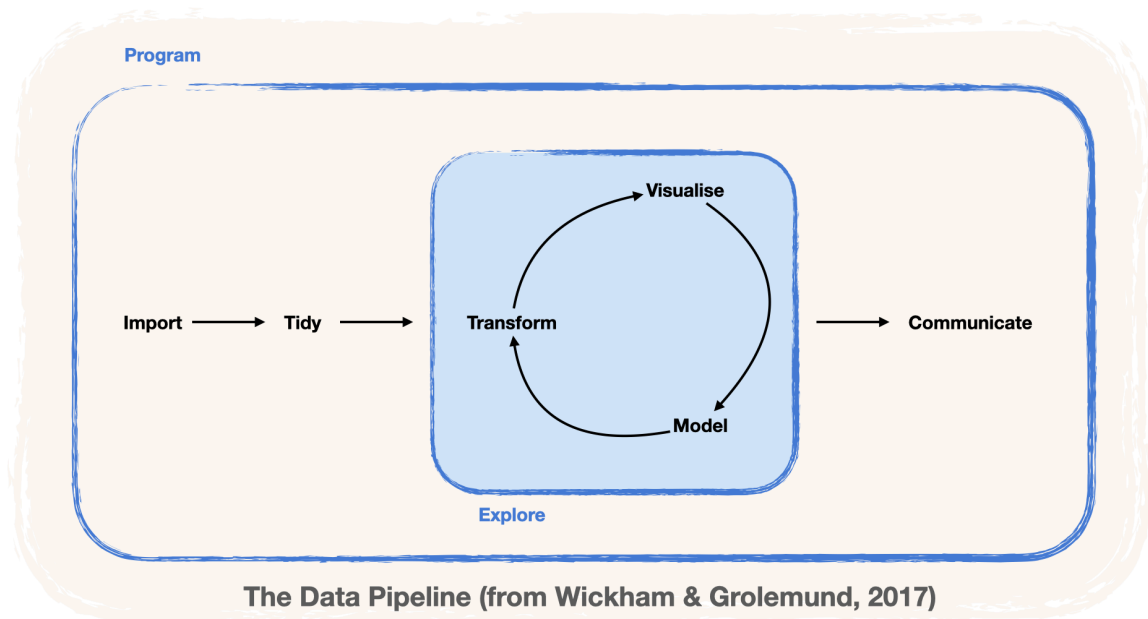


The Data Science Pipeline

Table of contents

1 Introduction

In this lesson we will explore how to carry the steps of the Data Science pipeline to produce an analysis



To illustrate, we will study the stocks of **five tech companies**

- First we'll focus on a single company
- Then we'll see how all these companies have fared over time.

2 Importing the data

- We've imported the data using the package `{tidyquant}`.

```
get_what <- "stock.prices"
companies <- c("AAPL",
               "MSFT",
               "GOOG",
               "AMZN",
               "TSLA")

stocks_data <-
  companies |>
  map(tq_get, get = get_what) |>
  bind_rows()
```

- The dataset contains the information on stocks for 6 companies: Apple Inc. (**AAPL**), Microsoft Corporation (**MSFT**), Alphabet Inc. (**GOOG**), Amazon Inc. (**AMZN**) and Tesla Inc. (**TSLA**).

```
stocks_data |>
  slice_head(n=6) |>
  kable()
```

Table 2.1: The first six lines of the whole dataset

symbol	date	open	high	low	close	volume	adjusted
AAPL	2014-01-02	19.84571	19.89393	19.71500	19.75464	234684800	17.29666
AAPL	2014-01-03	19.74500	19.77500	19.30107	19.32071	392467600	16.91672
AAPL	2014-01-06	19.19464	19.52857	19.05714	19.42607	412610800	17.00897
AAPL	2014-01-07	19.44000	19.49857	19.21143	19.28714	317209200	16.88733
AAPL	2014-01-08	19.24321	19.48429	19.23893	19.40929	258529600	16.99427
AAPL	2014-01-09	19.52857	19.53071	19.11964	19.16143	279148800	16.77725

3 Tidying the data

3.1 Choosing a Company

- A ticker can be stored in the variable `company_ticker`.
- Later, this will help us *parametrise* our report

```
company_ticker <- "AAPL"
```

3.2 Choosing a Timeframe

- To narrow our focus, we restrain our analysis to a given timeframe.
- We will focus on the performance of these stocks since the beginning of the COVID-19 pandemic (March 11 2020) until now.
- The `start_date` variable can store this information.

```
start_date <- ymd("2020-03-11")
```

