# SCRAPING
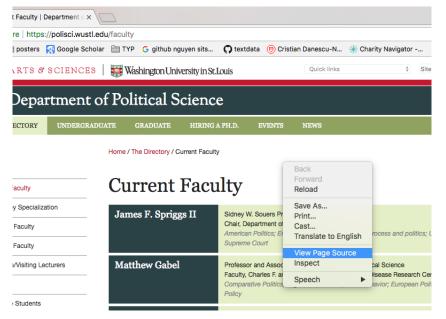
Erin Rossiter

- I use what we go over today (and this week) a lot.
- Scraping, APIs, text handling is where Python beats R.
- The frustrating stuff from last week (exception handling!) is very useful here.

# First things first: Page Source

```
<!DOCTYPE html>
<html>
<head>
<title>Page Title</title>
</head>
<body>

<h1>My First Heading</h1>
<p>My first paragraph.</p>

</body>
</html>
```

posters   Google Scholar   TYP   github nguyen sits...   textdata   Cristian Danescu-N...   Charity Navigator -...

ARTS & SCIENCES | Washington University in St.Louis

Quick links   Site

# Department of Political Science

| ECTORY | UNDERGRADUATE | GRADUATE | HIRING A PH.D. | EVENTS | NEWS |

Home / The Directory / Current Faculty

## Current Faculty

aculty

y Specialization

Faculty

Faculty

s/Visiting Lecturers

Students

| **James F. Spriggs II** | Sidney W. Souers Pr<br>Chair, Department of<br>*American Politics; E*<br>*Supreme Court* | | *rocess and politics; l* |

| **Matthew Gabel** | Professor and Assoc<br>Faculty, Charles F. a<br>*Comparative Politics*<br>*Policy* | cal Science<br>isease Research Cer<br>*avior; European Poli* |

| | Back |
| | Forward |
| | Reload |
| | Save As... |
| | Print... |
| | Cast... |
| | Translate to English |
| | **View Page Source** |
| | Inspect |
| | Speech ▶ |

https://polisci.wustl.edu/faculty

← → C ⟳ ⓘ view-source:https://polisci.wustl.edu/faculty ☆ :

