



Workshop 'Computational Text Analysis'

Session 1: Web-Scraping

Mirko Wegemann

24 March 2025



About us I

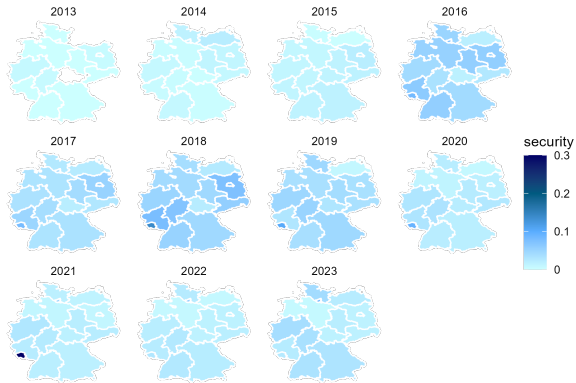
Mirko Wegemann (he/him)

- Postdoctoral Researcher at the University of Münster, Germany
- Research Agenda
 - Political Parties
 - Gender & Politics
 - Political Communication



About us III

- classification of gender appeals by the radical right using a multilingual transformer model (joint work with Leon Küstermann)





About you

- your name (if you want pronoun)
- research topic
- what's text analysis doing for your studies?
- do you have any experiences with text analysis?
- do you have any experiences with images-as-data?



Goals of this workshop

1. Automation of data collection
 2. Analysis of textual data
 - unsupervised approaches (e.g., topic models)
 - supervised approaches (e.g. text classification)
 3. Analysis of images-as-data
- understanding and dealing with key challenges, implementing first own analyses, getting a 'glimpse' into what's possible



What is web-scraping?



The process of systematically collecting data from the Internet to store it in a structured data format.



How can we utilize it in research? I



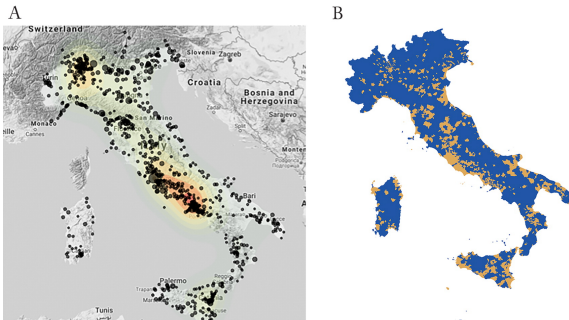
Figure 1.1 Location of UNESCO World Heritage Sites in danger (as of March 2014). Cultural sites are marked with triangles, natural sites with dots

World Heritage in Danger? (Munzert 2015, p. 5)

From a table on a [Wikipedia page](#) to a nice visualization



How can we utilize it in research? II



Bischof and Kurer (2023) show how local campaigning of M5S mobilizes voters



How can we utilize it in research? III

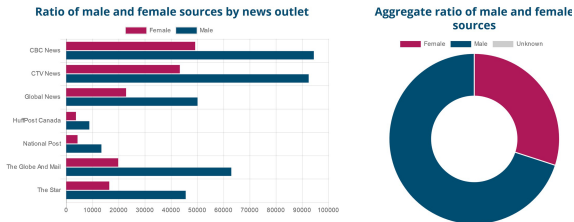


Fig 1. The Gender Gap Tracker online dashboard page. Reprinted from <https://gendergaptracker.informedopinions.org/> under a CC BY license, with permission from Informed Opinions, original copyright 2018.

Asr et al. (2021) on how Canadian newspaper quote men more often than women ([still online!](#))



Types of web-scraping

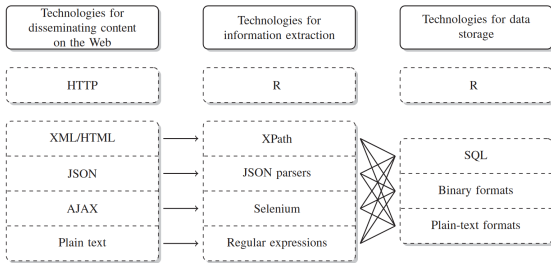


Figure 1.4 Technologies for disseminating, extracting, and storing web data

Munzert (2015, p. 10)



What we'll cover...

