

Package ‘HighFreq’

June 2, 2022

Type Package

Title High Frequency Time Series Management

Version 0.1

Date 2018-09-12

Author Jerzy Pawlowski (algoquant)

Maintainer Jerzy Pawlowski <jp3900@nyu.edu>

Description Functions for chaining and joining time series, scrubbing bad data, managing time zones and aligning time indices, converting TAQ data to OHLC format, aggregating data to lower frequency, estimating volatility, skew, and higher moments.

License MPL-2.0

Depends xts,
quantmod,
rutils

Imports xts,
quantmod,
rutils,
RcppRoll,
Rcpp

LinkingTo Rcpp, RcppArmadillo

SystemRequirements GNU make, C++11

Remotes github::algoquant/rutils,

VignetteBuilder knitr

LazyData true

ByteCompile true

Repository GitHub

URL <https://github.com/algoquant/HighFreq>

RoxygenNote 7.1.2

Encoding UTF-8

R topics documented:

agg_ohlc	3
agg_stats_r	4
back_test	5
calc_covar	7
calc_cvar	9
calc_eigen	10
calc_endpoints	11
calc_hurst	12
calc_hurst_ohlc	13
calc_inv	15
calc_kurtosis	16
calc_lm	18
calc_mean	19
calc_ranks	20
calc_reg	21
calc_scaled	23
calc_skew	24
calc_startpoints	26
calc_var	27
calc_varvec	28
calc_var_ag	29
calc_var_ohlc	30
calc_var_ohlc_ag	32
calc_var_ohlc_r	34
calc_weights	35
diffit	38
diff_vec	39
hf_data	40
lagit	41
lag_vec	42
lik_garch	43
mult_mat	44
mult_mat_ref	45
ohlc_returns	47
ohlc_sharpe	48
ohlc_skew	49
ohlc_variance	49
random_ohlc	51
random_taq	52
remove_jumps	53
roll_apply	54
roll_backtest	55
roll_conv	57
roll_count	58
roll_fun	59
roll_hurst	61
roll_kurtosis	62
roll_mean	63
roll_ohlc	65
roll_reg	66

roll_scale	68
roll_sharpe	69
roll_skew	70
roll_stats	71
roll_sum	72
roll_sumep	73
roll_var	75
roll_varvec	76
roll_var_ohlc	77
roll_vec	80
roll_vecw	81
roll_vwap	82
roll_wsum	83
roll_zscores	85
run_covar	86
run_max	88
run_mean	89
run_min	91
run_reg	92
run_var	94
run_var_ohlc	95
run_zscores	96
save_rets	98
save_rets_ohlc	99
save_scrub_agg	100
save_taq	101
scrub_agg	101
scrub_taq	103
seasonality	103
sim_ar	104
sim_df	105
sim_garch	106
sim_ou	108
sim_schwartz	109
which_extreme	110
which_jumps	111

agg_ohlc

*Aggregate a time series of data into a single bar of OHLC data.***Description**

Aggregate a time series of data into a single bar of *OHLC* data.

Usage

```
agg_ohlc(tseries)
```

Arguments

tseries *A time series or a matrix with multiple columns of data.*

Details

The function `agg_ohlc()` aggregates a time series of data into a single bar of *OHLC* data. It can accept either a single column of data or four columns of *OHLC* data. It can also accept an additional column containing the trading volume.

The function `agg_ohlc()` calculates the *open* value as equal to the *open* value of the first row of *tseries*. The *high* value as the maximum of the *high* column of *tseries*. The *low* value as the minimum of the *low* column of *tseries*. The *close* value as the *close* of the last row of *tseries*. The *volume* value as the sum of the *volume* column of *tseries*.

For a single column of data, the *open*, *high*, *low*, and *close* values are all the same.

Value

A *matrix* containing a single row, with the *open*, *high*, *low*, and *close* values, and also the total *volume* (if provided as either the second or fifth column of *tseries*).

Examples

```
## Not run:
# Define matrix of OHLC data
ohlc <- coredata(rutils::etfenv$VTI[, 1:5])
# Aggregate to single row matrix
ohlcagg <- HighFreq::agg_ohlc(ohlc)
# Compare with calculation in R
all.equal(drop(ohlcagg),
  c(ohlc[1, 1], max(ohlc[, 2]), min(ohlc[, 3]), ohlc[NROW(ohlc), 4], sum(ohlc[, 5])),
  check.attributes=FALSE)

## End(Not run)
```

agg_stats_r	<i>Calculate the aggregation (weighted average) of a statistical estimator over a OHLC time series using R code.</i>
-------------	--

Description

Calculate the aggregation (weighted average) of a statistical estimator over a *OHLC* time series using R code.

Usage

```
agg_stats_r(ohlc, calcBars = "ohlc_variance", weighted = TRUE, ...)
```

Arguments

...	additional parameters to the function <code>calcBars</code> .
ohlc	An <i>OHLC</i> time series of prices and trading volumes, in <i>xts</i> format.
calcBars	A <i>character</i> string representing a function for calculating statistics for individual <i>OHLC</i> bars.
weighted	<i>Boolean</i> argument: should estimate be weighted by the trading volume? (default is TRUE)

Details

The function `agg_stats_r()` calculates a single number representing the volume weighted average of statistics of individual *OHLC* bars. It first calls the function `calcBars` to calculate a vector of statistics for the *OHLC* bars. For example, the statistic may simply be the difference between the *High* minus *Low* prices. In this case the function `calcBars` would calculate a vector of *High* minus *Low* prices. The function `agg_stats_r()` then calculates a trade volume weighted average of the vector of statistics.

The function `agg_stats_r()` is implemented in R code.

Value

A single *numeric* value equal to the volume weighted average of an estimator over the time series.

Examples

```
# Calculate weighted average variance for SPY (single number)
variance <- agg_stats_r(ohlc=HighFreq::SPY, calcBars="ohlc_variance")
# Calculate time series of daily skew estimates for SPY
skew_daily <- apply.daily(x=HighFreq::SPY, FUN=agg_stats_r, calcBars="ohlc_skew")
```

back_test	<i>Simulate (backtest) a rolling portfolio optimization strategy, using RcppArmadillo.</i>
-----------	--

Description

Simulate (backtest) a rolling portfolio optimization strategy, using RcppArmadillo.

Usage

```
back_test(
  excess,
  returns,
  startp,
  endp,
  lambda = 0,
  method = "sharpen",
  eigen_thresh = 1e-05,
  eigen_max = 0L,
  confl = 0.1,
  alpha = 0,
  rankw = FALSE,
  centerw = FALSE,
  scalew = "voltarget",
  vol_target = 0.01,
  coeff = 1,
  bid_offer = 0
)
```

Arguments

returns	A <i>time series</i> or a <i>matrix</i> of asset returns data.
excess	A <i>time series</i> or a <i>matrix</i> of excess returns data (the returns in excess of the risk-free rate).
startp	An <i>integer vector</i> of start points.
endp	An <i>integer vector</i> of end points.
lambda	A <i>numeric</i> decay factor to multiply the past portfolio weights. (The default is $\lambda = 0$ - no memory.)
coeff	A <i>numeric</i> multiplier of the weights. (The default is 1)
bid_offer	A <i>numeric</i> bid-offer spread (the default is 0)
method	A <i>string</i> specifying the method for calculating the weights (see Details) (the default is <code>method = "sharpen"</code>)
eigen_thresh	A <i>numeric</i> threshold level for discarding small singular values in order to regularize the inverse of the covariance matrix (the default is $1e-5$).
eigen_max	An <i>integer</i> equal to the number of singular values used for calculating the regularized inverse of the returns matrix (the default is 0 - equivalent to <code>eigen_max</code> equal to the number of columns of returns).
conf1	The confidence level for calculating the quantiles of returns (the default is <code>conf1 = 0.75</code>).
alpha	The shrinkage intensity between 0 and 1. (the default is 0).
rankw	A <i>Boolean</i> specifying whether the weights should be ranked (the default is <code>rankw = FALSE</code>).
centerw	A <i>Boolean</i> specifying whether the weights should be centered (the default is <code>centerw = FALSE</code>).
scalew	A <i>string</i> specifying the method for scaling the weights (the default is <code>scalew = "voltarget"</code>).
vol_target	A <i>numeric</i> volatility target for scaling the weights (the default is $1e-5$)

Details

The function `back_test()` performs a backtest simulation of a rolling portfolio optimization strategy over a *vector* of the end points `endp`.

It performs a loop over the end points `endp`, and subsets the *matrix* of the excess asset returns `excess` along its rows, between the corresponding *start point* and the *end point*. It passes the subset matrix of excess returns into the function `calc_weights()`, which calculates the optimal portfolio weights at each *end point*. The arguments `eigen_max`, `alpha`, `method`, `rankw`, `centerw`, and `scalew` are also passed to the function `calc_weights()`.

It then recursively averages the weights w_i at the *end point* $= i$ with the weights w_{i-1} from the previous *end point* $= (i-1)$, using the decay factor $\lambda = \lambda$:

$$w_i = (1 - \lambda)w_i + \lambda w_{i-1}$$

The purpose of averaging the weights is to reduce their variance, to improve their out-of-sample performance. It is equivalent to extending the portfolio holding period beyond the time interval between neighboring *end points*.

The function `back_test()` then calculates the out-of-sample strategy returns by multiplying the average weights times the future asset returns.

The function `back_test()` multiplies the out-of-sample strategy returns by the coefficient `coeff` (with default equal to 1), which allows simulating either a trending strategy (if `coeff = 1`), or a reverting strategy (if `coeff = -1`).

The function `back_test()` calculates the transaction costs by multiplying the bid-offer spread `bid_offer` times the absolute difference between the current weights minus the weights from the previous period. Then it subtracts the transaction costs from the out-of-sample strategy returns.

The function `back_test()` returns a *time series* (column vector) of strategy returns, of the same length as the number of rows of returns.

Value

A column *vector* of strategy returns, with the same length as the number of rows of returns.

Examples

```
## Not run:
# Calculate the ETF daily excess returns
returns <- na.omit(rutils::etfenv$returns[, 1:16])
# riskf is the daily risk-free rate
riskf <- 0.03/260
excess <- returns - riskf
# Define monthly end points without initial warmup period
endp <- rutils::calc_endpoints(returns, interval="months")
endp <- endp[endp > 0]
nrows <- NROW(endp)
# Define 12-month look-back interval and start points over sliding window
look_back <- 12
startp <- c(rep_len(1, look_back-1), endp[1:(nrows-look_back+1)])
# Define return shrinkage and regularization intensities
alpha <- 0.5
eigen_max <- 3
# Simulate a monthly rolling portfolio optimization strategy
pnls <- HighFreq::back_test(excess, returns,
                           startp-1, endp-1,
                           eigen_max = eigen_max,
                           alpha = alpha)
pnls <- xts::xts(pnls, index(returns))
colnames(pnls) <- "strat_rets"
# Plot dygraph of strategy
dygraphs::dygraph(cumsum(pnls),
  main="Cumulative Returns of Max Sharpe Portfolio Strategy")

## End(Not run)
```

calc_covar

Calculate the covariance matrix of the columns of a time series using RcppArmadillo.

Description

Calculate the covariance matrix of the columns of a *time series* using RcppArmadillo.

Usage

```
calc_covar(tseries, method = "moment", conf1 = 0.75)
```

Arguments

tseries	A <i>time series</i> or a <i>matrix</i> of data.
method	A <i>string</i> specifying the type of the covariance model (the default is method = "moment" - see Details).
conf1	The confidence level for calculating the quantiles of returns (the default is conf1 = 0.75).

Details

The function `calc_covar()` calculates the covariance matrix of the columns of a *time series* or a *matrix* of data using RcppArmadillo C++ code. The covariance is a measure of the codependency of the data.

If method = "moment" (the default) then `calc_covar()` calculates the covariance as the second co-moment:

$$\sigma_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

Then `calc_covar()` performs the same calculation as the R function `stats::cov()`.

If method = "quantile" then it calculates the covariance as the difference between the quantiles as follows:

$$\mu = q_{\alpha} - q_{1-\alpha}$$

Where α is the confidence level for calculating the quantiles.

If method = "nonparametric" then it calculates the covariance as the Median Absolute Deviation (*MAD*):

$$MAD = \text{median}(\text{abs}(x - \text{median}(x)))$$

It also multiplies the *MAD* by a factor of 1.4826, to make it comparable to the standard deviation.

If method = "nonparametric" then `calc_covar()` performs the same calculation as the function `stats::mad()`, but it's much faster because it uses RcppArmadillo C++ code.

If the number of rows of `tseries` is less than 3 then it returns zeros.

Value

A square matrix with the covariance coefficients of the columns of the *time series* `tseries`.

Examples

```
## Not run:
# Calculate VTI and XLF returns
returns <- na.omit(rutils::etfenv$returns[, c("VTI", "XLF")])
# Compare HighFreq::calc_covar() with standard var()
all.equal(drop(HighFreq::calc_covar(returns)),
  cov(returns), check.attributes=FALSE)
# Compare the speed of RcppArmadillo with matrixStats and with R code
library(microbenchmark)
summary(microbenchmark(
  Rcpp=HighFreq::calc_covar(returns),
  Rcode=cov(returns),
```



```

times=10))[, c(1, 4, 5)] # end microbenchmark summary
# Compare HighFreq::calc_cvar() with stats::mad()
all.equal(drop(HighFreq::calc_cvar(returns, method="nonparametric")),
  sapply(returns, mad), check.attributes=FALSE)
# Compare the speed of RcppArmadillo with stats::mad()
summary(microbenchmark(
  Rcpp=HighFreq::calc_cvar(returns, method="nonparametric"),
  Rcode=sapply(returns, mad),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary

## End(Not run)

```

calc_cvar	<i>Calculate the Value at Risk (VaR) or the Conditional Value at Risk (CVaR) of an xts time series of returns, using R code.</i>
-----------	--

Description

Calculate the Value at Risk (*VaR*) or the Conditional Value at Risk (*CVaR*) of an *xts time series* of returns, using R code.

Usage

```
calc_cvar(tseries, method = "var", confi = pnorm(-2))
```

Arguments

tseries	An <i>xts time series</i> of returns with multiple columns.
method	A <i>string</i> specifying the type of risk measure (the default is method = "var" - see Details).
confi	The confidence level for calculating the quantile (the default is confi = pnorm(-2) = 0.02275).

Details

The function `calc_cvar()` calculates the Value at Risk (*VaR*) or the Conditional Value at Risk (*CVaR*) of an *xts time series* of returns, using R

The Value at Risk (*VaR*) and the Conditional Value at Risk (*CVaR*) are measures of the tail risk of returns.

If method = "var" then `calc_cvar()` calculates the Value at Risk (*VaR*) as the quantile of the returns as follows:

$$\alpha = \int_{-\infty}^{\text{VaR}(\alpha)} f(r) \, dr$$

Where α is the confidence level for calculating the quantile, and $f(r)$ is the probability density (distribution) of returns.

If method = "cvar" then `calc_cvar()` calculates the Value at Risk (*VaR*) as the Expected Tail Loss (*ETL*) of the returns as follows:

$$\text{CVaR} = \frac{1}{\alpha} \int_0^\alpha \text{VaR}(p) \, dp$$

Where α is the confidence level for calculating the quantile.

Value

A vector with the risk measures of the columns of the input *time series* tseries.

Examples

```
## Not run:
# Calculate VTI and XLF returns
returns <- na.omit(rutils::etfenv$returns[, c("VTI", "XLF")])
# Calculate VaR
all.equal(HighFreq::calc_cvar(returns),
  sapply(returns, quantile, probs=pnorm(-2)), check.attributes=FALSE)
# Calculate CVaR
all.equal(HighFreq::calc_cvar(returns, method="cvar", confi=0.02),
  sapply(returns, function(x) mean(x[x < quantile(x, 0.02)])),
  check.attributes=FALSE)

## End(Not run)
```

calc_eigen

Calculate the eigen decomposition of the covariance matrix of returns data using RcppArmadillo.

Description

Calculate the eigen decomposition of the *covariance matrix* of returns data using RcppArmadillo.

Usage

```
calc_eigen(tseries)
```

Arguments

tseries A *time series* or *matrix* of returns data.

Details

The function calc_eigen() first calculates the *covariance matrix* of tseries, and then calculates the eigen decomposition of the *covariance matrix*.

Value

A list with two elements: a *vector* of eigenvalues (named "values"), and a *matrix* of eigenvectors (named "vectors").

Examples

```
## Not run:
# Create matrix of random data
datav <- matrix(rnorm(5e6), nc=5)
# Calculate eigen decomposition
eigend <- HighFreq::calc_eigen(scale(datav, scale=FALSE))
# Calculate PCA
```

```

pcad <- prcomp(datav)
# Compare PCA with eigen decomposition
all.equal(pcad$sdev^2, drop(eigend$values))
all.equal(abs(unname(pcad$rotation)), abs(eigend$vectors))
# Compare the speed of Rcpp with R code
summary(microbenchmark(
  Rcpp=HighFreq::calc_eigen(datav),
  Rcode=prcomp(datav),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary

## End(Not run)

```

calc_endpoints	<i>Calculate a vector of end points that divides a vector into equal intervals.</i>
----------------	---

Description

Calculate a vector of end points that divides a vector into equal intervals.

Usage

```
calc_endpoints(length, step = 1L, stub = 0L)
```

Arguments

length	An <i>integer</i> equal to the length of the vector to be divided into equal intervals.
step	The number of elements in each interval between neighboring end points.
stub	An <i>integer</i> value equal to the first end point for calculating the end points.

Details

The end points are a vector of integers which divide a vector of length equal to length into equally spaced intervals. If a whole number of intervals doesn't fit over the vector, then calc_endpoints() adds a stub interval at the end.

The first end point is equal to the argument step, unless the argument stub is provided, and then it becomes the first end point.

For example, consider the end points for a vector of length 20 divided into intervals of length step=5: 0,5,10,15,20. In order for all the differences between neighboring end points to be equal to 5, the first end point is set equal to 0. But 0 doesn't correspond to any vector element, so calc_endpoints() doesn't include it and it only retains the non-zero end points equal to: 5,10,15,20.

Since indexing in C++ code starts at 0, then calc_endpoints() shifts the end points by -1 and returns the vector equal to 4,9,14,19.

If stub = 1 then the first end point is equal to 1 and the end points are equal to: 1,6,11,16,20. The extra stub interval at the end is equal to 4 = 20 -16. And calc_endpoints() returns 0,5,10,15,19. The first value is equal to 0 which is the index of the first element in C++ code.

If stub = 2 then the first end point is equal to 2, with an extra stub interval at the end, and the end points are equal to: 2,7,12,17,20. And calc_endpoints() returns 1,6,11,16,19.

The function `calc_endpoints()` is similar to the function `rutils::calc_endpoints()` from package **rutils**.

But the end points are shifted by -1 compared to R code because indexing starts at 0 in C++ code, while it starts at 1 in R code. So if `calc_endpoints()` is used in R code then 1 should be added to it.

This works in R code because the vector element corresponding to index 0 is empty. For example, the R code: `(4:1)[c(0,1)]` produces `4`. So in R we can select vector elements using the end points starting at zero.

In C++ the end points must be shifted by -1 compared to R code, because indexing starts at 0 : $-1, 4, 9, 14, 19$. But there is no vector element corresponding to index -1 . So in C++ we cannot select vector elements using the end points starting at -1 . The solution is to drop the first placeholder end point.

Value

A vector of equally spaced *integers* representing the end points.

Examples

```
# Calculate end points without a stub interval
HighFreq::calc_endpoints(length=20, step=5)
# Calculate end points with a final stub interval
HighFreq::calc_endpoints(length=23, step=5)
# Calculate end points with initial and final stub intervals
HighFreq::calc_endpoints(length=20, step=5, stub=2)
# Calculate end points with initial and final stub intervals
HighFreq::calc_endpoints(length=20, step=5, stub=24)
```

calc_hurst

Calculate the Hurst exponent from the volatility ratio of aggregated returns.

Description

Calculate the Hurst exponent from the volatility ratio of aggregated returns.

Usage

```
calc_hurst(tseries, step = 1L)
```

Arguments

tseries	A <i>time series</i> or a <i>matrix</i> of prices.
step	The number of periods in each interval between neighboring end points.

Details

The function `calc_hurst()` calculates the Hurst exponent from the ratios of the volatilities of aggregated returns.

The aggregated volatility σ_t scales (increases) with the length of the aggregation interval Δt raised to the power of the *Hurst exponent* H :

$$\sigma_t = \sigma \Delta t^H$$

Where σ is the daily return volatility.

The *Hurst exponent* H is equal to the logarithm of the ratio of the volatilities divided by the logarithm of the time interval Δt :

$$H = \frac{\log \sigma_t - \log \sigma}{\log \Delta t}$$

The function `calc_hurst()` calls the function `calc_var_ag()` to calculate the aggregated volatility σ_t .

Value

The Hurst exponent calculated from the variance of aggregated returns.

Examples

```
## Not run:
# Calculate the log prices
prices <- na.omit(rutils::etfenv$prices[, c("XLP", "VTI")])
prices <- log(prices)
# Calculate the Hurst exponent from 21 day aggregations
calc_hurst(prices, step=21)

## End(Not run)
```

<code>calc_hurst_ohlc</code>	<i>Calculate the Hurst exponent from the volatility ratio of aggregated OHLC prices.</i>
------------------------------	--

Description

Calculate the Hurst exponent from the volatility ratio of aggregated *OHLC* prices.

Usage

```
calc_hurst_ohlc(
  ohlc,
  step = 1L,
  method = "yang_zhang",
  close_lag = 0L,
  scale = TRUE,
  index = 0L
)
```

Arguments

ohlc	A <i>time series</i> or a <i>matrix</i> of <i>OHLC</i> prices.
step	The number of periods in each interval between neighboring end points.
method	A <i>character</i> string representing the price range estimator for calculating the variance. The estimators include: <ul style="list-style-type: none"> • "close" close-to-close estimator, • "rogers_satchell" Rogers-Satchell estimator, • "garman_klass" Garman-Klass estimator, • "garman_klass_yz" Garman-Klass with account for close-to-open price jumps, • "yang_zhang" Yang-Zhang estimator, (The default is the method = "yang_zhang".)
close_lag	A <i>vector</i> with the lagged <i>close</i> prices of the <i>OHLC time series</i> . This is an optional argument. (The default is close_lag = 0).
scale	<i>Boolean</i> argument: Should the returns be divided by the time index, the number of seconds in each period? (The default is scale = TRUE).
index	A <i>vector</i> with the time index of the <i>time series</i> . This is an optional argument (the default is index = 0).

Details

The function `calc_hurst_ohlc()` calculates the Hurst exponent from the ratios of the volatilities of aggregated *OHLC* prices.

The aggregated volatility σ_t scales (increases) with the length of the aggregation interval Δt raised to the power of the *Hurst exponent* H :

$$\sigma_t = \sigma \Delta t^H$$

Where σ is the daily return volatility.

The *Hurst exponent* H is equal to the logarithm of the ratio of the volatilities divided by the logarithm of the time interval Δt :

$$H = \frac{\log \sigma_t - \log \sigma}{\log \Delta t}$$

The function `calc_hurst_ohlc()` calls the function `calc_var_ohlc_ag()` to calculate the aggregated volatility σ_t .

Value

The Hurst exponent calculated from the variance ratio of aggregated *OHLC* prices.

Examples

```
## Not run:
# Calculate the log ohlc prices
ohlc <- log(rutils::etfenv$VTI)
# Calculate the Hurst exponent from 21 day aggregations
calc_hurst_ohlc(ohlc, step=21)

## End(Not run)
```

calc_inv	Calculate the regularized inverse of a matrix of data using Singular Value Decomposition (SVD).
----------	---

Description

Calculate the regularized inverse of a *matrix* of data using Singular Value Decomposition (SVD).

Usage

```
calc_inv(tseries, eigen_thresh = 0.01, eigen_max = 0L)
```

Arguments

tseries	A <i>time series</i> or <i>matrix</i> of data.
eigen_thresh	A <i>numeric</i> threshold level for discarding small singular values in order to regularize the inverse of the matrix tseries (the default is 0.01).
eigen_max	An <i>integer</i> equal to the number of singular values used for calculating the regularized inverse of the matrix tseries (the default is eigen_max = 0 - equivalent to eigen_max equal to the number of columns of tseries).

Details

The function `calc_inv()` calculates the regularized inverse of the matrix `tseries` using Singular Value Decomposition (SVD).

The function `calc_inv()` first performs Singular Value Decomposition (SVD) of the matrix `tseries`. The SVD of a matrix C is defined as the factorization:

$$C = U \Sigma V^T$$

Where U and V are the left and right *singular matrices*, and Σ is a diagonal matrix of *singular values* $\Sigma = \{\sigma_i\}$.

The inverse C^{-1} of the matrix C can be calculated from the SVD matrices as:

$$C^{-1} = V \Sigma^{-1} U^T$$

The *regularized inverse* of the matrix C is given by:

$$C^{-1} = V_n \Sigma_n^{-1} U_n^T$$

Where U_n , V_n and Σ_n are the SVD matrices with the rows and columns corresponding to zero *singular values* removed.

The function `calc_inv()` performs regularization by discarding the smallest singular values σ_i that are less than the threshold level `eigen_thresh` times the sum of all the singular values:

$$\sigma_i < \text{eigen_thresh} \cdot \left(\sum \sigma_i \right)$$

It then discards additional singular values so that only the largest `eigen_max` singular values remain. It calculates the regularized inverse from the SVD matrices using only the largest singular values up to `eigen_max`. For example, if `eigen_max` = 3 then it only uses the 3 largest singular values. This has the effect of dimension shrinkage.

If the matrix `tseries` has a large number of small singular values, then the number of remaining singular values may be less than `eigen_max`.

Value

A *matrix* equal to the regularized inverse of the matrix *tseries*.

