

FRE7241 Algorithmic Portfolio Management

Lecture#1, Spring 2023

Jerzy Pawlowski jp3900@nyu.edu

NYU Tandon School of Engineering

March 21, 2023



NYU

**TANDON SCHOOL
OF ENGINEERING**

Welcome Students!

My name is Jerzy Pawlowski jp3900@nyu.edu

I'm an adjunct professor at NYU Tandon because I love teaching and I want to share my professional knowledge with young, enthusiastic students.

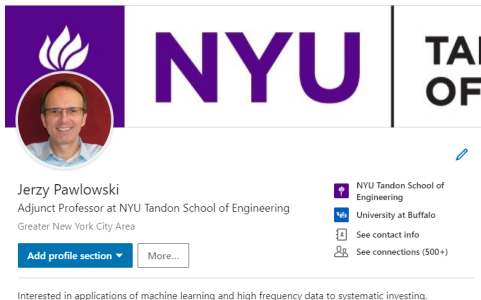
I'm interested in applications of *machine learning* to *systematic investing*.

I'm an advocate of *open-source software*, and I share it on GitHub:

[My GitHub account](#)

In my finance career, I have worked as a hedge fund *portfolio manager*, *CLO structurer* (banker), and *quant analyst*.

[My LinkedIn profile](#)

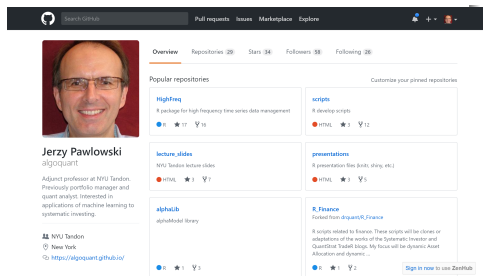


Jerzy Pawlowski
Adjunct Professor at NYU Tandon School of Engineering
Greater New York City Area

[Add profile section](#) [More...](#)

[NYU Tandon School of Engineering](#)
[University at Buffalo](#)
[See contact info](#)
[See connections \(500+\)](#)

Interested in applications of machine learning and high frequency data to systematic investing.



Jerzy Pawlowski
algoquant

Adjunct professor at NYU Tandon. Previously portfolio manager and quant analyst. Interested in applications of machine learning to systematic investing.

[NYU Tandon](#)
[New York](#)
<https://algoquant.github.io/>

Overview Repositories 20 Stars 34 Followers 58 Following 26

Popular repositories

Repository	Stars	Followers	Following
HighFreq A package for high frequency time series data management	17	17	17
lecture_slides NYU Tandon lecture slides	3	3	3
alphanlib alphanlib library	1	1	1
scripts A develop scripts	12	12	12
presentations A presentation files (pdfs, shpg, etc.)	5	5	5
R_Finance R scripts related to Finance. These scripts will be clones or adaptations of the works of the Systematic Investor and QuantGest Trade bings. My focus will be dynamic Asset Allocation and dynamic ...	1	1	1

[Sign in now to use ZenHub](#)

FRE7241 Course Description and Objectives

Course Description

The course will apply the R programming language to *trend following*, *momentum trading*, *statistical arbitrage* (pairs trading), and other active portfolio management strategies. The course will implement volatility and price *forecasting models*, asset pricing and *factor models*, and *portfolio optimization*. The course will apply *machine learning* techniques, such as *parameter regularization* (shrinkage), *bagging* and *backtesting* (cross-validation).

FRE7241 Course Description and Objectives

Course Description

The course will apply the R programming language to *trend following*, *momentum trading*, *statistical arbitrage* (pairs trading), and other active portfolio management strategies. The course will implement volatility and price *forecasting models*, asset pricing and *factor models*, and *portfolio optimization*. The course will apply *machine learning* techniques, such as *parameter regularization* (shrinkage), *bagging* and *backtesting* (cross-validation).

Course Objectives

Students will learn through R coding exercises how to:

- download data from external sources, and to scrub and format it.
- estimate time series parameters, and fit models such as *ARIMA*, *GARCH*, and factor models.
- optimize portfolios under different constraints and risk-return objectives.
- backtest active portfolio management strategies and evaluate their performance.

FRE7241 Course Description and Objectives

Course Description

The course will apply the R programming language to *trend following*, *momentum trading*, *statistical arbitrage* (pairs trading), and other active portfolio management strategies. The course will implement volatility and price forecasting models, asset pricing and *factor models*, and *portfolio optimization*. The course will apply *machine learning* techniques, such as *parameter regularization* (shrinkage), *bagging* and *backtesting* (cross-validation).

Course Objectives

Students will learn through R coding exercises how to:

- download data from external sources, and to scrub and format it.
- estimate time series parameters, and fit models such as *ARIMA*, *GARCH*, and factor models.
- optimize portfolios under different constraints and risk-return objectives.
- backtest active portfolio management strategies and evaluate their performance.

Course Prerequisites

FRE6123 Financial Risk Management and Asset Pricing. The R language is considered to be challenging, so this course requires programming experience with other languages such as C++ or Python. Students with less programming experience are encouraged to first take *FRE6871 R in Finance*, and also *FRE6883 Financial Computing* by prof. Song Tang. Students should also have knowledge of basic statistics (random variables, estimators, hypothesis testing, regression, etc.)

Homeworks and Tests

Homeworks and Tests

Grading will be based on homeworks and tests. There will be no final exam.

The tests will require writing code, which should run directly when pasted into an R session, and should produce the required output, without any modifications.

Students will be allowed to consult lecture slides, and to copy code from them, and to copy from books or any online sources, but they will be required to provide references to those external sources (such as links or titles and page numbers).

The tests will be closely based on code contained in the lecture slides, so students are encouraged to become very familiar with those slides.

Students will submit their homework and test files only through *Brightspace* (not emails).

Students will be required to bring their laptop computers to class and run the R Interpreter, and the RStudio Integrated Development Environment (*IDE*), during the lecture.

Homeworks will also include reading assignments designed to help prepare for tests.

Homeworks and Tests

Homeworks and Tests

Grading will be based on homeworks and tests. There will be no final exam.

The tests will require writing code, which should run directly when pasted into an R session, and should produce the required output, without any modifications.

Students will be allowed to consult lecture slides, and to copy code from them, and to copy from books or any online sources, but they will be required to provide references to those external sources (such as links or titles and page numbers).

The tests will be closely based on code contained in the lecture slides, so students are encouraged to become very familiar with those slides.

Students will submit their homework and test files only through *Brightspace* (not emails).

Students will be required to bring their laptop computers to class and run the R Interpreter, and the RStudio Integrated Development Environment (*IDE*), during the lecture.

Homeworks will also include reading assignments designed to help prepare for tests.

Graduate Assistant

The graduate assistant (GA) will be Wenxin Li wl2570@nyu.edu.

The GA will answer questions during office hours, or via *Brightspace* forums, not via emails. Please send emails regarding lecture matters from *Brightspace* (not personal emails).

Tips for Solving Homeworks and Tests

Tips for Solving Homeworks and Tests

The tests will require mostly copying code samples from the lecture slides, making some modifications to them, and combining them with other code samples.

Partial credit will be given even for code that doesn't produce the correct output, but that has elements of code that can be useful for producing the right answer.

So don't leave test assignments unanswered, and instead copy any code samples from the lecture slides that are related to the solution and make sense.

Contact the GA during office hours via text or phone, and submit questions to the GA or to me via *Brightspace*.

Tips for Solving Homeworks and Tests

Tips for Solving Homeworks and Tests

The tests will require mostly copying code samples from the lecture slides, making some modifications to them, and combining them with other code samples.

Partial credit will be given even for code that doesn't produce the correct output, but that has elements of code that can be useful for producing the right answer.

So don't leave test assignments unanswered, and instead copy any code samples from the lecture slides that are related to the solution and make sense.

Contact the GA during office hours via text or phone, and submit questions to the GA or to me via *Brightspace*.

Please Submit *Minimal Working Examples* With Your Questions

When submitting questions, please provide a *minimal working example* that produces the error in R, with the following items:

- The *complete* R code that produces the error, including the seed value for random numbers,
- The version of R (output of command: `sessionInfo()`), and the versions of R packages,
- The type and version of your operating system (Windows or OSX),
- The dataset file used by the R code,
- The text or screenshots of error messages,

You can read more about producing *minimal working examples* here: <http://stackoverflow.com/help/mcve>
<http://www.jaredknowles.com/journal/2013/5/27/writing-a-minimal-working-example-mwe-in-r>

Course Grading Policies

Numerical Scores

Tests will be graded and assigned numerical scores. Each part of the tests will be graded separately and assigned a numerical score.

Maximum scores will be given only for complete code, that produces the correct output when it's pasted into an R session, without any modifications. As long as the R code uses the required functions and produces the correct output, it will be given full credit.

Partial credit will be given even for code that doesn't produce the correct output, but that has elements of code that can be useful for producing the right answer.

Course Grading Policies

Numerical Scores

Tests will be graded and assigned numerical scores. Each part of the tests will be graded separately and assigned a numerical score.

Maximum scores will be given only for complete code, that produces the correct output when it's pasted into an R session, without any modifications. As long as the R code uses the required functions and produces the correct output, it will be given full credit.

Partial credit will be given even for code that doesn't produce the correct output, but that has elements of code that can be useful for producing the right answer.

Letter Grades

Letter grades for the course will be derived from the cumulative scores obtained for all the tests. Very high numerical scores close to the maximum won't guarantee an A letter grade, since grading will also depend on the difficulty of the assignments.

Course Grading Policies

Numerical Scores

Tests will be graded and assigned numerical scores. Each part of the tests will be graded separately and assigned a numerical score.

Maximum scores will be given only for complete code, that produces the correct output when it's pasted into an R session, without any modifications. As long as the R code uses the required functions and produces the correct output, it will be given full credit.

Partial credit will be given even for code that doesn't produce the correct output, but that has elements of code that can be useful for producing the right answer.

Letter Grades

Letter grades for the course will be derived from the cumulative scores obtained for all the tests. Very high numerical scores close to the maximum won't guarantee an A letter grade, since grading will also depend on the difficulty of the assignments.

Plagiarism

Plagiarism (copying from other students) and cheating will be punished.

But copying code from lecture slides, books, or any online sources is allowed and encouraged.

Students must provide references to any external sources from which they copy code (such as links or titles and page numbers).

FRE7241 Course Materials

Lecture Slides

The course will be mostly self-contained, using detailed lecture slides containing extensive, working R code examples.

The course will also utilize data and tutorials which are freely available on the internet.

FRE7241 Course Materials

Lecture Slides

The course will be mostly self-contained, using detailed lecture slides containing extensive, working R code examples.

The course will also utilize data and tutorials which are freely available on the internet.

FRE7241 Recommended Textbooks

- *Advances in Financial Machine Learning* by Marcos Lopez de Prado, great overview of machine learning techniques applied to quant trading.
- *Financial Data and Models Using R* by Clifford Ang, provides a good introduction to time series, portfolio optimization, and performance measures.
- *Systematic Trading* by Rob Carver, explains practical systematic trading rules.
- *Automated Trading* by Chris Conlan, explains how to implement a practical computer trading system.
- *Statistics and Data Analysis for Financial Engineering* by David Ruppert, introduces regression, cointegration, multivariate time series analysis, *ARIMA*, *GARCH*, *CAPM*, and factor models, with examples in R.
- *Financial Risk Modelling and Portfolio Optimization with R* by Bernhard Pfaff, introduces volatility models, portfolio optimization, and tactical asset allocation, with a great review of R packages and examples in R.

Many textbooks can be downloaded in electronic format from the [NYU Library](#).

FRE7241 Supplementary Books

- *Introduction to Statistical Learning* by Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani, introduces machine learning techniques using R, but without deep learning.
- *Quantitative Risk Management* by Alexander J. McNeil, Rudiger Frey, and Paul Embrechts: review of Value at Risk, factor models, ARMA and GARCH, extreme value theory, and credit risk models.
- *Applied Econometrics with R* by Christian Kleiber and Achim Zeileis, introduces advanced statistical models and econometrics.
- *The Art of R Programming* by Norman Matloff, contains a good introduction to R and to statistical models.
- *Advanced R* by Hadley Wickham, is the best book for learning the advanced features of R.
- *Numerical Recipes in C++* by William Press, Saul Teukolsky, William Vetterling, and Brian Flannery, is a great reference for linear algebra and numerical methods, implemented in working C++ code.
- The books *R in Action* by Robert Kabacoff and *R for Everyone* by Jared Lander, are good introductions to R and to statistical models.
- *Quant Finance books* by Jerzy Pawlowski.
- *Quant Trading books* by Jerzy Pawlowski.

FRE7241 Supplementary Materials

Robert Carver's trading blog

Great blog about practical systematic trading and investments, with Python code: <http://qoppac.blogspot.com/>

Introduction to Computational Finance with R

Good course by prof. Eric Zivot, with lots of R examples:

<https://www.datacamp.com/community/open-courses/computational-finance-and-financial-econometrics-with-r>

Notepad++ is a free source code editor for MS Windows, that supports several programming languages, including R.

Notepad++ has a very convenient and fast *search and replace* function, that allows *search and replace* in multiple files.

<http://notepad-plus-plus.org/>



Internal R Help and Documentation

The function `help()` displays documentation on a function or subject.

Preceding the keyword with a single "?" is equivalent to calling `help()`.

```
> # Display documentation on function "getwd"
> help(getwd)
> # Equivalent to "help(getwd)"
> ?getwd
```

The function `help.start()` displays a page with links to internal documentation.

```
> # Open the hypertext documentation
> help.start()
```

R documentation is also available in RGui under the help tab.

The *pdf* files with R documentation are also available directly under:

<C:/Program Files/R/R-3.1.2/doc/manual/>
(the exact path will depend on the R version.)



Introduction to R by Venables and R Core Team.

R Style Guides

DataCamp R style guide

The DataCamp R style guide is very close to what I have adopted:
[DataCamp R style guide](#)

Google R style guide

The Google R style guide is similar to DataCamp's:
[Google R style guide](#)

Stack Exchange

Stack Overflow

Stack Overflow is a Q&A forum for computer programming, and is part of Stack Exchange

<http://stackoverflow.com>

<http://stackoverflow.com/questions/tagged/r>

<http://stackoverflow.com/tags/r/info>

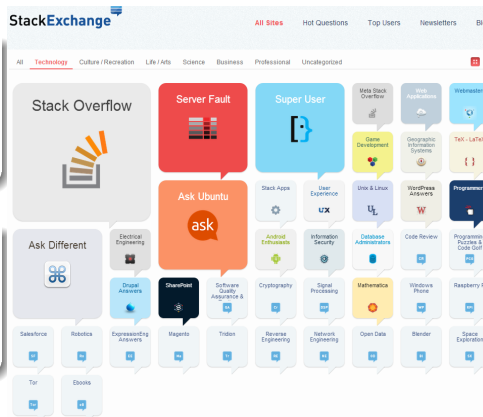
Stack Exchange

Stack Exchange is a family of Q&A forums in a variety of fields

<http://stackexchange.com/>

<http://stackexchange.com/sites#technology>

<http://quant.stackexchange.com/>



RStudio Support

RStudio has extensive online help, Q&A database, and documentation

<https://support.rstudio.com/hc/en-us>

<https://support.rstudio.com/hc/en-us/sections/200107586-Using-RStudio>

<https://support.rstudio.com/hc/en-us/sections/200148796-Advanced-Topics>

R Online Books and References

Hadley Wickham book *Advanced R*

The best book for learning the advanced features of R: <http://adv-r.had.co.nz/>

Cookbook for R by Winston Chang from *RStudio*

Good plotting, but not interactive: <http://www.cookbook-r.com/>

Efficient R programming by Colin Gillespie and Robin Lovelace

Good tips for fast R programming: <https://csgillespie.github.io/efficientR/programming.html>

Endmemo web book

Good, but not interactive: <http://www.endmemo.com/program/R/>

Quick-R by Robert Kabacoff

Good, but not interactive: <http://www.statmethods.net/>

R for Beginners by Emmanuel Paradis

Good, basic introduction to R: http://cran.r-project.org/doc/contrib/Paradis-rdebuts_en.pdf

R Online Interactive Courses

Datacamp Interactive Courses

Datacamp introduction to R: <https://www.datacamp.com/courses/introduction-to-r/>

Datacamp list of free courses: <https://www.datacamp.com/community/open-courses>

Datacamp basic statistics in R: <https://www.datacamp.com/community/open-courses/basic-statistics>

Datacamp computational finance in R:

<https://www.datacamp.com/community/open-courses/computational-finance-and-financial-econometrics-with-r>

Datacamp machine learning in R:

<https://www.datacamp.com/community/open-courses/kaggle-r-tutorial-on-machine-learning>

Try R

Interactive R tutorial, but rather basic: <http://tryr.codeschool.com/>

R Blogs and Experts

R-Bloggers

R-Bloggers is an aggregator of blogs dedicated to R

<http://www.r-bloggers.com/>

Tal Galili is the author of R-Bloggers and has his own excellent blog

<http://www.r-statistics.com/>

Dirk Eddebuettel

Dirk is a *Top Answerer* for R questions on Stackoverflow, the author of the Rcpp package, and the CRAN Finance View

<http://dirk.eddebuettel.com/>

<http://dirk.eddebuettel.com/code/>

<http://dirk.eddebuettel.com/blog/>

<http://www.rinfinance.com/>

Romain Francois

Romain is an R Enthusiast and Rcpp Hero

<http://romainfrancois.blog.free.fr/>

<http://romainfrancois.blog.free.fr/index.php?tag/graphgallery>

<http://blog.r-enthusiasts.com/>

More R Blogs and Experts

Revolution Analytics Blog

R blog by Revolution Analytics software vendor

<http://blog.revolutionanalytics.com/>

RStudio Blog

R blog by *RStudio*

<http://blog.rstudio.org/>

GitHub for Hosting Software Projects Online

GitHub is an internet-based online service for hosting repositories of software projects.

GitHub provides version control using *git* (desved by Linus Torvalds).

Most R projects are now hosted on *GitHub*.

Google uses *GitHub* to host its *tensorflow* library for machine learning:

<https://github.com/tensorflow/tensorflow>

All the *FRE-7241* and *FRE-6871* lectures are hosted on *GitHub*:

https://github.com/algoquant/lecture_slides

<https://github.com/algoquant>

Hosting projects on *Google* is a great way to advertize your skills and network with experts.

The screenshot shows the GitHub profile of Jerzy Pawlowski (username: algoquant). The profile includes a bio stating he is an adjunct professor at NYU Tandon, previously a portfolio manager and quant analyst, and is interested in machine learning for systematic investing. His location is New York, and his website is <https://algoquant.github.io/>. The 'Popular repositories' section lists several projects:

- HighFreq**: R package for high-frequency time series data management. 17 stars, 15 forks.
- lecture_slides**: NYU Tandon lecture slides. 3 stars, 1 fork.
- alphaHub**: alphahub library. 1 star, 1 fork.
- scripts**: R develop scripts. 12 stars, 1 fork.
- presentations**: R presentation files (pdf, shap, etc.). 5 stars, 1 fork.
- R_Finance**: Scripts related to finance. Forked from [dqquantR_Finance](#). 1 star, 1 fork.

What is R?

- An open-source software environment for statistical computing and graphics.
- An interpreted language, that allows interactive code development.
- A functional language where every operator is an R function.
- A very expressive language that can perform complex operations with very few lines of code.
- A language with metaprogramming facilities that allow programming on the language.
- A language written in C/C++, which can easily call other C/C++ programs.
- Can be easily extended with *packages* (function libraries), providing the latest developments like *Machine Learning*.
- Supports object-oriented programming with *classes* and *methods*.
- Vectorized functions written in C/C++, allow very fast execution of loops over vector elements.



Why is R More Difficult Than Other Languages?

R is more difficult than other languages because:

- R is a *functional* language, which makes its syntax unfamiliar to users of procedural languages like C/C++.
- The huge number of user-created *packages* makes it difficult to tell which are the best for particular applications.
- R can produce very cryptic *warning* and *error* messages, because it's a programming environment, so it performs many operations quietly, but those can sometimes fail.
- Fixing errors usually requires analyzing the complex structure of the R programming environment.

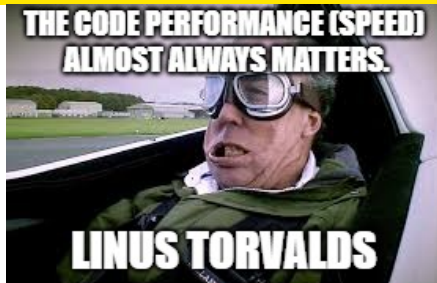


This course is designed to teach the most useful elements of R for financial analysis, through case studies and examples,

What are the Best Ways to Use R?

