# Statistical Programming
# and Open Science Methods

## Tidy data scraping

Joachim Gassen
Humboldt-Universität zu Berlin

September 02, 2022

# Time table Monday, February 17

| When? | What? |
|-------|-------|
| 10:00 | Welcome and Coffee |
| 10:30 | Tidy data scraping |
| 12:00 | Lunch |
| 13:30 | Code along: Unit testing in function development |
| 15:00 | Coffee |
| 15:30 | Group work presentations |
| 16:30 | End of Day |

# Time table Tuesday, February 18

| When? | What? |
|-------|-------|
| 09:00 | Explore your researcher degrees of freedom |
| 10:30 | Coffee |
| 11:00 | Providing data access via RESTful APIs |
| 12:30 | Lunch |
| 13:30 | Group Work Presentations |
| 15:00 | Coffee and Wrap Up |
| 15:30 | End of Event |

# Parsing an HTML table I

# Parsing an HTML table II

# Parsing an HTML table III

```r
library(tidyverse)
library(rvest)

url_sp500_const <- paste0(
  "https://en.wikipedia.org/wiki/",
  "List_of_S%26P_500_companies"
)

url_sp500_const %>%
  read_html() %>%
  html_node(xpath = '//*[@id="constituents"]') %>%
  html_table() -> sp500_constituents_raw
```

# Retrieving local URLs from within table I

# Retrieving local URLs from within table II

```r
url_sp500_const %>%
  read_html() %>%
  html_node(xpath = '//*[@id="constituents"]') %>%
  html_nodes("td:nth-child(2) a") %>%
  html_attr("href")-> links
```

# Scraping ill-structured tables I

# Scraping ill-structured tables II

```
xml_data <- read_html(url) %>%
  html_node('#mw-content-text div table.infobox.vcard')

xml_data  %>%
  html_table(fill = TRUE) %>%
  rename(tag = X1, content = X2) %>%
  filter(tag != "")
```

# Another approach to scraping

# This won't work

```r
library(rvest)

url_bt_open_data <- "https://www.bundestag.de/services/opendata"
url_bt_open_data %>%
  read_html() %>%
  html_node(
    xpath = '//*[@id="bt-collapse-543410"]/div[1]/div/div/div[1]/table'
  ) %>% html_table() -> pp_table
```

```
## Error in UseMethod("html_table"): no applicable method for 'html_tab
```

Reason: The web page is being created dynamically by JavaScript (or similar)

# The idea of headless browsing: meet Selenium



▶ Selenium offers a way to script a web browser so that data can be scraped using the navigation that a web page provides

▶ Allows for various web browser and all sorts of user web browser interaction

▶ Uses a docker container holding the actual browser environment

▶ See `code/btag_open_data_scrape_data.R` for a demonstration on how to use it

# Parsing XML Data



Gaston Sanchez, https://github.com/gastonstat/tutorial-R-web-data
CC BY-NC-SA 4.0

# My task . . .

- ▶ Develop a function that parses the XML files of the Plenarprotokolle of the 19th Wahlperiode to extract all speaches into a tidy data structure
- ▶ Implement some basic test routines verifying that the code works
- ▶ See `code/btag_open_data_scrape_data.R` and `code/test/test_btag_open_data_scrape_data.R`