

FRE7241 Algorithmic Portfolio Management

Lecture#1, Fall 2023

Jerzy Pawlowski jp3900@nyu.edu

NYU Tandon School of Engineering

September 5, 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.

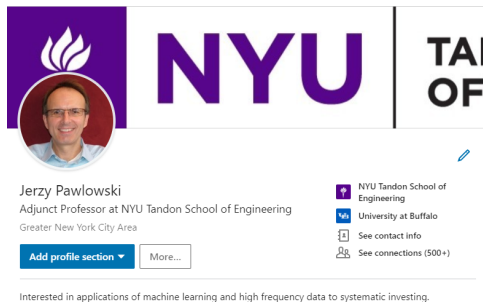
I'm interested in applications of *machine learning* to *systematic investing*.

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](#)

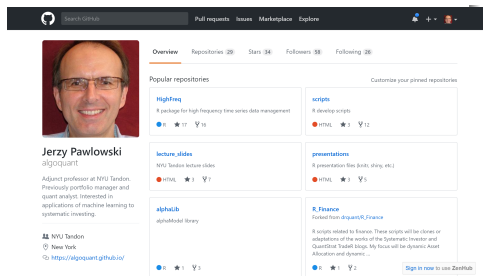


Jerzy Pawlowski
Adjunct Professor at NYU Tandon School of Engineering
Greater New York City Area

[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	Language
HighFreq A package for high-frequency time series data management	17	Python
lecture_slides NYU Tandon lecture slides	7	HTML
alphanlib alphanlib library	3	Python
scripts A develop scripts	12	HTML
presentations A presentation files (pdfs, shps, etc.)	5	HTML
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 ...	2	Python

[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 be announced several days in advance.

The homeworks and 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 be announced several days in advance.

The homeworks and 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 Raunak Bhupal rb4986@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

Homeworks and tests will be graded and assigned numerical scores. Each part of homeworks and 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

Homeworks and tests will be graded and assigned numerical scores. Each part of homeworks and 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

Homeworks and tests will be graded and assigned numerical scores. Each part of homeworks and 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 - Machine learning techniques applied to trading and portfolio management.
- *Systematic Trading* by Robert Carver - Practical trading knowledge by an experienced portfolio manager.
- *Systematic Trading* by Robert Carver - Practical investment knowledge by a successful investor.
- *Quantitative Trading* by Xin Guo, Tze Leung Lai, Howard Shek, Samuel Po-Shing Wong - Advanced topics in quantitative trading by academic experts.
- *Financial Data and Models Using R* by Clifford Ang - Good introduction to time series, portfolio optimization, and performance measures.
- *Automated Trading* by Chris Conlan - How to implement practical computer trading systems.
- *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>

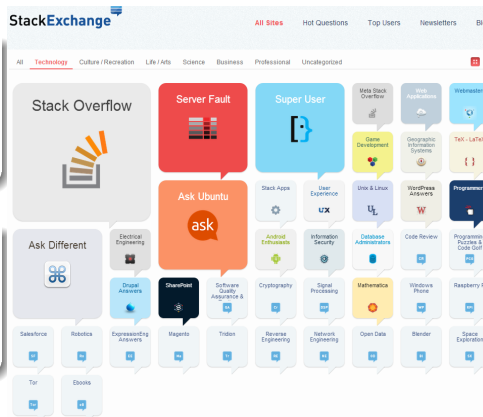
Stack Exchange

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 Eddelbuettel

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

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

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

<http://dirk.eddelbuettel.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 (algoquant). The profile includes a bio: "Adjunct professor at NYU Tandon. Previously portfolio manager and quant analyst. Interested in applications of machine learning to systematic investing." and a location of "New York". Below the profile, there are several repositories listed:

- HighFreq**: R package for high-frequency time series data management. 17 stars, 15 forks.
- scripts**: R develop scripts. 3 stars, 12 forks.
- lecture_slides**: NYU Tandon lecture slides. 3 stars, 1 fork.
- presentations**: R presentation files (pdf, shap, etc.). 3 stars, 5 forks.
- alphatub**: alphatub library. 1 star, 3 forks.
- R_Finance**: R scripts related to finance. These scripts will be clones or adaptations of the works of the Systematic Investor and Quantitative Trading blogs. My focus will be dynamic Asset Allocation and dynamic... 1 star, 2 forks.

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.

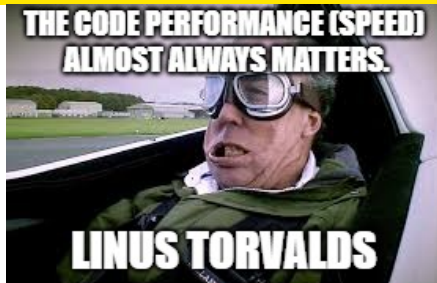


This course is designed to teach the most useful elements of R for financial analysis, through case studies and examples,

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 outputs of the two methods
> all.equal(cumsumv, cumsumv2)
> # Microbenchmark the two methods
> library(microbenchmark)
> summary(microbenchmark(
+   cumsum=cumsum(vectorv), # Vectorized
+   loop_alloc={cumsumv2 <- vectorv # Allocate memory to cumsumv3
+     for (i in 2:NROW(vectorv))
+       cumsumv2[i] <- (vectorv[i] + cumsumv2[i-1])
+   },
+   loop_nalloc={cumsumv3 <- vectorv[1] # Doesn't allocate memory to
+     for (i in 2:NROW(vectorv))
+       cumsumv3[i] <- (vectorv[i] + cumsumv3[i-1])
```

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 environment with the following components:

- Source Editor:** Contains an R script with code for data manipulation and portfolio optimization. The code includes comments and uses functions like `tail`, `update.alphaModel`, `recalc.alphaModel`, `as.vector`, `write.csv`, `read.csv`, `xts`, `diff`, and `library`.
- Console:** Shows output from the `install.packages()` function, including warnings about internet connectivity and the successful installation of the `PerformanceAnalytics` package from R-Forge.
- Workspace/History Pane:** Displays the `?MASS` help page for the `MASS` package, showing its description, installation details, and usage instructions.

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] "2023-09-09 14:09:31 EDT"
>
> Sys.Date() # Get date only
[1] "2023-09-09"
```

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 4.3.0 (2023-04-21)
Platform: aarch64-apple-darwin20 (64-bit)
Running under: macOS Ventura 13.3.1

Matrix products: default
BLAS:   /Library/Frameworks/R.framework/Versions/4.3-arm64/Resources
LAPACK: /Library/Frameworks/R.framework/Versions/4.3-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

time zone: America/New_York
tzcode source: internal

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

other attached packages:
[1] knitr_1.42      HighFreq_0.1    rutils_0.2      dygraphs_1.1
[5] quantmod_0.4.22 TTR_0.24.3      xts_0.13.1      zoo_1.8-12

