# Lecture 6: Data Acquisition

James Sears*

AFRE 891/991 SS 25

Michigan State University

*Parts of these slides are adapted from **"Data Science for Economists"** by Grant McDermott and **"Advanced Data Analytics"** by Nick Hagerty.

# Table of Contents

# Prologue

# Prologue

Let's load in the packages that we'll need for today.

```
pacman :: p_load(lubridate, rvest, stringr, tidyverse)
```

# Data Acquisition

In order to wrangle, clean, or visualize data, we first need... data.

# Where Do Data Come from?

There is a whole spectrum, from DIY to plug-and-chug.

**1. Pre-cleaned datasets posted on secondary repositories**

- General and journal replication repositories
  - e.g. **Harvard Dataverse**
- Community-based repositories
  - e.g. **Kaggle**
- Public GitHub repositories
  - e.g. **Johns Hopkins COVID-19 caseloads, deaths, and vaccination data**

# Where Do Data Come from?

There is a whole spectrum, from DIY to plug-and-chug.

**2. Open data libraries**

- US Gov't: **Data.gov**
- Institutions, non-profits, and thinktanks
    - e.g. **World Bank Open Data**, **Pew Research**, **NBER Public Use Data**, **Economic Policy Institute**, **AEA Data Sources**
- Tech and other sites
    - e.g. **data.world**, **Our World in Data**, **Stanford Large Network Datasets**
- **Map of 2600+ Open Data Portals Worldwide**

# Where Do Data Come from?

There is a whole spectrum, from DIY to plug-and-chug.

## 3. Websites of primary data providers

- Government agencies; some private companies and NGOs; scientific researchers.
    - i.e. **EPA + USGS Water Quality Portal**

## 4. Programmatic data access through Application Program Interfaces (APIs)

- We'll talk more about these next lecture

# Where Do Data Come from?

There is a whole spectrum, from DIY to plug-and-chug.

**5. Liberate previously inaccessible data**

- Build relationships with people in government or the private sector.
- Find the right person, cold-email them and ask politely.
  - i.e. my master's thesis!
- File a Freedom of Information Act (FOIA) request (a last resort; very aggressive move).

# Where Do Data Come from?

**6. Compile data yourself**

- Assemble systematic information from many disparate sources.
- E.g. historical archives, websites, PDF reports.

**7. Collect your own primary data**

- Run surveys or experiments.

# Where to Look for Data?

There is **no "one-stop shop"**. Where to look entirely depends on your topic.

For economics research:

- **Search the literature:** Find papers related to your topic and check the Data section.
  - Good for learning the "standard" sources for common things (e.g., weather data).
- **Find your own data** that the literature hasn't used yet.
  - It's hard to find a novel use for an already widely-used dataset.
  - Cross-disciplinary arbitrage
- **Combine data in new ways:** most new projects will require joining data from 2+ sources.
  - E.g. state-level policy changes + household-level outcome data.

# Where to Look for Data?

There is **no "one-stop shop"**. Where to look entirely depends on your topic.

A few useful starting points:

- **"How to Find Data: Tips for Finding Data"** (Davidson College Library).
- **"Data Sets for Quantitative Research"** (University of Missouri Libraries).
- **Google Dataset Search**, **AWS Data Exchange**

The majority of data sources described above have the data easily accessible once found.

If the data aren't already machine readable, then we can take advantage of...

# Intro to Web Scraping

# Web Scraping

**Web scraping** is the process of **extracting semi-structured web data** and **converting into a structured dataset**

- Useful when information is already online but not available in a handy format.

# Structure of Webpages

Webpages are largely made out of **two types of files** that we have to parse:

**HTML**

- **H**yper**T**ext **M**arkup **L**anguage is a **markup language**, Like Markdown
- It specifies the **structure** of a webpage.

**CSS**

- **C**ascading **S**tyle **S**heets is a language for **formatting the appearance of a webpage**
  - CSS **properties** specify **how** to format: what font, what color, how wide
  - CSS **selectors** specify **what** to format: which structural elements get what rule

# How Websites Render Content

It's also worth realising that there are **two ways** that web content gets rendered in a browser:

1. **Server-side (back-end)**
2. **Client-side (front-end)**

You can read **here** for more details (including example scripts), but for our purposes the essential features are as follows...