Examples

```
## Not run:
# Calculate ETF returns
returns <- na.omit(rutils::etfenv$returns)
# Calculate covariance matrix
covmat <- cov(returns)
# Calculate regularized inverse using RcppArmadillo
invmat <- HighFreq::calc_inv(covmat, eigen_max=3)
# Calculate regularized inverse from SVD in R
svdec <- svd(covmat)
eigen_max <- 1:3
invsvd <- svdec$v[, eigen_max] %*% (t(svdec$u[, eigen_max]) / svdec$d[eigen_max])
# Compare RcppArmadillo with R
all.equal(invmat, invsvd)

## End(Not run)
```

calc_kurtosis	Calculate the kurtosis of the columns of a time series or a matrix using RcppArmadillo.
---------------	---

Description

Calculate the kurtosis of the columns of a *time series* or a *matrix* using RcppArmadillo.

Usage

```
calc_kurtosis(tseries, method = "moment", confl = 0.75)
```

Arguments

tseries	A <i>time series</i> or a <i>matrix</i> of data.
method	A <i>string</i> specifying the type of the kurtosis model (the default is method = "moment" - see Details).
confl	The confidence level for calculating the quantiles of returns (the default is confl = 0.75).

Details

The function `calc_kurtosis()` calculates the kurtosis of the columns of the *matrix* *tseries* using RcppArmadillo C++ code.

If `method = "moment"` (the default) then `calc_kurtosis()` calculates the fourth moment of the data. But it doesn't de-mean the columns of *tseries* because that requires copying the matrix *tseries*, so it's time-consuming.

If method = "quantile" then it calculates the skewness κ from the differences between the quantiles of the data as follows:

$$\kappa = \frac{q_{\alpha} - q_{1-\alpha}}{q_{0.75} - q_{0.25}}$$

Where α is the confidence level for calculating the quantiles.

If method = "nonparametric" then it calculates the kurtosis as the difference between the mean of the data minus its median, divided by the standard deviation.

If the number of rows of tseries is less than 3 then it returns zeros.

The code examples below compare the function calc_kurtosis() with the kurtosis calculated using R code.

Value

A single-row matrix with the kurtosis of the columns of tseries.

Examples

```
## Not run:
# Define a single-column time series of returns
returns <- na.omit(rutils::etfenv$returns$VTI)
# Calculate the moment kurtosis
HighFreq::calc_kurtosis(returns)
# Calculate the moment kurtosis in R
calc_kurtr <- function(x) {
  x <- (x-mean(x))
  sum(x^4)/var(x)^2/NROW(x)
} # end calc_kurtr
all.equal(HighFreq::calc_kurtosis(returns),
  calc_kurtr(returns), check.attributes=FALSE)
# Compare the speed of RcppArmadillo with R code
library(microbenchmark)
summary(microbenchmark(
  Rcpp=HighFreq::calc_kurtosis(returns),
  Rcode=calc_kurtr(returns),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary
# Calculate the quantile kurtosis
HighFreq::calc_kurtosis(returns, method="quantile", conf=0.9)
# Calculate the quantile kurtosis in R
calc_kurtq <- function(x, a=0.9) {
  quantiles <- quantile(x, c(1-a, 0.25, 0.75, a), type=5)
  (quantiles[4] - quantiles[1])/(quantiles[3] - quantiles[2])
} # end calc_kurtq
all.equal(drop(HighFreq::calc_kurtosis(returns, method="quantile", conf=0.9)),
  calc_kurtq(returns, a=0.9), check.attributes=FALSE)
# Compare the speed of RcppArmadillo with R code
summary(microbenchmark(
  Rcpp=HighFreq::calc_kurtosis(returns, method="quantile"),
  Rcode=calc_kurtq(returns),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary
# Calculate the nonparametric kurtosis
HighFreq::calc_kurtosis(returns, method="nonparametric")
# Compare HighFreq::calc_kurtosis() with R nonparametric kurtosis
all.equal(drop(HighFreq::calc_kurtosis(returns, method="nonparametric")),
  (mean(returns)-median(returns))/sd(returns),
  check.attributes=FALSE)
```

```
# Compare the speed of RcppArmadillo with R code
summary(microbenchmark(
  Rcpp=HighFreq::calc_kurtosis(returns, method="nonparametric"),
  Rcode=(mean(returns)-median(returns))/sd(returns),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary

## End(Not run)
```

calc_lm	<i>Perform multivariate linear regression using least squares and return a named list of regression coefficients, their t-values, and p-values.</i>
---------	---

Description

Perform multivariate linear regression using least squares and return a named list of regression coefficients, their t-values, and p-values.

Usage

```
calc_lm(response, predictor)
```

Arguments

response	A single-column <i>time series</i> or a <i>vector</i> of response data.
predictor	A <i>time series</i> or a <i>matrix</i> of predictor data.

Details

The function `calc_lm()` performs the same calculations as the function `lm()` from package *stats*. It uses RcppArmadillo C++ code so it's several times faster than `lm()`. The code was inspired by this article (but it's not identical to it): <http://gallery.rcpp.org/articles/fast-linear-model-with-armadillo/>

Value

A named list with three elements: a *matrix* of coefficients (named "*coefficients*"), the *z-score* of the last residual (named "*zscore*"), and a *vector* with the R-squared and F-statistic (named "*stats*"). The numeric *matrix* of coefficients named "*coefficients*" contains the alpha and beta coefficients, and their *t-values* and *p-values*.

Examples

```
## Not run:
# Calculate historical returns
returns <- na.omit(rutils::etfenv$returns[, c("XLF", "VTI", "IEF")])
# Response equals XLF returns
response <- returns[, 1]
# Predictor matrix equals VTI and IEF returns
predictor <- returns[, -1]
# Perform multivariate regression using lm()
lmod <- lm(response ~ predictor)
lmodsum <- summary(lmod)
# Perform multivariate regression using calc_lm()
```

```

reg_arma <- HighFreq::calc_lm(response=response, predictor=predictor)
# Compare the outputs of both functions
all.equal(reg_arma$coefficients[, "coeff"], unname(coef(lmod)))
all.equal(unname(reg_arma$coefficients), unname(lmodsum$coefficients))
all.equal(unname(reg_arma$stats), c(lmodsum$r.squared, unname(lmodsum$fstatistic[1])))
# Compare the speed of RcppArmadillo with R code
summary(microbenchmark(
  Rcpp=HighFreq::calc_lm(response=response, predictor=predictor),
  Rcode=lm(response ~ predictor),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary

## End(Not run)

```

calc_mean	<i>Calculate the mean (location) of the columns of a time series or a matrix using RcppArmadillo.</i>
-----------	---

Description

Calculate the mean (location) of the columns of a *time series* or a *matrix* using RcppArmadillo.

Usage

```
calc_mean(tseries, method = "moment", conf1 = 0.75)
```

Arguments

tseries	A <i>time series</i> or a <i>matrix</i> of data.
method	A <i>string</i> specifying the type of the mean (location) model (the default is method = "moment" - see Details).
conf1	The confidence level for calculating the quantiles of returns (the default is conf1 = 0.75).

Details

The function `calc_mean()` calculates the mean (location) values of the columns of the *time series* `tseries` using RcppArmadillo C++ code.

If `method = "moment"` (the default) then `calc_mean()` calculates the location as the mean - the first moment of the data.

If `method = "quantile"` then it calculates the location μ as the sum of the quantiles as follows:

$$\mu = q_{\alpha} + q_{1-\alpha}$$

Where α is the confidence level for calculating the quantiles.

If `method = "nonparametric"` then it calculates the location as the median.

The code examples below compare the function `calc_mean()` with the mean (location) calculated using R code.

Value

A single-row matrix with the mean (location) of the columns of `tseries`.

Examples

```
## Not run:
# Calculate historical returns
returns <- na.omit(rutils::etfenv$returns[, c("XLP", "VTI")])
# Calculate the column means in RcppArmadillo
HighFreq::calc_mean(returns)
# Calculate the column means in R
sapply(returns, mean)
# Compare the values
all.equal(drop(HighFreq::calc_mean(returns)),
  sapply(returns, mean), check.attributes=FALSE)
# Compare the speed of RcppArmadillo with R code
library(microbenchmark)
summary(microbenchmark(
  Rcpp=HighFreq::calc_mean(returns),
  Rcode=sapply(returns, mean),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary
# Calculate the quantile mean (location)
HighFreq::calc_mean(returns, method="quantile", confl=0.9)
# Calculate the quantile mean (location) in R
colSums(sapply(returns, quantile, c(0.9, 0.1), type=5))
# Compare the values
all.equal(drop(HighFreq::calc_mean(returns, method="quantile", confl=0.9)),
  colSums(sapply(returns, quantile, c(0.9, 0.1), type=5)),
  check.attributes=FALSE)
# Compare the speed of RcppArmadillo with R code
summary(microbenchmark(
  Rcpp=HighFreq::calc_mean(returns, method="quantile", confl=0.9),
  Rcode=colSums(sapply(returns, quantile, c(0.9, 0.1), type=5)),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary
# Calculate the column medians in RcppArmadillo
HighFreq::calc_mean(returns, method="nonparametric")
# Calculate the column medians in R
sapply(returns, median)
# Compare the values
all.equal(drop(HighFreq::calc_mean(returns, method="nonparametric")),
  sapply(returns, median), check.attributes=FALSE)
# Compare the speed of RcppArmadillo with R code
summary(microbenchmark(
  Rcpp=HighFreq::calc_mean(returns, method="nonparametric"),
  Rcode=sapply(returns, median),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary

## End(Not run)
```

calc_ranks

Calculate the ranks of the elements of a single-column time series or a vector using RcppArmadillo.

Description

Calculate the ranks of the elements of a single-column *time series* or a *vector* using RcppArmadillo.

Usage

```
calc_ranks(tseries)
```

Arguments

tseries A single-column *time series* or a *vector*.

Details

The function `calc_ranks()` calculates the ranks of the elements of a single-column *time series* or a *vector*. It uses the RcppArmadillo function `arma::sort_index()`. The function `arma::sort_index()` calculates the permutation index to sort a given vector into ascending order.

Applying the function `arma::sort_index()` twice: `arma::sort_index(arma::sort_index())`, calculates the *reverse* permutation index to sort the vector from ascending order back into its original unsorted order. The permutation index produced by: `arma::sort_index(arma::sort_index())` is the *reverse* of the permutation index produced by: `arma::sort_index()`.

The ranks of the elements are equal to the *reverse* permutation index. The function `calc_ranks()` calculates the *reverse* permutation index.

Value

An *integer vector* with the ranks of the elements of the `tseries`.

Examples

```
## Not run:
# Create a vector of random data
datav <- round(runif(7), 2)
# Calculate the ranks of the elements in two ways
all.equal(rank(datav), drop(HighFreq::calc_ranks(datav)))
# Create a time series of random data
datav <- xts::xts(runif(7), seq.Date(Sys.Date(), by=1, length.out=7))
# Calculate the ranks of the elements in two ways
all.equal(rank(coredatav(datav)), drop(HighFreq::calc_ranks(datav)))
# Compare the speed of RcppArmadillo with R code
datav <- runif(7)
library(microbenchmark)
summary(microbenchmark(
  Rcpp=calc_ranks(datav),
  Rcode=rank(datav),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary

## End(Not run)
```

calc_reg

Perform multivariate regression using different methods, and return a vector of regression coefficients, their t-values, and the last residual z-score.

Description

Perform multivariate regression using different methods, and return a vector of regression coefficients, their t-values, and the last residual z-score.

Usage

```
calc_reg(
  response,
  predictor,
  intercept = TRUE,
  method = "least_squares",
  eigen_thresh = 1e-05,
  eigen_max = 0L,
  confl = 0.1,
  alpha = 0
)
```

Arguments

response	A single-column <i>time series</i> or a <i>vector</i> of response data.
predictor	A <i>time series</i> or a <i>matrix</i> of predictor data.
intercept	A <i>Boolean</i> specifying whether an intercept term should be added to the predictor (the default is intercept = TRUE).
method	A <i>string</i> specifying the type of the regression model the default is method = "least_squares" - see Details).
eigen_thresh	A <i>numeric</i> threshold level for discarding small singular values in order to regularize the inverse of the predictor matrix (the default is 1e-5).
eigen_max	An <i>integer</i> equal to the number of singular values used for calculating the regularized inverse of the predictor matrix (the default is 0 - equivalent to eigen_max equal to the number of columns of predictor).
confl	The confidence level for calculating the quantiles of returns (the default is confl = 0.75).
alpha	The shrinkage intensity between 0 and 1. (the default is 0).

Details

The function `calc_reg()` performs multivariate regression using different methods, and returns a vector of regression coefficients, their t-values, and the last residual z-score.

If `method = "least_squares"` (the default) then it performs the standard least squares regression, the same as the function `calc_lm()`, and the function `lm()` from the R package *stats*. But it uses RcppArmadillo C++ code so it's several times faster than `lm()`.

If `method = "regular"` then it performs shrinkage regression. It calculates the regularized inverse of the predictor matrix from its singular value decomposition. It performs regularization by selecting only the largest singular values equal in number to `eigen_max`.

If `method = "quantile"` then it performs quantile regression (not implemented yet).

If `intercept = TRUE` then an extra intercept column (unit column) is added to the predictor matrix (the default is `intercept = FALSE`).

The length of the return vector depends on the number of columns of the predictor matrix (including the intercept column, if it's added). The length of the return vector is equal to the number of

regression coefficients, plus their t-values, plus the z-score. The number of regression coefficients is equal to the number of columns of the predictor matrix (including the intercept column, if it's added). The number of t-values is equal to the number of coefficients.

For example, if the number of columns of the predictor matrix is equal to n , and if `intercept = TRUE` (the default), then `calc_reg()` returns a vector with $2n+3$ elements: $n+1$ regression coefficients (including the intercept coefficient), $n+1$ corresponding t-values, and 1 z-score value.

If `intercept = FALSE`, then `calc_reg()` returns a vector with $2n+1$ elements: n regression coefficients (without the intercept coefficient), n corresponding t-values, and 1 z-score value.

Value

A single-row matrix with A vector with the regression coefficients, their t-values, and the last residual z-score.

Examples

```
## Not run:
# Calculate historical returns
returns <- na.omit(rutils::etfenv$returns[, c("XLF", "VTI", "IEF")])
# Response equals XLF returns
response <- returns[, 1]
# Predictor matrix equals VTI and IEF returns
predictor <- returns[, -1]
# Perform multivariate regression using lm()
lmod <- lm(response ~ predictor)
lmodsum <- summary(lmod)
coeff <- lmodsum$coefficients
# Perform multivariate regression using calc_reg()
reg_arma <- drop(HighFreq::calc_reg(response=response, predictor=predictor))
# Compare the outputs of both functions
all.equal(reg_arma[1:(2*(1+NCOL(predictor)))],
  c(coeff[, "Estimate"], coeff[, "t value"]), check.attributes=FALSE)
# Compare the speed of RcppArmadillo with R code
library(microbenchmark)
summary(microbenchmark(
  Rcpp=HighFreq::calc_reg(response=response, predictor=predictor),
  Rcode=lm(response ~ predictor),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary

## End(Not run)
```

calc_scaled	Scale (standardize) the columns of a matrix of data using RcppArmadillo.
-------------	--

Description

Scale (standardize) the columns of a *matrix* of data using RcppArmadillo.

Usage

```
calc_scaled(tseries, use_median = FALSE)
```

Arguments

tseries	A <i>time series</i> or <i>matrix</i> of data.
use_median	A <i>Boolean</i> argument: if TRUE then the centrality (central tendency) is calculated as the <i>median</i> and the dispersion is calculated as the <i>median absolute deviation</i> (<i>MAD</i>). If use_median = FALSE then the centrality is calculated as the <i>mean</i> and the dispersion is calculated as the <i>standard deviation</i> (the default is FALSE)

Details

The function `calc_scaled()` scales (standardizes) the columns of the `tseries` argument using RcppArmadillo.

If the argument `use_median` is FALSE (the default), then it performs the same calculation as the standard R function `scale()`, and it calculates the centrality (central tendency) as the *mean* and the dispersion as the *standard deviation*.

If the argument `use_median` is TRUE, then it calculates the centrality as the *median* and the dispersion as the *median absolute deviation* (*MAD*).

If the number of rows of `tseries` is less than 3 then it returns `tseries` unscaled.

The function `calc_scaled()` uses RcppArmadillo C++ code and is about 5 times faster than function `scale()`, for a *matrix* with 1,000 rows and 20 columns.

Value

A *matrix* with the same dimensions as the input argument `tseries`.

Examples

```
## Not run:
# Create a matrix of random data
returns <- matrix(rnorm(20000), nc=20)
scaled <- calc_scaled(tseries=returns, use_median=FALSE)
scaled2 <- scale(returns)
all.equal(scaled, scaled2, check.attributes=FALSE)
# Compare the speed of Rcpp with R code
library(microbenchmark)
summary(microbenchmark(
  Rcpp=calc_scaled(tseries=returns, use_median=FALSE),
  Rcode=scale(returns),
  times=100))[, c(1, 4, 5)] # end microbenchmark summary

## End(Not run)
```

calc_skew	Calculate the skewness of the columns of a <i>time series</i> or a <i>matrix</i> using RcppArmadillo.
-----------	---

Description

Calculate the skewness of the columns of a *time series* or a *matrix* using RcppArmadillo.

Usage

```
calc_skew(tseries, method = "moment", confl = 0.75)
```

Arguments

tseries	A <i>time series</i> or a <i>matrix</i> of data.
method	A <i>string</i> specifying the type of the skewness model (the default is method = "moment" - see Details).
confl	The confidence level for calculating the quantiles of returns (the default is confl = 0.75).

Details

The function `calc_skew()` calculates the skewness of the columns of a *time series* or a *matrix* of data using RcppArmadillo C++ code.

If method = "moment" (the default) then `calc_skew()` calculates the skewness as the third moment of the data.

If method = "quantile" then it calculates the skewness ς from the differences between the quantiles of the data as follows:

$$\varsigma = \frac{q_{\alpha} + q_{1-\alpha} - 2 * q_{0.5}}{q_{\alpha} - q_{1-\alpha}}$$

Where α is the confidence level for calculating the quantiles.

If method = "nonparametric" then it calculates the skewness as the difference between the mean of the data minus its median, divided by the standard deviation.

If the number of rows of tseries is less than 3 then it returns zeros.

The code examples below compare the function `calc_skew()` with the skewness calculated using R code.

Value

A single-row matrix with the skewness of the columns of tseries.

Examples

```
## Not run:
# Define a single-column time series of returns
returns <- na.omit(rutils::etfenv$returns$VTI)
# Calculate the moment skewness
HighFreq::calc_skew(returns)
# Calculate the moment skewness in R
calc_skewr <- function(x) {
  x <- (x-mean(x))
  sum(x^3)/var(x)^1.5/NROW(x)
} # end calc_skewr
all.equal(HighFreq::calc_skew(returns),
  calc_skewr(returns), check.attributes=FALSE)
# Compare the speed of RcppArmadillo with R code
library(microbenchmark)
summary(microbenchmark(
  Rcpp=HighFreq::calc_skew(returns),
  Rcode=calc_skewr(returns),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary
```

```

# Calculate the quantile skewness
HighFreq::calc_skew(returns, method="quantile", confl=0.9)
# Calculate the quantile skewness in R
calc_skewq <- function(x, a = 0.75) {
  quantiles <- quantile(x, c(1-a, 0.5, a), type=5)
  (quantiles[3] + quantiles[1] - 2*quantiles[2])/(quantiles[3] - quantiles[1])
} # end calc_skewq
all.equal(drop(HighFreq::calc_skew(returns, method="quantile", confl=0.9)),
  calc_skewq(returns, a=0.9), check.attributes=FALSE)
# Compare the speed of RcppArmadillo with R code
summary(microbenchmark(
  Rcpp=HighFreq::calc_skew(returns, method="quantile"),
  Rcode=calc_skewq(returns),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary
# Calculate the nonparametric skewness
HighFreq::calc_skew(returns, method="nonparametric")
# Compare HighFreq::calc_skew() with R nonparametric skewness
all.equal(drop(HighFreq::calc_skew(returns, method="nonparametric")),
  (mean(returns)-median(returns))/sd(returns),
  check.attributes=FALSE)
# Compare the speed of RcppArmadillo with R code
summary(microbenchmark(
  Rcpp=HighFreq::calc_skew(returns, method="nonparametric"),
  Rcode=(mean(returns)-median(returns))/sd(returns),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary

## End(Not run)

```

calc_startpoints	<i>Calculate a vector of start points by lagging (shifting) a vector of end points.</i>
------------------	---

Description

Calculate a vector of start points by lagging (shifting) a vector of end points.

Usage

```
calc_startpoints(endp, look_back)
```

Arguments

endp	An <i>integer</i> vector of end points.
look_back	The length of the look-back interval, equal to the lag (shift) applied to the end points.

Details

The start points are equal to the values of the vector endp lagged (shifted) by an amount equal to look_back. In addition, an extra value of 1 is added to them, to avoid data overlaps. The lag operation requires appending a beginning warmup interval containing zeros, so that the vector of start points has the same length as the endp.

For example, consider the end points for a vector of length 25 divided into equal intervals of length 5: 4, 9, 14, 19, 24. (In C++ the vector indexing starts at 0 not 1, so it's shifted by -1.) Then the start points for look_back = 2 are equal to: 0, 0, 5, 10, 15. The differences between the end points minus the corresponding start points are equal to 9, except for the warmup interval.

Value

An *integer* vector with the same number of elements as the vector endp.

Examples

```
# Calculate end points
endp <- HighFreq::calc_endpoints(25, 5)
# Calculate start points corresponding to the end points
startp <- HighFreq::calc_startpoints(endp, 2)
```

calc_var	<i>Calculate the dispersion (variance) of the columns of a time series or a matrix using RcppArmadillo.</i>
----------	---

Description

Calculate the dispersion (variance) of the columns of a *time series* or a *matrix* using RcppArmadillo.

Usage

```
calc_var(tseries, method = "moment", conf1 = 0.75)
```

Arguments

tseries	A <i>time series</i> or a <i>matrix</i> of data.
method	A <i>string</i> specifying the type of the dispersion model (the default is method = "moment" - see Details).
conf1	The confidence level for calculating the quantiles of returns (the default is conf1 = 0.75).

Details

The function `calc_var()` calculates the dispersion of the columns of a *time series* or a *matrix* of data using RcppArmadillo C++ code.

The dispersion is a measure of the variability of the data. Examples of dispersion are the variance and the Median Absolute Deviation (*MAD*).

If method = "moment" (the default) then `calc_var()` calculates the dispersion as the second moment of the data σ^2 (the variance). Then `calc_var()` performs the same calculation as the function `colVars()` from package `matrixStats`, but it's much faster because it uses RcppArmadillo C++ code.

If method = "quantile" then it calculates the dispersion as the difference between the quantiles as follows:

$$\mu = q_{\alpha} - q_{1-\alpha}$$

Where α is the confidence level for calculating the quantiles.

If method = "nonparametric" then it calculates the dispersion as the Median Absolute Deviation (*MAD*):

$$MAD = \text{median}(\text{abs}(x - \text{median}(x)))$$

It also multiplies the *MAD* by a factor of 1.4826, to make it comparable to the standard deviation.

If method = "nonparametric" then calc_var() performs the same calculation as the function stats::mad(), but it's much faster because it uses RcppArmadillo C++ code.

If the number of rows of tseries is less than 3 then it returns zeros.

Value

A row vector equal to the dispersion of the columns of the matrix tseries.

Examples

```
## Not run:
# Calculate VTI and XLF returns
returns <- na.omit(rutils::etfenv$returns[, c("VTI", "XLF")])
# Compare HighFreq::calc_var() with standard var()
all.equal(drop(HighFreq::calc_var(returns)),
  apply(returns, 2, var), check.attributes=FALSE)
# Compare HighFreq::calc_var() with matrixStats
all.equal(drop(HighFreq::calc_var(returns)),
  matrixStats::colVars(returns), check.attributes=FALSE)
# Compare the speed of RcppArmadillo with matrixStats and with R code
library(microbenchmark)
summary(microbenchmark(
  Rcpp=HighFreq::calc_var(returns),
  matrixStats=matrixStats::colVars(returns),
  Rcode=apply(returns, 2, var),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary
# Compare HighFreq::calc_var() with stats::mad()
all.equal(drop(HighFreq::calc_var(returns, method="nonparametric")),
  sapply(returns, mad), check.attributes=FALSE)
# Compare the speed of RcppArmadillo with stats::mad()
summary(microbenchmark(
  Rcpp=HighFreq::calc_var(returns, method="nonparametric"),
  Rcode=sapply(returns, mad),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary

## End(Not run)
```

calc_varvec

Calculate the variance of a single-column time series or a vector using RcppArmadillo.

Description

Calculate the variance of a single-column *time series* or a *vector* using RcppArmadillo.

Usage

```
calc_varvec(tseries)
```

Arguments

tseries A single-column *time series* or a *vector*.

Details

The function `calc_varvec()` calculates the variance of a *vector* using RcppArmadillo C++ code, so it's significantly faster than the R function `var()`.

Value

A *numeric* value equal to the variance of the *vector*.

Examples

```
## Not run:
# Create a vector of random returns
returns <- rnorm(1e6)
# Compare calc_varvec() with standard var()
all.equal(HighFreq::calc_varvec(returns),
  var(returns))
# Compare the speed of RcppArmadillo with R code
library(microbenchmark)
summary(microbenchmark(
  Rcpp=HighFreq::calc_varvec(returns),
  Rcode=var(returns),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary

## End(Not run)
```

calc_var_ag

Calculate the variance of returns aggregated over end points.

Description

Calculate the variance of returns aggregated over end points.

Usage

```
calc_var_ag(tseries, step = 1L)
```

Arguments

tseries A *time series* or a *matrix* of prices.

step The number of periods in each interval between neighboring end points.

Details

The function `calc_var_ag()` calculates the variance of returns aggregated over end points.

It first calculates the end points spaced apart by the number of periods equal to the argument `step`. Then it calculates the aggregated returns by differencing the prices `tseries` calculated at the end points. Finally it calculates the variance of the returns.