If used properly, R can be fast and interactive:

- Use R as an interface to libraries written in C++, Java, and JavaScript.
- Avoid using too many R function calls (every command in R is a function).
- Avoid using `apply()` and `for()` loops for large datasets.
- Use R functions which are *compiled* C++ code, instead of using interpreted R code.
- Use package *data.table* for high performance data management.
- Use package *shiny* for interactive charts of live models running in R.
- Use package *dygraphs* for interactive time series plots.
- Use package *knitr* for *RMarkdown* documents.
- Pre-allocate memory for new objects.
- Write C++ functions in *Rcpp* and *RcppArmadillo*.



```
> # Calculate cumulative sum of a vector
> vectorv <- runif(1e5)
> # Use compiled function
> cumsumv <- cumsum(vectorv)
> # Use for loop
> cumsumv2 <- vectorv
> for (i in 2:NROW(vectorv))
+   cumsumv2[i] <- (vectorv[i] + cumsumv2[i-1])
> # Compare the two methods
> all.equal(cumsumv, cumsumv2)
> # Microbenchmark the two methods
> library(microbenchmark)
> summary(microbenchmark(
+   cumsum=cumsum(vectorv),
+   loop_alloc={
+     cumsumv2 <- vectorv
+     for (i in 2:NROW(vectorv))
+       cumsumv2[i] <- (vectorv[i] + cumsumv2[i-1])
+   },
+   loop_nalloc={
+     # Doesn't allocate memory to cumsumv3
```

The R License

R is open-source software released under the GNU General Public License:

<http://www.r-project.org/Licenses>



Some other R packages are released under the Creative Commons Attribution-ShareAlike License:

<http://creativecommons.org>



Installing R and *RStudio*

Students will be required to bring their laptop computers to all the lectures, and to run the R Interpreter and **RStudio** RStudio during the lecture.

Laptop computers will be necessary for following the lectures, and for performing tests.

Students will be required to install and to become proficient with the R Interpreter.

Students can download the R Interpreter from CRAN (Comprehensive R Archive Network):

<http://cran.r-project.org/>

To invoke the RGui interface, click on:

<C:/Program Files/R/R-3.1.2/bin/x64/RGui.exe>



Students will be required to install and to become proficient with the *RStudio* Integrated Development Environment (*IDE*),

<http://www.rstudio.com/products/rstudio/>



Using RStudio

The screenshot displays the RStudio IDE interface. The top menu bar includes File, Edit, Code, View, Plots, Session, Project, Build, Tools, and Help. Below the menu is a toolbar with icons for file operations and running code. The main editor window shows a script with R code for portfolio optimization, including data loading, model fitting, and saving results. The console at the bottom shows warnings about internet connectivity and successful package installation. The right-hand pane displays the 'Workspace' and 'History' tabs, with the 'Packages' tab active, showing the 'library' package and its description.

```

2087 # Run quasi-CEP mode
2088 cep.ticks <- 0:100 # number of ticks cut off from tail
2089 n.buffer <- 500 # buffer size of ticks fed into model
2090 model.cep <- model.test
2091 ts.prices <- model.test$prices
2092 cep.signals <- sapply(cep.ticks, function(cep.tick)
2093 {
2094   cep.prices <- tail(last(ts.prices, cep.tick), n.buffer)
2095   model.cep <- update.alphaModel(model=model.cep, ts.prices=cep.prices)
2096   model.cep <- recalc.alphaModel(model.cep)
2097   as.vector(last(model.cep$signals))
2098 })
2099 write.csv(cep.signals, "S:/Data/R_Data/signals.cep.csv")
2100 write.csv(model.test$signals, "S:/Data/R_Data/signals.csv")
2101
2102
2103
2104 #####
2105 ## Portfolio Optimization ##
2106 #####
2107 library(DEoptim)
2108
2109 ## Load data
2110 stock.sectors.prices <- read.csv(paste(alpha.dir, "stock_sectors.csv", sep=""), stringsAsFactors = FALSE)
2111 stock.sectors.prices <- xts(stock.sectors.prices[, -1], order.by=as.POSIXt(stock.sectors.prices[, 1]))
2112 ts.rets <- diff(stock.sectors.prices, lag=1)
2113 ts.rets[1,] <- ts.rets[2,]
2114 <

```

```

Warning in install.packages :
  InternetOpenUrl failed: 'A connection with the server could not be established'
Warning in install.packages :
  InternetOpenUrl failed: 'A connection with the server could not be established'
Warning in install.packages :
  unable to access index for repository http://www.stats.ox.ac.uk/pub/Rwin/bin/windows/contrib/3.0
Installing package into 'C:/Users/Jerzy/Documents/R/win-library/3.0'
(as 'lib' is unspecified)
trying URL 'http://R-Forge.R-project.org/bin/windows/contrib/3.0/PerformanceAnalytics_1.1.2.zip'
Content type 'application/zip' length 2205138 bytes (2.1 MB)
opened URL
downloaded 2.1 Mb

```

Workspace **History**

```

??MASS
install.packages()
packageDescription("MASS")
?unloadNamespace
?library
?data
install.packages("PerformanceAnalytics", repos="http://R-Forge.R-project.org")
R.home
R.home
R.home("home")
R.home()
?Startup

```

Files **Plots** **Packages** **Help**

R: Loading and Listing of Packages

library (base)

Loading and Listing of Packages

Description

library and require load add-on packages.

Usage

```

library(package, help, pos = 2, lib.loc = NULL,
character.only = FALSE, logical.return = FALSE,
warn.conflicts = TRUE, quietly = FALSE,
verbose = getOption("verbose"))

```

Arguments

package, help the name of a package, given as a [name](#) or literal character string, or a character vector of package names, or a [package file](#) (e.g., a .zip or .tar.gz file).

A First R Session

Variables are created by an assignment operation, and they don't have to be declared.

The standard assignment operator in R is the arrow symbol "`<=`".

R interprets text in quotes ("`\"`") as character strings.

Text that is not in quotes ("`\"`") is interpreted as a *symbol* or *expression*.

Typing a *symbol* or *expression* evaluates it.

R uses the hash "`#`" sign to mark text as comments.

All text after the hash "`#`" sign is treated as a comment, and is not executed as code.

```
> # "<=" and "=" are valid assignment operators
> myvar <- 3
>
> # Typing a symbol or expression evaluates it
> myvar
[1] 3
>
> # Text in quotes is interpreted as a string
> myvar <- "Hello World!"
>
> # Typing a symbol or expression evaluates it
> myvar
[1] "Hello World!"
>
> myvar # Text after hash is treated as comment
[1] "Hello World!"
```

Exploring an R Session

The function `getwd()` returns a vector of length 1, with the first element containing a string with the name of the current working directory (`cwd`).

The function `setwd()` accepts a character string as input (the name of the directory), and sets the working directory to that string.

R is a functional language, and R commands are functions, so they must be followed by parentheses `()`.

```
> getwd() # Get cwd
> setwd("/Users/jerzy/Develop/R") # Set cwd
> getwd() # Get cwd
```

Get system date and time

Just the date

```
> Sys.time() # Get date and time
[1] "2023-03-21 14:23:51 EDT"
>
> Sys.Date() # Get date only
[1] "2023-03-21"
```

The R Workspace

The workspace is the current R working environment, which includes all user-defined objects and the command history.

The function `ls()` returns names of objects in the R workspace.

The function `rm()` removes objects from the R workspace.

The workspace can be saved into and loaded back from an `.RData` file (compressed binary file format).

The function `save.image()` saves the whole workspace.

The function `save()` saves just the selected objects.

The function `load()` reads data from `.RData` files, and *invisibly* returns a vector of names of objects created in the workspace.

```
> var1 <- 3 # Define new object
> ls() # List all objects in workspace
> # List objects starting with "v"
> ls(pattern=glob2rx("v*"))
> # Remove all objects starting with "v"
> rm(list=ls(pattern=glob2rx("v*")))
> save.image() # Save workspace to file .RData in cwd
> rm(var1) # Remove object
> ls() # List objects
> load(".RData")
> ls() # List objects
> var2 <- 5 # Define another object
> save(var1, var2, # Save selected objects
+       file="/Users/jerzy/Develop/lecture_slides/data/my_data.RData")
> rm(list=ls()) # Remove all objects
> ls() # List objects
> loadobj <- load(file="/Users/jerzy/Develop/lecture_slides/data/my_data.RData")
> loadobj
> ls() # List objects
```

The R Workspace (cont.)

When you quit R you'll be prompted "Save workspace image?"

If you answer *YES* then the workspace will be saved into the `.RData` file in the `cwd`.

When you start R again, the workspace will be automatically loaded from the existing `.RData` file.

```
> q() # quit R session
```

The function `history()` displays recent commands.

You can also save and load the command history from a file.

```
> history(5) # Display last 5 commands
> savehistory(file="myfile") # Default is ".Rhistory"
> loadhistory(file="myfile") # Default is ".Rhistory"
```

R Session Info

The function `sessionInfo()` returns information about the current R session.

- R version,
- OS platform,
- locale settings,
- list of packages that are loaded and attached to the search path,
- list of packages that are loaded, but *not* attached to the search path,

```
> sessionInfo() # Get R version and other session info
R version 4.2.1 (2022-06-23)
Platform: aarch64-apple-darwin20 (64-bit)
Running under: macOS Monterey 12.5.1

Matrix products: default
BLAS: /Library/Frameworks/R.framework/Versions/4.2-arm64/Resources
LAPACK: /Library/Frameworks/R.framework/Versions/4.2-arm64/Resources

locale:
[1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8

attached base packages:
[1] graphics grDevices utils datasets stats methods base

other attached packages:
[1] knitr_1.40 HighFreq_0.1 rutils_0.2 dygraphs_1.1
[5] quantmod_0.4.20 TTR_0.24.3 xts_0.12.1 zoo_1.8-10

loaded via a namespace (and not attached):
[1] Rcpp_1.0.9 rstudioapi_0.13 magrittr_2.0.3 lattice_
[5] rlang_1.0.6 fastmap_1.1.0 highr_0.9 stringr_
[9] tools_4.2.1 grid_4.2.1 xfun_0.32 cli_3.6.0
[13] htmltools_0.5.3 digest_0.6.29 lifecycle_1.0.3 htmlwidg
[17] vctr_0.5.2 curl_4.3.2 evaluate_0.16 glue_1.6
[21] stringi_1.7.8 compiler_4.2.1
```

Global *Options* Settings

R uses a list of global *options* which affect how R computes and displays results.

The function `options()` either sets or displays the values of global *options*.

`options("globop")` displays the current value of option "globop".

`getOption("globop")` displays the current value of option "globop".

`options(globop=value)` sets the option "globop" equal to "value".

```
> # ?options # Long list of global options
> # Interpret strings as characters, not factors
> getOption("stringsAsFactors") # Display option
> options("stringsAsFactors") # Display option
> options(stringsAsFactors=FALSE) # Set option
> # Number of digits printed for numeric values
> options(digits=3)
> # Control exponential scientific notation of print method
> # Positive "scipen" values bias towards fixed notation
> # Negative "scipen" values bias towards scientific notation
> options(scipen=100)
> # Maximum number of items printed to console
> options(max.print=30)
> # Warning levels options
> # Negative - warnings are ignored
> options(warn=-1)
> # zero - warnings are stored and printed after top-confl function
> options(warn=0)
> # One - warnings are printed as they occur
> options(warn=1)
> # 2 or larger - warnings are turned into errors
> options(warn=2)
> # Save all options in variable
> optionv <- options()
> # Restore all options from variable
> options(optionv)
```

Environments in R

Environments consist of a *frame* (a set of symbol-value pairs) and an *enclosure* (a pointer to an enclosing environment).

There are three system environments:

- `globalenv()` the user's workspace,
- `baseenv()` the environment of the base package,
- `emptyenv()` the only environment without an enclosure,

Environments form a tree structure of successive enclosures, with the empty environment at its root.

Packages have their own environments.

The enclosure of the base package is the empty environment.

```
> rm(list=ls())
> # Get base environment
> baseenv()
> # Get global environment
> globalenv()
> # Get current environment
> environment()
> # Get environment class
> class(environment())
> # Define variable in current environment
> globv <- 1
> # Get objects in current environment
> ls(environment())
> # Create new environment
> new_env <- new.env()
> # Get calling environment of new environment
> parent.env(new_env)
> # Assign Value to Name
> assign("new_var1", 3, envir=new_env)
> # Create object in new environment
> new_env$new_var2 <- 11
> # Get objects in new environment
> ls(new_env)
> # Get objects in current environment
> ls(environment())
> # Environments are subset like listv
> new_env$new_var1
> # Environments are subset like listv
> new_env[["new_var1"]]
```

The R Search Path

R evaluates variables using the search path, a series of environments:

- global environment,
- package environments,
- base environment,

The function `search()` returns the search path for R objects.

The function `attach()` attaches objects to the search path.

Using `attach()` allows referencing object components by their names alone, rather than as components of objects.

The function `detach()` detaches objects from the search path.

The function `find()` finds where objects are located on the search path.

Rule of Thumb

Be very careful with using `attach()`.

Make sure to `detach()` objects once they're not needed.

```
> search() # Get search path for R objects
[1] ".GlobalEnv"      "package:knitr"      "package:graphics"
[4] "package:grDevices" "package:utils"      "package:datasets"
[7] "package:HighFreq"  "package:rutils"     "package:dygraphs"
[10] "package:quantmod"  "package:TTR"        "package:xts"
[13] "package:zoo"       "package:stats"      "package:methods"
[16] "Autoloads"        "package:base"

> my_list <- list(flowers=c("rose", "daisy", "tulip"),
+               trees=c("pine", "oak", "maple"))
> my_list$trees
[1] "pine" "oak"  "maple"
> attach(my_list)
> trees
[1] "pine" "oak"  "maple"
> search() # Get search path for R objects
[1] ".GlobalEnv"      "my_list"            "package:knitr"
[4] "package:graphics" "package:grDevices"  "package:utils"
[7] "package:datasets" "package:HighFreq"   "package:rutils"
[10] "package:dygraphs" "package:quantmod"   "package:TTR"
[13] "package:xts"      "package:zoo"        "package:stats"
[16] "package:methods"  "Autoloads"          "package:base"
> detach(my_list)
> head(trees) # "trees" is in datasets base package
  Girth Height Volume
1   8.3    70   10.3
2   8.6    65   10.3
3   8.8    63   10.2
4  10.5    72   16.4
5  10.7    81   18.8
6  10.8    83   19.7
```


Extracting Time Series from Environments

The function `mget()` accepts a vector of strings and returns a list of the corresponding objects extracted from an *environment*.

The extractor (accessor) functions from package *quantmod*: `C1()`, `Vo()`, etc., extract columns from *OHLC* data.

A list of *xts* series can be flattened into a single *xts* series using the function `do.call()`.

The function `do.call()` executes a function call using a function name and a list of arguments.

`do.call()` passes the list elements individually, instead of passing the whole list as one argument.

The function `eapply()` is similar to `lapply()`, and applies a function to objects in an *environment*, and returns a list.

Time series can also be extracted from an *environment* by coercing it into a list, and then subsetting and merging it into an *xts* series using the function `do.call()`.

```
> library(rutils) # Load package rutils
> # Define ETF symbols
> symbolv <- c("VTI", "VEU", "IEF", "VNQ")
> # Extract symbolv from rutils::etfenv
> pricev <- mget(symbolv, envir=rutils::etfenv)
> # pricev is a list of xts series
> class(pricev)
> class(pricev[[1]])
> # Extract Close prices
> pricev <- lapply(pricev, quantmod::C1)
> # Collapse list into time series the hard way
> xts1 <- cbind(pricev[[1]], pricev[[2]], pricev[[3]], pricev[[4]])
> class(xts1)
> dim(xts1)
> # Collapse list into time series using do.call()
> pricev <- do.call(cbind, pricev)
> all.equal(xts1, pricev)
> class(pricev)
> dim(pricev)
> # Extract and cbind in single step
> pricev <- do.call(cbind, lapply(
+   mget(symbolv, envir=rutils::etfenv), quantmod::C1))
> # Or
> # Extract and bind all data, subset by symbolv
> pricev <- lapply(symbolv, function(symbol) {
+   quantmod::C1(get(symbol, envir=rutils::etfenv))
+ }) # end lapply
> # Same, but loop over etfenv without anonymous function
> pricev <- do.call(cbind,
+   lapply(as.list(rutils::etfenv)[symbolv], quantmod::C1))
> # Same, but works only for OHLC series - produces error
> pricev <- do.call(cbind,
+   eapply(rutils::etfenv, quantmod::C1)[symbolv])
```

Managing Time Series

Time series columns can be renamed, and then saved into .csv files.

The function `strsplit()` splits the elements of a character vector.

The package *zoo* contains functions `write.zoo()` and `read.zoo()` for writing and reading *zoo* time series from .txt and .csv files.

The function `eapply()` is similar to `lapply()`, and applies a function to objects in an *environment*, and returns a list.

The function `assign()` assigns a value to an object in a specified *environment*, by referencing it using a character string (name).

The function `save()` writes objects to compressed binary .RData files.

```
> # Drop ".Close" from column names
> colnames(pricew[, 1:4])
> do.call(rbind, strsplit(colnames(pricew[, 1:4]), split=".[.])", 1
> colnames(pricew) <- do.call(rbind, strsplit(colnames(pricew), spli
> # Or
> colnames(pricew) <- unname(sapply(colnames(pricew),
+   function(colname) strsplit(colname, split=".[.])[[1]][1]))
> tail(pricew, 3)
> # Which objects in global environment are class xts?
> unlist(eapply(globalenv(), is.xts))
> # Save xts to csv file
> write.zoo(pricew,
+   file="/Users/jerzy/Develop/lecture_slides/data/etf_series.csv",
> # Copy prices into etfenv
> etfenv$etf_list <- etf_list
> # Or
> assign("prices", pricew, envir=etfenv)
> # Save to .RData file
> save(etfenv, file="etf_data.RData")
```

Referencing Object Components Using with()

The function `with()` evaluates an expression in an environment constructed from the data.

`with()` allows referencing object components by their names alone.

It's often better to use `with()` instead of `attach()`.

```
> # "trees" is in datasets base package
> head(trees, 3)
  Girth Height Volume
1   8.3    70   10.3
2   8.6    65   10.3
3   8.8    63   10.2
> colnames(trees)
[1] "Girth" "Height" "Volume"
> mean(Girth)

Error in mean(Girth): object 'Girth' not found

> mean(trees$Girth)
[1] 13.2
> with(trees,
+       c(mean(Girth), mean(Height), mean(Volume)))
[1] 13.2 76.0 30.2
```

R Packages

Types of R Packages

R can run libraries of functions called packages,

R packages can also contain data,

Most packages need to be *loaded* into R before they can be used,

R includes a number of base packages that are already installed and loaded,

There's also a special package called the base package, which is responsible for all the basic R functionality, datasets is a base package containing various datasets, for example EuStockMarkets,

The *base* Packages

R includes a number of packages that are pre-installed (often called *base* packages),

Some *base* packages:

- *base* - basic R functionality,
- *stats* - statistical functions and random number generation,
- *graphics* - basic graphics,
- *utils* - utility functions,
- *datasets* - popular datasets,
- *parallel* - support for parallel computation,

Very popular packages:

- *MASS* - functions and datasets for "Modern Applied Statistics with S",
- *ggplot2* - grammar of graphics plots,
- *shiny* - interactive web graphics from R,
- *slidify* - HTML5 slide shows from R,
- *devtools* - create R packages,
- *roxygen2* - document R packages,
- *Rcpp* - integrate C++ code with R,
- *RcppArmadillo* - interface to Armadillo linear algebra library,
- *forecast* - linear models and forecasting,
- *tseries* - time series analysis and computational finance,
- *zoo* - time series and ordered objects,
- *xts* - advanced time series objects,
- *quantmod* - quantitative financial modeling framework,
- *caTools* - moving window statistics for graphics and time series objects,

CRAN Package Views

CRAN view for package **AER**:

<http://cran.r-project.org/web/packages/AER/>

Note:

- Authors,
- Version number,
- Reference manual,
- Vignettes,
- Dependencies on other packages.

