# A Hands-on Session for Web Scraping

Md. Izhar Ashraf ashraf@imsc.res.in

March 17, 2025

## Contents

1	Introduction to Web Scraping			
	1.1	Definition and Purpose of Web Scraping	2	
	1.2	Importance of Web Scraping in Epigraphy	2	
	1.3	Legal and Ethical Considerations	2	
<b>2</b>	Understanding Web Basics			
	2.1	Building Blocks of a Webpage	3	
	2.2	Understanding the Structure of a Webpage	4	
	2.3	Introduction to the Document Object Model (DOM)	6	
3	Tools and Libraries for Web Scraping			
	3.1	BeautifulSoup	7	
		3.1.1 Requests	7	
		3.1.2 XPath (XML Path Language) Syntax	9	
	3.2	Scrapy	11	
	3.3	Selenium	12	
4	Set	ting & installation of required packages	13	

## 1 Introduction to Web Scraping

### 1.1 Definition and Purpose of Web Scraping

- Web scraping is the automated collection of data from websites, using software to extract information for analysis or use in various applications.
- Imagine you're an archaeologist, sifting through sand to uncover hidden treasures. Web scraping is similar! Websites hold a vast amount of information, but it's often coded and not readily accessible.
- Think of websites like libraries with countless books. Scraping allows you to efficiently target and collect specific information from those books, without reading everything cover to cover.

## 1.2 Importance of Web Scraping in Epigraphy

- Epigraphy, the study of inscriptions or epigraphs as writing, often involves the collection of data from various online databases, digital libraries, and archives.
- By employing web scraping techniques, researchers can automate the extraction of valuable epigraphic data, including texts, translations, metadata, and images.
- This not only accelerates the research process but also enables the creation of comprehensive datasets that can be used for further study.

## 1.3 Legal and Ethical Considerations

- While web scraping is a powerful tool for researchers, it's crucial to navigate the legal and ethical considerations involved.
- There are special files called "robots.txt" that tell web scrapers which parts of a website are okay to access. It's like a map showing where you can "dig" for information.
- We also need to consider legal aspects like copyright. The information itself might be freely available, but how it's presented might be protected.
- Finally, websites can get overloaded if too many requests come in at once. It's important to be respectful and avoid overwhelming them.
- By following these guidelines, web scraping becomes a valuable tool for research and knowledge sharing in epigraphy and beyond!

## 2 Understanding Web Basics

### 2.1 Building Blocks of a Webpage

- Websites are built with three key ingredients: HTML, CSS, and JavaScript [1].
- Think of them like the bricks, mortar, and electrical wiring of a house.
- HTML (HyperText Markup Language): The foundation, like bricks. It defines the structure and content of a webpage, like headings, paragraphs, and images.

```
<!DOCTYPE html>
1
  <html>
2
  <head>
3
           <title>Page Title</title>
4
  </head>
5
  <body>
6
           <h1>This is a Heading</h1>
7
           This is a paragraph.
8
  </body>
9
  </html>
10
```

• CSS (Cascading Style Sheets): The decorator, like mortar. It defines the visual style of the webpage, like fonts, colors, and layout.

```
<! DOCTYPE html>
1
 <html>
 <head>
3
    <title>Page Title</title>
4
 </head>
5
 <body>
6
    <h1 style="color:blue;_text-align:_center;">This is a
      \hookrightarrow Heading</h1>
    This is a
8
      \hookrightarrow paragraph.
 </body>
9
 </html>
```

• JavaScript (JS): The electrician, like wiring. It adds interactivity to webpages, like animations or forms that respond to your clicks.

```
<!DOCTYPE html>
1
  <html>
2
  <head>
3
    <title>Change Text Color</title>
4
    <script>
5
      function changeColor() {
6
        // Get a reference to the heading element
        var heading = document.getElementById("myHeading");
8
9
        // Change the heading's text color
```

```
heading.style.color = "green";
11
      }
    </script>
13
  </head>
14
  <body>
    <h1 id="myHeading" style="text-align:_center;">This is a
          Heading</h1>
      \hookrightarrow
    This is a
17

→ paragraph.

    <button onclick="changeColor()">Magic Button</button>
18
  </body>
19
  </html>
20
```

