_t = 2 * D_t * \left( \text{PRICE}_t - \frac{(\text{BID}_{t+300} + \text{OFR}_{t+300})}{2} \right),$$

where  $D_t$  is 1 (-1) if  $\text{trade}_t$  was buy (sell) (see Boehmer (2005), Bessembinder (2003)). Note that the time indication of BID and OFR refers to the registered time of the quote in seconds.

- value\_trade: trade value

$$\text{trade value}_t = \text{SIZE}_t * \text{PRICE}_t.$$

- signed\_value\_trade: signed trade value

$$\text{signed trade value}_t = D_t * (\text{SIZE}_t * \text{PRICE}_t),$$

where  $D_t$  is 1 (-1) if  $\text{trade}_t$  was buy (sell) (see Boehmer (2005), Bessembinder (2003)).

- di\_diff: depth imbalance (as a difference)

$$\text{depth imbalance (as difference)}_t = \frac{D_t * (\text{OFRSIZ}_t - \text{BIDSIZ}_t)}{(\text{OFRSIZ}_t + \text{BIDSIZ}_t)},$$

where  $D_t$  is 1 (-1) if  $\text{trade}_t$  was buy (sell) (see Boehmer (2005), Bessembinder (2003)). Note that the input of this function consists of the matched trades and quotes, so this is where the time indication refers to (and thus not to the registered quote timestamp).

- di\_div: depth imbalance (as ratio)

$$\text{depth imbalance (as ratio)}_t = \left( \frac{D_t * \text{OFRSIZ}_t}{\text{BIDSIZ}_t} \right)^{D_t},$$

where  $D_t$  is 1 (-1) if  $\text{trade}_t$  was buy (sell) (see Boehmer (2005), Bessembinder (2003)). Note that the input of this function consists of the matched trades and quotes, so this is where the time indication refers to (and thus not to the registered quote timestamp).

- pes: proportional effective spread

$$\text{proportional effective spread}_t = \frac{\text{effective spread}_t}{(\text{OFR}_t + \text{BID}_t)/2}$$

(Venkataraman, 2001).

Note that the input of this function consists of the matched trades and quotes, so this is where the time indication refers to (and thus not to the registered quote timestamp).

- prs: proportional realized spread

$$\text{proportional realized spread}_t = \frac{\text{realized spread}_t}{(\text{OFR}_t + \text{BID}_t)/2}$$

(Venkataraman, 2001).

Note that the input of this function consists of the matched trades and quotes, so this is where the time indication refers to (and thus not to the registered

- price\_impact: price impact

$$\text{price impact}_t = \frac{\text{effective spread}_t - \text{realized spread}_t}{2}$$

(see Boehmer (2005), Bessembinder (2003)).

- prop\_price\_impact: proportional price impact

$$\text{proportional price impact}_t = \frac{(\text{effective spread}_t - \text{realized spread}_t)}{\frac{\text{OFR}_t + \text{BID}_t}{2}}$$

(Venkataraman, 2001). Note that the input of this function consists of the matched trades and quotes, so this is where the time indication refers to (and thus not to the registered quote timestamp).

- tspread: half traded spread

$$\text{half traded spread}_t = D_t * (\text{PRICE}_t - \frac{(\text{BID}_t + \text{OFR}_t)}{2}),$$

where  $D_t$  is 1 (-1) if  $\text{trade}_t$  was buy (sell) (see Boehmer (2005), Bessembinder (2003)). Note that the input of this function consists of the matched trades and quotes, so this is where the time indication refers to (and thus not to the registered quote timestamp).

- pts: proportional half traded spread

$$\text{proportional half traded spread}_t = \frac{\text{half traded spread}_t}{\frac{\text{OFR}_t + \text{BID}_t}{2}}.$$

Note that the input of this function consists of the matched trades and quotes, so this is where the time indication refers to (and thus not to the registered quote timestamp).

- p\_return\_sqr: squared log return on trade prices

$$\text{squared log return on Trade prices}_t = (\log(\text{PRICE}_t) - \log(\text{PRICE}_{t-1}))^2.$$

- p\_return\_abs: absolute log return on trade prices

$$\text{absolute log return on Trade prices}_t = |\log(\text{PRICE}_t) - \log(\text{PRICE}_{t-1})|.$$

- qs: quoted spread

$$\text{quoted spread}_t = \text{OFR}_t - \text{BID}_t$$

Note that the input of this function consists of the matched trades and quotes, so this is where the time indication refers to (and thus not to the registered quote timestamp).

