

# Support notebook for Financial Risk Forecasting

This should serve to accompany the [Book Code](#) from Jon Danielsson's *Financial Risk Forecasting* and the Seminar notebooks for the class. All the code used is in R.

Thanks to Alvaro Aguirre for developing the notebook and material, and to Jia Rong Fan for the information on financial data sources.

## References

If you want further material, the following list provides information on various topics:

- [A very short introduction to R](#): Concise tutorial on the basics of R
- [The official introduction to R](#): Official documentation of the language which covers its usage
- [R Programming tutorial by freeCodeCamp.org](#): Complete video for learning R
- [R Markdown tutorial](#): Using R to create documents
- [R Markdown cheat sheet](#): Practical guide to use R Markdown
- [Dplyr cheat sheet](#): Guide on data wrangling using dplyr
- [ggplot2 tutorial](#): Complete tutorial on ggplot2

## Table of Contents

- [1 Financial data sources](#)
- [2 Control Flow](#)
- [3 Data frames handling](#)
- [4 Time Series, Dates and Zoo](#)
- [5 OLS and Maximum Likelihood](#)
- [6 GARCH Maximum Likelihood](#)
- [7 Laplace and Laplace usage](#)
- [8 VaR probabilities](#)
- [9 Simulated VaR with simple and continuous returns](#)
- [10 Statistical testing](#)
- [11 The Apply family of functions](#)
- [12 Introduction to dplyr](#)
- [13 dplyr](#)
- [14 R Markdown and Jupyter Notebooks](#)

## 1 Financial data sources

While financial markets generate a vast amount of data (sometimes with new observations every nanosecond), it is generally difficult and costly to get this data. Most data providers are commercial with complicated interfaces, and free data providers tend to have erratic access and errors in data.

Online data providers generally provide an *application programming interface* (API), an interface which allows users to submit data requests and receive responses. For most languages, there are often open-source or proprietary packages available which allow end-users to retrieve data using a relatively simple function call.

Today, many universities maintain subscriptions to data providers — this makes obtaining complete and accurate data much easier. Below we discuss some of the data providers that can be of use.

While this information is accurate at the time of writing (July 2020), given the rapid developments in the space, it is likely to become quickly outdated.

While we have used each of these data sources below, we have no specific recommendation.

### Common Sources of Financial Data

#### Bloomberg

One of the most ubiquitous data sources in finance is the [Bloomberg Terminal](#). Due to its pervasiveness throughout the industry, there are numerous packages in practically every language to access its APIs.

LSE students have access to a number of Bloomberg terminals in the library and the Masters students common rooms.

#### Wind

The [Wind Financial Terminal](#) (WFT) also provides market data like the Bloomberg Terminal, but with a specific focus on the Chinese financial markets. It supports APIs for MATLAB, R, C++ and Python, among others.

#### WRDS

The Wharton business school at University of Pennsylvania provides service called [Wharton Research Data Services](#) (WRDS) that many universities subscribe to. This provides a common interface to a number of databases, including CRSP and TAC high-frequency data. WRDS and many of its databases are available to LSE students and staff via wrds.wharton.upenn.edu.

#### CRSP

For historical US stock market data, one of the major sources is [The Center for Research in Security Prices](#) (CRSP, pronounced "crisp"), headquartered at the University of Chicago. LSE students can access CRSP via WRDS.

#### Yahoo Finance

The go-to place for many researchers requiring financial data has been [finance.yahoo.com](#). This data can be automatically downloaded into many software packages, including Matlab, R and Python.

There are three problems with Yahoo Finance.

1. Yahoo occasionally changes how the API works, requiring updates to software;
2. It often is unavailable for days or weeks;
3. There are errors in the data. For example, UK prices, quoted in pence by convention, sometimes appear in pounds for one or two days, reverting back to pence. On other occasions, numbers are simply wrong.

#### EOD Historical Data

[End of Day Historical Data](#) provides fundamental data API, live and end of day historical prices for stocks, ETFs and mutual funds from exchanges all around the world.

#### DBnomics

[DBnomics](#) is an open source data aggregator that provides access to dozens of data providers, including several countries' statistics.

#### Federal Reserve Economic Data (FRED)

The IFRED Economic Research ([https://fred.stlouisfed.org/](#)) is a good source for macroeconomic data including data for unemployment, GDP, interest rates, the money supply, etc. It can be accessed from DBnomics.

#### IEX

IEX provides access to US equity data via [https://iextrading.com/developer/](#).

#### ECB FX

The European Central Bank [Statistical Data Warehouse](#) and its corresponding [SDMX interface](#) allow for retrieval of daily Euro FX data.

The entire dataset is here [http://www.ecb.europa.eu/stats/eurofxref/eurofxref-hist.zip](#), and it can be accessed in R using

```
wget http://www.ecb.europa.eu/stats/eurofxref/eurofxref-hist.zip -O eurofxref-hist.zip
```

#### Alpha Vantage

[Alpha Vantage](#) provides free daily and realtime stock price data, and API access is available in both R and Python. Its data source appears to be the same as Yahoo Finance and hence subject to the same errors.

#### Quandl

[Quandl](#) provides a common API access for R and Python to a large number of commercial databases, and some that are free. While comprehensive, one may need to subscribe to data from a number of providers.

#### Fama-French Data Library:

The [Fama-French Data Library](#) provides a large amount of historical data compiled by Eugene Fama and Kenneth French. The data is updated regularly, and the Fama-French 3-factor data is especially useful for analyzing fund and portfolio performance.

## Preliminary Considerations

### Downloading and importing data

Data from this sources can either be downloaded directly, as a .csv or .xls file, or in most cases be imported into R using an API or package.

### Symbols and Names

All stocks that are trading in the market are associated with a *ticker symbol*, serving as an identifier for the specific security. These are often specific to the particular exchange or a country of listing. This can often lead to confusion or ambiguity when a company has cross-listings — for instance, the Japanese car manufacturer Toyota is listed as [7203](#) on the Tokyo Stock Exchange, [TXX](#) on the London Stock Exchange, and [TM](#) on the New York Stock Exchange.

Some identical securities have different names. Depending on the source, searching for data on the S&P-500 (a major American stock market index) might require the ticker names [SPX](#) (Bloomberg), [%SPC](#) (Yahoo Finance), [TSE](#) (Google Finance) and so on.

When searching for a particular security from the data source, it is good practice to verify that the ticker symbol used indeed corresponds to the correct data series by checking the data description.

### Tickers, ISIN and PERMNO

It is important to understand that Tickers can change over time (for example due to a merger or name change) and be recycled. This can create problems if not taken into consideration when querying financial data. There are other securities identifiers that do not change over time. The [ISIN](#) is the International Securities Identification Number, which is an internationally recognized 12-characters code that is unique for each stock. Unlike Tickers, the ISIN for a stock is the same regardless of the market. Another important identifier is the [PERMNO](#), which is the permanent issue identifier of the CRSP dataset.

### Dates

A date generally has three components, year (YYYY), month (MM) and day (DD). Depending on the granularity of data, there may also be a time component: hours, minutes, seconds and fractions of a second. Software packages often define some arbitrary date as origin. For example, R's origin date/time is set as January 1 1970 00:00:00, and all dates are relative to this date.

Dates often present problems when using multiple languages or programs to carry out numerical work. Excel has two different conventions for dates depending on version — the origin year can either be 1900 or 1904 — this requires verification before use. Furthermore, Excel does not allow dates before 1900.

A date can be represented numerically in many different ways. Consider the date 13 September 2018:

Format	Example
DD-MM-YYYY	13-09-2018
MM-DD-YYYY	09-13-2018
YYYY-MM-DD	2018-09-13
YYMMDD	20180913

The best way is to use the YYMMDD convention for two reasons:

1. It can be represented as an integer, not as a string, making data handling more convenient
2. It sorts naturally (in chronological order) even as integers

### Best practice in data downloading

Given the need to ensure that accurate data is used, it may be best to have a separate step in downloading and running the data. Get data from a provider, save it as a CSV (or Excel) file and then have a separate run for processing the data.

There are a few reasons why this is preferable:

1. Limits for API calls:
  - Many data providers place limits on the amount/speed of data downloads
2. Minimize time spent importing data
  - Especially true for large datasets, since CSV/xlsx reading is relatively quick
3. Quick checking for data omissions
  - Useful to have a pre-analysis check of data quality using plots and tables

Never manipulate your data in Excel. Perform any analysis in R or the programming language of your choice. This will make it tractable and can prevent mistakes.

## 2 Control Flow in R

In many cases we will want to repeat a piece of code several times, or execute part of our code only if some conditions are met.

Programming languages allow us to do that quite easily with loops and conditionals, which are specific statements that give instructions to computers. In this section we will cover the most important loops and conditionals in R.

For more details you can check [the official introduction to R](#) (Section

### For statements

A *for loop* evaluates a piece of code for different values. Usually, we think of *for loops* in two types, depending on the values we are evaluating: Loops over sequences, and loops over elements.

#### Looping over sequences

Probably the most common way to build a *for* statement is to loop over a numeric sequence. For example, if we wanted to `print()` the integers from 1 to 5:

```
In [1]: for (i in 1:5) {
# print(i)
}

[1] 1
[1] 2
[1] 3
[1] 4
[1] 5
```

In line 1, we are telling R that we want to create an integer sequence from 1 to 5, and use the name `i` to refer to values from this sequence. Then, line 2 is executed for every possible value specified in line 1.

Basically we are telling R:

1. Create a sequence from a to b
2. Take the first element of the sequence and assign it to "i"
3. Execute the code between {} using the current value of "i"
4. When done, check if there are more elements in the sequence
  - 4.1 If there are, take the next element and assign that value to "i", then go back to step 3
  - 4.2 If there are not, exit the loop

Recall the syntax `1:5` means "Create a sequence from 1 to 5", and if we don't specify the step, it takes 1 by default:

```
In [2]: 1:5

1 2 3 4 5
```

This is equivalent to doing:

```
In [3]: seq(1,5,1)

1 2 3 4 5
```

So our *for loop* could also be written as:

```
In [4]: for (i in seq(1,5,1)) {
# print(i)
}

[1] 1
[1] 2
[1] 3
[1] 4
[1] 5
```

It is a convention to create integer sequences using `first_element:last_element`.

The main usage of this type of loops is working with **indexes**. For example, consider we have a vector of 500 elements filled with `NA`:

```
In [5]: vec <- rep(NA, 500)

# Let's say we want to replace all odd elements with a 1, and keep the even elements as NA. We could build a for loop over the sequence
seq(1, 500, by = 2)
```

```
In [6]: for (i in seq(1, 500, by = 2)) { # Loop over i = 1, 3, 5, ...
vec[i] <- 1 # Replace the i-th element in vec for i
}
```

An interesting application for using *for loops* over sequences is constructing the Fibonacci sequence: 1, 1, 2, 3, 5, 8, 13, 21, ...

Recall the *i*-th term of the Fibonacci sequence is the sum of the *i*-1-th and *i*-2-th terms.

Let's say we want to create a vector with the first 50 terms of the Fibonacci sequence:

```
In [7]: fibonacci <- rep(NA, 50) # Create the empty vector
fibonacci[1:2] <- 1 # The first two values are 1
for (i in 3:50){ # Loop over the indexes 3 to 100
  fibonacci[i] <- fibonacci[i-1] + fibonacci[i-2] # Replace the i-th term by the sum of the two prev
  ious terms
}

In [8]: fibonacci

1 1 2 3 5 8 13 21 34 55 89 144 233 377 610 987 1597 2584 4181 6765 10946 17711 28657 46368
75025 121393 196418 317811 514229 832040 1346289 2178309 3524578 5702887 9227465 14930352 24157817
48058169 63245986 102334155 165580141 267914296 433949437 701408733 1134903170 1836311903 2971250793
7867152675 778742049 12586269025
```

We could put this in to a function that gives us the *i*-th fibonacci number:

```
In [9]: fibonacci <- function(N) {
  fibo <- rep(NA, N)
  fibo[1:2] <- 1
  for (i in 3:N){
    fibo[i] <- fibo[i-1] + fibo[i-2]
  }
  paste0("The ", N, "-th Fibonacci number is: ", fibo[N])
}

fibonacci(10)
fibonacci(50)

'The 10-th Fibonacci number is: 55'
'The 50-th Fibonacci number is: 12586269025'
```

Or, for example, if you have a vector with different model names:

```
In [10]: model <- c("GARCH", "EWMA", "HS")

# And another vector with a value associated to each model respectively:

In [11]: violation_ratio <- c(1.06, 0.95, 0.98)
```

We can output this as:

```
In [12]: for (i in 1:length(model)){
  cat("The value for the", model[i], "model is:", violation_ratio[i], "\n")
}

The value for the GARCH model is: 1.06
The value for the EWMA model is: 0.95
The value for the HS model is: 0.98
```

### Looping over elements

The other common way to think of *for loops* although technically analogous, is to loop over elements of a vector, which don't necessarily have to be numbers. The syntax is similar.

For example, consider the `model` I defined above:

```
In [13]: model

[1] "GARCH" "EWMA" "HS"
```

You can loop over the vector `model` by doing `for (element in model)`. Note that `element` is just a placeholder, it can be anything else. Let's say you want to say hello to everyone on the list:

```
In [14]: for (model_type in model){
  cat("My favorite model is: ", model_type, "\n")
}

My favorite model is GARCH
My favorite model is EWMA
My favorite model is HS

Another example:

In [15]: estimation_windows = c(300, 500, 1000, 2000) # Vector of different "windows"
n = 5000
rand <- rnorm(n) # Create a vector of random numbers of length 5k
for (window in estimation_windows){ # Loop over estimation windows
  # Get last window-length elements of rand
  tmp <- rand[(n-window+1):n]
  cat("Average of last", window, "elements: ", mean(tmp), "\n")
}

Average of last 300 elements: -0.05571649
Average of last 500 elements: -0.03766916
Average of last 1000 elements: -0.01084457
Average of last 2000 elements: 0.01729325
```

## If statements

A conditional, or *if*, statement, evaluates a piece of code that can be either `TRUE` or `FALSE` and takes a different action depending on the value. If the condition is `TRUE`, then it will execute the piece of code inside the brackets {}, if it is `FALSE`, it will not. You can then specify a different action to be performed if the condition is `FALSE`, which would go in an *else* statement. It can even include *else if* statements to evaluate another condition if the first one was not met.

