Web Scraping Recap
○○○○○○○

Web Scraping Methods
○○○○○○○○○○○○○○○

Scraping in R: Core Functions w. Examples
○○○○○○○○○○○○○○○○○○○○○○○○○

# Class 7: Web Scraping in R
## MAST5953: Web Scraping and Text Mining

Dr. Miriam Sorace

www.miriamsorace.net

25 January 2023

Web Scraping Recap
0000000

Web Scraping Methods
000000000000000

Scraping in R: Core Functions w. Examples
0000000000000000000000000

## Outline of Today's Class

Web Scraping Recap

Web Scraping Methods

Scraping in R: Core Functions w. Examples

# Web Scraping Recap

# Web Scraping

- ▶ Scrapers: programmes that grab *specific* content from webpages
- ▶ Spiders (or Crawlers): programmes that index entire pages and follow every link.
- ▶ Web scraping involves some spidering when - for example - parsing entire web-pages. But scrapers are *selective/purposive*

# Web Scraping: Managing Expectations

▶ It is rarely a one-shot affair: websites change constantly

▶ Difficult to generalise: each website is different

▶ Writing a scraping code that works well takes a long time

# Web Scraping: Netiquette!

▶ Search for & respect terms of service
  ▶ add "/robot.txt" at the end of URL
▶ Careful with personal/private data
▶ Check if there are APIs or ready-to-download files already
▶ Stay identifiable (do not mask your IP)
▶ Reduce traffic as much as possible
  ▶ light/efficient scrapers
  ▶ use pauses (Sys.sleep())
  ▶ avoid over-scraping

# Web Scraping Steps

1. Analyse the structure of the webpage (e.g. HTML/XML trees)
    - ▶ Identify relevant nodes
    - ▶ www.selectorgadget.com
2. Test: write scraper and debug
3. Execute the download

# First Task: Install SelectorGadget!

# First Task: Install SelectorGadget!

Web Scraping Recap
0000000

Web Scraping Methods
●000000000000000

Scraping in R: Core Functions w. Examples
00000000000000000000000000

# Web Scraping Methods

Web Scraping Recap
0000000

Web Scraping Methods
0●000000000000000

Scraping in R: Core Functions w. Examples
000000000000000000000000000

## Extraction Methods

1. Regular Expressions
2. Node queries
3. APIs (*Application Programming Interfaces*)
   ▶ We'll explore them in the next class on social media data

# Regular Expressions

▶ Scraper works by extracting from meaningful alphanumeric patterns in the webpage text

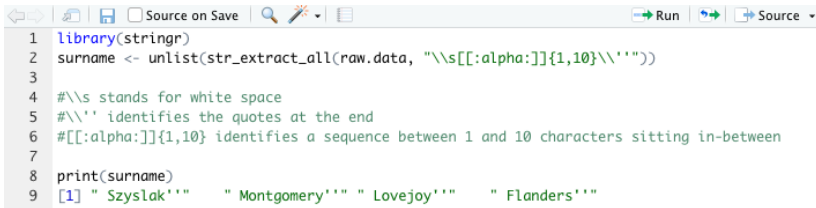    ▶ E.g.: grab all text in-between quotes ("").

# Regular Expressions

```
[1] "Moe Szyslak"        "Burns, C. Montgomery" "Rev. Timothy Lovejoy"
[4] "Ned Flanders"
```

▶ Suppose from the web content above, you wanted to extract the surnames

▶ You can exploit the regular pattern of surnames sitting between a space and closing with ''

## Regular Expressions

```
[1] "Moe Szyslak"           "Burns, C. Montgomery" "Rev. Timothy Lovejoy"
[4] "Ned Flanders"
```

▶ You would write a regular expression scraper as the below

```
      Source on Save                                    Run       Source
  1   library(stringr)
  2   surname <- unlist(str_extract_all(raw.data, "\\s[[:alpha:]]{1,10}\\''"))
  3
  4   #\\s stands for white space
  5   #\\'' identifies the quotes at the end
  6   #[[:alpha:]]{1,10} identifies a sequence between 1 and 10 characters sitting in-between
  7
  8   print(surname)
  9   [1] " Szyslak''"    " Montgomery''" " Lovejoy''"    " Flanders''"
```