⦂⦂ Apps ✉ email ▣ posters 🎓 Google Scholar ▶ TYP ⚙ github nguyen sits... ⬤ textdata ⦿ Cristian Danescu-N... ♥ Charity Navigator -... ▶ writeup ● Python »

```
1   <!DOCTYPE html>
2   <html xmlns="http://www.w3.org/1999/xhtml" lang="en" dir="ltr"
3     xmlns:content="http://purl.org/rss/1.0/modules/content/"
4     xmlns:dc="http://purl.org/dc/terms/"
5     xmlns:foaf="http://xmlns.com/foaf/0.1/"
6     xmlns:og="http://ogp.me/ns#"
7     xmlns:rdfs="http://www.w3.org/2000/01/rdf-schema#"
8     xmlns:sioc="http://rdfs.org/sioc/ns#"
9     xmlns:sioct="http://rdfs.org/sioc/types#"
10    xmlns:skos="http://www.w3.org/2004/02/skos/core#"
11    xmlns:xsd="http://www.w3.org/2001/XMLSchema#" version="HTML+RDFa 1.1">
12    <head>
13      <meta charset="utf-8" />
14  <link rel="apple-touch-icon" href="/sites/all/themes/awesomesauce/images/iOS-57.png" />
15  <link rel="apple-touch-icon" href="/sites/all/themes/awesomesauce/images/iOS-114.png" sizes="114x114" />
16  <link rel="apple-touch-icon" href="/sites/all/themes/awesomesauce/images/iOS-72.png" sizes="72x72" />
17  <link rel="shortcut icon" href="https://polisci.wustl.edu/sites/polisci.wustl.edu/themes/sitetheme/favicon.ico"
      type="image/vnd.microsoft.icon" />
18  <!--[if lt IE 9]><script type="text/javascript" src="/sites/all/themes/innovate/scripts/html5.js?p7rays"></script><![endif]--><meta
      name="viewport" content="width=device-width, initial-scale=1.0, maximum-scale=1.0" />
19  <meta name="Generator" content="Drupal 7 (http://drupal.org)" />
20  <link rel="canonical" href="/faculty" />
21  <link rel="shortlink" href="/node/4" />
22  <meta http-equiv="X-UA-Compatible" content="IE=Edge,chrome=1" />
23      <title>Current Faculty | Department of Political Science</title>
24  <link type="text/css" rel="stylesheet" href="https://polisci.wustl.edu/files/polisci/css/css_xE-rWrJf-fncB6ztZfd2huxqgxu4WO-
      qwma6Xer30m4.css" media="all" />
25  <link type="text/css" rel="stylesheet" href="https://polisci.wustl.edu/files/polisci/css/css_OOi-hABmd-SMR1h8-obE-
      0sOAEBXpoLcCXmeWHEFYrA.css" media="all" />
26  <link type="text/css" rel="stylesheet" href="https://polisci.wustl.edu/files/polisci/css/css_OAaTCSbhUuU8GVBlD1vHzOHt_tgz3N-
      F0_OqXjxw4.css" media="all" />
27  <link type="text/css" rel="stylesheet"
      href="https://polisci.wustl.edu/files/polisci/css/css_DWPYTaZP2gGao6omeTINDslfzQowIJrrplASKETTOGk.css" media="all" />
28  <link type="text/css" rel="stylesheet" href="https://polisci.wustl.edu/files/polisci/css/css_TcoAtVfl5vJt1FjzFWhK9VXLceW8r3ZfCRavLyghXuw.css" media="screen" />
29  <link type="text/css" rel="stylesheet" href="https://polisci.wustl.edu/files/polisci/css/css_OF-
      ro9xGfLiKzLxUFuAhYPuYCsdII2_jbt2RGCO6WY.css" media="handheld, only screen and (max-width: 767px)" />
30  <link type="text/css" rel="stylesheet" href="https://polisci.wustl.edu/files/polisci/css/css_NbnqFpRdg8MhD-MbJZmXr3i7Cl2-
      RxRSbC4Im43WfDY.css" media="handheld, only screen and (max-device-width: 480px)" />
31  <link type="text/css" rel="stylesheet" href="https://polisci.wustl.edu/files/polisci/css/css_Tuq8NnOU2u2oOjPX8LEIDKpgjiCBdgR-
      COOYw_xpfMY.css" media="print" />
32
33  <!--[if IE 8]>
34  <link type="text/css" rel="stylesheet"
      href="https://polisci.wustl.edu/files/polisci/css/css_KHVZdciTQ9NLR5qC9yFi34NDUDVt77RAYvQyNmD6VF4.css" media="all" />
35  <![endif]-->
36
```

## Scraping HTML

`urllib2`

- "crawler"; in the background
- navigate to url

`BeautifulSoup`

- download HTML
- parse it for info

## Great approach when...

- Info is contained in HTML (older websites are great!)
- Webpage names have a pattern
- Example:
  `http://www.presidency.ucsb.edu/press_briefings.php?year=2018`

## Not helpful when...

- Info comes from javascript call, or info is not in HTML for some other reason (fancier websites)
- Example:
  https://www.oyez.org/cases/2017/17-586
- Your task is not structured (e.g., using a search engine)
- Example:
  http://www.spencerdailyreporter.com/

- Beautiful Soup and Urllib are "crawlers" in the "background".

## A second approach: Remote driver

- Beautiful Soup and Urllib are "crawlers" in the "background".
  - They are incredibly fast...

