# Web Scraping for Sports Data with R

Yaqiong Yao

10/1/2020

# Outline

- Introduction

- Using R
  - Import files downloaded from websites
  - Static data
  - Dynamic data

- Summary

# Introduction

- Web scraping technique is used for capturing data from websites.

- Motivation of Web Scraping
  - Need to extract data from websites
  - A reproducible way of capturing data online

- Prerequisite
  - Having experience with R
  - A laptop with R and R studio installed

# Example

College basketball school index

- These data can be obtained by copy and paste manully.

- Web scraping technique helps capture the data efficiently.

# Web Scraping Using R

- Different web scraping techniques are required to deal with different situations of data in R.

- Data have been organized into files.
  - Directly download it and read it in R

- Data are contained in HTML pages.
  - Static data
  - Dynamic data

# Import Data Files from Websites

- These files that can be read by **read.csv** or related functions.

- They can be directly imported from a URL.

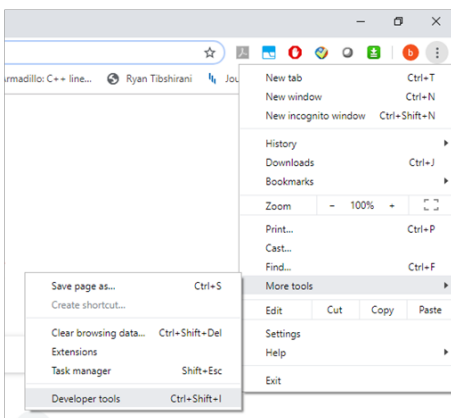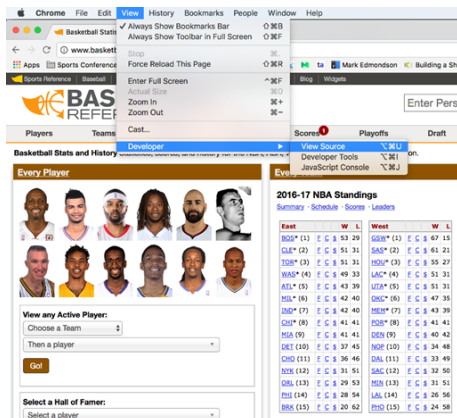- Example: we extract the most recent Australian Open Tennis Championships match (AUS Open):

```
url <- "http://www.tennis-data.co.uk/2020/ausopen.csv"
tennis_aus <- read.csv(url)
str(tennis_aus)
```

# Static Data and Dynamic Data

- Most of data in the web are not organized into files, which can be directly imported into R.

- Before we capture these data, we need to determined whether the data is static or dynamic based on the source code.

- Static data is the data that can be seen in the source code.

- We cannot see the dynamic data in the source code.

# Static Data and Dynamic Data

- The source code can be accessed by View → Developer → View Source in Chrome. Or right click the website and choose "View Page Source".

# Static Data and Dynamic Data

Exerciese: Determine what kind of the data are in the following examples, static or dynamic.

- http://tennisabstract.com/reports/atp_elo_ratings.html

- https://www.espncricinfo.com/ci/content/stats/

# Static Data and Dynamic Data



This is static data.

# Static Data and Dynamic Data



This is dynamic data.

# Web Scraping for Static Data in R

R provides several approaches for web scraping static data. Two of them will be discussed in this workshop.

- **readLines** function: Read the source code of the HTML page.

- **rvest** package: Capture useful data by identifying the elements contains the data in the source code.

# Web Scraping for Static Data in R

Use **readLines** function for College basketball school index.

```
web_page <- readLines("https://www.sports-reference.com/cbb/schools/")
head(web_page, n = 10L)
```

```
## [1] ""
## [2] "<!DOCTYPE html>"
## [3] "<html data-version=\"klecko-\" data-root=\"/home/cbb/build\" itemscope itemtype=\"https://schema.org/W
## [4] "<head>"
## [5] "    <meta charset=\"utf-8\">"
## [6] "    <meta http-equiv=\"x-ua-compatible\" content=\"ie=edge\">"
## [7] "    <meta name=\"viewport\" content=\"width=device-width, initial-scale=1.0, maximum-scale=2.0\" />"
## [8] "    <link rel=\"dns-prefetch\" href=\"https://d2p3bygnnzw9w3.cloudfront.net/req/202009101\" />"
## [9] ""
## [10] "    <title>School Index | College Basketball at Sports-Reference.com</title>"
```

- Gives the source code.