### 2.2 Understanding the Structure of a Webpage

- Imagine a webpage as a well-organized document. HTML tags create a hierarchy, like headings and subheadings, to structure the content.
- Each element has tags (like labels) that define its purpose (e.g., heading, paragraph, image).
- By understanding this structure, we can pinpoint the specific data we want to scrape.

```
<!DOCTYPE html>
  <html>
  <head>
3
     <title>Change Text Color</title>
4
     <style>
5
       .center-content {
6
         text-align: center;
7
       }
8
       table {
9
         margin: auto; /* Centers the table */
10
       }
11
     </style>
12
     <script>
13
       function changeColor() {
14
         // Get a reference to the heading element
         var heading = document.getElementById("myHeading");
16
17
         // Change the heading's text color
18
         heading.style.color = "green";
19
       }
20
     </script>
21
  </head>
22
  <body>
23
     <h1 id="myHeading" style="text-align:_center;">This is a
24
        ↔ Heading</h1>
```

```
This is a
25

→ paragraph.

    <br>>
26
27
    <div class="center-content">
28
      <button onclick="changeColor()">Magic Button</button>
29
      \langle br \rangle
30
      \langle br \rangle
31
      <!-- Adding a simple table -->
32
      33
       <caption>Demo Table 1</caption>
34
       \langle tr \rangle
35
         Column 1
36
         Column 2
37
         Column 3
38
       39
       \langle tr \rangle
40
         Row 1, Cell 1
41
         Row 1, Cell 2
42
         Row 1, Cell 3
43
       44
       \langle tr \rangle
45
         Row 2, Cell 1
46
         Row 2, Cell 2
47
         Row 2, Cell 3
48
        49
      50
      <br>>
53
      <!-- Adding another simple table -->
54
      <caption>Demo Table 2</caption>
56
       \langle tr \rangle
57
         ID
58
         >Site Name
       60
       \langle tr \rangle
61
         101
62
         Harappa
       64
       \langle tr \rangle
65
         102
66
         Mohenjo-daro
67
       68
       \langle tr \rangle
69
         103
70
         Rakhigarhi
71
```

### 2.3 Introduction to the Document Object Model (DOM)

- The Document Object Model (DOM) is a way to represent a webpage's structure as a tree.
- Each element on the page becomes a node in this tree, with parent-child relationships reflecting the HTML hierarchy.
- Understanding the DOM is crucial for web scraping tools like BeautifulSoup to navigate and extract specific data from a webpage.

```
document
   |- html
2
      - head
3
         |- title: "Change_Text_Color"
      4
         |- style
      L
5
           `- (.center-content and table styles)
         6
         `- script: function changeColor()
      I
      `- body
8
         |- h1#myHeading (style="text-align:")
9
         |- p (style="color:red;_text-align:_center")
10
         |- br
11
          `- div.center-content
12
             |- button (onclick="changeColor()")
13
             |- br
14
             |- br
             |- table#table1 (border="1")
16
                |- caption: "DemouTableu1"
             L
17
                |- tr
18
                    |- th: "Column⊔1"
19
                    |- th: "Column_{\sqcup}2"
20
                    `- th: "Column⊔3"
21
                |- tr
22
                    |- td: "Rowu1,uCellu1"
23
                    |- td: "Rowul,uCellu2"
24
                    `- td: "Row⊔1,⊔Cell⊔3"
25
                `- tr
26
                    |- td: "Row⊔2,⊔Cell⊔1"
27
                    |- td: "Rowu2,uCellu2"
28
                    `- td: "Rowu2,uCellu3"
29
             |- br
30
               table#table2 (border="1")
```

```
|- caption: "Demo<sub>u</sub>Table<sub>u</sub>2"
32
                    |- tr
33
                       |- th: "ID"
34
                       `- th: "SiteuName"
35
                    - tr
36
                       |- td: "101"
37
                       `- td: "Harappa"
38
                    - tr
39
                       |- td: "102"
40
                        `- td: "Mohenjo-daro"
41
                    `- tr
42
                       |- td: "103"
43
                        `- td: "Rakhigarhi"
44
```

## 3 Tools and Libraries for Web Scraping

Python is ideal for web scraping with its simple syntax and powerful libraries like BeautifulSoup, Requests, and Scrapy make data extraction efficient and accessible for all skill levels.

### 3.1 BeautifulSoup

Beautiful Soup [2], started by Leonard Richardson, is a Python package used for parsing HTML and XML documents, including those with malformed markup. It creates a parse tree that facilitates data extraction from HTML, making it highly useful for web scraping. However, it lacks the capability to download web pages on its own, a function complemented by the Requests library.