loaded via a namespace (and not attached):
 [1] digest_0.6.31    fastmap_1.1.1    xfun_0.39        lattice_0.20-45
 [5] magrittr_2.0.3    htmltools_0.5.5  cli_3.6.1         grid_4.3-0
 [9] compiler_4.3.0    highr_0.10       tools_4.3.0      rstudioapi_0.14
[13] curl_5.0.0        evaluate_0.20    Rcpp_1.0.10      rlang_1.1.1
[17] htmlwidgets_1.6.2
```

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(pricev)
> do.call(rbind, strsplit(colnames(pricev), split=".[.]"))[, 1]
> colnames(pricev) <- do.call(rbind, strsplit(colnames(pricev), split=".[.]"))[, 1]
> # Or
> colnames(pricev) <- unname(sapply(colnames(pricev),
+   function(colname) strsplit(colname, split=".[.]")[[1]][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$etf_list <- etf_list
> # Or
> assign("prices", pricev, 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 eval(expr, envir, enclos): 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), needs compilation (no), citation (AER citation info), materials (NEWS), in views (Econometrics, Survival, TimeSeries), and CRAN checks (AER results). There is a 'Downloads:' section with links for 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). At the bottom, it shows 'Reverse dependencies:' with reverse depends (lpack, rdd) and reverse suggests (censReg, glmx, lme4, micEconCES, mlogit, plm, REEMtree, 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](#), [glmx](#), [lme4](#), [micEconCES](#), [mlogit](#), [plm](#), [REEMtree](#), [sandwich](#)

CRAN Task Views

CRAN Finance Task View

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

Note:

- Maintainer,
- Topics,
- List of packages.

← → ↻ 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 [CRAN Task Views](#) for [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 [nlme](#) package). Many other suitable methods are available in 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 [havesGARCH](#) package provides a variety of 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 [series](#) 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 a variety of 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](#), [waveletComp2](#), [waveletComp3](#), [waveletComp4](#), [waveletComp5](#), [waveletComp6](#), [waveletComp7](#), [waveletComp8](#), [waveletComp9](#), [waveletComp10](#), [waveletComp11](#), [waveletComp12](#), [waveletComp13](#), [waveletComp14](#), [waveletComp15](#), [waveletComp16](#), [waveletComp17](#), [waveletComp18](#), [waveletComp19](#), [waveletComp20](#), [waveletComp21](#), [waveletComp22](#), [waveletComp23](#), [waveletComp24](#), [waveletComp25](#), [waveletComp26](#), [waveletComp27](#), [waveletComp28](#), [waveletComp29](#), [waveletComp30](#), [waveletComp31](#), [waveletComp32](#), [waveletComp33](#), [waveletComp34](#), [waveletComp35](#), [waveletComp36](#), [waveletComp37](#), [waveletComp38](#), [waveletComp39](#), [waveletComp40](#), [waveletComp41](#), [waveletComp42](#), [waveletComp43](#), [waveletComp44](#), [waveletComp45](#), [waveletComp46](#), [waveletComp47](#), [waveletComp48](#), [waveletComp49](#), [waveletComp50](#), [waveletComp51](#), [waveletComp52](#), [waveletComp53](#), [waveletComp54](#), [waveletComp55](#), [waveletComp56](#), [waveletComp57](#), [waveletComp58](#), [waveletComp59](#), [waveletComp60](#), [waveletComp61](#), [waveletComp62](#), [waveletComp63](#), [waveletComp64](#), [waveletComp65](#), [waveletComp66](#), [waveletComp67](#), [waveletComp68](#), [waveletComp69](#), [waveletComp70](#), [waveletComp71](#), [waveletComp72](#), [waveletComp73](#), [waveletComp74](#), [waveletComp75](#), [waveletComp76](#), [waveletComp77](#), [waveletComp78](#), [waveletComp79](#), [waveletComp80](#), [waveletComp81](#), [waveletComp82](#), [waveletComp83](#), [waveletComp84](#), [waveletComp85](#), [waveletComp86](#), [waveletComp87](#), [waveletComp88](#), [waveletComp89](#), [waveletComp90](#), [waveletComp91](#), [waveletComp92](#), [waveletComp93](#), [waveletComp94](#), [waveletComp95](#), [waveletComp96](#), [waveletComp97](#), [waveletComp98](#), [waveletComp99](#), [waveletComp100](#), [waveletComp101](#), [waveletComp102](#), [waveletComp103](#), [waveletComp104](#), [waveletComp105](#), [waveletComp106](#), [waveletComp107](#), [waveletComp108](#), [waveletComp109](#), [waveletComp110](#), [waveletComp111](#), [waveletComp112](#), [waveletComp113](#), [waveletComp114](#), [waveletComp115](#), [waveletComp116](#), [waveletComp117](#), [waveletComp118](#), [waveletComp119](#), [waveletComp120](#), [waveletComp121](#), [waveletComp122](#), [waveletComp123](#), [waveletComp124](#), [waveletComp125](#), [waveletComp126](#), [waveletComp127](#), [waveletComp128](#), [waveletComp129](#), [waveletComp130](#), [waveletComp131](#), [waveletComp132](#), [waveletComp133](#), [waveletComp134](#), [waveletComp135](#), [waveletComp136](#), [waveletComp137](#), [waveletComp138](#), [waveletComp139](#), [waveletComp140](#), [waveletComp141](#), [waveletComp142](#), [waveletComp143](#), [waveletComp144](#), [waveletComp145](#), [waveletComp146](#), [waveletComp147](#), [waveletComp148](#), [waveletComp149](#), [waveletComp150](#), [waveletComp151](#), [waveletComp152](#), [waveletComp153](#), [waveletComp154](#), [waveletComp155](#), [waveletComp156](#), [waveletComp157](#), [waveletComp158](#), [waveletComp159](#), [waveletComp160](#), [waveletComp161](#), [waveletComp162](#), [waveletComp163](#), [waveletComp164](#), [waveletComp165](#), [waveletComp166](#), [waveletComp167](#), [waveletComp168](#), [waveletComp169](#), [waveletComp170](#), [waveletComp171](#), [waveletComp172](#), [waveletComp173](#), [waveletComp174](#), [waveletComp175](#), [waveletComp176](#), [waveletComp177](#), [waveletComp178](#), [waveletComp179](#), [waveletComp180](#), [waveletComp181](#), [waveletComp182](#), [waveletComp183](#), [waveletComp184](#), [waveletComp185](#), [waveletComp186](#), [waveletComp187](#), [waveletComp188](#), [waveletComp189](#), [waveletComp190](#), [waveletComp191](#), [waveletComp192](#), [waveletComp193](#), [waveletComp194](#), [waveletComp195](#), [waveletComp196](#), [waveletComp197](#), [waveletComp198](#), [waveletComp199](#), [waveletComp200](#), [waveletComp201](#), [waveletComp202](#), [waveletComp203](#), [waveletComp204](#), [waveletComp205](#), [waveletComp206](#), [waveletComp207](#), [waveletComp208](#), [waveletComp209](#), [waveletComp210](#), [waveletComp211](#), [waveletComp212](#), [waveletComp213](#), [waveletComp214](#), [waveletComp215](#), [waveletComp216](#), [waveletComp217](#), [waveletComp218](#), [waveletComp219](#), [waveletComp220](#), [waveletComp221](#), [waveletComp222](#), [waveletComp223](#), [waveletComp224](#), [waveletComp225](#), [waveletComp226](#), [waveletComp227](#), [waveletComp228](#), [waveletComp229](#), [waveletComp230](#), [waveletComp231](#), [waveletComp232](#), [waveletComp233](#), [waveletComp234](#), [waveletComp235](#), [waveletComp236](#), [waveletComp237](#), [waveletComp238](#), [waveletComp239](#), [waveletComp240](#), [waveletComp241](#), [waveletComp242](#), [waveletComp243](#), [waveletComp244](#), [waveletComp245](#), [waveletComp246](#), [waveletComp247](#), [waveletComp248](#), [waveletComp249](#), [waveletComp250](#), [waveletComp251](#), [waveletComp252](#), [waveletComp253](#), [waveletComp254](#), [waveletComp255](#), [waveletComp256](#), [waveletComp257](#), [waveletComp258](#), [waveletComp259](#), [waveletComp260](#), [waveletComp261](#), [waveletComp262](#), [waveletComp263](#), [waveletComp264](#), [waveletComp265](#), [waveletComp266](#), [waveletComp267](#), [waveletComp268](#), [waveletComp269](#), [waveletComp270](#), [waveletComp271](#), [waveletComp272](#), [waveletComp273](#), [waveletComp274](#), [waveletComp275](#), [waveletComp276](#), [waveletComp277](#), [waveletComp278](#), [waveletComp279](#), [waveletComp280](#), [waveletComp281](#), [waveletComp282](#), [waveletComp283](#), [waveletComp284](#), [waveletComp285](#), [waveletComp286](#), [waveletComp287](#), [waveletComp288](#), [waveletComp289](#), [waveletComp290](#), [waveletComp291](#), [waveletComp292](#), [waveletComp293](#), [waveletComp294](#), [waveletComp295](#), [waveletComp296](#), [waveletComp297](#), [waveletComp298](#), [waveletComp299](#), [waveletComp300](#), [waveletComp301](#), [waveletComp302](#), [waveletComp303](#), [waveletComp304](#), [waveletComp305](#), [waveletComp306](#), [waveletComp307](#), [waveletComp308](#), [waveletComp309](#), [waveletComp310](#), [waveletComp311](#), [waveletComp312](#), [waveletComp313](#), [waveletComp314](#), [waveletComp315](#), [waveletComp316](#), [waveletComp317](#), [waveletComp318](#), [waveletComp319](#), [waveletComp320](#), [waveletComp321](#), [waveletComp322](#), [waveletComp323](#), [waveletComp324](#), [waveletComp325](#), [waveletComp326](#), [waveletComp327](#), [waveletComp328](#), [waveletComp329](#), [waveletComp330](#), [waveletComp331](#), [waveletComp332](#), [waveletComp333](#), [waveletComp334](#), [waveletComp335](#), [waveletComp336](#), [waveletComp337](#), [waveletComp338](#), [waveletComp339](#), [waveletComp340](#), [waveletComp341](#), [waveletComp342](#), [waveletComp343](#), [waveletComp344](#), [waveletComp345](#), [waveletComp346](#), [waveletComp347](#), [waveletComp348](#), [waveletComp349](#), [waveletComp350](#), [waveletComp351](#), [waveletComp352](#), [waveletComp353](#), [waveletComp354](#), [waveletComp355](#), [waveletComp356](#), [waveletComp357](#), [waveletComp358](#), [waveletComp359](#), [waveletComp360](#), [waveletComp361](#), [waveletComp362](#), [waveletComp363](#), [waveletComp364](#), [waveletComp365](#), [waveletComp366](#), [waveletComp367](#), [waveletComp368](#), [waveletComp369](#), [waveletComp370](#), [waveletComp371](#), [waveletComp372](#), [waveletComp373](#), [waveletComp374](#), [waveletComp375](#), [waveletComp376](#), [waveletComp377](#), [waveletComp378](#), [waveletComp379](#), [waveletComp380](#), [waveletComp381](#), [waveletComp382](#), [waveletComp383](#), [waveletComp384](#), [waveletComp385](#), [waveletComp386](#), [waveletComp387](#), [waveletComp388](#), [waveletComp389](#), [waveletComp390](#), [waveletComp391](#), [waveletComp392](#), [waveletComp393](#), [waveletComp394](#), [waveletComp395](#), [waveletComp396](#), [waveletComp397](#), [waveletComp398](#), [waveletComp399](#), [waveletComp400](#), [waveletComp401](#), [waveletComp402](#), [waveletComp403](#), [waveletComp404](#), [waveletComp405](#), [waveletComp406](#), [waveletComp407](#), [waveletComp408](#), [waveletComp409](#), [waveletComp410](#), [waveletComp411](#), [waveletComp412](#), [waveletComp413](#), [waveletComp414](#), [waveletComp415](#), [waveletComp416](#), [waveletComp417](#), [waveletComp418](#), [waveletComp419](#), [waveletComp420](#), [waveletComp421](#), [waveletComp422](#), [waveletComp423](#), [waveletComp424](#), [waveletComp425](#), [waveletComp426](#), [waveletComp427](#), [waveletComp428](#), [waveletComp429](#), [waveletComp430](#), [waveletComp431](#), [waveletComp432](#), [waveletComp433](#), [waveletComp434](#), [waveletComp435](#), [waveletComp436](#), [waveletComp437](#), [waveletComp438](#), [waveletComp439](#), [waveletComp440](#), [waveletComp441](#), [waveletComp442](#), [waveletComp443](#), [waveletComp444](#), [waveletComp445](#), [waveletComp446](#), [waveletComp447](#), [waveletComp448](#), [waveletComp449](#), [waveletComp450](#), [waveletComp451](#), [waveletComp452](#), [waveletComp453](#), [waveletComp454](#), [waveletComp455](#), [waveletComp456](#), [waveletComp457](#), [waveletComp458](#), [waveletComp459](#), [waveletComp460](#), [waveletComp461](#), [waveletComp462](#), [waveletComp463](#), [waveletComp464](#), [waveletComp465](#), [waveletComp466](#), [waveletComp467](#), [waveletComp468](#), [waveletComp469](#), [waveletComp470](#), [waveletComp471](#), [waveletComp472](#), [waveletComp473](#), [waveletComp474](#), [waveletComp475](#), [waveletComp476](#), [waveletComp477](#), [waveletComp478](#), [waveletComp479](#), [waveletComp480](#), [waveletComp481](#), [waveletComp482](#), [waveletComp483](#), [waveletComp484](#), [waveletComp485](#), [waveletComp486](#), [waveletComp487](#), [waveletComp488](#), [waveletComp489](#), [waveletComp490](#), [waveletComp491](#), [waveletComp492](#), [waveletComp493](#), [waveletComp494](#), [waveletComp495](#), [waveletComp496](#), [waveletComp497](#), [waveletComp498](#), [waveletComp499](#), [waveletComp500](#), [waveletComp501](#), [waveletComp502](#), [waveletComp503](#), [waveletComp504](#), [waveletComp505](#), [waveletComp506](#), [waveletComp507](#), [waveletComp508](#), [waveletComp509](#), [waveletComp510](#), [waveletComp511](#), [waveletComp512](#), [waveletComp513](#), [waveletComp514](#), [waveletComp515](#), [waveletComp516](#), [waveletComp517](#), [waveletComp518](#), [waveletComp519](#), [waveletComp520](#), [waveletComp521](#), [waveletComp522](#), [waveletComp523](#), [waveletComp524](#), [waveletComp525](#), [waveletComp526](#), [waveletComp527](#), [waveletComp528](#), [waveletComp529](#), [waveletComp530](#), [waveletComp531](#), [waveletComp532](#), [waveletComp533](#), [waveletComp534](#), [waveletComp535](#), [waveletComp536](#), [waveletComp537](#), [waveletComp538](#), [waveletComp539](#), [waveletComp540](#), [waveletComp541](#), [waveletComp542](#), [waveletComp543](#), [waveletComp544](#), [waveletComp545](#), [waveletComp546](#), [waveletComp547](#), [waveletComp548](#), [waveletComp549](#), [waveletComp550](#), [waveletComp551](#), [waveletComp552](#), [waveletComp553](#), [waveletComp554](#), [waveletComp555](#), [waveletComp556](#), [waveletComp557](#), [waveletComp558](#), [waveletComp559](#), [waveletComp560](#), [waveletComp561](#), [waveletComp562](#), [waveletComp563](#), [waveletComp564](#), [waveletComp565](#), [waveletComp566](#), [waveletComp567](#), [waveletComp568](#), [waveletComp569](#), [waveletComp570](#), [waveletComp571](#), [waveletComp572](#), [waveletComp573](#), [waveletComp574](#), [waveletComp575](#), [waveletComp576](#), [waveletComp577](#), [waveletComp578](#), [waveletComp579](#), [waveletComp580](#), [waveletComp581](#), [waveletComp582](#), [waveletComp583](#), [waveletComp584](#), [waveletComp585](#), [waveletComp586](#), [waveletComp587](#), [waveletComp588](#), [waveletComp589](#), [waveletComp590](#), [waveletComp591](#), [waveletComp592](#), [waveletComp593](#), [waveletComp594](#), [waveletComp595](#), [waveletComp596](#), [waveletComp597](#), [waveletComp598](#), [waveletComp599](#), [waveletComp600](#), [waveletComp601](#), [waveletComp602](#), [waveletComp603](#), [waveletComp604](#), [waveletComp605](#), [waveletComp606](#), [waveletComp607](#), [waveletComp608](#), [waveletComp609](#), [waveletComp610](#), [waveletComp611](#), [waveletComp612](#), [waveletComp613](#), [waveletComp614](#), [waveletComp615](#), [waveletComp616](#), [waveletComp617](#), [waveletComp618](#), [waveletComp619](#), [waveletComp620](#), [waveletComp621](#), [waveletComp622](#), [waveletComp623](#), [waveletComp624](#), [waveletComp625](#), [waveletComp626](#), [waveletComp627](#), [waveletComp628](#), [waveletComp629](#), [waveletComp630](#), [waveletComp631](#), [waveletComp632](#), [waveletComp633](#), [waveletComp634](#), [waveletComp635](#), [waveletComp636](#), [waveletComp637](#), [waveletComp638](#), [waveletComp639](#), [waveletComp640](#), [waveletComp641](#), [waveletComp642](#), [waveletComp643](#), [waveletComp644](#), [waveletComp645](#), [waveletComp646](#), [waveletComp647](#), [waveletComp648](#), [waveletComp649](#), [waveletComp650](#), [waveletComp651](#), [waveletComp652](#), [waveletComp653](#), [waveletComp654](#), [waveletComp655](#), [waveletComp656](#), [waveletComp657](#), [waveletComp658](#), [waveletComp659](#), [waveletComp660](#), [waveletComp661](#), [waveletComp662](#), [waveletComp663](#), [waveletComp664](#), [waveletComp665](#), [waveletComp666](#), [waveletComp667](#), [waveletComp668](#), [waveletComp669](#), [waveletComp670](#), [waveletComp671](#), [waveletComp672](#), [waveletComp673](#), [waveletComp674](#), [waveletComp675](#), [waveletComp676](#), [waveletComp677](#), [waveletComp678](#), [waveletComp679](#), [waveletComp680](#), [waveletComp681](#), [waveletComp682](#), [waveletComp683](#), [waveletComp684](#), [waveletComp685](#), [waveletComp686](#), [waveletComp687](#), [waveletComp688](#), [waveletComp689](#), [waveletComp690](#), [waveletComp691](#), [waveletComp692](#), [waveletComp693](#), [waveletComp694](#), [waveletComp695](#), [waveletComp696](#), [waveletComp697](#), [waveletComp698](#), [waveletComp699](#), [waveletComp700](#), [waveletComp701](#), [waveletComp702](#), [waveletComp703](#), [waveletComp704](#), [waveletComp705](#), [waveletComp706](#), [waveletComp707](#), [waveletComp708](#), [waveletComp709](#), [waveletComp710](#), [waveletComp711](#), [waveletComp712](#), [waveletComp713](#), [waveletComp714](#), [waveletComp715](#), [waveletComp716](#), [waveletComp717](#), [waveletComp718](#), [waveletComp719](#), [waveletComp720](#), [waveletComp721](#), [waveletComp722](#), [waveletComp723](#), [waveletComp724](#), [waveletComp725](#), [waveletComp726](#), [waveletComp727](#), [waveletComp728](#), [waveletComp729](#), [waveletComp730](#), [waveletComp731](#), [waveletComp732](#), [waveletComp733](#), [waveletComp734](#), [waveletComp735](#), [waveletComp736](#), [waveletComp737](#), [waveletComp738](#), [waveletComp739](#), [waveletComp740](#), [waveletComp741](#), [waveletComp742](#), [waveletComp743](#), [waveletComp744](#), [waveletComp745](#), [waveletComp746](#), [waveletComp747](#), [waveletComp748](#), [waveletComp749](#), [waveletComp750](#), [waveletComp751](#), [waveletComp752](#), [waveletComp753](#), [waveletComp754](#), [waveletComp755](#), [waveletComp756](#), [waveletComp757](#), [waveletComp758](#), [waveletComp759](#), [waveletComp760](#), [waveletComp761](#), [waveletComp762](#), [waveletComp763](#), [waveletComp764](#), [waveletComp765](#), [waveletComp766](#), [waveletComp767](#), [waveletComp768](#), [waveletComp769](#), [waveletComp770](#), [waveletComp771](#), [waveletComp772](#), [waveletComp773](#), [waveletComp774](#), [waveletComp775](#), [waveletComp776</](#)

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/4.3-arm64/Resources/1"
[2] "/Library/Frameworks/R.framework/Versions/4.3-arm64/Resources/2"
[3] "/Library/Frameworks/R.framework/Versions/4.3-arm64/Resources/3"
[4] "/Library/Frameworks/R.framework/Versions/4.3-arm64/Resources/4"
[5] "/Library/Frameworks/R.framework/Versions/4.3-arm64/Resources/5"
[6] "/Library/Frameworks/R.framework/Versions/4.3-arm64/Resources/6"
[7] "/Library/Frameworks/R.framework/Versions/4.3-arm64/Resources/7"
[8] "/Library/Frameworks/R.framework/Versions/4.3-arm64/Resources/8"
[9] "/Library/Frameworks/R.framework/Versions/4.3-arm64/Resources/9"
[10] "/Library/Frameworks/R.framework/Versions/4.3-arm64/Resources/10"
[11] "/Library/Frameworks/R.framework/Versions/4.3-arm64/Resources/11"
[12] "/Library/Frameworks/R.framework/Versions/4.3-arm64/Resources/12"
[13] "/Library/Frameworks/R.framework/Versions/4.3-arm64/Resources/13"
[14] "/Library/Frameworks/R.framework/Versions/4.3-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
+   vapply = vapply(1:NROW(matrixv), function(indeks)
+     sum(matrixv[indeks, ]),
+     FUN.VALUE = c(sum=0)), # end vapply
+   sapply = 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.

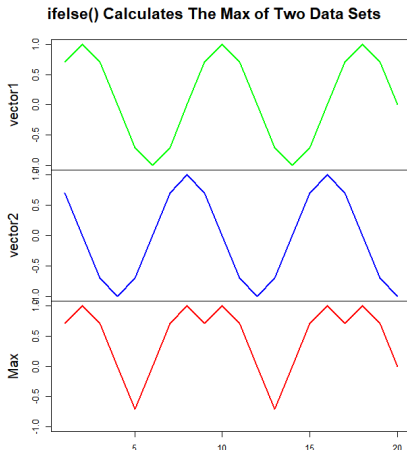
```
> # rnorm() vectorized with respect to "mean" and "sd"
> vec_rnorm <- function(n, mean=0, sd=1) {
+   if (NROW(mean)==1 && NROW(sd)==1)
+     rnorm(n=n, mean=mean, sd=sd)
+   else
+     mapply(rnorm, n=n, mean=mean, sd=sd)
+ } # end vec_rnorm
> # Call vec_rnorm() on vector of "sd"
> vec_rnorm(n=2, 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(vectorv))
+   cumsumv2[i] <- (vectorv[i] + cumsumv2[i-1])
> # Compare the outputs of the two methods
> all.equal(cumsumv, cumsumv2)
> # Microbenchmark the two methods
> library(microbenchmark)
> summary(microbenchmark(
+   cumsum=cumsum(vectorv), # Vectorized
+   loop_alloc={cumsumv2 <- vectorv # Allocate memory to cumsumv3
+     for (i in 2:NROW(vectorv))
+       cumsumv2[i] <- (vectorv[i] + cumsumv2[i-1])
+   },
+   loop_nalloc={cumsumv3 <- vectorv[1] # Doesn't allocate memory to
+     for (i in 2:NROW(vectorv))
+       cumsumv3[i] <- (vectorv[i] + cumsumv3[i-1])
```

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.

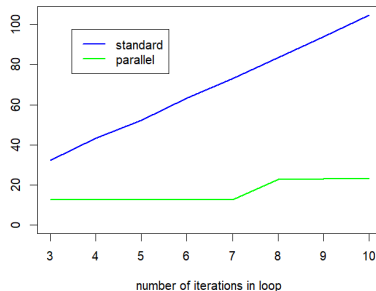
The function `plot()` by default plots a scatterplot, but can also plot lines using the argument `type="l"`.

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],
+   quantv = 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

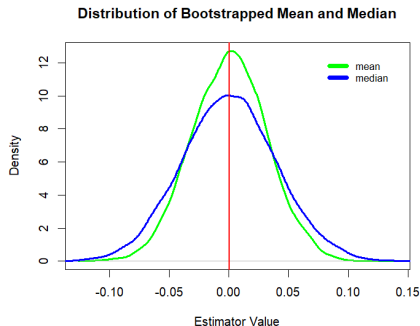
The standard errors of estimators can be calculated using a *bootstrap* simulation.

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:

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

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)
> nrows <- 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, nrows) {
+     samplev <- retp[sample.int(nrows, replace=TRUE)]
+     c(sd=sd(samplev), mad=mad(samplev))
+   }, retp=retp, nrows=nrows) # end parLapply
> # Bootstrap sd and MAD under Mac-OSX or Linux
> bootd <- mclapply(1:nboot, function(x) {
+   samplev <- retp[sample.int(nrows, 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
> predm <- rnorm(nrows, mean=2)
> noise <- rnorm(nrows)
> respv <- (-3 + 2*predictor + noise)
> desv <- cbind(respv, predm)
> # 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 ~ predm, 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
```