# Server-Side Content Rendering

The scripts that build a **server-side** website aren't run on our computer

- Rather, the script that "builds" the site is **run on a host server**
  - All information is directly embedded in the webpage's HTML
- e.g. Wikipedia tables are already populated with all the info we see in our browser

| No. | Peak | Range (or island) | Location | Coordinates [1] | Prominence (m) | Height (m) | Col (m) |
|---|---|---|---|---|---|---|---|
| 1. | Mount Everest | Himalayas | China Nepal | 27°59'17.6500"N 86°55'30.0652"E | 8,848.86 | 8,848.86 | 0 |
| 2. | Aconcagua | Andes | Argentina | 32°39'11"S 70°0'42"W | 6,960.8 | 6,960.8 | 0 |
| 3. | Denali / Mount McKinley* | Alaska Range | United States | 63°4'10"N 151°0'26"W | 6,155 | 6,191 | 47 |
| 4. | Mount Kilimanjaro* | Eastern Rift mountains | Tanzania | 3°4'0"S 37°21'33"E | 5,885 | 5,895 | 10 |
| 5. | Pico Cristóbal Colón | Sierra Nevada de Santa Marta | Colombia | 10°50'18"N 73°41'12"W | 5,509 | 5,700 | 191 |

The 125 most topographically prominent summits on Earth

# Server-Side Content Rendering

The scripts that build a **server-side** website aren't run on our computer

- Rather, the script that "builds" the site is **run on a host server**
  - All information is .hi-blue[directly embedded] in the webpage's HTML
- e.g. Wikipedia tables are already populated with all the info we see in our browser

- **Webscraping process:** finding the correct selectors (CSS or Xpath), iterating through (dynamic) webpages (e.g. "Next page" and "Show more" tabs)
- **Key concepts:** CSS, Xpath, HTML

# Client-Side Content Rendering

The scripts that build a **client-side** website aren't run on our computer

- Website contains an empty template of HTML and CSS
    - May contain a "skeleton" table without any values
- When we visit the page URL, our browser sends a **request** to the host server
- If everything is okay (e.g. our request is valid), then the server sends a **response** script, which our browser executes and uses to populate the HTML template with the specific information that we want.

# Client-Side Content Rendering

- If everything is okay (e.g. our request is valid), then the server sends a **response** script, which our browser executes and uses to populate the HTML template with the specific information that we want.

# Web Scraping

Over the next two lectures we'll cover the main differences between the two approaches and general workflows.

However, I want to forewarn that web scraping typically involves a fair bit of **detective work**, iterating and adjusting steps

- According to the type of data you want
- To match the specifics of a given website

In short, web scraping involves **as much art as it does science.**

The good news, though: **if you can see it, you can scrape it.**

# Ethical + Legal Considerations

The last line brings up an important consideration: just because you **can** scrape it, doesn't mean you **should**.

## Legality

- In short, it is **currently legal** to scrape data from the web using automated tools as long as the data are **publicly available** (hiQ Labs vs. Linkedin)
  - We'll chat more later on about what this means for "hidden" APIs
- May get blocked due to violating a site's Terms of Service preventing scraping

## Ethicality

- Need to consider the impact your scraper will have on the host server
  - Easy to write a function that can overwhelm a website's host with rapid requests
  - We'll return to the "be nice" mantra later on

# Scraping Static Websites

# Static Scraping: Preliminaries

Today we'll be using **SelectorGadget**, which is a Chrome extension that makes it easy to discover CSS selectors.

- Install the extension directly **here**.

Please note that SelectorGadget is only available for **Chrome**. If you prefer using **Firefox**, then you can try **ScrapeMate**.

# Static Scraping: Preliminaries

The primary R package that we'll be using today is **rvest**, a simple webscraping library inspired by Python's **Beautiful Soup**, but with extra tidyverse functionality.

**rvest** is designed to work with webpages that are built server-side and thus requires knowledge of the relevant CSS selectors...

Which means that now is probably a good time for us to cover what these are in more detail.

# CSS Selectors and SelectorGadget

CSS **selectors** are the **"what"** of the display rules. They identify which rules should be applied to which elements.

- E.g. Text elements that are selected as **".h1"** (i.e. top line headers) are usually larger and displayed more prominently than text elements selected as **".h2"** (i.e. sub-headers).

The key point is that if you can **identify the CSS selector(s)** of the content you want, then you can **isolate it from the rest** of the webpage content that you don't want.

# CSS Selectors and SelectorGadget

This where SelectorGadget comes in. We'll work through an extended example (with a twist!) below, but I highly recommend looking over this **quick vignette** soon.