If there are extra periods that don't fit over the length of `tseries`, then `calc_var_ag()` loops over all possible stub intervals, then it calculates all the corresponding variance values, and averages them.

For example, if the number of rows of `tseries` is equal to 20, and `step=3` then 6 end points fit over the length of `tseries`, and there are 2 extra periods that must fit into stubs, either at the beginning or at the end (or both).

The aggregated volatility σ_t scales (increases) with the length of the aggregation interval Δt raised to the power of the *Hurst exponent* H :

$$\sigma_t = \sigma \Delta t^H$$

Where σ is the daily return volatility.

The function `calc_var_ag()` can therefore be used to calculate the *Hurst exponent* from the volatility ratio.

Value

The variance of aggregated returns.

Examples

```
## Not run:
# Calculate the log prices
prices <- na.omit(rutils::etfenv$prices[, c("XLP", "VTI")])
prices <- log(prices)
# Calculate the daily variance of percentage returns
calc_var_ag(prices, step=1)
# Calculate the daily variance using R
sapply(rutils::diffit(prices), var)
# Calculate the variance of returns aggregated over 21 days
calc_var_ag(prices, step=21)
# The variance over 21 days is approximately 21 times the daily variance
21*calc_var_ag(prices, step=1)

## End(Not run)
```

calc_var_ohlc

Calculate the variance of returns from OHLC prices using different price range estimators.

Description

Calculate the variance of returns from *OHLC* prices using different price range estimators.

Usage

```
calc_var_ohlc(
  ohlc,
  method = "yang_zhang",
  close_lag = 0L,
  scale = TRUE,
  index = 0L
)
```

Arguments

ohlc	A <i>time series</i> or a <i>matrix</i> of <i>OHLC</i> prices.
method	A <i>character</i> string representing the price range estimator for calculating the variance. The estimators include: <ul style="list-style-type: none"> • "close" close-to-close estimator, • "rogers_satchell" Rogers-Satchell estimator, • "garman_klass" Garman-Klass estimator, • "garman_klass_yz" Garman-Klass with account for close-to-open price jumps, • "yang_zhang" Yang-Zhang estimator, (The default is the method = "yang_zhang".)
close_lag	A <i>vector</i> with the lagged <i>close</i> prices of the <i>OHLC time series</i> . This is an optional argument. (The default is close_lag = 0).
scale	<i>Boolean</i> argument: Should the returns be divided by the time index, the number of seconds in each period? (The default is scale = TRUE).
index	A <i>vector</i> with the time index of the <i>time series</i> . This is an optional argument (the default is index = 0).

Details

The function `calc_var_ohlc()` calculates the variance from all the different intra-day and day-over-day returns (defined as the differences of *OHLC* prices), using several different variance estimation methods.

The function `calc_var_ohlc()` does not calculate the logarithm of the prices. So if the argument `ohlc` contains dollar prices then `calc_var_ohlc()` calculates the dollar variance. If the argument `ohlc` contains the log prices then `calc_var_ohlc()` calculates the percentage variance.

The default method is "yang_zhang", which theoretically has the lowest standard error among unbiased estimators. The methods "close", "garman_klass_yz", and "yang_zhang" do account for *close-to-open* price jumps, while the methods "garman_klass" and "rogers_satchell" do not account for *close-to-open* price jumps.

If `scale` is `TRUE` (the default), then the returns are divided by the differences of the time index (which scales the variance to the units of variance per second squared). This is useful when calculating the variance from minutely bar data, because dividing returns by the number of seconds decreases the effect of overnight price jumps. If the time index is in days, then the variance is equal to the variance per day squared.

If the number of rows of `ohlc` is less than 3 then it returns zero.

The optional argument `index` is the time index of the *time series* `ohlc`. If the time index is in seconds, then the differences of the index are equal to the number of seconds in each time period. If the time index is in days, then the differences are equal to the number of days in each time period.

The optional argument `close_lag` are the lagged *close* prices of the *OHLC time series*. Passing in the lagged *close* prices speeds up the calculation, so it's useful for rolling calculations.

The function `calc_var_ohlc()` is implemented in RcppArmadillo C++ code, and it's over 10 times faster than `calc_var_ohlc_r()`, which is implemented in R code.

Value

A single *numeric* value equal to the variance of the *OHLC time series*.

Examples

```
## Not run:
# Extract the log OHLC prices of SPY
ohlc <- log(HighFreq::SPY)
# Extract the time index of SPY prices
indeks <- c(1, diff(xts::.index(ohlc)))
# Calculate the variance of SPY returns, with scaling of the returns
HighFreq::calc_var_ohlc(ohlc,
  method="yang_zhang", scale=TRUE, index=indeks)
# Calculate variance without accounting for overnight jumps
HighFreq::calc_var_ohlc(ohlc,
  method="rogers_satchell", scale=TRUE, index=indeks)
# Calculate the variance without scaling the returns
HighFreq::calc_var_ohlc(ohlc, scale=FALSE)
# Calculate the variance by passing in the lagged close prices
close_lag <- HighFreq::lagit(ohlc[, 4])
all.equal(HighFreq::calc_var_ohlc(ohlc),
  HighFreq::calc_var_ohlc(ohlc, close_lag=close_lag))
# Compare with HighFreq::calc_var_ohlc_r()
all.equal(HighFreq::calc_var_ohlc(ohlc, index=indeks),
  HighFreq::calc_var_ohlc_r(ohlc))
# Compare the speed of Rcpp with R code
library(microbenchmark)
summary(microbenchmark(
  Rcpp=HighFreq::calc_var_ohlc(ohlc),
  Rcode=HighFreq::calc_var_ohlc_r(ohlc),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary

## End(Not run)
```

<code>calc_var_ohlc_ag</code>	<i>Calculate the variance of aggregated OHLC prices using different price range estimators.</i>
-------------------------------	---

Description

Calculate the variance of aggregated *OHLC* prices using different price range estimators.

Usage

```
calc_var_ohlc_ag(
  ohlc,
  step = 1L,
```



```

    method = "yang_zhang",
    close_lag = 0L,
    scale = TRUE,
    index = 0L
)

```

Arguments

ohlc	A <i>time series</i> or a <i>matrix</i> of <i>OHLC</i> prices.
step	The number of periods in each interval between neighboring end points.
method	A <i>character</i> string representing the price range estimator for calculating the variance. The estimators include: <ul style="list-style-type: none"> • "close" close-to-close estimator, • "rogers_satchell" Rogers-Satchell estimator, • "garman_klass" Garman-Klass estimator, • "garman_klass_yz" Garman-Klass with account for close-to-open price jumps, • "yang_zhang" Yang-Zhang estimator, (The default is the method = "yang_zhang".)
close_lag	A <i>vector</i> with the lagged <i>close</i> prices of the <i>OHLC time series</i> . This is an optional argument. (The default is close_lag = 0).
scale	<i>Boolean</i> argument: Should the returns be divided by the time index, the number of seconds in each period? (The default is scale = TRUE).
index	A <i>vector</i> with the time index of the <i>time series</i> . This is an optional argument (the default is index = 0).

Details

The function `calc_var_ohlc_ag()` calculates the variance of *OHLC* prices aggregated over end points.

It first calculates the end points spaced apart by the number of periods equal to the argument `step`. Then it aggregates the *OHLC* prices to the end points. Finally it calculates the variance of the aggregated *OHLC* prices.

If there are extra periods that don't fit over the length of `ohlc`, then `calc_var_ohlc_ag()` loops over all possible stub intervals, it calculates all the corresponding variance values, and it averages them.

For example, if the number of rows of `ohlc` is equal to 20, and `step=3` then 6 end points fit over the length of `ohlc`, and there are 2 extra periods that must fit into stubs, either at the beginning or at the end (or both).

The aggregated volatility σ_t scales (increases) with the length of the aggregation interval Δt raised to the power of the *Hurst exponent* H :

$$\sigma_t = \sigma \Delta t^H$$

Where σ is the daily return volatility.

The function `calc_var_ohlc_ag()` can therefore be used to calculate the *Hurst exponent* from the volatility ratio.

Value

The variance of aggregated *OHLC* prices.

Examples

```
## Not run:
# Calculate the log ohlc prices
ohlc <- log(rutils::etfenv$VTI)
# Calculate the daily variance of percentage returns
calc_var_ohlc_ag(ohlc, step=1)
# Calculate the variance of returns aggregated over 21 days
calc_var_ohlc_ag(ohlc, step=21)
# The variance over 21 days is approximately 21 times the daily variance
21*calc_var_ohlc_ag(ohlc, step=1)

## End(Not run)
```

calc_var_ohlc_r	<i>Calculate the variance of an OHLC time series, using different range estimators for variance.</i>
-----------------	--

Description

Calculate the variance of an *OHLC* time series, using different range estimators for variance.

Usage

```
calc_var_ohlc_r(ohlc, method = "yang_zhang", scalit = TRUE)
```

Arguments

ohlc	An <i>OHLC</i> time series of prices in <i>xts</i> format.
method	A <i>character</i> string representing the method for estimating variance. The methods include: <ul style="list-style-type: none"> "close" close to close, "garman_klass" Garman-Klass, "garman_klass_yz" Garman-Klass with account for close-to-open price jumps, "rogers_satchell" Rogers-Satchell, "yang_zhang" Yang-Zhang, (default is "yang_zhang")
scalit	<i>Boolean</i> argument: should the returns be divided by the number of seconds in each period? (default is TRUE)

Details

The function `calc_var_ohlc_r()` calculates the variance from all the different intra-day and day-over-day returns (defined as the differences of *OHLC* prices), using several different variance estimation methods.

The default method is "yang_zhang", which theoretically has the lowest standard error among unbiased estimators. The methods "close", "garman_klass_yz", and "yang_zhang" do account for close-to-open price jumps, while the methods "garman_klass" and "rogers_satchell" do not account for close-to-open price jumps.

If `scalit` is `TRUE` (the default), then the returns are divided by the differences of the time index (which scales the variance to the units of variance per second squared.) This is useful when calculating the variance from minutely bar data, because dividing returns by the number of seconds decreases the effect of overnight price jumps. If the time index is in days, then the variance is equal to the variance per day squared.

The function `calc_var_ohlc_r()` is implemented in R code.

Value

A single *numeric* value equal to the variance.

Examples

```
# Calculate the variance of SPY returns
HighFreq::calc_var_ohlc_r(HighFreq::SPY, method="yang_zhang")
# Calculate variance without accounting for overnight jumps
HighFreq::calc_var_ohlc_r(HighFreq::SPY, method="rogers_satchell")
# Calculate the variance without scaling the returns
HighFreq::calc_var_ohlc_r(HighFreq::SPY, scalit=FALSE)
```

calc_weights	<i>Calculate the optimal portfolio weights using a variety of different objective functions.</i>
--------------	--

Description

Calculate the optimal portfolio weights using a variety of different objective functions.

Usage

```
calc_weights(
  returns,
  method = "maxsharpe",
  eigen_thresh = 1e-05,
  eigen_max = 0L,
  confl = 0.1,
  alpha = 0,
  rankw = FALSE,
  centerw = FALSE,
  scalew = "voltarget",
  vol_target = 0.01
)
```

Arguments

returns	A <i>time series</i> or a <i>matrix</i> of returns data (the returns in excess of the risk-free rate).
method	A <i>string</i> specifying the method for calculating the weights (see Details) (the default is <code>method = "sharpe"</code>)
eigen_thresh	A <i>numeric</i> threshold level for discarding small singular values in order to regularize the inverse of the covariance matrix of returns (the default is <code>1e-5</code>).

eigen_max	An <i>integer</i> equal to the number of singular values used for calculating the regularized inverse of the covariance matrix of returns (the default is 0 - equivalent to eigen_max equal to the number of columns of returns).
confl	The confidence level for calculating the quantiles of returns (the default is confl = 0.75).
alpha	The shrinkage intensity of returns. (values between 0 and 1 - the default is 0).
rankw	A <i>Boolean</i> specifying whether the weights should be ranked (the default is rankw = FALSE).
centerw	A <i>Boolean</i> specifying whether the weights should be centered (the default is centerw = FALSE).
scalew	A <i>string</i> specifying the method for scaling the weights (the default is scalew = "voltarget").
vol_target	A <i>numeric</i> volatility target for scaling the weights (the default is 0.001)

Details

The function `calc_weights()` calculates the optimal portfolio weights using a variety of different objective functions.

If `method = "maxsharpe"` (the default) then `calc_weights()` calculates the weights of the maximum Sharpe portfolio, by multiplying the regularized inverse of the *covariance matrix* C^{-1} times the mean column returns μ :

$$w = C^{-1}\mu$$

If `method = "maxsharpemed"` then `calc_weights()` uses the medians instead of the means.

If `method = "minvarlin"` then it calculates the weights of the minimum variance portfolio under linear constraint, by multiplying the regularized inverse of the *covariance matrix* times the unit vector:

$$w = C^{-1}1$$

If `method = "minvarquad"` then it calculates the weights of the minimum variance portfolio under quadratic constraint (which is the highest order principal component).

If `method = "sharpem"` then it calculates the momentum weights equal to the Sharpe ratios (the returns divided by their standard deviations):

$$w = \frac{\mu}{\sigma}$$

If `method = "kellym"` then it calculates the momentum weights equal to the Kelly ratios (the returns divided by their variance):

$$w = \frac{\mu}{\sigma^2}$$

`calc_weights()` calls the function `calc_inv()` to calculate the regularized inverse of the *covariance matrix* of returns. It performs regularization by selecting only the largest eigenvalues equal in number to `eigen_max`.

In addition, `calc_weights()` applies shrinkage to the columns of returns, by shrinking their means to their common mean value:

$$r'_i = (1 - \alpha)r_i + \alpha\bar{r}$$

Where r_i is the mean of column i and \bar{r} is the mean of all the columns. The shrinkage intensity alpha determines the amount of shrinkage that is applied, with `alpha = 0` representing no shrinkage

(with the column means r_i unchanged), and $\alpha = 1$ representing complete shrinkage (with the column means all equal to the single mean of all the columns: $r_i = \bar{r}$).

After the weights are calculated, they are scaled, depending on several arguments.

If `rankw = TRUE` then the weights are converted into their ranks. The default is `rankw = FALSE`.

If `centerw = TRUE` then the weights are centered so that their sum is equal to 0. The default is `centerw = FALSE`.

If `scalew = "voltarget"` (the default) then the weights are scaled (multiplied by a factor) so that the weighted portfolio has an in-sample volatility equal to `vol_target`.

If `scalew = "voleqw"` then the weights are scaled so that the weighted portfolio has the same volatility as the equal weight portfolio.

If `scalew = "sumone"` then the weights are scaled so that their sum is equal to 1. If `scalew = "sumsq"` then the weights are scaled so that their sum of squares is equal to 1. If `scalew = "none"` then the weights are not scaled.

The function `calc_weights()` is written in RcppArmadillo C++ code.

Value

A column *vector* of the same length as the number of columns of returns.

Examples

```
## Not run:
# Calculate covariance matrix and eigen decomposition of ETF returns
returns <- na.omit(rutils::etfenv$returns[, 1:16])
ncols <- NCOL(returns)
eigend <- eigen(cov(returns))
# Calculate regularized inverse of covariance matrix
eigen_max <- 3
eigenvec <- eigend$vectors[, 1:eigen_max]
eigenval <- eigend$values[1:eigen_max]
invmat <- eigenvec %*% (t(eigenvec) / eigenval)
# Define shrinkage intensity and apply shrinkage to the mean returns
alpha <- 0.5
colmeans <- colMeans(returns)
colmeans <- ((1-alpha)*colmeans + alpha*mean(colmeans))
# Calculate weights using R
weightsr <- drop(invmat %*% colmeans)
weightsr <- weightsr*sd(rowMeans(returns))/sd(returns %*% weightsr)
weightsr <- 0.01*weightsr/sd(returns %*% weightsr)
weightsr <- weightsr/sqrt(sum(weightsr^2))
# Calculate weights using RcppArmadillo
weightcpp <- drop(HighFreq::calc_weights(returns, eigen_max=eigen_max, alpha=alpha, scalew="sumsq"))
all.equal(weightcpp, weightsr)

## End(Not run)
```

diffit	Calculate the row differences of a time series or a matrix using RcppArmadillo.
--------	---

Description

Calculate the row differences of a *time series* or a *matrix* using *RcppArmadillo*.

Usage

```
diffit(tseries, lagg = 1L, pad_zeros = TRUE)
```

Arguments

tseries	A <i>time series</i> or a <i>matrix</i> .
lagg	An <i>integer</i> equal to the number of rows (time periods) to lag when calculating the differences (the default is lagg = 1).
pad_zeros	<i>Boolean</i> argument: Should the output <i>matrix</i> be padded (extended) with zero values, in order to return a <i>matrix</i> with the same number of rows as the input? (the default is pad_zeros = TRUE)

Details

The function `diffit()` calculates the differences between the rows of the input *matrix* `tseries` and its lagged version.

The argument `lagg` specifies the number of lags applied to the rows of the lagged version of `tseries`. For positive `lagg` values, the lagged version of `tseries` has its rows shifted *forward* (down) by the number equal to `lagg` rows. For negative `lagg` values, the lagged version of `tseries` has its rows shifted *backward* (up) by the number equal to `-lagg` rows. For example, if `lagg=3` then the lagged version will have its rows shifted down by 3 rows, and the differences will be taken between each row minus the row three time periods before it (in the past). The default is `lagg = 1`.

The argument `pad_zeros` specifies whether the output *matrix* should be padded (extended) with zero values in order to return a *matrix* with the same number of rows as the input `tseries`. The default is `pad_zeros = TRUE`. If `pad_zeros = FALSE` then the return *matrix* has a smaller number of rows than the input `tseries`. The padding operation can be time-consuming, because it requires the copying of data.

The function `diffit()` is implemented in *RcppArmadillo* C++ code, which makes it much faster than R code.

Value

A *matrix* containing the differences between the rows of the input *matrix* `tseries`.

Examples

```
## Not run:
# Create a matrix of random data
datav <- matrix(sample(15), nc=3)
# Calculate differences with lagged rows
HighFreq::diffit(datav, lagg=2)
```

```
# Calculate differences with advanced rows
HighFreq::diffit(datav, lagg=-2)
# Compare HighFreq::diffit() with rutils::diffit()
all.equal(HighFreq::diffit(datav, lagg=2),
  rutils::diffit(datav, lagg=2),
  check.attributes=FALSE)
# Compare the speed of RcppArmadillo with R code
library(microbenchmark)
summary(microbenchmark(
  Rcpp=HighFreq::diffit(datav, lagg=2),
  Rcode=rutils::diffit(datav, lagg=2),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary

## End(Not run)
```

diff_vec	<i>Calculate the differences between the neighboring elements of a single-column time series or a vector.</i>
----------	---

Description

Calculate the differences between the neighboring elements of a single-column *time series* or a *vector*.

Usage

```
diff_vec(tseries, lagg = 1L, pad_zeros = TRUE)
```

Arguments

tseries	A single-column <i>time series</i> or a <i>vector</i> .
lagg	An <i>integer</i> equal to the number of time periods to lag when calculating the differences (the default is lagg = 1).
pad_zeros	<i>Boolean</i> argument: Should the output <i>vector</i> be padded (extended) with zeros, in order to return a <i>vector</i> of the same length as the input? (the default is pad_zeros = TRUE)

Details

The function `diff_vec()` calculates the differences between the input *time series* or *vector* and its lagged version.

The argument `lagg` specifies the number of lags. For example, if `lagg=3` then the differences will be taken between each element minus the element three time periods before it (in the past). The default is `lagg = 1`.

The argument `pad_zeros` specifies whether the output *vector* should be padded (extended) with zeros at the front, in order to return a *vector* of the same length as the input. The default is `pad_zeros = TRUE`. The padding operation can be time-consuming, because it requires the copying of data.

The function `diff_vec()` is implemented in RcppArmadillo C++ code, which makes it several times faster than R code.

Value

A column *vector* containing the differences between the elements of the input vector.

Examples

```
## Not run:
# Create a vector of random returns
returns <- rnorm(1e6)
# Compare diff_vec() with rutils::diffit()
all.equal(drop(HighFreq::diff_vec(returns, lagg=3, pad=TRUE)),
  rutils::diffit(returns, lagg=3))
# Compare the speed of RcppArmadillo with R code
library(microbenchmark)
summary(microbenchmark(
  Rcpp=HighFreq::diff_vec(returns, lagg=3, pad=TRUE),
  Rcode=rutils::diffit(returns, lagg=3),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary

## End(Not run)
```

hf_data

High frequency data sets

Description

hf_data.RData is a file containing the datasets:

SPY an xts time series containing 1-minute OHLC bar data for the SPY etf, from 2008-01-02 to 2014-05-19. SPY contains 625,425 rows of data, each row contains a single minute bar.

TLT an xts time series containing 1-minute OHLC bar data for the TLT etf, up to 2014-05-19.

VXX an xts time series containing 1-minute OHLC bar data for the VXX etf, up to 2014-05-19.

Usage

```
data(hf_data) # not required - data is lazy load
```

Format

Each xts time series contains OHLC data, with each row containing a single minute bar:

Open Open price in the bar

High High price in the bar

Low Low price in the bar

Close Close price in the bar

Volume trading volume in the bar

Source

<https://wrds-web.wharton.upenn.edu/wrds/>

References

Wharton Research Data Service ([WRDS](#))

Examples

```
# data(hf_data) # not required - data is lazy load
head(SPY)
chart_Series(x=SPY["2009"])
```

lagit	<i>Apply a lag to the rows of a time series or a matrix using RcppArmadillo.</i>
-------	--

Description

Apply a lag to the rows of a *time series* or a *matrix* using RcppArmadillo.

Usage

```
lagit(tseries, lagg = 1L, pad_zeros = TRUE)
```

Arguments

tseries	<i>A time series or a matrix.</i>
lagg	<i>An integer equal to the number of periods to lag (the default is lagg = 1).</i>
pad_zeros	<i>Boolean argument: Should the output be padded with zeros? (The default is pad_zeros = TRUE.)</i>

Details

The function `lagit()` applies a lag to the input *matrix* by shifting its rows by the number equal to the argument `lagg`. For positive `lagg` values, the rows are shifted *forward* (down), and for negative `lagg` values they are shifted *backward* (up).

The output *matrix* is padded with either zeros (the default), or with rows of data from `tseries`, so that it has the same dimensions as `tseries`. If the `lagg` is positive, then the first row is copied and added upfront. If the `lagg` is negative, then the last row is copied and added to the end.

As a rule, if `tseries` contains returns data, then the output *matrix* should be padded with zeros, to avoid data snooping. If `tseries` contains prices, then the output *matrix* should be padded with the prices.

Value

A *matrix* with the same dimensions as the input argument `tseries`.

Examples

```
## Not run:
# Create a matrix of random returns
returns <- matrix(rnorm(5e6), nc=5)
# Compare lagit() with rutils::lagit()
all.equal(HighFreq::lagit(returns), rutils::lagit(returns))
# Compare the speed of RcppArmadillo with R code
library(microbenchmark)
summary(microbenchmark(
  Rcpp=HighFreq::lagit(returns),
  Rcode=rutils::lagit(returns),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary

## End(Not run)
```

lag_vec	Apply a lag to a single-column <i>time series</i> or a <i>vector</i> using RcppArmadillo.
---------	---

Description

Apply a lag to a single-column *time series* or a *vector* using RcppArmadillo.

Usage

```
lag_vec(tseries, lagg = 1L, pad_zeros = TRUE)
```

Arguments

tseries	A single-column <i>time series</i> or a <i>vector</i> .
lagg	An <i>integer</i> equal to the number of periods to lag. (The default is lagg = 1.)
pad_zeros	<i>Boolean</i> argument: Should the output be padded with zeros? (The default is pad_zeros = TRUE.)

Details

The function `lag_vec()` applies a lag to the input *time series* `tseries` by shifting its elements by the number equal to the argument `lagg`. For positive `lagg` values, the elements are shifted forward in time (down), and for negative `lagg` values they are shifted backward (up).

The output *vector* is padded with either zeros (the default), or with data from `tseries`, so that it has the same number of element as `tseries`. If the `lagg` is positive, then the first element is copied and added upfront. If the `lagg` is negative, then the last element is copied and added to the end.

As a rule, if `tseries` contains returns data, then the output *matrix* should be padded with zeros, to avoid data snooping. If `tseries` contains prices, then the output *matrix* should be padded with the prices.

Value

A column *vector* with the same number of elements as the input time series.

Examples

```
## Not run:
# Create a vector of random returns
returns <- rnorm(1e6)
# Compare lag_vec() with rutils::lagit()
all.equal(drop(HighFreq::lag_vec(returns)),
  rutils::lagit(returns))
# Compare the speed of RcppArmadillo with R code
library(microbenchmark)
summary(microbenchmark(
  Rcpp=HighFreq::lag_vec(returns),
  Rcode=rutils::lagit(returns),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary

## End(Not run)
```

lik_garch	<i>Calculate the log-likelihood of a time series of returns assuming a GARCH(1,1) process.</i>
-----------	--

Description

Calculate the log-likelihood of a time series of returns assuming a *GARCH(1,1)* process.

Usage

```
lik_garch(omega, alpha, beta, returns, minval = 1e-06)
```

Arguments

omega	Parameter proportional to the long-term average level of variance.
alpha	The weight associated with recent realized variance updates.
beta	The weight associated with the past variance estimates.
returns	A single-column <i>matrix</i> of returns.
minval	The floor value applied to the variance, to avoid zero values. (The default is minval = 0.000001.)

Details

The function `lik_garch()` calculates the log-likelihood of a time series of returns assuming a *GARCH(1,1)* process.

It first estimates the rolling variance of the returns argument using function `sim_garch()`:

$$\sigma_i^2 = \omega + \alpha r_i^2 + \beta \sigma_{i-1}^2$$

Where r_i is the time series of returns, and σ_i^2 is the estimated rolling variance. And ω , α , and β are the *GARCH* parameters. It applies the floor value `minval` to the variance, to avoid zero values. So the minimum value of the variance is equal to `minval`.