1. 'static' HTML structures (`rvest`)
 - all content is available by parsing the HTML
2. 'dynamic' webpages (`RSelenium`)
 - content is not visible in HTML source code but becomes available after interacting with the webpage



Disclaimer: Legal considerations I

Is web-scraping even allowed?

It depends...

- violation against terms of service?



Disclaimer: Legal considerations IV



- hiQ scraped data from LinkedIn public profiles
- LinkedIn tried to prevent that
- initial ruling: scraping allowed; later rulings: hiQ breached User Agreement → settlement

In the end, “[l]egalities depend a lot on where you live. However, as a general principle, if the data is public, non-personal, and factual, you’re likely to be ok” (Wickham et al. 2023)



HTML Basics

Websites are based on **HyperText Markup Language**

- HTML contains information about the structure of a web-page
- it's responsible for how content is graphically displayed

An example



HTML Elements and Attributes

- HTML consists of elements, tags and attributes
 - elements are the different components of a webpage (e.g. headlines, text, images)
 - elements are mostly embedded into tags (`<element>content</element>`) but some come without start and end tags
 - attributes are additional information of an element (e.g., the size of an image, the font type etc.)

We will only cover HTML in a superficial way, but try this [tutorial](#)



Important attributes

- *href*: weblink, always comes with a element
- *src*: source of an image



How to use HTML for scraping?

We need to identify the CSS selector of an element of interest.

There are **two** options:

- manual approach: cursor on element of interest > right-click > Inspect
- semi-automatic approach: download [SelectorGadget](#) or save as bookmark



Basic pipeline

Setup

- install SelectorGadget
- R library: `rvest`



Step 1

Download HTML source with `rvest`

```
1 > library(rvest)
2 > url <- "https://press.un.org/en/content/press-release"
3 > html <- read_html(url)
4 > html
5 {html_document}
6 <html lang="en" dir="ltr">
7 [1] <head>\n<meta http-equiv="Content-Type"
    content="text/html; charset=UTF-8">\n<meta
    charset="utf-8">\n<link rel="canonical" href="https:
    ...
8 [2] <body class="layout-one-sidebar
    layout-sidebar-first page-view-home-press
    path-content">\n    <div
    class="visually-hidden-focusable bg- ..
```



Step 2

Example: **UN** press releases

The screenshot shows the United Nations website at the URL `press.un.org/en/content/press-release`. The page features the United Nations logo and the text "United Nations Meetings Coverage and Press Releases". A navigation menu includes links for Home, Secretary-General, General Assembly, Security Council, Economic and Social Council, and International Court of Justice. The main content area is titled "Press Release" and displays a press release dated 22 May 2024 with the ID SC/15705. The title of the press release is "Security Council 2140 Committee Considers Panel of Experts' Midterm Update". The text of the press release states: "On 3 May 2024, the Security Council Committee established pursuant to resolution 2140 (2014) held informal consultations to consider the midterm update of the Panel of Experts, submitted in accordance with paragraph 3 of resolution 2707 (2023)." To the right of the main content, there is a "Highlights" section with two items: "Highlights of the Security Council Activities" and "Highlights of General Assembly 78th Session New York, 5 Sep. - 22 Dec. 2023". In the top right corner of the page, there is a language selector with "Français" and a search icon. A red arrow points to the "Français" link.



Step 3

The screenshot shows the United Nations website's 'Meetings Coverage and Press Releases' section. A green box highlights the 'Press Release' header. A red arrow points to the title of a press release: 'Security Council 2140 Committee Considers Panel of Experts' Midterm Update'. A grey box with the text 'click on element of interest' is positioned over the title. Another red arrow points to a search bar at the bottom right, with the text 'copy CSS selector' next to it. The search bar contains the text 'H:' and has buttons for 'Clear (1)', 'Toggle Position', 'XPath', and 'X'.

Try it out yourself!