A classic example:

```
In [16]: x = 5

if (x > 0) {
  print("Positive")
} else if (x < 0) {
  print("Negative")
} else {
  print("Zero")
}

[1] "Positive"
```

With this we can easily define a function that checks if a number is even:

```
In [17]: is_even <- function(x){
  if (x %% 2 == 0) {
    cat(x, "is even", "\n")
  } else {
    cat(x, "is odd", "\n")
  }
}

is_even(4)
is_even(7)

4 is even
7 is odd
```

## While loops

A while loop will evaluate a condition and as long as it is `TRUE`, it will execute a piece of code. **NOTE:** you need to make sure that the loop will not run indefinitely. To do this, at some point in the loop the condition needs to be set to `FALSE`.

A basic usage is:

```
In [18]: a <- 1

while (a < 6) {
  print(a)
  a = a + 1
}

[1] 1
[1] 2
[1] 3
[1] 4
[1] 5
```

Normally, you use *while loops* when the number of times the loop has to be executed is not pre-determined. For example, you can create a simple function that uses a *while loop* to tell you how many times you need to have a number before its absolute value is less than 0.5:

```
In [19]: less_half <- function(x){
  steps <- 0
  while (abs(x) > 0.5) {
    x = x/2
    steps = steps + 1
  }
  cat("Number of steps:", steps, "\n")
}

less_half(10)
less_half(0.2)
less_half(-500)

Number of steps: 5
Number of steps: 0
Number of steps: 10
```

## 3 Data frames handling

Objects in R can be of different classes. For example, we can have a vector, which is an ordered array of observations of data of the same type (numbers, characters, logicals). We can also have a matrix, which is a rectangular arrangement of elements of the same data type. You can check what class an object is by running `class(object)`.

A data frame is a two-dimensional structure in which each column contains values of one variable and each row contains one set of values, or "observation" from each column. It is perhaps the most common way of storing data in R and the one we will use most of the time.

One of the main advantages of a data frame in comparison to a matrix, is that each column can have a different data type. For example, you can have one column with numbers, one with text, one with dates, and one with logicals, whereas a matrix limits you to only one data type. Keep in mind that a data frame needs all its columns to be of the same length.

## Creating Data Frames

There are different ways to create a data frame. We will focus on three:

### 1. Reading a file

In finance it is very common to have a .csv file with the data you want to analyze. For example, this data could be the output of any of the financial data sources we mentioned above. We can import directly a .csv file into R as a `data.frame` object using the function `read.csv()`.

The file `stocks.csv` holds daily returns for 6 different companies for the years 2015 to 2019. It has been downloaded from the CRSP database. Let's import it:

```
In [20]: stocks <- read.csv("stocks.csv") # Importing data to a dataframe

# After importing a file, it is recommended to see how it looks like before starting to perform any type of analysis. For this, we can use the function head() to see the first elements of our object:

In [21]: head(stocks)

      MSFT      XOM      GE      JPM      INTC      C      date
0.006651827 0.004101577 -0.0083447205 0.004941789 0.001927142 0.002768165 2015-01-02
-0.009346543 -0.027743311 -0.0185265618 -0.031537108 -0.01340056 -0.030222284 2015-01-05
-0.014678200 -0.005330180 -0.0217804820 -0.026271083 -0.018812857 -0.035839635 2015-01-06
0.012624989 0.005303180 0.0004149139 0.001524837 0.020758053 0.009227297 2015-01-07
0.028989594 0.016507990 0.0119710801 0.022098988 0.018430117 0.014935902 2015-01-08
-0.008440521 -0.001410995 -0.014050240 -0.017539929 0.001906182 -0.022586158 2015-01-09
```

It is **very** important to check the documentation of the functions we are using to make sure we are doing things right. For example, we need to be careful if:

- The CSV file has headers or not - A *header* is the title of a column. We can have a csv file where the first row includes the column titles, or a csv with no titles
- The separator is a comma (,), a semicolon (;), a tab (\t) or another symbol
- We are getting data from a country where the decimal symbol is a comma (,) instead of a period (.)

By checking the documentation of a function, we will see how to deal with each of these specific cases and more. To do so, run `?read.csv`

```
In [22]: ?read.csv

Now let's check that our variable df is actually a data frame:

In [23]: class(stocks)

[1] "data.frame"
```

We can check the *structure* of a data frame using the function `str()`. This can be quite useful to see what every column of the data frame holds, and what type of data it has:

```
In [24]: str(stocks)

'data.frame':   1258 obs. of  7 variables:
 $ MSFT: num  0.00665 -0.00935 -0.01468 0.01262 0.02899 ...
 $ XOM : num  -0.009345 -0.027743 -0.0185265618 -0.031537108 -0.01340056 -0.030222284 2015-01-05 ...
 $ GE : num  -0.0083447205 -0.0217804820 -0.026271083 -0.018812857 -0.035839635 2015-01-06 ...
 $ JPM : num  0.00494 -0.03154 -0.02627 0.0152 0.0221 ...
 $ INTC: num  0.001931 -0.01134 -0.01881 0.02076 0.01843 ...
 $ C : num  0.00277 -0.03202 -0.03584 0.00923 0.01494 ...
 $ date: Factor w/ 1258 levels "2015-01-02","2015-01-05","2015-01-06",...: 1 2 3 4 5 6 7 8 9 10 ...

We see the data frame has 8 columns and 1258 rows. The structure function shows us the names of each column, along with the data type it holds, and some observations. We see there is a column called X that serves as an index, with variable type integer, a column for every stock, which holds numerical data, and finally a column for dates, which is of type factor.
We can also check parts of the structure with different functions:

In [25]: dim(stocks)
nrow(stocks)
ncol(stocks)
colnames(stocks)

1258 7

1258

7

'MSFT' 'XOM' 'GE' 'JPM' 'INTC' 'C' 'date'
```

### 2. Creating a Data Frame from scratch

In some cases you might want to create a data frame from a list of vectors. This can easily be done with the `data.frame()` function:

```
In [26]: df <- data.frame(col1 = 1:3,
# col2 = c("A", "B", "C"),
col3 = c(TRUE, TRUE, FALSE),
col4 = c(1,0, 2,2, 3,3))

You have
```



```
[40]: jpm <- stocks[, c("JPM", "date")]
class(jpm)
dim(jpm)
```

JPM	date
0.004941769	2015-01-02
-0.031537108	2015-01-05
-0.026271083	2015-01-06
0.001524837	2015-01-07
0.022099986	2015-01-08
-0.017539029	2015-01-09
'data.frame'	
1258	2

Note that if we had only taken the returns, the output would have been a vector and not a data frame:

```
In [41]: jpm2 <- stocks[, c("JPM", "date")]
head(jpm2)
class(jpm2)
```

0.00494176931922958 -0.0315371077982067 -0.0262710827191262 0.00152483684516575 0.0220999863448632  
-0.0175399291295206

'numeric'

What if we wanted to create another column called negative, that tells us if the returns for a particular day where negative?

We can easily add a new variable by using the syntax: `data_frame$new_variable`

```
In [42]: # We can use an ifelse to assign values conditional on another value
jpm$negative <- ifelse(jpm$JPM < 0, TRUE, FALSE)
head(jpm)
```

JPM	date	negative
0.004941769	2015-01-02	FALSE
-0.031537108	2015-01-05	TRUE
-0.026271083	2015-01-06	TRUE
0.001524837	2015-01-07	FALSE
0.022099986	2015-01-08	FALSE
-0.017539029	2015-01-09	TRUE

## Merging data frames

If we have datasets with a common column, we can use the `merge()` function and R will create a new dataset using that common column as a joiner. For example:

```
In [43]: MSFT <- stocks[, c("date", "MSFT")]
INTC <- stocks[, c("date", "INTC")]
head(MSFT)
head(INTC)
```

date	MSFT
2015-01-02	0.006651827
2015-01-05	-0.009346543
2015-01-06	-0.014678200
2015-01-07	0.012624969
2015-01-08	0.028989394
2015-01-09	-0.008440521

date	INTC
2015-01-02	0.001927142
2015-01-05	-0.011340056
2015-01-06	-0.018812857
2015-01-07	0.020758053
2015-01-08	0.018430117
2015-01-09	0.001906182

In case there are several common variables, we can specify how to join the data frames (for more information run `?merge`), but in a simple case as this, R will be smart enough to know what we mean:

```
In [44]: new_df <- merge(MSFT, INTC)
head(new_df)
```

date	MSFT	INTC
2015-01-02	0.006651827	0.001927142
2015-01-05	-0.009346543	-0.011340056
2015-01-06	-0.014678200	-0.018812857
2015-01-07	0.012624969	0.020758053
2015-01-08	0.028989394	0.018430117
2015-01-09	-0.008440521	0.001906182

## Reshaping

The package `reshape2` provides various useful functions to transform a dataset. Consider the following raw output from the CRSP database, including the same stocks' returns and time period:

```
In [45]: crsp <- read_csv("crsp.csv")
head(crsp)
```

PERMNO	date	TICKER	COMMON	PPC	RET	CFACPR
10107	20150102	MSFT	MICROSOFT CORP	46.76	0.006674	1
11850	20150102	XOM	EXXON MOBIL CORP	92.83	0.004410	1
12060	20150102	GE	GENERAL ELECTRIC CO	25.06	-0.006310	1
47896	20150102	JPM	JPMORGAN CHASE & CO	62.49	0.004954	1
59328	20150102	INTC	INTEL CORP	36.36	0.001929	1
75919	20150102	C	CITIGROUP INC	54.28	0.002772	1

If we want to create a data frame with the returns for every stock, we can use the `dcast()` function

```
In [46]: library(reshape2)
RET <- dcast(crsp, date ~ PERMNO, value.var = "RET")
head(RET)
```

date	10107	11850	12060	47896	59328	70519
20150102	0.006674	0.004410	-0.006310	0.004954	0.001929	0.002772
20150105	-0.009303	-0.027382	-0.018356	-0.031045	-0.011278	-0.031515
20150106	-0.014571	-0.005316	-0.021545	-0.025929	-0.018637	-0.035205
20150107	0.012705	0.010133	0.000415	0.001526	0.020875	0.009270
20150108	0.029418	0.016645	0.012043	0.022346	0.018601	0.015048
20150109	-0.008405	-0.001410	-0.013952	-0.017387	0.001908	-0.022333

Since we used the `PERMNO` as the variable to separate the `RET` observations, this is now the name of the columns. We can easily rename them:

```
In [47]: names(RET) <- c("date", "MSFT", "XOM", "GE", "JPM", "INTC", "C")
head(RET)
```

date	MSFT	XOM	GE	JPM	INTC	C
20150102	0.006674	0.004410	-0.006310	0.004954	0.001929	0.002772
20150105	-0.009303	-0.027382	-0.018356	-0.031045	-0.011278	-0.031515
20150106	-0.014571	-0.005316	-0.021545	-0.025929	-0.018637	-0.035205
20150107	0.012705	0.010133	0.000415	0.001526	0.020875	0.009270
20150108	0.029418	0.016645	0.012043	0.022346	0.018601	0.015048
20150109	-0.008405	-0.001410	-0.013952	-0.017387	0.001908	-0.022333

## Saving the data frame

Once we have finished cleaning and handling our dataset, we can save it so we don't have to repeat the procedure next time we want to use it. The most common file formats to save our data frame are either a `.RData` or a `.csv` file.

`.RData` files are specific to R and can store as many objects as you'd like within a single file. You can save your entire environment (all your variables) in a single file and then easily load it and continue working where you left it. To save our data frame in a `.RData` file we can run:

```
In [48]: save(RET, file = "RET.RData")
```

Then it can be directly loaded into R using `load("RET.RData")`.

To save your dataframe as a `.csv` file you can use the function `write.csv()`:

```
In [49]: write.csv(RET, file = "RET.csv")
```

## 4 | Time Series

Most of financial applications involve working with dates. It can be monthly, weekly, daily, or even intra-day data. Storing data as text is not helpful since we cannot order or subset it easily. R has a specific data type called `Date`. In this section we will explore some packages that will help us to work with `Date` objects.

### The lubridate package

It provides functions to work with date-times and time-spans. The full documentation is here: <https://cran.r-project.org/web/packages/lubridate/lubridate.pdf>

The main two functions we will use are `ymd()` and `dmy()`, which stand for year-month-day, and day-month-year respectively. These are quite powerful functions that transform a string object into date. As long as the object we pass on follow the structure of year-month-day for `ymd()`, the function will recognize it and transform it into date. For example:

```
In [50]: library(lubridate)
ymd("20200110")
class(ymd("20200110"))
```

Attaching package: 'lubridate'

The following object is masked from 'package:base':

date

2020-01-10

'Date'

```
In [51]: ymd("2015JAN11")
class(ymd("20200110"))
```

2015-01-11

'Date'

```
In [52]: ymd("04-MAR-25")
class(ymd("04MAR5"))
```

2004-03-05

'Date'

```
In [53]: dmy("1/june/2019")
class(dmy("1/june/2019"))
```

2019-06-01

'Date'

```
In [54]: dmy("28-december-14")
class(dmy("28-december-14"))
```

2014-12-28

'Date'

### The zoo package

This package functions to work with ordered indexed observations. It is very commonly used in time series analysis. The complete documentation is available here: <https://cran.r-project.org/web/packages/zoo/zoo.pdf>

```
In [55]: ret <- c(0.18, 0.02, -0.29, 0.00, 0.15)
class(ret)
```

'numeric'

```
In [56]: dates <- c("20190121", "20190122", "20190123", "20190124", "20190125")
class(dates)
```

'character'

```
In [57]: library(zoo)
ret.ts <- zoo(ret, order.by = ymd(dates))
class(ret.ts)
```

Attaching package: 'zoo'

The following objects are masked from 'package:base':

as.Date, as.Date.numeric

'zoo'

```
In [58]: ret.ts
```

2019-01-21 2019-01-22 2019-01-23 2019-01-24 2019-01-25  
0.18 0.02 -0.29 0.00 0.15

### Lag (and lead) functions

This function allows us to take the lag or leads of a time series object. The syntax is:

`lag(x, k, na.pad = F)` where:

- `x` = a time series object to lag
- `k` = number of lags (in units of observations); could be positive or negative
- `na.pad` = if TRUE, add NAs for missing observations

```
In [59]: ret.ts
lag(ret.ts, k = 2)
```

2019-01-21 2019-01-22 2019-01-23 2019-01-24 2019-01-25  
0.18 0.02 -0.29 0.00 0.15

2019-01-21 2019-01-22 2019-01-23  
-0.29 0.00 0.15

```
In [60]: ret.ts
lag(ret.ts, k = 2, na.pad = TRUE)
```

2019-01-21 2019-01-22 2019-01-23 2019-01-24 2019-01-25  
0.18 0.02 -0.29 0.00 0.15

2019-01-21 2019-01-22 2019-01-23 2019-01-24 2019-01-25  
-0.29 0.00 0.15 NA NA

```
In [61]: ret.ts
lag(ret.ts, k = -2)
```