The package source code can be downloaded by clicking on the **package source** link,



The screenshot shows the CRAN web page for the 'AER' package. The browser address bar displays 'cran.us.r-project.org/web/packages/AER/'. The page title is 'AER: Applied Econometrics with R'. Below the title, it states 'Functions, data sets, examples, demos, and vignettes for the book Christian Kleiber and Achim Zeileis (2008), Applie'. The page lists various details about the package, including its version (1.2-1), dependencies (R (≥ 2.13.0), car (≥ 2.0-1), lme4, sandwich, survival, zoo), imports (stats, Formula (≥ 0.2-0)), suggests (boot, dymlm, effects, foreign, ineq, KernSmooth, lattice, MASS, mlogit, nlme, rnet, np, plm, pscl), published date (2013-11-07), author (Christian Kleiber [aut], Achim Zeileis [aut, cre]), maintainer (Achim Zeileis <Achim.Zeileis@R-project.org>), license (GPL-2), and needs compilation (no). It also provides citation information, news, and in-view links for 'Econometrics', 'Survival', and 'TimeSeries'. The CRAN checks section shows 'AER results'. The downloads section lists the reference manual (AER.pdf), vignettes (Applied Econometrics with R: Package Vignette and Errata, Sweave Example: Linear Regression for Economics Journals Data), package source (AER_1.2-1.tar.gz), MacOS X binary (AER_1.2-1.tgz), Windows binary (AER_1.2-1.zip), and old sources (AER archive). The reverse dependencies section lists 'lpack' and 'rdd'. The reverse suggests section lists 'censReg', 'glm', 'lme4', 'micEconCES', 'mlogit', 'plm', 'REEMtree', and 'sandwich'.

cran.us.r-project.org/web/packages/AER/

AER: Applied Econometrics with R

Functions, data sets, examples, demos, and vignettes for the book Christian Kleiber and Achim Zeileis (2008), Applie

Version: 1.2-1

Depends: R (≥ 2.13.0), [car](#) (≥ 2.0-1), [lme4](#), [sandwich](#), [survival](#), [zoo](#)

Imports: stats, [Formula](#) (≥ 0.2-0)

Suggests: [boot](#), [dymlm](#), [effects](#), [foreign](#), [ineq](#), [KernSmooth](#), [lattice](#), [MASS](#), [mlogit](#), [nlme](#), [rnet](#), [np](#), [plm](#), [pscl](#)

Published: 2013-11-07

Author: Christian Kleiber [aut], Achim Zeileis [aut, cre]

Maintainer: Achim Zeileis <Achim.Zeileis@R-project.org>

License: [GPL-2](#)

NeedsCompilation: no

Citation: [AER citation info](#)

Materials: [NEWS](#)

In views: [Econometrics](#), [Survival](#), [TimeSeries](#)

CRAN checks: [AER results](#)

Downloads:

Reference manual: [AER.pdf](#)

Vignettes: [Applied Econometrics with R: Package Vignette and Errata](#)
[Sweave Example: Linear Regression for Economics Journals Data](#)

Package source: [AER_1.2-1.tar.gz](#)

MacOS X binary: [AER_1.2-1.tgz](#)

Windows binary: [AER_1.2-1.zip](#)

Old sources: [AER archive](#)

Reverse dependencies:

Reverse depends: [lpack](#), [rdd](#)

Reverse suggests: [censReg](#), [glm](#), [lme4](#), [micEconCES](#), [mlogit](#), [plm](#), [REEMtree](#), [sandwich](#)

CRAN Task Views

CRAN Finance Task View

<http://cran.r-project.org/>

Note:

- Maintainer,
- Topics,
- List of packages.

← → C cran.us.r-project.org



CRAN

[Mirrors](#)

[What's new?](#)

[Task Views](#)

[Search](#)

About R

[R Homepage](#)

[The R Journal](#)

Software

[R Sources](#)

[R Binaries](#)

[Packages](#)

[Other](#)

Documentation

[Manuals](#)

[FAQs](#)

[Contributed](#)

CRAN Task View: Empirical Finance

Maintainer: Dirk Eddelbuettel

Contact: Dirk Eddelbuettel at R-project.org

Version: 2014-01-16

This CRAN Task View contains a list of packages useful for empirical work in Finance,

Besides these packages, a very wide variety of functions suitable for empirical work in Finance are available in R packages on the Comprehensive R Archive Network (CRAN). Consequently, several of the following packages are also available on the [Optimization](#), [Robust](#), [SocialSciences](#) and [TimeSeries](#) Task Views.

Please send suggestions for additions and extensions for this task view to the [task view maintainer](#).

Standard regression models

- A detailed overview of the available regression methodologies is provided by the [lm](#) package.
- Linear models such as ordinary least squares (OLS) can be estimated by `lm()` (from the [stats](#) package) or `lmfit()` (from the [lmfit](#) package). Many other suitable methods are available, such as `nls()` from the [nlme](#) package.
- For the linear model, a variety of regression diagnostic tests are provided by the [car](#) package, which may be of interest as well.

Time series

- A detailed overview of tools for time series analysis can be found in the [TimeSeries](#) Task View.
- Classical time series functionality is provided by the [arima](#)() and [KalmanLike](#)() functions.
- The [dse](#) and [limsac](#) packages provide a variety of more advanced estimation methods.
- For volatility modeling, the standard GARCH(1,1) model can be estimated with the [rugarch](#) package. The [rugarch](#) package can be used to model a variety of univariate GARCH processes. The [rugarch](#) package provides methods for fit, forecast, simulation, inference and plotting as well. The [rugarch](#) package can also estimate and simulate the Beta-t-EGARCH model by Harvey. The [hugarch](#) package provides methods for estimating and simulating multivariate GARCH processes. The [ccgarch](#) package can estimate (multivariate) Conditional Correlation GARCH processes. The [AutoSEARCH](#) package provides automated general-to-specific model selection.
- Unit root and cointegration tests are provided by [tseries](#) and [urca](#). The [Rmetrics](#) package provides unit roots and more. The [CADTest](#) package implements the Hansen unit root tests.
- [MSBVAR](#) provides Bayesian estimation of vector autoregressive models. The [dlm](#) package provides dynamic linear models.
- The [vars](#) package offers estimation, diagnostics, forecasting and error decomposition for vector autoregressive models.
- The [dyn](#) and [dynlm](#) are suitable for dynamic (linear) regression models.
- Several packages provide wavelet analysis functionality: [rwt](#), [wavelets](#), [waveslim](#), [wavelet](#), [waveletComp](#), [waveletR](#), [waveletTools](#), [waveletVis](#), [waveletVis2](#), [waveletVis3](#), [waveletVis4](#), [waveletVis5](#), [waveletVis6](#), [waveletVis7](#), [waveletVis8](#), [waveletVis9](#), [waveletVis10](#), [waveletVis11](#), [waveletVis12](#), [waveletVis13](#), [waveletVis14](#), [waveletVis15](#), [waveletVis16](#), [waveletVis17](#), [waveletVis18](#), [waveletVis19](#), [waveletVis20](#), [waveletVis21](#), [waveletVis22](#), [waveletVis23](#), [waveletVis24](#), [waveletVis25](#), [waveletVis26](#), [waveletVis27](#), [waveletVis28](#), [waveletVis29](#), [waveletVis30](#), [waveletVis31](#), [waveletVis32](#), [waveletVis33](#), [waveletVis34](#), [waveletVis35](#), [waveletVis36](#), [waveletVis37](#), [waveletVis38](#), [waveletVis39](#), [waveletVis40](#), [waveletVis41](#), [waveletVis42](#), [waveletVis43](#), [waveletVis44](#), [waveletVis45](#), [waveletVis46](#), [waveletVis47](#), [waveletVis48](#), [waveletVis49](#), [waveletVis50](#), [waveletVis51](#), [waveletVis52](#), [waveletVis53](#), [waveletVis54](#), [waveletVis55](#), [waveletVis56](#), [waveletVis57](#), [waveletVis58](#), [waveletVis59](#), [waveletVis60](#), [waveletVis61](#), [waveletVis62](#), [waveletVis63](#), [waveletVis64](#), [waveletVis65](#), [waveletVis66](#), [waveletVis67](#), [waveletVis68](#), [waveletVis69](#), [waveletVis70](#), [waveletVis71](#), [waveletVis72](#), [waveletVis73](#), [waveletVis74](#), [waveletVis75](#), [waveletVis76](#), [waveletVis77](#), [waveletVis78](#), [waveletVis79](#), [waveletVis80](#), [waveletVis81](#), [waveletVis82](#), [waveletVis83](#), [waveletVis84](#), [waveletVis85](#), [waveletVis86](#), [waveletVis87](#), [waveletVis88](#), [waveletVis89](#), [waveletVis90](#), [waveletVis91](#), [waveletVis92](#), [waveletVis93](#), [waveletVis94](#), [waveletVis95](#), [waveletVis96](#), [waveletVis97](#), [waveletVis98](#), [waveletVis99](#), [waveletVis100](#).

Installing Packages

Most packages need to be *installed* before they can be loaded and used.

Some packages like *MASS* are installed with base R (but not loaded).

Installing a package means downloading and saving its files to a local computer directory (hard disk), so they can be *loaded* by the R system.

The function `install.packages()` installs packages from the R command line.

Most widely used packages are available on the *CRAN* repository:

<http://cran.r-project.org/web/packages/>

Or on *R-Forge* or *GitHub*:

<https://r-forge.r-project.org/>

<https://github.com/>

Packages can also be installed in *RStudio* from the menu (go to **Tools** and then **Install packages**),

Packages residing on GitHub can be installed using the devtools packages.

```
> getOption("repos") # get default package source
> .libPaths() # get package save directory
> install.packages("AER") # install "AER" from CRAN
> # install "PerformanceAnalytics" from R-Forge
> install.packages(
+   pkgs="PerformanceAnalytics", # name
+   lib="C:/Users/Jerzy/Downloads", # directory
+   repos="http://R-Forge.R-project.org") # source
> # install devtools from CRAN
> install.packages("devtools")
> # load devtools
> library(devtools)
> # install package "babynamesv" from GitHub
> install_github(repo="hadley/babynamesv")
```


Installing Packages From Source

Sometimes packages aren't available in compiled form, so it's necessary to install them from their source code.

To install a package from source, the user needs to first install compilers and development tools:

For Windows install Rtools:

<https://cran.r-project.org/bin/windows/Rtools/>

For Mac OSX install XCode developer tools:

<https://developer.apple.com/xcode/downloads/>

The function `install.packages()` with argument `type="source"` installs a package from source.

The function `download.packages()` downloads the package's installation files (compressed tar format) to a local directory.

The function `install.packages()` can then be used to install the package from the downloaded files.

```
> # install package "PortfolioAnalytics" from source
> install.packages("PortfolioAnalytics",
+   type="source",
+   repos="http://r-forge.r-project.org")
> # download files for package "PortfolioAnalytics"
> download.packages(pkgs = "PortfolioAnalytics",
+   destdir = ".", # download to cwd
+   type = "source",
+   repos="http://r-forge.r-project.org")
> # install "PortfolioAnalytics" from local tar source
> install.packages(
+   "C:/Users/Jerzy/Downloads/PortfolioAnalytics_0.9.3598.tar.gz",
+   repos=NULL, type="source")
```

Installed Packages

`defaultPackages` contains a list of packages loaded on startup by default.

The function `installed.packages()` returns a matrix of all packages installed on the system.

```
> getOption("defaultPackages")
> # matrix of installed package information
> pack_info <- installed.packages()
> dim(pack_info)
> # get all installed package names
> sort(unname(pack_info[, "Package"]))
> # get a few package names and their versions
> pack_info[sample(x=1:100, 5), c("Package", "Version")]
> # get info for package "xts"
> t(pack_info["xts", ])
```

Package Files and Directories

Package installation files are organized into multiple directories, including some of the following:

- `~/R` containing R source code files,
- `~/src` containing C++ and Fortran source code files,
- `~/data` containing datasets,
- `~/man` containing documentation files,

```
> # list directories in "PortfolioAnalytics" sub-directory
> gsub(
+   "C:/Users/Jerzy/Documents/R/win-library/3.1",
+   "~",
+   list.dirs(
+     file.path(
+       .libPaths()[1],
+       "PortfolioAnalytics")))
[1] "/Library/Frameworks/R.framework/Versions/4.2-arm64/Resources/1"
[2] "/Library/Frameworks/R.framework/Versions/4.2-arm64/Resources/2"
[3] "/Library/Frameworks/R.framework/Versions/4.2-arm64/Resources/3"
[4] "/Library/Frameworks/R.framework/Versions/4.2-arm64/Resources/4"
[5] "/Library/Frameworks/R.framework/Versions/4.2-arm64/Resources/5"
[6] "/Library/Frameworks/R.framework/Versions/4.2-arm64/Resources/6"
[7] "/Library/Frameworks/R.framework/Versions/4.2-arm64/Resources/7"
[8] "/Library/Frameworks/R.framework/Versions/4.2-arm64/Resources/8"
[9] "/Library/Frameworks/R.framework/Versions/4.2-arm64/Resources/9"
[10] "/Library/Frameworks/R.framework/Versions/4.2-arm64/Resources/10"
[11] "/Library/Frameworks/R.framework/Versions/4.2-arm64/Resources/11"
[12] "/Library/Frameworks/R.framework/Versions/4.2-arm64/Resources/12"
[13] "/Library/Frameworks/R.framework/Versions/4.2-arm64/Resources/13"
[14] "/Library/Frameworks/R.framework/Versions/4.2-arm64/Resources/14"
```

Loading Packages

Most packages need to be *loaded* before they can be used in an R session.

Loading a package means attaching the package *namespace* to the *search path*, which allows R to call the package functions and data.

The functions `library()` and `require()` load packages, but in slightly different ways.

`library()` produces an *error* (halts execution) if the package can't be loaded.

`require()` returns `TRUE` if the package is loaded successfully, and `FALSE` otherwise.

Therefore `library()` is usually used in script files that might be sourced, while `require()` is used inside functions.

```
> # load package, produce error if can't be loaded
> library(MASS)
> # load package, return TRUE if loaded successfully
> require(MASS)
> # load quietly
> library(MASS, quietly=TRUE)
> # load without any messages
> suppressMessages(library(MASS))
> # remove package from search path
> detach(MASS)
> # install package if it can't be loaded successfully
> if (!require("xts")) install.packages("xts")
```

Referencing Package Objects

After a package is *loaded*, the package functions and data can be accessed by name.

Package objects can also be accessed without *loading* the package, by using the double-colon ":" reference operator.

For example, `TTR::VWAP()` references the function `VWAP()` from the package *TTR*.

This way users don't have to load the package *TTR* (with `library(TTR)`) to use functions from the package *TTR*.

Using the ":" operator displays the source of objects, and makes R code easier to analyze.

```
> # calculate VTI volume-weighted average price
> vwapv <- TTR::VWAP(
+   price=quantmod::Cl(rutils::etfenv$VTI),
+   volume=quantmod::Vo(rutils::etfenv$VTI), n=10)
```

Exploring Packages

The package *Ecdat* contains data sets for econometric analysis.

The data frame *Garch* contains daily currency prices.

The function `data()` loads external data or listv data sets in a package.

Some packages provide *lazy loading* of their data sets, which means they automatically load their data sets when they're needed (when they are called by some operation).

The package's data isn't loaded into R memory when the package is *loaded*, so it's not listed using `ls()`, but the package data is available without calling the function `data()`.

The function `data()` isn't required to load data sets that are set up for *lazy loading*.

```
> library() # list all packages installed on the system
> search() # list all loaded packages on search path
>
> # get documentation for package "Ecdat"
> packageDescription("Ecdat") # get short description
> help(package="Ecdat") # load help page
> library(Ecdat) # load package "Ecdat"
> data(package="Ecdat") # list all datasets in "Ecdat"
> ls("package:Ecdat") # list all objects in "Ecdat"
> browseVignettes("Ecdat") # view package vignette
> detach("package:Ecdat") # remove Ecdat from search path
```

```
> library(Ecdat) # load econometric data sets
> class(Garch) # Garch is a data frame from "Ecdat"
> dim(Garch) # daily currency prices
> head(Garch[, -2]) # col 'dm' is Deutsch Mark
> detach("package:Ecdat") # remove Ecdat from search path
```

Package Namespaces

Package *namespaces*:

- Provide a mechanism for calling objects from a package,
- Hide functions and data internal to the package,
- Prevent naming conflicts between user and package names,

When a package is loaded using `library()` or `require()`, its *namespace* is attached to the search path.

```
> search() # get search path for R objects
> library(MASS) # load package "MASS"
> head(ls("package:MASS")) # list some objects in "MASS"
> detach("package:MASS") # remove "MASS" from search path
```

Package Namespaces and the Search Path

Packages may be loaded without their *namespace* being attached to the search path.

When packages are loaded, then packages they depend on are also loaded, but their *namespaces* aren't necessarily attached to the search path.

The function `loadedNamespaces()` lists all loaded *namespaces*, including those that aren't on the search path.

The function `search()` returns the current search path for R objects.

`search()` returns many package *namespaces*, but not all the loaded *namespaces*.

```
> loadedNamespaces() # get names of loaded namespaces
>
> search() # get search path for R objects
```


Not Attached Namespaces

The function `sessionInfo()` returns information about the current R session, including packages that are loaded, but *not attached* to the search path.

`sessionInfo()` lists those packages as "loaded via a namespace (and not attached)"

```
> # get session info,  
> # including packages not attached to the search path  
> sessionInfo()
```

Non-Visible Objects

Non-visible objects (variables or functions) are either:

- objects from *not attached namespaces*,
- objects *not exported* outside a package,

Objects from packages that aren't attached can be accessed using the double-colon ":" reference operator.

Objects that are *not exported* outside a package can be accessed using the triple-colon ":::" reference operator.

Colon operators automatically load the associated package.

Non-visible objects in namespaces often use the ".*" name syntax.

```
> plot.xts # package xts isn't loaded and attached
> head(xts::plot.xts, 3)
> methods("cbind") # get all methods for function "cbind"
> stats::cbind.ts # cbind isn't exported from package stats
> stats:::cbind.ts # view the non-visible function
> getAnywhere("cbind.ts")
> library(MASS) # load package 'MASS'
> select # code of primitive function from package 'MASS'
```

Exploring Namespaces and Non-Visible Objects

The function `getAnywhere()` displays information about R objects, including non-visible objects.

Objects referenced *within* packages have different search paths than other objects:

Their search path starts in the package *namespace*, then the global environment and then finally the regular search path.

This way references to objects from *within* a package are resolved to the package, and they're not masked by objects of the same name in other environments.

```
> getAnywhere("cbind.ts")
```

Benchmarking the Speed of R Code

The function `system.time()` calculates the execution time (in seconds) used to evaluate a given expression.

`system.time()` returns the "*user time*" (execution time of user instructions), the "*system time*" (execution time of operating system calls), and "*elapsed time*" (total execution time, including system latency waiting).

The function `microbenchmark()` from package `microbenchmark` calculates and compares the execution time of R expressions (in milliseconds), and is more accurate than `system.time()`.

The time it takes to execute an expression is not always the same, since it depends on the state of the processor, caching, etc.

`microbenchmark()` executes the expression many times, and returns the distribution of total execution times.

```
> library(microbenchmark)
> vectorv <- runif(1e6)
> # sqrt() and "^0.5" are the same
> all.equal(sqrt(vectorv), vectorv^0.5)
> # sqrt() is much faster than "^0.5"
> system.time(vectorv^0.5)
> microbenchmark(
+   power = vectorv^0.5,
+   sqrt = sqrt(vectorv),
+   times=10)
```

The "*times*" parameter is the number of times the expression is evaluated.

The choice of the "*times*" parameter is a tradeoff between the time it takes to run `microbenchmark()`, and the desired accuracy,

Using apply() Instead of for() and while() Loops

All the different R loops have similar speed, with `apply()` the fastest, then `vapply()`, `lapply()` and `sapply()` slightly slower, and `for()` loops the slowest.

More importantly, the `apply()` syntax is more readable and concise, and fits the functional language paradigm of R, so it's preferred over `for()` loops.

Both `vapply()` and `lapply()` are *compiled (primitive)* functions, and therefore can be faster than other `apply()` functions.

```
> # Calculate matrix of random data with 5,000 rows
> matrixv <- matrix(rnorm(10000), ncol=2)
> # Allocate memory for row sums
> rowsumv <- numeric(NROW(matrixv))
> summary(microbenchmark(
+   rowsumv = rowSums(matrixv), # end rowsumv
+   applyloop = apply(matrixv, 1, sum), # end apply
+   applyloop = lapply(1:NROW(matrixv), function(indeks)
+     sum(matrixv[indeks, ])), # end lapply
+   v_apply = vapply(1:NROW(matrixv), function(indeks)
+     sum(matrixv[indeks, ]),
+     FUN.VALUE = c(sum=0)), # end vapply
+   s_apply = sapply(1:NROW(matrixv), function(indeks)
+     sum(matrixv[indeks, ])), # end sapply
+   forloop = for (i in 1:NROW(matrixv)) {
+     rowsumv[i] <- sum(matrixv[i,])
+   }, # end for
+   times=10))[, c(1, 4, 5)] # end microbenchmark summary
```

Increasing Speed of Loops by Pre-allocating Memory

R performs automatic memory management as users assign values to objects.

R doesn't require allocating the full memory for vectors or lists, and allows appending new data to existing objects ("growing" them).

For example, R allows assigning a value to a vector element that doesn't exist yet.

This forces R to allocate additional memory for that element, which carries a small speed penalty.

But when data is appended to an object using the functions `c()`, `append()`, `cbind()`, or `rbind()`, then R allocates memory for the whole new object and copies all the existing values, which is very memory intensive and slow.

It is therefore preferable to pre-allocate memory for large objects before performing loops.