The function `lik_garch()` calculates the log-likelihood assuming a normal distribution of returns conditional on the variance σ_{i-1}^2 in the previous period, as follows:

$$likelihood = - \sum_{i=1}^n \left(\frac{r_i^2}{\sigma_{i-1}^2} + \log(\sigma_{i-1}^2) \right)$$

Value

The log-likelihood value.

Examples

```
## Not run:
# Define the GARCH model parameters
alpha <- 0.79
betav <- 0.2
om_ega <- 1e-4*(1-alpha-betav)
# Calculate historical VTI returns
returns <- na.omit(rutils::etfenv$returns$VTI)
# Calculate the log-likelihood of VTI returns assuming GARCH(1,1)
HighFreq::lik_garch(omega=om_ega, alpha=alpha, beta=betav, returns=returns)

## End(Not run)
```

mult_mat

Multiply the rows or columns of a matrix times a vector, element-wise.

Description

Multiply the rows or columns of a *matrix* times a *vector*, element-wise.

Usage

```
mult_mat(vector, matrix, byrow = TRUE)
```

Arguments

vector	A <i>numeric vector</i> .
matrix	A <i>numeric matrix</i> .
byrow	A <i>Boolean</i> argument: if TRUE then multiply the rows of matrix by vector, otherwise multiply the columns (the default is byrow = TRUE.)

Details

The function `mult_mat()` multiplies the rows or columns of a *matrix* times a *vector*, element-wise.

If `byrow = TRUE` (the default), then function `mult_mat()` multiplies the rows of the argument *matrix* times the argument *vector*. Otherwise it multiplies the columns of *matrix*.

In R, *matrix* multiplication is performed by columns. Performing multiplication by rows is often required, for example when multiplying stock returns by portfolio weights. But performing multiplication by rows requires explicit loops in R, or it requires *matrix* transpose. And both are slow.

The function `mult_mat()` uses RcppArmadillo C++ code, so when multiplying large *matrix* columns it's several times faster than vectorized R code, and it's even much faster compared to R when multiplying the *matrix* rows.

The function `mult_mat()` performs loops over the *matrix* rows and columns using the *Armadillo* operators `each_row()` and `each_col()`, instead of performing explicit `for()` loops (both methods are equally fast).

Value

A *matrix* equal to the product of the arguments *matrix* times *vector*, with the same dimensions as the argument *matrix*.

Examples

```
## Not run:
# Create vector and matrix data
matrixv <- matrix(round(runif(25e4), 2), nc=5e2)
vectorv <- round(runif(5e2), 2)

# Multiply the matrix rows using R
matrixr <- t(vectorv*t(matrixv))
# Multiply the matrix rows using C++
matrixp <- HighFreq::mult_mat(vectorv, matrixv, byrow=TRUE)
all.equal(matrixr, matrixp)
# Compare the speed of Rcpp with R code
library(microbenchmark)
summary(microbenchmark(
  Rcpp=HighFreq::mult_mat(vectorv, matrixv, byrow=TRUE),
  Rcode=t(vectorv*t(matrixv)),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary

# Multiply the matrix columns using R
matrixr <- vectorv*matrixv
# Multiply the matrix columns using C++
matrixp <- HighFreq::mult_mat(vectorv, matrixv, byrow=FALSE)
all.equal(matrixr, matrixp)
# Compare the speed of Rcpp with R code
library(microbenchmark)
summary(microbenchmark(
  Rcpp=HighFreq::mult_mat(vectorv, matrixv, byrow=FALSE),
  Rcode=vectorv*matrixv,
  times=10))[, c(1, 4, 5)] # end microbenchmark summary

## End(Not run)
```

mult_mat_ref

Multiply the rows or columns of a matrix times a vector, element-wise and in place (without copying).

Description

Multiply the rows or columns of a *matrix* times a *vector*, element-wise and in place (without copying).

Usage

```
mult_mat_ref(vector, matrix, byrow = TRUE)
```

Arguments

vector	A <i>numeric vector</i> .
matrix	A <i>numeric matrix</i> .
byrow	A <i>Boolean</i> argument: if TRUE then multiply the rows of matrix by vector, otherwise multiply the columns (the default is byrow = TRUE.)

Details

The function `mult_mat_ref()` multiplies the rows or columns of a *matrix* times a *vector*, element-wise and in place (without copying).

It accepts a *pointer* to the argument *matrix*, and replaces the old *matrix* values with the new values. It performs the calculation in place, without copying the *matrix* in memory, which can significantly increase the computation speed for large matrices.

If `byrow = TRUE` (the default), then function `mult_mat_ref()` multiplies the rows of the argument *matrix* times the argument vector. Otherwise it multiplies the columns of *matrix*.

In R, *matrix* multiplication is performed by columns. Performing multiplication by rows is often required, for example when multiplying stock returns by portfolio weights. But performing multiplication by rows requires explicit loops in R, or it requires *matrix* transpose. And both are slow.

The function `mult_mat_ref()` uses RcppArmadillo C++ code, so when multiplying large *matrix* columns it's several times faster than vectorized R code, and it's even much faster compared to R when multiplying the *matrix* rows.

The function `mult_mat_ref()` performs loops over the *matrix* rows and columns using the *Armadillo* operators `each_row()` and `each_col()`, instead of performing explicit `for()` loops (both methods are equally fast).

Value

Void (no return value).

Examples

```
## Not run:
# Create vector and matrix data
matrixv <- matrix(round(runif(25e4), 2), nc=5e2)
vectorv <- round(runif(5e2), 2)

# Multiply the matrix rows using R
matrixr <- t(vectorv*t(matrixv))
# Multiply the matrix rows using C++
HighFreq::mult_mat_ref(vectorv, matrixv, byrow=TRUE)
all.equal(matrixr, matrixv)
# Compare the speed of Rcpp with R code
library(microbenchmark)
summary(microbenchmark(
  Rcpp=HighFreq::mult_mat_ref(vectorv, matrixv, byrow=TRUE),
  Rcode=t(vectorv*t(matrixv)),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary

# Multiply the matrix columns using R
matrixr <- vectorv*matrixv
# Multiply the matrix columns using C++
HighFreq::mult_mat_ref(vectorv, matrixv, byrow=FALSE)
```

```

all.equal(matrixr, matrixv)
# Compare the speed of Rcpp with R code
library(microbenchmark)
summary(microbenchmark(
  Rcpp=HighFreq::mult_mat_ref(vectorv, matrixv, byrow=FALSE),
  Rcode=vectorv*matrixv,
  times=10))[, c(1, 4, 5)] # end microbenchmark summary

## End(Not run)

```

ohlc_returns	Calculate single period percentage returns from either <i>TAQ</i> or <i>OHLC</i> prices.
--------------	--

Description

Calculate single period percentage returns from either *TAQ* or *OHLC* prices.

Usage

```
ohlc_returns(xtsv, lagg = 1, colnum = 4, scalit = TRUE)
```

Arguments

xtsv	An <i>xts</i> time series of either <i>TAQ</i> or <i>OHLC</i> data.
lagg	An integer equal to the number of time periods of lag. (default is 1)
colnum	The column number to extract from the <i>OHLC</i> data. (default is 4, or the <i>Close</i> prices column)
scalit	<i>Boolean</i> argument: should the returns be divided by the number of seconds in each period? (default is TRUE)

Details

The function `ohlc_returns()` calculates the percentage returns for either *TAQ* or *OHLC* data, defined as the difference of log prices. Multi-period returns can be calculated by setting the `lag` parameter to values greater than 1 (the default).

If `scalit` is TRUE (the default), then the returns are divided by the differences of the time index (which scales the returns to units of returns per second.)

The time index of the `xtsv` time series is assumed to be in *POSIXct* format, so that its internal value is equal to the number of seconds that have elapsed since the *epoch*.

If `scalit` is TRUE (the default), then the returns are expressed in the scale of the time index of the `xtsv` time series. For example, if the time index is in seconds, then the returns are given in units of returns per second. If the time index is in days, then the returns are equal to the returns per day.

The function `ohlc_returns()` identifies the `xtsv` time series as *TAQ* data when it has six columns, otherwise assumes it's *OHLC* data. By default, for *OHLC* data, it differences the *Close* prices, but can also difference other prices depending on the value of `colnum`.

Value

A single-column *xts* time series of returns.

Examples

```
# Calculate secondly returns from TAQ data
returns <- HighFreq::ohlc_returns(xtsv=HighFreq::SPY_TAQ)
# Calculate close to close returns
returns <- HighFreq::ohlc_returns(xtsv=HighFreq::SPY)
# Calculate open to open returns
returns <- HighFreq::ohlc_returns(xtsv=HighFreq::SPY, colnum=1)
```

ohlc_sharpe	<i>Calculate time series of point Sharpe-like statistics for each row of a OHLC time series.</i>
-------------	--

Description

Calculate time series of point Sharpe-like statistics for each row of a *OHLC* time series.

Usage

```
ohlc_sharpe(ohlc, method = "close")
```

Arguments

ohlc	An <i>OHLC</i> time series of prices in <i>xts</i> format.
method	A <i>character</i> string representing method for estimating the Sharpe-like exponent.

Details

The function `ohlc_sharpe()` calculates Sharpe-like statistics for each row of a *OHLC* time series. The Sharpe-like statistic is defined as the ratio of the difference between *Close* minus *Open* prices divided by the difference between *High* minus *Low* prices. This statistic may also be interpreted as something like a *Hurst exponent* for a single row of data. The motivation for the Sharpe-like statistic is the notion that if prices are trending in the same direction inside a given time bar of data, then this statistic is close to either 1 or -1.

Value

An *xts* time series with the same number of rows as the argument `ohlc`.

Examples

```
# Calculate time series of running Sharpe ratios for SPY
sharpe_running <- ohlc_sharpe(HighFreq::SPY)
```

ohlc_skew	<i>Calculate time series of point skew estimates from a OHLC time series, assuming zero drift.</i>
-----------	--

Description

Calculate time series of point skew estimates from a *OHLC* time series, assuming zero drift.

Usage

```
ohlc_skew(ohlc, method = "rogers_satchell")
```

Arguments

ohlc	An <i>OHLC</i> time series of prices in <i>xts</i> format.
method	A character string representing method for estimating skew.

Details

The function `ohlc_skew()` calculates a time series of skew estimates from *OHLC* prices, one for each row of *OHLC* data. The skew estimates are expressed in the time scale of the index of the *OHLC* time series. For example, if the time index is in seconds, then the skew is given in units of skew per second. If the time index is in days, then the skew is equal to the skew per day.

Currently only the "close" skew estimation method is correct (assuming zero drift), while the "rogers_satchell" method produces a skew-like indicator, proportional to the skew. The default method is "rogers_satchell".

Value

A time series of point skew estimates.

Examples

```
# Calculate time series of skew estimates for SPY
skew <- HighFreq::ohlc_skew(HighFreq::SPY)
```

ohlc_variance	<i>Calculate a time series of point estimates of variance for an OHLC time series, using different range estimators for variance.</i>
---------------	---

Description

Calculates the point variance estimates from individual rows of *OHLC* prices (rows of data), using the squared differences of *OHLC* prices at each point in time, without averaging them over time.

Usage

```
ohlc_variance(ohlc, method = "yang_zhang", scalit = TRUE)
```

Arguments

<code>ohlc</code>	An <i>OHLC</i> time series of prices in <i>xts</i> format.
<code>method</code>	A <i>character</i> string representing the method for estimating variance. The methods include: <ul style="list-style-type: none"> "close" close to close, "garman_klass" Garman-Klass, "garman_klass_yz" Garman-Klass with account for close-to-open price jumps, "rogers_satchell" Rogers-Satchell, "yang_zhang" Yang-Zhang, (default is "yang_zhang")
<code>scalit</code>	<i>Boolean</i> argument: should the returns be divided by the number of seconds in each period? (default is TRUE)

Details

The function `ohlc_variance()` calculates a time series of point variance estimates of percentage returns, from *OHLC* prices, without averaging them over time. For example, the method "close" simply calculates the squares of the differences of the log *Close* prices.

The other methods calculate the squares of other possible differences of the log *OHLC* prices. This way the point variance estimates only depend on the price differences within individual rows of data (and possibly from the neighboring rows.) All the methods are implemented assuming zero drift, since the calculations are performed only for a single row of data, at a single point in time.

The user can choose from several different variance estimation methods. The methods "close", "garman_klass_yz", and "yang_zhang" do account for close-to-open price jumps, while the methods "garman_klass" and "rogers_satchell" do not account for close-to-open price jumps. The default method is "yang_zhang", which theoretically has the lowest standard error among unbiased estimators.

The point variance estimates can be passed into function `roll_vwap()` to perform averaging, to calculate rolling variance estimates. This is appropriate only for the methods "garman_klass" and "rogers_satchell", since they don't require subtracting the rolling mean from the point variance estimates.

The point variance estimates can also be considered to be technical indicators, and can be used as inputs into trading models.

If `scalit` is TRUE (the default), then the variance is divided by the squared differences of the time index (which scales the variance to units of variance per second squared.) This is useful for example, when calculating intra-day variance from minutely bar data, because dividing returns by the number of seconds decreases the effect of overnight price jumps.

If `scalit` is TRUE (the default), then the variance is expressed in the scale of the time index of the *OHLC* time series. For example, if the time index is in seconds, then the variance is given in units of variance per second squared. If the time index is in days, then the variance is equal to the variance per day squared.

The time index of the `ohlc` time series is assumed to be in *POSIXct* format, so that its internal value is equal to the number of seconds that have elapsed since the *epoch*.

The function `ohlc_variance()` performs similar calculations to the function `volatility()` from package **TTR**, but it assumes zero drift, and doesn't calculate a running sum using `runSum()`. It's also a little faster because it performs less data validation.

Value

An *xts* time series with a single column and the same number of rows as the argument *ohlc*.

Examples

```
# Create minutely OHLC time series of random prices
ohlc <- HighFreq::random_ohlc()
# Calculate variance estimates for ohlc
var_running <- HighFreq::ohlc_variance(ohlc)
# Calculate variance estimates for SPY
var_running <- HighFreq::ohlc_variance(HighFreq::SPY, method="yang_zhang")
# Calculate SPY variance without overnight jumps
var_running <- HighFreq::ohlc_variance(HighFreq::SPY, method="rogers_satchell")
```

random_ohlc	<i>Calculate a random OHLC time series of prices and trading volumes, in xts format.</i>
-------------	--

Description

Calculate a random *OHLC* time series either by simulating random prices following geometric Brownian motion, or by randomly sampling from an input time series.

Usage

```
random_ohlc(
  ohlc = NULL,
  reducit = TRUE,
  volat = 6.5e-05,
  drift = 0,
  indeks = seq(from = as.POSIXct(paste(Sys.Date() - 3, "09:30:00")), to =
    as.POSIXct(paste(Sys.Date() - 1, "16:00:00")), by = "1 sec"),
  ...
)
```

Arguments

ohlc	An <i>OHLC</i> time series of prices and trading volumes, in <i>xts</i> format (default is <i>NULL</i>).
volat	The volatility per period of the <i>indeks</i> time index (default is 6.5e-05 per second, or about 0.01=1.0% per day).
drift	The drift per period of the <i>indeks</i> time index (default is 0.0).
indeks	The time index for the <i>OHLC</i> time series.
reducit	<i>Boolean</i> argument: should <i>ohlc</i> time series be transformed to reduced form? (default is TRUE)

Details

If the input `ohlcv` time series is *NULL* (the default), then the function `random_ohlcv()` simulates a minutely *OHLC* time series of random prices following geometric Brownian motion, over the two previous calendar days.

If the input `ohlcv` time series is not *NULL*, then the rows of `ohlcv` are randomly sampled, to produce a random time series.

If `reducit` is *TRUE* (the default), then the `ohlcv` time series is first transformed to reduced form, then randomly sampled, and finally converted to standard form.

Note: randomly sampling from an intraday time series over multiple days will cause the overnight price jumps to be re-arranged into intraday price jumps. This will cause moment estimates to become inflated compared to the original time series.

Value

An *xts* time series with the same dimensions and the same time index as the input `ohlcv` time series.

Examples

```
# Create minutely synthetic OHLC time series of random prices
ohlcv <- HighFreq::random_ohlcv()
# Create random time series from SPY by randomly sampling it
ohlcv <- HighFreq::random_ohlcv(ohlcv=HighFreq::SPY["2012-02-13/2012-02-15"])
```

random_taq	<i>Calculate a random TAQ time series of prices and trading volumes, in xts format.</i>
------------	---

Description

Calculate a *TAQ* time series of random prices following geometric Brownian motion, combined with random trading volumes.

Usage

```
random_taq(
  volat = 6.5e-05,
  drift = 0,
  indeks = seq(from = as.POSIXct(paste(Sys.Date() - 3, "09:30:00")), to =
    as.POSIXct(paste(Sys.Date() - 1, "16:00:00")), by = "1 sec"),
  bid_offer = 0.001,
  ...
)
```

Arguments

<code>bid_offer</code>	The bid-offer spread expressed as a fraction of the prices (default is 0.001=10bps).
<code>volat</code>	The volatility per period of the <code>indeks</code> time index (default is 6.5e-05 per second, or about 0.01=1.0% per day).
<code>drift</code>	The drift per period of the <code>indeks</code> time index (default is 0.0).
<code>indeks</code>	The time index for the <i>TAQ</i> time series.

Details

The function `random_taq()` calculates an *xts* time series with four columns containing random prices following geometric Brownian motion: the bid, ask, and trade prices, combined with random trade volume data. If `indeks` isn't supplied as an argument, then by default it's equal to the secondly index over the two previous calendar days.

Value

An *xts* time series, with time index equal to the input `indeks` time index, and with four columns containing the bid, ask, and trade prices, and the trade volume.

Examples

```
# Create secondly TAQ time series of random prices
taq <- HighFreq::random_taq()
# Create random TAQ time series from SPY index
taq <- HighFreq::random_taq(indeks=index(HighFreq::SPY["2012-02-13/2012-02-15"]))
```

<code>remove_jumps</code>	<i>Remove overnight close-to-open price jumps from an OHLC time series, by adding adjustment terms to its prices.</i>
---------------------------	---

Description

Remove overnight close-to-open price jumps from an *OHLC* time series, by adding adjustment terms to its prices.

Usage

```
remove_jumps(ohlc)
```

Arguments

`ohlc` An *OHLC* time series of prices and trading volumes, in *xts* format.

Details

The function `remove_jumps()` removes the overnight close-to-open price jumps from an *OHLC* time series, by adjusting its prices so that the first *Open* price of the day is equal to the last *Close* price of the previous day.

The function `remove_jumps()` adds adjustment terms to all the *OHLC* prices, so that intra-day returns and volatilities are not affected.

The function `remove_jumps()` identifies overnight periods as those that are greater than 60 seconds. This assumes that intra-day periods between neighboring rows of data are 60 seconds or less.

The time index of the *ohlc* time series is assumed to be in *POSIXct* format, so that its internal value is equal to the number of seconds that have elapsed since the *epoch*.

Value

An *OHLC* time series with the same dimensions and the same time index as the input *ohlc* time series.

Examples

```
# Remove overnight close-to-open price jumps from SPY data
ohlcv <- remove_jumps(HighFreq::SPY)
```

roll_apply	<i>Apply an aggregation function over a rolling look-back interval and the end points of an OHLC time series, using R code.</i>
------------	---

Description

Apply an aggregation function over a rolling look-back interval and the end points of an *OHLC* time series, using R code.

Usage

```
roll_apply(
  xtsv,
  agg_fun,
  look_back = 2,
  endpoints = seq_along(xtsv),
  by_columns = FALSE,
  out_xts = TRUE,
  ...
)
```

Arguments

...	additional parameters to the function <code>agg_fun</code> .
xtsv	An <i>OHLC</i> time series of prices and trading volumes, in <i>xts</i> format.
agg_fun	The name of the aggregation function to be applied over a rolling look-back interval.
look_back	The number of end points in the look-back interval used for applying the aggregation function (including the current row).
by_columns	<i>Boolean</i> argument: should the function <code>agg_fun()</code> be applied column-wise (individually), or should it be applied to all the columns combined? (default is <code>FALSE</code>)
out_xts	<i>Boolean</i> argument: should the output be coerced into an <i>xts</i> series? (default is <code>TRUE</code>)
endpoints	An integer vector of end points.

Details

The function `roll_apply()` applies an aggregation function over a rolling look-back interval attached at the end points of an *OHLC* time series.

The function `roll_apply()` is implemented in R code.

`HighFreq::roll_apply()` performs similar operations to the functions `rollapply()` and `period.apply()` from package `xts`, and also the function `apply.rolling()` from package `PerformanceAnalytics`. (The function `rollapply()` isn't exported from the package `xts`.)

But `HighFreq::roll_apply()` is faster because it performs less type-checking and skips other overhead. Unlike the other functions, `roll_apply()` doesn't produce any leading *NA* values.

The function `roll_apply()` can be called in two different ways, depending on the argument `endpoints`. If the argument `endpoints` isn't explicitly passed to `roll_apply()`, then the default value is used, and `roll_apply()` performs aggregations over overlapping intervals at each point in time.

If the argument `endpoints` is explicitly passed to `roll_apply()`, then `roll_apply()` performs aggregations over intervals attached at the endpoints. If `look_back=2` then the aggregations are performed over non-overlapping intervals, otherwise they are performed over overlapping intervals.

If the argument `out_xts` is `TRUE` (the default) then the output is coerced into an *xts* series, with the number of rows equal to the length of argument `endpoints`. Otherwise a list is returned, with the length equal to the length of argument `endpoints`.

If `out_xts` is `TRUE` and the aggregation function `agg_fun()` returns a single value, then `roll_apply()` returns an *xts* time series with a single column. If `out_xts` is `TRUE` and if `agg_fun()` returns a vector of values, then `roll_apply()` returns an *xts* time series with multiple columns, equal to the length of the vector returned by the aggregation function `agg_fun()`.

Value

Either an *xts* time series with the number of rows equal to the length of argument `endpoints`, or a list the length of argument `endpoints`.

Examples

```
# extract a single day of SPY data
ohlc <- HighFreq::SPY["2012-02-13"]
interval <- 11 # number of data points between end points
look_back <- 4 # number of end points in look-back interval
# Calculate the rolling sums of ohlc columns over a rolling look-back interval
agg_regations <- roll_apply(ohlc, agg_fun=sum, look_back=look_back, by_columns=TRUE)
# Apply a vector-valued aggregation function over a rolling look-back interval
agg_function <- function(ohlc) c(max(ohlc[, 2]), min(ohlc[, 3]))
agg_regations <- roll_apply(ohlc, agg_fun=agg_function, look_back=look_back)
# Define end points at 11-minute intervals (HighFreq::SPY is minutely bars)
endpoints <- rutils::endpoints(ohlc, interval=interval)
# Calculate the sums of ohlc columns over endpoints using non-overlapping intervals
agg_regations <- roll_apply(ohlc, agg_fun=sum, endpoints=endpoints, by_columns=TRUE)
# Apply a vector-valued aggregation function over the endpoints of ohlc
# using overlapping intervals
agg_regations <- roll_apply(ohlc, agg_fun=agg_function,
                           look_back=5, endpoints=endpoints)
```

roll_backtest

Perform a backtest simulation of a trading strategy (model) over a vector of end points along a time series of prices.

Description

Perform a backtest simulation of a trading strategy (model) over a vector of end points along a time series of prices.

Usage

```
roll_backtest(
  xtsv,
  train_func,
  trade_func,
  look_back = look_forward,
  look_forward,
  endpoints = rutils::calc_endpoints(xtsv, look_forward),
  ...
)
```

Arguments

...	additional parameters to the functions <code>train_func()</code> and <code>trade_func()</code> .
<code>xtsv</code>	A time series of prices, asset returns, trading volumes, and other data, in <i>xts</i> format.
<code>train_func</code>	The name of the function for training (calibrating) a forecasting model, to be applied over a rolling look-back interval.
<code>trade_func</code>	The name of the trading model function, to be applied over a rolling look-forward interval.
<code>look_back</code>	The size of the look-back interval, equal to the number of rows of data used for training the forecasting model.
<code>look_forward</code>	The size of the look-forward interval, equal to the number of rows of data used for trading the strategy.
<code>endpoints</code>	A vector of end points along the rows of the <code>xtsv</code> time series, given as either integers or dates.

Details

The function `roll_backtest()` performs a rolling backtest simulation of a trading strategy over a vector of end points. At each end point, it trains (calibrates) a forecasting model using past data taken from the `xtsv` time series over the look-back interval, and applies the forecasts to the `trade_func()` trading model, using out-of-sample future data from the look-forward interval.

The function `trade_func()` should simulate the trading model, and it should return a named list with at least two elements: a named vector of performance statistics, and an *xts* time series of out-of-sample returns. The list returned by `trade_func()` can also have additional elements, like the in-sample calibrated model statistics, etc.

The function `roll_backtest()` returns a named list containing the listv returned by function `trade_func()`. The list names are equal to the *endpoints* dates. The number of list elements is equal to the number of *endpoints* minus two (because the first and last end points can't be included in the backtest).

Value

An *xts* time series with the number of rows equal to the number of end points minus two.

Examples

```
## Not run:
# Combine two time series of prices
```



```

prices <- cbind(rutils::etfenv$XLU, rutils::etfenv$XLP)
look_back <- 252
look_forward <- 22
# Define end points
endpoints <- rutils::calc_endpoints(prices, look_forward)
# Perform back-test
back_test <- roll_backtest(endpoints=endpoints,
  look_forward=look_forward,
  look_back=look_back,
  train_func = train_model,
  trade_func = trade_model,
  model_params = model_params,
  trading_params = trading_params,
  xtsv=prices)

## End(Not run)

```

roll_conv	<i>Calculate the rolling convolutions (weighted sums) of a time series with a single-column matrix of weights.</i>
-----------	--

Description

Calculate the rolling convolutions (weighted sums) of a *time series* with a single-column *matrix* of weights.