# Regular Expressions

| | |
|---|---|
| [:digit:] | Digits: 0 1 2 3 4 5 6 7 8 9 |
| [:lower:] | Lowercase characters: a–z |
| [:upper:] | Uppercase characters: A–Z |
| [:alpha:] | Alphabetic characters: a–z and A–Z |
| [:alnum:] | Digits and alphabetic characters |
| [:punct:] | Punctuation characters: . , ; etc. |
| [:graph:] | Graphical characters: [:alnum:] and [:punct:] |
| [:blank:] | Blank characters: Space and tab |
| [:space:] | Space characters: Space, tab, newline, and other space characters |
| [:print:] | Printable characters: [:alnum:], [:punct:] and [:space:] |

*Source:* Adapted from http://stat.ethz.ch/R-manual/R-patched/library/base/html/regex.html

# Regular Expressions

| | |
|---|---|
| ? | The preceding item is optional and will be matched at most once |
| * | The preceding item will be matched zero or more times |
| + | The preceding item will be matched one or more times |
| {n} | The preceding item is matched exactly *n* times |
| {n, } | The preceding item is matched *n* or more times |
| {n,m} | The preceding item is matched at least *n* times, but not more than *m* times |

*Source:* Adapted from http://stat.ethz.ch/R-manual/R-patched/library/base/html/regex.html

**Table 8.3** Selected symbols with special meaning

| | |
|---|---|
| \w | Word characters: [[:alnum:]_] |
| \W | No word characters: [^[:alnum:]_] |
| \s | Space characters: [[:blank:]] |
| \S | No space characters: [^[:blank:]] |
| \d | Digits: [[:digit:]] |
| \D | No digits: [^[:digit:]] |
| \b | Word edge |
| \B | No word edge |
| \< | Word beginning |
| \> | Word end |

# Regular Expressions

| Function | Description | Output |
|---|---|---|
| *Functions using regular expressions* | | |
| str_extract() | Extracts first string that matches pattern | Character vector |
| str_extract_all() | Extracts all strings that match pattern | List of character vectors |
| str_locate() | Returns position of first pattern match | Matrix of start/end positions |
| str_locate_all() | Returns positions of all pattern matches | List of matrices |
| str_replace() | Replaces first pattern match | Character vector |
| str_replace_all() | Replaces all pattern matches | Character vector |
| str_split() | Splits string at pattern | List of character vectors |
| str_split_fixed() | Splits string at pattern into fixed number of pieces | Matrix of character vectors |
| str_detect() | Detects patterns in string | Boolean vector |
| str_count() | Counts number of pattern occurrences in string | Numeric vector |
| *Further functions* | | |
| str_sub() | Extracts strings by position | Character vector |
| str_dup() | Duplicates strings | Character vector |
| str_length() | Returns length of string | Numeric vector |
| str_pad() | Pads a string | Character vector |
| str_trim() | Discards string padding | Character vector |
| str_c() | Concatenates strings | Character vector |

Web Scraping Recap
0000000

Web Scraping Methods
00000000●0000000

Scraping in R: Core Functions w. Examples
0000000000000000000000000

# Regular Expressions
Advantages

- ▶ Only useful when HTML/XML are malformed **and** clear patterns are retrievable in the page
- ▶ Mostly used in creating URL lists to scrape multiple pages
- ▶ Mostly used in data cleaning/variable creation

# Regular Expressions
Disadvantages

▶ Building and interpreting a regular expression scraper is challenging: hard to test and debug

▶ Limited applicability: not all web pages have consistent repeated patterns that *uniquely* identify a specific content