- pqs: proportional quoted spread

$$\text{proportional quoted spread}_t = \frac{\text{quoted spread}_t}{\frac{\text{OFR}_t + \text{BID}_t}{2}}$$

(Venkataraman, 2001). Note that the input of this function consists of the matched trades and quotes, so this is where the time indication refers to (and thus not to the registered quote timestamp).

- logqs: log quoted spread

$$\log \text{ quoted spread}_t = \log\left(\frac{\text{OFR}_t}{\text{BID}_t}\right)$$

(Hasbrouck and Seppi, 2001). Note that the input of this function consists of the matched trades and quotes, so this is where the time indication refers to (and thus not to the registered quote timestamp).

- logsize: log quoted size

$$\log \text{ quoted size}_t = \log(\text{OFRSIZ}_t) - \log(\text{BIDSIZ}_t)$$

(Hasbrouck and Seppi, 2001). Note that the input of this function consists of the matched trades and quotes, so this is where the time indication refers to (and thus not to the registered quote timestamp).

- qslope: quoted slope

$$\text{quoted slope}_t = \frac{\text{quoted spread}_t}{\log \text{ quoted size}_t}$$

(Hasbrouck and Seppi, 2001).

- logqslope: log quoted slope

$$\log \text{ quoted slope}_t = \frac{\log \text{ quoted spread}_t}{\log \text{ quoted size}_t}.$$

- mq\_return\_sqr: midquote squared return

$$\text{midquote squared return}_t = (\log(\text{midquote}_t) - \log(\text{midquote}_{t-1}))^2,$$

where  $\text{midquote}_t = \frac{\text{BID}_t + \text{OFR}_t}{2}$ .

- mq\_return\_abs: midquote absolute return

$$\text{midquote absolute return}_t = |\log(\text{midquote}_t) - \log(\text{midquote}_{t-1})|,$$

where  $\text{midquote}_t = \frac{\text{BID}_t + \text{OFR}_t}{2}$ .

- signed\_trade\_size: signed trade size

$$\text{signed trade size}_t = D_t * \text{SIZE}_t,$$

where  $D_t$  is 1 (-1) if  $\text{trade}_t$  was buy (sell).

## Author(s)

Jonathan Cornelissen and Kris Boudt

## References

- Bessembinder, H. (2003). Issues in assessing trade execution costs. *Journal of Financial Markets*, 223-257.
- Boehmer, E. (2005). Dimensions of execution quality: Recent evidence for US equity markets. *Journal of Financial Economics* 78 (3), 553-582.
- Hasbrouck, J. and D. J. Seppi (2001). Common factors in prices, order flows and liquidity. *Journal of Financial Economics*, 383-411.
- Venkataraman, K. (2001). Automated versus floor trading: An analysis of execution costs on the paris and new york exchanges. *The Journal of Finance*, 56, 1445-1485.

## Examples

```
#load data samples
data("sample_tdata");
data("sample_qdata");
tdata = sample_tdata;
qdata = sample_qdata;
#match the trade and quote data
tqdata = matchTradesQuotes(tdata,qdata);

#calculate the proportional realized spread:
prs = tqLiquidity(tqdata,tdata,qdata,type="prs");

#calculate the effective spread:
es = tqLiquidity(tqdata,type="es");
```

---

tradesCleanup	<i>Cleans trade data</i>
---------------	--------------------------

---

## Description

This is a wrapper function for cleaning the trade data of all stocks in "ticker" over the interval [from,to]. The result is saved in the folder datadestination. The function returns a vector indicating how many trades were removed at each cleaning step.

In case you supply the argument "rawtdata", the on-disk functionality is ignored and the function returns a list with the cleaned trades as xts object (see examples).

The following cleaning functions are performed sequentially: [noZeroPrices](#), [selectExchange](#), [salesCondition](#), [mergeTradesSameTimestamp](#).

Since the function [rmTradeOutliers](#) also requires cleaned quote data as input, it is not incorporated here and there is a separate wrapper called [tradesCleanupFinal](#).

## Usage

```
tradesCleanup(from,to,datasource,datadestination,ticker,exchanges,
              tdataraw,report,selection,...)
```

## Arguments

from	character indicating first date to clean, e.g. "2008-01-30".
to	character indicating last date to clean, e.g. "2008-01-31".
datasource	character indicating the folder in which the original data is stored.
datadestination	character indicating the folder in which the cleaned data is stored.
ticker	vector of tickers for which the data should be cleaned, e.g. ticker = c("AAPL","AIG")
exchanges	vector of stock exchange symbols for all tickers in vector "ticker". It thus should have the same length as the vector ticker. Only data from one exchanges will be retained for each stock respectively, e.g. exchanges = c("Q","N") The possible exchange symbols are: <ul style="list-style-type: none"> <li>• A: AMEX</li> <li>• N: NYSE</li> </ul>