Usage

```
roll_conv(tseries, weights)
```

Arguments

tseries	A <i>time series</i> or a <i>matrix</i> of data.
weights	A single-column <i>matrix</i> of weights.

Details

The function `roll_conv()` calculates the convolutions of the *matrix* columns with a single-column *matrix* of weights. It performs a loop over the *matrix* rows and multiplies the past (higher) values by the weights. It calculates the rolling weighted sums of the past values.

The function `roll_conv()` uses the RcppArmadillo function `arma::conv2()`. It performs a similar calculation to the standard R function `filter(x=tseries, filter=weights, method="convolution", sides=1)`, but it's over 6 times faster, and it doesn't produce any leading NA values.

Value

A *matrix* with the same dimensions as the input argument `tseries`.

Examples

```
## Not run:
# First example
# Calculate a time series of returns
returns <- na.omit(rutils::etfenv$returns[, c("IEF", "VTI")])
# Create simple weights equal to a 1 value plus zeros
weights <- matrix(c(1, rep(0, 10)), nc=1)
# Calculate rolling weighted sums
weighted <- HighFreq::roll_conv(returns, weights)
# Compare with original
all.equal(coredata(returns), weighted, check.attributes=FALSE)
# Second example
# Calculate exponentially decaying weights
weights <- exp(-0.2*(1:11))
weights <- matrix(weights/sum(weights), nc=1)
# Calculate rolling weighted sums
weighted <- HighFreq::roll_conv(returns, weights)
# Calculate rolling weighted sums using filter()
filtered <- filter(x=returns, filter=weights, method="convolution", sides=1)
# Compare both methods
all.equal(filtered[-(1:11), ], weighted[-(1:11), ], check.attributes=FALSE)

## End(Not run)
```

roll_count	<i>Count the number of consecutive TRUE elements in a Boolean vector, and reset the count to zero after every FALSE element.</i>
------------	--

Description

Count the number of consecutive TRUE elements in a Boolean vector, and reset the count to zero after every FALSE element.

Usage

```
roll_count(tseries)
```

Arguments

tseries	<i>A Boolean vector of data.</i>
---------	----------------------------------

Details

The function `roll_count()` calculates the number of consecutive TRUE elements in a Boolean vector, and it resets the count to zero after every FALSE element.

For example, the Boolean vector FALSE, TRUE, TRUE, FALSE, FALSE, TRUE, TRUE, TRUE, TRUE, TRUE, FALSE, is translated into 0, 1, 2, 0, 0, 1, 2, 3, 4, 5, 0.

Value

An *integer vector* of the same length as the argument `tseries`.

Examples

```
## Not run:
# Calculate the number of consecutive TRUE elements
drop(HighFreq::roll_count(c(FALSE, TRUE, TRUE, FALSE, FALSE, TRUE, TRUE, TRUE, TRUE, TRUE, FALSE)))

## End(Not run)
```

roll_fun	<i>Calculate a matrix of estimator values over a rolling look-back interval attached at the end points of a time series or a matrix.</i>
----------	--

Description

Calculate a *matrix* of estimator values over a rolling look-back interval attached at the end points of a *time series* or a *matrix*.

Usage

```
roll_fun(
  tseries,
  fun = "calc_var",
  startp = 0L,
  endp = 0L,
  step = 1L,
  look_back = 1L,
  stub = 0L,
  method = "moment",
  confl = 0.75
)
```

Arguments

tseries	A <i>time series</i> or a <i>matrix</i> of data.
fun	A <i>string</i> specifying the estimator function (the default is fun = "calc_var".)
startp	An <i>integer</i> vector of start points (the default is startp = 0).
endp	An <i>integer</i> vector of end points (the default is endp = 0).
step	The number of time periods between the end points (the default is step = 1).
look_back	The number of end points in the look-back interval (the default is look_back = 1).
stub	An <i>integer</i> value equal to the first end point for calculating the end points (the default is stub = 0).
method	A <i>string</i> specifying the type of the model for the estimator (the default is method = "moment".)
confl	The confidence level for calculating the quantiles of returns (the default is confl = 0.75).

Details

The function `roll_fun()` calculates a *matrix* of estimator values, over rolling look-back intervals attached at the end points of the *time series* `tseries`.

The function `roll_fun()` performs a loop over the end points, and at each end point it subsets the time series `tseries` over a look-back interval equal to `look_back` number of end points.

It passes the subset time series to the function specified by the argument `fun`, which calculates the statistic. See the functions `calc_*`() for a description of the different estimators.

If the arguments `endp` and `startp` are not given then it first calculates a vector of end points separated by step time periods. It calculates the end points along the rows of `tseries` using the function `calc_endpoints()`, with the number of time periods between the end points equal to step time periods.

For example, the rolling variance at 25 day end points, with a 75 day look-back, can be calculated using the parameters `step = 25` and `look_back = 3`.

The function `roll_fun()` is implemented in RcppArmadillo C++ code, so it's many times faster than the equivalent R code.

Value

A *matrix* with the same number of columns as the input time series `tseries`, and the number of rows equal to the number of end points.

Examples

```
## Not run:
# Define time series of returns using package rutils
returns <- na.omit(rutils::etfenv$returns$VTI)
# Calculate the rolling variance at 25 day end points, with a 75 day look-back
var_rollfun <- HighFreq::roll_fun(returns, fun="calc_var", step=25, look_back=3)
# Calculate the rolling variance using roll_var()
var_roll <- HighFreq::roll_var(returns, step=25, look_back=3)
# Compare the two methods
all.equal(var_rollfun, var_roll, check.attributes=FALSE)
# Define end points and start points
endp <- HighFreq::calc_endpoints(NROW(returns), step=25)
startp <- HighFreq::calc_startpoints(endp, look_back=3)
# Calculate the rolling variance using RcppArmadillo
var_rollfun <- HighFreq::roll_fun(returns, fun="calc_var", startp=startp, endp=endp)
# Calculate the rolling variance using R code
var_roll <- sapply(1:NROW(endp), function(it) {
  var(returns[startp[it]:endp[it]+1, ])
}) # end sapply
# Compare the two methods
all.equal(drop(var_rollfun), var_roll, check.attributes=FALSE)
# Compare the speed of RcppArmadillo with R code
library(microbenchmark)
summary(microbenchmark(
  Rcpp=HighFreq::roll_fun(returns, fun="calc_var", startp=startp, endp=endp),
  Rcode=sapply(1:NROW(endp), function(it) {
    var(returns[startp[it]:endp[it]+1, ])
  }),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary

## End(Not run)
```

roll_hurst	Calculate a time series of <i>Hurst</i> exponents over a rolling look-back interval.
------------	--

Description

Calculate a time series of *Hurst* exponents over a rolling look-back interval.

Usage

```
roll_hurst(ohlc, look_back = 11)
```

Arguments

ohlc	An <i>OHLC</i> time series of prices in <i>xts</i> format.
look_back	The size of the look-back interval, equal to the number of rows of data used for aggregating the <i>OHLC</i> prices.

Details

The function `roll_hurst()` calculates a time series of *Hurst* exponents from *OHLC* prices, over a rolling look-back interval.

The *Hurst* exponent is defined as the logarithm of the ratio of the price range, divided by the standard deviation of returns, and divided by the logarithm of the interval length.

The function `roll_hurst()` doesn't use the same definition as the rescaled range definition of the *Hurst* exponent. First, because the price range is calculated using *High* and *Low* prices, which produces bigger range values, and higher *Hurst* exponent estimates. Second, because the *Hurst* exponent is estimated using a single aggregation interval, instead of multiple intervals in the rescaled range definition.

The rationale for using a different definition of the *Hurst* exponent is that it's designed to be a technical indicator for use as input into trading models, rather than an estimator for statistical analysis.

Value

An *xts* time series with a single column and the same number of rows as the argument `ohlc`.

Examples

```
# Calculate rolling Hurst for SPY in March 2009
hurst_rolling <- roll_hurst(ohlc=HighFreq::SPY["2009-03"], look_back=11)
chart_Series(hurst_rolling["2009-03-10/2009-03-12"], name="SPY hurst_rolling")
```

roll_kurtosis	Calculate a <i>matrix</i> of kurtosis estimates over a rolling look-back interval attached at the end points of a time series or a matrix.
---------------	--

Description

Calculate a *matrix* of kurtosis estimates over a rolling look-back interval attached at the end points of a *time series* or a *matrix*.

Usage

```
roll_kurtosis(
  tseries,
  startp = 0L,
  endp = 0L,
  step = 1L,
  look_back = 1L,
  stub = 0L,
  method = "moment",
  confl = 0.75
)
```

Arguments

tseries	A <i>time series</i> or a <i>matrix</i> of data.
startp	An <i>integer</i> vector of start points (the default is startp = 0).
endp	An <i>integer</i> vector of end points (the default is endp = 0).
step	The number of time periods between the end points (the default is step = 1).
look_back	The number of end points in the look-back interval (the default is look_back = 1).
stub	An <i>integer</i> value equal to the first end point for calculating the end points (the default is stub = 0).
method	A <i>string</i> specifying the type of the kurtosis model (the default is method = "moment" - see Details).
confl	The confidence level for calculating the quantiles of returns (the default is confl = 0.75).

Details

The function roll_kurtosis() calculates a *matrix* of kurtosis estimates over rolling look-back intervals attached at the end points of the *time series* tseries.

The function roll_kurtosis() performs a loop over the end points, and at each end point it subsets the time series tseries over a look-back interval equal to look_back number of end points.

It passes the subset time series to the function calc_kurtosis(), which calculates the kurtosis. See the function calc_kurtosis() for a description of the kurtosis methods.

If the arguments endp and startp are not given then it first calculates a vector of end points separated by step time periods. It calculates the end points along the rows of tseries using the function

calc_endpoints(), with the number of time periods between the end points equal to step time periods.

For example, the rolling kurtosis at 25 day end points, with a 75 day look-back, can be calculated using the parameters step = 25 and look_back = 3.

The function roll_kurtosis() is implemented in RcppArmadillo C++ code, so it's many times faster than the equivalent R code.

Value

A *matrix* of kurtosis estimates with the same number of columns as the input time series tseries, and the number of rows equal to the number of end points.

Examples

```
## Not run:
# Define time series of returns using package rutils
returns <- na.omit(rutils::etfenv$returns$VTI)
# Define end points and start points
endp <- 1 + HighFreq::calc_endpoints(NROW(returns), step=25)
startp <- HighFreq::calc_startpoints(endp, look_back=3)
# Calculate the rolling kurtosis at 25 day end points, with a 75 day look-back
kurtosisv <- HighFreq::roll_kurtosis(returns, step=25, look_back=3)
# Calculate the rolling kurtosis using R code
kurt_r <- sapply(1:NROW(endp), function(it) {
  HighFreq::calc_kurtosis(returns[startp[it]:endp[it], ])
}) # end sapply
# Compare the kurtosis estimates
all.equal(drop(kurtosisv), kurt_r, check.attributes=FALSE)
# Compare the speed of RcppArmadillo with R code
library(microbenchmark)
summary(microbenchmark(
  Rcpp=HighFreq::roll_kurtosis(returns, step=25, look_back=3),
  Rcode=sapply(1:NROW(endp), function(it) {
    HighFreq::calc_kurtosis(returns[startp[it]:endp[it], ])
  }),
  times=100))[, c(1, 4, 5)] # end microbenchmark summary

## End(Not run)
```

roll_mean

Calculate a matrix of mean (location) estimates over a rolling look-back interval attached at the end points of a time series or a matrix.

Description

Calculate a *matrix* of mean (location) estimates over a rolling look-back interval attached at the end points of a *time series* or a *matrix*.

Usage

```
roll_mean(
  tseries,
  startp = 0L,
```

```

endp = 0L,
step = 1L,
look_back = 1L,
stub = 0L,
method = "moment",
confl = 0.75
)

```

Arguments

tseries	A <i>time series</i> or a <i>matrix</i> of data.
startp	An <i>integer</i> vector of start points (the default is startp = 0).
endp	An <i>integer</i> vector of end points (the default is endp = 0).
step	The number of time periods between the end points (the default is step = 1).
look_back	The number of end points in the look-back interval (the default is look_back = 1).
stub	An <i>integer</i> value equal to the first end point for calculating the end points (the default is stub = 0).
method	A <i>character</i> string representing the type of mean measure of (the default is method = "moment").

Details

The function `roll_mean()` calculates a *matrix* of mean (location) estimates over rolling look-back intervals attached at the end points of the *time series* tseries.

The function `roll_mean()` performs a loop over the end points, and at each end point it subsets the time series tseries over a look-back interval equal to look_back number of end points.

It passes the subset time series to the function `calc_mean()`, which calculates the mean (location). See the function `calc_mean()` for a description of the mean methods.

If the arguments endp and startp are not given then it first calculates a vector of end points separated by step time periods. It calculates the end points along the rows of tseries using the function `calc_endpoints()`, with the number of time periods between the end points equal to step time periods.

For example, the rolling mean at 25 day end points, with a 75 day look-back, can be calculated using the parameters step = 25 and look_back = 3.

The function `roll_mean()` with the parameter step = 1 performs the same calculation as the function `roll_mean()` from package **RcppRoll**, but it's several times faster because it uses RcppArmadillo C++ code.

The function `roll_mean()` is implemented in RcppArmadillo C++ code, so it's many times faster than the equivalent R code.

If only a simple rolling mean is required (not the median) then other functions like `roll_sum()` or `roll_vec()` may be even faster.

Value

A *matrix* of mean (location) estimates with the same number of columns as the input time series tseries, and the number of rows equal to the number of end points.

Examples

```
## Not run:
# Define time series of returns using package rutils
returns <- na.omit(rutils::etfenv$returns$VTI)
# Calculate the rolling means at 25 day end points, with a 75 day look-back
means <- HighFreq::roll_mean(returns, step=25, look_back=3)
# Compare the mean estimates over 11-period look-back intervals
all.equal(HighFreq::roll_mean(returns, look_back=11)[-(1:10)],
  drop(RcppRoll::roll_mean(returns, n=11)), check.attributes=FALSE)
# Define end points and start points
endp <- HighFreq::calc_endpoints(NROW(returns), step=25)
startp <- HighFreq::calc_startpoints(endp, look_back=3)
# Calculate the rolling means using RcppArmadillo
means <- HighFreq::roll_mean(returns, startp=startp, endp=endp)
# Calculate the rolling medians using RcppArmadillo
medianscpp <- HighFreq::roll_mean(returns, startp=startp, endp=endp, method="nonparametric")
# Calculate the rolling medians using R
medians = sapply(1:NROW(endp), function(i) {
  median(returns[startp[i]:endp[i] + 1])
}) # end sapply
all.equal(medians, drop(medianscpp))
# Compare the speed of RcppArmadillo with R code
library(microbenchmark)
summary(microbenchmark(
  Rcpp=HighFreq::roll_mean(returns, startp=startp, endp=endp, method="nonparametric"),
  Rcode=sapply(1:NROW(endp), function(i) {median(returns[startp[i]:endp[i] + 1])}),
  times=10))[, c(1, 4, 5)]

## End(Not run)
```

roll_ohlc

Aggregate a time series to an OHLC time series with lower periodicity.

Description

Given a time series of prices at a higher periodicity (say seconds), it calculates the *OHLC* prices at a lower periodicity (say minutes).

Usage

```
roll_ohlc(tseries, endp)
```

Arguments

tseries	A <i>time series</i> or a <i>matrix</i> with multiple columns of data.
endp	An <i>integer vector</i> of end points.

Details

The function `roll_ohlc()` performs a loop over the end points *endp*, along the rows of the data *tseries*. At each end point, it selects the past rows of the data *tseries*, starting at the first bar after the previous end point, and then calls the function `agg_ohlc()` on the selected data *tseries* to calculate the aggregations.

The function `roll_ohlc()` can accept either a single column of data or four columns of *OHLC* data. It can also accept an additional column containing the trading volume.

The function `roll_ohlc()` performs a similar aggregation as the function `to.period()` from package *xts*.

Value

A *matrix* with *OHLC* data, with the number of rows equal to the number of *endp* minus one.

Examples

```
## Not run:
# Define matrix of OHLC data
ohlc <- rutils::etfenv$VTI[, 1:5]
# Define end points at 25 day intervals
endp <- HighFreq::calc_endpoints(NROW(ohlc), step=25)
# Aggregate over endp:
ohlcagg <- HighFreq::roll_ohlc(tseries=ohlc, endp=endp)
# Compare with xts::to.period()
ohlcagg_xts <- .Call("toPeriod", ohlc, as.integer(endp+1), TRUE, NCOL(ohlc), FALSE, FALSE, colnames(ohlc), PACKAGE="xts")
all.equal(ohlcagg, coredata(ohlcagg_xts), check.attributes=FALSE)

## End(Not run)
```

roll_reg	<i>Calculate a matrix of regression coefficients, their t-values, and z-scores, at the end points of the predictor matrix.</i>
----------	--

Description

Calculate a *matrix* of regression coefficients, their t-values, and z-scores, at the end points of the predictor matrix.

Usage

```
roll_reg(
  response,
  predictor,
  intercept = TRUE,
  startp = 0L,
  endp = 0L,
  step = 1L,
  look_back = 1L,
  stub = 0L,
  method = "least_squares",
  eigen_thresh = 1e-05,
  eigen_max = 0L,
  confl = 0.1,
  alpha = 0
)
```

Arguments

response	A single-column <i>time series</i> or a <i>vector</i> of response data.
predictor	A <i>time series</i> or a <i>matrix</i> of predictor data.
startp	An <i>integer</i> vector of start points (the default is startp = 0).
endp	An <i>integer</i> vector of end points (the default is endp = 0).
step	The number of time periods between the end points (the default is step = 1).
look_back	The number of end points in the look-back interval (the default is look_back = 1).
stub	An <i>integer</i> value equal to the first end point for calculating the end points (the default is stub = 0).
intercept	A <i>Boolean</i> specifying whether an intercept term should be added to the predictor (the default is intercept = TRUE).
method	A <i>string</i> specifying the type of the regression model the default is method = "least_squares" - see Details).
eigen_thresh	A <i>numeric</i> threshold level for discarding small singular values in order to regularize the inverse of the predictor matrix (the default is 1e-5).
eigen_max	An <i>integer</i> equal to the number of singular values used for calculating the regularized inverse of the predictor matrix (the default is 0 - equivalent to eigen_max equal to the number of columns of predictor).
confl	The confidence level for calculating the quantiles of returns (the default is confl = 0.75).
alpha	The shrinkage intensity between 0 and 1. (the default is 0).

Details

The function `roll_reg()` calculates a *matrix* of regression coefficients, their t-values, and z-scores at the end points of the predictor matrix.

The function `roll_reg()` performs a loop over the end points, and at each end point it subsets the time series predictor over a look-back interval equal to `look_back` number of end points.

If the arguments `endp` and `startp` are not given then it first calculates a vector of end points separated by `step` time periods. It calculates the end points along the rows of predictor using the function `calc_endpoints()`, with the number of time periods between the end points equal to `step` time periods.

For example, the rolling regression at 25 day end points, with a 75 day look-back, can be calculated using the parameters `step = 25` and `look_back = 3`.

It passes the subset time series to the function `calc_reg()`, which calculates the regression coefficients, their t-values, and the z-score.

If `intercept = TRUE` (the default) then an extra intercept column (unit column) is added to the predictor matrix.

The number of columns of the return matrix depends on the number of columns of the predictor matrix (including the intercept column, if it's added). The number of columns of the return matrix is equal to the number of regression coefficients, plus their t-values, plus the z-score column. The number of regression coefficients is equal to the number of columns of the predictor matrix (including the intercept column, if it's added). The number of t-values is equal to the number of coefficients. For example, if the number of columns of the predictor matrix is equal to n , and if `intercept = TRUE` (the default), then `roll_reg()` returns a matrix with $2n+3$ columns: $n+1$ regression coefficients (including the intercept coefficient), $n+1$ corresponding t-values, and 1 z-score column.

Value

A *matrix* with the regression coefficients, their t-values, and z-scores, and with the same number of rows as predictor a number of columns equal to $2n+3$, where n is the number of columns of predictor.

Examples

```
## Not run:
# Calculate historical returns
returns <- na.omit(rutils::etfenv$returns[, c("XLP", "VTI")])
# Define monthly end points and start points
endp <- xts::endpoints(returns, on="months")[-1]
look_back <- 12
startp <- c(rep(1, look_back), endp[1:(NROW(endp)-look_back)])
# Calculate rolling betas using RcppArmadillo
reg_stats <- HighFreq::roll_reg(response=returns[, 1], predictor=returns[, 2], endp=(endp-1), startp=(startp-1))
betas <- reg_stats[, 2]
# Calculate rolling betas in R
betas_r <- sapply(1:NROW(endp), FUN=function(ep) {
  datav <- returns[startp[ep]:endp[ep], ]
  drop(cov(datav[, 1], datav[, 2])/var(datav[, 2]))
}) # end sapply
# Compare the outputs of both functions
all.equal(betas, betas_r, check.attributes=FALSE)

## End(Not run)
```

roll_scale	<i>Perform a rolling scaling (standardization) of the columns of a matrix of data using RcppArmadillo.</i>
------------	--

Description

Perform a rolling scaling (standardization) of the columns of a *matrix* of data using RcppArmadillo.

Usage

```
roll_scale(matrix, look_back, use_median = FALSE)
```

Arguments

use_median	A <i>Boolean</i> argument: if TRUE then the centrality (central tendency) is calculated as the <i>median</i> and the dispersion is calculated as the <i>median absolute deviation (MAD)</i> . If use_median is FALSE then the centrality is calculated as the <i>mean</i> and the dispersion is calculated as the <i>standard deviation</i> (the default is use_median = FALSE)
matrix	A <i>matrix</i> of data.
look_back	The length of the look-back interval, equal to the number of rows of data used in the scaling.

Details

The function `roll_scale()` performs a rolling scaling (standardization) of the columns of the `matrix` argument using `RcppArmadillo`. The function `roll_scale()` performs a loop over the rows of `matrix`, subsets a number of previous (past) rows equal to `look_back`, and scales the subset matrix. It assigns the last row of the scaled subset *matrix* to the return matrix.

If the argument `use_median` is `FALSE` (the default), then it performs the same calculation as the function `roll::roll_scale()`. If the argument `use_median` is `TRUE`, then it calculates the centrality as the *median* and the dispersion as the *median absolute deviation (MAD)*.

Value

A *matrix* with the same dimensions as the input argument `matrix`.

Examples

```
## Not run:
matrixv <- matrix(rnorm(20000), nc=2)
look_back <- 11
rolled_scaled <- roll::roll_scale(data=matrixv, width = look_back, min_obs=1)
rolled_scaled2 <- roll_scale(matrix=matrixv, look_back = look_back, use_median=FALSE)
all.equal(rolled_scaled[-1, ], rolled_scaled2[-1, ])

## End(Not run)
```

roll_sharpe	<i>Calculate a time series of Sharpe ratios over a rolling look-back interval for an OHLC time series.</i>
-------------	--

Description

Calculate a time series of Sharpe ratios over a rolling look-back interval for an *OHLC* time series.

Usage

```
roll_sharpe(ohlc, look_back = 11)
```

Arguments

<code>ohlc</code>	An <i>OHLC</i> time series of prices in <i>xts</i> format.
<code>look_back</code>	The size of the look-back interval, equal to the number of rows of data used for aggregating the <i>OHLC</i> prices.

Details

The function `roll_sharpe()` calculates the rolling Sharpe ratio defined as the ratio of percentage returns over the look-back interval, divided by the average volatility of percentage returns.

Value

An *xts* time series with a single column and the same number of rows as the argument `ohlc`.

Examples

```
# Calculate rolling Sharpe ratio over SPY
sharpe_rolling <- roll_sharpe(ohlc=HighFreq::SPY, look_back=11)
```

roll_skew	<i>Calculate a matrix of skewness estimates over a rolling look-back interval attached at the end points of a time series or a matrix.</i>
-----------	--

Description

Calculate a *matrix* of skewness estimates over a rolling look-back interval attached at the end points of a *time series* or a *matrix*.

Usage

```
roll_skew(
  tseries,
  startp = 0L,
  endp = 0L,
  step = 1L,
  look_back = 1L,
  stub = 0L,
  method = "moment",
  confl = 0.75
)
```

Arguments

tseries	A <i>time series</i> or a <i>matrix</i> of data.
startp	An <i>integer</i> vector of start points (the default is startp = 0).
endp	An <i>integer</i> vector of end points (the default is endp = 0).
step	The number of time periods between the end points (the default is step = 1).
look_back	The number of end points in the look-back interval (the default is look_back = 1).
stub	An <i>integer</i> value equal to the first end point for calculating the end points (the default is stub = 0).
method	A <i>string</i> specifying the type of the skewness model (the default is method = "moment" - see Details).
confl	The confidence level for calculating the quantiles of returns (the default is confl = 0.75).

Details

The function roll_skew() calculates a *matrix* of skewness estimates over rolling look-back intervals attached at the end points of the *time series* tseries.

The function roll_skew() performs a loop over the end points, and at each end point it subsets the time series tseries over a look-back interval equal to look_back number of end points.

It passes the subset time series to the function calc_skew(), which calculates the skewness. See the function calc_skew() for a description of the skewness methods.

If the arguments `endp` and `startp` are not given then it first calculates a vector of end points separated by step time periods. It calculates the end points along the rows of `tseries` using the function `calc_endpoints()`, with the number of time periods between the end points equal to step time periods.

For example, the rolling skewness at 25 day end points, with a 75 day look-back, can be calculated using the parameters `step = 25` and `look_back = 3`.

The function `roll_skew()` is implemented in RcppArmadillo C++ code, so it's many times faster than the equivalent R code.

Value

A *matrix* of skewness estimates with the same number of columns as the input time series `tseries`, and the number of rows equal to the number of end points.