▶ Ignores the useful and meaningful hierarchical structure of web pages

▶ Strongly discouraged for webscraping, best for data cleaning

Web Scraping Recap
0000000

Web Scraping Methods
0000000000●00000

Scraping in R: Core Functions w. Examples
000000000000000000000000

# Node Queries

- ▶ This is the most common procedure
- ▶ The scraper is built by exploiting the tree structure of the web page
    - ▶ Inspected via web developer tools & SelectorGadget
- ▶ Utilises query functions from R library such as `rvest` to locate nodes and extract information

# Node Queries
Advantages

▶ Intuitive to write and read

▶ Uniquely identifies the relevant content

Web Scraping Recap
○○○○○○○

Web Scraping Methods
○○○○○○○○○○○○●○○○

Scraping in R: Core Functions w. Examples
○○○○○○○○○○○○○○○○○○○○○○○○○

# Node Queries
Disadvantages

▶ Costly when the HTML/XML structure is malformed
▶ Costly when web pages get constantly updated

## Application Programming Interfaces

▶ Scraper is built by using APIs 'wrappers', exploiting the webpage's API, which provides the website's data in pure/cleaned form

▶ We'll see concrete examples of this in the next class where we'll look at social media data

# Application Programming Interfaces
## Advantages

▶ This is the scraping gold standard.

▶ Very easy: clean, standardised data

▶ APIs codes are robust and regularly updated

# Application Programming Interfaces
Disadvantages

▶ Not many websites provide their APIs

▶ Websites that do, often have licences/fees/limits attached to API use

▶ API owners often limit their functionality from one day to the next, rendering the scraper useless.

Web Scraping Recap
OOOOOOO

Web Scraping Methods
OOOOOOOOOOOOOOO

Scraping in R: Core Functions w. Examples
●OOOOOOOOOOOOOOOOOOOOOOOO

Scraping in R: Core Functions w. Examples

# Parsing the Page

### read_html

#### Read In .Html Content

Read in the content from a .html file. This is generalized, reading in all body text. For finer control the user should utilize the xml2 and rvest packages.

**Keywords**    html

#### Usage

```
read_html(file, skip = 0, remove.empty = TRUE, trim = TRUE, ...)
read_xml(file, skip = 0, remove.empty = TRUE, trim = TRUE, ...)
```

#### Arguments

| | |
|---|---|
| **file** | The path to the .html file. |
| **skip** | The number of lines to skip. |
| **remove.empty** | logical. If TRUE empty elements in the vector are removed. |
| **trim** | logical. If TRUE the leading/training white space is removed. |
| **...** | Other arguments passed to xml2::read_html(). |

#### Value

Returns a character vector.

# Extracting Tables

### html_table

**Parse An Html Table Into A Data Frame.**

Parse an html table into a data frame.

#### Usage

```
html_table(x, header = NA, trim = TRUE, fill = FALSE, dec = ".")
```

#### Arguments

| | |
|---|---|
| **x** | A node, node set or document. |
| **header** | Use first row as header? If `NA`, will use first row if it consists of `<th>` tags. |
| **trim** | Remove leading and trailing whitespace within each cell? |
| **fill** | If `TRUE`, automatically fill rows with fewer than the maximum number of columns with `NA`s. |
| **dec** | The character used as decimal mark. |

#### Assumptions

`html_table` currently makes a few assumptions:

- No cells span multiple rows
- Headers are in the first row

# Extracting Tables
Example

# Extracting Tables
Example

- ▶ url<-"https:
  //finance.yahoo.com/quote/AAPL/history?p=AAPL"
- ▶ apple finance<- read html(url)
- ▶ tableaslist<-html table(apple finance)
- ▶ table<-as.data.frame(tableaslist)

# Extracting Tables
Example

# Extracting Text

## html_text

From rvest v0.3.6
by Hadley Wickham

**Extract Attributes, Text And Tag Name From Html.**

Extract attributes, text and tag name from html.

**Usage**

```
html_text(x, trim = FALSE)

html_name(x)

html_children(x)

html_attrs(x)

html_attr(x, name, default = NA_character_)
```

**Arguments**

| | |
|---|---|
| **x** | A document, node, or node set. |
| **trim** | If `TRUE` will trim leading and trailing spaces. |
| **name** | Name of attribute to retrieve. |
| **default** | A string used as a default value when the attribute does not exist in every node. |

**Value**

`html_attr` , `html_tag` and `html_text` , a character vector; `html_attrs` , a list.