The function `numeric(k)` returns a numeric vector of zeros of length `k`, while `numeric(0)` returns an empty (zero length) numeric vector (not to be confused with a `NULL` object).

```
> vectorv <- rnorm(5000)
> summary(microbenchmark(
+ # Allocate full memory for cumulative sum
+   forloop = {cumsumv <- numeric(NROW(vectorv))
+     cumsumv[1] <- vectorv[1]
+     for (i in 2:NROW(vectorv)) {
+       cumsumv[i] <- cumsumv[i-1] + vectorv[i]
+     }}, # end for
+ # Allocate zero memory for cumulative sum
+   grow_vec = {cumsumv <- numeric(0)
+     cumsumv[1] <- vectorv[1]
+     for (i in 2:NROW(vectorv)) {
+       # Add new element to "cumsumv" ("grow" it)
+       cumsumv[i] <- cumsumv[i-1] + vectorv[i]
+     }}, # end for
+ # Allocate zero memory for cumulative sum
+   com_bine = {cumsumv <- numeric(0)
+     cumsumv[1] <- vectorv[1]
+     for (i in 2:NROW(vectorv)) {
+       # Add new element to "cumsumv" ("grow" it)
+       cumsumv <- c(cumsumv, vectorv[i])
+     }}, # end for
+   times=10))[, c(1, 4, 5)]
```

Vectorized Functions for Vector Computations

Vectorized functions accept vectors as their arguments, and return a vector of the same length as their value.

Many *vectorized* functions are also *compiled* (they pass their data to compiled C++ code), which makes them very fast.

The following *vectorized compiled* functions calculate cumulative values over large vectors:

- `cummax()`
- `cummin()`
- `cumsum()`
- `cumprod()`

Standard arithmetic operations ("`+`", "`-`", etc.) can be applied to vectors, and are implemented as *vectorized compiled* functions.

`ifelse()` and `which()` are *vectorized compiled* functions for logical operations.

But many *vectorized* functions perform their calculations in R code, and are therefore slow, but convenient to use.

```
> vector1 <- rnorm(1000000)
> vector2 <- rnorm(1000000)
> big_vector <- numeric(1000000)
> # Sum two vectors in two different ways
> summary(microbenchmark(
+   # Sum vectors using "for" loop
+   rloop = (for (i in 1:NROW(vector1)) {
+     big_vector[i] <- vector1[i] + vector2[i]
+   }),
+   # Sum vectors using vectorized "+"
+   vectorvized = (vector1 + vector2),
+   times=10))[, c(1, 4, 5)] # end microbenchmark summary
> # Allocate memory for cumulative sum
> cumsumv <- numeric(NROW(big_vector))
> cumsumv[1] <- big_vector[1]
> # Calculate cumulative sum in two different ways
> summary(microbenchmark(
+   # Cumulative sum using "for" loop
+   rloop = (for (i in 2:NROW(big_vector)) {
+     cumsumv[i] <- cumsumv[i-1] + big_vector[i]
+   }),
+   # Cumulative sum using "cumsum"
+   vectorvized = cumsum(big_vector),
+   times=10))[, c(1, 4, 5)] # end microbenchmark summary
```

Vectorized Functions for Matrix Computations

`apply()` loops are very inefficient for calculating statistics over rows and columns of very large matrices.

R has very fast *vectorized compiled* functions for calculating sums and means of rows and columns:

- `rowSums()`
- `colSums()`
- `rowMeans()`
- `colMeans()`

These *vectorized* functions are also *compiled* functions, so they're very fast because they pass their data to compiled C++ code, which performs the loop calculations.

```
> # Calculate matrix of random data with 5,000 rows
> matrixv <- matrix(rnorm(10000), ncol=2)
> # Calculate row sums two different ways
> all.equal(rowSums(matrixv),
+   apply(matrixv, 1, sum))
> summary(microbenchmark(
+   rowsumv = rowSums(matrixv),
+   applyloop = apply(matrixv, 1, sum),
+   times=10))[, c(1, 4, 5)] # end microbenchmark summary
```


Fast R Code for Matrix Computations

The functions `pmax()` and `pmin()` calculate the "parallel" maxima (minima) of multiple vector arguments.

`pmax()` and `pmin()` return a vector, whose n -th element is equal to the maximum (minimum) of the n -th elements of the arguments, with shorter vectors recycled if necessary.

`pmax.int()` and `pmin.int()` are methods of generic functions `pmax()` and `pmin()`, designed for atomic vectors.

`pmax()` can be used to quickly calculate the maximum values of rows of a matrix, by first converting the matrix columns into a list, and then passing them to `pmax()`.

`pmax.int()` and `pmin.int()` are very fast because they are *compiled* functions (compiled from C++ code).

```
> library(microbenchmark)
> str(pmax)
> # Calculate row maximums two different ways
> summary(microbenchmark(
+   pmax=do.call(pmax.int,
+   lapply(seq_along(matrixv[1, ]),
+     function(indeks) matrixv[, indeks])),
+   applyloop=unlist(lapply(seq_along(matrixv[, 1]),
+     function(indeks) max(matrixv[indeks, ]))),
+   times=10))[, c(1, 4, 5)]
```

Package matrixStats for Fast Matrix Computations

The package *matrixStats* contains functions for calculating aggregations over matrix columns and rows, and other matrix computations, such as:

- estimating location and scale: `rowRanges()`, `colRanges()`, and `rowMaxs()`, `rowMins()`, etc.,
- testing and counting values: `colAnyMissings()`, `colAnys()`, etc.,
- cumulative functions: `colCumsums()`, `colCummins()`, etc.,
- binning and differencing: `binCounts()`, `colDiffs()`, etc.,

A summary of *matrixStats* functions can be found under:

<https://cran.r-project.org/web/packages/matrixStats/vignettes/matrixStats-methods.html>

The *matrixStats* functions are very fast because they are *compiled* functions (compiled from C++ code).

```
> install.packages("matrixStats") # Install package matrixStats
> library(matrixStats) # Load package matrixStats
> # Calculate row min values three different ways
> summary(microbenchmark(
+   rowmins = rowMins(matrixv),
+   pmin =
+     do.call(pmin.int,
+       lapply(seq_along(matrixv[1, ]),
+         function(indeks)
+           matrixv[, indeks])),
+   as_dframe =
+     do.call(pmin.int,
+       as.data.frame.matrix(matrixv)),
+   times=10))[, c(1, 4, 5)] # end microbenchmark summary
```

Package Rfast for Fast Matrix and Numerical Computations

The package *Rfast* contains functions for fast matrix and numerical computations, such as:

- `colMedians()` and `rowMedians()` for matrix column and row medians,
- `colCumSums()`, `colCumMins()` for cumulative sums and min/max,
- `eigen.sym()` for performing eigenvalue matrix decomposition,

The Rfast functions are very fast because they are *compiled* functions (compiled from C++ code).

```
> install.packages("Rfast") # Install package Rfast
> library(Rfast) # Load package Rfast
> # Benchmark speed of calculating ranks
> vectorv <- 1e3
> all.equal(rank(vectorv), Rfast::Rank(vectorv))
> library(microbenchmark)
> summary(microbenchmark(
+   rcode = rank(vectorv),
+   Rfast = Rfast::Rank(vectorv),
+   times=10))[, c(1, 4, 5)] # end microbenchmark summary
> # Benchmark speed of calculating column medians
> matrixv <- matrix(1e4, nc=10)
> all.equal(matrixStats::colMedians(matrixv), Rfast::colMedians(matrixv))
> summary(microbenchmark(
+   matrixStats = matrixStats::colMedians(matrixv),
+   Rfast = Rfast::colMedians(matrixv),
+   times=10))[, c(1, 4, 5)] # end microbenchmark summary
```

Writing Fast R Code Using Vectorized Operations

R-style code is code that relies on *vectorized compiled* functions, instead of `for()` loops.

`for()` loops in R are slow because they call functions multiple times, and individual function calls are compute-intensive and slow.

The brackets "`[]`" operator is a *vectorized compiled* function, and is therefore very fast.

Vectorized assignments using brackets "`[]`" and Boolean or integer vectors to subset vectors or matrices are therefore preferable to `for()` loops.

R code that uses *vectorized compiled* functions can be as fast as C++ code.

R-style code is also very *expressive*, i.e. it allows performing complex operations with very few lines of code.

```
> summary(microbenchmark( # Assign values to vector three different
+ # Fast vectorized assignment loop performed in C using brackets "
+   brackets = {vectorv <- numeric(10)
+     vectorv[] <- 2},
+ # Slow because loop is performed in R
+   forloop = {vectorv <- numeric(10)
+     for (indeks in seq_along(vectorv))
+       vectorv[indeks] <- 2},
+   times=10))[, c(1, 4, 5)] # end microbenchmark summary
> summary(microbenchmark( # Assign values to vector two different v
+ # Fast vectorized assignment loop performed in C using brackets "
+   brackets = {vectorv <- numeric(10)
+     vectorv[4:7] <- rnorm(4)},
+ # Slow because loop is performed in R
+   forloop = {vectorv <- numeric(10)
+     for (indeks in 4:7)
+       vectorv[indeks] <- rnorm(1)},
+   times=10))[, c(1, 4, 5)] # end microbenchmark summary
```

Vectorized Functions

Functions which use vectorized operations and functions are automatically *vectorized* themselves.

Functions which only call other compiled C++ vectorized functions, are also very fast.

But not all functions are vectorized, or they're not vectorized with respect to their *parameters*.

Some *vectorized* functions perform their calculations in R code, and are therefore slow, but convenient to use.

```
> # Define function vectorized automatically
> my_fun <- function(input, param) {
+   param*input
+ } # end my_fun
> # "input" is vectorized
> my_fun(input=1:3, param=2)
> # "param" is vectorized
> my_fun(input=10, param=2:4)
> # Define vectors of parameters of rnorm()
> stdevs <- structure(1:3, names=paste0("sd=", 1:3))
> means <- structure(-1:1, names=paste0("mean=", -1:1))
> # "sd" argument of rnorm() isn't vectorized
> rnorm(1, sd=stdevs)
> # "mean" argument of rnorm() isn't vectorized
> rnorm(1, mean=means)
```

Performing sapply() Loops Over Function Parameters

Many functions aren't vectorized with respect to their *parameters*.

Performing sapply() loops over a function's parameters produces vector output.

```
> # Loop over stdevs produces vector output
> set.seed(1121)
> sapply(stdevs, function(stdev) rnorm(n=2, sd=stdev))
> # Same
> set.seed(1121)
> sapply(stdevs, rnorm, n=2, mean=0)
> # Loop over means
> set.seed(1121)
> sapply(means, function(meanv) rnorm(n=2, mean=meanv))
> # Same
> set.seed(1121)
> sapply(means, rnorm, n=2)
```

Creating Vectorized Functions

In order to *vectorize* a function with respect to one of its *parameters*, it's necessary to perform a loop over it.

The function `Vectorize()` performs an `apply()` loop over the arguments of a function, and returns a vectorized version of the function.

`Vectorize()` vectorizes the arguments passed to "vectorize.args".

`Vectorize()` is an example of a *higher order* function: it accepts a function as its argument and returns a function as its value.

Functions that are vectorized using `Vectorize()` or `apply()` loops are just as slow as `apply()` loops, but convenient to use.

```
> # rnorm() vectorized with respect to "stdev"
> vec_rnorm <- function(n, mean=0, sd=1) {
+   if (NROW(sd)==1)
+     rnorm(n=n, mean=mean, sd=sd)
+   else
+     sapply(sd, rnorm, n=n, mean=mean)
+ } # end vec_rnorm
> set.seed(1121)
> vec_rnorm(n=2, sd=stdevs)
> # rnorm() vectorized with respect to "mean" and "sd"
> vec_rnorm <- Vectorize(FUN=rnorm,
+   vectorize.args=c("mean", "sd")
+ ) # end Vectorize
> set.seed(1121)
> vec_rnorm(n=2, sd=stdevs)
> set.seed(1121)
> vec_rnorm(n=2, mean=means)
```

The mapply() Functional

The `mapply()` functional is a multivariate version of `sapply()`, that allows calling a non-vectorized function in a vectorized way.

`mapply()` accepts a multivariate function passed to the "FUN" argument and any number of vector arguments passed to the dots "...".

`mapply()` calls "FUN" on the vectors passed to the dots "...", one element at a time:

$$\begin{aligned} \text{mapply}(\text{FUN} = \text{fun}, \text{vec1}, \text{vec2}, \dots) = \\ [\text{fun}(\text{vec1}_{1,1}, \text{vec2}_{1,1}, \dots), \dots, \\ \text{fun}(\text{vec1}_{i,i}, \text{vec2}_{i,i}, \dots), \dots] \end{aligned}$$

`mapply()` passes the first vector to the first argument of "FUN", the second vector to the second argument, etc.

The first element of the output vector is equal to "FUN" called on the first elements of the input vectors, the second element is "FUN" called on the second elements, etc.

```
> str(sum)
> # na.rm is bound by name
> mapply(sum, 6:9, c(5, NA, 3), 2:6, na.rm=TRUE)
> str(rnorm)
> # mapply vectorizes both arguments "mean" and "sd"
> mapply(rnorm, n=5, mean=means, sd=stdevs)
> mapply(function(input, e_xp) input^e_xp,
+ 1:5, seq(from=1, by=0.2, length.out=5))
```

The output of `mapply()` is a vector of length equal to the longest vector passed to the dots "...", with the elements of the other vectors recycled if necessary,

Vectorizing Functions Using mapply()

The mapply() functional is a multivariate version of sapply(), that allows calling a non-vectorized function in a vectorized way.

mapply() can be used to vectorize several function arguments simultaneously.

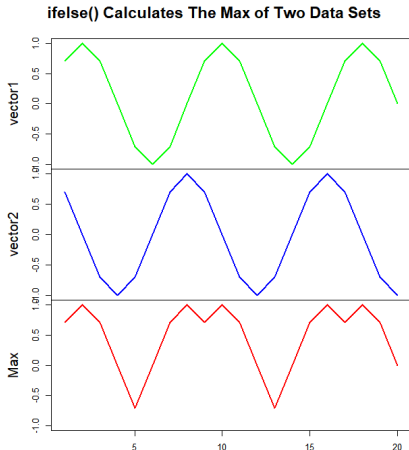
```
> # rnorm() vectorized with respect to "mean" and "sd"
> vec_rnorm <- function(n, mean=0, sd=1) {
+   if (NROW(mean)==1 && NROW(sd)==1)
+     rnorm(n=n, mean=mean, sd=sd)
+   else
+     mapply(rnorm, n=n, mean=mean, sd=sd)
+ } # end vec_rnorm
> # Call vec_rnorm() on vector of "sd"
> vec_rnorm(n=2, sd=stdevs)
> # Call vec_rnorm() on vector of "mean"
> vec_rnorm(n=2, mean=means)
```

Vectorized if-else Statements Using Function ifelse()

The function `ifelse()` performs *vectorized* if-else statements on vectors.

`ifelse()` is much faster than performing an element-wise loop in R.

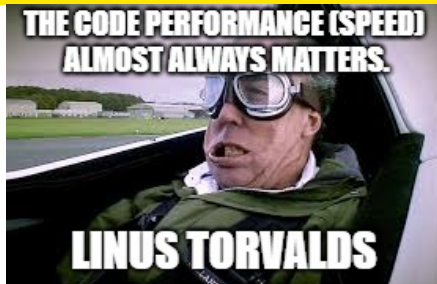
```
> # Create two numeric vectors
> vector1 <- sin(0.25*pi*1:20)
> vector2 <- cos(0.25*pi*1:20)
> # Create third vector using 'ifelse'
> vector3 <- ifelse(vector1 > vector2, vector1, vector2)
> # cbind all three together
> vector3 <- cbind(vector1, vector2, vector3)
> colnames(vector3)[3] <- "Max"
> # Set plotting parameters
> x11(width=6, height=7)
> par(oma=c(0, 1, 1, 1), mar=c(0, 2, 2, 1),
+     mgp=c(2, 1, 0), cex.lab=0.5, cex.axis=1.0, cex.main=1.8, cex.
> # Plot matrix
> zoo::plot.zoo(vector3, lwd=2, ylim=c(-1, 1),
+   xlab="", col=c("green", "blue", "red"),
+   main="ifelse() Calculates The Max of Two Data Sets")
```



It's *Always* Important to Write Fast R Code

How to write fast R code:

- Avoid using `apply()` and `for()` loops for large datasets.
- Use R functions which are *compiled* C++ code, instead of using interpreted R code.
- Avoid using too many R function calls (every command in R is a function).
- Pre-allocate memory for new objects, instead of appending to them ("growing" them).
- Write C++ functions in *Rcpp* and *RcppArmadillo*.
- Use *function methods* directly instead of using *generic functions*.
- Create specialized functions by extracting only the essential R code from *function methods*.
- *Byte-compile* R functions using the *byte compiler* in package *compiler*.



```
> # Calculate cumulative sum of a vector
> vectorv <- runif(1e5)
> # Use compiled function
> cumsumv <- cumsum(vectorv)
> # Use for loop
> cumsumv2 <- vectorv
> for (i in 2:NROW(cumsumv2))
+   cumsumv2[i] <- (cumsumv2[i] + cumsumv2[i-1])
> # Compare the two methods
> all.equal(cumsumv, cumsumv2)
> # Microbenchmark the two methods
> library(microbenchmark)
> summary(microbenchmark(
+   cumsum=cumsum(vectorv),
+   loop_alloc={
+     cumsumv2 <- vectorv
+     for (i in 2:NROW(cumsumv2))
+       cumsumv2[i] <- (cumsumv2[i] + cumsumv2[i-1])
+   },
+   loop_nalloc={
+     # Doesn't allocate memory to cumsumv3
```

Parallel Computing in R

Parallel Computing in R

Parallel computing means splitting a computing task into separate sub-tasks, and then simultaneously computing the sub-tasks on several computers or CPU cores.

There are many different packages that allow parallel computing in R, most importantly package *parallel*, and packages *foreach*, *doParallel*, and related packages:

<http://cran.r-project.org/web/views/HighPerformanceComputing.html>

<http://blog.revolutionanalytics.com/high-performance-computing/>

<http://gforge.se/2015/02/how-to-go-parallel-in-r-basics-tips/>

R Base Package *parallel*

The package *parallel* provides functions for parallel computing using multiple cores of CPUs,

The package *parallel* is part of the standard R distribution, so it doesn't need to be installed.

<http://adv-r.had.co.nz/Profiling.html#parallelise>

<https://github.com/tobiothub/R-parallel/wiki/R-parallel-package-overview>

Packages *foreach*, *doParallel*, and Related Packages

<http://blog.revolutionanalytics.com/2015/10/updates-to-the-foreach-package-and-its-friends.html>

Parallel Computing Using Package *parallel*

The package *parallel* provides functions for parallel computing using multiple cores of CPUs.

The package *parallel* is part of the standard R distribution, so it doesn't need to be installed.

Different functions from package *parallel* need to be called depending on the operating system (*Windows*, *Mac-OSX*, or *Linux*).

Parallel computing requires additional resources and time for distributing the computing tasks and collecting the output, which produces a computing overhead.

Therefore parallel computing can actually be slower for small computations, or for computations that can't be naturally separated into sub-tasks.

```
> library(parallel) # Load package parallel
> # Get short description
> packageDescription("parallel")
> # Load help page
> help(package="parallel")
> # List all objects in "parallel"
> ls("package:parallel")
```

Performing Parallel Loops Using Package *parallel*

Some computing tasks naturally lend themselves to parallel computing, like for example performing loops.

Different functions from package *parallel* need to be called depending on the operating system (*Windows*, *Mac-OSX*, or *Linux*).

The function `mclapply()` performs loops (similar to `lapply()`) using parallel computing on several CPU cores under *Mac-OSX* or *Linux*.

Under *Windows*, a cluster of R processes (one per each CPU core) need to be started first, by calling the function `makeCluster()`.

Mac-OSX and *Linux* don't require calling the function `makeCluster()`.

The function `parLapply()` is similar to `lapply()`, and performs loops under *Windows* using parallel computing on several CPU cores.

```
> # Define function that pauses execution
> paws <- function(x, sleep_time=0.01) {
+   Sys.sleep(sleep_time)
+   x
+ } # end paws
> library(parallel) # Load package parallel
> # Calculate number of available cores
> ncores <- detectCores() - 1
> # Initialize compute cluster under Windows
> cluster <- makeCluster(ncores)
> # Perform parallel loop under Windows
> outv <- parLapply(cluster, 1:10, paws)
> # Perform parallel loop under Mac-OSX or Linux
> outv <- mclapply(1:10, paws, mc.cores=ncores)
> library(microbenchmark) # Load package microbenchmark
> # Compare speed of lapply versus parallel computing
> summary(microbenchmark(
+   standard = lapply(1:10, paws),
+   parallel = parLapply(cluster, 1:10, paws),
+   times=10)
+ )[, c(1, 4, 5)]
```

Computing Advantage of Parallel Computing

Parallel computing provides an increasing advantage for larger number of loop iterations.

The function `stopCluster()` stops the R processes running on several CPU cores.

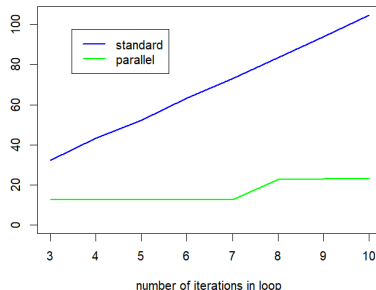
The function `plot()` by default plots a scatterplot, but can also plot lines using the argument `type="l"`.