Examples

```
## Not run:
# Define time series of returns using package rutils
returns <- na.omit(rutils::etfenv$returns$VTI)
# Define end points and start points
endp <- 1 + HighFreq::calc_endpoints(NROW(returns), step=25)
startp <- HighFreq::calc_startpoints(endp, look_back=3)
# Calculate the rolling skewness at 25 day end points, with a 75 day look-back
skewv <- HighFreq::roll_skew(returns, step=25, look_back=3)
# Calculate the rolling skewness using R code
skewr <- sapply(1:NROW(endp), function(it) {
  HighFreq::calc_skew(returns[startp[it]:endp[it], ])
}) # end sapply
# Compare the skewness estimates
all.equal(drop(skewv), skewr, check.attributes=FALSE)
# Compare the speed of RcppArmadillo with R code
library(microbenchmark)
summary(microbenchmark(
  Rcpp=HighFreq::roll_skew(returns, step=25, look_back=3),
  Rcode=sapply(1:NROW(endp), function(it) {
    HighFreq::calc_skew(returns[startp[it]:endp[it], ])
  }),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary

## End(Not run)
```

roll_stats

Calculate a vector of statistics over an OHLC time series, and calculate a rolling mean over the statistics.

Description

Calculate a vector of statistics over an *OHLC* time series, and calculate a rolling mean over the statistics.

Usage

```
roll_stats(
  ohlc,
  calc_stats = "ohlc_variance",
  look_back = 11,
  weighted = TRUE,
  ...
)
```

Arguments

...	additional parameters to the function <code>calc_stats</code> .
ohlc	An <i>OHLC</i> time series of prices and trading volumes, in <i>xts</i> format.
calc_stats	The name of the function for estimating statistics of a single row of <i>OHLC</i> data, such as volatility, skew, and higher moments.
look_back	The size of the look-back interval, equal to the number of rows of data used for calculating the rolling mean.
weighted	<i>Boolean</i> argument: should statistic be weighted by trade volume? (default TRUE)

Details

The function `roll_stats()` calculates a vector of statistics over an *OHLC* time series, such as volatility, skew, and higher moments. The statistics could also be any other aggregation of a single row of *OHLC* data, for example the *High* price minus the *Low* price squared. The length of the vector of statistics is equal to the number of rows of the argument `ohlc`. Then it calculates a trade volume weighted rolling mean over the vector of statistics over and calculate statistics.

Value

An *xts* time series with a single column and the same number of rows as the argument `ohlc`.

Examples

```
# Calculate time series of rolling variance and skew estimates
var_rolling <- roll_stats(ohlc=HighFreq::SPY, look_back=21)
skew_rolling <- roll_stats(ohlc=HighFreq::SPY, calc_stats="ohlc_skew", look_back=21)
skew_rolling <- skew_rolling/(var_rolling)^(1.5)
skew_rolling[1, ] <- 0
skew_rolling <- rutils::na_locf(skew_rolling)
```

roll_sum

Calculate the rolling sums over a time series or a matrix using Rcpp.

Description

Calculate the rolling sums over a *time series* or a *matrix* using *Rcpp*.

Usage

```
roll_sum(tseries, look_back = 1L)
```


Arguments

tseries	A <i>time series</i> or a <i>matrix</i> .
look_back	The length of the look-back interval, equal to the number of data points included in calculating the rolling sum (the default is look_back = 1).

Details

The function roll_sum() calculates the rolling sums over the columns of the data tseries.

The function roll_sum() returns a *matrix* with the same dimensions as the input argument tseries.

The function roll_sum() uses the fast RcppArmadillo function arma::cumsum(), without explicit loops. The function roll_sum() is several times faster than rutils::roll_sum() which uses vectorized R code.

Value

A *matrix* with the same dimensions as the input argument tseries.

Examples

```
## Not run:
# Calculate historical returns
returns <- na.omit(rutils::etfenv$returns[, c("VTI", "IEF")])
# Define parameters
look_back <- 22
# Calculate rolling sums and compare with rutils::roll_sum()
c_sum <- HighFreq::roll_sum(returns, look_back)
r_sum <- rutils::roll_sum(returns, look_back)
all.equal(c_sum, coredata(r_sum), check.attributes=FALSE)
# Calculate rolling sums using R code
r_sum <- apply(zoo::coredata(returns), 2, cumsum)
lag_sum <- rbind(matrix(numeric(2*look_back), nc=2), r_sum[1:(NROW(r_sum) - look_back), ])
r_sum <- (r_sum - lag_sum)
all.equal(c_sum, r_sum, check.attributes=FALSE)

## End(Not run)
```

roll_sumep	<i>Calculate the rolling sums at the end points of a time series or a matrix.</i>
------------	---

Description

Calculate the rolling sums at the end points of a *time series* or a *matrix*.

Usage

```
roll_sumep(
  tseries,
  startp = 0L,
  endp = 0L,
  step = 1L,
```

```

    look_back = 1L,
    stub = 0L
  )

```

Arguments

<code>tseries</code>	A <i>time series</i> or a <i>matrix</i> .
<code>startp</code>	An <i>integer</i> vector of start points (the default is <code>startp = 0</code>).
<code>endp</code>	An <i>integer</i> vector of end points (the default is <code>endp = 0</code>).
<code>step</code>	The number of time periods between the end points (the default is <code>step = 1</code>).
<code>look_back</code>	The number of end points in the look-back interval (the default is <code>look_back = 1</code>).
<code>stub</code>	An <i>integer</i> value equal to the first end point for calculating the end points.

Details

The function `roll_sumep()` calculates the rolling sums at the end points of the *time series* `tseries`.

The function `roll_sumep()` is implemented in RcppArmadillo C++ code, which makes it several times faster than R code.

Value

A *matrix* with the same number of columns as the input time series `tseries`, and the number of rows equal to the number of end points.

Examples

```

## Not run:
# Calculate historical returns
returns <- na.omit(rutils::etfenv$returns[, c("VTI", "IEF")])
# Define end points at 25 day intervals
endp <- HighFreq::calc_endpoints(NROW(returns), step=25)
# Define start points as 75 day lag of end points
startp <- HighFreq::calc_startpoints(endp, look_back=3)
# Calculate rolling sums using Rcpp
c_sum <- HighFreq::roll_sumep(returns, startp=startp, endp=endp)
# Calculate rolling sums using R code
r_sum <- sapply(1:NROW(endp), function(ep) {
  colSums(returns[(startp[ep]+1):(endp[ep]+1)], )
}) # end sapply
r_sum <- t(r_sum)
all.equal(c_sum, r_sum, check.attributes=FALSE)

## End(Not run)

```

roll_var	Calculate a <i>matrix</i> of dispersion (variance) estimates over a rolling look-back interval attached at the end points of a <i>time series</i> or a <i>matrix</i> .
----------	--

Description

Calculate a *matrix* of dispersion (variance) estimates over a rolling look-back interval attached at the end points of a *time series* or a *matrix*.

Usage

```
roll_var(
  tseries,
  startp = 0L,
  endp = 0L,
  step = 1L,
  look_back = 1L,
  stub = 0L,
  method = "moment",
  confl = 0.75
)
```

Arguments

tseries	A <i>time series</i> or a <i>matrix</i> of data.
startp	An <i>integer</i> vector of start points (the default is startp = 0).
endp	An <i>integer</i> vector of end points (the default is endp = 0).
step	The number of time periods between the end points (the default is step = 1).
look_back	The number of end points in the look-back interval (the default is look_back = 1).
stub	An <i>integer</i> value equal to the first end point for calculating the end points (the default is stub = 0).
method	A <i>character</i> string representing the type of the measure of dispersion (the default is method = "moment").

Details

The function roll_var() calculates a *matrix* of dispersion (variance) estimates over rolling look-back intervals attached at the end points of the *time series* tseries.

The function roll_var() performs a loop over the end points, and at each end point it subsets the time series tseries over a look-back interval equal to look_back number of end points.

It passes the subset time series to the function calc_var(), which calculates the dispersion. See the function calc_var() for a description of the dispersion methods.

If the arguments endp and startp are not given then it first calculates a vector of end points separated by step time periods. It calculates the end points along the rows of tseries using the function calc_endpoints(), with the number of time periods between the end points equal to step time periods.

For example, the rolling variance at 25 day end points, with a 75 day look-back, can be calculated using the parameters `step = 25` and `look_back = 3`.

The function `roll_var()` with the parameter `step = 1` performs the same calculation as the function `roll_var()` from package **RcppRoll**, but it's several times faster because it uses RcppArmadillo C++ code.

The function `roll_var()` is implemented in RcppArmadillo C++ code, so it's many times faster than the equivalent R code.

Value

A *matrix* dispersion (variance) estimates with the same number of columns as the input time series *tseries*, and the number of rows equal to the number of end points.

Examples

```
## Not run:
# Define time series of returns using package rutils
returns <- na.omit(rutils::etfenv$returns$VTI)
# Calculate the rolling variance at 25 day end points, with a 75 day look-back
variance <- HighFreq::roll_var(returns, step=25, look_back=3)
# Compare the variance estimates over 11-period look-back intervals
all.equal(HighFreq::roll_var(returns, look_back=11)[-1:10], ,
  drop(RcppRoll::roll_var(returns, n=11)), check.attributes=FALSE)
# Compare the speed of HighFreq::roll_var() with RcppRoll::roll_var()
library(microbenchmark)
summary(microbenchmark(
  Rcpp=HighFreq::roll_var(returns, look_back=11),
  RcppRoll=RcppRoll::roll_var(returns, n=11),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary
# Compare the speed of HighFreq::roll_var() with TTR::runMAD()
summary(microbenchmark(
  Rcpp=HighFreq::roll_var(returns, look_back=11, method="quantile"),
  TTR=TTR::runMAD(returns, n = 11),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary

## End(Not run)
```

roll_varvec

Calculate a vector of variance estimates over a rolling look-back interval for a single-column time series or a single-column matrix, using RcppArmadillo.

Description

Calculate a *vector* of variance estimates over a rolling look-back interval for a single-column *time series* or a single-column *matrix*, using RcppArmadillo.

Usage

```
roll_varvec(tseries, look_back = 1L)
```

Arguments

tseries	A single-column <i>time series</i> or a single-column <i>matrix</i> .
look_back	The length of the look-back interval, equal to the number of <i>vector</i> elements used for calculating a single variance estimate (the default is look_back = 1).

Details

The function roll_varvec() calculates a *vector* of variance estimates over a rolling look-back interval for a single-column *time series* or a single-column *matrix*, using RcppArmadillo C++ code.

The function roll_varvec() uses an expanding look-back interval in the initial warmup period, to calculate the same number of elements as the input argument tseries.

The function roll_varvec() performs the same calculation as the function roll_var() from package **RcppRoll**, but it's several times faster because it uses RcppArmadillo C++ code.

Value

A single-column *matrix* with the same number of elements as the input argument tseries.

Examples

```
## Not run:
# Create a vector of random returns
returns <- rnorm(1e6)
# Compare the variance estimates over 11-period look-back intervals
all.equal(drop(HighFreq::roll_varvec(returns, look_back=11))[-(1:10)],
  RcppRoll::roll_var(returns, n=11))
# Compare the speed of RcppArmadillo with RcppRoll
library(microbenchmark)
summary(microbenchmark(
  Rcpp=HighFreq::roll_varvec(returns, look_back=11),
  RcppRoll=RcppRoll::roll_var(returns, n=11),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary

## End(Not run)
```

roll_var_ohlc	<i>Calculate a vector of variance estimates over a rolling look-back interval attached at the end points of a time series or a matrix with OHLC price data.</i>
---------------	---

Description

Calculate a *vector* of variance estimates over a rolling look-back interval attached at the end points of a *time series* or a *matrix* with *OHLC* price data.

Usage

```
roll_var_ohlc(
  ohlc,
  startp = 0L,
  endp = 0L,
```

```

    step = 1L,
    look_back = 1L,
    stub = 0L,
    method = "yang_zhang",
    scale = TRUE,
    index = 0L
)

```

Arguments

ohlc	A <i>time series</i> or a <i>matrix</i> with <i>OHLC</i> price data.
startp	An <i>integer</i> vector of start points (the default is startp = 0).
endp	An <i>integer</i> vector of end points (the default is endp = 0).
step	The number of time periods between the end points (the default is step = 1).
look_back	The number of end points in the look-back interval (the default is look_back = 1).
stub	An <i>integer</i> value equal to the first end point for calculating the end points (the default is stub = 0).
method	A <i>character</i> string representing the price range estimator for calculating the variance. The estimators include: <ul style="list-style-type: none"> • "close" close-to-close estimator, • "rogers_satchell" Rogers-Satchell estimator, • "garman_klass" Garman-Klass estimator, • "garman_klass_yz" Garman-Klass with account for close-to-open price jumps, • "yang_zhang" Yang-Zhang estimator, (The default is the "yang_zhang" estimator.)
scale	<i>Boolean</i> argument: Should the returns be divided by the time index, the number of seconds in each period? (The default is scale = TRUE.)
index	A <i>vector</i> with the time index of the <i>time series</i> . This is an optional argument (the default is index=0).

Details

The function `roll_var_ohlc()` calculates a *vector* of variance estimates over a rolling look-back interval attached at the end points of the *time series* `ohlc`.

The input *OHLC time series* `ohlc` is assumed to contain the log prices.

The function `roll_var_ohlc()` performs a loop over the end points, subsets the previous (past) rows of `ohlc`, and passes them into the function `calc_var_ohlc()`.

At each end point, the variance is calculated over a look-back interval equal to `look_back` number of end points. In the initial warmup period, the variance is calculated over an expanding look-back interval.

If the arguments `endp` and `startp` are not given then it first calculates a vector of end points separated by `step` time periods. It calculates the end points along the rows of `ohlc` using the function `calc_endpoints()`, with the number of time periods between the end points equal to `step` time periods.

For example, the rolling variance at daily end points with an 11 day look-back, can be calculated using the parameters `step = 1` and `look_back = 1` (Assuming the `ohlc` data has daily frequency.)

Similarly, the rolling variance at 25 day end points with a 75 day look-back, can be calculated using the parameters `step = 25` and `look_back = 3` (because $3 \times 25 = 75$).

The function `roll_var_ohlc()` calculates the variance from all the different intra-day and day-over-day returns (defined as the differences between *OHLC* prices), using several different variance estimation methods.

The default method is *"yang_zhang"*, which theoretically has the lowest standard error among unbiased estimators. The methods *"close"*, *"garman_klass_yz"*, and *"yang_zhang"* do account for *close-to-open* price jumps, while the methods *"garman_klass"* and *"rogers_satchell"* do not account for *close-to-open* price jumps.

If `scale` is `TRUE` (the default), then the returns are divided by the differences of the time index (which scales the variance to the units of variance per second squared.) This is useful when calculating the variance from minutely bar data, because dividing returns by the number of seconds decreases the effect of overnight price jumps. If the time index is in days, then the variance is equal to the variance per day squared.

The optional argument `index` is the time index of the *time series* *ohlc*. If the time index is in seconds, then the differences of the index are equal to the number of seconds in each time period. If the time index is in days, then the differences are equal to the number of days in each time period.

The function `roll_var_ohlc()` is implemented in RcppArmadillo C++ code, so it's many times faster than the equivalent R code.

Value

A column *vector* of variance estimates, with the number of rows equal to the number of end points.

Examples

```
## Not run:
# Extract the log OHLC prices of SPY
ohlc <- log(HighFreq::SPY)
# Extract the time index of SPY prices
indeks <- c(1, diff(xts::index(ohlc)))
# Rolling variance at minutely end points, with a 21 minute look-back
var_rolling <- HighFreq::roll_var_ohlc(ohlc,
                                     step=1, look_back=21,
                                     method="yang_zhang",
                                     index=indeks, scale=TRUE)

# Daily OHLC prices
ohlc <- rutils::etfenv$VTI
indeks <- c(1, diff(xts::index(ohlc)))
# Rolling variance at 5 day end points, with a 20 day look-back (20=4*5)
var_rolling <- HighFreq::roll_var_ohlc(ohlc,
                                     step=5, look_back=4,
                                     method="yang_zhang",
                                     index=indeks, scale=TRUE)

# Same calculation in R
nrows <- NROW(ohlc)
close_lag = HighFreq::lagit(ohlc[, 4])
endp <- drop(HighFreq::calc_endpoints(nrows, 3)) + 1
startp <- drop(HighFreq::calc_startpoints(endp, 2))
n_pts <- NROW(endp)
var_rollingr <- sapply(2:n_pts, function(it) {
  rangev <- startp[it]:endp[it]
  sub_ohlc = ohlc[rangev, ]
  sub_close = close_lag[rangev]
```

```

    sub_index = indeks[rangev]
    HighFreq::calc_var_ohlc(sub_ohlc, close_lag=sub_close, scale=TRUE, index=sub_index)
  }) # end sapply
  var_rollingr <- c(0, var_rollingr)
  all.equal(drop(var_rolling), var_rollingr)

## End(Not run)

```

roll_vec	<i>Calculate the rolling sums over a single-column time series or a single-column matrix using Rcpp.</i>
----------	--

Description

Calculate the rolling sums over a single-column *time series* or a single-column *matrix* using *Rcpp*.

Usage

```
roll_vec(tseries, look_back)
```

Arguments

tseries	A single-column <i>time series</i> or a single-column <i>matrix</i> .
look_back	The length of the look-back interval, equal to the number of elements of data used for calculating the sum.

Details

The function `roll_vec()` calculates a single-column *matrix* of rolling sums, over a single-column *matrix* of data, using fast *Rcpp* C++ code. The function `roll_vec()` is several times faster than `rutils::roll_sum()` which uses vectorized R code.

Value

A single-column *matrix* of the same length as the argument `tseries`.

Examples

```

## Not run:
# Define a single-column matrix of returns
returns <- zoo::coredata(na.omit(rutils::etfenv$returns$VTI))
# Calculate rolling sums over 11-period look-back intervals
sum_rolling <- HighFreq::roll_vec(returns, look_back=11)
# Compare HighFreq::roll_vec() with rutils::roll_sum()
all.equal(HighFreq::roll_vec(returns, look_back=11),
          rutils::roll_sum(returns, look_back=11),
          check.attributes=FALSE)
# Compare the speed of Rcpp with R code
library(microbenchmark)
summary(microbenchmark(
  Rcpp=HighFreq::roll_vec(returns, look_back=11),
  Rcode=rutils::roll_sum(returns, look_back=11),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary

```



```
## End(Not run)
```

roll_vecw	Calculate the rolling weighted sums over a single-column <i>time series</i> or a single-column <i>matrix</i> using RcppArmadillo.
-----------	---

Description

Calculate the rolling weighted sums over a single-column *time series* or a single-column *matrix* using RcppArmadillo.

Usage

```
roll_vecw(tseries, weights)
```

Arguments

tseries	A single-column <i>time series</i> or a single-column <i>matrix</i> .
weights	A single-column <i>matrix</i> of weights.

Details

The function `roll_vecw()` calculates the rolling weighted sums of a single-column *matrix* over its past values (a convolution with the single-column *matrix* of weights), using RcppArmadillo. It performs a similar calculation as the standard R function `stats::filter(x=series, filter=weights, method="convolution", sides=1)`, but it's over 6 times faster, and it doesn't produce any NA values.

Value

A single-column *matrix* of the same length as the argument `tseries`.

Examples

```
## Not run:
# First example
# Define a single-column matrix of returns
returns <- zoo::coredata(na.omit(rutils::etfenv$returns$VTI))
# Create simple weights
weights <- matrix(c(1, rep(0, 10)))
# Calculate rolling weighted sums
weighted <- HighFreq::roll_vecw(tseries=returns, weights=weights)
# Compare with original
all.equal(zoo::coredata(returns), weighted, check.attributes=FALSE)
# Second example
# Create exponentially decaying weights
weights <- matrix(exp(-0.2*1:11))
weights <- weights/sum(weights)
# Calculate rolling weighted sums
weighted <- HighFreq::roll_vecw(tseries=returns, weights=weights)
# Calculate rolling weighted sums using filter()
```

```

filtered <- stats::filter(x=returns, filter=weights, method="convolution", sides=1)
# Compare both methods
all.equal(filtered[-(1:11)], weighted[-(1:11)], check.attributes=FALSE)
# Compare the speed of Rcpp with R code
library(microbenchmark)
summary(microbenchmark(
  Rcpp=HighFreq::roll_vecw(tseries=returns, weights=weights),
  Rcode=stats::filter(x=returns, filter=weights, method="convolution", sides=1),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary

## End(Not run)

```

roll_vwap	<i>Calculate the volume-weighted average price of an OHLC time series over a rolling look-back interval.</i>
-----------	--

Description

Performs the same operation as function `VWAP()` from package **TTR**, but using vectorized functions, so it's a little faster.

Usage

```
roll_vwap(ohlc, close = ohlc[, 4, drop = FALSE], look_back)
```

Arguments

ohlc	An <i>OHLC</i> time series of prices in <i>xts</i> format.
close	A time series of close prices.
look_back	The size of the look-back interval, equal to the number of rows of data used for calculating the average price.

Details

The function `roll_vwap()` calculates the volume-weighted average closing price, defined as the sum of the prices multiplied by trading volumes in the look-back interval, divided by the sum of trading volumes in the interval. If the argument `close` is passed in explicitly, then its volume-weighted average value over time is calculated.

Value

An *xts* time series with a single column and the same number of rows as the argument `ohlc`.

Examples

```

# Calculate and plot rolling volume-weighted average closing prices (VWAP)
prices_rolling <- roll_vwap(ohlc=HighFreq::SPY["2013-11"], look_back=11)
chart_Series(HighFreq::SPY["2013-11-12"], name="SPY prices")
add_TA(prices_rolling["2013-11-12"], on=1, col="red", lwd=2)
legend("top", legend=c("SPY prices", "VWAP prices"),
bg="white", lty=c(1, 1), lwd=c(2, 2),
col=c("black", "red"), bty="n")

```

```
# Calculate running returns
returns_running <- ohlc_returns(xtsv=HighFreq::SPY)
# Calculate the rolling volume-weighted average returns
roll_vwap(ohlc=HighFreq::SPY, close=returns_running, look_back=11)
```

roll_wsum	<i>Calculate the rolling weighted sums over a time series or a matrix using Rcpp.</i>
-----------	---

Description

Calculate the rolling weighted sums over a *time series* or a *matrix* using *Rcpp*.

Usage

```
roll_wsum(tseries, endp = NULL, look_back = 1L, stub = NULL, weights = NULL)
```

Arguments

tseries	A <i>time series</i> or a <i>matrix</i> .
endp	An <i>integer</i> vector of end points (the default is endp = NULL).
look_back	The length of the look-back interval, equal to the number of data points included in calculating the rolling sum (the default is look_back = 1).
stub	An <i>integer</i> value equal to the first end point for calculating the end points (the default is stub = NULL).
weights	A single-column <i>matrix</i> of weights (the default is weights = NULL).

Details

The function roll_wsum() calculates the rolling weighted sums over the columns of the data tseries.

The function roll_wsum() calculates the rolling weighted sums as convolutions of the columns of tseries with the *column vector* of weights using the RcppArmadillo function arma::conv2(). It performs a similar calculation to the standard R function stats::filter(x=returns, filter=weights, method="convolution", sides=1), but it can be many times faster, and it doesn't produce any leading NA values.

The function roll_wsum() returns a *matrix* with the same dimensions as the input argument tseries.

The arguments weights, endp, and stub are optional.

If the argument weights is not supplied, then simple sums are calculated, not weighted sums.

If either the stub or endp arguments are supplied, then the rolling sums are calculated at the end points.

If only the argument stub is supplied, then the end points are calculated from the stub and look_back arguments. The first end point is equal to stub and the end points are spaced look_back periods apart.

If the arguments weights, endp, and stub are not supplied, then the sums are calculated over a number of data points equal to look_back.

The function roll_wsum() is also several times faster than rutils::roll_sum() which uses vectorized R code.

Technical note: The function roll_wsum() has arguments with default values equal to NULL, which are implemented in Rcpp code.

Value

A *matrix* with the same dimensions as the input argument *tseries*.

Examples

```
## Not run:
# First example
# Calculate historical returns
returns <- na.omit(rutils::etfenv$returns[, c("VTI", "IEF")])
# Define parameters
look_back <- 22
# Calculate rolling sums and compare with rutils::roll_sum()
c_sum <- HighFreq::roll_sum(returns, look_back)
r_sum <- rutils::roll_sum(returns, look_back)
all.equal(c_sum, coredata(r_sum), check.attributes=FALSE)
# Calculate rolling sums using R code
r_sum <- apply(zoo::coredata(returns), 2, cumsum)
lag_sum <- rbind(matrix(numeric(2*look_back), nc=2), r_sum[1:(NROW(r_sum) - look_back), ])
r_sum <- (r_sum - lag_sum)
all.equal(c_sum, r_sum, check.attributes=FALSE)

# Calculate rolling sums at end points
stu_b <- 21
c_sum <- HighFreq::roll_wsum(returns, look_back, stub=stu_b)
endp <- (stu_b + look_back*(0:(NROW(returns) %% look_back)))
endp <- endp[endp < NROW(returns)]
r_sum <- apply(zoo::coredata(returns), 2, cumsum)
r_sum <- r_sum[endp+1, ]
lag_sum <- rbind(numeric(2), r_sum[1:(NROW(r_sum) - 1), ])
r_sum <- (r_sum - lag_sum)
all.equal(c_sum, r_sum, check.attributes=FALSE)

# Calculate rolling sums at end points - pass in endp
c_sum <- HighFreq::roll_wsum(returns, endp=endp)
all.equal(c_sum, r_sum, check.attributes=FALSE)

# Create exponentially decaying weights
weights <- exp(-0.2*(1:11))
weights <- matrix(weights/sum(weights), nc=1)
# Calculate rolling weighted sum
c_sum <- HighFreq::roll_wsum(returns, weights=weights)
# Calculate rolling weighted sum using filter()
filtered <- filter(x=returns, filter=weights, method="convolution", sides=1)
all.equal(c_sum[-(1:11), ], filtered[-(1:11), ], check.attributes=FALSE)

# Calculate rolling weighted sums at end points
c_sum <- HighFreq::roll_wsum(returns, endp=endp, weights=weights)
all.equal(c_sum, filtered[endp+1, ], check.attributes=FALSE)

# Create simple weights equal to a 1 value plus zeros
weights <- matrix(c(1, rep(0, 10)), nc=1)
# Calculate rolling weighted sum
weighted <- HighFreq::roll_wsum(returns, weights=weights)
# Compare with original
all.equal(coredata(returns), weighted, check.attributes=FALSE)
```

```
## End(Not run)
```

roll_zscores	<i>Calculate a vector of z-scores of the residuals of rolling regressions at the end points of the predictor matrix.</i>
--------------	--

Description

Calculate a *vector* of z-scores of the residuals of rolling regressions at the end points of the predictor matrix.