**Examples**

```
# NOT RUN {
movie <- read_html("https://en.wikipedia.org/wiki/The_Lego_Movie")
cast <- html_nodes(movie, "tr:nth-child(8) .plainlist a")
html_text(cast)
```

# Extracting Text
## Example

# Extracting Text
## Example

# Extracting Text
Example

- url<-"https:
  //www.europarl.europa.eu/doceo/document/
  CRE-9-2020-06-17-INT-1-161-0000_EN.html"

- v_speech1<-read_html(url)

- text <- v_speech1 %>%
    html_nodes(".contents") %>%
    html_text() %>%
    as.data.frame()

# Extracting Text
## Example

1

        Guy Verhofstadt,   on behalf of the Renew Group. - Madam President, the reason for
this resolution is very simple: it is not a resolution to the Commission, it is not a resolution
to the colleagues in the Parliament, it is a resolution that is in fact made for the Council. I
hope that the Council this time can reach an agreement so that we can launch this conference, be
cause let's be honest, it becomes more and more like the monster of Loch Ness. From time to time
it appears, then it disappears, it is more like an illusion, like a fantasy and we must avoid th
at.
2                          This is serious business, this conference is not one or other li
ttle instrument for the European Parliament. This conference is vital, it is crucial for the futu
re of our European Union, because let's face it, this Covid crisis is another illustration of it.
This Covid crisis will change the world. It is a world that will be completely different, a worl
d that will no longer be dominated, I think, by the United States of America who is withdrawing f
rom the international scene for the moment. It may well be dominated by China who wants to becom
e, and who already is, a world power and who wants to dominate the world and international scene.

3 Between America and China there will be Europe stuck between them and if we don't want to be st
uck between the US at the one hand and China on the other hand, and if we want really the instrum
ents to defend the interests of our citizens, we need another European Union. This European Union
is not fit for purpose. Not fit for the future and we all know it. The Council knows it, the Com
mission knows it and certainly our citizens know it. So it's time now to start really this exerci
se because we have not a lot of time, maybe we can tell that to your colleagues in the Council. T
here is the Covid crisis, there are precedent crises, tomorrow there could be another challenge.

# Extracting Links

▶ Use `html_attr("href")` after identifying the relevant node

```
html_attrs(x)

html_attr(x, name, default = NA_character_)
```

### Arguments

**x**        A document, node, or node set.

**trim**     If `TRUE` will trim leading and trailing spaces.

**name**     Name of attribute to retrieve.

**default**  A string used as a default value when the attribute does not exist in every node.

### Value

`html_attr` , `html_tag` and `html_text` , a character vector; `html_attrs` , a list.

### Examples

```
# NOT RUN {
movie <- read_html("https://en.wikipedia.org/wiki/The_Lego_Movie")
cast <- html_nodes(movie, "tr:nth-child(8) .plainlist a")
html_text(cast)
html_name(cast)
html_attrs(cast)
html_attr(cast, "href")
# }
```

# Extracting Links
## Example

# Extracting Links
Example

# Extracting Links
Example

▶ url<-"https://www.europarl.europa.eu/meps/en/
97058/GUY_VERHOFSTADT/main-activities/
plenary-speeches#detailedcardmep"

▶ verhof_speech <- read_html(url)

▶ speech_links <- verhof_speech %>%
  html_nodes(".erpl_search-results-list-expandable-bloc
a") %>%
  html_attr("href") %>%
  as.data.frame()

# Extracting Links
Example

| ▲ | . |
|---|---|
| 1 | https://www.europarl.europa.eu/doceo/document/C... |
| 2 | https://www.europarl.europa.eu/doceo/document/C... |
| 3 | https://www.europarl.europa.eu/doceo/document/C... |
| 4 | https://www.europarl.europa.eu/doceo/document/C... |
| 5 | https://www.europarl.europa.eu/doceo/document/C... |
| 6 | https://www.europarl.europa.eu/doceo/document/C... |
| 7 | https://www.europarl.europa.eu/doceo/document/C... |
| 8 | https://www.europarl.europa.eu/doceo/document/C... |
| 9 | https://www.europarl.europa.eu/doceo/document/C... |
| 10 | https://www.europarl.europa.eu/doceo/document/C... |

## Looping over list of URLs

▶ This is useful in situations where urls have patterns

▶ E.g. In EP website, MEPs urls are identical barring from MEP number and name.

