# Session 15
# Gathering Electronic Data 1

R for Social Data Science

Jeffrey Ziegler, PhD

Assistant Professor in Political Science & Data Science
Trinity College Dublin

Fall 2022

# Road map for today

Last week:

- Data visualization

This time:

- Online data sources

- Data collection

- Web technologies

- HTML fundamentals

# Online data sources

- Data downloadable in tabular format (E.g. CSV/TSV, XLS, DTA, etc.)

- Data available online as a table (E.g. webpages with rendered tables)

- Unstructured data available online (E.g. simple webpages)

- Interactive webpages with user-input (E.g. webpages with logins, dropdown menus)

- Web APIs (special interfaces for querying, e.g. Twitter, Google)

# Online data collection

- Tabular format: download single or multiple files (automate with 'download.file()' in R, 'wget' in Python/Terminal)

- Online tables and unstructured data: simple web scraping (HTML with XPath, 'rvest' in R, 'beautifulsoup' in Python)

- Interactive webpages: web scraping with headless browser (Selenium, 'RSelenium' in R, 'selenium' in Python)

- Web API: sending requests and processing responses (HTTP queries, 'httr' in R, 'requests' in Python)

# WEB TABLES



## Members of the 1st Dáil

From Wikipedia, the free encyclopedia

Members by constituency   [ edit ]

| Members of the 1st Dáil[4] | | |
|---|---|---|
| **Constituency** | **Name** | **Party** |
| Antrim East | Robert McCalmont | Irish Unionist |
| Antrim Mid | Hugh O'Neill | Irish Unionist |
| Antrim North | Peter Kerr-Smiley | Irish Unionist |
| Antrim South | Charles Curtis Craig | Irish Unionist |
| Armagh Mid | James Rolston Lonsdale | Irish Unionist |
| Armagh North | William Allen | Irish Unionist |
| Armagh South | Patrick Donnelly | Irish Parliamentary |
| Belfast Cromac | William Arthur Lindsay | Irish Unionist |
| Belfast Duncairn | Edward Carson | Irish Unionist |
| Belfast Falls | Joseph Devlin | Irish Parliamentary |
| Belfast Ormeau | Thomas Moles | Irish Unionist |
| Belfast Pottinger | Herbert Dixon | Irish Unionist |

Source: Wikipedia

# UNSTRUCTURED DATA



Source: Eur-Lex

# INTERACTIVE WEBPAGES



Source: Izbori.ba

- Manual scraping (copy-pasting) can be:
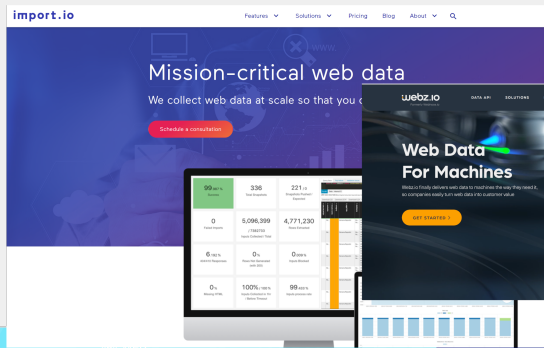  - Extremely laborious and time-consuming
  - Very error-prone
  - Often impossible to reproduce exactly
  - Automated data collection
  - Easy to scale up (computer time is cheap)
  - Less error-prone
  - Usually, perfectly reproducible
  - There is a trade-off (time invested in automation vs time saved)
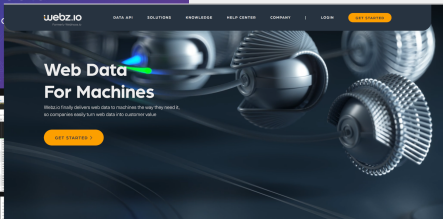  - However, it is good to err on the side of automation

# COMMERCIAL SOLUTIONS

# Web technologies

- Key technologies used to disseminate content on the Web:

  - XML/HTML (E**x**tensible **M**arkup **L**anguage/**H**yper**t**ext **M**arkup **L**anguage)

  - CSS (**C**ascading **S**tyle **S**heets)

  - JavaScript

  - API (**A**pplication **P**rogramming **I**nterface)

  - JSON (**J**ava**S**cript **O**bject **N**otation)

# Static vs dynamic websites

- The critical feature of a website which determines approach to scraping its content

- Static websites all have prebuild source code which is served at user's request

  - ▶ No real-time processing of user's input

  - ▶ Can contain elements that change appearance of a website

  - ▶ Example: Course website

  - ▶ Dynamic websites render websites in real-time as a response to user's input

  - ▶ They can use a range of technologies to achieve it (JavaScript, Python Django, PHP)

  - ▶ Example: Google Maps

# HTML: Hypertext Markup Language

- HTML (**H**yper**t**ext **M**arkup **L**anguage) is a mark-up language for webpages

- Forms the basis of static websites

- Your browser renders (interprets) HTML for viewing

- Current version is HTML5

Extra - W3Schools: Try HTML

# HTML BASICS

- Basic unit of HTML is an *element* (aka *node*)

- Elements, typically, begin with an start tag (e.g. '<h1>')

- And finish with an *end tag* (e.g. '</h1>')