Distribution of Bootstrapped Regression Coefficients

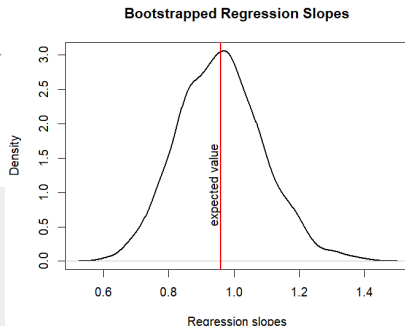
The *bootstrapped* coefficient values can be used to calculate the probability distribution of the coefficients and their standard errors,

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

`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), stdererror=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).

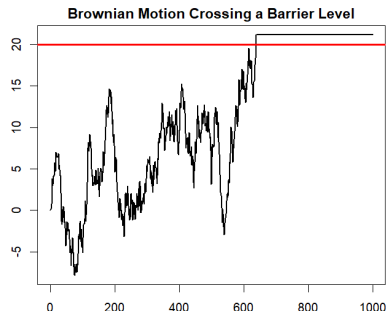
```
> # Collapse the bootstrap list into a vector
> class(bootd)
> bootd <- unlist(bootd)
> # Mean and standard error of beta regression coefficient
> c(mean=mean(bootd), stderror=sd(bootd))
> # Plot density of bootstrapped beta coefficients
> plot(density(bootd),
+      lwd=2, xlab="Regression slopes",
+      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] <- rnorm(1) # 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)
```

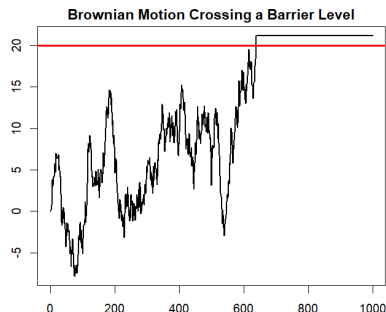


Simulating Brownian Motion Using Vectorized Functions

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

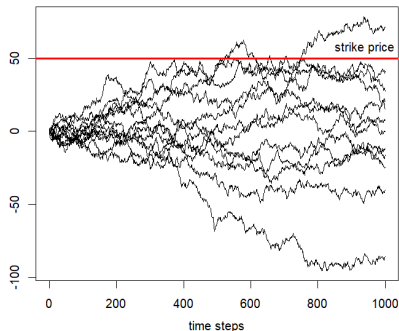
The statistics of Brownian motion can be estimated by simulating multiple paths.

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
> # Perform parallel bootstrap under Windows
> clusterSetRNGStream(cluster, 1121) # Reset random number generato
> 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 Window
> 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
> # Perform parallel bootstrap under Windows
> clusterSetRNGStream(cluster, 1121) # Reset random number generat
> 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 Window
> 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", "QQQ")
> # 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 Database for Investment Strategies

The database 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.

The *SPY ETF* owns the *S&P500* index constituents. *SPY* is the biggest, the most liquid, and the oldest ETF. *SPY* has over \$400 billion of shares outstanding, and trades over \$20 billion per day, at a bid-ask spread of only one tick (\$0.01, or about 0.0022%).

The *QQQ ETF* owns the *Nasdaq-100* index constituents.

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.

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
IVW	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 ETF
SVXY	Short VIX Futures	Commodity Based ETF
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 ETF
QQQ	QQQ Trust	US Equity ETF

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.

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*.

Stock Databases And Survivorship Bias

The file `sp500_constituents.csv` contains a *data frame* of over 700 present (and also some past) *S&P500* index constituents.

The file `sp500_constituents.csv` is updated with stocks recently added to the *S&P500* index by downloading the *SPY ETF Holdings*.

But the file `sp500_constituents.csv` doesn't include companies that have gone bankrupt. For example, it doesn't include Enron, which was in the *S&P500* index before it went bankrupt in 2001.

Most databases of stock prices don't include companies that have gone bankrupt or have been liquidated.

This introduces a *survivorship bias* to the data, which can skew portfolio simulations and strategy backtests.

Accurate strategy simulations require starting with a portfolio of companies at a "point in time" in the past, and tracking them over time.

Research databases like the *WRDS* database provide stock prices of companies that are no longer traded.

The stock tickers are stored in the column "Ticker" of the `sp500 data frame`.

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"
```

Wharton Research Data Services WRDS

Wharton Research Data Services (*WRDS*) is a distributor of premium third party data for the academic and research communities.

WRDS provides time series of security prices and fundamental company data, and other financial, econometric, and social datasets.

WRDS provides stock price, options and implied volatilities, stock fundamentals, financial ratios, zoo::indexes, earnings estimates, analyst ratings, etc.

WRDS redistributes fundamental company data from *Compustat*, *S&P Capital IQ*, *Thomson Reuters*, *FactSet*, *Hedge Fund Research*, *Markit*, etc.

NYU students can obtain user accounts for *WRDS* data.

Subscriptions

- Bank Regulatory
- Beta Suite by WRDS
- Blockholders
- CBOE Indexes
- Compustat - Capital IQ
- Contributed Data
- CRSP
- CSMAR
- DMEF Academic Data
- Dow Jones
- Financial Ratios Suite by WRDS
- IBES
- IHS Global Insight
- Infotrip
- Institutional Shareholder Services (ISS)
- Intraday Indicators by WRDS
- IRI
- Linking Suite by WRDS
- OTC Markets
- Penn World Tables
- Peters and Taylor Total Q
- PHLX
- Public
- Research Quotient
- SEC Order Execution
- TAQ
- Thomson Reuters
- TRACE

Analytics by WRDS

SEC Analytics Suite

Researchers will gain an overview of WRDS' powerful SEC research platform. Some key features they will learn about include:

- Access to 15 million filings from a single index
- Full-text searching over 3.5 million filings
- Sentiment analysis

U.S. Daily Event Study

Investigate abnormal stock returns/volumes around event dates by uploading your own "events" file, or analyzing reaction to firm-specific events from Capital IQ's Key Development database.

- Instant visualization of the effect of events on U.S. equities
- Output includes statistics, plots

Intraday Indicators

Stock specific daily and intraday (5 min, 15 min and 30 min) indicators created from the TAQ intraday datasets.

- Stock specific daily and intraday (5 min, 15 min and 30 min) indicators.
- Requires subscription to TAQ dataset.

Classroom by WRDS

Kernel Density of Asset Returns

The kernel density is proportional to the number of data points close to a given point.

The kernel density is analogous to a histogram, but it provides more detailed information about the distribution of the data.

The smoothing kernel $K(x)$ is a symmetric function which decreases with the distance x .

The kernel density d_r at a point r is equal to the sum over the kernel function $K(x)$:

$$d_r = \sum_{j=1}^n K(r - r_j)$$

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

The parameter *smoothing bandwidth* is the standard deviation of the smoothing kernel $K(x)$.

The function `density()` returns a vector of densities at equally spaced points, not for the original data points.

The function `approx()` interpolates a vector of data into another vector.

```
> library(rutils) # Load package rutils
> # Calculate VTI percentage returns
> retp <- rutils::etfenv$returns$VTI
> retp <- drop(coredata(na.omit(retp)))
> nrow <- NROW(retp)
> # Mean and standard deviation of returns
> c(mean(retp), sd(retp))
> # Calculate the smoothing bandwidth as the MAD of returns 10 points
> retp <- sort(retp)
> bwidh <- 10*mad(rutils::diffit(retp, lagg=10))
> # Calculate the kernel density
> densv <- sapply(1:nrow, function(it) {
+   sum(dnorm(retp-retp[it], sd=bwidh))
+ }) # end sapply
> madv <- mad(retp)
> plot(retp, densv, xlim=c(-5*madv, 5*madv),
+   t="l", col="blue", lwd=3,
+   xlab="returns", ylab="density",
+   main="Density of VTI Returns")
> # Calculate the kernel density using density()
> densv <- density(retp, bw=bwidh)
> NROW(densv$y)
> x11(width=6, height=5)
> plot(densv, xlim=c(-5*madv, 5*madv),
+   xlab="returns", ylab="density",
+   col="blue", lwd=3, main="Density of VTI Returns")
> # Interpolate the densv vector into returns
> densv <- approx(densv$x, densv$y, xout=retp)
> all.equal(densv$x, retp)
> plot(densv, xlim=c(-5*madv, 5*madv),
+   xlab="returns", ylab="density",
+   t="l", col="blue", lwd=3,
+   main="Density of VTI Returns")
```

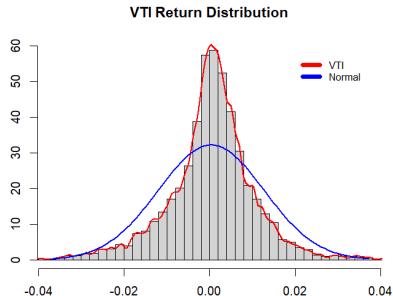
Distribution of Asset Returns

Asset returns are usually not normally distributed and they exhibit *leptokurtosis* (large kurtosis, or fat tails).

The function `hist()` calculates and plots a histogram, and returns its data *invisibly*.

The parameter `breaks` is the number of cells of the histogram.

The function `lines()` draws a line through specified points.