▶ A scraper can exploit this and loop over a series of webpages with identical content

▶ This 'simply' requires creating a list of univocal url codes, pasting those to the unchanging - generic - url and running a standard `for` loop

▶ The web scraper we built will be inside the `for` loop

## Looping over list of URLs
Example

- ▶ Consider the following urls:
    1. https://www.europarl.europa.eu/meps/en/97058/GUY_
       VERHOFSTADT/main-activities/plenary-speeches#
       detailedcardmep/
    2. https://www.europarl.europa.eu/meps/en/197395/
       ALICE_KUHNKE/main-activities/plenary-speeches#
       detailedcardmep
    3. https://www.europarl.europa.eu/meps/en/124846/
       PINA_PICIERNO/main-activities/plenary-speeches#
       detailedcardmep

# Looping over list of URLs
Example

```
#### loop through URLs

MEPcodes <- c("97058/GUY_VERHOFSTADT", "197395/ALICE_KUHNKE", "124846/PINA_PICIERNO")

urls <- paste0("https://www.europarl.europa.eu/meps/en/", MEPcodes , "/main-activities/plenary-speeches#detailedcardmep" )

catcherlist<-list()

for (i in urls) {
 page <- read_html(i)
 Name <- page %>% html_nodes("#presentationmep .erpl_title-h1") %>% html_text() %>% as.character()
 Title <- page %>% html_nodes(".erpl_search-results-list-expandable-block .t-item") %>% html_text() %>% as.character()
 Link <- page %>% html_nodes(".erpl_search-results-list-expandable-block a") %>% html_attr("href") %>% as.character()
 temp <- list(Name, Title, Link)
  catcherlist <- rbind(catcherlist,temp)
}

df<-as.data.frame(catcherlist)
```

# Looping over list of links
Example

# Looping over list of links
Example

```
##follow links

#Setup empty data frame

catcher_text <- data.frame(Name=character(),Text=character())


for (i in my_data$links) {

  page <- read_html(i)
  Name<-page %>% html_nodes(".doc_subtitle_level1_bis .bold") %>% html_text() %>% as.character()
  Text<-page %>% html_nodes(".contents") %>% html_text() %>% as.character()
  temp_text <- data.frame(Name, Text) #fill temporary repository
  catcher_text <- rbind(catcher_text,temp_text) #convert into dataframe
}
```

# Clicking Buttons
E.g. "Load More"

# RSelenium

- ▶ *Dynamic* interaction (filling fields, clicking buttons) with the website is not possible through rvest
- ▶ You need RSelenium for this: but **outside of the scopes of this course**
- ▶ RSelenium requires a bit of environment preparation:
  - ▶ It requires selenium initiation, by opening a browser session (not simply the specification of a url)
  - ▶ Only after setting up the 'Remote Server/Driver' it can navigate the specific web page ... v. time consuming
- ▶ The core functions to perform a button click are RSelenium's $findElement and $clickElement
- ▶ If you want to know more: http://joshuamccrain.com/ tutorials/web_scraping_R_selenium.html

## What we have learnt today ...

▶ SelectorGadget can make the first step of scraping (web-page analysis) so much easier

▶ There are 3 core strategies in scraping
  ▶ Regular expression (worst case)
  ▶ Node queries (typical case)
  ▶ APIs (gold standard but rare)

▶ The core functions in R to perform **node query** scraping - using the package rvest

▶ Awareness of RSelenium and the problems in scraping dynamic pages

▶ Next week: more on scraping via APIs: we're going to learn how to scrape Twitter data