Here are two helpful links if you're interested in reading more about **CSS** (i.e Cascading Style Sheets) and **SelectorGadget**

27 / 48

# Application 1: Wikipedia

Let's say you watched the U.S. Olympic Marathon Trials last year and now want to scrape the wikipedia page on **marathon world record progression**

- Women's record: 3:40:22 in 1926 to 2:11:53 last year!

First, open up this page in your browser. Take a look at its structure: What type of objects does it contain? How many tables does it have? Do these tables all share the same columns? What row- and columns-spans? Etc.

# Application 1: Wikipedia

# Application 1: Wikipedia

Once you've familiarised yourself with the structure, read the whole page
into R using the `rvest::read_html()` function.

```
mthn = read_html("https://en.wikipedia.org/wiki/Marathon_world_record_prog
mthn
```

```
## {html_document}
## <html class="client-nojs vector-feature-language-in-header-enabled vector-fe
## [1] <head>\n<meta http-equiv="Content-Type" content="text/html; charset=UTF-
## [2] <body class="skin--responsive skin-vector skin-vector-search-vue mediawi
```

# Application 1: Wikipedia

As you can see, this is an **XML** document that contains **everything** needed to render the Wikipedia page.[1]

It's kind of like viewing someone's entire LaTeX document (preamble, syntax, etc.) when all we want are the data from some tables in their paper.

[1] XML stands for Extensible Markup Language and is one of the primary languages used for encoding and formatting web pages.

# First Table: Women's Records

Let's start by scraping our first table from the page, which documents the **women's record progression**.

The first thing we need to do is identify the table's unique CSS selector using SelectorGadget.

**Note:** this *will* require trial and error - and a lot of clicking.

Start by activating SelectorGadget and clicking on a chart element.

Clicking on another chart element expands the selection (but too much!)

Click the elements you *don't* want until they turn red:

# First Table: Women's Records

Working through this iterative process yields

`"div+ .wikitable :nth-child(1)"`.

We can use this unique CSS selector to isolate the women's record table content from the rest of the HTML document.

- **Extract the table content** with `html_element()`
- Parse the **HTML table into an R data frame** with `html_table()`

# First Table: Women's Records

- **Extract the table content** with `html_element()`
- Parse the **HTML table into an R data frame** with `html_table()`

```
women ← mthn %>%
  html_element("div+ .wikitable :nth-child(1)") %>% ## select table eleme
  html_table()                                       ## convert to data fr

women
```

```
## # A tibble: 45 × 7
##    Time      Name                 Nationality    Date    Event/Place    Source No
##    <chr>     <chr>                <chr>          <chr>   <chr>          <chr>    <
##  1 5:40:xx   Marie-Louise Ledru   France         Septe…  Tour de Pari…  ARRS[…  "
##  2 3:40:22   Violet Piercy        United Kingdom Octob…  London [nb 7]  IAAF[…  "
##  3 3:37:07   Merry Lepper         United States  Decem…  Culver City,…  IAAF[…  "
##  4 3:27:45   Dale Greig           United Kingdom May 2…  Ryde           IAAF,…  "
##  5 3:19:33   Mildred Sampson      New Zealand    July …  Auckland, Ne…  IAAF[…  "
##  6 3:14:23   Maureen Wilton       Canada         May 6…  Toronto, Can…  IAAF,…  "
##  7 3:07:27.2 Anni Pede-Erdkamp    West Germany   Septe…  Waldniel, We…  IAAF,…  "
##  8 3:02:53   Caroline Walker      United States  Febru…  Seaside, OR    IAAF,…  "
```

# First Table: Women's Records

Great, it worked! Now convert the date string to a format that R actually understands.[2]

```
women %>%
  mutate(Date = mdy(Date)) %>% ## convert string to date format
  head(4)
```

```
## # A tibble: 4 × 7
##   Time    Name                Nationality    Date       Event/Place    Source No
##   <chr>   <chr>               <chr>          <date>     <chr>          <chr>    <
## 1 5:40:xx Marie-Louise Ledru  France         1918-09-29 Tour de Pari…  ARRS[… "
## 2 3:40:22 Violet Piercy       United Kingd…  1926-10-03 London [nb 7]  IAAF[… "
## 3 3:37:07 Merry Lepper        United States  NA         Culver City,…  IAAF[… "
## 4 3:27:45 Dale Greig          United Kingd…  1964-05-23 Ryde           IAAF,… "
```

[2] *Note:* If column name had spaces or capital letters, we could use the `janitor::clean_names()` convenience function to clean them. (Q: How else could we have done this?)

