

Lecture 4

*Data Collection II: Web-scrapping Primer; Scrapping Data with **selenium***

Byeong-Hak Choe

bchoe@geneseo.edu

SUNY Geneseo

February 13, 2026



Premier on Web- scrapping



Data Collection via Web-scraping

- Web pages can be a rich data source, but **web scraping is powerful**.
 - Careless scraping can **harm websites, violate rules, or compromise privacy**.
- Our goal in this module:
 - Learn the **web fundamentals** (client/server, HTTPS, URL, HTML/DOM),
 - Understand **ethical, responsible scraping**



“Legal” Is Not the Same as “Ethical”

“If you can see things in your web browser, you can scrape them.”

- *Legally (U.S.):* **publicly available** data may sometimes be scraped using automated tools in US (e.g., **hiQ Labs vs. LinkedIn Corp.**)
- *But legality \neq permission or responsibility:*
 - *Technically:* it may be possible.
 - *Ethically:* you still must consider terms or service (ToS), robots.txt, privacy, and data minimization.
 - *Practically:* you can trigger blocks or harm service quality (e.g., overloading servers, ToS/privacy issues).

Warning

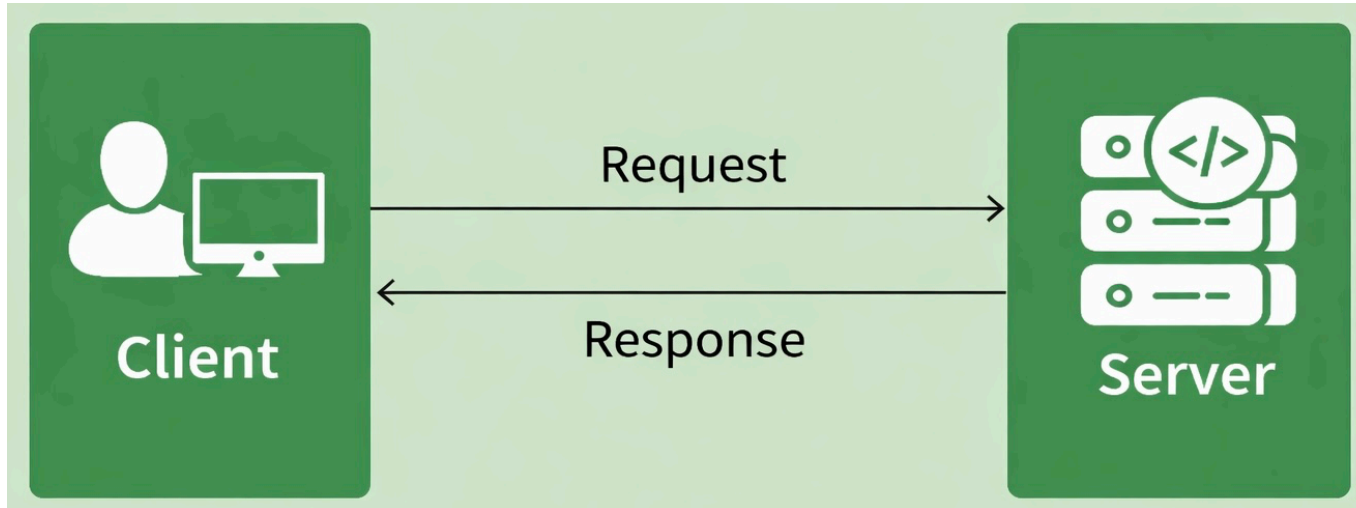
Legal \neq ethical. Even if data is “public,” ToS, privacy expectations, and platform blocks still matter.



Web Basics: Clients and Servers



Clients and Servers



- Devices on the web act as **clients** and **servers**.
- Your browser is a **client**; websites and data live on **servers**.
 - **Client**: your computer/phone + a browser (Chrome/Firefox/Safari).
 - **Server**: a computer that stores webpages/data and sends them when requested.
- When you load a page, your browser sends a **request** and the server sends back a **response** (the page content).



Hypertext Transfer Protocol Secure (HTTPS)

- **HTTP** is how clients and servers communicate.
- **HTTPS** is encrypted HTTP (safer).

When we type a URL starting with [https://](#):

1. Browser finds the server.
2. Browser and server establish a secure connection.
3. Browser sends a request for content.
4. Server responds (e.g., **200 OK**) and sends data.
5. Browser decrypts and displays the page.