3.3 Filtering the data

- We're interested in the company APPLE INC, we can use the ticker AAPL
- We'll use the `filter()` to subset the data for the company in the given timeframe.

```
company_data <-  
stocks_data |>  
filter(symbol == company_ticker,  
       date > start_date)
```

- We can explore the data by printing its first 6 lines:

```
head(company_data) |>  
  kable()
```

Table 3.1: The first six lines of the apple dataset

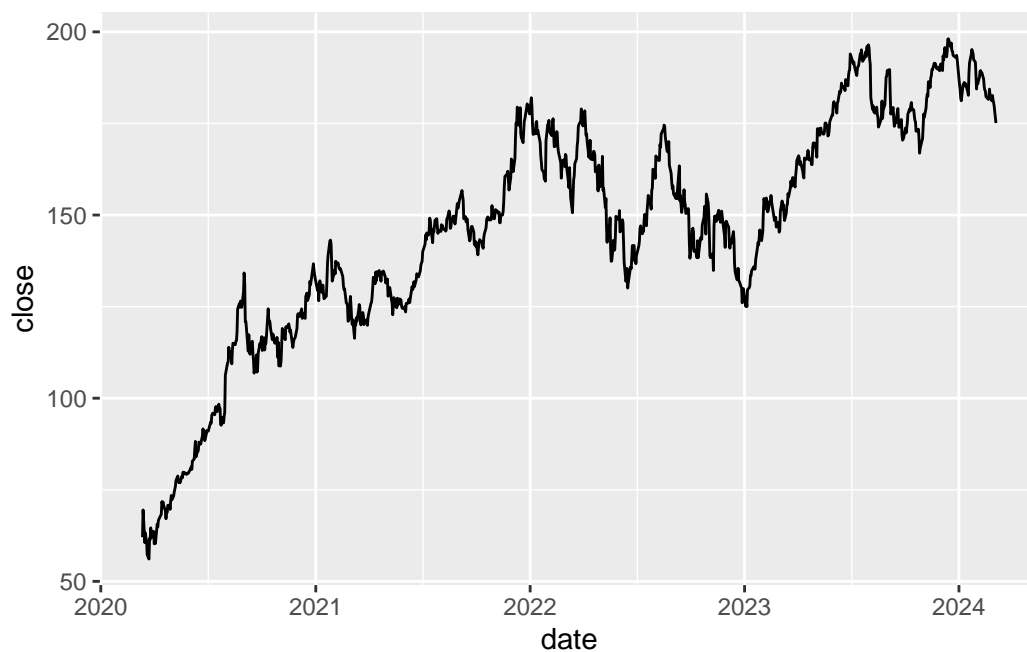
symbol	date	open	high	low	close	volume	adjusted
AAPL	2020-03-12	63.9850	67.5000	62.0000	62.0575	418474000	60.52466
AAPL	2020-03-13	66.2225	69.9800	63.2375	69.4925	370732000	67.77601
AAPL	2020-03-16	60.4875	64.7700	60.0000	60.5525	322423600	59.05684
AAPL	2020-03-17	61.8775	64.4025	59.6000	63.2150	324056000	61.65357
AAPL	2020-03-18	59.9425	62.5000	59.2800	61.6675	300233600	60.14429
AAPL	2020-03-19	61.8475	63.2100	60.6525	61.1950	271857200	59.68346

4 Understanding the data

4.1 Visualise

- With this data and the functions in `ggplot()`, we can create a first visualisation of the closing stock price (`close`).
- We set the dates on the x axis and the `close` price in the y axis.

```
company_data |>  
  ggplot(aes(x = date))+  
  geom_line(aes(y = close))
```



4.2 Transform

- The visualisation offers a first glance. We can transform again to ask the questions on the returns.

- Remember that the difference in log-price are approximations of the returns, i.e. the **percentage gain after selling the stock**.

Definition 4.1. Let p_t denote the closing price of the stock, the log-return r_t can be defined as:

$$r_t = \log(p_t) - \log(p_{t-1}) \approx \frac{p_t - p_{t-1}}{p_{t-1}} \quad (4.1)$$

- We use the function `mutate()`, alongside `lag()` to create a column with the daily (log) returns and the definition in equation Definition ??

```
company_data <- company_data |>
  mutate(daily_log_returns = log(close)-log(lag(close)))
```

4.3 Visualize again: log-returns

We can construct a visualisation with this. Additionally, we can add layers to our visualisation to decorate it at will.

```
company_data |>
  ggplot(aes(x = date))+
  geom_line(aes(y = daily_log_returns), alpha = 0.5, color = "#555555") +
  geom_hline(yintercept = 0, lty = 3)+
  labs(
    title = str_glue("Daily Returns of the stock for {company_ticker}"),
    subtitle = str_glue("Close stock prices since {start_date |> format('%d %B, %Y')}"),
    x = "Date",
    y = "Returns"
  ) +
  theme_minimal()
```

Warning: Removed 1 row containing missing values or values outside the scale range (``geom_line()``).