The function `lines()` adds lines to a plot.

The function `legend()` adds a legend to a plot.

```
> # Compare speed of lapply with parallel computing
> iterations <- 3:10
> compute_times <- sapply(iterations,
+   function(max_iterations) {
+     summary(microbenchmark(
+       standard = lapply(1:max_iterations, paws),
+       parallel = parLapply(cluster, 1:max_iterations, paws),
+       times=10))[, 4]
+   }) # end sapply
> compute_times <- t(compute_times)
> colnames(compute_times) <- c("standard", "parallel")
> rownames(compute_times) <- iterations
> # Stop R processes over cluster under Windows
> stopCluster(cluster)
```

Compute times



```
> x11(width=6, height=5)
> plot(x=rownames(compute_times),
+   y=compute_times[, "standard"],
+   type="l", lwd=2, col="blue",
+   main="Compute times",
+   xlab="number of iterations in loop", ylab="",
+   ylim=c(0, max(compute_times[, "standard"])))
> lines(x=rownames(compute_times),
+   y=compute_times[, "parallel"], lwd=2, col="green")
> legend(x="topleft", legend=colnames(compute_times),
+   inset=0.1, cex=1.0, bg="white",
+   lwd=2, lty=1, col=c("blue", "green"))
```

Parallel Computing Over Matrices

Very often we need to perform time consuming calculations over columns of matrices.

The function `parCapply()` performs an apply loop over columns of matrices using parallel computing on several CPU cores.

```
> # Calculate matrix of random data
> matrixv <- matrix(rnorm(1e5), ncol=100)
> # Define aggregation function over column of matrix
> aggfun <- function(column) {
+   output <- 0
+   for (indeks in 1:NROW(column))
+     output <- output + column[indeks]
+   output
+ } # end aggfun
> # Perform parallel aggregations over columns of matrix
> aggs <- parCapply(cluster, matrixv, aggfun)
> # Compare speed of apply with parallel computing
> summary(microbenchmark(
+   applyloop=apply(matrixv, MARGIN=2, aggfun),
+   parapplyloop=parCapply(cluster, matrixv, aggfun),
+   times=10)
+ ), c(1, 4, 5))
> # Stop R processes over cluster under Windows
> stopCluster(cluster)
```


Initializing Parallel Clusters Under *Windows*

Under *Windows* the child processes in the parallel compute cluster don't inherit data and objects from their parent process.

Therefore the required data must be either passed into `parLapply()` via the dots `"..."` argument, or by calling the function `clusterExport()`.

Objects from packages must be either referenced using the double-colon operator `::`, or the packages must be loaded in the child processes.

```
> basep <- 2
> # Fails because child processes don't know basep:
> parLapply(cluster, 2:4,
+   function(exponent) basep^exponent)
> # basep passed to child via dots ... argument:
> parLapply(cluster, 2:4,
+   function(exponent, basep) basep^exponent,
+   basep=basep)
> # basep passed to child via clusterExport:
> clusterExport(cluster, "basep")
> parLapply(cluster, 2:4,
+   function(exponent) basep^exponent)
> # Fails because child processes don't know zoo::index():
> parSapply(cluster, c("VTI", "IEF", "DBC"),
+   function(symbol)
+     NROW(zoo::index(get(symbol, envir=rutils::etfenv))))
> # zoo function referenced using "::" in child process:
> parSapply(cluster, c("VTI", "IEF", "DBC"),
+   function(symbol)
+     NROW(zoo::index(get(symbol, envir=rutils::etfenv))))
> # Package zoo loaded in child process:
> parSapply(cluster, c("VTI", "IEF", "DBC"),
+   function(symbol) {
+     stopifnot("package:zoo" %in% search() || require("zoo", quiet=TRUE))
+     NROW(zoo::index(get(symbol, envir=rutils::etfenv)))
+   }) # end parSapply
> # Stop R processes over cluster under Windows
> stopCluster(cluster)
```

Reproducible Parallel Simulations Under *Windows*

Simulations use pseudo-random number generators, and in order to perform reproducible results, they must set the *seed* value, so that the number generators produce the same sequence of pseudo-random numbers.

The function `set.seed()` initializes the random number generator by specifying the *seed* value, so that the number generator produces the same sequence of numbers for a given *seed* value.

But under *Windows* `set.seed()` doesn't initialize the random number generators of child processes, and they don't produce the same sequence of numbers.

The function `clusterSetRNGStream()` initializes the random number generators of child processes under *Windows*.

The function `set.seed()` does initialize the random number generators of child processes under *Mac-OSX* and *Linux*.

```
> library(parallel) # Load package parallel
> # Calculate number of available cores
> ncores <- detectCores() - 1
> # Initialize compute cluster under Windows
> cluster <- makeCluster(ncores)
> # Set seed for cluster under Windows
> # Doesn't work: set.seed(1121)
> clusterSetRNGStream(cluster, 1121)
> # Perform parallel loop under Windows
> output <- parLapply(cluster, 1:70, rnorm, n=100)
> sum(unlist(output))
> # Stop R processes over cluster under Windows
> stopCluster(cluster)
> # Perform parallel loop under Mac-OSX or Linux
> output <- mclapply(1:10, rnorm, mc.cores=ncores, n=100)
```

Monte Carlo Simulation

Monte Carlo simulation consists of generating random samples from a given probability distribution.

The *Monte Carlo* data samples can then be used to calculate different parameters of the probability distribution (moments, quantiles, etc.), and its functionals.

The *quantile* of a probability distribution is the value of the *random variable* x , such that the probability of values less than x is equal to the given *probability* p .

The *quantile* of a data sample can be calculated by first sorting the sample, and then finding the value corresponding closest to the given *probability* p .

The function `quantile()` calculates the sample quantiles. It uses interpolation to improve the accuracy. Information about the different interpolation methods can be found by typing `?quantile`.

The function `sort()` returns a vector sorted into ascending order.

```
> set.seed(1121) # Reset random number generator
> # Sample from Standard Normal Distribution
> nrows <- 1000
> datav <- rnorm(nrows)
> # Sample mean - MC estimate
> mean(datav)
> # Sample standard deviation - MC estimate
> sd(datav)
> # Monte Carlo estimate of cumulative probability
> pnorm(-2)
> sum(datav < (-2))/nrows
> # Monte Carlo estimate of quantile
> confl <- 0.02
> qnorm(confl) # Exact value
> cutoff <- confl*nrows
> datav <- sort(datav)
> datav[cutoff] # Naive Monte Carlo value
> quantile(datav, probs=confl)
> # Analyze the source code of quantile()
> stats:::quantile.default
> # Microbenchmark quantile
> library(microbenchmark)
> summary(microbenchmark(
+   monte_carlo = datav[cutoff],
+   quantilev = quantile(datav, probs=confl),
+   times=100))[, c(1, 4, 5)] # end microbenchmark summary
```

Standard Errors of Estimators Using Bootstrap Simulation

The *bootstrap* procedure uses *Monte Carlo* simulation to generate a distribution of estimator values.

The *bootstrap* procedure generates new data by randomly sampling with replacement from the observed (empirical) data set.

If the original data consists of simulated random numbers then we simply simulate another set of these random numbers.

The *bootstrapped* datasets are used to recalculate the estimator many times, to provide a distribution of the estimator and its standard error.

```
> # Sample from Standard Normal Distribution
> nrows <- 1000; datav <- rnorm(nrows)
> # Sample mean and standard deviation
> mean(datav); sd(datav)
> # Bootstrap of sample mean and median
> nboot <- 10000
> bootd <- sapply(1:nboot, function(x) {
+   # Sample from Standard Normal Distribution
+   samplev <- rnorm(nrows)
+   c(mean=mean(samplev), median=median(samplev))
+ }) # end sapply
> bootd[, 1:3]
> bootd <- t(bootd)
> # Standard error from formula
> sd(datav)/sqrt(nrows)
> # Standard error of mean from bootstrap
> sd(bootd[, "mean"])
> # Standard error of median from bootstrap
> sd(bootd[, "median"])
```

The Distribution of Estimators Using Bootstrap Simulation

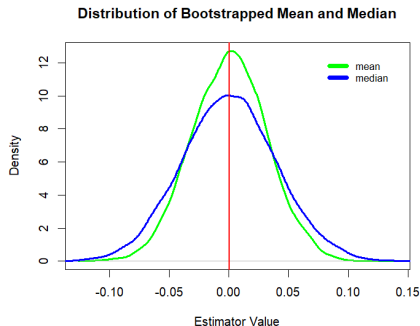
The standard errors of estimators can be calculated using a *bootstrap* simulation.

The *bootstrap* procedure generates new data by randomly sampling with replacement from the observed (empirical) data set.

The *bootstrapped* dataset is used to recalculate the estimator many times.

The *bootstrapped* estimator values are then used to calculate the probability distribution of the estimator and its standard error.

The function `density()` calculates a kernel estimate of the probability density for a sample of data.



```
> # Plot the densities of the bootstrap data
> x11(width=6, height=5)
> plot(density(boot[, "mean"]), lwd=3, xlab="Estimator Value",
+      main="Distribution of Bootstrapped Mean and Median", col="green",
+      lwd=6, bg="white", col=c("green", "blue"))
> lines(density(boot[, "median"]), lwd=3, col="blue")
> abline(v=mean(boot[, "mean"]), lwd=2, col="red")
> legend("topright", inset=0.05, cex=0.8, title=NULL,
+      leg=c("mean", "median"), bty="n",
+      lwd=6, bg="white", col=c("green", "blue"))
```

Bootstrapping Using Vectorized Operations

Bootstrap simulations can be accelerated by using vectorized operations instead of R loops.

But using vectorized operations requires calculating a matrix of random data, instead of calculating random vectors in a loop.

This is another example of the tradeoff between speed and memory usage in simulations.

Faster code often requires more memory than slower code.

```
> set.seed(1121) # Reset random number generator
> nrows <- 1000
> # Bootstrap of sample mean and median
> nboot <- 100
> bootd <- sapply(1:nboot, function(x) median(rnorm(nrows)))
> # Perform vectorized bootstrap
> set.seed(1121) # Reset random number generator
> # Calculate matrix of random data
> samplev <- matrix(rnorm(nboot*nrows), ncol=nboot)
> bootv <- Rfast::colMedians(samplev)
> all.equal(bootd, bootv)
> # Compare speed of loops with vectorized R code
> library(microbenchmark)
> summary(microbenchmark(
+   loop = sapply(1:nboot, function(x) median(rnorm(nrows))),
+   cpp = {
+     samplev <- matrix(rnorm(nboot*nrows), ncol=nboot)
+     Rfast::colMedians(samplev)
+   },
+   times=10))[, c(1, 4, 5)] # end microbenchmark summary
```

Bootstrapping Standard Errors Using Parallel Computing

The *bootstrap* procedure performs a loop, which naturally lends itself to parallel computing.

Different functions from package *parallel* need to be called depending on the operating system (*Windows*, *Mac-OSX*, or *Linux*).

The function `makeCluster()` starts running R processes on several CPU cores under *Windows*.

The function `parLapply()` is similar to `lapply()`, and performs loops under *Windows* using parallel computing on several CPU cores.

The R processes started by `makeCluster()` don't inherit any data from the parent R process.

Therefore the required data must be either passed into `parLapply()` via the dots "... " argument, or by calling the function `clusterExport()`.

The function `mclapply()` performs loops using parallel computing on several CPU cores under *Mac-OSX* or *Linux*.

The function `stopCluster()` stops the R processes running on several CPU cores.

```
> library(parallel) # Load package parallel
> ncores <- detectCores() - 1 # Number of cores
> cluster <- makeCluster(ncores) # Initialize compute cluster under
> set.seed(1121) # Reset random number generator
> # Sample from Standard Normal Distribution
> nrows <- 1000
> # Bootstrap mean and median under Windows
> nboot <- 10000
> bootd <- parLapply(cluster, 1:nboot,
+   function(x, datav, nrows) {
+     samplev <- rnorm(nrows)
+     c(mean=mean(samplev), median=median(samplev))
+   }, datav=datav, nrows=nrows) # end parLapply
> # Bootstrap mean and median under Mac-OSX or Linux
> bootd <- mclapply(1:nboot,
+   function(x) {
+     samplev <- rnorm(nrows)
+     c(mean=mean(samplev), median=median(samplev))
+   }, mc.cores=ncores) # end mclapply
> bootd <- rutils::do_call(rbind, bootd)
> # Means and standard errors from bootstrap
> apply(bootd, MARGIN=2, function(x)
+   c(mean=mean(x), stderor=sd(x)))
> # Standard error from formula
> sd(datav)/sqrt(nrows)
> stopCluster(cluster) # Stop R processes over cluster under Windows
```

Parallel Bootstrapping of the *Median Absolute Deviation*

The *Median Absolute Deviation* (*MAD*) is a robust measure of dispersion (variability), defined using the median instead of the mean:

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

The advantage of *MAD* is that it's always well defined, even for data that has infinite variance.

For normally distributed data the *MAD* has a larger standard error than the standard deviation.

But for distributions with fat tails (like asset returns), the standard deviation has a larger standard error than the *MAD*.

The *MAD* for normally distributed data is equal to $\Phi^{-1}(0.75) \cdot \hat{\sigma} = 0.6745 \cdot \hat{\sigma}$.

The function `mad()` calculates the *MAD* and divides it by $\Phi^{-1}(0.75)$ to make it comparable to the standard deviation.

```
> nrows <- 1000
> datav <- rnorm(nrows)
> sd(datav); mad(datav)
> median(abs(datav - median(datav)))
> median(abs(datav - median(datav)))/qnorm(0.75)
> # Bootstrap of sd and mad estimators
> nboot <- 10000
> bootd <- sapply(1:nboot, function(x) {
+   samplev <- rnorm(nrows)
+   c(sd=sd(samplev), mad=mad(samplev))
+ }) # end sapply
> bootd <- t(bootd)
> # Analyze bootstrapped variance
> head(bootd)
> sum(is.na(bootd))
> # Means and standard errors from bootstrap
> apply(bootd, MARGIN=2, function(x)
+   c(mean=mean(x), stdev=sd(x)))
> # Parallel bootstrap under Windows
> library(parallel) # Load package parallel
> ncores <- detectCores() - 1 # Number of cores
> cluster <- makeCluster(ncores) # Initialize compute cluster
> bootd <- parLapply(cluster, 1:nboot,
+   function(x, datav) {
+     samplev <- rnorm(nrows)
+     c(sd=sd(samplev), mad=mad(samplev))
+   }, datav=datav) # end parLapply
> # Parallel bootstrap under Mac-OSX or Linux
> bootd <- mclapply(1:nboot, function(x) {
+   samplev <- rnorm(nrows)
+   c(sd=sd(samplev), mad=mad(samplev))
+ }, mc.cores=ncores) # end mclapply
> stopCluster(cluster) # Stop R processes over cluster
> bootd <- rutils::do_call(rbind, bootd)
> # Means and standard errors from bootstrap
> apply(bootd, MARGIN=2, function(x)
+   c(mean=mean(x), stdev=sd(x)))
```


Resampling From Empirical Datasets

Resampling is randomly selecting data from an existing dataset, to create a new dataset with similar properties to the existing dataset.

Resampling is usually performed with replacement, so that each draw is independent from the others.

Resampling is performed when it's not possible or convenient to obtain another set of empirical data, so we simulate a new data set by randomly sampling from the existing data.

The function `sample()` selects a random sample from a vector of data elements.

The function `sample.int()` is a *method* that selects a random sample of *integers*.

The function `sample.int()` with argument `replace=TRUE` selects a sample with replacement (the *integers* can repeat).

The function `sample.int()` is a little faster than `sample()`.

```
> # Calculate time series of VTI returns
> library(rutils)
> retp <- rutils::etfenv$returns$VTI
> retp <- na.omit(retp)
> nrow <- NROW(retp)
> # Sample from VTI returns
> samplev <- retp[sample.int(nrow, replace=TRUE)]
> c(sd=sd(samplev), mad=mad(samplev))
> # sample.int() is a little faster than sample()
> library(microbenchmark)
> summary(microbenchmark(
+   sample.int = sample.int(1e3),
+   sample = sample(1e3),
+   times=10))[, c(1, 4, 5)]
```

Bootstrapping From Empirical Datasets

Bootstrapping is usually performed by resampling from an observed (empirical) dataset.

Resampling consists of randomly selecting data from an existing dataset, with replacement.

Resampling produces a new *bootstrapped* dataset with similar properties to the existing dataset.

The *bootstrapped* dataset is used to recalculate the estimator many times.

The *bootstrapped* estimator values are then used to calculate the probability distribution of the estimator and its standard error.

Bootstrapping shows that for asset returns, the *Median Absolute Deviation (MAD)* has a smaller relative standard error than the standard deviation.

Bootstrapping doesn't provide accurate estimates for estimators which are sensitive to the ordering and correlations in the data.

```
> # Sample from time series of VTI returns
> library(rutils)
> retp <- rutils::etfenv$returns$VTI
> retp <- na.omit(retp)
> nrow <- NROW(retp)
> # Bootstrap sd and MAD under Windows
> library(parallel) # Load package parallel
> ncores <- detectCores() - 1 # Number of cores
> cluster <- makeCluster(ncores) # Initialize compute cluster under Windows
> clusterSetRNGStream(cluster, 1121) # Reset random number generator
> nboot <- 10000
> bootd <- parLapply(cluster, 1:nboot,
+   function(x, retp, nrow) {
+     samplev <- retp[sample.int(nrow, replace=TRUE)]
+     c(sd=sd(samplev), mad=mad(samplev))
+   }, retp=retp, nrow=nrow) # end parLapply
> # Bootstrap sd and MAD under Mac-OSX or Linux
> bootd <- mclapply(1:nboot, function(x) {
+   samplev <- retp[sample.int(nrow, replace=TRUE)]
+   c(sd=sd(samplev), mad=mad(samplev))
+ }, mc.cores=ncores) # end mclapply
> stopCluster(cluster) # Stop R processes over cluster under Windows
> bootd <- rutils::do_call(rbind, bootd)
> # Standard error assuming normal distribution of returns
> sd(retp)/sqrt(nboot)
> # Means and standard errors from bootstrap
> stderrors <- apply(bootd, MARGIN=2,
+   function(x) c(mean=mean(x), stdev=sd(x)))
> stderrors
> # Relative standard errors
> stderrors[2, ]/stderrors[1, ]
```

Standard Errors of Regression Coefficients Using Bootstrap

The standard errors of the regression coefficients can be calculated using a *bootstrap* simulation.

The *bootstrap* procedure creates new design matrices by randomly sampling with replacement from the regression design matrix.

Regressions are performed on the *bootstrapped* design matrices, and the regression coefficients are saved into a matrix of *bootstrapped* coefficients.

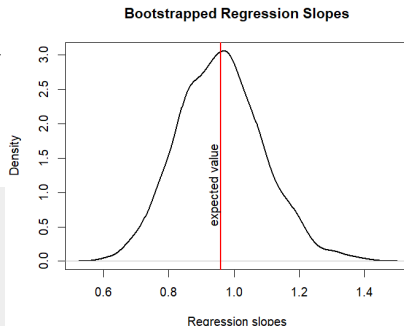
```
> # Initialize random number generator
> set.seed(1121)
> # Define explanatory and response variables
> nrows <- 100
> predv <- rnorm(nrows, mean=2)
> noise <- rnorm(nrows)
> respv <- (-3 + 2*predictor + noise)
> desv <- cbind(respv, predv)
> # Calculate alpha and beta regression coefficients
> betav <- cov(desv[, 1], desv[, 2])/var(desv[, 2])
> alpha <- mean(desv[, 1]) - betav*mean(desv[, 2])
> x11(width=6, height=5)
> plot(respv ~ predv, data=desv)
> abline(a=alpha, b=betav, lwd=3, col="blue")
> # Bootstrap of beta regression coefficient
> nboot <- 100
> bootd <- sapply(1:nboot, function(x) {
+   samplev <- sample.int(nrows, replace=TRUE)
+   desv <- desv[samplev, ]
+   cov(desv[, 1], desv[, 2])/var(desv[, 2])
+ }) # end sapply
```

The *bootstrapped* coefficient values can be used to calculate the probability distribution of the coefficients and their standard errors.

`abline()` plots a straight line on the existing plot.

The function `text()` draws text on a plot, and can be used to draw plot labels.

```
> # Mean and standard error of beta regression coefficient
> c(mean=mean(bootd), stdev=sqrt(sd(bootd)))
> # Plot density of bootstrapped beta coefficients
> plot(density(bootd), lwd=2, xlab="Regression slopes",
+      main="Bootstrapped Regression Slopes")
> # Add line for expected value
> abline(v=mean(bootd), lwd=2, col="red")
> text(x=mean(bootd)-0.01, y=1.0, labels="expected value",
+      lwd=2, srt=90, pos=3)
```



Bootstrapping Regressions Using Parallel Computing

The *bootstrap* procedure performs a loop, which naturally lends itself to parallel computing.

Different functions from package *parallel* need to be called depending on the operating system (*Windows*, *Mac-OSX*, or *Linux*).