HTTP Status Codes

```
1 # library for making HTTPS requests in Python
2 import requests
```

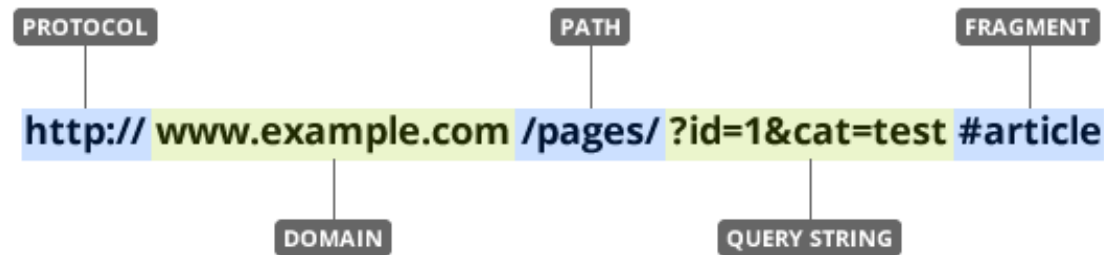
```
1 p = 'https://bcdanl.github.io/210'
2 response = requests.get(p)
3 print(response.status_code)
4 print(response.reason)
```

- **200 OK** → success; content returned.

```
1 p = 'https://bcdanl.github.io/2100'
2 response = requests.get(p)
3 print(response.status_code)
4 print(response.reason)
```

- **404 Not Found** → URL/page doesn't exist (typo, removed page, broken link).

URL (what you're actually requesting)



- A **URL** is a location for a resource on the internet.
- Often includes:
 - Protocol (`https`)
 - Domain (`example.com`)
 - Path (`/products`)
 - **Query string** (`?id=...&cat=...`) ← common in data pages
 - **Fragment** (`#section`) ← in-page reference



HTML Basics



The Big Idea: Scraping = Selecting from HTML

- **HTML** (HyperText Markup Language) is the markup that defines the **structure** of a web page (headings, paragraphs, links, tables, etc.).
- When you “scrape,” you usually:
 1. Load a page
 2. Examine the **HTML**
 3. Extract specific elements (title, price, table, links, etc.)
- **If you don’t understand HTML, you can’t reliably target the right data.**
- Selenium is not “magic”—it automates a browser, but you still need to:
 - Inspect the HTML to identify and target the right elements



HTML in Browser vs. HTML Source Code

DANL 210: Data Preparation and Management, Spring 2026

</> Class Code

Home

Syllabus

Brightspace

Google Colab

Lecture (PDF)

Lecture

Classwork

Homework

Exams

Project

Weeks

DANL 210: Data Preparation and Management, Spring 2026

Instructor: Byeong-Hak Choe ([Email](#))

Welcome! 🙌

— Explore, Learn, and Grow with Data Analytics! 🌟

Lecture

	Title	Subtitle	Date
Lecture 1	Syllabus and Course Outline		January 21, 2026
Lecture 2	Python Fundamentals		January 23, 2026
Lecture 3	Data Collection I: <code>DataFrame</code> ; Spyder IDE; Scrapping Web-tables with <code>pd.read_html()</code>		February 9, 2026
Lecture 4	Data Collection II: Web-scrapping Primer; Scrapping Data with <code>selenium</code>		February 13, 2026