Text

Here, we retrieve every paragraph on the webpage.

```
1 > library(rvest)
2 > (paragraphs <- read_html(url) %>%
3 +   html_elements("p") %>%
4 +   html_text())
5 [1] "On 3 May 2024, the Security Council Committee
     established pursuant to resolution 2140 (2014) held
     informal consultations to consider the midterm
     update of the Panel of Experts, submitted in
     accordance with paragraph 3 of resolution 2707
     (2023)."
```

```
6 [2] "United Nations Secretary-General Antonio Guterres
     flew from Muscat, Oman, to Manama, Bahrain, in the
     early evening of Wednesday, 15 May."
```




Tables

`rvest` has a pre-defined function `read_table()` to extract information from HTML tables

```
1 > html <- read_html(url2)
2 > table <- html %>%
3 +   html_element(".wikitable:nth-child(4)") %>%
4 +   html_table()
```



Images

For images, it is a bit more complicated.

1. open a session
2. retrieve image source link
3. download the image to your directory



Images II

```
1 > session <- session(url)
2 >
3 > # access links for image sources
4 > imgsrc <- session %>%
5 +   read_html() %>%
6 +   html_nodes("img") %>%
7 +   html_attr("src")
8 >
9 > # access page of source image (only retrieve the
  > # first image here)
10 > img <- session_jump_to(session, paste0(root_url,
  >   imgsrc[[1]]))
11 >
12 > # write to our project's directory
13 > writeBin(img$response$content, basename(imgsrc[1]))
```

Let's do it in \mathbb{R}



Loops I

Often, we want to automatize these steps for multiple pages.

Two options:

1. creating empty objects and filling them in a `for-loop`
2. defining a function, apply it and retrieve objects of interest from list

→ usually, functions are more versatile and can be easier run in parallel



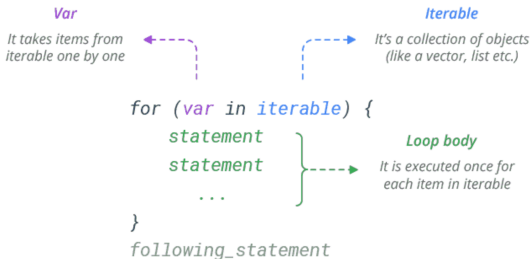
Loops II

Before creating a loop/function

1. check the pagination structure of a website (e.g. [United Nations](#) use `'?page=#'` to list results)
2. make sure which elements you need to retrieve (often, you just want links but some information, such as dates, may not be available on sub-pages, so you want to gather these as well)
3. test the pipeline on a single element before using it in a loop



For-loops



Tutorial on for-loops
Graph and guide to for-loops



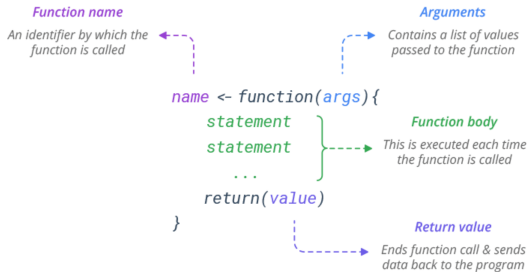
For-loops for automation

```
1
2 > urls <- c("https://www.eui.eu/en/home",
3             "https://www.ceu.edu/")
4 > links <- c()
5 > for(url in 1:length(urls)){
6 +     html <- read_html(urls[[url]])
7 +     links[url] <- html %>%
8 +       html_node("h1") %>%
9 +       html_text()
10 + }
11 > links
[1] "\nOpening doors to the future\n" "Apply Now for
    2024-2025!"
```

Let's do it in \mathbb{R}



Functions in R



Tutorial on Functions Graph and guide to functions



Functions for automation

```
1 > h1_scrape <- function(url){
2   +   html <- read_html(urls[[url]])
3   +   links[url] <- html %>%
4   +     html_node("h1") %>%
5   +     html_text()
6   + }
7 >
8 > (links <- sapply(1:length(urls), h1_scrape))
9 [1] "\nOpening doors to the future\n" "Apply Now for
    2024-2025!"
```

