# Extracting Data From the Web Using Python

## Web Data Extraction Workshop
## D-Lab
## UC Berkeley

Berkeley
UNIVERSITY OF CALIFORNIA

# Learning Objectives

- The purpose of web scraping and APIs when to use either
- HTML & CSS and their use in web scraping.
- Use the requests and BeautifulSoup libraries to acquire and parse data from websites.
- Use a 3rd-party Python library to make API calls.
- Working with data pulled from APIs.

# Why Webscrape

- Tons of web data useful for social scientists and humanists
  - social media
  - news media
  - government publications
  - organizational records

- Two kinds of ways to get data off the web
  - Webscraping – i.e. user-facing websites for humans
  - APIs – i.e. application-facing, for computers

# Webscraping v. APIs

- **Webscraping Benefits**
  - Any content that can be viewed on a webpage can be scraped. Period
  - No API needed
  - No rate-limiting or authentication (usually). Your IP can be blocked (403)
- **Webscraping Challenges**
  - Rarely tailored for researchers
  - Messy, unstructured, inconsistent
  - Entirely site-dependent
- **Rule of thumb:**
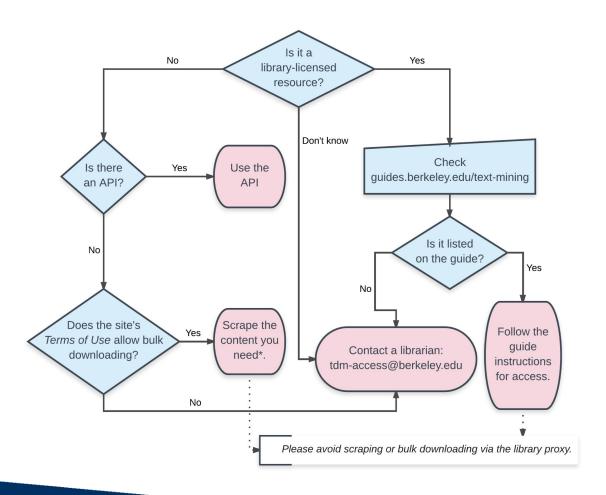  - Check for API first. If not available, scrape.

Berkeley
UNIVERSITY OF CALIFORNIA

# Some Disclaimers!!

- Check a site's terms and conditions before scraping.

- Be nice – don't hammer the site's server.
    - Add a time delay in between batches of scrapping.

- Sites change their layout all the time. Your scraper will break.
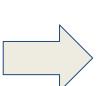
# Workflow

# What's a website

- Some combination of codebase, database.
- The "front end" product is HTML + CSS stylesheets + javascript.
- Browser turns the left image into the right.

# Webscraping returns HTML

- It's easy to pull HTML from a website

- It's much more difficult to find the information you want from that HTML

- So we have to learn how to **parse** HTML to find the data we want



## Berkeley
UNIVERSITY OF CALIFORNIA

# Basic strategy of webscraping:

1.  Find out what kind of HTML element your data is in. (Use your browser's "inspector")

2.  Think about how you can differentiate those elements from other, similar elements in the webpage using HTML/CSS anatomy.

3.  Use Python and add-on modules like BeautifulSoup to extract just that data.

# HTML: Basic structure

```html
<!DOCTYPE html>
<html>
    <head>
        <title>Page title</title>
    </head>
    <body>
        <p>Hello world!</p>
    </body>
</html>
```

Berkeley
UNIVERSITY OF CALIFORNIA

# HTML as a Tree



Each branch of the tree is called an **element**

# HTML Elements

Generally speaking, an HTML element has three components:

1. Tags (starting and ending the element)
2. Attributes (giving information about the element)
3. Text, or Content (the text inside the element)



UNIVERSITY OF CALIFORNIA

# HTML: Tags



THIS SAYS "BEGIN ITALICS NOW."
THIS IS THE ACTUAL TEXT
THIS SAYS "END ITALICS NOW."