The function `makeCluster()` starts running R processes on several CPU cores under *Windows*.

The function `parLapply()` is similar to `lapply()`, and performs loops under *Windows* using parallel computing on several CPU cores.

The R processes started by `makeCluster()` don't inherit any data from the parent R process.

Therefore the required data must be passed into `parLapply()` via the dots `"..."` argument.

The function `mclapply()` performs loops using parallel computing on several CPU cores under *Mac-OSX* or *Linux*.

The function `stopCluster()` stops the R processes running on several CPU cores.

```
> library(parallel) # Load package parallel
> ncores <- detectCores() - 1 # Number of cores
> cluster <- makeCluster(ncores) # Initialize compute cluster under Windows
> # Bootstrap of regression under Windows
> bootd <- parLapply(cluster, 1:1000,
+   function(x, desv) {
+     samplev <- sample.int(nrows, replace=TRUE)
+     desv <- desv[samplev, ]
+     cov(desv[, 1], desv[, 2])/var(desv[, 2])
+   }, design=desv) # end parLapply
> # Bootstrap of regression under Mac-OSX or Linux
> bootd <- mclapply(1:1000,
+   function(x) {
+     samplev <- sample.int(nrows, replace=TRUE)
+     desv <- desv[samplev, ]
+     cov(desv[, 1], desv[, 2])/var(desv[, 2])
+   }, mc.cores=ncores) # end mclapply
> stopCluster(cluster) # Stop R processes over cluster under Windows
```

Analyzing the Bootstrap Data

The *bootstrap* loop produces a *list* which can be collapsed into a vector.

The function `unlist()` collapses a list with atomic elements into a vector (which can cause type coercion).

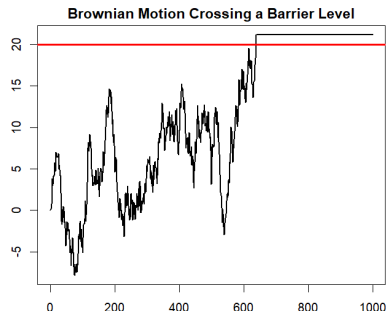
```
> # Collapse the bootstrap list into a vector
> class(bootd)
> bootd <- unlist(bootd)
> # Mean and standard error of beta regression coefficient
> c(mean=mean(bootd), stderror=sd(bootd))
> # Plot density of bootstrapped beta coefficients
> plot(density(bootd),
+      lwd=2, xlab="Regression slopes",
+      main="Bootstrapped Regression Slopes")
> # Add line for expected value
> abline(v=mean(bootd), lwd=2, col="red")
> text(x=mean(bootd)-0.01, y=1.0, labels="expected value",
+      lwd=2, srt=90, pos=3)
```

Simulating Brownian Motion Using while() Loops

while() loops are often used in simulations, when the number of required loops is unknown in advance.

Below is an example of a simulation of the path of *Brownian Motion* crossing a barrier level.

```
> set.seed(1121) # Reset random number generator
> barl <- 20 # Barrier level
> nrows <- 1000 # Number of simulation steps
> pathv <- numeric(nrows) # Allocate path vector
> pathv[1] <- 0 # Initialize path
> it <- 2 # Initialize simulation index
> while ((it <= nrows) && (pathv[it - 1] < barl)) {
+ # Simulate next step
+   pathv[it] <- pathv[it - 1] + rnorm(1)
+   it <- it + 1 # Advance index
+ } # end while
> # Fill remaining path after it crosses barl
> if (it <= nrows)
+   pathv[it:nrows] <- pathv[it - 1]
> # Plot the Brownian motion
> x11(width=6, height=5)
> par(mar=c(3, 3, 2, 1), oma=c(1, 1, 1, 1))
> plot(pathv, type="l", col="black",
+       lty="solid", lwd=2, xlab="", ylab="")
> abline(h=barl, lwd=3, col="red")
> title(main="Brownian Motion Crossing a Barrier Level", line=0.5)
```

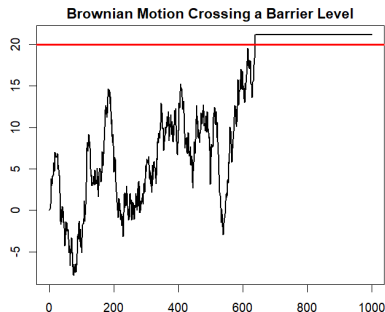


Simulating Brownian Motion Using Vectorized Functions

Simulations in R can be accelerated by pre-computing a vector of random numbers, instead of generating them one at a time in a loop.

Vectors of random numbers allow using *vectorized* functions, instead of inefficient (slow) `while()` loops.

```
> set.seed(1121) # Reset random number generator
> barl <- 20 # Barrier level
> nrows <- 1000 # Number of simulation steps
> # Simulate path of Brownian motion
> pathv <- cumsum(rnorm(nrows))
> # Find index when path crosses barl
> crossp <- which(pathv > barl)
> # Fill remaining path after it crosses barl
> if (NROW(crossp)>0) {
+   pathv[(crossp[1]+1):nrows] <- pathv[crossp[1]]
+ } # end if
> # Plot the Brownian motion
> x11(width=6, height=5)
> par(mar=c(3, 3, 2, 1), oma=c(1, 1, 1, 1))
> plot(pathv, type="l", col="black",
+      lty="solid", lwd=2, xlab="", ylab="")
> abline(h=barl, lwd=3, col="red")
> title(main="Brownian Motion Crossing a Barrier Level", line=0.5)
```



The tradeoff between speed and memory usage: more memory may be used than necessary, since the simulation may stop before all the pre-computed random numbers are used up.

But the simulation is much faster because the path is simulated using *vectorized* functions,

Estimating the Statistics of Brownian Motion

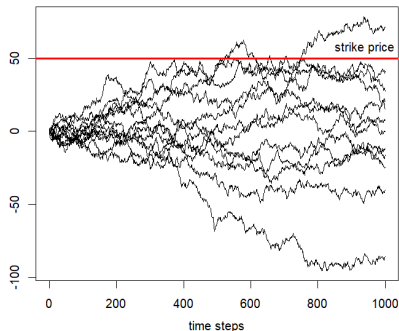
The statistics of Brownian motion can be estimated by simulating multiple paths.

An example of a statistic is the expected value of Brownian motion at a fixed time horizon, which is the option payout for the strike price k : $\mathbb{E}[(p_t - k)_+]$.

Another statistic is the probability of Brownian motion crossing a boundary (barrier) b : $\mathbb{E}[\mathbb{1}(p_t - b)]$.

```
> # Define Brownian motion parameters
> sigmav <- 1.0 # Volatility
> drift <- 0.0 # Drift
> nrows <- 1000 # Number of simulation steps
> nsimu <- 100 # Number of simulations
> # Simulate multiple paths of Brownian motion
> set.seed(1121)
> pathm <- rnorm(nsimu*nrows, mean=drift, sd=sigmav)
> pathm <- matrix(pathm, nc=nsimu)
> pathm <- matrixStats::colCumsums(pathm)
> # Final distribution of paths
> mean(pathm[nrows, ]) ; sd(pathm[nrows, ])
> # Calculate option payout at maturity
> strikep <- 50 # Strike price
> payouts <- (pathm[nrows, ] - strikep)
> sum(payouts[payouts > 0])/nsimu
> # Calculate probability of crossing the barrier at any point
> bar1 <- 50
> crossi <- (colSums(pathm > bar1) > 0)
> sum(crossi)/nsimu
```

Paths of Brownian Motion



```
> # Plot in window
> x11(width=6, height=5)
> par(mar=c(4, 3, 2, 2), oma=c(0, 0, 0, 0), mgp=c(2.5, 1, 0))
> # Select and plot full range of paths
> ordern <- order(pathm[nrows, ])
> pathm[nrows, ordern]
> indeks <- ordern[seq(1, 100, 9)]
> zoo::plot.zoo(pathm[, indeks], main="Paths of Brownian Motion",
+   xlab="time steps", ylab=NA, plot.type="single")
> abline(h=strikep, col="red", lwd=3)
> text(x=(nrows-60), y=strikep, labels="strike price", pos=3, cex=1.5)
```

Bootstrapping From Time Series of Prices

Bootstrapping from a time series of prices requires first converting the prices to *percentage* returns, then bootstrapping the returns, and finally converting them back to prices.

Bootstrapping from *percentage* returns ensures that the bootstrapped prices are not negative.

Below is a simulation of the frequency of bootstrapped prices crossing a barrier level.

```
> # Calculate percentage returns from VTI prices
> library(rutils)
> pricev <- quantmod::Cl(rutils::etfenv$VTI)
> startd <- as.numeric(pricev[1, ])
> retp <- rutils::diffit(log(pricev))
> class(retp); head(retp)
> sum(is.na(retp))
> nrows <- NROW(retp)
> # Define barrier level with respect to prices
> barl <- 1.5*max(pricev)
> # Calculate single bootstrap sample
> samplev <- retp[sample.int(nrows, replace=TRUE)]
> # Calculate prices from percentage returns
> samplev <- startd*exp(cumsum(samplev))
> # Calculate if prices crossed barrier
> sum(samplev > barl) > 0
```

```
> library(parallel) # Load package parallel
> ncores <- detectCores() - 1 # Number of cores
> cluster <- makeCluster(ncores) # Initialize compute cluster under Windows
> # Perform parallel bootstrap under Windows
> clusterSetRNGStream(cluster, 1121) # Reset random number generator
> clusterExport(cluster, c("startd", "barl"))
> nboot <- 10000
> bootd <- parLapply(cluster, 1:nboot,
+   function(x, retp, nrows) {
+     samplev <- retp[sample.int(nrows, replace=TRUE)]
+     # Calculate prices from percentage returns
+     samplev <- startd*exp(cumsum(samplev))
+     # Calculate if prices crossed barrier
+     sum(samplev > barl) > 0
+   }, retp=retp, nrows=nrows) # end parLapply
> # Perform parallel bootstrap under Mac-OSX or Linux
> bootd <- mclapply(1:nboot, function(x) {
+   samplev <- retp[sample.int(nrows, replace=TRUE)]
+   # Calculate prices from percentage returns
+   samplev <- startd*exp(cumsum(samplev))
+   # Calculate if prices crossed barrier
+   sum(samplev > barl) > 0
+ }, mc.cores=ncores) # end mclapply
> stopCluster(cluster) # Stop R processes over cluster under Windows
> bootd <- rutils::do.call(rbind, bootd)
> # Calculate frequency of crossing barrier
> sum(bootd)/nboot
```

Bootstrapping From OHLC Prices

Bootstrapping from OHLC prices requires updating all the price columns, not just the *Close* prices.

The *Close* prices are bootstrapped first, and then the other columns are updated using the differences of the OHLC price columns.

Below is a simulation of the frequency of the *High* prices crossing a barrier level.

```
> # Calculate percentage returns from VTI prices
> library(rutils)
> ohlc <- rutils::etfenv$VTI
> pricev <- as.numeric(ohlc[, 4])
> startd <- pricev[1]
> retp <- rutils::diffit(log(pricev))
> nrows <- NROW(retp)
> # Calculate difference of OHLC price columns
> ohlc_diff <- ohlc[, 1:3] - pricev
> class(retp); head(retp)
> # Calculate bootstrap prices from percentage returns
> datav <- sample.int(nrows, replace=TRUE)
> boot_pricev <- startd*exp(cumsum(retp[datav]))
> boot_ohlc <- ohlc_diff + boot_prices
> boot_ohlc <- cbind(boot_ohlc, boot_pricev)
> # Define barrier level with respect to prices
> barl <- 1.5*max(pricev)
> # Calculate if High bootstrapped prices crossed barrier level
> sum(boot_ohlc[, 2] > barl) > 0
```

```
> library(parallel) # Load package parallel
> ncores <- detectCores() - 1 # Number of cores
> cluster <- makeCluster(ncores) # Initialize compute cluster under Windows
> # Perform parallel bootstrap under Windows
> clusterSetRNGStream(cluster, 1121) # Reset random number generator
> clusterExport(cluster, c("startd", "barl", "ohlc_diff"))
> nboot <- 10000
> bootd <- parLapply(cluster, 1:nboot,
+   function(x, retp, nrows) {
+     # Calculate OHLC prices from percentage returns
+     datav <- sample.int(nrows, replace=TRUE)
+     boot_pricev <- startd*exp(cumsum(retp[datav]))
+     boot_ohlc <- ohlc_diff + boot_prices
+     boot_ohlc <- cbind(boot_ohlc, boot_pricev)
+     # Calculate statistic
+     sum(boot_ohlc[, 2] > barl) > 0
+   }, retp=retp, nrows=nrows) # end parLapply
> # Perform parallel bootstrap under Mac-OSX or Linux
> bootd <- mclapply(1:nboot, function(x) {
+   # Calculate OHLC prices from percentage returns
+   datav <- sample.int(nrows, replace=TRUE)
+   boot_pricev <- startd*exp(cumsum(retp[datav]))
+   boot_ohlc <- ohlc_diff + boot_prices
+   boot_ohlc <- cbind(boot_ohlc, boot_pricev)
+   # Calculate statistic
+   sum(boot_ohlc[, 2] > barl) > 0
+ }, mc.cores=ncores) # end mclapply
> stopCluster(cluster) # Stop R processes over cluster under Windows
> bootd <- rutils::do_call(rbind, bootd)
> # Calculate frequency of crossing barrier
> sum(bootd)/nboot
```

The *ETF* Database

Exchange-traded Funds (*ETFs*) are funds which invest in portfolios of assets, such as stocks, commodities, or bonds.

ETFs are shares in portfolios of assets, and they are traded just like stocks.

ETFs provide investors with convenient, low cost, and liquid instruments to invest in various portfolios of assets.

The file `etf_list.csv` contains a database of exchange-traded funds (*ETFs*) and exchange traded notes (*ETNs*).

We will select a portfolio of *ETFs* for illustrating various investment strategies.

```
> # Select ETF symbols for asset allocation
> symbolv <- c("VTI", "VEU", "EEM", "XLY", "XLP", "XLE", "XLF",
+ "XLV", "XLI", "XLB", "XLK", "XLU", "VYM", "IVW", "IWB", "IWD",
+ "IWF", "IEF", "TLT", "VNQ", "DBC", "GLD", "USO", "VXX", "SVXY",
+ "MTUM", "IVE", "VLUE", "QUAL", "VTV", "USMV", "AIEQ")
> # Read etf database into data frame
> etflist <- read.csv(file="/Users/jerzy/Develop/lecture_slides/data/etf_list.csv")
> rownames(etflist) <- etflist$Symbol
> # Select from etflist only those ETF's in symbolv
> etflist <- etflist[symbolv, ]
> # Shorten names
> etfnames <- sapply(etflist$Name, function(name) {
+   namesplit <- strsplit(name, split=" ")[1]
+   namesplit <- namesplit[c(-1, -NROW(namesplit))]
+   name_match <- match("Select", namesplit)
+   if (!is.na(name_match))
+     namesplit <- namesplit[-name_match]
+   paste(namesplit, collapse=" ")
+ }) # end sapply
> etflist$Name <- etfnames
> etflist["IEF", "Name"] <- "10 year Treasury Bond Fund"
> etflist["TLT", "Name"] <- "20 plus year Treasury Bond Fund"
> etflist["XLY", "Name"] <- "Consumer Discr. Sector Fund"
> etflist["EEM", "Name"] <- "Emerging Market Stock Fund"
> etflist["MTUM", "Name"] <- "Momentum Factor Fund"
> etflist["SVXY", "Name"] <- "Short VIX Futures"
> etflist["VXX", "Name"] <- "Long VIX Futures"
> etflist["DBC", "Name"] <- "Commodity Futures Fund"
> etflist["USO", "Name"] <- "WTI Oil Futures Fund"
> etflist["GLD", "Name"] <- "Physical Gold Fund"
```

ETF Portfolio for Investment Strategies

The portfolio contains *ETFs* representing different *industry sectors* and *investment styles*.

The *ETFs* with names *X** represent *industry sector funds* (energy, financial, etc.)

The *ETFs* with names *I** represent *style funds* (value, growth, size).

IWB is the Russell 1000 small-cap fund.

MTUM is an *ETF* which owns a stock portfolio representing the *momentum factor*.

DBC is an *ETF* providing the total return on a portfolio of commodity futures.

VXX is an *ETN* providing the total return of *long VIX* futures contracts (specifically the *S&P VIX Short-Term Futures Index*).

VXX is *bearish* because it's *long VIX* futures, and the *VIX* *rises* when stock prices *drop*.

SVXY is an *ETF* providing the total return of *short VIX* futures contracts.

SVXY is *bullish* because it's *short VIX* futures, and the *VIX* *drops* when stock prices *rise*.

Symbol	Name	Fund.Type
VTI	Total Stock Market	US Equity ETF
VEU	FTSE All World Ex US	Global Equity ETF
EEM	Emerging Market Stock Fund	Global Equity ETF
XLY	Consumer Discr. Sector Fund	US Equity ETF
XLP	Consumer Staples Sector Fund	US Equity ETF
XLE	Energy Sector Fund	US Equity ETF
XLF	Financial Sector Fund	US Equity ETF
XLV	Health Care Sector Fund	US Equity ETF
XLI	Industrial Sector Fund	US Equity ETF
XLB	Materials Sector Fund	US Equity ETF
XLK	Technology Sector Fund	US Equity ETF
XLU	Utilities Sector Fund	US Equity ETF
VYM	Large-cap Value	US Equity ETF
IWV	S&P 500 Growth Index Fund	US Equity ETF
IWB	Russell 1000	US Equity ETF
IWD	Russell 1000 Value	US Equity ETF
IWF	Russell 1000 Growth	US Equity ETF
IEF	10 year Treasury Bond Fund	US Fixed Income ETF
TLT	20 plus year Treasury Bond Fund	US Fixed Income ETF
VNQ	REIT ETF - DNQ	US Equity ETF
DBC	Commodity Futures Fund	Commodity Based ETF
GLD	Physical Gold Fund	Commodity Based ETF
USO	WTI Oil Futures Fund	Commodity Based ETF
VXX	Long VIX Futures	Commodity Based ETN
SVXY	Short VIX Futures	Commodity Based ETN
MTUM	Momentum Factor Fund	US Equity ETF
IVE	S&P 500 Value Index Fund	US Equity ETF
VLUE	MSCI USA Value Factor	US Equity ETF
QUAL	MSCI USA Quality Factor	US Equity ETF
VTV	Value	US Equity ETF
USMV	MSCI USA Minimum Volatility Fund	US Equity ETF
AIEQ	AI Powered Equity	US Asset Allocation ETN

Exchange Traded Notes (ETNs)

ETNs are similar to *ETFs*, with the difference that *ETFs* are shares in a fund which owns the underlying assets, while *ETNs* are notes from issuers which promise payouts according to a formula tied to the underlying asset.

ETFs are similar to mutual funds, while *ETNs* are similar to corporate bonds.

ETNs are technically unsecured corporate debt, but instead of fixed coupons, they promise to provide returns on a market index or futures contract.

The *ETN* issuer promises the payout and is responsible for tracking the index.

The *ETN* investor has counterparty credit risk to the *ETN* issuer.

Downloading ETF Prices Using Package *quantmod*

The function `getSymbols()` downloads time series data into the specified *environment*.

`getSymbols()` downloads the daily *OHLC* prices and trading volume (Open, High, Low, Close, Adjusted, Volume).

`getSymbols()` creates objects in the specified *environment* from the input strings (names), and assigns the data to those objects, without returning them as a function value, as a *side effect*.

If the argument "auto.assign" is set to `FALSE`, then `getSymbols()` returns the data, instead of assigning it silently.

Yahoo data quality deteriorated significantly in 2017, and *Google* data quality is also poor, leaving *Tiingo* and *Alpha Vantage* as the only major providers of free daily *OHLC* stock prices.

But *Quandl* doesn't provide free *ETF* prices, leaving *Alpha Vantage* as the best provider of free daily *ETF* prices.

```
> # Select ETF symbols for asset allocation
> symbolv <- c("VTI", "VEU", "EEM", "XLY", "XLP", "XLE", "XLF",
+ "XLV", "XLI", "XLB", "XLK", "XLU", "VYM", "IVW", "IWB", "IWD",
+ "IWF", "IEF", "TLT", "VNQ", "DBC", "GLD", "USO", "VXX", "SVXY",
+ "MTUM", "IVE", "VLUE", "QUAL", "VTI", "USMV", "AIEQ")
> library(rutils) # Load package rutils
> etfenv <- new.env() # New environment for data
> # Boolean vector of symbols already downloaded
> isdownloaded <- symbolv %in% ls(etfenv)
> # Download data for symbolv using single command - creates pacing
> getSymbols.av(symbolv, adjust=TRUE, env=etfenv,
+ output.size="full", api.key="T7JPW54ES8G75310")
> # Download data from Alpha Vantage using while loop
> n attempts <- 0 # number of download attempts
> while ((sum(!isdownloaded) > 0) & (n attempts < 10)) {
+ # Download data and copy it into environment
+ n attempts <- n attempts + 1
+ cat("Download attempt = ", n attempts, "\n")
+ for (symbol in na.omit(symbolv[!isdownloaded][1:5])) {
+ cat("Processing: ", symbol, "\n")
+ tryCatch( # With error handler
+ quantmod::getSymbols.av(symbol, adjust=TRUE, env=etfenv, auto.assign=
+ # Error handler captures error condition
+ error=function(error_cond) {
+ print(paste("error handler: ", error_cond))
+ }, # end error handler
+ finally=print(paste("symbol=", symbol))
+ ) # end tryCatch
+ } # end for
+ # Update vector of symbols already downloaded
+ isdownloaded <- symbolv %in% ls(etfenv)
+ cat("Pausing 1 minute to avoid pacing...\n")
+ Sys.sleep(65)
+ } # end while
> # Download all symbolv using single command - creates pacing error
> # quantmod::getSymbols.av(symbolv, env=etfenv, adjust=TRUE, from=
```

Inspecting ETF Prices in an Environment

The function `get()` retrieves objects that are referenced using character strings, instead of their names.

