# Accessing and Managing Financial Data with Tidy Finance
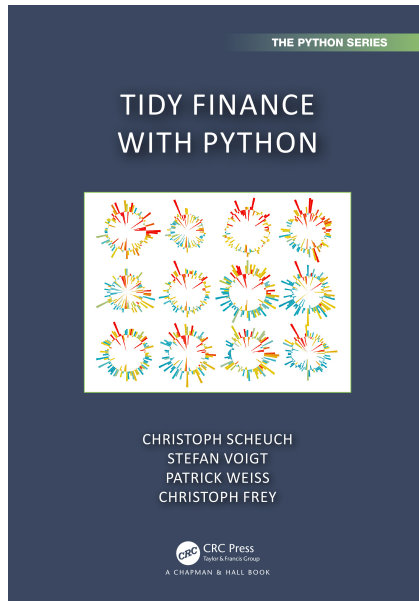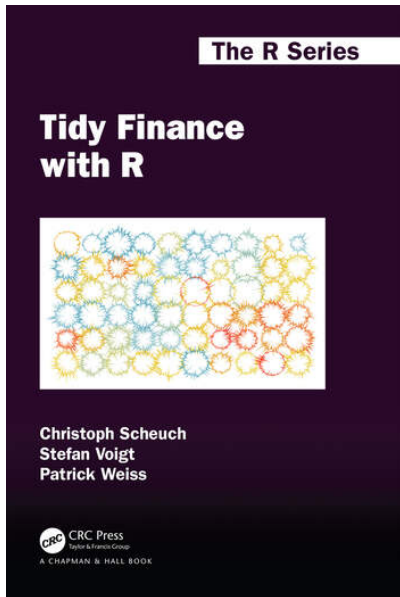## useR! Conference

Christoph Scheuch

2024-07-09

# A few years ago …

Two PhD students in Vienna
- ▶ Hardly any public code or data
- ▶ Hard to reproduce papers
- ▶ 80% of time spent preparing data

## Since then ...

# What is Tidy Finance?

A **transparent**, **open-source** approach to research in financial economics, featuring **multiple programming languages**
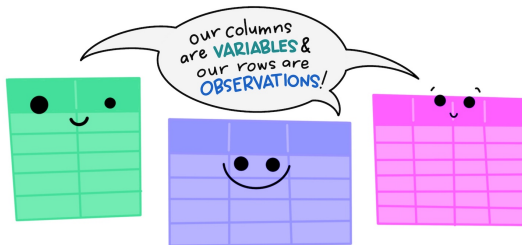
The `tidyfinance` packages is a simple way to:

- ▶ Load our approach into R
- ▶ Use helper functions to download & process data
- ▶ Easily compile multiple data sources

# Why tidy?

1. Write code that is **easy to read** for humans
2. **Compose simple functions** to solve complex problems
3. **Embrace functional programming** for reproducible results
4. **Reuse data structures** across applications
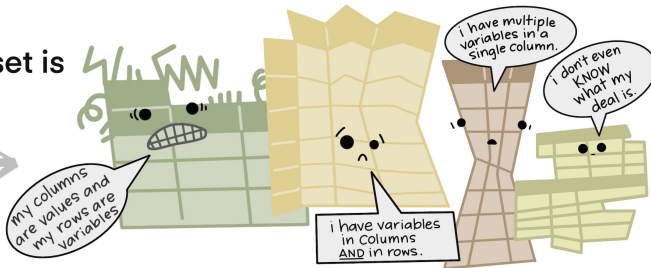
# Recap: what is *tidy data*?

The standard structure of tidy data means that **"tidy datasets are all alike..."**

"...but every messy dataset is messy in its own way."