- Content of element is found between start and end tags

- *Attributes* are special words used within a start tag to control element's behaviour (e.g. 'style="color:Red;"')

- Some HTML tag examples:

  - ▶ Document structure: '<html>', '<body>', '<header>'

  - ▶ Document components: '<h1>', '<title>', '<div>'

  - ▶ Text style: '<b>', '<i>'

  - ▶ Hyperlinks: '<a>'

# HTML tree relationships

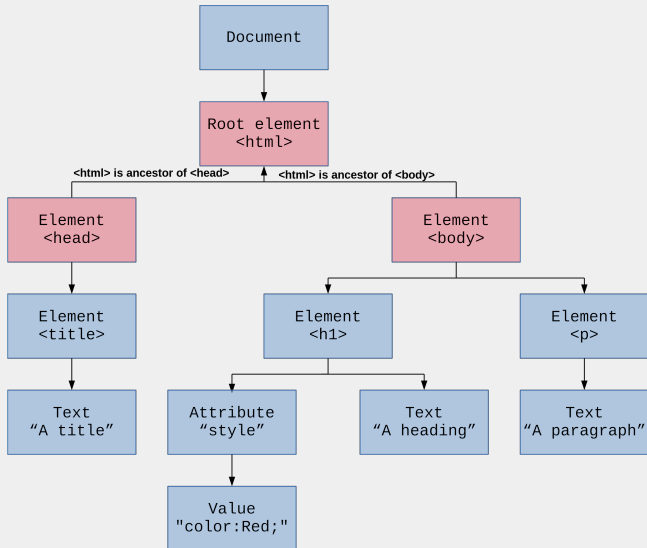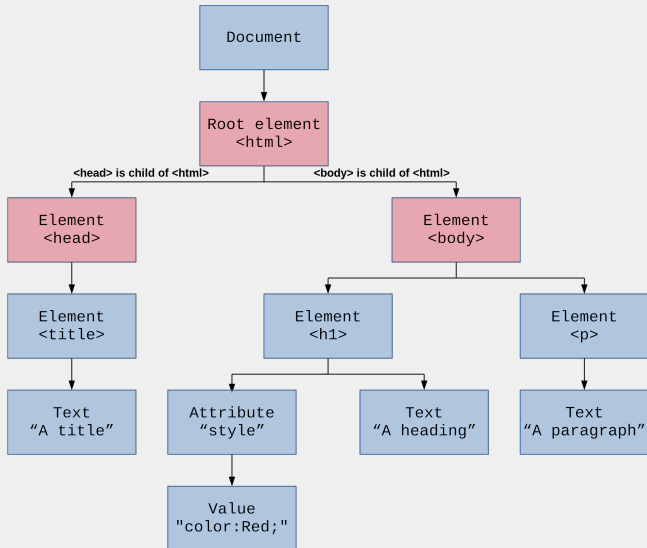- All elements (nodes) in HTML tree are connected by relationships

- These relationship can be of the following types:
  - ▶ Ancestors (parents)
  - ▶ Descendants (children)
  - ▶ Siblings

# HTML PARENT/ANCESTOR

# HTML CHILDREN/DESCENDANTS

# Ex: Parsing HTML tree

```r
1  html_txt <- "\n
2  <!DOCTYPE html> \n
3  <html>\n
4      <head>\n
5          <title >A title </title > \n
6      </head>\n
7      <body>\n
8          <h1 style ='color:Red;'>A heading</h1> \n
9          <p>A paragraph.</p> \n
10     </body>\n
11 </html>"
12 html <- rvest::read_html(html_txt)
13 str(html)
```

```
List of 2
$ node:<externalptr>
$ doc :<externalptr>
- attr(, "class")= chr [1:2] "xml_document" "xml_node"
```

# Ex: Parsing HTML tree

```
1 children <- rvest::html_children(html)
2 children
```

```
{xml_nodeset (2)}
[1] <head>\n<meta http-equiv="Content-Type" content="text/html; charset=UTF-8">\n<title>A title
[2] <body>\n\n        <h1 style="color:Red;">A heading</h1> \n\n        <p>A paragraph.</p> \n\
```

```
1 body <- children[2]
2 rvest::html_name(body)
```

```
[1] "body"
```

# Ex: Parsing HTML tree

```
1  children2 <- rvest::html_children(body)
2  children2

   {xml_nodeset (2)}
   [1] <h1 style="color:Red;">A heading</h1>
   [2] <p>A paragraph.</p>

1  rvest::html_attrs(children2[1])

   [[1]]
   style
   "color:Red;"

1  rvest::html_text(children2[1])

   [1] "A heading"
```

# Tutorial - HTML basics and scraping tables

- We will extract the table of countries with their GDP from a Wikipedia article

- Start by loading in the webpage using 'rvest''s 'read_html()' function

- Go the webpage of the article and locate the elements that would be helpful for table extraction

- Extract the '<table>' node that correponds to the main table

- Extract '<tbody>' element as a child of this element

- Extract the table of with data using 'rvest''s 'html_table()' function

- Tidy up the extracted table

## Overview

This time:

- Online data sources

- Data collection

- Web technologies

- HTML fundamentals

Next time:

- XML, XPath

- APIs