## A second approach: Remote driver

- Beautiful Soup and Urllib are "crawlers" in the "background".
  - They are incredibly fast...
  - ... but also easier to detect and block

## A second approach: Remote driver

- Beautiful Soup and Urllib are "crawlers" in the "background".
  - They are incredibly fast...
  - ... but also easier to detect and block
  - Plus, sometimes they aren't amenable to your task

## A second approach: Remote driver

- Beautiful Soup and Urllib are "crawlers" in the "background".
  - They are incredibly fast...
  - ... but also easier to detect and block
  - Plus, sometimes they aren't amenable to your task
- Selenium is a "remote driver" of your favorite browser

## A second approach: Remote driver

- Beautiful Soup and Urllib are "crawlers" in the "background".
  - They are incredibly fast...
  - ... but also easier to detect and block
  - Plus, sometimes they aren't amenable to your task
- Selenium is a "remote driver" of your favorite browser
  - Simulates behavior of human (e.g., point, click, type)

## A second approach: Remote driver

- Beautiful Soup and Urllib are "crawlers" in the "background".
    - They are incredibly fast...
    - ... but also easier to detect and block
    - Plus, sometimes they aren't amenable to your task
- Selenium is a "remote driver" of your favorite browser
    - Simulates behavior of human (e.g., point, click, type)
    - Less likely to be tracked and blocked

## A second approach: Remote driver

- Beautiful Soup and Urllib are "crawlers" in the "background".
  - They are incredibly fast...
  - ... but also easier to detect and block
  - Plus, sometimes they aren't amenable to your task

- Selenium is a "remote driver" of your favorite browser
  - Simulates behavior of human (e.g., point, click, type)
  - Less likely to be tracked and blocked
  - Flexible to undefined tasks

## A second approach: Remote driver

- Beautiful Soup and Urllib are "crawlers" in the "background".
  - They are incredibly fast...
  - ... but also easier to detect and block
  - Plus, sometimes they aren't amenable to your task

- Selenium is a "remote driver" of your favorite browser
  - Simulates behavior of human (e.g., point, click, type)
  - Less likely to be tracked and blocked
  - Flexible to undefined tasks
  - Cons:

## A second approach: Remote driver

- Beautiful Soup and Urllib are "crawlers" in the "background".
    - They are incredibly fast...
    - ... but also easier to detect and block
    - Plus, sometimes they aren't amenable to your task

- Selenium is a "remote driver" of your favorite browser
    - Simulates behavior of human (e.g., point, click, type)
    - Less likely to be tracked and blocked
    - Flexible to undefined tasks
    - Cons:
        - slow

## A second approach: Remote driver

- Beautiful Soup and Urllib are "crawlers" in the "background".
  - They are incredibly fast...
  - ... but also easier to detect and block
  - Plus, sometimes they aren't amenable to your task

- Selenium is a "remote driver" of your favorite browser
  - Simulates behavior of human (e.g., point, click, type)
  - Less likely to be tracked and blocked
  - Flexible to undefined tasks
  - Cons:
    - slow
    - dependent on Internet connection

## A second approach: Remote driver

- Beautiful Soup and Urllib are "crawlers" in the "background".
    - They are incredibly fast...
    - ... but also easier to detect and block
    - Plus, sometimes they aren't amenable to your task

- Selenium is a "remote driver" of your favorite browser
    - Simulates behavior of human (e.g., point, click, type)
    - Less likely to be tracked and blocked
    - Flexible to undefined tasks
    - Cons:
        - slow
        - dependent on Internet connection
        - start-up costs (time, patience) in getting it set up on your machine

- Always check terms of service first.

# Some tricks and advice to become a pro-scraper

- Always check terms of service first.
- Inspect source code and structure of website... think about what approach will work best for your task

# Some tricks and advice to become a pro-scraper

- Always check terms of service first.
- Inspect source code and structure of website... think about what approach will work best for your task
- I prefer Google Chrome to view page source and insect elements

# Some tricks and advice to become a pro-scraper

- Always check terms of service first.
- Inspect source code and structure of website... think about what approach will work best for your task
- I prefer Google Chrome to view page source and insect elements
- Use time breaks to avoid being blocked

# Some tricks and advice to become a pro-scraper

- Always check terms of service first.
- Inspect source code and structure of website... think about what approach will work best for your task
- I prefer Google Chrome to view page source and insect elements
- Use time breaks to avoid being blocked
- Expect your code to break–websites change :(

1. How to read and write files
2. How to navigate to url & parse HTML – approach #1
3. Example of selenium – approach #2
4. Lab to practice approach #1