```
> # Plot histogram
> histp <- hist(retp, breaks=100, freq=FALSE,
+   xlim=c(-5*madv, 5*madv), xlab="", ylab="",
+   main="VTI Return Distribution")
> # Draw kernel density of histogram
> lines(densv, col="red", lwd=2)
> # Add density of normal distribution
> curve(expr=dnorm(x, mean=mean(retp), sd=sd(retp)),
+   add=TRUE, lwd=2, col="blue")
> # Add legend
> legend("topright", inset=0.05, cex=0.8, title=NULL,
+   leg=c("VTI", "Normal"), bty="n",
+   lwd=6, bg="white", col=c("red", "blue"))
```

The Quantile-Quantile Plot

A *Quantile-Quantile* (*Q-Q*) plot is a plot of points with the same *quantiles*, from two probability distributions.

If the two distributions are similar then all the points in the *Q-Q* plot lie along the diagonal.

The *VTI Q-Q* plot shows that the *VTI* return distribution has fat tails.

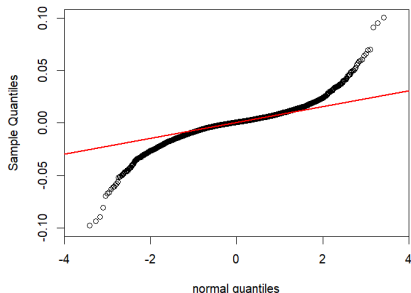
The *p*-value of the *Shapiro-Wilk* test is very close to zero, which shows that the *VTI* returns are very unlikely to be normal.

The function `shapiro.test()` performs the *Shapiro-Wilk* test of normality.

The function `qqnorm()` produces a normal *Q-Q* plot.

The function `qqline()` fits a line to the normal quantiles.

VTI Q-Q Plot



```
> # Create normal Q-Q plot
> qqnorm(retp, ylim=c(-0.1, 0.1), main="VTI Q-Q Plot",
+   xlab="Normal Quantiles")
> # Fit a line to the normal quantiles
> qqline(retp, col="red", lwd=2)
> # Perform Shapiro-Wilk test
> shapiro.test(retp[1:499])
```

Boxplots of Distributions of Values

Box-and-whisker plots (*boxplots*) are graphical representations of a distribution of values.

The bottom and top box edges (*hinges*) are equal to the first and third quartiles, and the *box* width is equal to the interquartile range (*IQR*).

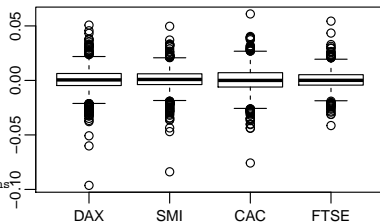
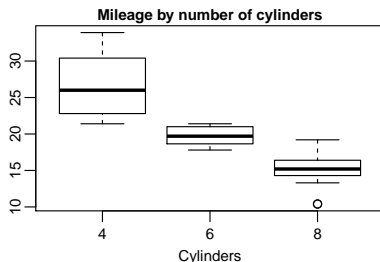
The nominal range is equal to 1.5 times the *IQR* above and below the box *hinges*.

The *whiskers* are dashed vertical lines representing values beyond the first and third quartiles, but within the nominal range.

The *whiskers* end at the last values within the nominal range, while the open circles represent outlier values beyond the nominal range.

The function `boxplot()` has two methods: one for formula objects (for categorical variables), and another for data frames.

```
> # Boxplot method for formula
> boxplot(formula=mpg ~ cyl, data=mtcars,
+   main="Mileage by number of cylinders",
+   xlab="Cylinders", ylab="Miles per gallon")
> # Boxplot method for data frame of EuStockMarkets percentage returns
> boxplot(x=diff(log(EuStockMarkets)))
```



Higher Moments of Asset Returns

The estimators of moments of a probability distribution are given by:

Sample mean: $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

Sample variance: $\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$

With their expected values equal to the population mean and standard deviation:

$\mathbb{E}[\bar{x}] = \mu$ and $\mathbb{E}[\hat{\sigma}] = \sigma$

The sample skewness (third moment):

$$\varsigma = \frac{n}{(n-1)(n-2)} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\hat{\sigma}} \right)^3$$

The sample kurtosis (fourth moment):

$$\kappa = \frac{n(n+1)}{(n-1)(n-2)(n-3)} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\hat{\sigma}} \right)^4$$

The normal distribution has skewness equal to 0 and kurtosis equal to 3.

Stock returns typically have negative skewness and kurtosis much greater than 3.

```
> # Calculate VTI percentage returns
> retp <- na.omit(rutils::etfenv$returns$VTI)
> # Number of observations
> nrow <- NROW(retp)
> # Mean of VTI returns
> retm <- mean(retp)
> # Standard deviation of VTI returns
> stdev <- sd(retp)
> # Skewness of VTI returns
> nrow/((nrow-1)*(nrow-2))*sum(((retp - retm)/stdev)^3)
> # Kurtosis of VTI returns
> nrow*(nrow+1)/((nrow-1)^3)*sum(((retp - retm)/stdev)^4)
> # Random normal returns
> retp <- rnorm(nrow, sd=stdev)
> # Mean and standard deviation of random normal returns
> retm <- mean(retp)
> stdev <- sd(retp)
> # Skewness of random normal returns
> nrow/((nrow-1)*(nrow-2))*sum(((retp - retm)/stdev)^3)
> # Kurtosis of random normal returns
> nrow*(nrow+1)/((nrow-1)^3)*sum(((retp - retm)/stdev)^4)
```

Functions for Calculating Skew and Kurtosis

R provides an easy way for users to write functions.

The function `calc_skew()` calculates the skew of returns, and `calc_kurt()` calculates the kurtosis.

Functions return the value of the last expression that is evaluated.

```
> # calc_skew() calculates skew of returns
> calc_skew <- function(retp) {
+   retp <- na.omit(retp)
+   sum(((retp - mean(retp))/sd(retp))^3)/NROW(retp)
+ } # end calc_skew
> # calc_kurt() calculates kurtosis of returns
> calc_kurt <- function(retp) {
+   retp <- na.omit(retp)
+   sum(((retp - mean(retp))/sd(retp))^4)/NROW(retp)
+ } # end calc_kurt
> # Calculate skew and kurtosis of VTI returns
> calc_skew(retp)
> calc_kurt(retp)
> # calcmom() calculates the moments of returns
> calcmom <- function(retp, moment=3) {
+   retp <- na.omit(retp)
+   sum(((retp - mean(retp))/sd(retp))^moment)/NROW(retp)
+ } # end calcmom
> # Calculate skew and kurtosis of VTI returns
> calcmom(retp, moment=3)
> calcmom(retp, moment=4)
```

Standard Errors of Estimators

Statistical estimators are functions of samples (which are random variables), and therefore are themselves *random variables*.

The *standard error* (SE) of an estimator is defined as its *standard deviation* (not to be confused with the *population standard deviation* of the underlying random variable).

For example, the *standard error* of the estimator of the mean is equal to:

$$\sigma_{\mu} = \frac{\sigma}{\sqrt{n}}$$

Where σ is the *population standard deviation* (which is usually unknown).

The *estimator* of this *standard error* is equal to:

$$SE_{\mu} = \frac{\hat{\sigma}}{\sqrt{n}}$$

where: $\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$ is the sample standard deviation (the estimator of the population standard deviation).

```
> set.seed(1121) # Reset random number generator
> # Sample from Standard Normal Distribution
> nrows <- 1000
> datav <- rnorm(nrows)
> # Sample mean
> mean(datav)
> # Sample standard deviation
> sd(datav)
> # Standard error of sample mean
> sd(datav)/sqrt(nrows)
```

Normal (Gaussian) Probability Distribution

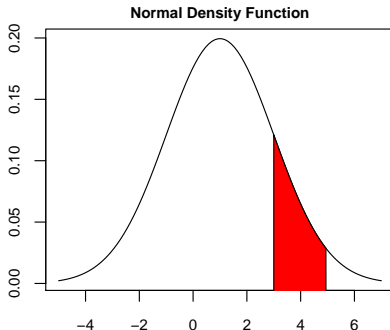
The *Normal (Gaussian)* probability density function is given by:

$$\phi(x, \mu, \sigma) = \frac{e^{-(x-\mu)^2/2\sigma^2}}{\sigma\sqrt{2\pi}}$$

The *Standard Normal* distribution $\phi(0, 1)$ is a special case of the *Normal* $\phi(\mu, \sigma)$ with $\mu = 0$ and $\sigma = 1$.

The function `dnorm()` calculates the *Normal* probability density.

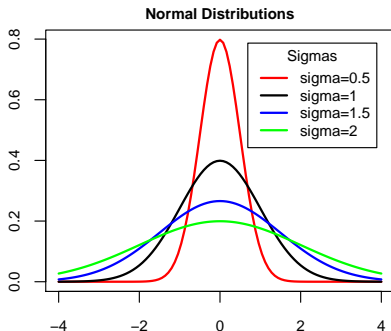
```
> xvar <- seq(-5, 7, length=100)
> yvar <- dnorm(xvar, mean=1.0, sd=2.0)
> plot(xvar, yvar, type="l", lty="solid", xlab="", ylab="")
> title(main="Normal Density Function", line=0.5)
> startp <- 3; endd <- 5 # Set lower and upper bounds
> # Set polygon base
> subv <- ((xvar >= startp) & (xvar <= endd))
> polygon(c(startp, xvar[subv], endd), # Draw polygon
+ c(-1, yvar[subv], -1), col="red")
```



Normal (Gaussian) Probability Distributions

Plots of several *Normal* distributions with different values of σ , using the function `curve()` for plotting functions given by their name.

```
> sigmavs <- c(0.5, 1, 1.5, 2) # Sigma values
> # Create plot colors
> colorv <- c("red", "black", "blue", "green")
> # Create legend labels
> labelv <- paste("sigma", sigmavs, sep="")
> for (it in 1:4) { # Plot four curves
+   curve(expr=dnorm(x, sd=sigmavs[it]),
+     xlim=c(-4, 4), xlab="", ylab="", lwd=2,
+     col=colorv[it], add=as.logical(it-1))
+ } # end for
> # Add title
> title(main="Normal Distributions", line=0.5)
> # Add legend
> legend("topright", inset=0.05, title="Sigmas",
+   labelv, cex=0.8, lwd=2, lty=1, bty="n", col=colorv)
```



Student's t -distribution

Let z_1, \dots, z_ν be independent standard normal random variables, with sample mean: $\bar{z} = \frac{1}{\nu} \sum_{i=1}^{\nu} z_i$ ($\mathbb{E}[\bar{z}] = \mu$) and sample variance:

$$\hat{\sigma}^2 = \frac{1}{\nu-1} \sum_{i=1}^{\nu} (z_i - \bar{z})^2$$

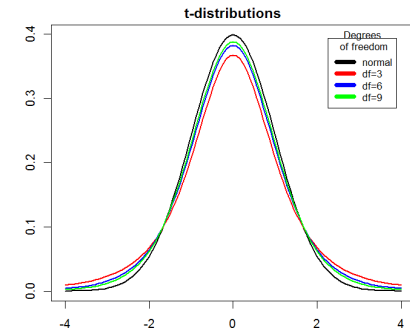
Then the random variable (t -ratio):

$$t = \frac{\bar{z} - \mu}{\hat{\sigma} / \sqrt{\nu}}$$

Follows the t -distribution with ν degrees of freedom, with the probability density function:

$$f(t) = \frac{\Gamma((\nu + 1)/2)}{\sqrt{\pi\nu} \Gamma(\nu/2)} (1 + t^2/\nu)^{-(\nu+1)/2}$$

```
> degf <- c(3, 6, 9) # Df values
> colorv <- c("black", "red", "blue", "green")
> labelv <- c("normal", paste("df", degf, sep=" "))
> # Plot a Normal probability distribution
> curve(expr=dnorm, xlim=c(-4, 4), xlab="", ylab="", lwd=2)
> for (it in 1:3) { # Plot three t-distributions
+   curve(expr=dt(x, df=degf[it]), xlab="", ylab="",
+   lwd=2, col=colorv[it+1], add=TRUE)
+ } # end for
```



```
> # Add title
> title(main="t-distributions", line=0.5)
> # Add legend
> legend("topright", inset=0.05, bty="n",
+       title="Degrees\n of freedom", labelv,
+       cex=0.8, lwd=6, lty=1, col=colorv)
```

Mixture Models of Returns

Mixture models are produced by randomly sampling data from different distributions.

The mixture of two normal distributions with different variances produces a distribution with *leptokurtosis* (large kurtosis, or fat tails).

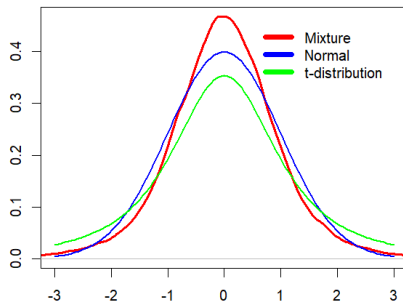
Student's *t-distribution* has fat tails because the sample variance in the denominator of the *t-ratio* is variable.

The time-dependent volatility of asset returns is referred to as *heteroskedasticity*.

Random processes with *heteroskedasticity* can be considered a type of mixture model.

The *heteroskedasticity* produces *leptokurtosis* (large kurtosis, or fat tails).

Mixture of Normal Returns