2019-01-21 2019-01-22 2019-01-23 2019-01-24 2019-01-25  
0.18 0.02 -0.29 0.00 0.15

2019-01-23 2019-01-24 2019-01-25  
0.18 0.02 -0.29

### The diff function

Takes the lagged difference of a time series. Syntax:

`diff(x, lag, differences, na.pad = F)` where:

- `x` = a time series object
- `lag` = number of lags (in units of observations)
- `differences` = the order of the difference

```
In [62]: ret.ts
```

2019-01-21 2019-01-22 2019-01-23 2019-01-24 2019-01-25  
0.18 0.02 -0.29 0.00 0.15

```
In [63]: diff(ret.ts, lag = 1, na.pad = TRUE)
```

2019-01-21 2019-01-22 2019-01-23 2019-01-24 2019-01-25  
NA -0.16 -0.31 0.29 0.15

```
In [64]: diff(ret.ts, lag = 1, differences = 2, na.pad = TRUE)
```

2019-01-21 2019-01-22 2019-01-23 2019-01-24 2019-01-25  
NA NA -0.15 0.60 -0.14

```
In [65]: diff(ret.ts, lag = 2, differences = 1, na.pad = TRUE)
```

2019-01-21 2019-01-22 2019-01-23 2019-01-24 2019-01-25  
NA NA -0.47 -0.02 0.44

### The rollapply function

`rollapply` is a function for applying another function to rolling margins of an array. For example, you can construct a cumulative sum or mean. Syntax is:

`rollapply(data, width, FUN, fill, align...)` where:

- `data` = a time series object
- `width` = the window width
- `FUN` = the function to apply
- `fill` = NA to replace missing observations
- `align` = either "left", "right", or "center" (default)

```
In [66]: # Example: Sum the last two consecutive returns
ret.ts
```

2019-01-21 2019-01-22 2019-01-23 2019-01-24 2019-01-25  
0.18 0.02 -0.29 0.00 0.15

```
In [67]: rollapply(ret.ts, width = 2, FUN = sum, align = "left")
```

2019-01-21 2019-01-22 2019-01-23 2019-01-24 2019-01-25  
0.20 -0.27 -0.29 0.15

```
In [68]: rollapply(ret.ts, width = 2, FUN = sum, align = "right", fill = NA)
```

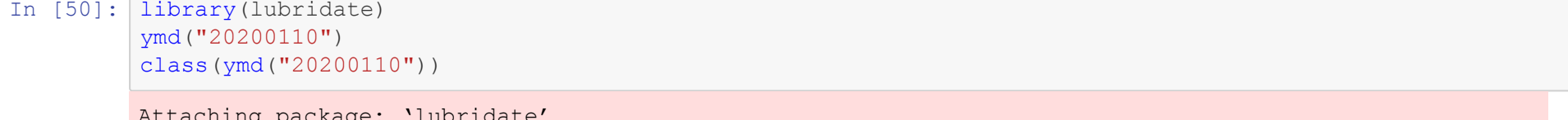
2019-01-21 2019-01-22 2019-01-23 2019-01-24 2019-01-25  
NA 0.20 -0.27 -0.29 0.15

## Plotting zoo objects

When plotting a zoo object, R will automatically use the dates as the x-axis. Let's create a series of fictitious returns to see this:

```
In [69]: ret <- rt(5000, df = 4) # Random draws from a Student t with 3 dof
dates <- seq(Sys.Date(), by = "-1 day", length.out = 2000) # 2000 days back from today
ret.ts <- zoo(ret, order.by = ymd(dates)) # Transform to zoo
```

```
In [70]: plot(ret.ts, main = "Returns from the past 2000 days",
xlab = "Date", ylab = "Returns", col = "mediumblue")
```



## Zooming into a time period

We can use the `window()` function to subset a zoo object to a given time period. For example, let's say we are interested in the returns of 1919. Since it is a zoo object, R will automatically adjust the x-axis to show months:

```
In [71]: sub.ret.ts <- window(ret.ts, start = "2019/1/1", end = "2019/12/31")
plot(sub.ret.ts, main = "Returns from 2019",
xlab = "Date", ylab = "Returns", col = "mediumblue")
```



## 5 | OLS and Maximum Likelihood

Consider the linear model:

$$y_i = \alpha + \beta_1 x_{i,1} + \beta_2 x_{i,2} + \dots + \beta_K x_{i,K} + \epsilon_i$$

Where  $E[\epsilon_i | X_i] = 0$

In matrix notation:

$$y = X\beta + \epsilon$$

Where  $X$  is a matrix containing one independent variable for each column,  $\beta$  is a vector of coefficients and  $\epsilon$  is a vector of errors.

The Ordinary Least Squares (OLS) estimator  $\hat{\beta}$  minimizes the sum of squared residuals:

$$\hat{\beta} = \arg \min \sum (y_i - \alpha - \beta_1 x_{i,1} - \beta_2 x_{i,2} - \dots - \beta_K x_{i,K})^2$$

Or in matrix notation:

$$\hat{\beta} = \arg \min (y - X\beta)'(y - X\beta)$$

Taking the first order conditions:

$$-X'y + yX + 2X'X\hat{\beta} = 0$$

Assuming invertibility of  $(X'X)$ , we have:

$$\hat{\beta} = (X'X)^{-1}X'y$$

Notation and definitions:

- Fitted values:  $\hat{y} = X\hat{\beta}$
- Residuals:  $\hat{\epsilon} = y - \hat{y}$
- Total Sum of Squares (SST):  $(y - \bar{y})'(y - \bar{y})$
- Sum of Squares of Residuals (SSR):  $\hat{\epsilon}'\hat{\epsilon}$
- R squared:  $1 - SSR/SST$

### OLS "by hand"

We can use matrix algebra in R to solve OLS. We will use the file "Shiller.csv". This dataset was developed by Robert Shiller and includes historical data on monthly stock prices, dividends, and earnings, consumer price index (to allow conversion to real values), all starting January 1871. More details on the dataset can be found here: <http://www.econ.yale.edu/~shiller/data.htm>

```
In [72]: # Data load and cleaning
df <- read.csv("Shiller.csv", header = T, na.strings = NA)
head(df)
```

```
df$date <- paste0(df$date, "01") # Add days to dates
# Create time series
ret.ts <- zoo(df$R, order.by = ymd(df$date))
PDIFF.ts <- diff(zoo(df$PD_r, order.by = ymd(df$date)), na.pad = T)
head(ret.ts)
head(PDIFF.ts)
```

```
# Adjust length
ret.ts <- ret.ts[2:length(ret.ts)]*100
PDIFF.ts <- PDIFF.ts[2:length(PDIFF.ts)]
head(ret.ts)
head(PDIFF.ts)
```

date	P	D	E	CPI	date.frac	long.int	P_r	D_r	E_r	P.E10	R_r	PD_r
187101	4.44	0.26	0.4	12.46	1871.04	5.32	77.17	4.52	6.85	NA	NA	17.07301
187102	4.50	0.26	0.4	12.84	1871.13	5.32	75.90	4.39	6.75	NA	-0.016457172	17.28929
187103	4.61	0.26	0.4	13.03	1871.21	5.33	76.62	4.32	6.65	NA	0.009486166	17.73611
187104	4.74	0.26	0.4	12.56	1871.29	5.33	81.76	4.48	6.90	NA	0.049168297	18.25001
187105												



```
[13]: # ARCH(1) MLE
LL_GARCH1 <- function(par, x) {
  omega <- par[1] # First element of par
  alpha <- par[2] # Second element of par
  T <- length(x) # Number of rows
  loglikelihood <- -(T-1)/2 * log(2*pi)
  # Loop (create the sum and add every instance to loglikelihood
  for(i in 2:T){
    loglikelihood = loglikelihood - 1/2 * (log(omega + alpha*(x[i-1]^2)*x[i]^2/(omega + alpha*x[i-1]^2)
  )
  }
  return(-loglikelihood) # output the log likelihood
}
res <- optim(c(0.1,0.5), LL_GARCH1, gr = NULL, ret, hessian = T) # Optimization
cat("Omega:", res$par[1], "\n",
  "Alpha:", res$par[2])

Omega: 0.9956736
Alpha: 0.3486696
```

## GARCH(1,1)

In a GARCH(1,1) model, the conditional volatility follows:

$$\sigma_t^2 = \omega + \alpha\gamma_{t-1}^2 + \beta\sigma_{t-1}^2$$

The conditional density is:

$$f(y_t|y_t) = \frac{1}{\sqrt{2\pi(\omega + \alpha\gamma_{t-1}^2 + \beta\sigma_{t-1}^2)}} \exp\left(-\frac{1}{2} \frac{y_t^2}{\omega + \alpha\gamma_{t-1}^2 + \beta\sigma_{t-1}^2}\right)$$

So the log-likelihood function is:

$$\log L = -\frac{T-1}{2} \log(2\pi) - \frac{1}{2} \sum_{t=2}^T \left( \log(\omega + \alpha\gamma_{t-1}^2 + \beta\sigma_{t-1}^2) + \frac{y_t^2}{\omega + \alpha\gamma_{t-1}^2 + \beta\sigma_{t-1}^2} \right)$$

Writing the function in R:

```
In [94]: # GARCH(1,1) MLE
LL_GARCH1_1 <- function(par, x) {
  omega <- par[1] # First element of par
  alpha <- par[2] # Second element of par
  beta <- par[3] # Third element of par
  T <- length(x) # Number of rows
  loglikelihood <- -(T-1)/2 * log(2*pi)
  sigma_2 <- rep(NA, T)
  sigma_2[1] <- var(x) # Initialize the sigma with unconditional volatility
  # Loop (create the sum and add every instance to loglikelihood
  for(i in 2:T){
    loglikelihood = loglikelihood - 1/2 * (log(omega + alpha*x[i-1]^2 + beta*sigma_2[i-1])
    + x[i]^2/(omega + alpha*x[i-1]^2 + beta*sigma_2[i-1]))
    sigma_2[i] <- omega + alpha*x[i-1]^2 + beta*sigma_2[i-1] # Update sigma
  }
  return(-loglikelihood) # output the log likelihood
}
res <- optim(c(0.1, 0, 0), LL_GARCH1_1, gr = NULL, ret, hessian = T,
  method = "L-BFGS-B", lower = c(0,0,0)) # Optimization
cat("Omega:", res$par[1], "\n",
  "Alpha:", res$par[2], "\n",
  "Beta:", res$par[3])
## Note that the values obtained in GARCH ML will depend on the optimisation method used

Omega: 0.01010722
Alpha: 0.1030158
Beta: 0.883244
```

Some important components of the optimization process, the first:

- Value of  $\sigma_t^2$ : We need to assign a value to  $\sigma_t^2$ , as a initialization of the conditional volatility. General, we choose the unconditional variance of the data for this
- Optimization method: We have chosen the L-BFGS-B algorithm, since it allows us to place bounds on the parameters. We have specified non-negativity by setting `lower = c(0,0,0)`, since this is desired in GARCH
- Initial values of the parameters: Optimization works best with. You need to assign values to initialize the parameters, making sure these are compliant with the restrictions of the function and the parameter bounds. This will have a small effect on the parameter values, but can have a large effect in computing time

## Optimization packages

In this section we will use the function `optim()` to fit the maximum likelihood estimation. This was enough for our purposes, but if you have to carry on more complex optimization problems, `optim()` can fail. We recommend the package [nlmin](#), which currently is the state of the art for minimizing functions in non-linear settings.

## 7 | rugarch and rmgarch usage

The packages that we use in the course for estimating GARCH models are `rugarch` for univariate models, and `rmgarch` for multivariate models. The complete documentation of the packages is:

- `rugarch`: <https://cran.r-project.org/web/packages/rugarch/rugarch.pdf>
- `rmgarch`: <https://cran.r-project.org/web/packages/rmgarch/rmgarch.pdf>

Both packages were developed by Alexios Galanos and are constantly maintained and updated. The development code can be found here: <https://github.com/alegalanos/>. If you are curious on learning more programming, you are encourage to browse the page to see how the package was created.

## Univariate GARCH models with rugarch

To estimate a univariate GARCH model you need to follow two steps:

1. Create an object of the `uGARCHspec` class which is the specification of the model you want to estimate. This includes the type of model (standard, asymmetric, power, etc), the GARCH order, the distribution, and model for the mean
2. Fit the specified model to the data

### 1. uGARCHspec()

First let's take a look at `uGARCHspec()`, which is the function to create an instance of the `uGARCHspec` class. The complete syntax with its default values is:

```
uGARCHspec(variance.model = list(model = "sGARCH", garchOrder = c(1, 1),
  submodel = NULL, external.regressors = NULL, variance.targeting = FALSE),
  mean.model = list(armaOrder = c(1, 1), include.mean = TRUE, arch = FALSE,
  archow = 1, arfima = FALSE, external.regressors = NULL, archex = FALSE),
  distribution.model = "norm", start.pars = list(), fixed.pars = list(), ...)
```

You can check the details of this function and what every argument means by running `?uGARCHspec`.

For the purposes of our course, we are going to focus on three arguments of the `uGARCHspec()` function:

#### variance.model

This argument takes a `list` with the specifications of the GARCH model. Its most important components are:

- `model`: Specify the type of model, currently implemented models are "sGARCH", "iGARCH", "eGARCH", "gjrGARCH", "apGARCH", and "GARCH" and "ccGARCH"
- `garchOrder`: Specify the ARCH(q) and GARCH(p) orders

Other components includes options to do `variance.targeting` or specify `external.regressors`.

#### mean.model

This argument takes a `list` with the specifications of the mean model, if assumed to have one. A traditional assumption is that the model has zero mean, in which case we specify `armaOrder = c(0,0)`, `include.mean = FALSE`. However, it is also common to assume that the mean follows a certain ARMA process.

#### distribution.model

Here you can specify the conditional density to use for innovations. The default is a normal distribution, but we can specify the use of a Student-t distribution by setting the value to `std`.

### 2. uGARCHfit()

Once the specification has been created, you can fit this to the data using `uGARCHfit(spec = ..., data = ...)`. The result will be an object of the class `uGARCHfit`, which is a list that contains useful information which will be shown below.

## GARCH(1,1)

Let's use `rugarch` to estimate a GARCH(1,1) model.

First, let's load the library, and the S&P500 data:

```
In [95]: library(rugarch)
sp <- read.csv("sp500.csv")
ret <- diff(log(Sp500))/100
ret <- ret[1:length(ret)]

Loading required package: parallel
Registered 83 method overwritten by 'xts':
 method from
as.xts.xts zoo

Attaching package: 'rugarch'

The following object is masked from 'package:stats':
  sigma
```

Now let's use `uGARCHspec()` to specify the type of model we want. In this case, we want to use a GARCH(1,1), so the `garchOrder` in `variance.model` will be set as `c(1,1)`. Note that even if we only have one component, we need to put it inside a `list`.