#### 3.1.1 Requests

Requests is a straightforward Python library for HTTP requests, essential for retrieving web content. Paired with BeautifulSoup, it efficiently manages web scraping tasks. Beautiful Soup is a Python library for parsing HTML and XML documents. It provides idiomatic ways of navigating, searching, and modifying the parse tree. Here's a concise cheat sheet of its basic syntax and functions:

#### Creating a Soup Object

To parse a document, you first need to create a BeautifulSoup object, which represents the document as a nested data structure:

```
1 from bs4 import BeautifulSoup
2 # Parse HTML from a string
3 soup = BeautifulSoup(html_doc, 'html.parser')
4 # Parse HTML from a file
5 with open("index.html") as file:
6 soup = BeautifulSoup(file, 'html.parser')
7 # Parses XML content
8 soup = BeautifulSoup(xml content, 'lxml')
```

#### Navigating the Tree

BeautifulSoup offers multiple ways to navigate and search the parse tree:

```
1 # Accessing tag names
2 soup.title
3 # Accessing tag attributes
4 soup.title.name
5 soup.title.string
6 soup.title['attribute']
7 # Navigating using tag names
8 soup.body.a
9 # Navigating using methods
10 soup.find_all('a')
11 soup.find(id='link3')
```

#### Searching the Tree

BeautifulSoup's search methods allow you to find elements based on their attributes, text content, or filter functions:

```
# Find all tags with a specific name
2 soup.find all('a')
3 # Find the first tag with a specific name
4 soup.find('title')
5 # Find tags using keyword arguments
6 soup.find_all(id='link2')
7 soup.find all(href=True)
8 # Searching by CSS class
 soup.find_all("a", class_="sister")
9
10 # Using string arguments
soup.find all(string="Elsie")
12 # Using regular expressions
13 import re
soup.find all(string=re.compile("Dormouse"))
15 # Using a list
 soup.find_all(["a", "b"])
16
```

#### Modifying the Tree

You can also modify the HTML or XML tree in various ways:

```
1 # Changing tag names and attributes
2 tag.name = 'blockquote'
3 tag['attribute'] = 'newuvalue'
4 # Adding and removing tags
5 new_tag = soup.new_tag('a', href="http://www.example.com")
6 soup.body.append(new_tag)
7 tag.extract() # removes a tag from the tree
```

This cheat sheet covers the basics to get started with BeautifulSoup for web scraping tasks. For more detailed information, refer to the official BeautifulSoup documentation.

#### Example code

```
import requests
1
  from bs4 import BeautifulSoup
2
3
  # The URL of the website you want to scrape
4
  url = 'http://example.com/news'
5
  # Use Requests to get the webpage content
6
  response = requests.get(url)
7
  # Create a BeautifulSoup object and specify the parser
8
  soup = BeautifulSoup(response.text, 'html.parser')
9
  # Find all 'h2' elements with a class 'article-title'
10
  article_titles = soup.find_all('h2', class_='article-title')
11
  # Loop through the list of titles and print them
12
 for title in article titles:
      print(title.text.strip())
14
```

#### 3.1.2 XPath (XML Path Language) Syntax

XPath (XML Path Language) is a query language designed for navigating and selecting nodes from an XML document. It allows for precise location of elements, attributes, text, and more within XML files using a path-like syntax.

XPath expressions can be used to navigate through elements and attributes in an XML document, making it a powerful tool for XML querying and transformation tasks.

#### Example code

```
from lxml import etree
  import requests
2
3 # Fetch the HTML content
4 url = 'https://www.example.com/'
 response = requests.get(url)
5
6 html content = response.text
 # Parse the HTML
7
 parser = etree.HTMLParser()
8
  tree = etree.fromstring(html content, parser)
9
 out list= tree.xpath("//div/text()")
10
```

Syntax	Description	
nodename	Selects nodes with the name nodename	
/	Selects from the root node	
//	Selects nodes from the current node that match the se- lection no matter where they are	
	Selects the current node	
	Selects the parent of the current node	
Q	Selects attributes	
*	Selects all elements/nodes regardless of their name	
@*	Selects all attributes of the current node	
node()	Selects all nodes (element, attribute, text, namespace, processing-instruction, comment, and document nodes)	
[n]	Selects the n-th node (1-based index)	
I	Combines two expressions, selecting nodes that match either expression	
//book[1]	Selects the first book element	
<pre>//book[last()]</pre>	Selects the last book element	
<pre>//book[position()&lt;3]</pre>	Selects the first two book elements	
//book[@attr]	Selects all book elements that have an attribute named ${\tt attr}$	
<pre>//book[@attr='value']</pre>	Selects all book elements where the attr attribute has the value 'value'	
//book[price>35.00]	Selects all book elements with price elements having a value greater than $35.00$	
<pre>//title[@lang='en']</pre>	Selects all $\mathtt{title}$ elements with a $\mathtt{lang}$ attribute value of 'en'	

Table 1: XPath syntax

#### 3.2 Scrapy

For more complex web scraping projects, Scrapy, a comprehensive open-source web crawling and scraping framework, comes into play. Unlike BeautifulSoup, which is more suited for simple, direct web page extraction, Scrapy is built to scrape and crawl at scale. It allows for the extraction of data from websites and the automation of web interactions, making it a formidable tool for gathering data from multiple pages or even entire websites.