Classwork

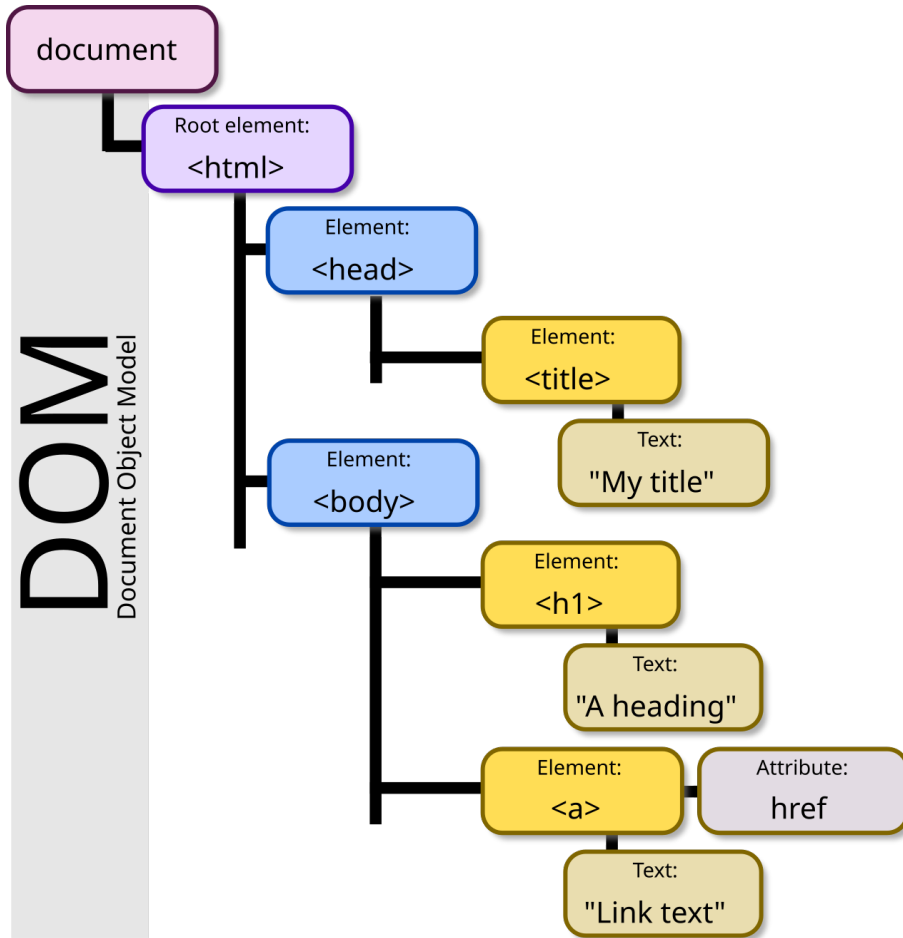
	Title	Subtitle	Date
--	-------	----------	------

```
...<!DOCTYPE html> == $0
<html xmlns="http://www.w3.org/1999/xhtml" lang="en" xml:lang="en">
  <head>
  </head>
  <body class="nav-sidebar docked nav-fixed fullcontent quarto-light" data-bs-offset="78" style="padding-top: 78px;">
    <div id="quarto-search-results"></div>
    <header id="quarto-header" class="headroom fixed-top">
      <!-- content -->
    </header>
    <div id="quarto-content" class="quarto-container page-columns page-rows-contents page-layout-article page-navbar" style="min-height: calc(-184px + 100vh);">
      <!-- sidebar -->
      <nav id="quarto-sidebar" class="sidebar collapse collapse-horizontal sidebar-navigation docked overflow-auto" style="top: 78px; max-height: calc(-78px + 100vh);">
        <div id="quarto-sidebar-glass" data-bs-toggle="collapse" data-bs-target="#quarto-sidebar, #quarto-sidebar-glass"></div>
        <!-- margin-sidebar -->
        <!-- main -->
      </nav>
      <main class="content" id="quarto-document-content">
        <header id="title-block-header" class="quarto-title-block default">
          <div style="display:block; margin:25px;">
          <p></p>
          <div style="display:block; margin:-10px;">
          <p></p>
          <div style="display:block; margin:5px;">
          <section id="lecture" class="level2">
          <section id="classwork" class="level2">
          <section id="homework" class="level2">
          <font size="5">
          <a onclick="window.scrollTo(0, 0); return false;" role="button" id="quarto-back-to-top">
          </a>
        </main>
        <font size="5">
        </div>
        <font size="5">
        </div>
      </body>
    </html>
```



Document Object Model (DOM)

The Browser's "Tree" of the Page



- The browser represents HTML as the **DOM** (Document Object Model).
- Selenium interacts with the **DOM**.
- Scraping often becomes:
 - "Find the node"
 - "Extract its text/attribute"



Inspecting HTML (your #1 web-scraping skill)

- Open a **Chrome** browser.
- Open DevTools:
 - **F12**, or right-click → **Inspect**
- Use it to find:
 - Element text
 - `id / class`
 - Attributes (like `href`, `data-*`)

HTML Elements (what you actually scrape)

- Most HTML is built from **elements** like:

```
1 <tagname>Content goes here...</tagname>
```

- Common ones you'll extract:
 - Headings: `<h1> ... </h1>`
 - Text blocks: `<p> ... </p>`
 - Links: ` ... `
 - Tables: `<table> ... </table>`
 - Containers: `<div> ... </div>`
 - Inline text: ` ... `



HTML Body: Links and Images

<a> (Link)

```
1 <a href="https://www.w3schools.com">This is a link</a>
```

- The **href** attribute is often what you scrape.

 (Image)

```
1 
```

- You may scrape **src** (image URL) or **alt** (description).



HTML Tables

```
1 <table style="width:100%">
2   <tr>
3     <th>Firstname</th>
4     <th>Lastname</th>
5     <th>Age</th>
6   </tr>
7   <tr>
8     <td>Eve</td>
9     <td>Jackson</td>
10    <td>94</td>
11  </tr>
12 </table>
```

- Table structure:
 - `<table>` table container
 - `<tr>` row
 - `<th>` header cell
 - `<td>` data cell

Lists you'll see in the wild

● Unordered List ()

```
1 <ul>
2   <li>Coffee</li>
3   <li>Tea</li>
4   <li>Milk</li>
5 </ul>
```

- Coffee
- Tea
- Milk

Ordered List ()

```
1 <ol>
2   <li>Coffee</li>
3   <li>Tea</li>
4   <li>Milk</li>
5 </ol>
```

1. Coffee
2. Tea
3. Milk

Containers you'll target a lot: `<div>` and ``

`<div>` – *block-level container*

```
1 <div style="background-color:black;color:white;padding:20px;">
2   <h2>Seoul</h2>
3   <p>Seoul is the capital city of South Korea...</p>
4 </div>
```

Seoul

Seoul is the capital city of South Korea...

- Often used to group major page sections.

`` – *inline container*

```
1 <p>My mother has <span style="color:blue;font-weight:bold">blue</span> eyes...</p>
```

My mother has **blue** eyes...