```
> # Mixture of two normal distributions with sd=1 and sd=2
> nrows <- 1e5
> retp <- c(rnorm(nrows/2), 2*rnorm(nrows/2))
> retp <- (retp-mean(retp))/sd(retp)
> # Kurtosis of normal
> calc_kurt(rnorm(nrows))
> # Kurtosis of mixture
> calc_kurt(retp)
> # Or
> nrows*sum(retp^4)/(nrows-1)^2
```

```
> # Plot the distributions
> plot(density(retp), xlab="", ylab="",
+      main="Mixture of Normal Returns",
+      xlim=c(-3, 3), type="l", lwd=3, col="red")
> curve(expr=dnorm, lwd=2, col="blue", add=TRUE)
> curve(expr=dt(x, df=3), lwd=2, col="green", add=TRUE)
> # Add legend
> legend("topright", inset=0.05, lty=1, lwd=6, bty="n",
+       legend=c("Mixture", "Normal", "t-distribution"),
+       col=c("red", "blue", "green"))
```

Non-standard Student's *t*-distribution

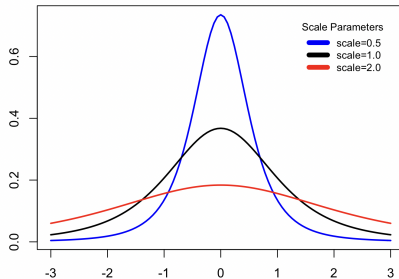
The non-standard Student's *t*-distribution has the probability density function:

$$f(t) = \frac{\Gamma((\nu + 1)/2)}{\sqrt{\pi\nu} \sigma \Gamma(\nu/2)} \left(1 + \left(\frac{t - \mu}{\sigma}\right)^2 / \nu\right)^{-(\nu+1)/2}$$

It has non-zero mean equal to the location parameter μ , and a standard deviation proportional to the scale parameter σ .

```
> dev.new(width=6, height=5, noRStudioGD=TRUE)
> # x11(width=6, height=5)
> # Define density of non-standard t-distribution
> tdistr <- function(x, dfree, locv=0, scalev=1) {
+   dt((x-locv)/scalev, df=dfree)/scalev
+ } # end tdistr
> # Or
> tdistr <- function(x, dfree, locv=0, scalev=1) {
+   gamma((dfree+1)/2)/(sqrt(pi*dfree)*gamma(dfree/2)*scalev)*
+   (1+((x-locv)/scalev)^2/dfree)^(-(dfree+1)/2)
+ } # end tdistr
> # Calculate vector of scale values
> scalev <- c(0.5, 1.0, 2.0)
> colorv <- c("blue", "black", "red")
> labelv <- paste("scale", format(scalev, digits=2), sep="")
> # Plot three t-distributions
> for (it in 1:3) {
+   curve(expr=tdistr(x, dfree=3, scalev=scalev[it]), xlim=c(-3, 3),
+   xlab="", ylab="", lwd=2, col=colorv[it], add=(it>1))
+ } # end for
```

t-distributions with Different Scale Parameters



```
> # Add title
> title(main="t-distributions with Different Scale Parameters", line=1)
> # Add legend
> legend("topright", inset=0.05, bty="n", title="Scale Parameters",
+       cex=0.8, lwd=6, lty=1, col=colorv)
```

The *Shapiro-Wilk* Test of Normality

The *Shapiro-Wilk* test is designed to test the *null hypothesis* that a sample: $\{x_1, \dots, x_n\}$ is from a normally distributed population.

The test statistic is equal to:

$$W = \frac{(\sum_{i=1}^n a_i x_{(i)})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Where the: $\{a_1, \dots, a_n\}$ are proportional to the *order statistics* of random variables from the normal distribution.

$x_{(k)}$ is the *k*-th *order statistic*, and is equal to the *k*-th smallest value in the sample: $\{x_1, \dots, x_n\}$.

The *Shapiro-Wilk* statistic follows its own distribution, and is less than or equal to 1.

The *Shapiro-Wilk* statistic is close to 1 for samples from normal distributions.

The *p*-value for *VTI* returns is extremely small, and we conclude that the *null hypothesis* is FALSE, and the *VTI* returns are not from a normally distributed population.

The *Shapiro-Wilk* test is not reliable for large sample sizes, so it's limited to less than 5000 sample size.

```
> # Calculate VTI percentage returns
> library(rutils)
> retp <- as.numeric(na.omit(rutils::etfenv$returns$VTI))[1:499]
> # Reduce number of output digits
> ndigits <- options(digits=5)
> # Shapiro-Wilk test for normal distribution
> nrow <- NROW(retp)
> shapiro.test(rnorm(nrow))
```

Shapiro-Wilk normality test

```
data:  rnorm(nrow)
W = 0.997, p-value = 0.46
> # Shapiro-Wilk test for VTI returns
> shapiro.test(retp)
```

Shapiro-Wilk normality test

```
data:  retp
W = 0.993, p-value = 0.022
> # Shapiro-Wilk test for uniform distribution
> shapiro.test(runif(nrow))
```

Shapiro-Wilk normality test

```
data:  runif(nrow)
W = 0.955, p-value = 3.6e-11
> # Restore output digits
> options(digits=ndigits$digits)
```

The Jarque-Bera Test of Normality

The *Jarque-Bera* test is designed to test the *null hypothesis* that a sample: $\{x_1, \dots, x_n\}$ is from a normally distributed population.

The test statistic is equal to:

$$JB = \frac{n}{6}(\varsigma^2 + \frac{1}{4}(\kappa - 3)^2)$$

Where the *skewness* and *kurtosis* are defined as:

$$\varsigma = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\hat{\sigma}} \right)^3 \quad \kappa = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\hat{\sigma}} \right)^4$$

The *Jarque-Bera* statistic asymptotically follows the *chi-squared* distribution with 2 degrees of freedom.

The *Jarque-Bera* statistic is small for samples from normal distributions.

The *p*-value for *VTI* returns is extremely small, and we conclude that the *null hypothesis* is FALSE, and the *VTI* returns are not from a normally distributed population.

```
> library(tseries) # Load package tseries
> # Jarque-Bera test for normal distribution
> jarque.bera.test(rnorm(nrows))
```

Jarque Bera Test

```
data:  rnorm(nrows)
X-squared = 0.8, df = 2, p-value = 0.7
> # Jarque-Bera test for VTI returns
> jarque.bera.test(retp)
```

Jarque Bera Test

```
data:  retp
X-squared = 2, df = 2, p-value = 0.4
> # Jarque-Bera test for uniform distribution
> jarque.bera.test(runif(NROW(retp)))
```

Jarque Bera Test

```
data:  runif(NROW(retp))
X-squared = 31, df = 2, p-value = 2e-07
```

The Kolmogorov-Smirnov Test for Probability Distributions

The *Kolmogorov-Smirnov* test *null hypothesis* is that two samples: $\{x_1, \dots, x_n\}$ and $\{y_1, \dots, y_n\}$ were obtained from the same probability distribution.

The *Kolmogorov-Smirnov* statistic depends on the maximum difference between two empirical cumulative distribution functions (cumulative frequencies):

$$D = \sup_i |P(x_i) - P(y_i)|$$

The function `ks.test()` performs the *Kolmogorov-Smirnov* test and returns the statistic and its *p*-value *invisibly*.

The second argument to `ks.test()` can be either a numeric vector of data values, or a name of a cumulative distribution function.

The *Kolmogorov-Smirnov* test can be used as a *goodness of fit* test, to test if a set of observations fits a probability distribution.

```
> # KS test for normal distribution
> ks_test <- ks.test(rnorm(100), pnorm)
> ks_test$p.value
> # KS test for uniform distribution
> ks.test(runif(100), pnorm)
> # KS test for two shifted normal distributions
> ks.test(rnorm(100), rnorm(100, mean=0.1))
> ks.test(rnorm(100), rnorm(100, mean=1.0))
> # KS test for two different normal distributions
> ks.test(rnorm(100), rnorm(100, sd=2.0))
> # KS test for VTI returns vs normal distribution
> retp <- as.numeric(na.omit(rutils::etfenv$returns$VTI))
> retp <- (retp - mean(retp))/sd(retp)
> ks.test(retp, pnorm)
```

Chi-squared Distribution

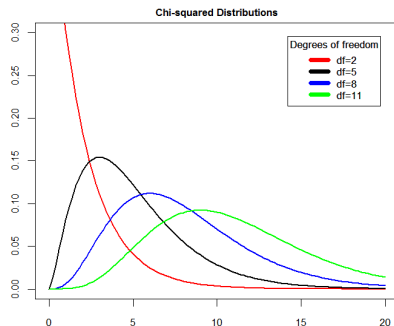
Let z_1, \dots, z_k be independent standard *Normal* random variables.

Then the random variable $X = \sum_{i=1}^k z_i^2$ is distributed according to the *Chi-squared* distribution with k degrees of freedom: $X \sim \chi_k^2$, and its probability density function is given by:

$$f(x) = \frac{x^{k/2-1} e^{-x/2}}{2^{k/2} \Gamma(k/2)}$$

The *Chi-squared* distribution with k degrees of freedom has mean equal to k and variance equal to $2k$.

```
> # Degrees of freedom
> degf <- c(2, 5, 8, 11)
> # Plot four curves in loop
> colorv <- c("red", "black", "blue", "green")
> for (it in 1:4) {
+   curve(dchisq(x, df=degf[it]),
+         xlim=c(0, 20), ylim=c(0, 0.3),
+         xlab="", ylab="", col=colorv[it],
+         lwd=2, add=as.logical(it-1))
+ } # end for
```



```
> # Add title
> title(main="Chi-squared Distributions", line=0.5)
> # Add legend
> labelv <- paste("df", degf, sep="=")
> legend("topright", inset=0.05, bty="n",
+       title="Degrees of freedom", labelv,
+       cex=0.8, lwd=6, lty=1, col=colorv)
```


The *Chi-squared* Test for the Goodness of Fit

Goodness of Fit tests are designed to test if a set of observations fits an assumed theoretical probability distribution.

The *Chi-squared* test tests if a frequency of counts fits the specified distribution.

The *Chi-squared* statistic is the sum of squared differences between the observed frequencies o_i and the theoretical frequencies p_i :

$$\chi^2 = N \sum_{i=1}^n \frac{(o_i - p_i)^2}{p_i}$$

Where N is the total number of observations.

The *null hypothesis* is that the observed frequencies are consistent with the theoretical distribution.

The function `chisq.test()` performs the *Chi-squared* test and returns the statistic and its *p*-value *invisibly*.

The parameter `breaks` in the function `hist()` should be chosen large enough to capture the shape of the frequency distribution.

```
> # Observed frequencies from random normal data
> histp <- hist(rnorm(1e3, mean=0), breaks=100, plot=FALSE)
> countsn <- histp$counts
> # Theoretical frequencies
> countst <- rutils::diffit(pnorm(histp$breaks))
> # Perform Chi-squared test for normal data
> chisq.test(x=countsn, p=countst, rescale.p=TRUE, simulate.p.value=TRUE)
> # Return p-value
> chisq_test <- chisq.test(x=countsn, p=countst, rescale.p=TRUE, simulate.p.value=TRUE)
> chisq_test$p.value
> # Observed frequencies from shifted normal data
> histp <- hist(rnorm(1e3, mean=2), breaks=100, plot=FALSE)
> countsn <- histp$counts/sum(histp$counts)
> # Theoretical frequencies
> countst <- rutils::diffit(pnorm(histp$breaks))
> # Perform Chi-squared test for shifted normal data
> chisq.test(x=countsn, p=countst, rescale.p=TRUE, simulate.p.value=TRUE)
> # Calculate histogram of VTI returns
> histp <- hist(retp, breaks=100, plot=FALSE)
> countsn <- histp$counts
> # Calculate cumulative probabilities and then difference them
> countst <- pt((histp$breaks-locv)/scalev, df=2)
> countst <- rutils::diffit(countst)
> # Perform Chi-squared test for VTI returns
> chisq.test(x=countsn, p=countst, rescale.p=TRUE, simulate.p.value=TRUE)
```

The Likelihood Function of Student's *t*-distribution

The non-standard Student's *t*-distribution is:

$$f(t) = \frac{\Gamma((\nu + 1)/2)}{\sqrt{\pi\nu} \sigma \Gamma(\nu/2)} \left(1 + \left(\frac{t - \mu}{\sigma}\right)^2 / \nu\right)^{-(\nu+1)/2}$$

It has non-zero mean equal to the location parameter μ , and a standard deviation proportional to the scale parameter σ .

The negative logarithm of the probability density is equal to:

$$-\log(f(t)) = -\log\left(\frac{\Gamma((\nu + 1)/2)}{\sqrt{\pi\nu} \Gamma(\nu/2)}\right) + \log(\sigma) + \frac{\nu + 1}{2} \log\left(1 + \left(\frac{t - \mu}{\sigma}\right)^2 / \nu\right)$$

The *likelihood* function $\mathcal{L}(\theta|\bar{x})$ is a function of the model parameters θ , given the observed values \bar{x} , under the model's probability distribution $f(x|\theta)$:

$$\mathcal{L}(\theta|x) = \prod_{i=1}^n f(x_i|\theta)$$

```
> # Objective function from function dt()
> likefun <- function(par, dfree, data) {
+   -sum(log(dt(x=(data-par[1])/par[2], df=dfree)/par[2]))
+ } # end likefun
> # Demonstrate equivalence with log(dt())
> likefun(c(1, 0.5), 2, 2:5)
> -sum(log(dt(x=2:5-1)/0.5, df=2)/0.5))
> # Objective function is negative log-likelihood
> likefun <- function(par, dfree, data) {
+   sum(-log(gamma((dfree+1)/2)/(sqrt(pi*dfree)*gamma(dfree/2))) +
+     log(par[2]) + (dfree+1)*2*log(1+((data-par[1])/par[2])^2/dfree)
+ } # end likefun
```

The *likelihood* function measures how *likely* are the parameters, given the observed values \bar{x} .

The *maximum-likelihood* estimate (MLE) of the parameters are those that maximize the *likelihood* function:

$$\theta_{MLE} = \arg \max_{\theta} \mathcal{L}(\theta|x)$$

In practice the logarithm of the *likelihood* $\log(\mathcal{L})$ is maximized, instead of the *likelihood* itself.

Fitting Asset Returns into Student's t -distribution

The function `fitdistr()` from package *MASS* fits a univariate distribution to a sample of data, by performing *maximum likelihood* optimization.

The function `fitdistr()` performs a *maximum likelihood* optimization to find the non-standardized Student's t -distribution location and scale parameters.