`<i>text</i>`

THIS IS WHAT SHOWS UP ON YOUR SCREEN → *text*

| Tag | Meaning |
|---|---|
| <head> | page header (metadata, etc |
| <body> | holds all of the content |
| <p> | regular text (paragraph) |
| <h1>,<h2>,<h3> | header text, levels 1, 2, 3 |
| ol,,<ul>,<li> | ordered list, unordered list, list item |
| <a href="page.html"> | link to "page.html" |
| <table>,<tr>,<td> | table, table row, table item |
| <div>,<span> | general containers |

Berkeley
UNIVERSITY OF CALIFORNIA

# HTML Attributes

- HTML elements can have attributes.
- Attributes provide additional information about an element.
- Attributes are always specified in the start tag.
- Attributes come in name/value pairs like: name="value"



**Start Tag**

`<Item optional="1">`

Tag Name — Attribute Name — Attribute Value

Attribute

**End Tag**

`</Item>`

# Finding HTML

- Sometimes we can find the data we want just by using HTML tags or attributes (e.g, all the <a> tags)
- More often, this isn't enough: There might be 1000 <a> tags on a page. But maybe we want only the <a> tags inside of a <p> tag.
- This is where CSS comes in.

# CSS (Cascading Style Sheet)

- CSS defines how HTML elements are to be displayed
- HTML came first. But it was only meant to define content, not format it. While HTML contains tags like <font> and <color>, this is a very inefficient way to develop a website.
- To solve this problem, CSS was created specifically to display content on a webpage. Now, one can change the look of an entire website just by changing one file.
- Most web designers litter the HTML markup with tons of classes and ids to provide "hooks" for their CSS.
- You can piggyback on these "hooks" to jump to the parts of the markup that contain the data you need.

Berkeley
UNIVERSITY OF CALIFORNIA

# CSS Selectors & Declarations

| Type | HTML | CSS Selector |
|------|------|-------------|
| Element | <a>, | a<br>p a |
| Class | <a class="blue"> | .blue<br>a.blue |
| ID | <a id="blue"> | #blue<br>a#blue |

- Selector: *a*
- Property: *background-color*
- Value: *yellow*

Berkeley
UNIVERSITY OF CALIFORNIA

# CSS Hooks

## Basic Anatomy of a CSS Rule

Element   Property   Value

Rule

h2 {
    color: #333;  _ _ _ _ _ _ _ _
    font-size: 150%;
    margin: 5px 0 10px 15px;
}

Declaration

Declaration Syntax:
Property then Colon
Value then Semicolon

Declaring a CSS Rule for a **Class** Attribute

the XHTML
<a class="pdf" href="brochure.pdf">Brochure</a>

the CSS
.pdf {background: url(images/pdf.gif) no-repeat left 50%;}

use a period when writing a rule for a class

Declaring a CSS Rule for an **Id** Attribute

the XHTML
<div id="wrapper">Main Content</div>

the CSS
#wrapper {width: 750px; margin: 0 auto;}

use a pound sign when writing a rule for a id

**Berkeley**
UNIVERSITY OF CALIFORNIA

# CSS meets HTML

What does the following HTML render to?
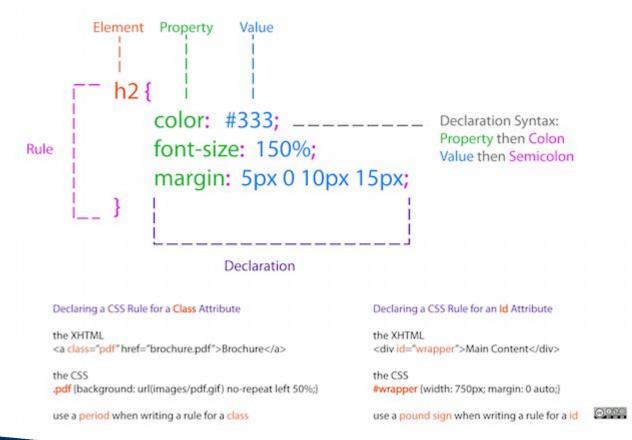
```html
<body>
    <table id="content">
        <tr class='kurtis'>
            <td class='firstname'>
                Kurtis
            </td>
            <td class='lastname'>
                McCoy
            </td>
        </tr>
        <tr class='leah'>
            <td class='firstname'>
                Leah
            </td>
            <td class='lastname'>
                Guerrero
             </td>
        </tr>
    </table>
```

Berkeley
UNIVERSITY OF CALIFORNIA

# Relevance to Webscrapping

- It's necessary to cover HTML & CSS because we'll be using their tags to parse scrapped content.
- In the next section we'll be using the BeautifulSoup library to scrape web data.
- As previously mentioned, that's the easy part.
- Now that we know how tags work and what they're used for, we'll have an easier time sifting through the pile of HTML data