—HADLEY WICKHAM

# A consistent interface to financial data

```
library(tidyfinance)

download_data(
  type = "factors_ff3_monthly",
  start_date = "2023-01-01",
  end_date = "2023-12-31"
) |>
  print(n = 5)
```

## A consistent interface to financial data

```r
library(tidyfinance)

download_data(
  type = "factors_ff3_monthly",
  start_date = "2023-01-01",
  end_date = "2023-12-31"
) |>
  print(n = 5)
```

```
# A tibble: 12 x 5
  date       risk_free mkt_excess     smb     hml
  <date>         <dbl>      <dbl>   <dbl>   <dbl>
1 2023-01-01    0.0035     0.0665  0.05   -0.0402
2 2023-02-01    0.0034    -0.0258  0.0117 -0.0081
3 2023-03-01    0.0036     0.0251 -0.0551 -0.0886
4 2023-04-01    0.0035     0.0061 -0.0335 -0.0004
5 2023-05-04    0.0036     0.0035  0.046  -0.0772
```

# Deep dive: download raw data

```
1  raw_data <- download_french_data(dataset)
2  raw_data <- raw_data$subsets$data[[1]]
```

# Deep dive: parse dates

```r
raw_data <- download_french_data(dataset)
raw_data <- raw_data$subsets$data[[1]]

if (grepl("monthly", type)) {
  processed_data <- raw_data |>
    mutate(date = floor_date(ymd(paste0(date, "01")), "month"))
} else {
  processed_data <- raw_data |>
    mutate(date = ymd(date))
}
```

# Deep dive: transform numeric columns

```
1  raw_data <- download_french_data(dataset)
2  raw_data <- raw_data$subsets$data[[1]]
3
4  if (grepl("monthly", type)) {
5    processed_data <- raw_data |>
6      mutate(date = floor_date(ymd(paste0(date, "01")), "month"))
7  } else {
8    processed_data <- raw_data |> mutate(date = ymd(date))
9  }
10
11 processed_data <- processed_data |>
12   mutate(across(-date, ~na_if(.,-99.99)),
13          across(-date, ~ . / 100))
```

## Deep dive: rename columns
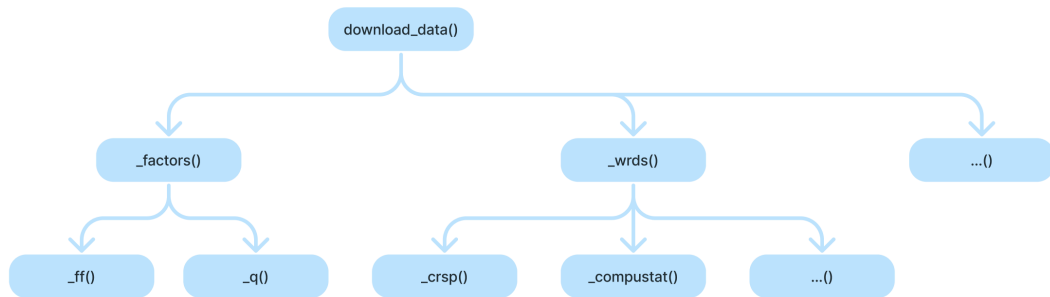
```
1  raw_data <- download_french_data(dataset)
2  raw_data <- raw_data$subsets$data[[1]]
3
4  if (grepl("monthly", type)) {
5    processed_data <- raw_data |>
6      mutate(date = floor_date(ymd(paste0(date, "01")), "month"))
7  } else {
8    processed_data <- raw_data |> mutate(date = ymd(date))
9  }
10
11 processed_data <- processed_data |>
12   mutate(across(-date, ~na_if(.,-99.99)),
13          across(-date, ~ . / 100)) |>
14   rename_with(tolower)
15
16 processed_data |>
17   rename_with(tolower) |>
```

# List of supported data sources

Currently 32 data sets supported from these domains:

▶ Fama-French factors
▶ Q factors
▶ Goyal-Welch macroeconomic predictors
▶ Wharton Research Database Service (WRDS)

# Easy to extend supported sources

# Example: load packages

```
library(tidyfinance)
library(dplyr)
```

# Example: load stock returns

```r
library(tidyfinance)
library(dplyr)

crsp <- download_data(
  "wrds_crsp_monthly", "2023-01-01", "2023-12-31"
)
```