	<ul style="list-style-type: none"> <li>• B: Boston</li> <li>• P: Arca</li> <li>• C: NSX</li> <li>• T/Q: NASDAQ</li> <li>• D: NASD ADF and TRF</li> <li>• X: Philadelphia</li> <li>• I: ISE</li> <li>• M: Chicago</li> <li>• W: CBOE</li> <li>• Z: BATS</li> </ul>
tdataraw	xts object containing (ONE day and for ONE stock only) raw trade data. This argument is NULL by default. Enabling it means the arguments from, to, data-source and datadestination will be ignored. (only advisable for small chunks of data)
report	boolean and TRUE by default. In case it is true the function returns (also) a vector indicating how many trades remained after each cleaning step.
selection	argument to be passed on to the cleaning routine <a href="#">mergeTradesSameTimestamp</a> . The default is "median".
...	additional arguments.

### Value

For each day an xts object is saved into the folder of that date, containing the cleaned data. This procedure is performed for each stock in "ticker". The function returns a vector indicating how many trades remained after each cleaning step.

In case you supply the argument "rawtdata", the on-disk functionality is ignored and the function returns a list with the cleaned trades as xts object (see examples).

### Author(s)

Jonathan Cornelissen and Kris Boudt

### References

Barndorff-Nielsen, O. E., P. R. Hansen, A. Lunde, and N. Shephard (2009). Realized kernels in practice: Trades and quotes. *Econometrics Journal* 12, C1-C32.

Brownlees, C.T. and Gallo, G.M. (2006). Financial econometric analysis at ultra-high frequency: Data handling concerns. *Computational Statistics & Data Analysis*, 51, pages 2232-2245.

Falkenberry, T.N. (2002). High frequency data filtering. Unpublished technical report.

### Examples

```
#Consider you have raw trade data for 1 stock for 1 day
data("sample_tdataraw");
head(sample_tdataraw);
dim(sample_tdataraw);
tdata_afterfirstcleaning = tradesCleanup(tdataraw=sample_tdataraw,exchanges="N");
tdata_afterfirstcleaning$report;
barplot(tdata_afterfirstcleaning$report);
dim(tdata_afterfirstcleaning$tdata);
```

```
#In case you have more data it is advised to use the on-disk functionality
#via "from","to","datasource",etc. arguments
```

---

tradesCleanupFinal	<i>Perform a final cleaning procedure on trade data</i>
--------------------	---

---

## Description

Function performs cleaning procedure [rmTradeOutliers](#) for the trades of all stocks in "ticker" over the interval [from,to] and saves the result in "datadestination". Note that preferably the input data for this function is trade and quote data cleaned by respectively e.g. [tradesCleanup](#) and [quotesCleanup](#).

## Usage

```
tradesCleanupFinal(from,to,datasource,datadestination,ticker,
                   tdata,qdata,...)
```

## Arguments

from	character indicating first date to clean, e.g. "2008-01-30".
to	character indicating last date to clean, e.g. "2008-01-31".
datasource	character indicating the folder in which the original data is stored.
datadestination	character indicating the folder in which the cleaned data is stored.
ticker	vector of tickers for which the data should be cleaned.
tdata	xts object containing (ONE day and for ONE stock only) trade data cleaned by <a href="#">tradesCleanup</a> . This argument is NULL by default. Enabling it, means the arguments from, to, datasource and datadestination will be ignored. (only advisable for small chunks of data)
qdata	xts object containing (ONE day and for ONE stock only) cleaned quote data. This argument is NULL by default. Enabling it means the arguments from, to, datasource, datadestination will be ignored. (only advisable for small chunks of data)
...	additional arguments.

## Value

For each day an xts object is saved into the folder of that date, containing the cleaned data. This procedure is performed for each stock in "ticker".

In case you supply the arguments "tdata" and "qdata", the on-disk functionality is ignored and the function returns a list with the cleaned trades as xts object (see examples).

## Author(s)

Jonathan Cornelissen and Kris Boudt



## References

Barndorff-Nielsen, O. E., P. R. Hansen, A. Lunde, and N. Shephard (2009). Realized kernels in practice: Trades and quotes. *Econometrics Journal* 12, C1-C32.

Brownlees, C.T. and Gallo, G.M. (2006). Financial econometric analysis at ultra-high frequency: Data handling concerns. *Computational Statistics & Data Analysis*, 51, pages 2232-2245.

## Examples

```
#Consider you have raw trade data for 1 stock for 1 day
#data("sample_qdata"); #load cleaned quote data
#data("sample_tdataraw"); #load raw trade data
#tdata_afterfirstcleaning = tradesCleanup(tdataraw=sample_tdataraw,
#exchange="N",report=FALSE);
#dim(tdata_afterfirstcleaning);
#tdata_afterfinalcleaning = tradesCleanupFinal(qdata=sample_qdata,
#tdata=tdata_afterfirstcleaning);
#dim(tdata_afterfinalcleaning);
#In case you have more data it is advised to use the on-disk functionality
#via "from","to","datasource",etc. arguments
```

# Index

## \*Topic **cleaning**

- autoSelectExchangeQuotes, 8
- autoSelectExchangeTrades, 9
- exchangeHoursOnly, 12
- mergeQuotesSameTimestamp, 21
- mergeTradesSameTimestamp, 22
- noZeroPrices, 23
- noZeroQuotes, 24
- quotesCleanup, 24
- rmLargeSpread, 38
- rmNegativeSpread, 38
- rmOutliers, 39
- rmTradeOutliers, 40
- salesCondition, 50
- selectExchange, 54
- tradesCleanup, 62
- tradesCleanupFinal, 64

## \*Topic **data manipulation**

- aggregatePrice, 3
- aggregateQuotes, 4
- aggregateTrades, 5
- aggregatets, 6
- makePsd, 17
- matchTradesQuotes, 19
- refreshTime, 32
- TAQLoad, 57

## \*Topic **datasets**

- lltc.xts, 17
- realized\_library, 32
- sample\_5minprices, 51
- sample\_5minprices\_jumps, 51
- sample\_qdata, 52
- sample\_qdataraw, 52
- sample\_real5minprices, 53
- sample\_tdata, 53
- sample\_tdataraw, 53
- sbux.xts, 54

## \*Topic **forecasting**

- harModel, 13
- heavyModel, 15

## \*Topic **liquidity**

- getTradeDirection, 12
- tqLiquidity, 58

## \*Topic **methods**

- rAccumulation, 26
- rCumSum, 31
- rMarginal, 36
- rScatterReturns, 45
- rZero, 49

## \*Topic **package**

- highfrequency-package, 3

## \*Topic **volatility**

- medRV, 20
- minRV, 22
- rAVGCov, 27
- rBPCov, 29
- rCov, 30
- rHYCov, 33
- rKernel.available, 34
- rKernelCov, 35
- rOWCov, 41
- rRTSCov, 42
- rThresholdCov, 46
- rTSCov, 47
- spotVol, 55

- aggregatePrice, 3
- aggregateQuotes, 4, 6
- aggregateTrades, 5, 6
- aggregatets, 6, 24
- autoSelectExchangeQuotes, 8
- autoSelectExchangeTrades, 9

- convert, 10, 57

- exchangeHoursOnly, 12

- getTradeDirection, 12

- harModel, 13
- heavyModel, 15
- highfrequency (highfrequency-package), 3
- highfrequency-package, 3

- lltc.xts, 17
- lm, 14

- makePsd, 17, 44, 48

makeReturns, 18  
matchTradesQuotes, 12, 19, 58  
medRV, 20  
mergeQuotesSameTimestamp, 21, 24, 25  
mergeTradesSameTimestamp, 22, 62, 63  
minRV, 22  
  
noZeroPrices, 23, 62  
noZeroQuotes, 24, 24  
  
previoustick, 24  
  
quotesCleanup, 24, 52, 64  
  
rAccumulation, 26, 37  
rAVGCov, 27  
rBPCov, 29, 42  
rCov, 30  
rCumSum, 31  
realized\_library, 32  
refreshTime, 32, 44, 48  
rHYCov, 33  
rKernel.available, 34, 36  
rKernelCov, 35, 35  
rMarginal, 27, 36  
rmLargeSpread, 24, 25, 38  
rmNegativeSpread, 38  
rmOutliers, 24, 25, 39  
rmTradeOutliers, 40, 62, 64  
rOWCov, 41, 42, 46  
rRTSCov, 42  
rScatterReturns, 45  
rThresholdCov, 42, 46  
rTSCov, 47  
rZero, 49  
  
salesCondition, 50, 62  
sample\_5minprices, 51  
sample\_5minprices\_jumps, 51  
sample\_qdata, 4, 52  
sample\_qdataraw, 52, 52  
sample\_real5minprices, 53  
sample\_tdata, 5, 53  
sample\_tdataraw, 53, 53  
sbux.xts, 54  
selectExchange, 24, 54, 62  
spotVol, 55  
  
TAQLoad, 10, 57  
tqLiquidity, 58  
tradesCleanup, 17, 53, 54, 62, 64  
tradesCleanupFinal, 62, 64