The function `eapply()` is similar to `lapply()`, and applies a function to objects in an *environment*, and returns a list.

```
> ls(etfenv) # List files in etfenv
> # Get class of object in etfenv
> class(get(x=symbolv[1], envir=etfenv))
> # Another way
> class(etfenv$VTI)
> colnames(etfenv$VTI)
> # Get first 3 rows of data
> head(etfenv$VTI, 3)
> # Get last 11 rows of data
> tail(etfenv$VTI, 11)
> # Get class of all objects in etfenv
> eapply(etfenv, class)
> # Get class of all objects in R workspace
> lapply(ls(), function(ob_jct) class(get(ob_jct)))
> # Get end dates of all objects in etfenv
> as.Date(sapply(etfenv, end))
```


Adjusting Stock Prices Using Package *quantmod*

Traded stock and bond prices experience jumps after splits and dividends, and must be adjusted to account for them.

The function `adjustOHLC()` adjusts *OHLC* prices.

The function `get()` retrieves objects that are referenced using character strings, instead of their names.

The function `assign()` assigns a value to an object in a specified *environment*, by referencing it using a character string (name).

The functions `get()` and `assign()` allow retrieving and assigning values to objects that are referenced using character strings.

The function `mget()` accepts a vector of strings and returns a list of the corresponding objects extracted from an *environment*.

If the argument "adjust" in function `getSymbols()` is set to `TRUE`, then `getSymbols()` returns adjusted data.

```
> # Check if object is an OHLC time series
> is.OHLC(etfenv$VTI)
> # Adjust single OHLC object using its name
> etfenv$VTI <- adjustOHLC(etfenv$VTI, use.Adjusted=TRUE)
>
> # Adjust OHLC object using string as name
> assign(symbolv[1], adjustOHLC(
+   get(x=symbolv[1], envir=etfenv), use.Adjusted=TRUE),
+   envir=etfenv)
>
> # Adjust objects in environment using vector of strings
> for (symbol in ls(etfenv)) {
+   assign(symbol,
+     adjustOHLC(get(symbol, envir=etfenv), use.Adjusted=TRUE),
+     envir=etfenv)
+ } # end for
```

Extracting Time Series from Environments

The function `mget()` accepts a vector of strings and returns a list of the corresponding objects extracted from an *environment*.

The extractor (accessor) functions from package *quantmod*: `C1()`, `Vo()`, etc., extract columns from *OHLC* data.

A list of *xts* series can be flattened into a single *xts* series using the function `do.call()`.

The function `do.call()` executes a function call using a function name and a list of arguments.

`do.call()` passes the list elements individually, instead of passing the whole list as one argument.

The function `eapply()` is similar to `lapply()`, and applies a function to objects in an *environment*, and returns a list.

Time series can also be extracted from an *environment* by coercing it into a list, and then subsetting and merging it into an *xts* series using the function `do.call()`.

```
> library(rutils) # Load package rutils
> # Define ETF symbols
> symbolv <- c("VTI", "VEU", "IEF", "VNQ")
> # Extract symbolv from rutils::etfenv
> pricev <- mget(symbolv, envir=rutils::etfenv)
> # pricev is a list of xts series
> class(pricev)
> class(pricev[[1]])
> tail(pricev[[1]])
> # Extract close prices
> pricev <- lapply(pricev, quantmod::C1)
> # Collapse list into time series the hard way
> prices2 <- cbind(pricev[[1]], pricev[[2]], pricev[[3]], pricev[[4]])
> class(prices2)
> dim(prices2)
> # Collapse list into time series using do.call()
> pricev <- do.call(cbind, pricev)
> all.equal(prices2, pricev)
> class(pricev)
> dim(pricev)
> # Or extract and cbind in single step
> pricev <- do.call(cbind, lapply(
+   mget(symbolv, envir=rutils::etfenv), quantmod::C1))
> # Or extract and bind all data, subset by symbolv
> pricev <- lapply(symbolv, function(symbol) {
+   quantmod::C1(get(symbol, envir=rutils::etfenv))
+ }) # end lapply
> # Or loop over etfenv without anonymous function
> pricev <- do.call(cbind,
+   lapply(as.list(rutils::etfenv)[symbolv], quantmod::C1))
> # Same, but works only for OHLC series - produces error
> pricev <- do.call(cbind,
+   eapply(rutils::etfenv, quantmod::C1)[symbolv])
```

Managing Time Series

Time series columns can be renamed, and then saved into .csv files.

The function `strsplit()` splits the elements of a character vector.

The package `zoo` contains functions `write.zoo()` and `read.zoo()` for writing and reading `zoo` time series from .txt and .csv files.

The function `eapply()` is similar to `lapply()`, and applies a function to objects in an *environment*, and returns a list.

The function `assign()` assigns a value to an object in a specified *environment*, by referencing it using a character string (name).

The function `save()` writes objects to compressed binary .RData files.

```
> # Column names end with ".Close"
> colnames(pricev)
> strsplit(colnames(pricev), split=".")
> do.call(rbind, strsplit(colnames(pricev), split="."))
> do.call(rbind, strsplit(colnames(pricev), split="."))[, 1]
> # Drop ".Close" from colnames
> colnames(pricev) <- rutils::get_name(colnames(pricev))
> # Or
> # colnames(pricev) <- do.call(rbind,
> #   strsplit(colnames(pricev), split="."))[, 1]
> tail(pricev, 3)
> # Which objects in global environment are class xts?
> unlist(eapply(globalenv(), is.xts))
> # Save xts to csv file
> write.zoo(pricev,
+   file="/Users/jerzy/Develop/lecture_slides/data/etf_series.csv",
> # Copy prices into etfenv
> etfenv$pricev <- pricev
> # Or
> assign("prices", pricev, envir=etfenv)
> # Save to .RData file
> save(etfenv, file="etf_data.RData")
```

Calculating Percentage Returns from Close Prices

The function `quantmod::dailyReturn()` calculates the percentage daily returns from the *Close* prices.

The `lapply()` and `sapply()` functionals perform a loop over the columns of *zoo* and *xts* series.

```
> # Extract VTI prices
> pricev <- etfenv$pricev[, "VTI"]
> pricev <- na.omit(pricev)
> # Calculate percentage returns "by hand"
> pricel <- as.numeric(pricev)
> pricel <- c(pricel[1], pricel[-NROW(pricel)])
> pricel <- xts(pricel, zoo::index(pricev))
> retp <- (pricev-pricel)/pricel
> # Calculate percentage returns using dailyReturn()
> retld <- quantmod::dailyReturn(pricev)
> head(cbind(retld, retp))
> all.equal(retld, retp, check.attributes=FALSE)
> # Calculate returns for all prices in etfenv$prices
> retp <- lapply(etfenv$pricev, function(xtsv) {
+   retld <- quantmod::dailyReturn(na.omit(xtsv))
+   colnames(retld) <- names(xtsv)
+   retld
+ }) # end lapply
> # "retp" is a list of xts
> class(retp)
> class(retp[[1]])
> # Flatten list of xts into a single xts
> retp <- do.call(cbind, retp)
> class(retp)
> dim(retp)
> # Copy retp into etfenv and save to .RData file
> # assign("retp", retp, envir=etfenv)
> etfenv$retp <- retp
> save(etfenv, file="/Users/jerzy/Develop/lecture_slides/data/etf_d
```

Managing Data Inside Environments

The function `as.environment()` coerces objects (list) into an environment.

The function `eapply()` is similar to `lapply()`, and applies a function to objects in an *environment*, and returns a list.

The function `mget()` accepts a vector of strings and returns a list of the corresponding objects extracted from an *environment*.

```
> library(rutils)
> startd <- "2012-05-10"; endd <- "2013-11-20"
> # Select all objects in environment and return as environment
> new_env <- as.environment(eapply(etfenv, "[",
+   paste(startd, endd, sep="/")))
> # Select only symbolv in environment and return as environment
> new_env <- as.environment(
+   lapply(as.list(etfenv)[symbolv], "[",
+     paste(startd, endd, sep="/")))
> # Extract and cbind Close prices and return to environment
> assign("prices", rutils::do_call(cbind,
+   lapply(ls(etfenv), function(symbol) {
+     xtsv <- quantmod::Cl(get(symbol, etfenv))
+     colnames(xtsv) <- symbol
+     xtsv
+   })), envir=new_env)
> # Get sizes of OHLC xts series in etfenv
> sapply(mget(symbolv, envir=etfenv), object.size)
> # Extract and cbind adjusted prices and return to environment
> colname <- function(xtsv)
+   strsplit(colnames(xtsv), split=".[.]"[1])[1]
> assign("prices", rutils::do_call(cbind,
+   lapply(mget(etfenv$symbolv, envir=etfenv),
+     function(xtsv) {
+       xtsv <- Ad(xtsv)
+       colnames(xtsv) <- colname(xtsv)
+       xtsv
+     })), envir=new_env)
```

Loading Stock Tickers

The file `sp500_constituents.csv` contains a *data frame* of *S&P500* constituents.

The stock tickers are stored in the column "Ticker".

The *data frame* contains duplicate tickers, which must be removed.

Some tickers (like "BRK.B" and "BF.B") are not valid symbols in *Tiingo*, so they must be renamed.

```
> # Load data frame of S&P500 constituents from CSV file
> sp500 <- read.csv(file="/Users/jerzy/Develop/lecture_slides/data/sp500.csv")
> # Inspect data frame of S&P500 constituents
> dim(sp500)
> colnames(sp500)
> # Extract tickers from the column Ticker
> symbolv <- sp500$Ticker
> # Get duplicate tickers
> tablev <- table(symbolv)
> duplicates <- tablev[tablev>1]
> duplicates <- names(duplicates)
> # Get duplicate records (rows) of sp500
> sp500[symbolv %in% duplicates, ]
> # Get unique tickers
> symbolv <- unique(symbolv)
> # Find index of ticker "BRK.B"
> which(symbolv=="BRK.B")
> # Rename "BRK.B" to "BRK-B" and "BF.B" to "BF-B"
> symbolv[which(symbolv=="BRK.B")] <- "BRK-B"
> symbolv[which(symbolv=="BF.B")] <- "BF-B"
```

Downloading Stock Time Series From *Tiingo*

Yahoo data quality deteriorated significantly in 2017, and *Google* data quality is also poor, leaving *Tiingo*, *Alpha Vantage*, and *Quandl* as the only major providers of free daily *OHLC* stock prices.

But *Quandl* doesn't provide free *ETF* prices, while *Tiingo* does.

The function `getSymbols()` has a *method* for downloading time series data from *Tiingo*, called `getSymbols.tiingo()`.

Users must first obtain a *Tiingo* API key, and then pass it in `getSymbols.tiingo()` calls:

<https://www.tiingo.com/>

Note that the data are downloaded as *xts* time series, with a date-time index of class *POSIXct* (not *Date*).

```
> # Load package rutils
> library(rutils)
> # Create new environment for data
> sp500env <- new.env()
> # Boolean vector of symbols already downloaded
> isdownloaded <- symbolv %in% ls(sp500env)
> # Download in while loop from Tiingo and copy into environment
> n attempts <- 0 # Number of download attempts
> while ((sum(!isdownloaded) > 0) & (n attempts < 3)) {
+   # Download data and copy it into environment
+   n attempts <- n attempts + 1
+   cat("Download attempt = ", n attempts, "\n")
+   for (symbol in symbolv[!isdownloaded]) {
+     cat("processing: ", symbol, "\n")
+     tryCatch( # With error handler
+       quantmod::getSymbols(symbol, src="tiingo", adjust=TRUE, auto.assign=
+         from="1990-01-01", env=sp500env, api.key="j84ac2b9c5bd
+     # Error handler captures error condition
+     error=function(error_cond) {
+       print(paste("error handler: ", error_cond))
+     }, # end error handler
+     finally=print(paste("symbol=", symbol))
+   ) # end tryCatch
+ } # end for
+ # Update vector of symbols already downloaded
+ isdownloaded <- symbolv %in% ls(sp500env)
+ Sys.sleep(2) # Wait 2 seconds until next attempt
+ } # end while
> class(sp500env$AAPL)
> class(zoo::index(sp500env$AAPL))
> tail(sp500env$AAPL)
> symbolv[!isdownloaded]
```

Coercing Date-time Indices

The date-time indices of the *OHLC* stock prices are in the `POSIXct` format suitable for intraday pricev, not daily prices.

The function `as.Date()` coerces `POSIXct` objects into `Date` objects.

The function `get()` retrieves objects that are referenced using character strings, instead of their names.

The function `assign()` assigns a value to an object in a specified *environment*, by referencing it using a character string (name).

The functions `get()` and `assign()` allow retrieving and assigning values to objects that are referenced using character strings.

```
> # The date-time index of AAPL is POSIXct
> class(zoo::index(sp500env$AAPL))
> # Coerce the date-time index of AAPL to Date
> zoo::index(sp500env$AAPL) <- as.Date(zoo::index(sp500env$AAPL))
> # Coerce all the date-time indices to Date
> for (symbol in ls(sp500env)) {
+   ohlc <- get(symbol, envir=sp500env)
+   zoo::index(ohlc) <- as.Date(zoo::index(ohlc))
+   assign(symbol, ohlc, envir=sp500env)
+ } # end for
```


Managing Exceptions in Stock Symbols

The column names for symbol "LOW" (Lowe's company) must be renamed for the extractor function `quantmod::Lo()` to work properly.

Tickers which contain a dot in their name (like "BRK.B") are not valid symbols in R, so they must be downloaded separately and renamed.

```
> # "LOW.Low" is a bad column name
> colnames(sp500env$LOW)
> strsplit(colnames(sp500env$LOW), split=".")
> do.call(cbind, strsplit(colnames(sp500env$LOW), split="."))
> do.call(cbind, strsplit(colnames(sp500env$LOW), split="."))[2, ]
> # Extract proper names from column names
> namesv <- rutils::get_name(colnames(sp500env$LOW), field=2)
> # Or
> # namesv <- do.call(rbind, strsplit(colnames(sp500env$LOW),
> #                                     split="."))[, 2]
> # Rename "LOW" colnames to "LOWES"
> colnames(sp500env$LOW) <- paste("LOWES", namesv, sep=".")
> sp500env$LOWES <- sp500env$LOW
> rm(LOW, envir=sp500env)
> # Rename BF-B colnames to "BFB"
> colnames(sp500env$BF-B) <- paste("BFB", namesv, sep=".")
> sp500env$BFB <- sp500env$BF-B
> rm("BF-B", envir=sp500env)
> # Rename BRK-B colnames
> sp500env$BRKB <- sp500env$BRK-B
> rm("BRK-B", envir=sp500env)
> colnames(sp500env$BRKB) <- gsub("BRK-B", "BRKB", colnames(sp500e
> # Save OHLC prices to .RData file
> save(sp500env, file="/Users/jerzy/Develop/lecture_slides/data/sp500.RData")
> # Download "BRK.B" separately with auto.assign=FALSE
> # BRKB <- quantmod::getSymbols("BRK-B", auto.assign=FALSE, src="tiingo", adjust=TRUE, from="1990-01-01", api.key="j84ac2b9c5bde2d68e3
> # colnames(BRKB) <- paste("BRKB", namesv, sep=".")
> # sp500env$BRKB <- BRKB
```



```
> # Plot OHLC candlestick chart for LOWES
> chart_Series(x=sp500env$LOWES["2019-12-/",
+   TA="add_Vo()", name="LOWES OHLC Stock Prices")
> # Plot dygraph
> dygraphs::dygraph(sp500env$LOWES["2019-12-/", -5], main="LOWES OHLC
+   dyCandlestick()
```

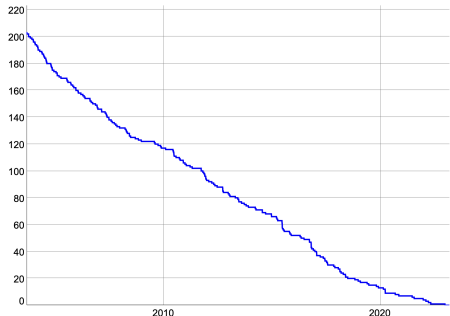
S&P500 Stock Index Constituent Prices

The file `sp500.RData` contains the *environment* `sp500.env` with *OHLC* prices and trading volumes of S&P500 stock index constituents.

The S&P500 stock index constituent data is of poor quality before 2000, so we'll mostly use the data after the year 2000.

```
> # Load S&P500 constituent stock prices
> load("/Users/jerzy/Develop/lecture_slides/data/sp500.RData")
> pricev <- eapply(sp500env, quantmod::C1)
> pricev <- rutils::do_call(cbind, pricev)
> # Carry forward non-NA prices
> pricev <- zoo::na.locf(pricev, na.rm=FALSE)
> # Drop ".Close" from column names
> colnames(pricev[, 1:4])
> colnames(pricev) <- rutils::get_name(colnames(pricev))
> # Or
> # colnames(pricev) <- do.call(rbind,
> #   strsplit(colnames(pricev), split=".[.]"))[, 1]
> # Calculate percentage returns of the S&P500 constituent stocks
> # retp <- xts::diff.xts(log(pricev))
> retp <- xts::diff.xts(pricev)/
+   rutils::lagit(pricev, pad_zeros=FALSE)
> set.seed(1121)
> samplev <- sample(NCOL(retp), s=100, replace=FALSE)
> prices100 <- pricev[, samplev]
> returns100 <- retp[, samplev]
> save(pricev, prices100,
+   file="/Users/jerzy/Develop/lecture_slides/data/sp500_prices.RData")
> save(retp, returns100,
+   file="/Users/jerzy/Develop/lecture_slides/data/sp500_returns.RData")
```

Number of S&P500 Constituents Without Prices

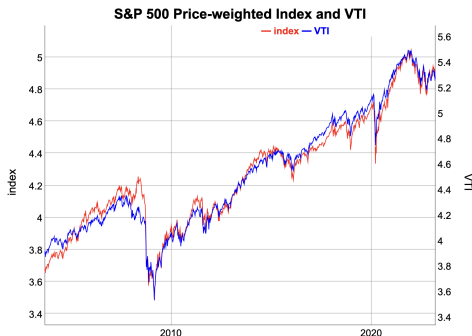


```
> # Calculate number of constituents without prices
> datav <- rowSums(is.na(pricev))
> datav <- xts::xts(datav, order.by=zoo::index(pricev))
> dygraphs::dygraph(datav, main="Number of S&P500 Constituents Witho
+   dyOptions(colors="blue", strokeWidth=2)
```

S&P500 Stock Portfolio Index

The price-weighted index of *S&P500* constituents closely follows the *VTI ETF*.

```
> # Calculate price weighted index of constituent
> ncols <- NCOL(pricev)
> pricev <- zoo::na.locf(pricev, fromLast=TRUE)
> indeks <- xts(rowSums(pricev)/ncols, zoo::index(pricev))
> colnames(indeks) <- "index"
> # Combine index with VTI
> datav <- cbind(indeks[zoo::index(etfenv$VTI)], etfenv$VTI[, 4])
> colnamev <- c("index", "VTI")
> colnames(datav) <- colnamev
> # Plot index with VTI
> endd <- rutils::calc_endpoints(datav, interval="weeks")
> dygraphs::dygraph(log(datav)[endd],
+   main="S&P 500 Price-weighted Index and VTI") %>%
+   dyAxis("y", label=colnamev[1], independentTicks=TRUE) %>%
+   dyAxis("y2", label=colnamev[2], independentTicks=TRUE) %>%
+   dySeries(name=colnamev[1], axis="y", col="red") %>%
+   dySeries(name=colnamev[2], axis="y2", col="blue")
```



Writing Time Series To Files

The data from *Tiingo* is downloaded as *xts* time series, with a date-time index of class *POSIXct* (not *Date*).

The function `save()` writes objects to compressed binary *.RData* files.

The easiest way to share data between R and Excel is through *.csv* files.

The package *zoo* contains functions `write.zoo()` and `read.zoo()` for writing and reading *zoo* time series from *.txt* and *.csv* files.