#### Scrapy Project Setup

Start a new Scrapy project:

#### scrapy startproject myproject

This creates a new Scrapy project with the name myproject. Generate a Spider:

scrapy genspider myspider example.com

This command generates a spider named myspider for the domain example.com.

#### **Scrapy Components**

**Spider**: Classes that define how a site will be scraped, including how to perform the crawl (i.e., follow links) and how to extract structured data from their pages. Basic spider example:

```
import scrapy
class MySpider(scrapy.Spider):
    name = 'example_spider'
    start_urls = ['http://example.com']

def parse(self, response):
    # Parsing code here
```

**Item**: Standard Python classes used to define the data structure for items you scrape. Example:

```
import scrapy
class MyItem(scrapy.Item):
    name = scrapy.Field()
    description = scrapy.Field()
```

**Item Pipeline**: Process and filter the items returned by spiders. Defined in ITEM\_PIPELINES setting.

Pipeline example:

```
1 class MyPipeline:
2 def process_item(self, item, spider):
3 # Process item here
4 return item
```

#### **Basic Commands**

Running a Spider:

```
scrapy crawl myspider
```

This command runs the spider named **myspider**. Shell:

```
scrapy shell 'http://example.com'
```

Opens the Scrapy shell for the given URL, allowing you to test your extraction code interactively.

#### Selectors

Scrapy uses selectors to extract data from HTML documents. There are two types of selectors: CSS and XPath. CSS:

```
response.css('title::text').get()
```

XPath:

```
response.xpath('//title/text()').get()
```

#### **Requests and Responses**

Following Links:

```
yield response.follow(next_page, self.parse)
```

Form Requests:

### 3.3 Selenium

Essential for dynamic web pages, Selenium automates browser actions, enabling interaction with JavaScript-driven content. Originally for testing, its real-user simulation is crucial for scraping JavaScript-heavy sites, making it invaluable for accessing dynamically loaded data, complementing BeautifulSoup, Requests, and Scrapy in the web scraping toolkit.

Import Selenium:

```
1 from selenium import webdriver
2 from selenium.webdriver.chrome.options import Options
3 from selenium.webdriver.common.by import By
4 5
5 chrome_options = Options()
6 driver = webdriver.Chrome(options=chrome_options)
```

Navigate to a Page:

driver.get('http://example.com')

Locate Elements:

- By ID:
- element = driver.find\_element(By.ID, 'elementId')
- By Name:
- element = driver.find\_element(By.NAME, 'name')
- By XPath:
- By CSS Selector:

• For multiple elements (returns a list):

```
elements = driver.find_elements_by(By.TAG, 'tag')
```

#### Interact with Elements:

- Click a button:
- button.click()
- Enter text in a text field:
- text\_field.send\_keys('text\_to\_enter')

#### Close the Browser:

```
driver.quit()
```

## 4 Setting & installation of required packages

To update and install Python 3 and pip on Ubuntu, use the following commands in your terminal:

```
sudo apt update
sudo apt upgrade -y
sudo apt install python3 python3-pip
```

Create a directory for your scraping project, set up a Python virtual environment, and activate it:

```
mkdir my_scraping_project
cd my_scraping_project
python3 -m venv demoVirEnv
source demoVirEnv/bin/activate
```

Install BeautifulSoup, Requests for fetching web content, Scrapy for more complex scraping, and Selenium for dynamic web pages:

```
    pip install beautifulsoup4
    pip install requests
    pip install lxml
    pip install scrapy
    pip install selenium
```

## References

- [1] W3School HTML. http://https://www.w3schools.com/
- [2] Beautiful Soup Documentation. https://beautiful-soup-4.readthedocs.io/en/ latest/#beautiful-soup-documentation
- [3] XPather. http://xpather.com/
- [4] Scrapy Tutorial. https://docs.scrapy.org/en/latest/intro/tutorial.html