

R CASE STUDY: WEB SCRAPING

John R Little
Duke University

DUKE UNIVERSITY: LAND ACKNOWLEDGEMENT

I want to take a moment to honor the land in Durham, NC. Duke University sits on the ancestral lands of the

DEMONSTRATION GOALS

- Building on earlier [Rfun workshops](#)
- Web scraping is fundamentally a deconstruction process
- Introduce just enough HTML/CSS
- Introduce the `library(rvest)` package for harvesting websites/HTML
- Tidyverse iteration with `purrr::map` - Point out useful documentation & resources

This is a demonstration of leveraging the Tidyverse. This is not a research design or HTML design

CAVEATS

- You will be as successful as the web author(s) were consistent
- Read and follow the *Terms of Use* for any target web host
- Read and honor the host's **robots.txt** | <https://www.robotstxt.org>
- Always **pause** to avoid the perception of a *Denial of Service* (DOS) attack

SCRAPING

Step one: Gather

ingest web page data for analysis

```
rvest::read_html()
```

Step two: Crawling

systematically (iterating) through a website, gathering data from more than one page (URL)

```
purrr::map()
```

Step three: Parsing

Separating the syntactic elements of a web page into meaningful data

```
rvest::html_nodes()  
rvest::html_text()  
rvest::html_attr()
```

HTML

Hyper Text Markup Language

```
1  <html>
2    <body>
3
4      <h1>My First Heading</h1>
5      <p>My first paragraph contains a
6      <a href="https://www.w3schools.com">link</a> to
7      W3schools.com
8    </p>
9
10   </body>
11 </html>
```

HTML + CSS

Cascading Style Sheets

```
1 <html>
2 <body>
3
4   <div class="abc"> ... </div>
5
6   <div id="xyz">
7     <span class="foo"> ... </span>
8   </div>
9
10  <span id="bar"> ... </span>
11
12 </body>
13 </html>
```

for example: <https://www.vondel.humanities.uva.nl/style.css>

PROCEDURE

The basic workflow of web scraping is

1. Development

- Import raw HTML of a single target page (page detail: a leaf)
- Parse the HTML of the test page and gather specific data
- Check *robots.txt* and *Terms Of Use* (TOU)
- In a web browser, manually browse and understand the target
- *Parse* the site navigation and develop an *iteration* plan
- *Iterate*: orchestrate/automate page crawling
- Perform a dry run with a limited subset of the target web site
- Construct pauses: avoid the posture of a DNS attack

2. Production

- Iterate/Crawl the site (navigation: branches)

SITE TREE