- Needs data cleaning and organization.

# Web Scraping for Static Data in R

Before we talk about web scraping by **rvest** package, we need to know how to locate the elements containing the data in the source code.

- Right click the page and choose "Inspect".

- Clike "Select an element in the page to inspect it".

- We can locate the element by CSS selector or XPATH.

# Web Scraping for Static Data in R

Use http://tennisabstract.com/reports/atp_elo_ratings.html as an example

- CSS selector: id = "reportable", class = "tablesorter"

# Web Scraping for Static Data in R

- XPATH: '//*[@id="reportable"]'

# Web Scraping for Static Data in R

Next, we are going to talk about how to use **rvest** for web scraping by using an example.

- Install **rvest** package from cran.

```
install.packages("rvest", repos = "http://cran.us.r-project.org")
require("rvest")
```

# Web Scraping for Static Data in R

- Web scraping data from
  http://tennisabstract.com/reports/atp_elo_ratings.html

```r
url_elo <- "http://tennisabstract.com/reports/atp_elo_ratings.html"
webpage <- read_html(url_elo)
elo_class <- webpage %>%
  html_nodes(".tablesorter") %>%
  html_table()
elo_id <- webpage %>%
  html_nodes("#reportable") %>%
  html_table()
identical(elo_class, elo_id)
```

```
## [1] TRUE
```

# Web Scraping for Static Data in R

```
elo_xpath <- webpage %>%
  html_nodes(xpath = '//*[@id="reportable"]') %>%
  html_table()
identical(elo_class, elo_xpath)
```

```
## [1] TRUE
head(elo_class[[1]])
```

```
##   Rank              Player  Age    Elo   HardRaw ClayRaw GrassRaw      hElo
## 1    1      Novak Djokovic 33.3 2255.4 NA  2142.9  2085.6   2013.9 NA 2199.1
## 2    2        Rafael Nadal 34.3 2185.0 NA  2045.2  2111.2   1677.9 NA 2115.1
## 3    3       Roger Federer 38.5 2170.0 NA  2051.7  1824.3   1933.8 NA 2110.9
## 4    4       Dominic Thiem 27.0 2079.8 NA  1989.8  2009.2   1614.3 NA 2034.8
## 5    5       Andrey Rublev 22.9 2023.5 NA  1910.8  1785.6   1516.4 NA 1967.2
## 6    6 Stefanos Tsitsipas 22.1 2022.2 NA  1939.0  1898.9   1573.1 NA 1980.6
##    cElo   gElo         Peak Match Peak Age Peak Elo
## 1 2170.5 2134.7 NA    2016 Miami F     28.8   2469.7
## 2 2148.1 1931.4 NA   2009 Madrid SF    22.9   2368.4
## 3 1997.1 2051.9 NA     2007 Dubai F    25.6   2379.4
## 4 2044.5 1847.0 NA    2016 Halle R16   22.8   2122.5
## 5 1904.5 1769.9 NA   2020 Hamburg F    22.9   2023.5
## 6 1960.5 1797.6 NA 2020 Cincinnati R16 22.0   2069.1
```

# Web Scraping for Static Data in R

- Except **html_nodes** and **html_table**, there are many other frequently used functions in **rvest**.
    - **html_node** : extract element
    - **html_text** : extract text
    - **html_attrs** : extract attributes
    - **html_form** : extract forms
- Please look up rvest cran for more information.
- SelectorGadget is a convenient tool to identify CSS selector.

# Web Scraping for Dynamic Data in R

- The dynamic data displayed in the website can be different because the website may provide user interaction.

- We need to automate the web browsing process in R for the dynamic data.

- **RSelenium** package helps this automating process by providing connection to Selenium Server.