```
> # Calculate VTI percentage returns
> retp <- as.numeric(na.omit(rutils::etfenv$returns$VTI))
> # Fit VTI returns using MASS::fitdistr()
> fitobj <- MASS::fitdistr(retp, densfun="t", df=3)
> summary(fitobj)
> # Fitted parameters
> fitobj$estimate
> locv <- fitobj$estimate[1]
> scalev <- fitobj$estimate[2]
> locv; scalev
> # Standard errors of parameters
> fitobj$sd
> # Log-likelihood value
> fitobj$value
> # Fit distribution using optim()
> initp <- c(mean=0, scale=0.01) # Initial parameters
> fitobj <- optim(par=initp,
+   fn=likefun, # Log-likelihood function
+   data=retp,
+   dfree=3, # Degrees of freedom
+   method="L-BFGS-B", # Quasi-Newton method
+   upper=c(1, 0.1), # Upper constraint
+   lower=c(-1, 1e-7)) # Lower constraint
> # Optimal parameters
> locv <- fitobj$par["mean"]
> scalev <- fitobj$par["scale"]
> locv; scalev
```

The Student's t -distribution Fitted to Asset Returns

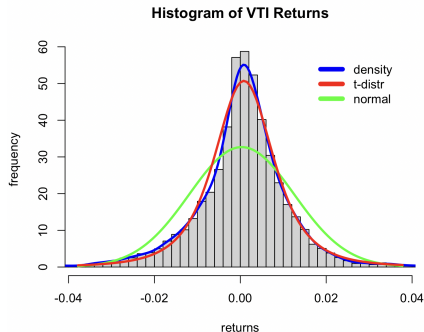
Asset returns typically exhibit *negative skewness* and *large kurtosis* (leptokurtosis), or fat tails.

Stock returns fit the non-standard t -distribution with 3 degrees of freedom quite well.

The function `hist()` calculates and plots a histogram, and returns its data *invisibly*.

The parameter `breaks` is the number of cells of the histogram.

```
> dev.new(width=6, height=5, noRStudioGD=TRUE)
> # x11(width=6, height=5)
> # Plot histogram of VTI returns
> madv <- mad(retp)
> histp <- hist(retp, col="lightgrey",
+   xlab="returns", breaks=100, xlim=c(-5*madv, 5*madv),
+   ylab="frequency", freq=FALSE, main="Histogram of VTI Returns")
> lines(density(retp, adjust=1.5), lwd=3, col="blue")
> # Plot the Normal probability distribution
> curve(expr=dnorm(x, mean=mean(retp),
+   sd=sd(retp)), add=TRUE, lwd=3, col="green")
> # Define non-standard t-distribution
> tdistr <- function(x, dfree, locv=0, scalev=1) {
+   dt((x-locv)/scalev, df=dfree)/scalev
+ } # end tdistr
> # Plot t-distribution function
> curve(expr=tdistr(x, dfree=3, locv=locv, scalev=scalev), col="red", lwd=3, add=TRUE)
> # Add legend
> legend("topright", inset=0.05, bty="n",
+   leg=c("density", "t-distr", "normal"),
+   lwd=6, lty=1, col=c("blue", "red", "green"))
```



Goodness of Fit of Student's *t*-distribution Fitted to Asset Returns

The Q-Q plot illustrates the relative distributions of two samples of data.

The Q-Q plot shows that stock returns fit the non-standard *t*-distribution with 3 degrees of freedom quite well.

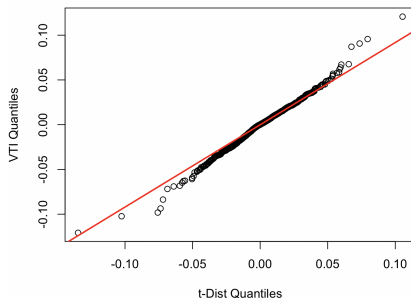
The function `qqplot()` produces a Q-Q plot for two samples of data.

The function `ks.test()` performs the *Kolmogorov-Smirnov* test for the similarity of two distributions.

The *null hypothesis* of the *Kolmogorov-Smirnov* test is that the two samples were obtained from the same probability distribution.

The *Kolmogorov-Smirnov* test rejects the *null hypothesis* that stock returns follow closely the non-standard *t*-distribution with 3 degrees of freedom.

Q-Q plot of VTI Returns vs Student's *t*-distribution



```
> # Calculate sample from non-standard t-distribution with df=3
> tdata <- scalev*rt(NROW(retp), df=3) + locv
> # Q-Q plot of VTI Returns vs non-standard t-distribution
> qqplot(tdata, retp, xlab="t-Dist Quantiles", ylab="VTI Quantiles"
+       main="Q-Q plot of VTI Returns vs Student's t-distribution")
> # Calculate quantiles of the distributions
> probs <- c(0.25, 0.75)
> qrets <- quantile(retp, probs)
> qtdata <- quantile(tdata, probs)
> # Calculate slope and plot line connecting quartiles
> slope <- diff(qrets)/diff(qtdata)
> intercept <- qrets[1]-slope*qtdata[1]
> abline(intercept,slope, lwd=2, col="red")
```

```
> # KS test for VTI returns vs t-distribution data
> ks.test(retp, tdata)
> # Define cumulative distribution of non-standard t-distribution
> pt distr <- function(x, dfree, locv=0, scalev=1) {
+   pt((x-locv)/scalev, df=dfree)
+ } # end pt distr
> # KS test for VTI returns vs cumulative t-distribution
> ks.test(sample(retp, replace=TRUE), pt distr, dfree=3, locv=locv, s
```

Leptokurtosis Fat Tails of Asset Returns

The probability under the *normal* distribution decreases exponentially for large values of x :

$$\phi(x) \propto e^{-x^2/2\sigma^2} \quad (\text{as } |x| \rightarrow \infty)$$

This is because a normal variable can be thought of as the sum of a large number of independent binomial variables of equal size.

So large values are produced only when all the contributing binomial variables are of the same sign, which is very improbable, so it produces extremely low tail probabilities (thin tails),

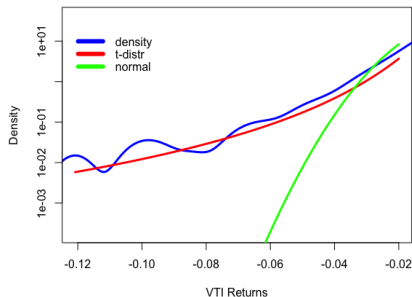
But in reality, the probability of large negative asset returns decreases much slower, as the negative power of the returns (fat tails).

The probability under Student's *t-distribution* decreases as a power for large values of x :

$$f(x) \propto |x|^{-(\nu+1)} \quad (\text{as } |x| \rightarrow \infty)$$

This is because a *t-variable* can be thought of as the sum of normal variables with different volatilities (different sizes).

Fat Left Tail of VTI Returns (density in log scale)



```
> # Plot log density of VTI returns
> plot(density(retp, adjust=4), xlab="VTI Returns", ylab="Density",
+      main="Fat Left Tail of VTI Returns (density in log scale)",
+      type="l", lwd=3, col="blue", xlim=c(min(retp), -0.02), log="y")
> # Plot t-distribution function
> curve(expr=dt((x-locv)/scalev, df=3)/scalev, lwd=3, col="red", add=TRUE)
> # Plot the Normal probability distribution
> curve(expr=dnorm(x, mean=mean(retp), sd=sd(retp)), lwd=3, col="green", add=TRUE)
> # Add legend
> legend("topleft", inset=0.01, bty="n", y.intersp=c(0.25, 0.25, 0.25),
+       leg=c("density", "t-distr", "normal"),
+       lwd=6, lty=1, col=c("blue", "red", "green"))
```

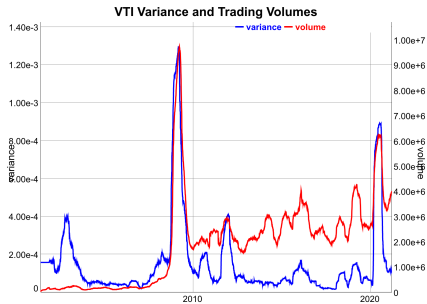
Trading Volumes

The average trading volumes have increased significantly since the 2008 crisis, mostly because of high frequency trading (HFT).

Higher levels of volatility coincide with higher *trading volumes*.

The time-dependent volatility of asset returns (*heteroskedasticity*) produces their fat tails (*leptokurtosis*).

```
> # Calculate VTI returns and trading volumes
> ohlc <- rutils::etfenv$VTI
> closep <- drop(coredata(quantmod::Cl(ohlc)))
> retp <- rutils::diffit(log(closep))
> volumv <- coredata(quantmod::Vo(ohlc))
> # Calculate trailing variance
> look_back <- 121
> varv <- HighFreq::roll_var_ohlc(log(ohlc), method="close", look_back=look_back, scale=FALSE)
> varv[1:look_back, ] <- varv[look_back+1, ]
> # Calculate trailing average volume
> volumr <- HighFreq::roll_var(volumv, look_back=look_back)/look_back
> # dygraph plot of VTI variance and trading volumes
> datav <- xts::xts(cbind(varv, volumr), zoo::index(ohlc))
> colnamev <- c("variance", "volume")
> colnames(datav) <- colnamev
> dygraphs::dygraph(datav, main="VTI Variance and Trading Volumes") %>%
+   dyAxis("y", label=colnamev[1], independentTicks=TRUE) %>%
+   dyAxis("y2", label=colnamev[2], independentTicks=TRUE) %>%
+   dySeries(name=colnamev[1], strokeWidth=2, axis="y", col="blue") %>%
+   dySeries(name=colnamev[2], strokeWidth=2, axis="y2", col="red") %>%
+   dyLegend(show="always", width=500)
```



Asset Returns in Trading Time

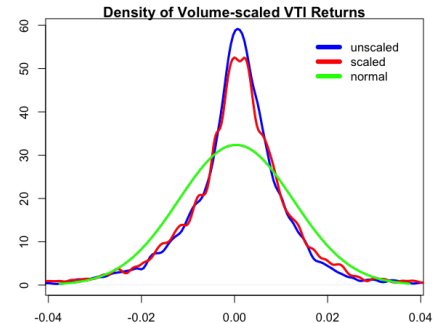
The time-dependent volatility of asset returns (*heteroskedasticity*) produces their fat tails (*leptokurtosis*).

If asset returns were measured at fixed intervals of *trading volumes* (*trading time* instead of clock time), then the volatility would be lower and less time-dependent.

The asset returns can be adjusted to *trading time* by dividing them by the *square root of the trading volumes*, to obtain scaled returns over equal trading volumes.

The scaled returns have a more positive *skewness* and a smaller *kurtosis* than unscaled returns.

```
> # Scale returns using volume (volume clock)
> retsc <- ifelse(volumv > 0, sqrt(volumv)*retp/sqrt(volumv), 0)
> retsc <- sd(retp)*retsc/sd(retsc)
> # retsc <- ifelse(volumv > 1e4, retp/volumv, 0)
> # Calculate moments of scaled returns
> nrow <- NROW(retp)
> sapply(list(retp=retp, retsc=retsc),
+   function(rets) {sapply(c(skew=3, kurt=4),
+     function(x) sum((rets/sd(rets))^x)/nrow)
+ }) # end sapply
```



```
> # x11(width=6, height=5)
> dev.new(width=6, height=5, noRStudioGD=TRUE)
> par(mar=c(3, 3, 2, 1), oma=c(1, 1, 1, 1))
> # Plot densities of SPY returns
> madv <- mad(retp)
> # bwidh <- mad(rutils::diffit(retp))
> plot(density(retp, bw=madv/10), xlim=c(-5*madv, 5*madv),
+   lwd=3, mgp=c(2, 1, 0), col="blue",
+   xlab="returns (standardized)", ylab="frequency",
+   main="Density of Volume-scaled VTI Returns")
> lines(density(retsc, bw=madv/10), lwd=3, col="red")
> curve(expr=dnorm(x, mean=mean(retp), sd=sd(retp)),
+   add=TRUE, lwd=3, col="green")
> # Add legend
> legend("topright", inset=0.05, bty="n"
```


Package *PerformanceAnalytics* for Risk and Performance Analysis

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

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

`managers` is an *xts* time series containing monthly percentage returns of six asset managers (HAM1 through HAM6), the EDHEC Long-Short Equity hedge fund index, the S&P 500, and US Treasury 10-year bond and 3-month bill total returns.

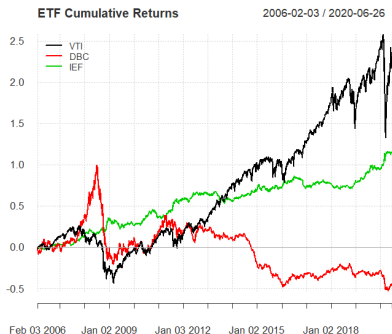
```
> # Load package PerformanceAnalytics
> library(PerformanceAnalytics)
> # Get documentation for package PerformanceAnalytics
> # Get short description
> packageDescription("PerformanceAnalytics")
> # Load help page
> help(package="PerformanceAnalytics")
> # List all objects in PerformanceAnalytics
> ls("package:PerformanceAnalytics")
> # List all datasets in PerformanceAnalytics
> data(package="PerformanceAnalytics")
> # Remove PerformanceAnalytics from search path
> detach("package:PerformanceAnalytics")
```

```
> perf_data <- unclass(data(
+   package="PerformanceAnalytics"))$results[, -(1:2)]
> apply(perf_data, 1, paste, collapse=" - ")
> # Load "managers" data set
> data(managers)
> class(managers)
> dim(managers)
> head(managers, 3)
```

Plots of Cumulative Returns

The function `chart.CumReturns()` from package *PerformanceAnalytics* plots the cumulative returns of a time series of returns.

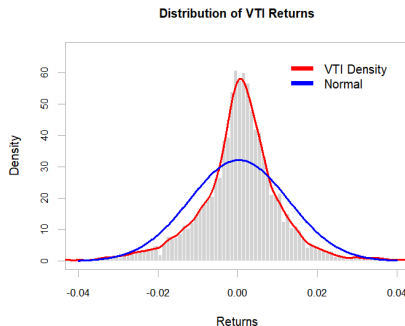
```
> # Load package "PerformanceAnalytics"
> library(PerformanceAnalytics)
> # Calculate ETF returns
> retp <- rutils::etfenv$returns[, c("VTI", "DBC", "IEF")]
> retp <- na.omit(retp)
> # Plot cumulative ETF returns
> x11(width=6, height=5)
> chart.CumReturns(retp, lwd=2, ylab="",
+   legend.loc="topleft", main="ETF Cumulative Returns")
```



The Distribution of Asset Returns

The function `chart.Histogram()` from package *PerformanceAnalytics* plots the histogram (frequency distribution) and the density of returns.

```
> retp <- na.omit(rutils::etfenv$returns$VTI)
> chart.Histogram(retp, xlim=c(-0.04, 0.04),
+   colorset = c("lightgray", "red", "blue"), lwd=3,
+   main=paste("Distribution of", colnames(retp), "Returns"),
+   methods = c("add.density", "add.normal"))
> legend("topright", inset=0.05, bty="n",
+   leg=c("VTI Density", "Normal"),
+   lwd=6, lty=1, col=c("red", "blue"))
```



The function `chart.Boxplot()` from package *PerformanceAnalytics* plots a box-and-whisker plot for a distribution of returns.

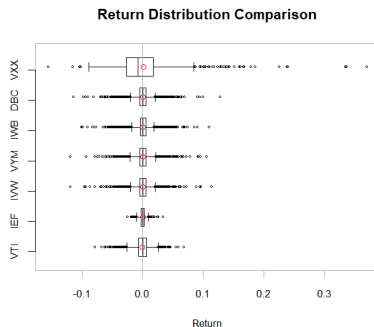
The function `chart.Boxplot()` is a wrapper and calls the function `graphics::boxplot()` to plot the box plots.

A *box plot* (box-and-whisker plot) is a graphical display of a distribution of data:

The *box* represents the upper and lower quartiles,

The vertical lines (whiskers) represent values beyond the quartiles.

Open circles represent values beyond the nominal range (outliers).



```
> retp <- rutils::etfenv$returns[,
+   c("VTII", "IEF", "IVW", "VYM", "IWB", "DBC", "VXX")]
> x11(width=6, height=5)
> chart.Boxplot(names=FALSE, retp)
> par(cex.lab=0.8, cex.axis=0.8)
> axis(side=2, at=(1:NCOL(retp))/7.5-0.05, labels=colnames(retp))
```

The Median Absolute Deviation Estimator of Dispersion

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

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

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

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.

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

```
> # Simulate normally distributed data
> 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
> bootd <- supply(1:10000, function(x) {
+   samplev <- datav[sample.int(nrows, replace=TRUE)]
+   c(sd=sd(samplev), mad=mad(samplev))
+ }) # end supply
> 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), stderror=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, 1:10000,
+   function(x, datav) {
+     samplev <- datav[sample.int(nrows, replace=TRUE)]
+     c(sd=sd(samplev), mad=mad(samplev))
+   }, datav=datav) # end parLapply
> # Parallel bootstrap under Mac-OSX or Linux
> bootd <- mclapply(1:10000, function(x) {
+   samplev <- datav[sample.int(nrows, replace=TRUE)]
+   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), stderror=sd(x)))
```

The Median Absolute Deviation of Asset Returns

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

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.

```
> # Calculate VTI returns
> retp <- na.omit(rutils::etfenv$returns$VTI)
> nrow <- NROW(retp)
> sd(retp)
> mad(retp)
> # Bootstrap of sd and mad estimators
> bootd <- sapply(1:10000, function(x) {
+   samplev <- retp[sample.int(nrow, replace=TRUE)]
+   c(sd=sd(samplev), mad=mad(samplev))
+ }) # end sapply
> bootd <- t(bootd)
> # Means and standard errors from bootstrap
> 100*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
> clusterExport(cluster, c("nrow", "returns"))
> bootd <- parLapply(cluster, 1:10000,
+   function(x) {
+     samplev <- retp[sample.int(nrow, replace=TRUE)]
+     c(sd=sd(samplev), mad=mad(samplev))
+   }) # end parLapply
> # Parallel bootstrap under Mac-OSX or Linux
> bootd <- mclapply(1:10000, 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
> 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)))
```

The Downside Deviation of Asset Returns

Some investors argue that positive returns don't represent risk, only those returns less than the target rate of return r_t .

The *Downside Deviation* (semi-deviation) σ_d is equal to the standard deviation of returns less than the target rate of return r_t :

$$\sigma_d = \sqrt{\frac{1}{n} \sum_{i=1}^n ([r_i - r_t]_-)^2}$$

The function `DownsideDeviation()` from package *PerformanceAnalytics* calculates the downside deviation, for either the full time series (`method="full"`) or only for the subseries less than the target rate of return r_t (`method="subset"`).

```
> library(PerformanceAnalytics)
> # Define target rate of return of 50 bps
> targetr <- 0.005
> # Calculate the full downside returns
> returns_sub <- (retp - targetr)
> returns_sub <- ifelse(returns_sub < 0, returns_sub, 0)
> nrows <- NROW(returns_sub)
> # Calculate the downside deviation
> all.equal(sqrt(sum(returns_sub^2)/nrows),
+   drop(DownsideDeviation(retp, MAR=targetr, method="full")))
> # Calculate the subset downside returns
> returns_sub <- (retp - targetr)
> returns_sub <- returns_sub[returns_sub < 0]
> nrows <- NROW(returns_sub)
> # Calculate the downside deviation
> all.equal(sqrt(sum(returns_sub^2)/nrows),
+   drop(DownsideDeviation(retp, MAR=targetr, method="subset")))
```

Drawdown Risk

The *drawdown* is the drop in prices from their historical peak, and is equal to the difference between the prices minus the cumulative maximum of the prices.

Drawdown risk determines the risk of liquidation due to stop loss limits.

```
> # Calculate time series of VTI drawdowns
> closep <- log(quantmod::Cl(rutils::etfenv$VTI))
> drawdns <- (closep - cummax(closep))
> # Extract the date index from the time series closep
> datev <- zoo::index(closep)
> # Calculate the maximum drawdown date and depth
> indexmin <- which.min(drawdns)
> datemin <- datev[indexmin]
> maxdd <- drawdns[datemin]
> # Calculate the drawdown start and end dates
> startd <- max(datev[(datev < datemin) & (drawdns == 0)])
> endd <- min(datev[(datev > datemin) & (drawdns == 0)])
> # dygraph plot of VTI drawdowns
> datav <- cbind(closep, drawdns)
> colnamev <- c("VTI", "Drawdowns")
> colnames(datav) <- colnamev
> dygraphs::dygraph(datav, main="VTI Drawdowns") %>%
+   dyAxis("y", label=colnamev[1], independentTicks=TRUE) %>%
+   dyAxis("y2", label=colnamev[2],
+     valueRange=(1.2*range(drawdns)+0.1), independentTicks=TRUE) %>%
+   dySeries(name=colnamev[1], axis="y", col="blue") %>%
+   dySeries(name=colnamev[2], axis="y2", col="red") %>%
+   dyEvent(startd, "start drawdown", col="blue") %>%
+   dyEvent(datemin, "max drawdown", col="red") %>%
+   dyEvent(endd, "end drawdown", col="green")
```



```
> # Plot VTI drawdowns using package quantmod
> plot_theme <- chart_theme()
> plot_theme$col$line.col <- c("blue")
> x11(width=6, height=5)
> quantmod::chart_Series(x=closep, name="VTI Drawdowns", theme=plot_theme)
> xval <- match(startd, datev)
> yval <- max(closep)
> abline(v=xval, col="blue")
> text(x=xval, y=0.95*yval, "start drawdown", col="blue", cex=0.9)
> xval <- match(datemin, datev)
> abline(v=xval, col="red")
> text(x=xval, y=0.9*yval, "max drawdown", col="red", cex=0.9)
> xval <- match(endd, datev)
> abline(v=xval, col="green")
> text(x=xval, y=0.85*yval, "end drawdown", col="green", cex=0.9)
```


Drawdown Risk Using PerformanceAnalytics::table.Drawdowns()

The function `table.Drawdowns()` from package *PerformanceAnalytics* calculates a data frame of drawdowns.

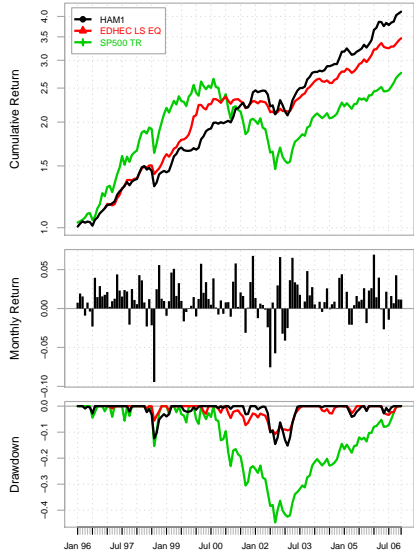
```
> library(xtable)
> library(PerformanceAnalytics)
> closep <- log(quantmod::Cl(rutils::etfenv$VTI))
> retp <- rutils::diffit(closep)
> # Calculate table of VTI drawdowns
> tablev <- PerformanceAnalytics::table.Drawdowns(retp, geometric=FALSE)
> # Convert dates to strings
> tablev <- cbind(sapply(tablev[, 1:3], as.character), tablev[, 4:7])
> # Print table of VTI drawdowns
> print(xtable(tablev), comment=FALSE, size="tiny", include.rownames=FALSE)
```

From	Trough	To	Depth	Length	To Trough	Recovery
2007-10-10	2009-03-09	2013-02-01	-0.57	1338.00	355.00	983.00
2020-02-20	2020-03-23	2020-08-21	-0.19	129.00	23.00	106.00
2022-01-04	2022-10-12		-0.12	419.00	195.00	
2018-09-21	2018-12-24	2019-07-02	-0.11	195.00	65.00	130.00
2015-06-24	2016-02-11	2016-07-18	-0.10	269.00	161.00	108.00

PerformanceSummary Plots

The function `charts.PerformanceSummary()` from package *PerformanceAnalytics* plots three charts: cumulative returns, return bars, and drawdowns, for time series of returns.

```
> data(managers)
> charts.PerformanceSummary(ham1,
+   main="", lwd=2, ylog=TRUE)
```



The Loss Distribution of Asset Returns

The distribution of returns has a long left tail of negative returns representing the risk of loss.

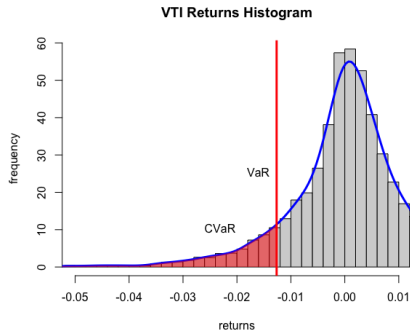
The *Value at Risk* (VaR) is equal to the quantile of returns corresponding to a given confidence level α .

The *Conditional Value at Risk* (CVaR) is equal to the average of negative returns less than the VaR.

The function `hist()` calculates and plots a histogram, and returns its data *invisibly*.

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

```
> # Calculate VTI percentage returns
> retp <- na.omit(rutils::etfenv$returns$VTI)
> confl <- 0.1
> varisk <- quantile(retp, confl)
> cvar <- mean(retp[retp <= varisk])
> # Plot histogram of VTI returns
> x11(width=6, height=5)
> par(mar=c(3, 2, 1, 0), oma=c(0, 0, 0, 0))
> histp <- hist(retp, col="lightgrey",
+   xlab="returns", ylab="frequency", breaks=100,
+   xlim=c(-0.05, 0.01), freq=FALSE, main="VTI Returns Histogram")
> # Calculate density
> densv <- density(retp, adjust=1.5)
```



```
> # Plot density
> lines(densv, lwd=3, col="blue")
> # Plot line for VaR
> abline(v=varisk, col="red", lwd=3)
> text(x=varisk, y=25, labels="VaR", lwd=2, pos=2)
> # Plot polygon shading for CVaR
> text(x=1.5*varisk, y=10, labels="CVaR", lwd=2, pos=2)
> varmax <- -0.06
> rangev <- (densv$x < varisk) & (densv$x > varmax)
> polygon(c(varmax, densv$x[rangev], varisk),
+   c(0, densv$y[rangev], 0), col=rgb(1, 0, 0, 0.5), border=NA)
```

Value at Risk (VaR)

The *Value at Risk* (VaR) is equal to the quantile of returns corresponding to a given confidence level α :

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

Where $f(r)$ is the probability density (distribution) of returns.

At a high confidence level, the value of VaR is subject to estimation error, and various numerical methods are used to approximate it.

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`.

A simpler but less accurate way of calculating the quantile is by sorting and selecting the data closest to the quantile.

The function `VaR()` from package *PerformanceAnalytics* calculates the *Value at Risk* using several different methods.

```
> # Calculate VTI percentage returns
> retp <- na.omit(rutils::etfenv$returns$VTI)
> nrow <- NROW(retp)
> confl <- 0.05
> # Calculate VaR approximately by sorting
> sortv <- sort(as.numeric(retp))
> cutoff <- round(confl*nrow)
> varisk <- sortv[cutoff]
> # Calculate VaR as quantile
> varisk <- quantile(retp, probs=confl)
> # PerformanceAnalytics VaR
> PerformanceAnalytics::VaR(retp, p=(1-confl), method="historical")
> all.equal(unname(varisk),
+   as.numeric(PerformanceAnalytics::VaR(retp,
+     p=(1-confl), method="historical")))
```

Conditional Value at Risk (CVaR)

The *Conditional Value at Risk* (CVaR) is equal to the average of negative returns less than the VaR:

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

The *Conditional Value at Risk* is also called the *Expected Shortfall* (ES), or the *Expected Tail Loss* (ETL).

The function `ETL()` from package *PerformanceAnalytics* calculates the *Conditional Value at Risk* using several different methods.

```
> # Calculate VaR as quantile
> varisk <- quantile(retp, confl)
> # Calculate CVaR as expected loss
> cvar <- mean(retp[retp <= varisk])
> # PerformanceAnalytics VaR
> PerformanceAnalytics::ETL(retp, p=(1-confl), method="historical")
> all.equal(unname(cvar),
+   as.numeric(PerformanceAnalytics::ETL(retp,
+     p=(1-confl), method="historical")))
```

Risk and Return Statistics

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

```
> # Calculate the risk-return statistics
> riskstats <-
+   PerformanceAnalytics::table.Stats(rutils::etfenv$returns)
> class(riskstats)
> # Transpose the data frame
> riskstats <- as.data.frame(t(riskstats))
> # Add Name column
> riskstats$Name <- rownames(riskstats)
> # Add Sharpe ratio column
> riskstats$Sharpe <- riskstats$"Arithmetic Mean"/riskstats$Stdev
> # Sort on Sharpe ratio
> riskstats <- riskstats[order(riskstats$Sharpe, decreasing=TRUE), ]
```

	Sharpe	Skewness	Kurtosis
QQQ	0.046	-0.507	6.62
USMV	0.041	-0.856	20.86
QUAL	0.035	-0.508	12.69
MTUM	0.033	-0.679	11.83
XLK	0.032	-0.217	9.83
IWF	0.030	-0.394	10.16
XLV	0.029	-0.324	11.12
IVW	0.028	-0.470	10.80
XLP	0.028	-0.428	11.71
GLD	0.026	-0.289	6.38
XLY	0.026	-0.562	8.53
IWB	0.025	-0.510	12.15
VTI	0.025	-0.477	12.28
XLI	0.023	-0.408	9.42
IVE	0.019	-0.547	11.81
XLU	0.018	-0.001	14.82
IWD	0.018	-0.470	13.44
VTV	0.018	-0.564	12.78
XLB	0.017	-0.372	7.98
VLUE	0.014	-0.956	16.41
XLE	0.014	-0.709	12.99
VYM	0.013	-0.496	11.56
SVXY	0.010	-18.060	649.08
AIEQ	0.010	-0.913	7.96
VNQ	0.006	-0.522	17.77
IEF	0.005	0.076	2.78
XLF	0.003	-0.439	18.18
TLT	0.002	0.009	3.73
VEU	0.001	-0.498	11.37
DBC	0.000	-0.503	3.25
EEM	-0.009	-24.012	944.15
VXX	-0.017	12.970	264.93
USO	-0.019	-1.144	14.21

Investor Risk and Return Preferences

Investors typically prefer larger *odd moments* of the return distribution (mean, skewness), and smaller *even moments* (variance, kurtosis).

But positive skewness is often associated with lower returns, which can be observed in the *VIX* volatility ETFs, *VXX* and *SVXY*.

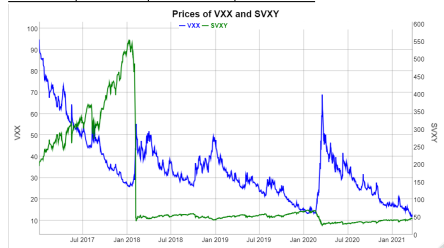
The *VXX* ETF is long the *VIX* index (effectively long an option), so it has positive skewness and small kurtosis, but negative returns (it's short market risk).

Since the *VXX* is effectively long an option, it pays option premiums so it has negative returns most of the time, with isolated periods of positive returns when markets drop.

The *SVXY* ETF is short the *VIX* index, so it has negative skewness and large kurtosis, but positive returns (it's long market risk).

Since the *SVXY* is effectively short an option, it earns option premiums so it has positive returns most of the time, but it suffers sharp losses when markets drop.

	Sharpe	Skewness	Kurtosis
VXX	-0.017	13.0	265
SVXY	0.010	-18.1	649



```
> # dygraph plot of VXX versus SVXY
> pricev <- na.omit(rutils::etfenv$pricev[, c("VXX", "SVXY")])
> pricev <- pricev["2017/"]
> colnamev <- c("VXX", "SVXY")
> colnames(pricev) <- colnamev
> dygraphs::dygraph(pricev, main="Prices of VXX and SVXY") %>%
+   dyAxis("y", label=colnamev[1], independentTicks=TRUE) %>%
+   dyAxis("y2", label=colnamev[2], independentTicks=TRUE) %>%
+   dySeries(name=colnamev[1], axis="y", strokeWidth=2, col="blue")
+   dySeries(name=colnamev[2], axis="y2", strokeWidth=2, col="green")
+   dyLegend(show="always", width=500) %>% dyLegend(show="always", width=500)
```

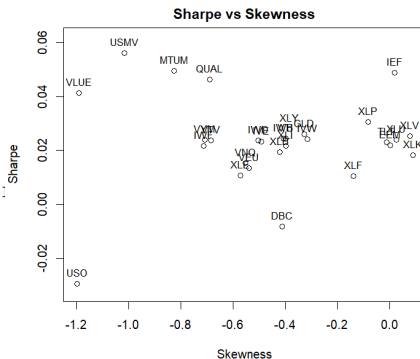
Skewness and Return Tradeoff

Similarly to the *VXX* and *SVXY*, for most other ETFs positive skewness is often associated with lower returns.

Some of the exceptions are bond ETFs (like *IEF*), which have both non-negative skewness and positive returns.

Another exception are commodity ETFs (like *USO* oil), which have both negative skewness and negative returns.

```
> # Remove VIX volatility ETF data
> riskstats <- riskstats[-match(c("VXX", "SVXY"), riskstats$Name), ]
> # Plot scatterplot of Sharpe vs Skewness
> plot(Sharpe ~ Skewness, data=riskstats,
+      ylim=1.1*range(riskstats$Sharpe),
+      main="Sharpe vs Skewness")
> # Add labels
> text(x=riskstats$Skewness, y=riskstats$Sharpe,
+      labels=riskstats$Name, pos=3, cex=0.8)
> # Plot scatterplot of Kurtosis vs Skewness
> x11(width=6, height=5)
> par(mar=c(4, 4, 2, 1), oma=c(0, 0, 0, 0))
> plot(Kurtosis ~ Skewness, data=riskstats,
+      ylim=c(1, max(riskstats$Kurtosis)),
+      main="Kurtosis vs Skewness")
> # Add labels
> text(x=riskstats$Skewness, y=riskstats$Kurtosis,
+      labels=riskstats$Name, pos=1, cex=0.8)
```



Risk-adjusted Return Measures

The *Sharpe ratio* S_r is equal to the excess returns (in excess of the risk-free rate r_f) divided by the standard deviation σ of the returns:

$$S_r = \frac{E[r - r_f]}{\sigma}$$

The *Sortino ratio* S_{Or} is equal to the excess returns divided by the *downside deviation* σ_d (standard deviation of returns that are less than a target rate of return r_t):

$$S_{Or} = \frac{E[r - r_t]}{\sigma_d}$$

The *Calmar ratio* C_r is equal to the excess returns divided by the *maximum drawdown* DD of the returns:

$$C_r = \frac{E[r - r_f]}{DD}$$

The *Dowd ratio* D_r is equal to the excess returns divided by the *Value at Risk* (VaR) of the returns:

$$D_r = \frac{E[r - r_f]}{VaR}$$

The *Conditional Dowd ratio* D_{c_r} is equal to the excess returns divided by the *Conditional Value at Risk* (CVaR) of the returns:

$$D_{c_r} = \frac{E[r - r_f]}{CVaR}$$

```
> library(PerformanceAnalytics)
> retp <- rutils::etfenv$returns[, c("VTI", "IEF")]
> retp <- na.omit(retp)
> # Calculate the Sharpe ratio
> confl <- 0.05
> PerformanceAnalytics::SharpeRatio(retp, p=(1-confl),
+   method="historical")
> # Calculate the Sortino ratio
> PerformanceAnalytics::SortinoRatio(retp)
> # Calculate the Calmar ratio
> PerformanceAnalytics::CalmarRatio(retp)
> # Calculate the Dowd ratio
> PerformanceAnalytics::SharpeRatio(retp, FUN="VaR",
+   p=(1-confl), method="historical")
> # Calculate the Dowd ratio from scratch
> varisk <- sapply(retp, quantile, probs=confl)
> ~sapply(retp, mean)/varisk
> # Calculate the Conditional Dowd ratio
> PerformanceAnalytics::SharpeRatio(retp, FUN="ES",
+   p=(1-confl), method="historical")
> # Calculate the Conditional Dowd ratio from scratch
> cvar <- sapply(retp, function(x) {
+   mean(x[x < quantile(x, confl)])
+ })
> ~sapply(retp, mean)/cvar
```

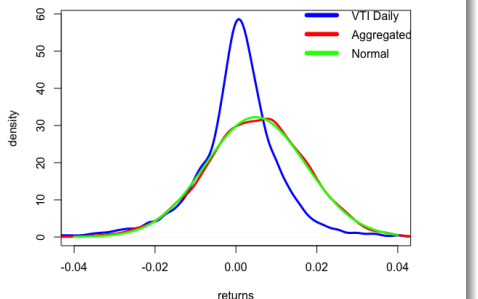
Risk of Aggregated Stock Returns

Stock returns aggregated over longer holding periods are closer to normally distributed, and their skewness, kurtosis, and tail risks are significantly lower than for daily returns.

Stocks become less risky over longer holding periods, so investors may choose to own a higher percentage of stocks, provided they hold them for a longer period of time.

```
> # Calculate VTI daily percentage returns
> retp <- na.omit(rutils::etfenv$returns$VTI)
> nrows <- NROW(retp)
> # Bootstrap aggregated annual VTI returns
> holdp <- 252
> reta <- sqrt(holdp)*sapply(1:nrows, function(x) {
+   mean(retp[sample.int(nrows, size=holdp, replace=TRUE)])
+ }) # end sapply
> # Calculate mean, standard deviation, skewness, and kurtosis
> datav <- cbind(retp, reta)
> colnames(datav) <- c("VTI", "Agg")
> sapply(datav, function(x) {
+   # Standardize the returns
+   meanv <- mean(x); stdev <- sd(x); x <- (x - meanv)/stdev
+   c(mean=meanv, stdev=stdev, skew=mean(x^3), kurt=mean(x^4))
+ }) # end sapply
> # Calculate the Sharpe and Dowd ratios
> confl <- 0.02
> ratiom <- sapply(datav, function(x) {
+   stdev <- sd(x)
+   varisk <- unname(quantile(x, probs=confl))
+   cvar <- mean(x[x < varisk])
+   mean(x)/c(Sharpe=stdev, Dowd=-varisk, DowdC=-cvar)
+ }) # end sapply
> # Annualize the daily risk
> ratiom[, 1] <- sqrt(252)*ratiom[, 1]
```

Distribution of Aggregated Stock Returns



```
> # Plot the densities of returns
> plot(density(retp), t="l", lwd=3, col="blue",
+   xlab="returns", ylab="density", xlim=c(-0.04, 0.04),
+   main="Distribution of Aggregated Stock Returns")
> lines(density(reta), t="l", col="red", lwd=3)
> curve(expr=dnorm(x, mean=mean(reta), sd=sd(reta)), col="green", lwd=3,
+   legend("topright", legend=c("VTI Daily", "Aggregated", "Normal"),
+   inset=-0.1, bg="white", lty=1, lwd=6, col=c("blue", "red", "green"))
```

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`,