# First Table: Women's Records

Alright that mostly worked, but there are a few hyperlink references in the dates leading to `NA`s.

- This is a case where using a regular expression is convenient: match the pattern `"[nb X]"` at the **end of the strings**

```r
women ← women %>%
  mutate(across(where(is.character), # mutate across all character string.
                ~str_replace(.x, "\\[nb [0-9]+\\]$", "")), # remove [nb 0
         Date = mdy(Date))  # convert string to date format

women
```

```
## # A tibble: 45 × 7
##    Time       Name           Nationality Date       Event/Place  Source No
##    <chr>      <chr>          <chr>       <date>      <chr>        <chr>  <
## 1 5:40:xx    Marie-Louise Led… France    1918-09-29 "Tour de Par… ARRS[… "
## 2 3:40:22    Violet Piercy   United Kin… 1926-10-03 "London "     IAAF[…
## 3 3:37:07    Merry Lepper    United Sta… 1963-12-16 "Culver City… IAAF[… "
## 4 3:27:45    Dale Greig      United Kin… 1964-05-23 "Ryde"        IAAF,…
```

# Table 2: Men's Records

We could stop here and plot the women's records, but while we're here let's grab the men's records as well.

**Challenge:** take a couple minutes to use SelectorGadget and find the CSS selector for the men's record table.

- Don't peek at next slide until you give it a shot!

```
men ← mthn %>%
  html_element("") %>% ## add the selector!
  html_table()    %>%
  mutate(across(where(is.character), # mutate across all character string.
               ~str_replace(.x, "\\[nb [0-9]+\\]$", "")), # remove [nb 0
         Date = mdy(Date))
```

```r
men ← mthn %>%
  html_element("p+ .wikitable :nth-child(1)") %>% ## add the selector!
  html_table()    %>%
  mutate(across(where(is.character), # mutate across all character string.
                ~str_replace(.x, "\\[nb [0-9]+\\]$", "")), # remove [nb 0
         Date = mdy(Date))
tail(men)
```

```
## # A tibble: 6 × 7
##   Time    Name          Nationality Date       Event/Place      Source    Not
##   <chr>   <chr>         <chr>       <date>     <chr>            <chr>       <
## 1 2:03:38 Patrick Makau   Kenya     2011-09-25 Berlin Marathon  IAAF,[82… "
## 2 2:03:23 Wilson Kipsang Kenya      2013-09-29 Berlin Marathon  IAAF[85]… "
## 3 2:02:57 Dennis Kimetto Kenya      2014-09-28 Berlin Marathon  IAAF[87]… "
## 4 2:01:39 Eliud Kipchoge Kenya      2018-09-16 Berlin Marathon  IAAF[89]  "
## 5 2:01:09 Eliud Kipchoge Kenya      2022-09-25 Berlin Marathon  World At… "
## 6 2:00:35 Kelvin Kiptum  Kenya      2023-10-08 Chicago Marathon World At… "
```

# Browser Inspection Tools

SelectorGadget is a great tool, but sometimes it **takes more work than necessary** and isn't available in all browsers.

**Alternate approach:** use the
**inspect web element**

- Chrome: Right-click > "Inspect" (`Ctrl + Shift + I`)
- Scroll over source elements until the table of interest is highlighted
- Use the selector that pops up over the element
  - i.e. `"table.wikitable"`

# Duplicate Selectors

If we look at the source code a bit more we can see that the selector `table.wikitable` **isn't unique**!

In cases like this, use `html_elements()` to retrieve all matching elements as a list

```
mthn_tabs ← mthn %>%
  html_elements("table.wikitable") ## select all elements matching the se

# first match is men's results
men2 ← mthn_tabs[[1]] %>% html_table()

# second match is women's results
women2 ← mthn_tabs[[2]] %>% html_table()
```

# Marathon Progression

**Challenge:** combine both tables into a single dataframe and plot the record progression

- Clean the `Time` variable (get rid of alphanumeric characters at the end)
- Separate the time variable into hours/minutes/seconds using `separate()`
- Use `hours(), minutes(), seconds()` from `lubridate` to create a time variable for the marathon time
- Create a plot with
  - Marathon time on the y-axis
    - Make sure to include a `scale_y_time()` layer!
  - Date on the x-axis
  - Linetype aesthetic mapping based on the group (men vs. women)
  - Use a nice theme

Answer in a few slides (no peeking until you've tried first)