The function `data.table::fread()` reads from *.csv* files over 6 times faster than the function `read.csv()`!

The function `data.table::fwrite()` writes to *.csv* files over 12 times faster than the function `write.csv()`, and 278 times faster than function `cat()`!

```
> # Save the environment to compressed .RData file
> dir_name <- "/Users/jerzy/Develop/lecture_slides/data/"
> save(sp500env, file=paste0(dir_name, "sp500.RData"))
> # Save the ETF prices into CSV files
> dir_name <- "/Users/jerzy/Develop/lecture_slides/data/SP500/"
> for (symbol in ls(sp500env)) {
+   zoo::write.zoo(sp500env$symbol, file=paste0(dir_name, symbol, ".csv"))
+ } # end for
> # Or using lapply()
> file_names <- lapply(ls(sp500env), function(symbol) {
+   xtsv <- get(symbol, envir=sp500env)
+   zoo::write.zoo(xtsv, file=paste0(dir_name, symbol, ".csv"))
+   symbol
+ }) # end lapply
> unlist(file_names)
> # Or using eapply() and data.table::fwrite()
> file_names <- eapply(sp500env, function(xtsv) {
+   file_name <- rutils::get_name(colnames(xtsv)[1])
+   data.table::fwrite(data.table::as.data.table(xtsv), file=paste0(
+     dir_name,
+     file_name
+   )) # end eapply
+ })
> unlist(file_names)
```

Reading Time Series from Files

The function `load()` reads data from `.RData` files, and *invisibly* returns a vector of names of objects created in the workspace.

The function `Sys.glob()` lists files matching names obtained from wildcard expansion.

The easiest way to share data between R and Excel is through `.csv` files.

The function `as.Date()` parses character strings, and coerces numeric and `POSIXct` objects into `Date` objects.

The function `data.table::setDF()` coerces a *data table* object into a *data frame* using a *side effect*, without making copies of data.

The function `data.table::fread()` reads from `.csv` files over 6 times faster than the function `read.csv()`!

```
> # Load the environment from compressed .RData file
> dir_name <- "/Users/jerzy/Develop/lecture_slides/data/"
> load(file=paste0(dir_name, "sp500.RData"))
> # Get all the .csv file names in the directory
> dir_name <- "/Users/jerzy/Develop/lecture_slides/data/SP500/"
> file_names <- Sys.glob(paste0(dir_name, "*.csv"))
> # Create new environment for data
> sp500env <- new.env()
> for (file_name in file_names) {
+   xtsv <- xts::as.xts(zoo::read.csv.zoo(file_name))
+   symbol <- rutils::get_name(colnames(xtsv)[1])
+   # symbol <- strsplit(colnames(xtsv), split=".")[[1]][1]
+   assign(symbol, xtsv, envir=sp500env)
+ } # end for
> # Or using fread()
> for (file_name in file_names) {
+   xtsv <- data.table::fread(file_name)
+   data.table::setDF(xtsv)
+   xtsv <- xts::xts(xtsv[, -1], as.Date(xtsv[, 1]))
+   symbol <- rutils::get_name(colnames(xtsv)[1])
+   assign(symbol, xtsv, envir=sp500env)
+ } # end for
```

Downloading Stock Time Series From *Alpha Vantage*

Yahoo data quality deteriorated significantly in 2017, and *Google* data quality is also poor, leaving *Tiingo*, *Alpha Vantage*, and *Quandl* as the only major providers of free daily *OHLC* stock prices.

But *Quandl* doesn't provide free *ETF* prices, while *Alpha Vantage* does.

The function `getSymbols()` has a *method* for downloading time series data from *Alpha Vantage*, called `getSymbols.av()`.

Users must first obtain an *Alpha Vantage* API key, and then pass it in `getSymbols.av()` calls:

<https://www.alphavantage.co/>

The function `adjustOHLC()` with argument `use.Adjusted=TRUE`, adjusts all the *OHLC* price columns, using the *Adjusted* price column.

```
> # Remove all files from environment(if necessary)
> rm(list=ls(sp500env), envir=sp500env)
> # Download in while loop from Alpha Vantage and copy into environment
> isdownloaded <- symbolv %in% ls(sp500env)
> n attempts <- 0
> while ((sum(!isdownloaded) > 0) & (n attempts < 10)) {
+   # Download data and copy it into environment
+   n attempts <- n attempts + 1
+   for (symbol in symbolv[!isdownloaded]) {
+     cat("processing: ", symbol, "\n")
+     tryCatch( # With error handler
+       quantmod::getSymbols(symbol, src="av", adjust=TRUE, auto.assign=TRUE,
+         output.size="full", api.key="T7JPW54ES8G75310"),
+     # error handler captures error condition
+     error=function(error_cond) {
+       print(paste("error handler: ", error_cond))
+     }, # end error handler
+     finally=print(paste("symbol=", symbol))
+   ) # end tryCatch
+ } # end for
+ # Update vector of symbols already downloaded
+ isdownloaded <- symbolv %in% ls(sp500env)
+ Sys.sleep(2) # Wait 2 seconds until next attempt
+ } # end while
> # Adjust all OHLC prices in environment
> for (symbol in ls(sp500env)) {
+   assign(symbol,
+     adjustOHLC(get(x=symbol, envir=sp500env), use.Adjusted=TRUE),
+     envir=sp500env)
+ } # end for
```

Downloading The *S&P500* Index Time Series From Yahoo

The *S&P500* stock market index is a capitalization-weighted average of the 500 largest U.S. companies, and covers about 80% of the U.S. stock market capitalization.

Yahoo provides daily *OHLC* prices for the *S&P500* index (symbol *^GSPC*), and for the *S&P500* total return index (symbol *^SP500TR*).

But special characters in some stock symbols, like "-" or "^" are not allowed in R names.

For example, the symbol *^GSPC* for the *S&P500* stock market index isn't a valid name in R.

The function `setSymbolLookup()` creates valid names corresponding to stock symbols, which are then used by the function `getSymbols()` to create objects with the valid names.

Yahoo data quality deteriorated significantly in 2017, and Google data quality is also poor, leaving *Alpha Vantage* and *Quandl* as the only major providers of free daily *OHLC* stock prices.

```
> # Assign name SP500 to ^GSPC symbol
> setSymbolLookup(SP500=list(name="^GSPC", src="yahoo"))
> getSymbolLookup()
> # view and clear options
> options("getSymbols.sources")
> options(getSymbols.sources=NULL)
> # Download S&P500 prices into etfenv
> quantmod::getSymbols("SP500", env=etfenv,
+   adjust=TRUE, auto.assign=TRUE, from="1990-01-01")
> chart_Series(x=etfenv$SP500["2016/"],
+   TA="add_Vo()", name="S&P500 index")
```

Downloading The *DJIA* Index Time Series From Yahoo

The Dow Jones Industrial Average (*DJIA*) stock market index is a price-weighted average of the 30 largest U.S. companies (same number of shares per company).

Yahoo provides daily *OHLC* prices for the *DJIA* index (symbol *^DJI*), and for the *DJITR* total return index (symbol *DJITR*).

But special characters in some stock symbols, like "-" or "^" are not allowed in R names.

For example, the symbol *^DJI* for the *DJIA* stock market index isn't a valid name in R.

The function `setSymbolLookup()` creates valid names corresponding to stock symbols, which are then used by the function `getSymbols()` to create objects with the valid names.

```
> # Assign name DJIA to ^DJI symbol
> setSymbolLookup(DJIA=list(name="^DJI", src="yahoo"))
> getSymbolLookup()
> # view and clear options
> options("getSymbols.sources")
> options(getSymbols.sources=NULL)
> # Download DJIA prices into etfenv
> quantmod::getSymbols("DJIA", env=etfenv,
+   adjust=TRUE, auto.assign=TRUE, from="1990-01-01")
> chart_Series(x=etfenv$DJIA["2016/"],
+   TA="add_Vo()", name="DJIA index")
```


Calculating Prices and Returns From *OHLC* Data

The function `na.locf()` from package *zoo* replaces NA values with the most recent non-NA values prior to it.

The function `na.locf()` with argument `fromLast=TRUE` replaces NA values with non-NA values in reverse order, starting from the end.

The function `rutils::get_name()` extracts symbol names (tickers) from a vector of character strings.

```
> pricev <- eapply(sp500env, quantmod::C1)
> pricev <- rutils::do_call(cbind, pricev)
> # Carry forward non-NA prices
> pricev <- zoo::na.locf(pricev, na.rm=FALSE)
> # Get first column name
> colnames(pricev[, 1])
> rutils::get_name(colnames(pricev[, 1]))
> # Modify column names
> colnames(pricev) <- rutils::get_name(colnames(pricev))
> # Or
> # colnames(pricev) <- do.call(rbind,
> #   strsplit(colnames(pricev), split=".[.]"))[, 1]
> # Calculate percentage returns
> retp <- xts::diff.xts(pricev)/
+   rutils::lagit(pricev, pad_zeros=FALSE)
> # Select a random sample of 100 prices and returns
> set.seed(1121)
> samplev <- sample(NCOL(retp), s=100, replace=FALSE)
> prices100 <- pricev[, samplev]
> returns100 <- retp[, samplev]
> # Save the data into binary files
> save(pricev, prices100,
+   file="/Users/jerzy/Develop/lecture_slides/data/sp500_prices.R")
> save(retp, returns100,
+   file="/Users/jerzy/Develop/lecture_slides/data/sp500_returns.R")
```

Downloading Stock Prices From Polygon

Polygon is a premium provider of live and historical stock price data, both daily and intraday (minutes).

Polygon provides 2 years of daily historical stock prices for free. But users must first obtain a *Polygon API key*.

Polygon provides the historical *OHLC* stock prices in *JSON* format.

JSON (JavaScript Object Notation) is a data format consisting of symbol-value pairs.

The package *jsonlite* contains functions for managing data in *JSON* format.

The functions `fromJSON()` and `toJSON()` convert data from *JSON* format to R objects, and vice versa.

The functions `read_json()` and `write_json()` read and write *JSON* format data in files.

The function `download.file()` downloads data from an internet website URL and writes it to a file.

```
> # Setup code
> symbol <- "SPY"
> startd <- as.Date("1990-01-01")
> todayd <- Sys.Date()
> tspan <- "day"
> # Replace below your own Polygon API key
> apikey <- "SEpnsBpiRyQNMJdl48r6d0o0_pjmCu5r"
> # Create url for download
> url1 <- paste0("https://api.polygon.io/v2/aggs/ticker/", symbol,
"> # Download SPY OHLC prices in JSON format from Polygon
> ohlc <- jsonlite::read_json(url1)
> class(ohlc)
> NROW(ohlc)
> names(ohlc)
> # Extract list of prices from json object
> ohlc <- ohlc$results
> # Coerce from list to matrix
> ohlc <- lapply(ohlc, unlist)
> ohlc <- do.call(rbind, ohlc)
> # Coerce time from milliseconds to dates
> dates <- ohlc[, "t"]/1e3
> dates <- as.POSIXct(dates, origin="1970-01-01")
> dates <- as.Date(dates)
> tail(dates)
> # Coerce from matrix to xts
> ohlc <- ohlc[, c("o","h","l","c","v","vw")]
> colnames(ohlc) <- c("Open", "High", "Low", "Close", "Volume", "VW
> ohlc <- xts::xts(ohlc, order.by=dates)
> tail(ohlc)
> # Save the xts time series to compressed RData file
> save(ohlc, file="/Users/jerzy/Data/spy_daily.RData")
> # Candlestick plot of SPY OHLC prices
> dygraphs::dygraph(ohlc[, 1:4], main=paste("Candlestick Plot of", s
+ dygraphs::dyCandlestick()
```

Downloading Multiple Stock Prices From Polygon

The stock prices for multiple stocks can be downloaded in a while() loop.

```
> # Select ETF symbols for asset allocation
```

```
> symbolv <- c("VTI", "VEU", "EEM", "XLY", "XLP", "XLI",
+ "XLV", "XLI", "XLB", "XLK", "XLU", "VYM", "IVW", "IWF",
+ "IEF", "TLT", "VNQ", "DBC", "GLD", "USO", "MTUM", "IVE",
+ "VLUE", "QUAL", "VTV", "USMV", "AIEI")
> # Setup code
> etfenv <- new.env() # New environment for data
> # Boolean vector of symbols already downloaded
> isdownloaded <- symbolv %in% ls(etfenv)
```

```
> # Download data from Polygon using while loop
> while (sum(!isdownloaded) > 0) {
+   for (symbol in symbolv[!isdownloaded]) {
+     cat("Processing:", symbol, "\n")
+     tryCatch({ # With error handler
+       # Download OHLC bars from Polygon into JSON format file
+       url1 <- paste0("https://api.polygon.io/v2/aggs/ticker/", symbol, "/range/1/", t
+       ohlc <- jsonlite::read_json(url1)
+       # Extract list of prices from json object
+       ohlc <- ohlc$results
+       # Coerce from list to matrix
+       ohlc <- lapply(ohlc, unlist)
+       ohlc <- do.call(rbind, ohlc)
+       # Coerce time from milliseconds to dates
+       dates <- ohlc[, "t"]/1e3
+       dates <- as.POSIXct(dates, origin="1970-01-01")
+       dates <- as.Date(dates)
+       # Coerce from matrix to xts
+       ohlc <- ohlc[, c("o", "h", "l", "c", "v", "vw")]
+       colnames(ohlc) <- paste0(symbol, ".", c("Open", "High", "Low", "Close", "Volume
+       ohlc <- xts::xts(ohlc, order.by=dates)
+       # Save to environment
+       assign(symbol, ohlc, envir=etfenv)
+       Sys.sleep(1)
+     },
+     error={function(error_cond) print(paste("Error handler:", error_cond))},
+     finally=print(paste0("symbol=", symbol))
+   ) # end tryCatch
+ } # end for
+ # Update vector of symbols already downloaded
+ isdownloaded <- symbolv %in% ls(etfenv)
+ } # end while
> save(etfenv, file="/Users/jerzy/Develop/lecture_slides/data/etf_data.RData")
```

Calculating the Stock Alphas, Betas, and Other Performance Statistics

The package *PerformanceAnalytics* contains functions for calculating risk and performance statistics, such as the *variance*, *skewness*, *kurtosis*, *beta*, *alpha*, etc.

The function `PerformanceAnalytics::table.CAPM()` calculates the *beta* β and *alpha* α values, the *Treynor* ratio, and other performance statistics.

The function `PerformanceAnalytics::table.Stats()` calculates a data frame of risk and return statistics of the return distributions.

```
> pricev <- eapply(etfenv, quantmod::C1)
> pricev <- do.call(cbind, pricev)
> # Drop ".Close" from colnames
> colnames(pricev) <- do.call(rbind, strsplit(colnames(pricev), split="."))
> # Calculate the log returns
> retp <- xts::diff.xts(log(pricev))
> # Copy prices and returns into etfenv
> etfenv$pricev <- pricev
> etfenv$retp <- retp
> # Copy symbolv into etfenv
> etfenv$symbolv <- symbolv
> # Calculate the risk-return statistics
> riskstats <- PerformanceAnalytics::table.Stats(retp)
> # Transpose the data frame
> riskstats <- as.data.frame(t(riskstats))
> # Add Name column
> riskstats$Name <- rownames(riskstats)
> # Copy riskstats into etfenv
> etfenv$riskstats <- riskstats
> # Calculate the beta, alpha, Treynor ratio, and other performance statistics
> capmstats <- PerformanceAnalytics::table.CAPM(Ra=retp[, symbolv], Rb=retp[, "VTI"], scale=1)
> colnamev <- strsplit(colnames(capmstats), split=" ")[,1]
> colnamev <- do.call(cbind, colnamev)[1, ]
> colnames(capmstats) <- colnamev
> capmstats <- t(capmstats)
> capmstats <- capmstats[, -1]
> colnamev <- colnames(capmstats)
> whichv <- match(c("Annualized Alpha", "Information Ratio", "Treynor Ratio", "Beta", "Alpha", "Information", "Treynor"), colnamev[whichv])
> colnames(capmstats) <- colnamev
> capmstats <- capmstats[order(capmstats[, "Alpha"], decreasing=TRUE), ]
> # Copy capmstats into etfenv
> etfenv$capmstats <- capmstats
> save(etfenv, file="/Users/jerzy/Develop/lecture_slides/data/etf_data.Rsave")
```

Scraping S&P500 Stock Index Constituents From Websites

The *S&P500* index constituents change over time, and *Standard & Poor's* replaces companies that have decreased in capitalization with ones that have increased.

The *S&P500* index may contain more than 500 stocks because some companies have several share classes of stock.

The *S&P500* index constituents may be scraped from websites like [Wikipedia](#), using dedicated packages.

The function `getURL()` from package *RCurl* downloads the *html* text data from an internet website URL.

The function `readHTMLTable()` from package *XML* extracts tables from *html* text data or from a remote URL, and returns them as a list of *data frames* or matrices.

`readHTMLTable()` can't parse secure URLs, so they must first be downloaded using function `getURL()`, and then parsed using `readHTMLTable()`.

```
> library(RCurl) # Load package RCurl
> library(XML) # Load package XML
> # Download text data from URL
> sp500 <- getURL(
+ "https://en.wikipedia.org/wiki/List_of_S%26P500_companies")
> # Extract tables from the text data
> sp500 <- readHTMLTable(sp500)
> str(sp500)
> # Extract colnames of data frames
> lapply(sp500, colnames)
> # Extract S&P500 constituents
> sp500 <- sp500[[1]]
> head(sp500)
> # Create valid R names from symbols containing "-" or "." character
> sp500$namesv <- gsub("-", "_", sp500$Ticker)
> sp500$namesv <- gsub("[.]", "_", sp500$names)
> # Write data frame of S&P500 constituents to CSV file
> write.csv(sp500,
+ file="/Users/jerzy/Develop/lecture_slides/data/sp500_Yahoo.csv",
+ row.names=FALSE)
```

Downloading S&P500 Time Series Data From Yahoo

Before time series data for the *S&P500* index constituents can be downloaded from *Yahoo*, it's necessary to create valid names corresponding to symbols containing special characters like "-".

The function `setSymbolLookup()` creates a lookup table for *Yahoo* symbols, using valid names in R.

For example *Yahoo* uses the symbol "BRK-B", which isn't a valid name in R, but can be mapped to "BRK_B", using the function `setSymbolLookup()`.

```
> library(rutils) # Load package rutils
> # Load data frame of S&P500 constituents from CSV file
> sp500 <- read.csv(file="/Users/jerzy/Develop/lecture_slides/data/sp500.csv")
> # Register symbols corresponding to R names
> for (indeks in 1:NROW(sp500)) {
+   cat("processing: ", sp500$Ticker[indeks], "\n")
+   setSymbolLookup(structure(
+     list(list(name=sp500$Ticker[indeks])),
+     names=sp500$names[indeks]))
+ } # end for
> sp500env <- new.env() # new environment for data
> # Remove all files (if necessary)
> rm(list=ls(sp500env), envir=sp500env)
> # Download data and copy it into environment
> rutils::get_data(sp500$names,
+   env_out=sp500env, startd="1990-01-01")
> # Or download in loop
> for (symbol in sp500$names) {
+   cat("processing: ", symbol, "\n")
+   rutils::get_data(symbol,
+     env_out=sp500env, startd="1990-01-01")
+ } # end for
> save(sp500env, file="/Users/jerzy/Develop/lecture_slides/data/sp500env.Rsave")
> chart_Series(x=sp500env$BRKB["2016/"],
+   TA="add_Vo()", name="BRK-B stock")
```

Downloading *FRED* Time Series Data

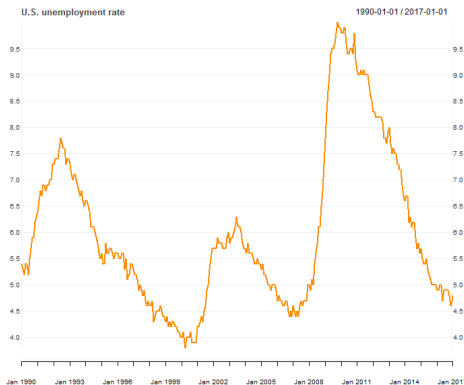
FRED is a database of economic time series maintained by the Federal Reserve Bank of St. Louis:

<http://research.stlouisfed.org/fred2/>

The function `getSymbols()` downloads time series data into the specified *environment*.

`getSymbols()` can download *FRED* data with the argument `"src"` set to *FRED*.

If the argument `"auto.assign"` is set to `FALSE`, then `getSymbols()` returns the data, instead of assigning it silently.



```
> # Download U.S. unemployment rate data
> unemp_rate <- quantmod::getSymbols("UNRATE",
+                                   auto.assign=FALSE, src="FRED")
> # Plot U.S. unemployment rate data
> chart_Series(unemp_rate["1990/"],
+              name="U.S. unemployment rate")
```

Homework Assignment

Required

- Study all the lecture slides in `FRE7241_Lecture_1.pdf`, and run all the code in `FRE7241_Lecture_1.R`,

Recommended

- Read the documentation for packages `rutils.pdf` and `HighFreq.pdf`,