We will assume zero mean, so we need to include this into `mean.model`:

```
In [96]: spec1 <- uGARCHspec(
  variance.model = list(garchOrder = c(1,1)),
  mean.model = list(armaOrder = c(0,0), include.mean = FALSE)
```

We can call `spec1` to see what is inside:

```
In [97]: spec1
-----+-----
*      GARCH Model Spec      *
-----+-----
Conditional Variance Dynamics
-----+-----
GARCH Model      : sGARCH(1,1)
Variance Targeting : FALSE
Conditional Mean Dynamics
-----+-----
Mean Model       : ARFIMA(0,0,0)
Include Mean     : FALSE
GARCH-in-Mean    : FALSE
Conditional Distribution
-----+-----
Distribution      : norm
Includes Skew     : FALSE
Includes Shape    : FALSE
Includes Lambda   : FALSE

Now we can fit the specified model to our data using the uGARCHfit() function:
```

```
In [98]: garch1_1 <- uGARCHfit(spec = spec1, data = ret)
```

We can check the class of this new object:

```
In [99]: class(garch1_1)
[1] "uGARCHfit"
```

Objects from the `uGARCHfit` class have two slots. You can think of every slot as a list of components. The two slots are:

1. `fit`
2. `fit`

To access a slot you need to use the syntax: `object$slot`.

Let's explore these to understand better:

#### @model

The `model` slot includes all the information that was needed to estimate the GARCH model. This includes both the model specification and the data. To see every that is included in this slot, you can use the `names()` function:

```
In [100]: names(garch1_1$model)
[1] "modeldesc" "modeldesc" "modeldata" "pars" "start.pars" "fixed.pars" "maxOrder" "pos.matrix" "fmodel" "pidx" "n.start"
```

And you can access each element with the `$` sign. For example, let's see what is inside `pars`:

```
In [101]: garch1_1$model$pars
      Level Fixed Include Estimate      LB      UB
mu      0.00000000    0    0    0      NA      NA
ar       0.00000000    0    0    0      NA      NA
ma       0.00000000    0    0    0      NA      NA
arma0    0.00000000    0    0    0      NA      NA
archm0   0.00000000    0    0    0      NA      NA
mxrvg    0.00000000    0    0    0      NA      NA
omega    0.01802185    0    1    1  2.220446e-18  1454.814
alpha1   0.10279939    0    1    1  0.000000e+00    1.000
beta1    0.88332744    0    1    1  0.000000e+00    1.000
gamma     0.00000000    0    0    0      NA      NA
eta1     0.00000000    0    0    0      NA      NA
delta    0.00000000    0    0    0      NA      NA
deta2    0.00000000    0    0    0      NA      NA
lambda   0.00000000    0    0    0      NA      NA
kxvrg    0.00000000    0    0    0      NA      NA
skew      0.00000000    0    0    0      NA      NA
shape    0.00000000    0    0    0      NA      NA
ghlmbda  0.00000000    0    0    0      NA      NA
xi        0.00000000    0    0    0      NA      NA
```

We have a matrix with all the possible parameters to be used in fitting a GARCH model. We see there is a value of 1 for the parameters that were included in our GARCH specification (omega, alpha1 and beta1). The matrix also includes what are the lower and upper bounds for these parameters, which has nothing to do with our particular data, but only with what is expected from the model.

We can also see the model description:

```
In [102]: garch1_1$model$modeldesc
Sdistribution
'norm'
Sdistno
1
Svmodel
sGARCH
```

And inside `garch1_1$model$modeldata` we will find the data vector we have used when fitting the model. You can think of the `@model` slot as everything that R needs to know in order to be able to estimate the model.

#### @fit

Inside the `fit` slot we will find the *estimated* model, including the coefficients, likelihood, fitted conditional variance, and more. Let's check everything that is included:

```
In [103]: names(garch1_1$fit)
[1] "hessian" "cvar" "var" "sigma" "condH" "z" "LLH" "log.likelihoods" "residuals" "coef" "robust.cvar" "A" "B" "scores"
"timer" "lpars" "solver"
```

Let's see the coefficients, along with their standard errors, in `matcoef`:

Note: You will see that these values are very similar to what we found using Maximum Likelihood. Since it is an optimization problem, we can expect to see some small differences

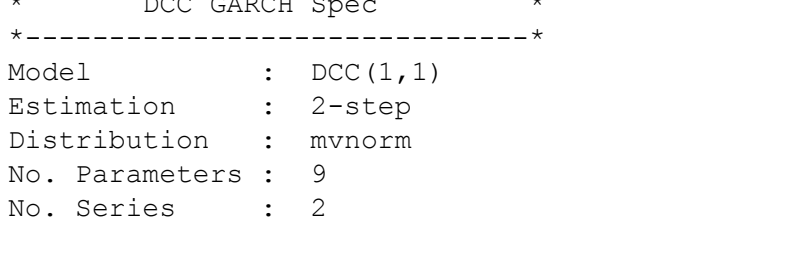
```
In [104]: garch1_1$fit$matcoef
      Estimate      Std. Error      t value      Pr(>|t|)
omega  0.01802185  0.00279939    6.45948  1.05987e-10
alpha1  0.10279939  0.00005806  11.298119  0.000000e+00
beta1   0.88332744  0.00972454   90.855407  0.000000e+00
```

We can also see the log-likelihood:

```
In [105]: garch1_1$fit$LLH
-6542.99177744059

This makes it easy to extract the estimated conditional volatility:
```

```
In [106]: plot(garch1_1$fit$sigma, type = "l")
```



We can also access the coefficients as likelihood using the following functions:

```
In [107]: coef(garch1_1)
      omega  0.0180218480714706
      alpha1 0.102799391800213
      beta1  0.8833274383512033
```

```
In [108]: likelihood(garch1_1)
-6542.99177744059
```

## Other specifications

Instead of using a conditionally normal GARCH(1,1) model specification, we can use a t-distribution or introduce asymmetric and power effects very easily.

Let's first fit a GARCH(1,1) using a Student-t as conditional distribution. We can see that the coefficients now include `shape`, which is the estimated degrees of freedom:

```
In [109]: # Student t
spec2 <- uGARCHspec(
  variance.model = list(garchOrder = c(1,1)),
  mean.model = list(armaOrder = c(0,0), include.mean = FALSE),
  distribution.model = "std")
tGARCH <- uGARCHfit(spec = spec2, data = ret)
tGARCH$fit$matcoef
      Estimate      Std. Error      t value      Pr(>|t|)
omega  0.00927054  0.002577062   3.597328  0.000215026
alpha1  0.09976745  0.010615749   9.398060  0.000000000
beta1   0.89872275  0.010206985  88.053343  0.000000000
shape  6.65154289  0.846320576  10.291389  0.000000000
```

Now let's fit an apGARCH model with fat tails. The parameter list will now include `gamma1` and `delta`, which are part of the apARCH model specification:

```
In [110]: # tapARCH model
spec3 <- uGARCHspec(
  variance.model = list(model = "apGARCH",
  mean.model = list(armaOrder = c(0,0), include.mean = FALSE),
  distribution.model = "std")
)
tapARCH <- uGARCHfit(spec = spec3, data = ret)
tapARCH$fit$matcoef
      Estimate      Std. Error      t value      Pr(>|t|)
omega  0.002231419  0.0015974472  13.96866  0
alpha1  0.086030565  0.0043436968  20.58583  0
beta1   0.90967493  0.0003901005  171.95797  0
gamma1  0.899897871  0.0003113951  3211.35291  0
delta   1.04334904  0.0904854633  11.28122  0
shape  7.72471637  0.8249865782  9.36458  0
```

## Issues

### Sample size

It is important to consider that the more parameter a model has, the more likely it is to run into estimation problems, and the more data we need. For example, as a rule of thumb, for a GARCH(1,1) we need at least 500 observations, while for a Student-t GARCH that minimum number is around 3,000 observations. In general, the more parameters are estimated, the more data is needed to be confident in the estimation.

### Optimization failure

`uGARCHfit()` essentially performs Maximum Likelihood to fit the specified model to our data. This optimization problem required a solver, which is the numerical method used to maximize the likelihood function. In some case, the default solver in `uGARCHfit()` can fail to converge. We recommend using the option `solver = "hybrid"` when fitting the model, since this will try different optimization methods in case the default one does not converge.

## Multivariate models with rmgarch

The package `rmgarch` allows to easily estimate models with multiple assets in the same fashion as `rugarch`. The procedure is analogous, we first need to specify the model we want to use, and then fit it to the data. However, specifying the model has extra steps. To explain this, we will focus on the estimation of a DCC model.

In a DCC model, we assume that each individual asset follows some type of univariate model, usually a GARCH. And then we model the correlation between the assets using an ARMA-like process.

The process for fitting a DCC model using `rmgarch` is then:

1. Specify the univariate volatility model each asset follows using `uGARCHspec()`
2. Use the `multispec()` function to create a multivariate specification. This is essentially a list of univariate specifications. If we are going to use the same for every asset, we can use `replicate()`
3. Then we need to create a `dccspec()` object, which takes in the list of univariate specifications for every asset, and the additional DCC part specifications like the `dccorder` and `distribution`
4. Fit the specification to the data

We will use the returns for JPM and Citigroup from the file `Y.Rdata` to go through this process:

```
In [111]: library(rmgarch)
load("RET.Rdata")
y <- RET[c("JPM", "C")]

We will assume a simple GARCH(1,1) with mean zero for each stock. We will create a single univariate specification, and then replicate it into multispec():
```

```
In [112]: # Create the univariate specification
uni_spec <- uGARCHspec(
  variance.model = list(garchOrder = c(1,1)),
  mean.model = list(armaOrder = c(0,0), include.mean = FALSE)
)
# Replicate it into a multispec element
mspec <- multispec(replicate(2, uni_spec))
```

Let's take a look inside the `mspec` object. We will see there are two univariate specifications, one per asset, and they are equal:

```
In [113]: mspec
-----+-----
*      GARCH Multi-Spec      *
-----+-----
Multiple Specifications : 2
Multi-Spec Type        : equal

You can check the specifications with mspec$spec
```

Now we proceed to create the specification for the DCC model using `dccspec()`:

```
In [114]: spec <- dccspec(
  # Univariate specifications - Needs to be multispec
  # DCC specification. We will assume an ARMA(1,1)-like process
  # Distribution, here multivariate normal
  distribution = "mvnorm"
)
```

We can call `spec` to see what is inside:

```
In [115]: spec
-----+-----
*      DCC GARCH Spec      *
-----+-----
Model      : DCC(1,1)
Estimation : 2-step
Distribution : mvnorm
No. Parameters : 9
No. Series : 2

We can again see more details in spec@model and spec@model$.
```

Now we can proceed to fit the specification to the data:

```
In [116]: res <- dccfit(spec, data = y)
res
-----+-----
*      DCC GARCH Fit      *
-----+-----
Distribution      : mvnorm
Model            : DCC(1,1)
No. Parameters    : 35
[VAR GARCH DCC UncQ] : 0+18+2+15]
No. Obs.          : 1258
Log-Likelihood    : 8156.696
Av.Log-Likelihood : 6.48

Optimal Parameters
-----+-----
Estimate      Std. Error      t value      Pr(>|t|)
[MSFT].omega  0.000033  0.000013  2.49435  0.012619
[MSFT].alpha1 0.000029  0.000016  1.8689  0.06633
[MSFT].beta1  0.141068  0.062791  2.2466  0.02463
[JPM].omega    0.000010  0.000001  0.03527  0.97495
[JPM].alpha1  0.000021  0.000013  1.6703  0.09464
[JPM].beta1    0.118458  0.056934  2.0799  0.037533
[C].omega      0.000026  0.000001  0.03527  0.97495
[C].alpha1     0.145713  0.06479  2.2466  0.02463
[C].beta1      0.118458  0.056934  2.0799  0.037533
[Joint]dcorr  0.100491  0.028630  3.5100  0.00048
[Joint]dcorr  0.685564  0.096975  7.1560  0.000000
```

Information Criteria
-----+-----
Akaike -12.953
Bayes -36.749
Shibata -12.954
Hannan-Quinn -12.940

Elapsed time : 0.933454

We can check what slots are inside:

```
In [117]: names(res@model)
[1] "modeldesc" "modeldesc" "modeldata" "varmodel" "pars" "start.pars" "fixed.pars" "maxgarchOrder" "maxdccOrder"
"pos.matrix" "pidx" "DCC" "mu" "residuals" "sigma" "lpars" "lpars" "midx" "idx" "model"
"convergence" "Nbar" "Qbar" "pik"
```

```
In [118]: # Coefficient matrix
res$fit$matcoef
      Estimate      Std. Error      t value      Pr(>|t|)
[JPM].omega  2.830088e-05  1.567794e-05  1.869924  6.163342e-02
[JPM].alpha1  1.410691e-01  6.279070e-02  2.246641  2.469300e-02
[JPM].beta1  6.913301e-01  1.265009e-01  5.101699  3.961728e-07
[C].omega    2.118301e-05  1.358032e-05  1.670279  9.486425e-02
[C].alpha1  1.184580e-01  5.89379e-02  2.078996  4.000000e-02
[C].beta1    7.952415e-01  9.828105e-02  8.091504  6.661338e-16
[Joint]dcorr1 1.004098e-01  2.862958e-02  3.510034  4.404007e-04
[Joint]dcorr1 8.855639e-01  9.067502e-02  7.560648  4.019070e-14

# Log Likelihood
res$fit$llh
8156.69553672898
```

The matrix `H` inside `res$fit` includes the covariances. It is 3-dimensional, since it includes the 2x2 covariance matrix for each of the T time periods:

```
In [120]: H <- res$fit$H
dim(H)
2 2 1258
```

```
In [121]: # First period's covariances
H[1,1]
0.0001754897 0.0001841418
0.0001841418 0.0002444920

We can extract the conditional correlation in two ways. One is computing it from H:
```

```
In [122]: # Initializing the vector
rhoDCC <- vector(length = dim(y)[1])
# Populate with the correlations
rhoDCC[1] <- H[1,2] / sqrt(H[1,1]*H[2,2])

Or we can directly extract it from res$fit$H:
```

```
In [123]: # Initializing the vector
rhoDCC2 <- vector(length = dim(y)[1])
for (i in 1:dim(y)[1]) {
  rhoDCC2[i] <- rcor(res)[i,2,i]
}
```

```
In [124]: plot(rhoDCC2, type = "l")
```



Now let's specify the VaR probability:

```
In [126]: load("RET.Rdata")
jpm <- RET.JPM
head(jpm)
0.004954 -0.031045 -0.025929 0.001526 0.022346 -0.017387
```

```
In [127]: # Plot the returns
plot(jpm, type = "l")
```

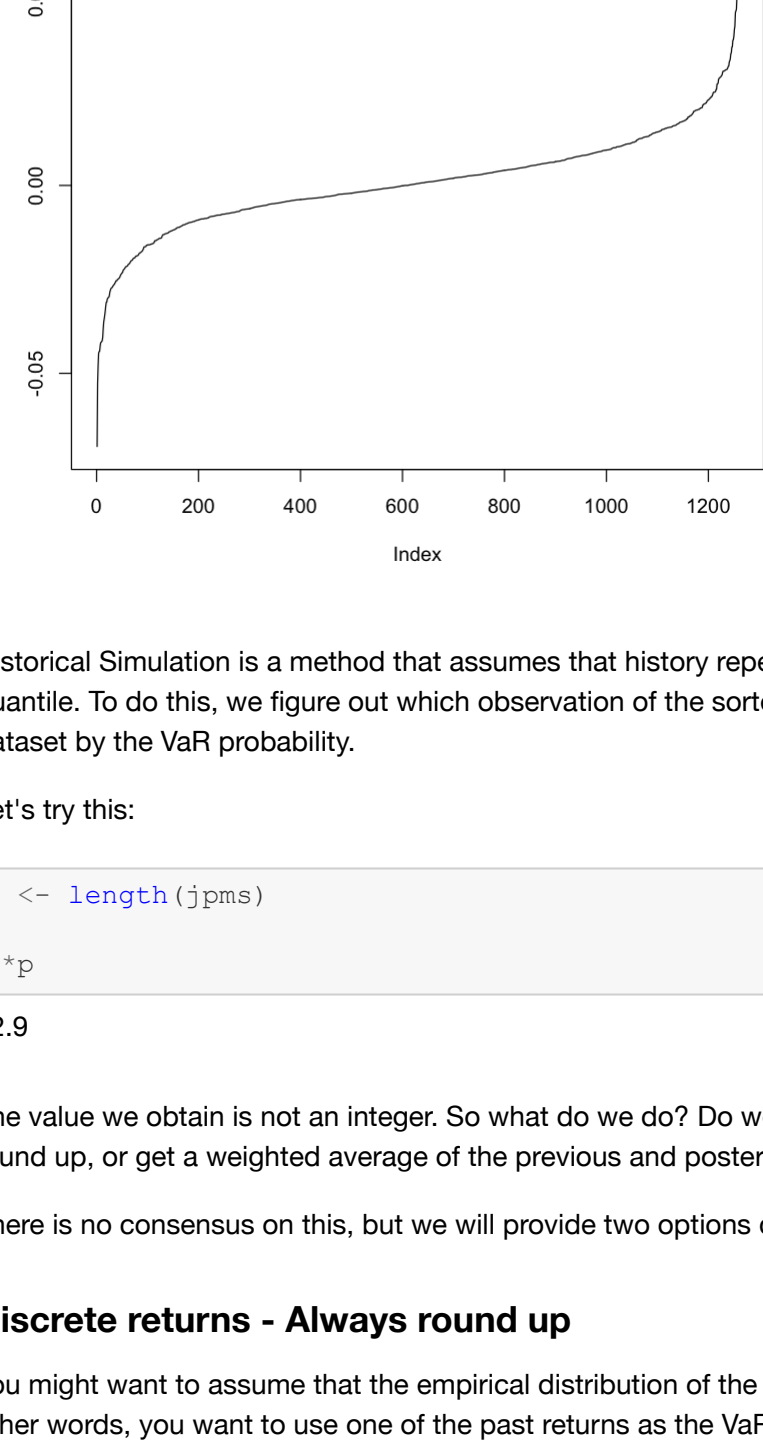




```
In [129]: p <- 0.05

Let's sort the values using sort() :
```

```
In [129]: jpmes <- sort(jpmes)
plot(jpmes, type = "l")
```



Historical Simulation is a method that assumes that history repeats itself. It uses the empirical distribution of the returns to find the p-th quantile. To do this, we figure out which observation of the sorted dataset p-th one. Practically, this involves multiplying the length of the dataset by the VaR probability.

Let's try this:

```
In [130]: n <- length(jpmes)
n*p
#p
62.9
```

The value we obtain is not an integer. So what do we do? Do we round this number to the closest integer, always round down, always round up, or get a weighted average of the previous and posterior observations?

There is no consensus on this, but we will provide two options depending on your assumption on the empirical distribution of the returns.

### Discrete returns - Always round up

You might want to assume that the empirical distribution of the returns is discrete. Meaning that you don't want to average observations. In other words, you want to use one of the past returns as the VaR estimation. In this case, you always need to **round up**. Imagine a case where  $n \cdot p = 80.1$ . Intuitively, it seems obvious that the VaR should round down to 80 and use the 80th observation. However, this would not be technically correct, because the quantile of the 80th observation is actually less than the required VaR probability  $p$ . So you would not be meeting the required probability. In many cases, VaR is computed as part of the financial regulation system, so it is quite important to meet the required probability. Using the 80th observation in this case might have given us the 4.9% VaR, whereas using the 81st observation might have given us the 5.6% VaR, but the latter does meet the required probability, while the first does not.

### Continuous returns - Weighted average

If you want to treat the empirical distribution of returns as continuous, you can compute the weighted average between the previous and posterior observations. For example, if  $n \cdot p = 80.1$ , you would use the  $0.9 \cdot (80th\ observation) + 0.1 \cdot (81st\ observation)$ .

For our purposes, we believe it is more intuitive to strictly use a past observation as the VaR estimation, so we will always round up. This can be easily done using the `ceiling()` function. Completing the HS estimation in our case:

```
In [131]: ceiling(n*p)
#p
63
```

```
In [132]: # VaR(5%)
jpmes[ceiling(n*p)]
#p
-0.020952
```

## 9 | Simulated VaR with simple and continuous returns

When working with simulations, we can choose if we want to use simple or continuous returns to simulate future prices.

### Simple returns

If we want to work with simple returns, we calculate one-day future prices in the following way:

1. Simulate one-day return:  $R_{t+1} \sim \mathcal{N}(0, \sigma^2)$
2. Calculate one-day future price:  $P_{t+1} = P_t \times (1 + R_{t+1})$

### Continuous returns

If we want to work with continuous returns, we calculate one-day future prices in the following way:

1. Simulate one-day return:  $y_{t+1} \sim \mathcal{N}(0, \sigma^2)$
2. Calculate one-day future price:  $P_{t+1} = P_t e^{(1/365) \times y_{t+1} \times e^{1/365}}$

### Comparison for VaR 5%

```
In [133]: # Simple returns
p <- 0.05
S <- 1e5
B <- 100
sigma <- 0.01
ret <- rnorm(S, 0, sigma^2)
Feim_simple <- P*(1+ret)
# VaR 5%
Feim_simple <- sort(Feim_simple - P)
Fe_simple[p*S]
#p
-0.0163162893227877
```

```
In [134]: # Compound returns
r <- 0.03 # Assuming risk free rate
Feim_comp <- P*exp(r*(1/365))*exp(ret)*exp(-0.5*sigma^2)
# VaR 5%
Feim_comp <- sort(Feim_comp - P)
Fe_comp[p*S]
#p
-0.013096236064345
```

## 10 | Statistical Testing

### Jarque-Bera

The Jarque-Bera (JB) test uses the third and fourth central moments of the sample data to check if it has a skewness and kurtosis matching a normal distribution. If the sample data comes from a normal distribution, the JB statistic has an asymptotic chi-square distribution with two degrees of freedom. So:

$$H_0: \text{The data is normally distributed}$$

$$H_1: \text{The data is NOT normally distributed}$$

We will generally use the JB test to check if our data follows a normal distribution. The function `jarque.bera.test()` from the `teststat` package will evaluate the data to compute the JB statistic and p-value. A small p-value will mean there is enough evidence to reject the null hypothesis that the third and fourth central moments of the sample data resemble the ones of a normal distribution.

Let's do a test on the normality of the S&P500 returns:

```
In [135]: library(teststat)
sp <- read.csv("sp500.csv")
ret <- diff(log(sp$price))
jarque.bera.test(ret)

Registered 93 method overwritten by 'quantmod':
  method from
as.zoo.data.frame zoo

Jarque Bera Test

data: ret
X-squared = 14521, df = 2, p-value < 2.2e-16
```

The p-value presented is the inverse of the JB statistic under the CDF of the asymptotic distribution. The p-value of this test is virtually zero, which means that we have enough evidence to reject the null hypothesis that our data has been drawn from a Normal distribution.

### Ljung-Box

The Ljung-Box (LB) test checks if the presented data exhibit serial correlation. It has as null-hypothesis that the data are independently distributed. The test statistic includes the sample size, and a sum of sample autocorrelations squared up to a specified lag, scaled by the sample size minus the lag. Under the null hypothesis, the statistic asymptotically follows a chi-square distribution with degrees of freedom equal to the number of lags being tested.

$$H_0: \text{The data exhibits serial correlation}$$

$$H_1: \text{The data does NOT exhibit serial correlation}$$

This type of test is commonly used for the residuals of ARMA models (Note: not used for the actuals values of the ARMA series). For our purposes, we normally use the LB test on returns and returns squared. Remember that it makes sense to do GARCH modeling when volatility exhibits autocorrelation, so we expect to see a rejection of the LB null hypothesis when applying the test to returns squared, even if we cannot reject the null when applied to returns.

We can implement the test using the `Box.test()` function and specifying the `type = "Ljung-Box"`:

```
In [136]: Box.test(ret, type = "Ljung-Box")
Box.test(ret~2, type = "Ljung-Box")

Box-Ljung test

data: ret
X-squared = 26.212, df = 1, p-value = 3.059e-07

Box-Ljung test

data: ret^2
X-squared = 212.66, df = 1, p-value < 2.2e-16
```

### Likelihood Ratio test

The Likelihood-Ratio (LR) test compares two models based on the ratio of their likelihoods, where one of the models was fitted over a restricted parameter space compared to the other one, so we can say it is nested. When we are working with log-likelihoods, we can naturally express the test statistic as a difference instead of a ratio.

Intuitively, you can think of the LR test as a check if a restriction imposed on a model makes a significant difference or not. The null hypothesis is that there is essentially no difference between the models. The LR statistic is asymptotically distributed as a chi-square with degrees of freedom equal to the difference in dimensionality in the parameter spaces that each model used in their maximum likelihood estimation.

$$H_0: \text{The two models are the same (the parameters excluded are jointly equal to zero)}$$

$$H_1: \text{The two models are not the same}$$

Let's give an application of this to make this clearer. Assume you want to build a model for the returns of the S&P500, and you are considering an ARCH(1) and a GARCH(1,1) model. Essentially, an ARCH(1) is a specific case of the GARCH family, in particular, it is a GARCH(1,0). We can see the ARCH maximization process as exactly the same, except that we are narrowing down the parameter space, since we are forcing  $\beta = 0$ . Hence, there is **one less dimension** in the ARCH parameter space MLE than in the GARCH MLE.

We can then perform a LR test to see if there is a gain in using a GARCH(1,1) instead of an ARCH(1) model.

Let's fit both models and build the LR test ourselves:

```
In [137]: library(rugarch)

# ARCH
spec_arch <- ugarchspec(
  variance.model = list(garchOrder = c(1,0,0)),
  mean.model = list(armaOrder = c(0,0), include.mean = FALSE)
)
ARCH <- ugarchfit(spec = spec_arch, data = ret, solver = "hybrid")

# GARCH
spec_garch <- ugarchspec(
  variance.model = list(garchOrder = c(1,1,1)),
  mean.model = list(armaOrder = c(0,0), include.mean = FALSE)
)
GARCH <- ugarchfit(spec = spec_garch, data = ret, solver = "hybrid")

# Perform the LR test
LR_statistic <- 2*(likelihood(GARCH)-likelihood(ARCH))
p_value <- 1 - pchisq(LR_statistic, df = 1)

cat("Likelihood of ARCH: ", round(likelihood(ARCH),2), "\n",
    "Likelihood of GARCH: ", round(likelihood(GARCH),2), "\n",
    "2 * (LH - LL) = ", round(LR_statistic,2), "\n",
    "p-value: ", p_value)
```

Likelihood of ARCH: 14598.73  
Likelihood of GARCH: 15449.95  
2 \* (LH - LL) = 1722.44  
p-value: 0

We find a p-value of 0, meaning that we have enough evidence to reject the null hypothesis that the two models are the same. The GARCH model has a significantly larger likelihood than the ARCH model.

## 11 | The Apply family of functions

The apply family of functions in R allows you to repetitively perform an action on multiple chunks of data. It is a good practice to get familiar with these functions since they are a direct way to perform operations that otherwise would require you to set up a loop. They are usually much faster too. The reason it is a family of functions instead of a single function, is that each one is specific to some data type and will provide the output in a different format.