# Marathon Progression

```r
men ← mutate(men, Group = "Men")
women ← mutate(women, Group = "Women")

mthn_prog ← rbind(men, women) %>%
  mutate(Time = str_replace_all(Time, "[:alpha:]", "")) %>%
   separate(col = "Time", into = c("Hours", "Minutes", "Seconds"), sep = "
  # add time variable using lubridate functions
  mutate(mthn_time = hours(Hours) + minutes(Minutes) + seconds(Seconds))

ggplot(mthn_prog, aes(x = Date, y = mthn_time)) +
  geom_line(aes(linetype = Group)) +
  scale_y_time() +
  labs(y = "Marathon World Record",
       x = "Year") +
  theme_minimal()
```
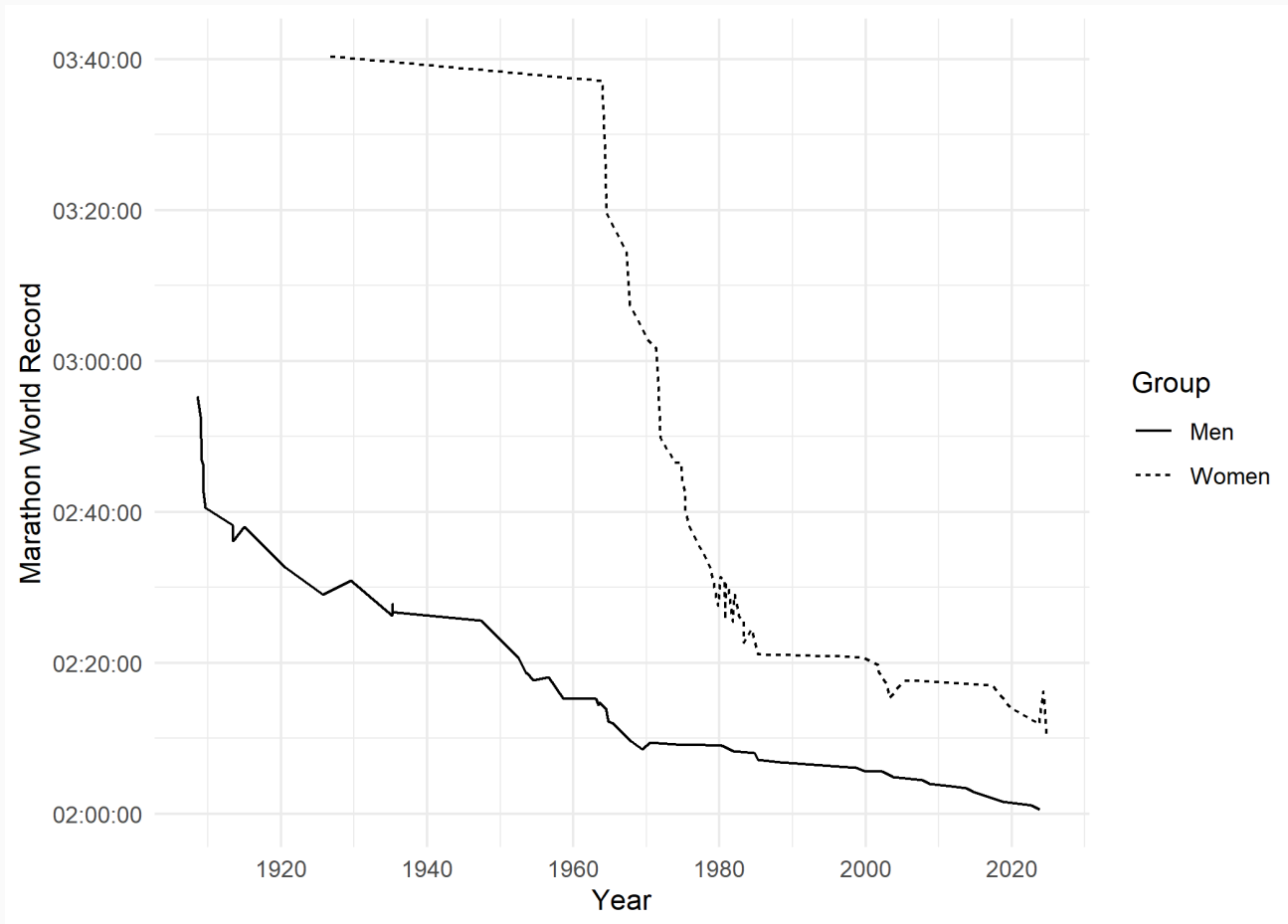
# Table of Contents