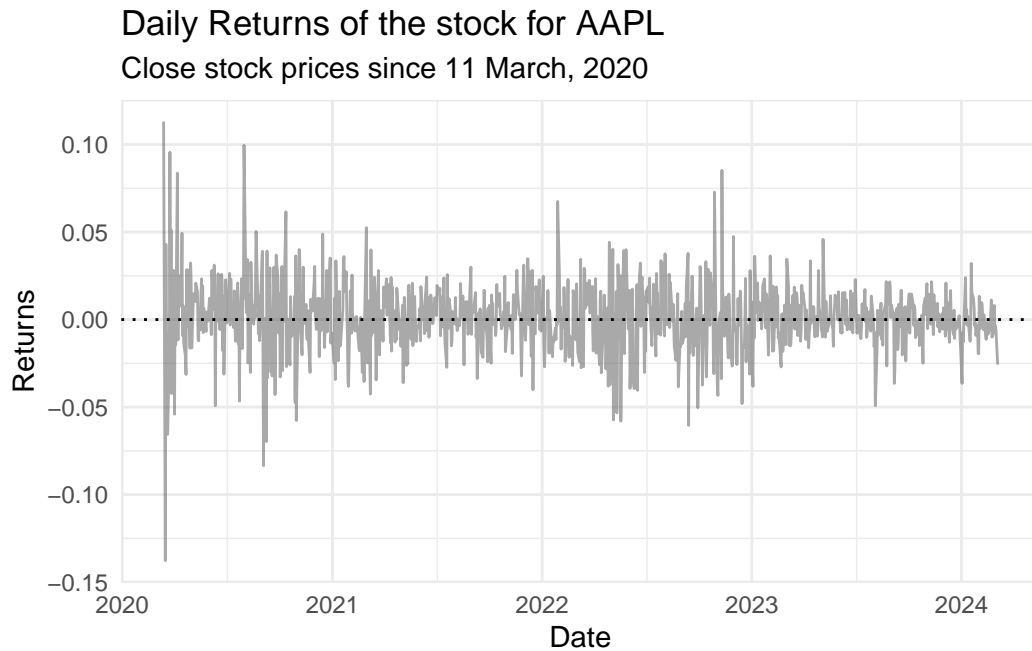


Figure 4.1: Log-returns for AAPL

! Look at the years!

What can you say about this plot?

- It appears that there is a high **variability** of the log-returns in the year 2020
- This increase in **variance** seems to stabilise in 2021 and **reappear** in 2022
- The year 2023 is also quite stable

All these point to signs of increased **variability** in times of **global crises**, which have added elements of uncertainty to the global supply chain.

4.4 Model: Create summary statistics

- We can obtain a summary table for some summary statistics with `summarise()`.
- We will compute the average return \bar{r} and estimate the standard deviation (SD) of the log-returns $\hat{\sigma}$

Definition 4.2. These statistics are defined as follows:

$$\bar{r} = \frac{1}{T} \sum_{t=1}^T r_t \quad (4.2)$$

$$\hat{\sigma} = \sqrt{\frac{1}{T-1} \sum_{t=1}^T (r_t - \bar{r})^2} \quad (4.3)$$

A word of warning

While the average return might not mean much in theoretical terms, the standard deviation might give some idea of the *risk* or *volatility*.

- We can show the values in the following table:

```
company_data |>
  summarise(`Average Return` = mean(daily_log_returns, na.rm = TRUE),
            `Average Risk (SD)` = sd(daily_log_returns, na.rm = TRUE)) |>
  gt() |>
  tab_header(
    title= "Summary statistics of Tech companies stocks",
    subtitle = str_glue("From {start_date |> format('%d %b, %Y')} to {Sys.Date() |> format(
  fmt_number(decimals = 4)
```

Table 4.1: Summary statistics of the log-returns for AAPL

Summary statistics of Tech companies stocks
From 11 Mar, 2020 to 05 Mar, 2024

Average Return	Average Risk (SD)
0.0010	0.0200

4.4.1 Summary statistics by year

And, seeing that the visualisation shows periods of high volatility in the year 2020, we can compute yearly measures of risk and volatility: