# WEB-SCRAPING WITH PYTHON

# Overview

- Web Scraping 101

- Basic HTML

- String Detection with Regex

- Python: *beautifulsoup*

  - *Example 1) Jeopardy Archive*

  - *Example 2) HackerNews Blog (If time allows)*

- Q&A / DataFest Ideas (If time allows)

# WEB-SCRAPING 101

# Web Scraping 101

*What is Web Scraping ?*

- Process of extracting data from websites by fetching them and extracting it

- Many different software options to use

- Fundamental process to creating datasets from online resources

*How does it work ?*

- Take unstructured data (HTML tags) → structured format (data frame)

- Unstructured data is estimated to make up 65%-85% of the World Wide Web (!)

# Web Scraping 101

*Examples of Unstructured Data:*

- **Human Generated:** Social Media, Websites, Photo Sharing, Court Reports, etc.

- **Machine Generated:** Weather Data, Surveillance Photos, Sensors, etc.

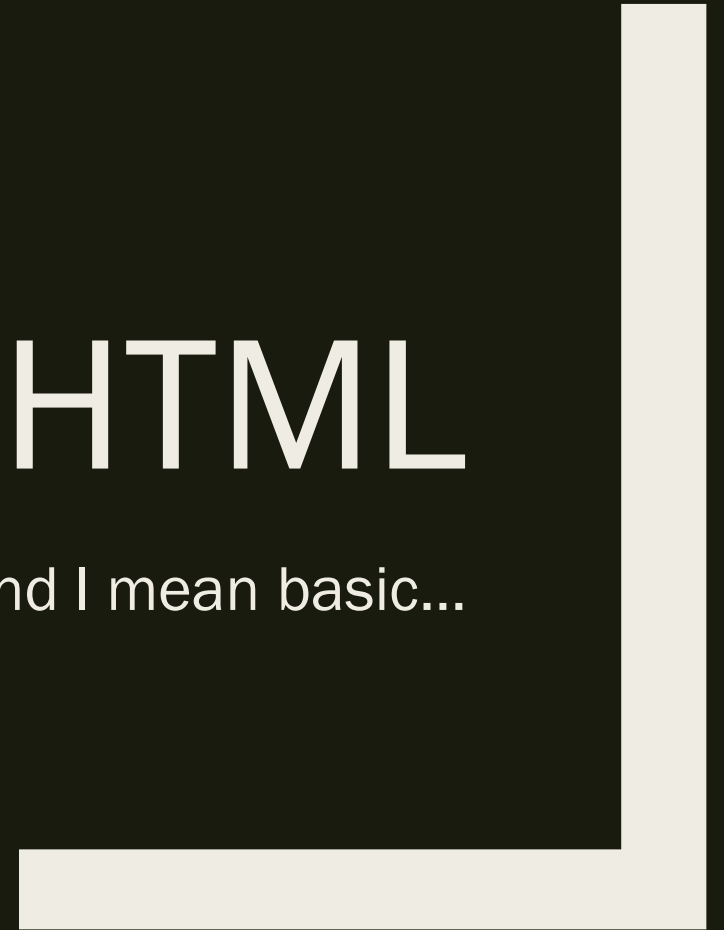| | Structured Data | Unstructured Data |
|---|---|---|
| **Characteristics** | • Pre-defined data models<br>• Usually text only<br>• Easy to search | • No pre-defined data model<br>• May be text, images, sound, video or other formats<br>• Difficult to search |
| **Resides in** | • Relational databases<br>• Data warehouses | • Applications<br>• NoSQL databases<br>• Data warehouses<br>• Data lakes |
| **Generated by** | Humans or machines | Humans or machines |
| **Typical applications** | • Airline reservation systems<br>• Inventory control<br>• CRM systems<br>• ERP systems | • Word processing<br>• Presentation software<br>• Email clients<br>• Tools for viewing or editing media |
| **Examples** | • Dates<br>• Phone numbers<br>• Social security numbers<br>• Credit card numbers<br>• Customer names<br>• Addresses<br>• Product names and numbers<br>• Transaction information | • Text files<br>• Reports<br>• Email messages<br>• Audio files<br>• Video files<br>• Images<br>• Surveillance imagery |

# Web Scraping 101

*Web Scraping Tools Available (For Free!)*

■ Python (Beautifulsoup Library)

– *https://pypi.org/project/beautifulsoup4/*

■ Selector Gadget

– *https://selectorgadget.com*

■ R (Rvest Package)

– *https://www.rdocumentation.org/packages/rvest/versions/0.3.5*

# BASIC HTML

And I mean basic...

# Basic HTML

- Websites are composed of a series of descriptive HTML tags
- **<tagname> \*content goes here \* </tagname>**
  - **<head> tags** contain meta information describing the site
  - **<title> tags** specify the document title
  - **<body> tags** contain page content that's visible
  - **<h1> tags** define large headings
  - **<p> tags** define paragraphs

```html
<html>
    <head>
        <title>Page title</title>
    </head>
    <body>
        <h1>This is a heading</h1>

        <p>This is a paragraph.</p>

        <p>This is another paragraph.</p>
    </body>
</html>
```

# Basic HTML Continued

- **<a> tags** define hyperlinks to URLs in the document
  - Important for iterating through pages
- **<td> tags** define a standard cell in an HTML table
  - Often where you will find data points
- **<img> tags** define images
- **<table> tags** define tables
- **<strong> tags** give text a strong emphasis

# But its never that simple...

# The Workaround: Selector Gadget

Allows you to directly extract the html elements of the contents on a web page you desire!

# Additional HTML tips …

■ Two of very common HTML attributes are **class** and **id.** They are used for grouping and identifying HTML tags.

# EXAMPLES & TIPS

Using Beautiful Soup to scrape the *Jeopardy* Archive and Hacker News

# Tip 1: Inspect your data source:

- Explore the website
  - *Which data points do you want to gather?*
  - *How should they be organized ?*
- Decipher the information of URLS!
  - *Base URL:*
    - Represent path to the website itself
  - *Query Parameters:*
    - Represent the additional values that can be declared on a page
      - *Starter parameters (start with ?)*
      - *Information (key = value)*
    - Parameters are often separated by an '&'

# Example URLs

What is the Base URL?

What are the Query Parameters?

- *http://www.j-archive.com/showgame.php?game_id=6527*

- https://news.ycombinator.com/submitted?id=whoishiring

- *https://laist.com/search_query.php?q=homelessness*

# Example URLs

What is the ==Base URL==?

What are the Query Parameters?

- ==http://www.j-archive.com/==showgame.php?game_id=6527

- ==https://news.ycombinator.com==/submitted?id=whoishiring

- ==https://laist.com/==search_query.php?q=homelessness

# Example URLs

What is the Base URL?

What are the Query Parameters?

- http://www.j-archive.com/showgame.php?game_id=6527

- https://news.ycombinator.com/submitted?id=whoishiring

- https://laist.com/search_query.php?q=homelessness

# Tip 2:
# Fetch the Page & Review Contents

```
from bs4 import BeautifulSoup
```

- ■ Using Beautiful Soup Package
  - – *Creates a 'soup' of html tags*

```
<!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.1//EN" "http://www.w3.org/TR/xhtml11/DTD/xhtml11.dtd">
<html xml:lang="en" xmlns="http://www.w3.org/1999/xhtml">
 <head>
  <title>
   J! Archive - Show #8215, aired 2020-05-01
  </title>
  <link href="http://www.j-archive.com/styles.css" rel="styleSheet" type="text/css"/>
  <link href="http://www.j-archive.com/favicon.ico" rel="shortcut icon"/>
  <meta content="Jeopardy Archive Quemars Ahmed Ashleigh McCord Jesse Laymon Final clue response question
uble" name="keywords"/>
  <meta content="An archive of clues and players for Jeopardy! show #8215." name="description"/>
 </head>
 <body>
  <div id="navbar">
   <span id="navbarlogo">
    <a href="http://www.j-archive.com">
     <img alt="J! Archive" height="22" src="http://www.j-archive.com/j-a.gif" width="100"/>
    </a>
   </span>
   <span id="navbartext">
```

- ■ Where is the content you want?

- ■ How is it stored on the site?

# Tip 3: Error Catching

- This will happen!!

- Important for ensuring the quality of the data captured.

```
>>> while True:
...     try:
...         x = int(input("Please enter a number: "))
...         break
...     except ValueError:
...         print("Oops!  That was no valid number.  Try again...")
...
```

- Try clause will be executed as normal

- If and only if try clause breaks is the exception raised. This is where you handle error management.

# Error handling example

```python
try:
    page_response.raise_for_status()
    pass
except:
    print("HTML ERROR CODE: " + str(page_response.status_code))
```

# ADDITIONAL TIPS / IDEAS

# R vs. Python for web scrapping?

- Use whatever you are comfortable with!

- Use the rvest package in R
  - *Similar to Beautiful Soup*

- At the end of the day, the Selector Gadget is the tool that will be the most help to you.

# A few websites that might be worth scraping!

- **■ *LAist***
  - *https://laist.com/search_query.php?q=homelessness*
  - *How has reporting on homelessness changed since quarantine?*

- **■ *Law360***
  - *https://www.law360.com/articles/1252836/coronavirus-the-latest-court-closures-and-restrictions*
  - *How fast did courts respond to quarantine measures?*

- **■ *Craigslist***
  - *https://losangeles.craigslist.org/*
  - *What are people selling / buying during a pandemic ?*

# Some final thoughts

■ When in doubt: <u>Stack Overflow</u>!

■ Additional Materials on HTML: https://www.w3schools.com/html/default.asp

■ Additional Materials on Regex: https://docs.python.org/3/library/re.html

■ Start simple and then expand your data collection!

   – *Don't try to capture too much data all at once. Get a single page working then add more dynamic features.*

   – *Dynamic sites? Check out the Selenium library in Python for additional tools.*

# THANK YOU!

Add me on LinkedIn!