Let's do it in \mathbb{R}



Post processing

Regular expressions

Often, the data you retrieve are not in the desired format. **Regular expressions** are handy for data processing, they can match/extract/remove a pattern in text

```
1 > library(stringr)
2 > date_time <- "2009/03/12 20:12:31"
3 > str_extract(date_time, "\\d{4}/\\d{2}/\\d{2}")
4 [1] "2009/03/12"
```



Regular expressions in R

- base R includes regular expressions in functions like `gsub` (replacement), `regmatches` (extraction) or `regexr` (detection)
- `stringr` is the `tidyverse` approach to regex and provides equivalent functions such as `str_replace`, `str_extract` or `str_detect`



More information on regular expressions

- You can try out regular expressions in online tools like [these](#)
- [StringR cheat sheet](#)
- [A guide to regular expressions](#)



Dynamic Webpages

Sometimes, webpages change only when we interact with them in a browser session (e.g. by clicking on specific objects). For these **webpages**, `rvest` is not applicable.

Usually, you can identify these pages by their usage of active classes

```
<html lang="en"> event scroll
▶ <head> ... </head>
▼ <body lang="en" data-jsactive="true" data-n
  <input id="dev_product_uri" type="hidden"
  ▼ <div id="website"> overflow
    ▶ <div id="website-header" class="ep-head
      <!--END OF REFACTORING-->
```



Functionality of RSelenium

Solution?



RSelenium

RSelenium was initially built to test webpages, we use it to send commands to a virtual browser window (e.g. clicking a button)



Before using RSelenium

Whenever possible, try to use the static approach, it's easier and less prone to error

- check other webpages which may contain the same information
- look for other sub-directories that may store the same content
- check the [Wayback Machine](#)
- use the 'Search' bar of a webpage to list results



Pipeline of RSelenium and wdman

A typical pipeline looks as follows:

- configuration of a session (now more complex, we need `wdman` for this)
- opening browser
- navigate to URL
- accept/decline cookies
- identify css selector for button, click on it (repeat if necessary)
- retrieve object of interest (usually links)



Setup

- configure driver
- extract client
- navigate to webpage

```
1 > url <- "https://www.europarl.europa.eu/news/en"  
2 > rd <- rsDriver(browser = "firefox",  
3 +               chromever = NULL,  
4 +               port = sample(1:65535, 1),  
5 +               verbose = FALSE)  
6 > browser <- rd[["client"]]  
7 > browser$navigate(url)
```




How to extract the text of a dynamic webpage

- identify elements by css selector
- elements contain lots of other information
- to extract plain text, loop through the elements and extract only the text

```
1 > paragraphs <- browser$findElements(using = 'css
  selector', value='p')
2 > text <- c()
3 > for(i in 1:length(paragraphs)){
4   +   text[i] <- paragraphs[[i]]$getElementText()[[1]]
5   + }
6 > text[1]
7 [1] "On Thursday afternoon, the Eurovision debate
  between the lead candidates for the presidency of
  the European Commission took place in the European
  Parliament."
```


Let's do it in \mathbb{R}



Concluding remarks

That's it for today. We covered the first aspect of text analysis on how to retrieve data.

- in the lab session: practice scraping on your own projects
- tomorrow: bags-of-words approaches
 - topic modelling
 - prediction tasks

Thank you for your attention!

References I

- Asr, F. T., Mazraeh, M., Lopes, A., Gautam, V., Gonzales, J., Rao, P., & Taboada, M. (2021). The Gender Gap Tracker: Using Natural Language Processing to measure gender bias in media. *PLOS ONE*, *16*(1), e0245533.
<https://doi.org/10.1371/journal.pone.0245533>
- Bischof, D., & Kurer, T. (2023). Place-Based Campaigning: The Political Impact of Real Grassroots Mobilization. *The Journal of Politics*, *85*(3), 984–1002.
<https://doi.org/10.1086/723985>
- Munzert, S. (2015). *Automated Data Collection with R: A Practical Guide to Web Scraping and Text Mining* (1st ed.). Wiley.
- Wickham, H., Çetinkaya-Rundel, M., & Grolemund, G. (2023). *R for Data Science*. O'Reilly Media, Inc.
<https://r4ds.hadley.nz/>