Usage

```
roll_zscores(
  response,
  predictor,
  startp = 0L,
  endp = 0L,
  step = 1L,
  look_back = 1L,
  stub = 0L
)
```

Arguments

response	A single-column <i>time series</i> or a <i>vector</i> of response data.
predictor	A <i>time series</i> or a <i>matrix</i> of predictor data.
startp	An <i>integer</i> vector of start points (the default is startp = 0).
endp	An <i>integer</i> vector of end points (the default is endp = 0).
step	The number of time periods between the end points (the default is step = 1).
look_back	The number of end points in the look-back interval (the default is look_back = 1).
stub	An <i>integer</i> value equal to the first end point for calculating the end points (the default is stub = 0).

Details

The function `roll_zscores()` calculates a *vector* of z-scores of the residuals of rolling regressions at the end points of the *time series* predictor.

The function `roll_zscores()` performs a loop over the end points, and at each end point it subsets the time series predictor over a look-back interval equal to `look_back` number of end points.

It passes the subset time series to the function `calc_lm()`, which calculates the regression data.

If the arguments `endp` and `startp` are not given then it first calculates a vector of end points separated by `step` time periods. It calculates the end points along the rows of predictor using the function `calc_endpoints()`, with the number of time periods between the end points equal to `step` time periods.

For example, the rolling variance at 25 day end points, with a 75 day look-back, can be calculated using the parameters `step = 25` and `look_back = 3`.

Value

A column *vector* of the same length as the number of rows of predictor.

Examples

```
## Not run:
# Calculate historical returns
returns <- na.omit(rutils::etfenv$returns[, c("XLF", "VTI", "IEF")])
# Response equals XLF returns
response <- returns[, 1]
# Predictor matrix equals VTI and IEF returns
predictor <- returns[, -1]
# Calculate Z-scores from rolling time series regression using RcppArmadillo
look_back <- 11
zscores <- HighFreq::roll_zscores(response=response, predictor=predictor, look_back)
# Calculate z-scores in R from rolling multivariate regression using lm()
zscoresr <- sapply(1:NROW(predictor), function(ro_w) {
  if (ro_w == 1) return(0)
  startpoint <- max(1, ro_w-look_back+1)
  responsi <- response[startpoint:ro_w]
  predicti <- predictor[startpoint:ro_w, ]
  lmod <- lm(responsi ~ predicti)
  residuals <- lmod$residuals
  residuals[NROW(residuals)]/sd(residuals)
}) # end sapply
# Compare the outputs of both functions
all.equal(zscores[-(1:look_back)], zscoresr[-(1:look_back)],
  check.attributes=FALSE)

## End(Not run)
```

run_covar

Calculate the running covariance of two streaming time series of returns.

Description

Calculate the running covariance of two streaming *time series* of returns.

Usage

```
run_covar(tseries, lambda)
```

Arguments

tseries A *time series* or a *matrix* with two columns of returns data.

lambda A *numeric* decay factor to multiply past estimates.

Details

The function `run_covar()` calculates the running covariance of two streaming *time series* of returns, by recursively weighing the products of their returns minus their means, with past covariance estimates σ_{t-1}^{cov} , using the decay factor λ :

$$\mu_t^1 = (1 - \lambda)r_t^1 + \lambda\mu_{t-1}^1$$

$$\mu_t^2 = (1 - \lambda)r_t^2 + \lambda\mu_{t-1}^2$$

$$\sigma_t^{cov} = (1 - \lambda)(r_t^1 - \mu_t^1)(r_t^2 - \mu_t^2) + \lambda\sigma_{t-1}^{cov}$$

Where σ_t^{cov} is the covariance estimate at time t , r_t^1 and r_t^2 are the two streaming returns data, and μ_t^1 and μ_t^2 are the means of the returns.

The above recursive formulas are convenient for processing live streaming data because they don't require maintaining a buffer of past data. The formulas are equivalent to a convolution with exponentially decaying weights, but they're faster to calculate.

The value of the decay factor λ should be in the range between 0 and 1. If λ is close to 1 then the decay is weak and past values have a greater weight, and the running covariance values have a stronger dependence on past values. This is equivalent to a long look-back interval. If λ is much less than 1 then the decay is strong and past values have a smaller weight, and the running covariance values have a weaker dependence on past values. This is equivalent to a short look-back interval.

The function `run_covar()` returns three columns of data: the covariance and the variances of the two columns of the argument `tseries`. This allows calculating the running correlation.

Value

A *matrix* with three columns of data: the covariance and the variances of the two columns of the argument `tseries`.

Examples

```
## Not run:
# Calculate historical returns
returns <- zoo::coredata(na.omit(rutils::etfenv$returns[, c("IEF", "VTI")]))
# Calculate the running covariance
lambda <- 0.9
covars <- HighFreq::run_covar(returns, lambda=lambda)
# Calculate running covariance using R code
filtered <- (1-lambda)*filter(returns[, 1]*returns[, 2],
  filter=lambda, init=as.numeric(returns[1, 1]*returns[1, 2])/(1-lambda),
  method="recursive")
all.equal(covars[, 1], unclass(filtered), check.attributes=FALSE)
# Calculate the running correlation
correl <- covars[, 1]/sqrt(covars[, 2]*covars[, 3])

## End(Not run)
```

run_max	Calculate the running maximum of streaming time series data.
---------	--

Description

Calculate the running maximum of streaming *time series* data.

Usage

```
run_max(tseries, lambda)
```

Arguments

tseries	A <i>time series</i> or a <i>matrix</i> .
lambda	A <i>numeric</i> decay factor to multiply past estimates.

Details

The function `run_max()` calculates the running maximum of streaming *time series* data by recursively weighing present and past values using the decay factor λ .

It first calculates the running mean of streaming data:

$$\mu_t = (1 - \lambda)p_t + \lambda\mu_{t-1}$$

Where μ_t is the mean value at time t , and p_t is the streaming data.

It then calculates the running maximums of streaming data, p_t^{max} :

$$p_t^{max} = \max(p_t, p_{t-1}^{max}) + (1 - \lambda)(\mu_{t-1} - p_{t-1}^{max})$$

The second term pulls the maximum value down to the mean value, allowing it to gradually "forget" the maximum value from the more distant past.

The above recursive formulas are convenient for processing live streaming data because they don't require maintaining a buffer of past data.

The value of the decay factor λ should be in the range between 0 and 1. If λ is close to 1 then the decay is weak and past values have a greater weight, and the running maximum values have a stronger dependence on past values. This is equivalent to a long look-back interval. If λ is much less than 1 then the decay is strong and past values have a smaller weight, and the running maximum values have a weaker dependence on past values. This is equivalent to a short look-back interval.

The function `run_max()` returns a *matrix* with the same dimensions as the input argument `tseries`.

Value

A *matrix* with the same dimensions as the input argument `tseries`.

Examples

```
## Not run:
# Calculate historical prices
prices <- zoo::coredata(quantmod::Cl(rutils::etfenv$VTI))
# Calculate the running maximums
lambda <- 0.9
maxs <- HighFreq::run_max(prices, lambda=lambda)
# Plot dygraph of VTI prices and running maximums
datav <- cbind(quantmod::Cl(rutils::etfenv$VTI), maxs)
colnames(datav) <- c("prices", "max")
colnamev <- colnames(datav)
dygraphs::dygraph(datav, main="VTI Prices and Running Maximums") %>%
  dySeries(name=colnamev[1], label=colnamev[1], strokeWidth=2, col="blue") %>%
  dySeries(name=colnamev[2], label=colnamev[2], strokeWidth=2, col="red")

## End(Not run)
```

run_mean

Calculate the running weighted means of streaming time series data.

Description

Calculate the running weighted means of streaming *time series* data.

Usage

```
run_mean(tseries, lambda, weights = 0L)
```

Arguments

tseries	A <i>time series</i> or a <i>matrix</i> .
lambda	A <i>numeric</i> decay factor to multiply past estimates.
weights	A single-column <i>matrix</i> of weights.

Details

The function `run_mean()` calculates the running weighted means of the streaming *time series* data p_t by recursively weighing present and past values using the decay factor λ . If the `weights` argument is omitted, then the function `run_mean()` simply calculates the exponentially weighted moving average value of the streaming *time series* data p_t :

$$\mu_t = (1 - \lambda)p_t + \lambda\mu_{t-1} = (1 - \lambda) \sum_{j=0}^n \lambda^j p_{t-j}$$

Some applications require applying additional weight factors, like for example the volume-weighted average price indicator. The streaming prices are multiplied by the streaming trading volumes.

If the `weights` argument is included, then the function `run_mean()` calculates the running weighted means in two steps.

First it calculates the running mean weights μ_t^w :

$$\mu_t^w = (1 - \lambda)w_t + \lambda\mu_{t-1}^w$$

Second it calculates the the running mean products μ_t^p of the weights w_t and the data p_t :

$$\mu_t^p = (1 - \lambda)w_t p_t + \lambda\mu_{t-1}^p$$

Where p_t is the streaming data, w_t are the streaming weights, μ_t^w are the running mean weights, and μ_t^p are the running mean products of the data and the weights.

The running mean weighted value μ_t is equal to the ratio of the data and weights products, divided by the mean weights:

$$\mu_t = \frac{\mu_t^p}{\mu_t^w}$$

The above recursive formulas are convenient for processing live streaming data because they don't require maintaining a buffer of past data. The formulas are equivalent to a convolution with exponentially decaying weights, but they're faster to calculate.

The value of the decay factor λ should be in the range between 0 and 1. If λ is close to 1 then the decay is weak and past values have a greater weight, and the running mean values have a stronger dependence on past values. This is equivalent to a long look-back interval. If λ is much less than 1 then the decay is strong and past values have a smaller weight, and the running mean values have a weaker dependence on past values. This is equivalent to a short look-back interval.

The function `run_mean()` performs the same calculation as the standard R function `stats::filter(x=series, filter=lambda, method="recursive")`, but it's several times faster.

The function `run_mean()` returns a *matrix* with the same dimensions as the input argument `tseries`.

Value

A *matrix* with the same dimensions as the input argument `tseries`.

Examples

```
## Not run:
# Calculate historical prices
ohlcv <- rutils::etfenv$VTI
closep <- quantmod::Cl(ohlcv)
# Calculate the running means
lambda <- 0.95
means <- HighFreq::run_mean(closep, lambda=lambda)
# Calculate running means using R code
filtered <- (1-lambda)*filter(prices,
  filter=lambda, init=as.numeric(prices[1, 1])/(1-lambda),
  method="recursive")
all.equal(drop(means), unclass(filtered), check.attributes=FALSE)

# Compare the speed of RcppArmadillo with R code
library(microbenchmark)
summary(microbenchmark(
  Rcpp=HighFreq::run_mean(prices, lambda=lambda),
  Rcode=filter(prices, filter=lambda, init=as.numeric(prices[1, 1])/(1-lambda), method="recursive"),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary

# Create weights equal to the trading volumes
```

```

weights <- quantmod::Vo(ohlc)
# Calculate the running weighted means
meanw <- HighFreq::run_mean(prices, lambda=lambda, weights=weights)
# Plot dygraph of the running weighted means
datav <- xts(cbind(means, meanw), zoo::index(ohlc))
colnames(datav) <- c("means running", "means weighted")
dygraphs::dygraph(datav, main="Running Means") %>%
  dyOptions(colors=c("blue", "red"), strokeWidth=1) %>%
  dyLegend(show="always", width=500)

## End(Not run)

```

run_min

*Calculate the running minimum of streaming time series data.***Description**

Calculate the running minimum of streaming *time series* data.

Usage

```
run_min(tseries, lambda)
```

Arguments

`tseries` *A time series or a matrix.*

`lambda` *A numeric decay factor to multiply past estimates.*

Details

The function `run_min()` calculates the running minimum of streaming *time series* data by recursively weighing present and past values using the decay factor λ .

It first calculates the running mean of streaming data:

$$\mu_t = (1 - \lambda)p_t + \lambda\mu_{t-1}$$

Where μ_t is the mean value at time t , and p_t is the streaming data.

It then calculates the running minimums of streaming data, p_t^{min} :

$$p_t^{min} = \min(p_t, p_{t-1}^{min}) + (1 - \lambda)(\mu_{t-1} - p_{t-1}^{min})$$

The second term pulls the minimum value up to the mean value, allowing it to gradually "forget" the minimum value from the more distant past.

The above recursive formulas are convenient for processing live streaming data because they don't require maintaining a buffer of past data.

The value of the decay factor λ should be in the range between 0 and 1. If λ is close to 1 then the decay is weak and past values have a greater weight, and the running minimum values have a stronger dependence on past values. This is equivalent to a long look-back interval. If λ is much less than 1 then the decay is strong and past values have a smaller weight, and the running minimum values have a weaker dependence on past values. This is equivalent to a short look-back interval.

The function `run_min()` returns a *matrix* with the same dimensions as the input argument `tseries`.

Value

A *matrix* with the same dimensions as the input argument `tseries`.

Examples

```
## Not run:
# Calculate historical prices
prices <- zoo::coredata(quantmod::Cl(rutils::etfenv$VTI))
# Calculate the running minimums
lambda <- 0.9
mins <- HighFreq::run_min(prices, lambda=lambda)
# Plot dygraph of VTI prices and running minimums
datav <- cbind(quantmod::Cl(rutils::etfenv$VTI), mins)
colnames(datav) <- c("prices", "min")
colnamev <- colnames(datav)
dygraphs::dygraph(datav, main="VTI Prices and Running Minimums") %>%
  dySeries(name=colnamev[1], label=colnamev[1], strokeWidth=2, col="blue") %>%
  dySeries(name=colnamev[2], label=colnamev[2], strokeWidth=2, col="red")

## End(Not run)
```

run_reg

Perform running regressions of streaming time series of response and predictor data, and calculate the alphas, betas, and the residuals.

Description

Perform running regressions of streaming *time series* of response and predictor data, and calculate the alphas, betas, and the residuals.

Usage

```
run_reg(response, predictor, lambda, method = "none")
```

Arguments

<code>response</code>	A single-column <i>time series</i> or a single-column <i>matrix</i> of response data.
<code>predictor</code>	A <i>time series</i> or a <i>matrix</i> of predictor data.
<code>lambda</code>	A <i>numeric</i> decay factor to multiply past estimates.
<code>method</code>	A <i>string</i> specifying the method for scaling the residuals (see Details) (the default is <code>method = "none"</code> - no scaling)

Details

The function `run_reg()` calculates the vectors of *alphas* α_t , *betas* β_t , and the *residuals* ϵ_t of running regressions, by recursively weighing the current estimates with past estimates, using the decay factor λ :

$$\mu_t^r = (1 - \lambda)r_t^r + \lambda\mu_{t-1}^r$$

$$\mu_t^p = (1 - \lambda)r_t^p + \lambda\mu_{t-1}^p$$

$$\begin{aligned}\sigma_t^2 &= (1 - \lambda)(r_t^{p2} - \mu_t^{p2}) + \lambda\sigma_{t-1}^2 \\ \sigma_t^{cov} &= (1 - \lambda)(r_t^r - \mu_t^r)(r_t^p - \mu_t^p) + \lambda\sigma_{t-1}^{cov} \\ \beta_t &= (1 - \lambda)\frac{\sigma_t^{cov}}{\sigma_t^2} + \lambda\beta_{t-1} \\ \epsilon_t &= (1 - \lambda)(r_t^r - \beta_t r_t^p) + \lambda\epsilon_{t-1}\end{aligned}$$

Where σ_t^{cov} are the covariances between the response and predictor data at time t ; σ_t^2 is the vector of predictor variances, and r_t^r and r_t^p are the streaming data of the response and predictor data.

The matrices σ^2 , σ^{cov} , and β have the same number of rows as the input argument predictor.

The above recursive formulas are convenient for processing live streaming data because they don't require maintaining a buffer of past data. The formulas are equivalent to a convolution with exponentially decaying weights, but they're faster to calculate.

The value of the decay factor λ should be in the range between 0 and 1. If λ is close to 1 then the decay is weak and past values have a greater weight, so the running values have a greater dependence on past values. This is equivalent to a long look-back interval. If λ is much less than 1 then the decay is strong and past values have a smaller weight, so the running values have a weaker dependence on past values. This is equivalent to a short look-back interval.

The *residuals* may be scaled by their volatilities to obtain the *z-scores*. The default is method = "none" - no scaling. If the argument method = "scale" then the *residuals* ϵ_t are divided by their volatilities σ^ϵ without subtracting their means:

$$\epsilon_t = \frac{\epsilon_t}{\sigma^\epsilon}$$

If the argument method = "standardize" then the means μ_ϵ are subtracted from the *residuals*, and then they are divided by their volatilities σ^ϵ :

$$\epsilon_t = \frac{\epsilon_t - \mu_\epsilon}{\sigma^\epsilon}$$

Which are equal to the *z-scores*.

The function run_reg() returns multiple columns of data. If the matrix predictor has n columns then run_reg() returns a matrix with n+2 columns. The first column contains the *residuals*, the second the *alphas*, and the remaining columns contain the *betas*.

Value

A *matrix* with the regression alphas, betas, and residuals.

Examples

```
## Not run:
# Calculate historical returns
returns <- na.omit(rutils::etfenv$returns[, c("XLF", "VTI", "IEF")])
# Response equals XLF returns
response <- returns[, 1]
# Predictor matrix equals VTI and IEF returns
predictor <- returns[, -1]
# Calculate the running regressions
lambda <- 0.9
regs <- HighFreq::run_reg(response=response, predictor=predictor, lambda=lambda)
# Plot the running residuals
datav <- cbind(cumsum(response), regs[, 1])
```

```

colnames(datav) <- c("XLF", "residuals")
colnamev <- colnames(datav)
dygraphs::dygraph(datav, main="Residuals of XLF Versus VTI and IEF") %>%
  dyAxis("y", label=colnamev[1], independentTicks=TRUE) %>%
  dyAxis("y2", label=colnamev[2], independentTicks=TRUE) %>%
  dySeries(name=colnamev[1], axis="y", label=colnamev[1], strokeWidth=1, col="blue") %>%
  dySeries(name=colnamev[2], axis="y2", label=colnamev[2], strokeWidth=1, col="red")

## End(Not run)

```

run_var

Calculate the running variance of streaming time series of returns.

Description

Calculate the running variance of streaming *time series* of returns.

Usage

```
run_var(tseries, lambda)
```

Arguments

tseries A *time series* or a *matrix* of returns.
lambda A *numeric* decay factor to multiply past estimates.

Details

The function `run_var()` calculates the running variance of a streaming *time series* of returns, by recursively weighing the squared returns r_t^2 minus the squared means μ_t^2 , with the past variance estimates σ_{t-1}^2 , using the decay factor λ :

$$\mu_t = (1 - \lambda)r_t + \lambda\mu_{t-1}$$

$$\sigma_t^2 = (1 - \lambda)(r_t^2 - \mu_t^2) + \lambda\sigma_{t-1}^2$$

Where σ_t^2 is the variance estimate at time t , and r_t are the streaming returns data.

The above recursive formulas are convenient for processing live streaming data because they don't require maintaining a buffer of past data. The formulas are equivalent to a convolution with exponentially decaying weights, but they're faster to calculate.

The value of the decay factor λ should be in the range between 0 and 1. If λ is close to 1 then the decay is weak and past values have a greater weight, and the running variance values have a stronger dependence on past values. This is equivalent to a long look-back interval. If λ is much less than 1 then the decay is strong and past values have a smaller weight, and the running variance values have a weaker dependence on past values. This is equivalent to a short look-back interval.

The function `run_var()` performs the same calculation as the standard R function `stats::filter(x=series, filter=weights, method="recursive")`, but it's several times faster.

The function `run_var()` returns a *matrix* with the same dimensions as the input argument `tseries`.

Value

A *matrix* with the same dimensions as the input argument `tseries`.

Examples

```
## Not run:
# Calculate historical returns
returns <- zoo::coredata(na.omit(rutils::etfenv$returns$VTI))
# Calculate the running variance
lambda <- 0.9
vars <- HighFreq::run_var(returns, lambda=lambda)
# Calculate running variance using R code
filtered <- (1-lambda)*filter(returns^2, filter=lambda,
  init=as.numeric(returns[1, 1])^2/(1-lambda),
  method="recursive")
all.equal(vars, unclass(filtered), check.attributes=FALSE)
# Compare the speed of RcppArmadillo with R code
library(microbenchmark)
summary(microbenchmark(
  Rcpp=HighFreq::run_var(returns, lambda=lambda),
  Rcode=filter(returns^2, filter=lambda, init=as.numeric(returns[1, 1])^2/(1-lambda), method="recursive"),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary

## End(Not run)
```

run_var_ohlc

Calculate the running variance of streaming OHLC price data.

Description

Calculate the running variance of streaming *OHLC* price data.

Usage

```
run_var_ohlc(ohlc, lambda)
```

Arguments

ohlc	A <i>time series</i> or a <i>matrix</i> with <i>OHLC</i> price data.
lambda	A <i>numeric</i> decay factor to multiply past estimates.

Details

The function `run_var_ohlc()` calculates a single-column *matrix* of variance estimates of streaming *OHLC* price data.

The function `run_var_ohlc()` calculates the variance from the differences between the *Open*, *High*, *Low*, and *Close* prices, using the *Yang-Zhang* range volatility estimator:

$$\sigma_t^2 = (1-\lambda)((O_t - C_{t-1})^2 + 0.134(C_t - O_t)^2 + 0.866((H_i - O_i)(H_i - C_i) + (L_i - O_i)(L_i - C_i))) + \lambda\sigma_{t-1}^2$$

It recursively weighs the current variance estimate with the past estimates σ_{t-1}^2 , using the decay factor λ .

The above recursive formula is convenient for processing live streaming data because it doesn't require maintaining a buffer of past data. The formula is equivalent to a convolution with exponentially decaying weights, but it's faster.

The function `run_var_ohlc()` does not calculate the logarithm of the prices. So if the argument `ohlc` contains dollar prices then `run_var_ohlc()` calculates the dollar variance. If the argument `ohlc` contains the log prices then `run_var_ohlc()` calculates the percentage variance.

The function `run_var_ohlc()` is implemented in RcppArmadillo C++ code, so it's many times faster than the equivalent R code.

Value

A single-column *matrix* of variance estimates, with the same number of rows as the input `ohlc` price data.

Examples

```
## Not run:
# Extract the log OHLC prices of VTI
ohlc <- log(rutils::etfenv$VTI)
# Calculate the running variance
var_running <- HighFreq::run_var_ohlc(ohlc, lambda=0.8)
# Calculate the rolling variance
var_rolling <- HighFreq::roll_var_ohlc(ohlc, look_back=5, method="yang_zhang", scale=FALSE)
datav <- cbind(var_running, var_rolling)
colnames(datav) <- c("running", "rolling")
colnamev <- colnames(datav)
datav <- xts::xts(datav, index(ohlc))
# dygraph plot of VTI running versus rolling volatility
dygraphs::dygraph(sqrt(datav[-(1:111), ]), main="Running and Rolling Volatility of VTI") %>%
  dyOptions(colors=c("red", "blue"), strokeWidth=1) %>%
  dyLegend(show="always", width=500)
# Compare the speed of running versus rolling volatility
library(microbenchmark)
summary(microbenchmark(
  running=HighFreq::run_var_ohlc(ohlc, lambda=0.8),
  rolling=HighFreq::roll_var_ohlc(ohlc, look_back=5, method="yang_zhang", scale=FALSE),
  times=10))[, c(1, 4, 5)]

## End(Not run)
```

run_zscores

Calculate the z-scores of running regressions of streaming time series of returns.

Description

Calculate the z-scores of running regressions of streaming *time series* of returns.

Usage

```
run_zscores(response, predictor, lambda, demean = TRUE)
```


Arguments

response	A single-column <i>time series</i> or a single-column <i>matrix</i> of response data.
predictor	A <i>time series</i> or a <i>matrix</i> of predictor data.
lambda	A <i>numeric</i> decay factor to multiply past estimates.
demean	A <i>Boolean</i> specifying whether the <i>z-scores</i> should be de-meanned (the default is <code>demean = TRUE</code>).

Details

The function `run_zscores()` calculates the vectors of *betas* β_t and the residuals ϵ_t of running regressions by recursively weighing the current estimates with past estimates, using the decay factor λ :

$$\begin{aligned}\sigma_t^2 &= (1 - \lambda)r_t^{p2} + \lambda\sigma_{t-1}^2 \\ \sigma_t^{cov} &= (1 - \lambda)r_t^r r_t^p + \lambda\sigma_{t-1}^{cov} \\ \beta_t &= (1 - \lambda)\frac{\sigma_t^{cov}}{\sigma_t^2} + \lambda\beta_{t-1} \\ \epsilon_t &= (1 - \lambda)(r_t^r - \beta_t r_t^p) + \lambda\epsilon_{t-1}\end{aligned}$$

Where σ_t^{cov} is the vector of covariances between the response and predictor returns, at time t ; σ_t^2 is the vector of predictor variances, and r_t^r and r_t^p are the streaming returns of the response and predictor data.

The above formulas for σ^2 and σ^{cov} are approximate because they don't subtract the means before squaring the returns. But they're very good approximations for daily returns.

The matrices σ^2 , σ^{cov} , and β have the same number of rows as the input argument predictor.