# Example: load factors

```
1  library(tidyfinance)
2  library(dplyr)
3
4  crsp <- download_data(
5    "wrds_crsp_monthly", "2023-01-01", "2023-12-31"
6  )
7
8  factors <- download_data(
9    "factors_ff3_monthly", "2023-01-01", "2023-12-31"
10  )
```

# Example: join data

```r
1  library(tidyfinance)
2  library(dplyr)
3
4  crsp <- download_data(
5    "wrds_crsp_monthly", "2023-01-01", "2023-12-31"
6  )
7
8  factors <- download_data(
9    "factors_ff3_monthly", "2023-01-01", "2023-12-31"
10 )
11
12 stock_returns <- crsp |>
13   left_join(factors, join_by(month == date))
```

## Example: winsorize column

```r
library(tidyfinance)
library(dplyr)

crsp <- download_data(
  "wrds_crsp_monthly", "2023-01-01", "2023-12-31"
)

factors <- download_data(
  "factors_ff3_monthly", "2023-01-01", "2023-12-31"
)

stock_returns <- crsp |>
  left_join(factors, join_by(month == date))

stock_returns <- stock_returns |>
  mutate(mktcap_winsorized = winsorize(mktcap, 0.05))
```

# Example: summary statistics

```
1  data |>
2    create_summary_statistics(mktcap, mktcap_winsorized)
3
```

# Example applications: summary statistics

```
1  data |>
2    create_summary_statistics(mktcap, mktcap_winsorized)
3
```

```
# A tibble: 2 x 7
  variable               n    mean     sd   min   q50        max
  <chr>              <int>   <dbl>  <dbl> <dbl> <dbl>      <dbl>
1 mktcap             45552  10287. 73406. 0.307  507. 3071345.
2 mktcap_winsorized  45552   4452.  9355.  8.84  507.   37332.
```

# Example: assign portfolios

```
1  data |>
2    group_by(date) |>
3    mutate(
4      portfolio = assign_portfolio(
5        pick(everything()), "mktcap_winsorized", n_portfolios = 10
6      )
7    )
8
```

## Example: assign portfolios

```
1  data |>
2    group_by(date) |>
3    mutate(
4      portfolio = assign_portfolio(
5        pick(everything()), "mktcap_winsorized", n_portfolios = 10
6      )
7    )
8
```

```
# A tibble: 45,552 x 4
  date        permno mktcap_winsorized portfolio
  <date>       <int>             <dbl>     <int>
1 2023-02-28   12591            37332.        10
2 2023-02-28   12592               98.9        3
3 2023-02-28   12615              383.         5
4 2023-02-28   12622            37332.        10
5 2023-02-28   12623             8577.         9
# i 45,547
```

# Example: assign portfolios

```
1  data |>
2    group_by(date) |>
3    mutate(
4      portfolio = assign_portfolio(
5        pick(everything()), "mktcap_winsorized", n_portfolios = 10
6      )
7    )
8
```

Currently working on `calculate_portfolio_returns()`

▶ Calculate value-weighted and equal-weighted returns for different sorting methods

# Example: estimate model

```
1  data |>
2    estimate_model("ret_excess ~ mkt_excess + smb + hml")
3
```

# Example: estimate model

```
1  data |>
2    estimate_model("ret_excess ~ mkt_excess + smb + hml")
3
```

```
  mkt_excess        smb       hml
1  0.9657947 0.9403686 0.2689844
```

# Example: estimate model

```
1  data |>
2    estimate_model("ret_excess ~ mkt_excess + smb + hml")
3
```

Curently working on `roll_capm_estimation()`:

▶ Estimate betas for different lookbacks and multiple factors

# Tidy approach to financial data

▶ Check out open source content at **tidy-finance.org**
▶ Get in touch for **teaching materials**
▶ Submit issues to **extend supported types**
▶ Follow me for news: linkedin.com/in/christophscheuch