Before getting started, let's take a look at a simple example, imagine you have the following dataset:

```
In [138]: load("RET.RData")
y <- subset(RET, select = c(date))
head(y)

      MSFT      XOM      GE      JPM      INTC      C
0.006674  0.004110 -0.008310  0.004954  0.001929  0.002772
-0.009303 -0.0027362 -0.018356 -0.031045 -0.01276 -0.031515
-0.014571 -0.000316 -0.021545 -0.025929 -0.018637 -0.032026
0.012755  0.010133  0.000415  0.001528  0.020975  0.009270
0.029418  0.016845  0.012043  0.022346  0.018601  0.015048
-0.008405 -0.001410 -0.013952 -0.017387  0.001908  0.022333
```

And you want to get the variance for each row. The first thing that comes to mind is either writing a line of code for each row:

```
In [139]: print(var(y$MSFT))
print(var(y$XOM))
print(var(y$GE))
print(var(y$JPM))
print(var(y$INTC))
print(var(y$C))

[1] 0.0002155043
[1] 0.0001436519
[1] 0.0003956245
[1] 0.0001749218
[1] 0.0002597558
[1] 0.0002448805
```

Or write a loop:

```
In [140]: for (stock in names(y)) {
  print(var(y[c(stock)]))
}

      MSFT      XOM      GE      JPM      INTC      C
MSFT 0.0002155043
XOM 0.0001436519
GE 0.0003956245
JPM 0.0001749218
INTC 0.0002597558
C 0.0002448805
```

The first approach is inefficient, especially with a large number of columns. The second approach is useful but a loop can seem unnecessary. We can use the apply family of function for this:

```
In [141]: apply(y, 2, var)

      MSFT      XOM      GE      JPM      INTC      C
0.000215504300024893
0.000143651860649884
0.00035262645651305
0.000174921788784438
0.000259755847689449
0.000244880468210118
```

### The apply function

The basic apply function has the syntax: `apply(X, MARGIN, FUN)` where:

- X is an array, matrix, or dataframe (the data you will be performing the function on)
- MARGIN specifies if you are going to perform the function across rows (1) or columns (2)
- FUN is the function you want to use, which can include options

For example, let's create a matrix:

```
In [142]: a <- matrix(rnorm(20), nrow = 5, ncol = 4)
a

      0.874044  -0.0293764  -0.24445109 -0.34984960
0.8525058  -0.0173984  -0.87864296 -0.21970682
1.7020600  -1.2409390  -0.54733383  1.22877162
-1.5881660 -0.1629312  -0.24592251  0.09574781
0.5411785  2.0624779  -0.04103685  1.04744810
```

Let's find the row and column sums using MARGIN 1 and 2 respectively:

```
In [143]: # Row sums
apply(a, 1, sum)

1.404903122431 -0.731143612430116 3.6304569031668 -1.9012718669659 3.6100676568428

In [144]: # Column sums
apply(a, 2, sum)

2.51191152274534 3.6552636957194 -1.95738723320548 1.80241110657793
```

Let's find the largest element of every row:

```
In [145]: # Index of every row
apply(a, 1, max)

1.02637644719398 0.884804567966169 1.70208003179984 0.0957478125732051 2.062477906833
```

However, `apply()` will not work on vectors:

```
In [146]: vec <- 1:20
apply(vec, 1, sum)

Error in apply(vec, 1, sum): dim(X) must have a positive length
Traceback:
1. apply(vec, 1, sum)
2. stop("dim(X) must have a positive length")
```

### These and apply

There are variations of `apply()` that work with other data types.

#### lapply()

This variation can be used for other objects like dataframes, lists or vectors. The output return is a **list**, which has the same number of elements as the object passed to it.

For example, let's create a list from three matrices:

```
In [147]: A <- matrix(1:19, nrow = 3, ncol = 3)
B <- matrix(1:14, nrow = 2, ncol = 2)
C <- matrix(rep(1, 100), nrow = 10, ncol = 10)

my_list <- list(A,B,C)

# We can extract the 2nd column from every element using the selector "["
lapply(my_list, "[", 2)

1.4 5 6
2.3 4
3.1 1 1 1 1 1 1 1 1 1
```

# Or the first row of all elements

```
In [148]: lapply(my_list, "[", 1, )

1.1 4 7
2.1 3
3.1 1 1 1 1 1 1 1 1 1

Note that the outputs are lists. To subset a list you need to use double brackets:
```

```
In [149]: a <- lapply(my_list, "[", 1, )
a[[1]]

1 4 7

We can use lapply() on a vector. Let's get the square root of every element from 1 to 10:
```

```
In [150]: vec = 1:10
a <- lapply(vec, sqrt)
class(a)

1
2
3
4
5
6
7
8
9
10
1.4142135623731
2.14142135623731
1.73205080756888
2.23606797749979
2.44948974278318
2.64575131106459
2.82842712474619
3
3.16227766016838

'list'
```

The `apply()` functions works very similar to `lapply()`, with the difference that it simplifies the output. While `lapply()` provides a list as output, `apply()` usually provides vectors, which are easier to work with:

```
In [151]: vec = 1:10
b <- lapply(vec, sqrt)
b

1
2
3
4
5
6
7
8
9
10
1.4142135623731 1.73205080756888 2.23606797749979 2.44948974278318 2.64575131106459 2.82842712474619
3.16227766016838

'numeric'
```

## 12 | Introduction to dplyr

The `dplyr` package is a fast, consistent tool for working with data frame like objects. It is the most common tool for data manipulation, and it can make your life easier. Here we will provide a very brief introduction to it. It is the next iteration of the `plyr` package (remember to load `plyr` before `dplyr`).

A complete cheat sheet on how to use `dplyr` can be found here: <https://rstudio.com/wp-content/uploads/2015/02/data-wrangling-cheat-sheet.pdf>

Some of its most useful functions are:

### filter

`filter` keeps rows meeting a condition. The syntax is `filter(data, ..., preserve)`, where:

- data = The data frame containing the rows you want to filter
- ... = The conditions to be evaluated. To use multiple conditions, use the logical operators & (AND), | (OR)
- preserve = F to recalculate the grouping structure based on the resulting data

### select

`select` keeps only the specified variables. Syntax is `select(data, ...)`, where:

- data = The data frame containing the variables you want to select
- ... = one or more unquoted expressions separated by commas; you can treat variable names like they are positions, so you can use expressions like `x:y` to select range of variables; positive values select variables, and negative values drop variables

### mutate

The `mutate` function add new variables and preserves existing ones. Syntax is `mutate(data, newvar = ...)` where:

- data = The data object you want to add the variable to
- ... = The function you want to use to create the new variable

### summarize

`summarize` creates one or more scalar variables summarizing the variables of an existing data object. Syntax is `summarize(data, newvar = ...)` where:

- data = The data object containing the variable you want to summarize
- ... = The function you want to use to create the new variable

### group\_by

`group_by` converts a data object into a grouped data object where operations are perform 'by group'. `ungroup()` removes grouping. Syntax is `group_by(data, ..., add, drop)`, where:

- data = The data object containing the variable you want to group by
- ... = The variables you want to group by
- add = F to override existing group (default)
- drop = T to drop empty groups

### join

Used to join data objects together. Syntax is `%>%join(x, y, by)` where:

- x, y = The data objects to join
- by = A character vector of variables to join by

There are different types of join:

- `inner_join` returns all rows from x where there are matching values in y, and all columns from x and y
- `left_join` returns all rows from x, and all columns from x and y. Rows in x with no match in y will have NA values
- `right_join` similar to `left_join`
- `outer_join` returns all rows and all columns from both x and y

### Pipes

`dplyr` allows to perform several operations at the same time by connecting different functions with the symbol `%>%`.

It is particularly useful if you want to execute a function by group.

### Tidyverse

`dplyr` is one of the packages included in [tidyverse](https://www.tidyverse.org/), which is a collection of R packages designed for data science.

### Examples

We will use `mpg`, one of the existing datasets in the `tidyverse` to demonstrate how `dplyr` works:

```
In [152]: # Load in this order
library(plyr)
library(dplyr)

# Tidyverse
library(tidyverse)

Attaching package: 'plyr'

The following object is masked from 'package:lubridate':
  here

Attaching package: 'dplyr'

The following objects are masked from 'package:plyr':
  arrange, count, desc, filter, id, mutate, rename, summarise,
  summarise

The following objects are masked from 'package:rmqarch':
  first, last

The following objects are masked from 'package:lubridate':
  intersect, setdiff, union

The following objects are masked from 'package:stats':
  filter, lag

The following objects are masked from 'package:base':
  intersect, setdiff, setequal, union

tidyverse 1.2.1 ---
✔ ggplot2 3.2.1      ✔ readr 1.3.1
✔ tidyr 1.0.0        ✔ purrr 0.3.3
✔ dplyr 1.1.0        ✔ stringr 1.4.0
✔ lubridate 3.2.1    ✔ forcats 0.4.0
```

--- Conflicts --- masks plyr::arrange() tidyverse\_conflicts() ---  
# lubridate::as\_datetime() masks base::as\_datetime()  
# purrr::compact() masks plyr::compact()  
# dplyr::count() masks plyr::count()  
# lubridate::date() masks base::date()  
# dplyr::fullwidth() masks plyr::fullwidth()  
# dplyr::filter() masks stats::filter()  
# dplyr::first() masks rmqarch::first()  
# plyr::here() masks lubridate::here()  
# dplyr::id() masks plyr::id()  
# lubridate::intersect() masks base::intersect()  
# dplyr::lag() masks stats::lag()  
# dplyr::last() masks rmqarch::last()  
# dplyr::mutate() masks plyr::mutate()  
# purrr::reduce() masks purrr::reduce()  
# dplyr::rename() masks plyr::rename()  
# lubridate::setdiff() masks base::setdiff()  
# dplyr::summarise() masks plyr::summarise()  
# dplyr::summarize() masks plyr::summarize()  
# lubridate::union() masks base::union()

```
In [153]: # First look at the dataset
data <- mpg
head(mpg)

  manufacturer model displ year cyl trans drv cty hwy fl class
1 audi a4 1.8 1999 4 auto(l5) f 18 29 p compact audi a4
2 audi a4 1.8 1999 4 manual(m5) f 21 29 p compact audi a4
3 audi a4 2.0 2008 4 manual(m5) f 20 31 p compact audi a4
4 audi a4 2.0 2008 4 auto(av) f 21 30 p compact audi a4
5 audi a4 2.8 1999 6 auto(l5) f 16 26 p compact audi a4
6 audi a4 2.8 1999 6 manual(m5) f 18 26 p compact audi a4
```

```
In [154]: # Select only manufacturer, model and class
data_select <- select(data, manufacturer, model, class)
head(data_select)

  manufacturer model class
1 audi a4 compact
2 audi a4 compact
3 audi a4 compact
4 audi a4 compact
5 audi a4 compact
6 audi a4 compact
```

```
In [155]: # Filter only the cars that are pickups
data_filter <- filter(data, class == "pickup")
head(data_filter)

  manufacturer model displ year cyl trans drv cty hwy fl class
1 dodge dakota pickup 4wd 3.7 2008 6 manual(m6) 4 15 19 r pickup
2 dodge dakota pickup 4wd 3.7 2008 6 auto(l4) 4 14 18 r pickup
3 dodge dakota pickup 4wd 3.9 1999 6 manual(m5) 4 13 17 r pickup
4 dodge dakota pickup 4wd 4.7 2008 8 auto(l5) 4 14 19 r pickup
5 dodge dakota pickup 4wd 4.7 2008 8 auto(l5) 4 14 19 r pickup
```

```
In [156]: # Select only the manufacturer model and class, and then filter
# Note that with pipes you don't have to specify "data" again
data_pipes <- data %>% select(manufacturer, model, class) %>% filter(class == "pickup")
head(data_pipes)

  manufacturer model class
1 dodge dakota pickup 4wd pickup
2 dodge dakota pickup 4wd pickup
3 dodge dakota pickup 4wd pickup
4 dodge dakota pickup 4wd pickup
5 dodge dakota pickup 4wd pickup
```

```
In [157]: # Create a new column that is a concatenation of manufacturer and model
data <- data %>% mutate(car = paste(manufacturer, model, sep = " - "))
head(data)

  manufacturer model displ year cyl trans drv cty hwy fl class car
1 audi a4 1.8 1999 4 auto(l5) f 18 29 p compact audi a4
2 audi a4 1.8 1999 4 manual(m5) f 21 29 p compact audi a4
3 audi a4 2.0 2008 4 manual(m5) f 20 31 p compact audi a4
4 audi a4 2.0 2008 4 auto(av) f 21 30 p compact audi a4
5 audi a4 2.8 1999 6 auto(l5) f 16 26 p compact audi a4
6 audi a4 2.8 1999 6 manual(m5) f 18 26 p compact audi a4
```

```
In [158]: # Get the mean for the hwy variable
avg_hwy <- data %>% summarize(avg_hwy = mean(hwy))
head(avg_hwy)

avg_hwy
23.44017
```

```
In [159]: # Get the mean of highway mileage for every manufacturer
avg_hwy2 <- data %>% group_by(manufacturer) %>% summarize(avg_hwy = mean(hwy))
head(avg_hwy2)

  manufacturer avg_hwy
1 audi 26.4444
2 chevrolet 21.89474
3 dodge 17.94595
4 ford 19.36000
5 honda 32.55556
6 hyundai 26.85714
```

```
In [160]: # For every manufacturer, get the average highway and city mileage
avg_hwy3 <- data %>% group_by(manufacturer) %>% summarize(avg_hwy = mean(hwy),
  avg_cty = mean(cty))
head(avg_hwy3)

  manufacturer avg_hwy avg_cty
1 audi 26.44444 17.611
```



```
[163]: library(ggplot2)
library(tidyverse)
```

We will use the `mpg` dataset:

```
In [164]: mpg <- mpg
          head(mpg)
```

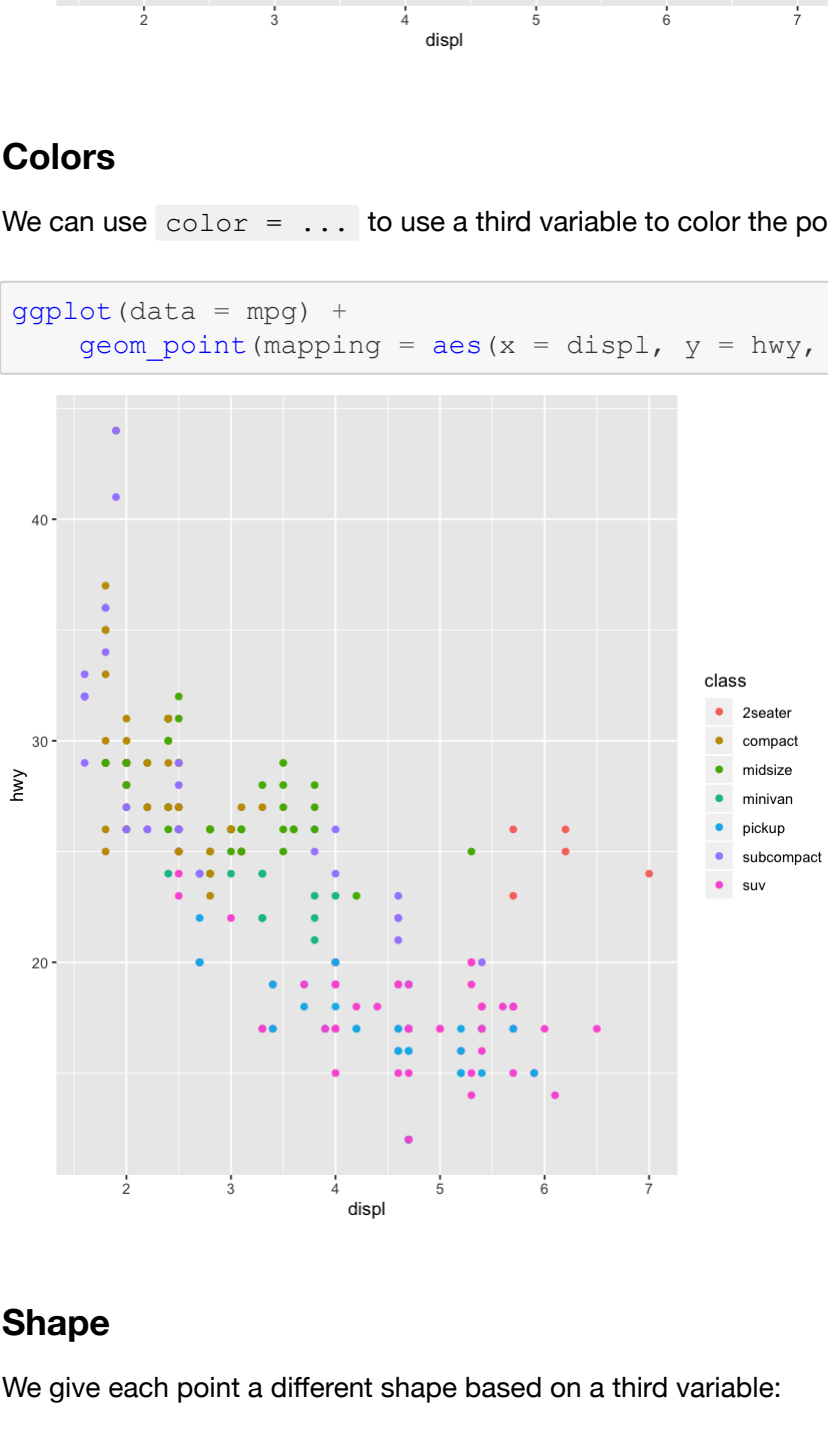
manufacturer	model	displ	year	cyl	trans	drv	cty	hwy	fl	class
audi	a4	1.8	1999	4	auto(l5)	f	18	29	p	compact
audi	a4	1.8	1999	4	manual(m5)	f	21	29	p	compact
audi	a4	2.0	2008	4	manual(m6)	f	20	31	p	compact
audi	a4	2.0	2008	4	auto(av)	f	21	30	p	compact
audi	a4	2.8	1999	6	auto(l5)	f	16	26	p	compact
audi	a4	2.8	1999	6	manual(m5)	f	18	26	p	compact

To create a `ggplot2` graph, we begin by using `ggplot(data = ...)` to specify the data that will feed the plot. On top of that we add the elements of the plot, like `geom_point(mapping = aes(x = ..., y = ...))` to add points. `aes()` stands for aesthetics, and it holds the values in the `x` and `y` axis. It allows for several options that we will explore:

## Scatter plot

We will build a scatter plot with engine displacement in the `x`-axis and highway consumption in the `y`-axis.

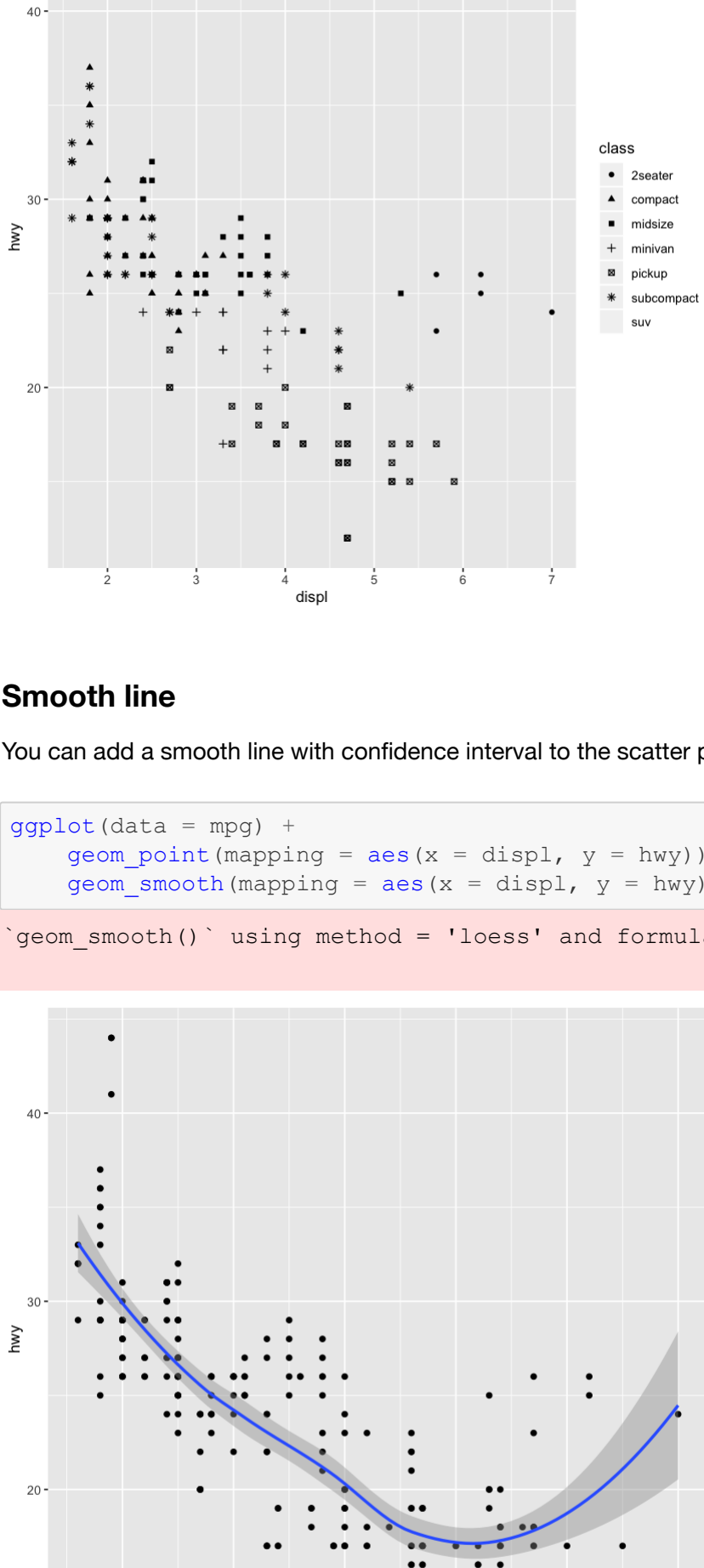
```
In [165]: # Scatter plot
          ggplot(data = mpg) +                # Specify the data
            geom_point(mapping = aes(x = displ, y = hwy)) # Add points
```



## Colors

We can use `color = ...` to use a third variable to color the points plotted. This also creates a legend automatically:

```
In [166]: ggplot(data = mpg) +
          geom_point(mapping = aes(x = displ, y = hwy, color = class))
```

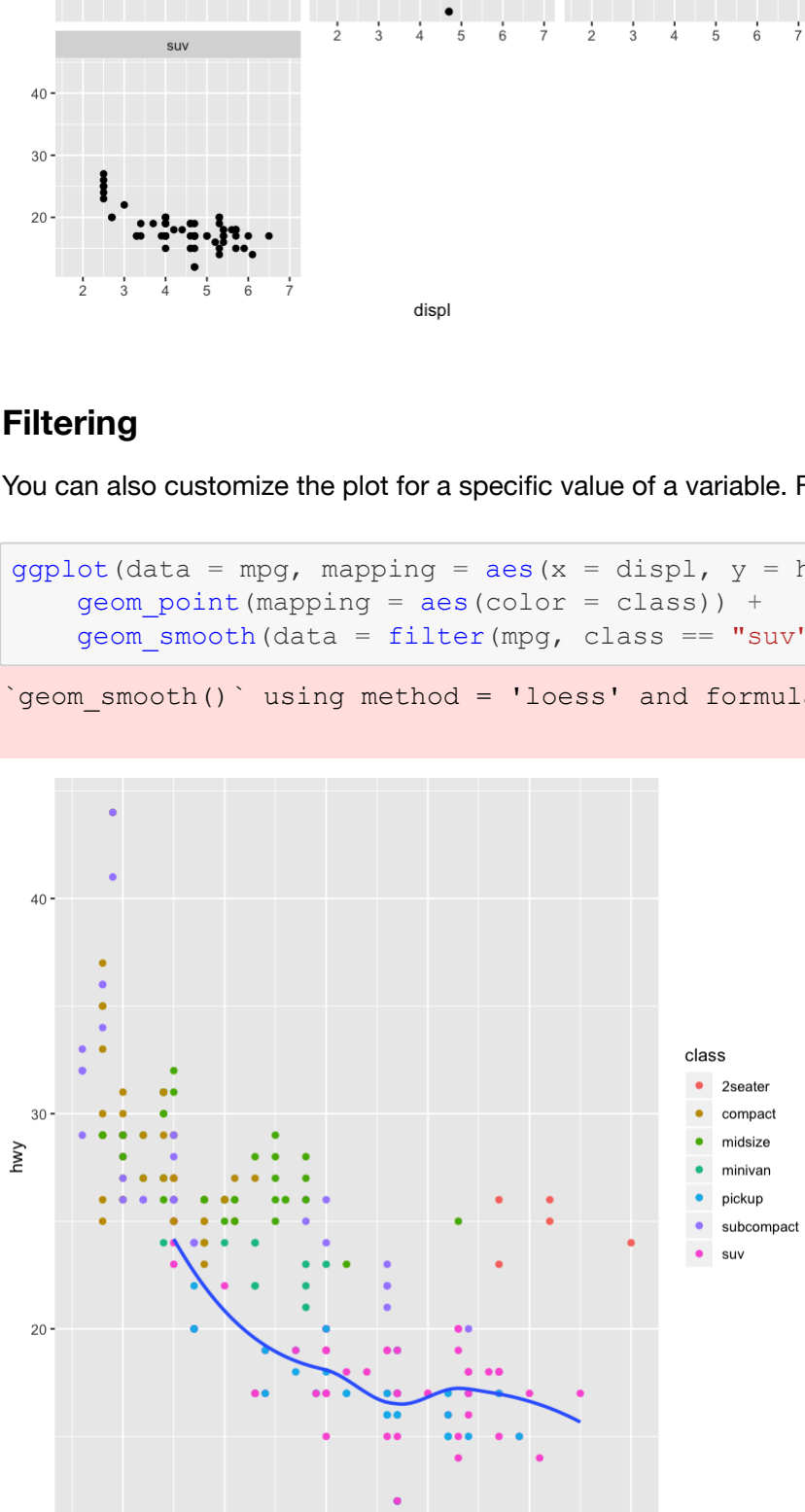


## Shape

We give each point a different shape based on a third variable:

```
In [167]: ggplot(data = mpg) +
          geom_point(mapping = aes(x = displ, y = hwy, shape = class))
```

Warning message:  
"The shape palette can deal with a maximum of 6 discrete values because more than 6 becomes difficult to discriminate; you have 7. Consider specifying shapes manually if you must have them."  
Warning message:  
"Removed 62 rows containing missing values (geom\_point)."

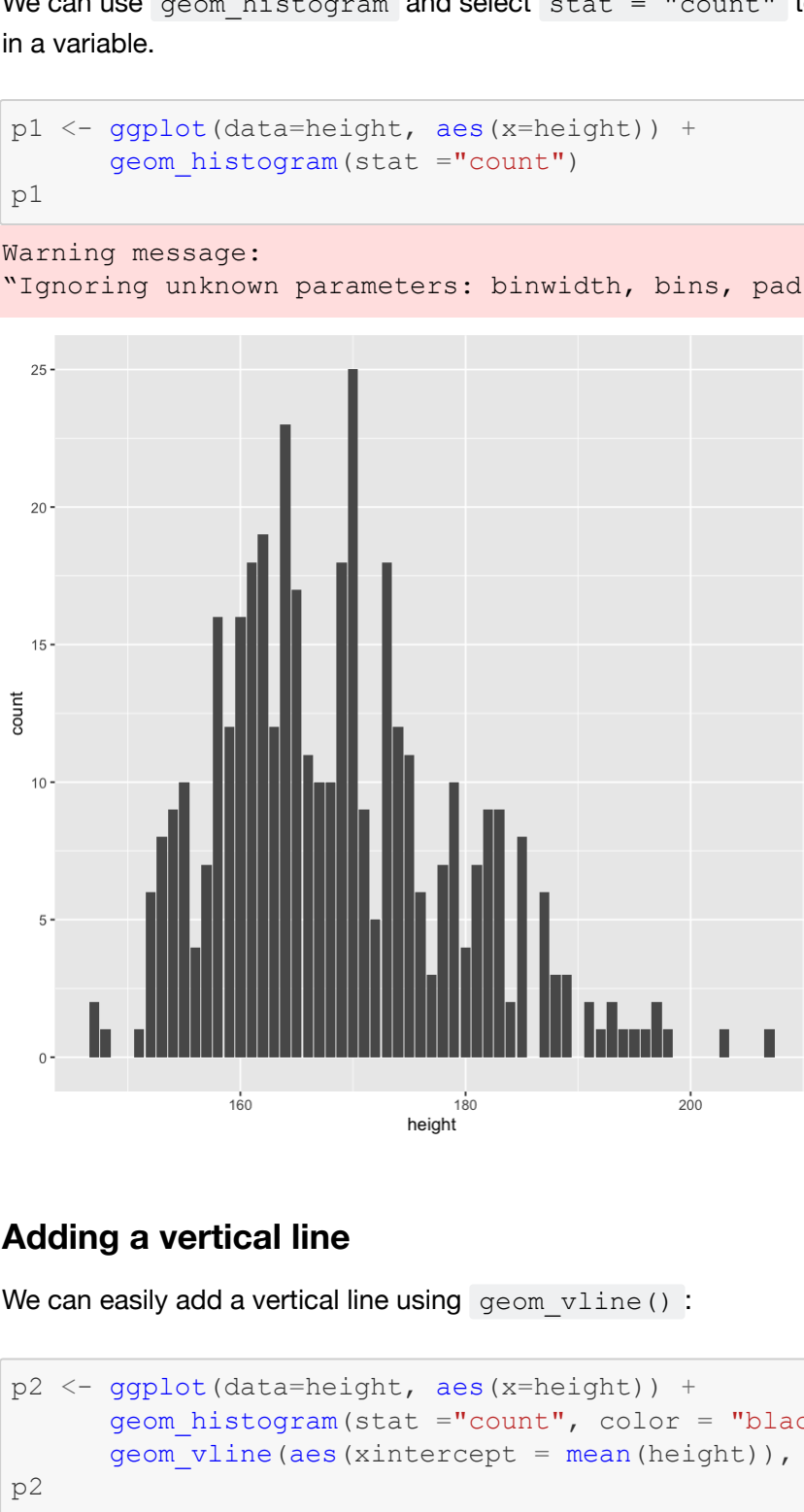


## Smooth line

You can add a smooth line with confidence interval to the scatter plot using `geom_smooth()` :

```
In [168]: ggplot(data = mpg) +
          geom_point(mapping = aes(x = displ, y = hwy)) +
          geom_smooth(mapping = aes(x = displ, y = hwy))
```