If the argument `demean = TRUE` (the default) then the *z-scores* z_t are calculated as equal to the residuals ϵ_t minus their means μ_ϵ , divided by their volatilities σ^ϵ :

$$z_t = \frac{\epsilon_t - \mu_\epsilon}{\sigma^\epsilon}$$

If the argument `demean = FALSE` then the *z-scores* are only divided by their volatilities without subtracting their means:

$$z_t = \frac{\epsilon_t}{\sigma^\epsilon}$$

The above recursive formulas are convenient for processing live streaming data because they don't require maintaining a buffer of past data. The formulas are equivalent to a convolution with exponentially decaying weights, but they're faster to calculate.

The value of the decay factor λ should be in the range between 0 and 1. If λ is close to 1 then the decay is weak and past values have a greater weight, and the running *z-score* values have a stronger dependence on past values. This is equivalent to a long look-back interval. If λ is much less than 1 then the decay is strong and past values have a smaller weight, and the running *z-score* values have a weaker dependence on past values. This is equivalent to a short look-back interval.

The function `run_zscores()` returns multiple columns of data. If the matrix predictor has n columns then `run_zscores()` returns a matrix with $2n+1$ columns. The first column contains the *z-scores*, and the remaining columns contain the *betas* and the *variances* of the predictor data.

Value

A *matrix* with the *z-scores*, *betas*, and the *variances* of the predictor data.

Examples

```
## Not run:
# Calculate historical returns
returns <- na.omit(rutils::etfenv$returns[, c("XLF", "VTI", "IEF")])
# Response equals XLF returns
response <- returns[, 1]
# Predictor matrix equals VTI and IEF returns
predictor <- returns[, -1]
# Calculate the running z-scores
lambda <- 0.9
zscores <- HighFreq::run_zscores(response=response, predictor=predictor, lambda=lambda)
# Plot the running z-scores
datav <- cbind(cumsum(response), zscores[, 1])
colnames(datav) <- c("XLF", "zscores")
colnamev <- colnames(datav)
dygraphs::dygraph(datav, main="Z-Scores of XLF Versus VTI and IEF") %>%
  dyAxis("y", label=colnamev[1], independentTicks=TRUE) %>%
  dyAxis("y2", label=colnamev[2], independentTicks=TRUE) %>%
  dySeries(name=colnamev[1], axis="y", label=colnamev[1], strokeWidth=1, col="blue") %>%
  dySeries(name=colnamev[2], axis="y2", label=colnamev[2], strokeWidth=1, col="red")

## End(Not run)
```

save_rets

Load, scrub, aggregate, and rbind multiple days of TAQ data for a single symbol. Calculate returns and save them to a single ‘.RData’ file.*

Description

Load, scrub, aggregate, and rbind multiple days of *TAQ* data for a single symbol. Calculate returns and save them to a single ‘*.RData’ file.

Usage

```
save_rets(
  symbol,
  data_dir = "E:/mktdata/sec/",
  output_dir = "E:/output/data/",
  look_back = 51,
  vol_mult = 2,
  period = "minutes",
  tzzone = "America/New_York"
)
```

Details

The function `save_rets` loads multiple days of *TAQ* data, then scrubs, aggregates, and rbinds them into a *OHLC* time series. It then calculates returns using function `ohlc_returns()`, and stores them in a variable named ‘`symbol.rets`’, and saves them to a file called ‘`symbol.rets.RData`’. The *TAQ* data files are assumed to be stored in separate directories for each ‘`symbol`’. Each ‘`symbol`’ has its own directory (named ‘`symbol`’) in the ‘`data_dir`’ directory. Each ‘`symbol`’ directory contains multiple daily ‘*.RData’ files, each file containing one day of *TAQ* data.

Value

A time series of returns and volume in *xts* format.

Examples

```
## Not run:
save_rets("SPY")

## End(Not run)
```

save_rets_ohlc

Load OHLC time series data for a single symbol, calculate its returns, and save them to a single '.RData' file, without aggregation.*

Description

Load *OHLC* time series data for a single symbol, calculate its returns, and save them to a single '*.RData' file, without aggregation.

Usage

```
save_rets_ohlc(
  symbol,
  data_dir = "E:/output/data/",
  output_dir = "E:/output/data/"
)
```

Details

The function `save_rets_ohlc()` loads *OHLC* time series data from a single file. It then calculates returns using function `ohlc_returns()`, and stores them in a variable named 'symbol.rets', and saves them to a file called 'symbol.rets.RData'.

Value

A time series of returns and volume in *xts* format.

Examples

```
## Not run:
save_rets_ohlc("SPY")

## End(Not run)
```

save_scrub_agg	<i>Load, scrub, aggregate, and rbind multiple days of TAQ data for a single symbol, and save the OHLC time series to a single ‘*.RData’ file.</i>
----------------	---

Description

Load, scrub, aggregate, and rbind multiple days of *TAQ* data for a single symbol, and save the *OHLC* time series to a single ‘*.RData’ file.

Usage

```
save_scrub_agg(
  symbol,
  data_dir = "E:/mktdata/sec/",
  output_dir = "E:/output/data/",
  look_back = 51,
  vol_mult = 2,
  period = "minutes",
  tzzone = "America/New_York"
)
```

Arguments

symbol	A <i>character</i> string representing symbol or ticker.
data_dir	A <i>character</i> string representing directory containing input ‘*.RData’ files.
output_dir	A <i>character</i> string representing directory containing output ‘*.RData’ files.

Details

The function `save_scrub_agg()` loads multiple days of *TAQ* data, then scrubs, aggregates, and rbinds them into a *OHLC* time series, and finally saves it to a single ‘*.RData’ file. The *OHLC* time series is stored in a variable named ‘symbol’, and then it’s saved to a file named ‘symbol.RData’ in the ‘output_dir’ directory. The *TAQ* data files are assumed to be stored in separate directories for each ‘symbol’. Each ‘symbol’ has its own directory (named ‘symbol’) in the ‘data_dir’ directory. Each ‘symbol’ directory contains multiple daily ‘*.RData’ files, each file containing one day of *TAQ* data.

Value

An *OHLC* time series in *xts* format.

Examples

```
## Not run:
# set data directories
data_dir <- "C:/Develop/data/hfreq/src/"
output_dir <- "C:/Develop/data/hfreq/scrub/"
symbol <- "SPY"
# Aggregate SPY TAQ data to 15-min OHLC bar data, and save the data to a file
save_scrub_agg(symbol=symbol, data_dir=data_dir, output_dir=output_dir, period="15 min")

## End(Not run)
```

save_taq	<i>Load and scrub multiple days of TAQ data for a single symbol, and save it to multiple '*.RData' files.</i>
----------	---

Description

Load and scrub multiple days of *TAQ* data for a single symbol, and save it to multiple '*.RData' files.

Usage

```
save_taq(
  symbol,
  data_dir = "E:/mktdata/sec/",
  output_dir = "E:/output/data/",
  look_back = 51,
  vol_mult = 2,
  tzzone = "America/New_York"
)
```

Details

The function `save_taq()` loads multiple days of *TAQ* data, scrubs it, and saves the scrubbed *TAQ* data to individual '*.RData' files. It uses the same file names for output as the input file names. The *TAQ* data files are assumed to be stored in separate directories for each 'symbol'. Each 'symbol' has its own directory (named 'symbol') in the 'data_dir' directory. Each 'symbol' directory contains multiple daily '*.RData' files, each file containing one day of *TAQ* data.

Value

a *TAQ* time series in *xts* format.

Examples

```
## Not run:
save_taq("SPY")

## End(Not run)
```

scrub_agg	<i>Scrub a single day of TAQ data, aggregate it, and convert to OHLC format.</i>
-----------	--

Description

Scrub a single day of *TAQ* data, aggregate it, and convert to *OHLC* format.

Usage

```
scrub_agg(  
  taq,  
  look_back = 51,  
  vol_mult = 2,  
  period = "minutes",  
  tzzone = "America/New_York"  
)
```

Arguments

period The aggregation period.

Details

The function `scrub_agg()` performs:

- index timezone conversion,
- data subset to trading hours,
- removal of duplicate time stamps,
- scrubbing of quotes with suspect bid-offer spreads,
- scrubbing of quotes with suspect price jumps,
- cbinding of mid prices with volume data,
- aggregation to OHLC using function `to.period()` from package *xts*,

Valid 'period' character strings include: "minutes", "3 min", "5 min", "10 min", "15 min", "30 min", and "hours". The time index of the output time series is rounded up to the next integer multiple of 'period'.

Value

A *OHLC* time series in *xts* format.

Examples

```
# Create random TAQ prices  
taq <- HighFreq::random_taq()  
# Aggregate to ten minutes OHLC data  
ohlc <- HighFreq::scrub_agg(taq, period="10 min")  
chart_Series(ohlc, name="random prices")  
# scrub and aggregate a single day of SPY TAQ data to OHLC  
ohlc <- HighFreq::scrub_agg(taq=HighFreq::SPY_TAQ)  
chart_Series(ohlc, name=symbol)
```

scrub_taq	<i>Scrub a single day of TAQ data in xts format, without aggregation.</i>
-----------	---

Description

Scrub a single day of *TAQ* data in *xts* format, without aggregation.

Usage

```
scrub_taq(taq, look_back = 51, vol_mult = 2, tzone = "America/New_York")
```

Arguments

taq	<i>TAQ</i> A time series in <i>xts</i> format.
tzone	The timezone to convert.

Details

The function `scrub_taq()` performs the same scrubbing operations as `scrub_agg`, except it doesn't aggregate, and returns the *TAQ* data in *xts* format.

Value

A *TAQ* time series in *xts* format.

Examples

```
taq <- HighFreq::scrub_taq(taq=HighFreq::SPY_TAQ, look_back=11, vol_mult=1)
# Create random TAQ prices and scrub them
taq <- HighFreq::random_taq()
taq <- HighFreq::scrub_taq(taq=taq)
taq <- HighFreq::scrub_taq(taq=taq, look_back=11, vol_mult=1)
```

season_ality	<i>Perform seasonality aggregations over a single-column xts time series.</i>
--------------	---

Description

Perform seasonality aggregations over a single-column *xts* time series.

Usage

```
season_ality(xtsv, indeks = format(zoo::index(xtsv), "%H:%M"))
```

Arguments

xtsv	A single-column <i>xts</i> time series.
indeks	A vector of <i>character</i> strings representing points in time, of the same length as the argument <i>xtsv</i> .

Details

The function `season_ality()` calculates the mean of values observed at the same points in time specified by the argument `indeks`. An example of a daily seasonality aggregation is the average price of a stock between 9:30AM and 10:00AM every day, over many days. The argument `indeks` is passed into function `tapply()`, and must be the same length as the argument `xtsv`.

Value

An *xts* time series with mean aggregations over the seasonality interval.

Examples

```
# Calculate running variance of each minutely OHLC bar of data
xtsv <- ohlc_variance(HighFreq::SPY)
# Remove overnight variance spikes at "09:31"
indeks <- format(index(xtsv), "%H:%M")
xtsv <- xtsv[!indeks=="09:31", ]
# Calculate daily seasonality of variance
var_seasonal <- season_ality(xtsv=xtsv)
chart_Series(x=var_seasonal, name=paste(colnames(var_seasonal),
  "daily seasonality of variance"))
```

sim_ar

Simulate autoregressive returns by recursively filtering a matrix of innovations through a matrix of autoregressive coefficients.

Description

Simulate *autoregressive* returns by recursively filtering a *matrix* of innovations through a *matrix* of *autoregressive* coefficients.

Usage

```
sim_ar(coeff, innov)
```

Arguments

`innov` A single-column *matrix* of innovations.
`coeff` A single-column *matrix* of *autoregressive* coefficients.

Details

The function `sim_ar()` recursively filters the *matrix* of innovations `innov` through the *matrix* of *autoregressive* coefficients `coeff`, using fast RcppArmadillo C++ code.

The function `sim_ar()` simulates an *autoregressive* process $AR(n)$ of order n :

$$r_i = \varphi_1 r_{i-1} + \varphi_2 r_{i-2} + \dots + \varphi_n r_{i-n} + \xi_i$$

Where r_i is the simulated output time series, φ_i are the *autoregressive* coefficients, and ξ_i are the standard normal *innovations*.

The order n of the *autoregressive* process $AR(n)$, is equal to the number of rows of the *autoregressive* coefficients `coeff`.

The function `sim_ar()` performs the same calculation as the standard R function `filter(x=innov, filter=coeff, method="recursive")`, but it's several times faster.

Value

A single-column *matrix* of simulated returns, with the same number of rows as the argument `innov`.

Examples

```
## Not run:
# Define AR coefficients
coeff <- matrix(c(0.1, 0.3, 0.5))
# Calculate matrix of innovations
innov <- matrix(rnorm(1e4, sd=0.01))
# Calculate recursive filter using filter()
filtered <- filter(innov, filter=coeff, method="recursive")
# Calculate recursive filter using RcppArmadillo
returns <- HighFreq::sim_ar(coeff, innov)
# Compare the two methods
all.equal(as.numeric(returns), as.numeric(filtered))
# Compare the speed of RcppArmadillo with R code
library(microbenchmark)
summary(microbenchmark(
  Rcpp=HighFreq::sim_ar(coeff, innov),
  Rcode=filter(innov, filter=coeff, method="recursive"),
  times=10))[, c(1, 4, 5)] # end microbenchmark summary

## End(Not run)
```

sim_df

*Simulate a Dickey-Fuller process using Rcpp.***Description**

Simulate a *Dickey-Fuller* process using *Rcpp*.

Usage

```
sim_df(init_price, eq_price, theta, coeff, innov)
```

Arguments

<code>init_price</code>	The initial price.
<code>eq_price</code>	The equilibrium price.
<code>theta</code>	The strength of mean reversion.
<code>coeff</code>	A single-column <i>matrix</i> of <i>autoregressive</i> coefficients.
<code>innov</code>	A single-column <i>matrix</i> of innovations (random numbers).

Details

The function `sim_df()` simulates the following *Dickey-Fuller* process:

$$r_i = \theta (\mu - p_{i-1}) + \varphi_1 r_{i-1} + \dots + \varphi_n r_{i-n} + \xi_i$$

$$p_i = p_{i-1} + r_i$$

Where r_i and p_i are the simulated returns and prices, θ and μ are the *Ornstein-Uhlenbeck* parameters, φ_i are the *autoregressive* coefficients, and ξ_i are the normal *innovations*. The recursion starts with: $r_1 = \xi_1$ and $p_1 = \text{init_price}$.

The *Dickey-Fuller* process is a combination of an *Ornstein-Uhlenbeck* process and an *autoregressive* process. The order n of the *autoregressive* process $AR(n)$, is equal to the number of rows of the *autoregressive* coefficients `coeff`.

The function `sim_df()` simulates the *Dickey-Fuller* process using fast *Rcpp* C++ code.

The function `sim_df()` returns a single-column *matrix* representing the *time series* of prices.

Value

A single-column *matrix* of simulated prices, with the same number of rows as the argument `innov`.

Examples

```
## Not run:
# Define the Ornstein-Uhlenbeck model parameters
init_price <- 1.0
eq_price <- 2.0
thetav <- 0.01
# Define AR coefficients
coeff <- matrix(c(0.1, 0.3, 0.5))
# Calculate matrix of standard normal innovations
innov <- matrix(rnorm(1e3, sd=0.01))
# Simulate Dickey-Fuller process using Rcpp
prices <- HighFreq::sim_df(init_price=init_price, eq_price=eq_price, theta=thetav, coeff=coeff, innov=innov)
plot(prices, t="l", main="Simulated Dickey-Fuller Prices")

## End(Not run)
```

sim_garch	<i>Simulate or estimate the rolling variance under a GARCH(1,1) process using Rcpp.</i>
-----------	---

Description

Simulate or estimate the rolling variance under a *GARCH(1,1)* process using *Rcpp*.

Usage

```
sim_garch(omega, alpha, beta, innov, is_random = TRUE)
```

Arguments

<code>omega</code>	Parameter proportional to the long-term average level of variance.
<code>alpha</code>	The weight associated with recent realized variance updates.
<code>beta</code>	The weight associated with the past variance estimates.
<code>innov</code>	A single-column <i>matrix</i> of innovations.
<code>is_random</code>	<i>Boolean</i> argument: Are the innovations random numbers or historical returns? (The default is <code>is_random = TRUE</code> .)

Details

The function `sim_garch()` simulates or estimates the rolling variance under a $GARCH(1,1)$ process using *Rcpp*.

If `is_random = TRUE` (the default) then the innovations `innov` are treated as random numbers ξ_i and the $GARCH(1,1)$ process is given by:

$$r_i = \sigma_{i-1} \xi_i$$

$$\sigma_i^2 = \omega + \alpha r_i^2 + \beta \sigma_{i-1}^2$$

Where r_i and σ_i^2 are the simulated returns and variance, and ω , α , and β are the $GARCH$ parameters, and ξ_i are standard normal *innovations*.

The long-term equilibrium level of the simulated variance is proportional to the parameter ω :

$$\sigma^2 = \frac{\omega}{1 - \alpha - \beta}$$

So the sum of α plus β should be less than 1, otherwise the volatility becomes explosive.

If `is_random = FALSE` then the function `sim_garch()` *estimates* the rolling variance from the historical returns. The innovations `innov` are equal to the historical returns r_i and the $GARCH(1,1)$ process is simply:

$$\sigma_i^2 = \omega + \alpha r_i^2 + \beta \sigma_{i-1}^2$$

Where σ_i^2 is the rolling variance.

The above should be viewed as a formula for *estimating* the rolling variance from the historical returns, rather than simulating them. It represents exponential smoothing of the squared returns with a decay factor equal to β .

The function `sim_garch()` simulates the $GARCH$ process using fast *Rcpp* C++ code.

Value

A *matrix* with two columns and with the same number of rows as the argument `innov`. The first column are the simulated returns and the second column is the variance.

Examples

```
## Not run:
# Define the GARCH model parameters
alpha <- 0.79
betav <- 0.2
om_ega <- 1e-4*(1-alpha-betav)
# Calculate matrix of standard normal innovations
innov <- matrix(rnorm(1e3))
# Simulate the GARCH process using Rcpp
garch_data <- HighFreq::sim_garch(omega=om_ega, alpha=alpha, beta=betav, innov=innov)
# Plot the GARCH rolling volatility and cumulative returns
plot(sqrt(garch_data[, 2]), t="l", main="Simulated GARCH Volatility", ylab="volatility")
plot(cumsum(garch_data[, 1]), t="l", main="Simulated GARCH Cumulative Returns", ylab="cumulative returns")
# Calculate historical VTI returns
returns <- na.omit(rutils::etfenv$returns$VTI)
# Estimate the GARCH volatility of VTI returns
garch_data <- HighFreq::sim_garch(omega=om_ega, alpha=alpha, beta=betav,
  innov=returns, is_random=FALSE)
# Plot dygraph of the estimated GARCH volatility
dygraphs::dygraph(xts::xts(sqrt(garch_data[, 2]), index(returns)),
  main="Estimated GARCH Volatility of VTI")
```

```
## End(Not run)
```

sim_ou

Simulate an Ornstein-Uhlenbeck process using Rcpp.

Description

Simulate an *Ornstein-Uhlenbeck* process using *Rcpp*.

Usage

```
sim_ou(init_price, eq_price, theta, innov)
```

Arguments

init_price	The initial price.
eq_price	The equilibrium price.
theta	The strength of mean reversion.
innov	A single-column <i>matrix</i> of innovations (random numbers).

Details

The function `sim_ou()` simulates the following *Ornstein-Uhlenbeck* process:

$$r_i = p_i - p_{i-1} = \theta (\mu - p_{i-1}) + \xi_i$$

$$p_i = p_{i-1} + r_i$$

Where r_i and p_i are the simulated returns and prices, θ , μ , and σ are the *Ornstein-Uhlenbeck* parameters, and ξ_i are the standard *innovations*. The recursion starts with: $r_1 = \xi_1$ and $p_1 = \text{init_price}$.

The function `sim_ou()` simulates the percentage returns as equal to the difference between the equilibrium price μ minus the latest price p_{i-1} , times the mean reversion parameter θ , plus a random normal innovation. The log prices are calculated as the sum of returns (not compounded), so they can become negative.

The function `sim_ou()` simulates the *Ornstein-Uhlenbeck* process using fast *Rcpp* C++ code.

The function `sim_ou()` returns a single-column *matrix* representing the *time series* of simulated prices.

Value

A single-column *matrix* of simulated prices, with the same number of rows as the argument `innov`.

Examples

```
## Not run:
# Define the Ornstein-Uhlenbeck model parameters
init_price <- 0.0
eq_price <- 1.0
sigmav <- 0.01
thetav <- 0.01
innov <- matrix(rnorm(1e3))
# Simulate Ornstein-Uhlenbeck process using Rcpp
prices <- HighFreq::sim_ou(init_price=init_price, eq_price=eq_price, volat=sigmav, theta=thetav, innov=innov)
plot(prices, t="l", main="Simulated Ornstein-Uhlenbeck Prices", ylab="prices")

## End(Not run)
```

sim_schwartz	<i>Simulate a Schwartz process using Rcpp.</i>
--------------	--

Description

Simulate a *Schwartz* process using *Rcpp*.

Usage

```
sim_schwartz(init_price, eq_price, theta, innov)
```

Arguments

init_price	The initial price.
eq_price	The equilibrium price.
theta	The strength of mean reversion.
innov	A single-column <i>matrix</i> of innovations (random numbers).

Details

The function `sim_schwartz()` simulates a *Schwartz* process using fast *Rcpp* C++ code.

The *Schwartz* process is the exponential of the *Ornstein-Uhlenbeck* process, and similar comments apply to it. The prices are calculated as the exponentially compounded returns, so they are never negative. The log prices can be obtained by taking the logarithm of the prices.

The function `sim_schwartz()` simulates the percentage returns as equal to the difference between the equilibrium price μ minus the latest price p_{i-1} , times the mean reversion parameter θ , plus a random normal innovation.

The function `sim_schwartz()` returns a single-column *matrix* representing the *time series* of simulated prices.

Value

A single-column *matrix* of simulated prices, with the same number of rows as the argument `innov`.

Examples

```
## Not run:
# Define the Schwartz model parameters
init_price <- 1.0
eq_price <- 2.0
thetav <- 0.01
innov <- matrix(rnorm(1e3, sd=0.01))
# Simulate Schwartz process using Rcpp
prices <- HighFreq::sim_schwartz(init_price=init_price, eq_price=eq_price, theta=thetav, innov=innov)
plot(prices, t="l", main="Simulated Schwartz Prices", ylab="prices")

## End(Not run)
```

which_extreme	<i>Calculate a Boolean vector that identifies extreme tail values in a single-column xts time series or vector, over a rolling look-back interval.</i>
---------------	--

Description

Calculate a *Boolean* vector that identifies extreme tail values in a single-column *xts* time series or vector, over a rolling look-back interval.

Usage

```
which_extreme(xtsv, look_back = 51, vol_mult = 2)
```

Arguments

xtsv	A single-column <i>xts</i> time series, or a <i>numeric</i> or <i>Boolean</i> vector.
look_back	The number of data points in rolling look-back interval for estimating rolling quantile.
vol_mult	The quantile multiplier.

Details

The function `which_extreme()` calculates a *Boolean* vector, with TRUE for values that belong to the extreme tails of the distribution of values.

The function `which_extreme()` applies a version of the Hampel median filter to identify extreme values, but instead of using the median absolute deviation (MAD), it uses the 0.9 quantile values calculated over a rolling look-back interval.

Extreme values are defined as those that exceed the product of the multiplier times the rolling quantile. Extreme values belong to the fat tails of the recent (trailing) distribution of values, so they are present only when the trailing distribution of values has fat tails. If the trailing distribution of values is closer to normal (without fat tails), then there are no extreme values.

The quantile multiplier `vol_mult` controls the threshold at which values are identified as extreme. Smaller quantile multiplier values will cause more values to be identified as extreme.

Value

A *Boolean* vector with the same number of rows as the input time series or vector.

Examples

```
# Create local copy of SPY TAQ data
taq <- HighFreq::SPY_TAQ
# scrub quotes with suspect bid-offer spreads
bid_offer <- taq[, "Ask.Price"] - taq[, "Bid.Price"]
sus_pect <- which_extreme(bid_offer, look_back=51, vol_mult=3)
# Remove suspect values
taq <- taq[!sus_pect]
```

which_jumps	<i>Calculate a Boolean vector that identifies isolated jumps (spikes) in a single-column xts time series or vector, over a rolling interval.</i>
-------------	--

Description

Calculate a *Boolean* vector that identifies isolated jumps (spikes) in a single-column *xts* time series or vector, over a rolling interval.

Usage

```
which_jumps(xtsv, look_back = 51, vol_mult = 2)
```

Details

The function `which_jumps()` calculates a *Boolean* vector, with TRUE for values that are isolated jumps (spikes).

The function `which_jumps()` applies a version of the Hampel median filter to identify jumps, but instead of using the median absolute deviation (MAD), it uses the 0.9 quantile of returns calculated over a rolling interval. This is in contrast to function `which_extreme()`, which applies a Hampel filter to the values themselves, instead of the returns. Returns are defined as simple differences between neighboring values.

Jumps (or spikes), are defined as isolated values that are very different from the neighboring values, either before or after. Jumps create pairs of large neighboring returns of opposite sign.

Jumps (spikes) must satisfy two conditions:

1. Neighboring returns both exceed a multiple of the rolling quantile,
2. The sum of neighboring returns doesn't exceed that multiple.

The quantile multiplier `vol_mult` controls the threshold at which values are identified as jumps. Smaller quantile multiplier values will cause more values to be identified as jumps.

Value

A *Boolean* vector with the same number of rows as the input time series or vector.

Examples

```
# Create local copy of SPY TAQ data
taq <- SPY_TAQ
# Calculate mid prices
mid_prices <- 0.5 * (taq[, "Bid.Price"] + taq[, "Ask.Price"])
# Replace whole rows containing suspect price jumps with NA, and perform locf()
taq[which_jumps(mid_prices, look_back=31, vol_mult=1.0), ] <- NA
taq <- xts::na.locf.xts(taq)
```