- Install **RSelenium** package.

```r
devtools::install_github("ropensci/RSelenium")
require("RSelenium")
```

# Web Scraping for Dynamic Data in R

- Use **RSelenium** to extract data on 2017 Australian Open Final

# Web Scraping for Dynamic Data in R

- Connect to a selenium server and open brower.

```r
rD <- rsDriver(port = 5560L, chromever = "85.0.4183.87")
remDr <- rD$client
```

- Extract Information and organize data.

```r
url <- "http://www.flashscore.com/match/Cj6I5iL9/#match-statistics;0"
remDr$navigate(url)
webElem <- remDr$findElements(using = 'class', "statBox")
webElem <- unlist(lapply(webElem, function(x){x$getElementText()}))[[1]]
head(unlist(strsplit(webElem, split = '\n')))
```

```
## [1] "Service"        "20"            "Aces"          "4"
## [5] "3"             "Double Faults"
```

```r
remDr$close()
```

# Web Scraping for Dynamic Data in R

- Frequently used functions of **RSelenium**:
  - ▶ rsDriver() : start a selenium server
  - ▶ navigate() : navigate web pages
  - ▶ findElements() : find elements by CSS seclector or XPATH
  - ▶ getPageSource() : get current page source
  - ▶ clickElement() : click element

- Please go to RSelenium cran for more details.

# Web Scraping for Dynamic Data in R

Exercise: Web Scraping for the history basketball recording of UConn

https://www.flashscore.com/team/connecticut-huskies/8rqVf3Tj/results/

- Start a selenium server and open web brower.

```
require("RSelenium")
rD <- rsDriver(port = 5533L, chromever = "85.0.4183.87")
remDr <- rD$client
url <- "https://www.flashscore.com/team/connecticut-huskies/8rqVf3Tj/result
remDr$navigate(url)
```

# Web Scraping for Dynamic Data in R

- Automate to click all "show more results".

```
repeat{
  x <- try(webElemMore <-
              remDr$findElement(using = 'xpath',
                                '//*[@id="live-table"]/div[1]/div/div/a'),
           silent=T)
  if (inherits(x, "try-error")) break
  webElemMore$clickElement()
}
```

- Extract data, such as time, home/away, score and result.

```
webElemTime <- remDr$findElements(using = 'xpath',
                                  '//*[@class="event__time"]')
webElemTime <-
  unlist(lapply(webElemTime, function(x){x$getElementText()}))
webElemTime <- gsub("\\n", " ", webElemTime)
```

# Web Scraping for Dynamic Data in R

```r
webElemHome <-
  remDr$findElements(using = 'class',
                     'event__participant')
webElemHome <-
  unlist(lapply(webElemHome, function(x){x$getElementText()}))

webElemScore <-
  remDr$findElements(using = 'class', 'event__score')
webElemScore <-
  unlist(lapply(webElemScore, function(x){x$getElementText()}))

webElemResult <-
  remDr$findElements(using = 'class', 'wld')
webElemResult <-
  unlist(lapply(webElemResult, function(x){x$getElementText()}))
```

# Web Scraping for Dynamic Data in R

- Organize dataset.

```
n <- length(webElemHome)
basketball <-
  data.frame(time = webElemTime,
             Home = webElemHome[seq(n) %% 2 == 1],
             Away = webElemHome[seq(n) %% 2 == 0],
             HomeS = webElemScore[seq(n) %% 2 == 1],
             AwayS = webElemScore[seq(n) %% 2 == 0],
             Result = webElemResult)

head(basketball)
remDr$close()
```

```
              time           Home          Away HomeS AwayS Result
1      08.03. 16:00         Tulane        UConn    76    80      W
2      05.03. 19:00          UConn      Houston    77    71      W
3      29.02. 14:00  East Carolina        UConn    63    84      W
4      26.02. 19:00          UConn  UCF Knights    81    65      W
5      23.02. 14:00          UConn South Florida    78    71      W
6 20.02. 19:00 AOT         Temple        UConn    93    89      L
```

# Summary

- For different kinds of data, we need to use different web scraping techiniques with R.

- One can simply use **read.csv** or related functions to directly import organized files from web pages.

- The static data can be extract with the help of **rvest**.

- We could use **RSelenium** to parse the dynamic data.

# Resources

- CSS and HTML crash course

- rvest

- RSelenium

- R task view: web technology