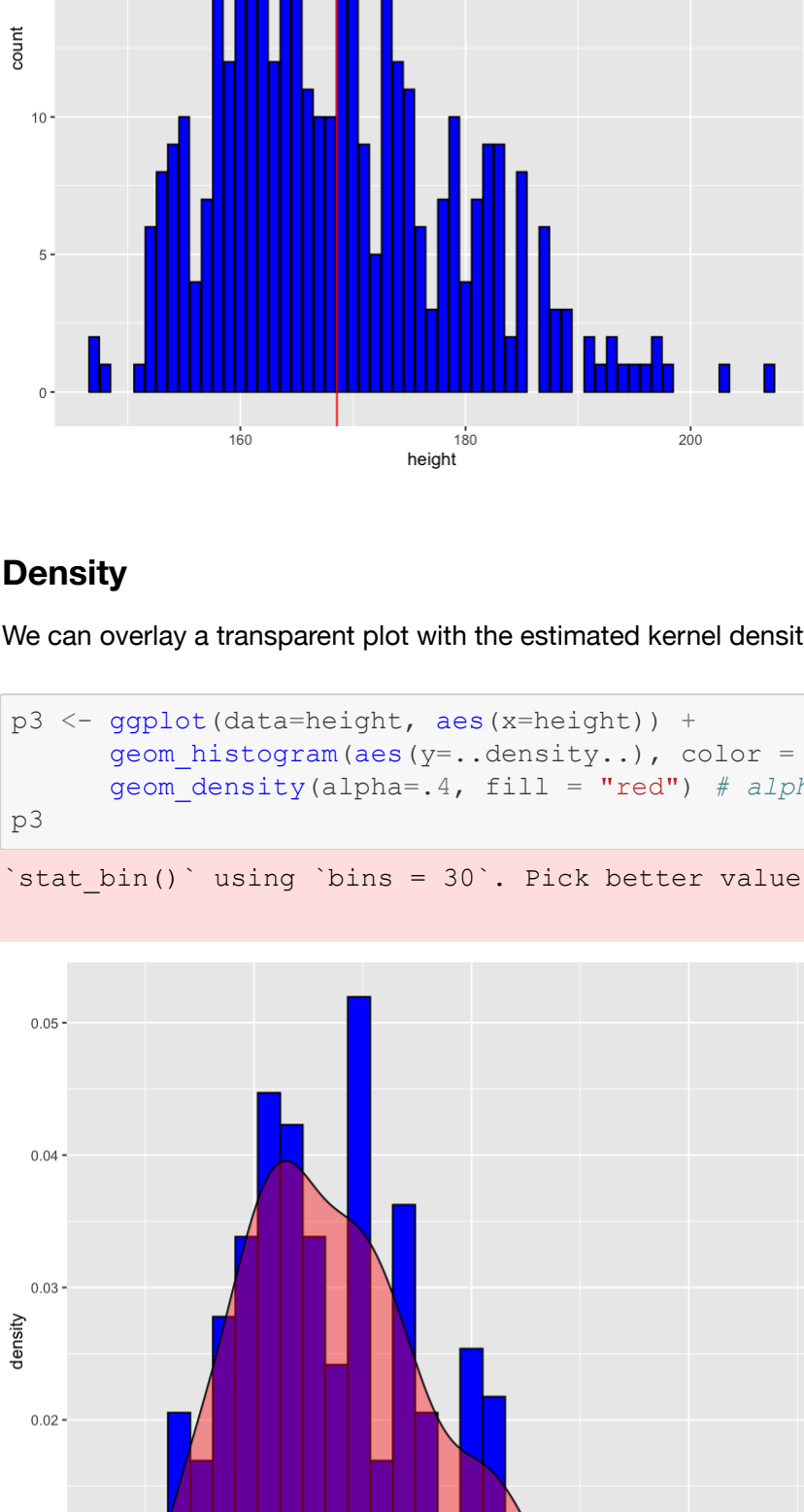
Warning message:  
"geom\_smooth() using method = 'loess' and formula 'y ~ x'"



## Separate plots by a variable

If you want to make a different plot for each element of a variable, you can do so with `facet_wrap()`:

```
In [169]: ggplot(data = mpg, mapping = aes(x = displ, y = hwy)) +
          facet_wrap(~ class, nrow = 3)
```

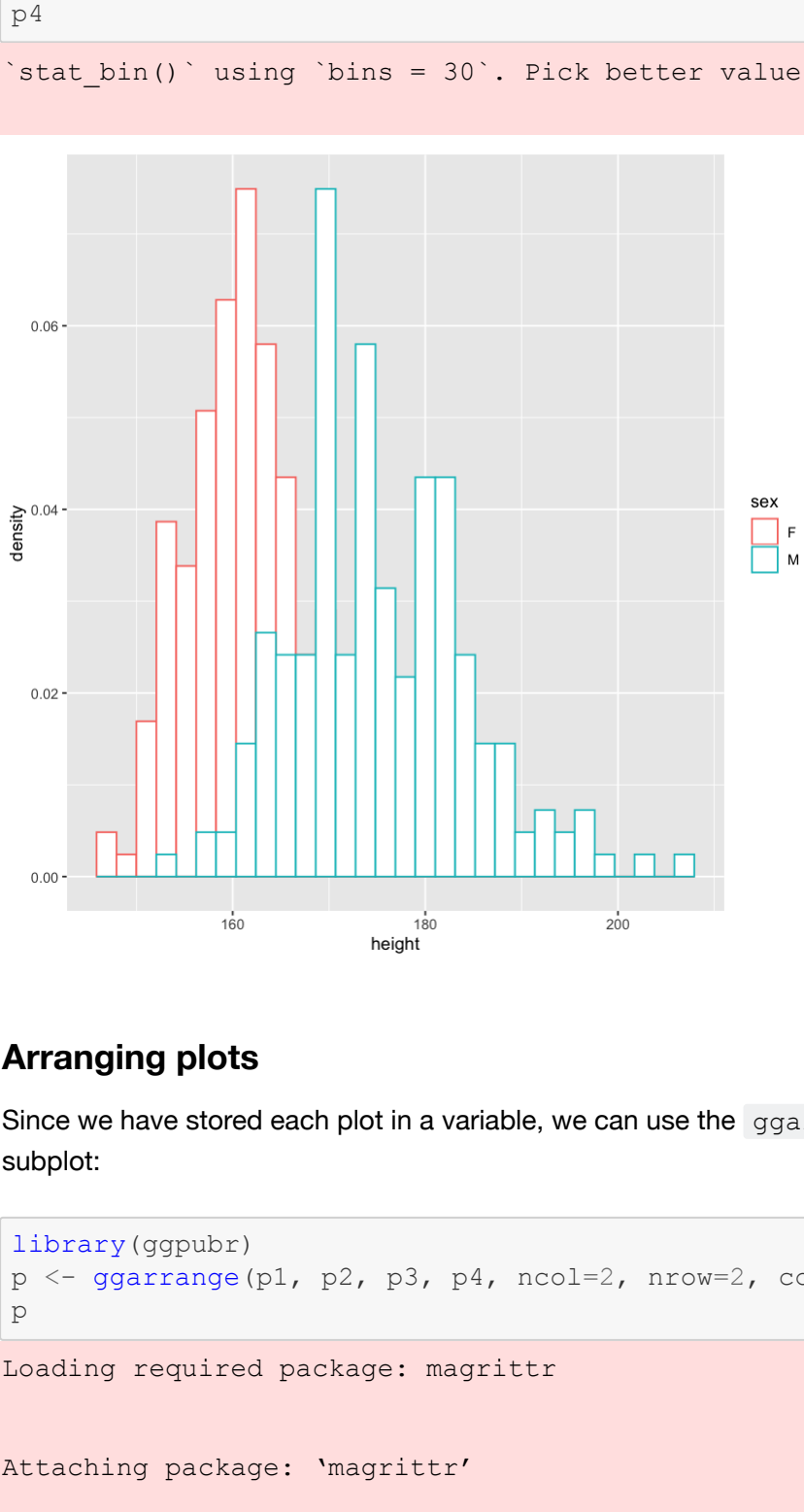


## Filtering

You can also customize the plot for a specific value of a variable. For example, let's add a smooth line only for SUVs:

```
In [170]: ggplot(data = mpg, mapping = aes(x = displ, y = hwy)) +
          geom_point(mapping = aes(color = class)) +
          geom_smooth(data = filter(mpg, class == "suv"), se = F)
```

Warning message:  
"geom\_smooth() using method = 'loess' and formula 'y ~ x'"



## Histograms and bar plots

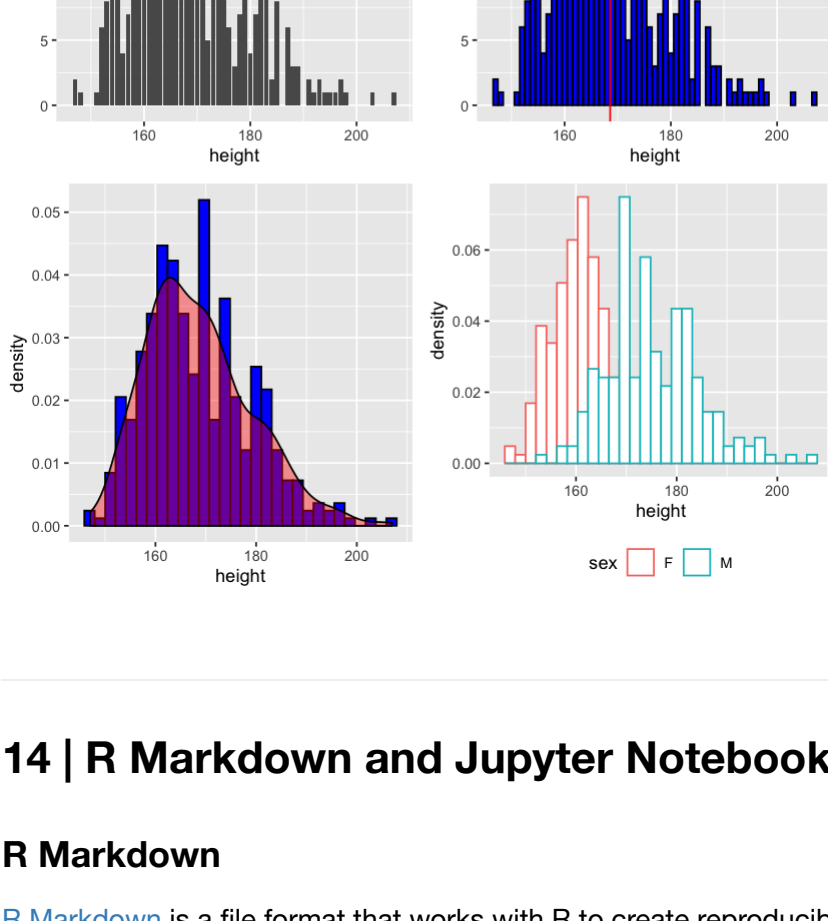
Let's create a dataframe with information on sex and height to build some more plots:

```
In [171]: height <- data.frame(
          sex = factor(rep(c("F", "M"), each=200)),
          height = round(rnorm(200, mean=161.5, sd=6), xnorm(200, mean=175.4, sd=9)))
          }
```

## Bar Plot

We can use `geom_histogram` and select `stat = "count"` to build a simple bar plot. Notice that you can store the `ggplot` object in a variable.

```
In [172]: p1 <- ggplot(data=height, aes(x=height)) +
          geom_histogram(stat = "count")
          p1
```



## Adding a vertical line

We can easily add a vertical line using `geom_vline()`:

```
In [173]: p2 <- ggplot(data=height, aes(x=height)) +
          geom_histogram(stat = "count", color = "black", fill="blue") +
          geom_vline(aes(xintercept = mean(height)), se = "red")
          p2
```

Warning message:  
"Ignoring unknown parameters: binwidth, bins, pad"



## Density

We can overlay a transparent plot with the estimated kernel density with `geom_density()`:

```
In [174]: p3 <- ggplot(data=height, aes(x=height)) +
          geom_histogram(aes(y=..density..), color = "black", fill="blue") +
          geom_density(alpha=.4, fill = "red") # alpha determines how transparent it is
          p3
```

Warning message:  
"stat\_bin() using 'bins = 30'. Pick better value with 'binwidth'."



## Overlapping histograms

With the option `position = identity` we can plot two histograms, one in front of the other:

```
In [175]: p4 <- ggplot(data=height, aes(x=height, color = sex)) +
          geom_histogram(aes(y=..density..), fill = "white", position="identity") # Position identity
          for overlapping histograms
          p4
```

Warning message:  
"stat\_bin() using 'bins = 30'. Pick better value with 'binwidth'."



## Arranging plots

Since we have stored each plot in a variable, we can use the `ggarrange()` function from the `ggpubr` package to plot each in a subplot:

```
In [176]: library(ggpubr)
          p <- ggarrange(p1, p2, p3, p4, ncol=2, nrow=2, common.legend = F, legend="bottom")
          p
```

Loading required package: magrittr

Attaching package: 'magrittr'

The following object is masked from 'package:purrr':

set\_names

The following object is masked from 'package:tidyr':

extract

Attaching package: 'ggpubr'

The following object is masked from 'package:plyr':

mutate

Warning message:  
"stat\_bin() using 'bins = 30'. Pick better value with 'binwidth'."

Warning message:  
"stat\_bin() using 'bins = 30'. Pick better value with 'binwidth'."



## 14 | R Markdown and Jupyter Notebooks

### R Markdown

**R Markdown** is a file format that works with R to create reproducible, dynamic reports. You can use it to easily embed your code and results into slideshows, pdf, html, Word, PowerPoint and more.

You can create a Markdown file from RStudio by choosing **File > New File > R Markdown...** You will be asked what type of output you would like to generate (html, pdf, word). Then you can start writing your document following a markdown syntax (check the [cheat sheet](#) in *Other references* for guidance). R Markdown allows you to embed code in the report that will be ran when the report is rendered and include the results.

### Jupyter notebooks

Jupyter notebooks are file types with extension `.ipynb` produced by [Project Jupyter](#) which contain both computer code (e.g. R, Python, Julia), and rich text elements (paragraphs, equations, images, links, etc). They are growing in popularity and are a great way to elaborate reports and documents. They can be edited on a web server and exported as html, pdf via LaTeX, or other file formats. They are interactive and you can independently run pieces of code. As an example, this set of notes has been created on a Jupyter notebook.