

Summary Report Web Scraping with BeautifulSoup and Requests

Mahir Hashimov
12695120
Unit 3
Deciphering Big Data
University of Essex Online
16 June 2024

Contents

Web Scraping.....	3
Objective	3
Steps Taken.....	3
Conclusion	5
Output	6

Web Scraping

Objective

The objective of this task was to write a web scraping script in Python to extract job listings for the keyword "Data Scientist" from Indeed.com and parse this data into an XML file. The process involved using the BeautifulSoup4 and requests Python libraries.

Steps Taken

1. *Setting Up the Environment*

- Installed necessary libraries using pip:

```
pip install requests
```

```
pip install BeautifulSoup4
```

2. *Identifying the Webpage and Data for Scraping*

- Target URL: <https://www.indeed.com/jobs?q=Data+Scientist&l=>
- Identified that the required data includes job titles, company names, locations, and job links.

3. *Writing the Initial Code*

- Sent an HTTP GET request to the URL.
- Parsed the HTML response using BeautifulSoup.
- Located job postings using HTML element classes and tags.

4. *Handling Errors and Challenges*

- HTTP 403 Error (Forbidden): Initial requests to the website returned a 403 error, indicating forbidden access. This was resolved by adding headers to mimic a browser request:

```
headers = {
```

```
    'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/91.0.4472.124 Safari/537.36'
```

```
}
```

```
response = requests.get(url, headers=headers)
```

5. *Parsing HTML and Extracting Data*

- Used BeautifulSoup to locate the job elements on the page:

```
job_elements = soup.find_all('div', class_='slider_container')
```

- Extracted job titles, company names, locations, and job links from the elements.

6. *Storing Data in XML Format*

- Used the xml.etree.ElementTree library to create an XML structure and save the data:

```
import xml.etree.ElementTree as ET
```

```
jobs = ET.Element('jobs')
```

```
for job_element in job_elements:
```

```
    title_element = job_element.find('h2', class_='jobTitle')
```

```
    company_element = job_element.find('span', class_='companyName')
```

```
    location_element = job_element.find('div', class_='companyLocation')
```

```
    link_element = job_element.find('a', class_='jcs-JobTitle')['href']
```

```
    job = ET.SubElement(jobs, 'job')
```

```
    title = ET.SubElement(job, 'title')
```

```
    title.text = title_element.text.strip()
```

```
    location = ET.SubElement(job, 'location')
```

```
    location.text = company_element.text.strip() +  
location_element.text.strip()
```

```
    link = ET.SubElement(job, 'link')
```

```
    link.text = 'https://www.indeed.com' + link_element
```

```
tree = ET.ElementTree(jobs)
```

```
tree.write('jobs.xml', encoding='utf-8', xml_declaration=True)
```

7. Saving and Reviewing the Output

- Saved the XML file and reviewed it to ensure the correct data was extracted and formatted.

8. Challenges and Resolutions

- 403 Forbidden Error: Initially faced a 403 error which was resolved by adding appropriate headers to the request to mimic a browser.
- Locating HTML Elements: Had to inspect the webpage to find the correct HTML tags and classes to extract the required data.
- XML Formatting: Ensured the extracted data was correctly formatted in the XML structure.

Conclusion

The project involved a combination of web scraping techniques, handling HTTP errors, parsing HTML content, and structuring data in XML format. Despite facing initial challenges with access restrictions and HTML parsing, these were overcome through methodical troubleshooting and adjusting the code to handle different scenarios. The final script successfully extracts job listings from Indeed.com and saves them in an XML file, which was then added to a GitHub repository for easy access and sharing.

Output

This XML file does not appear to have any style information associated with it. The document tree is shown below.

```
<?xml>

</jobs>

<?xml>
<title>Data Scientist</title>
<location>CVS HealthRemote in New York, NY</location>
<link>https://www.indeed.com/rc/clk?jk=61d7bcae075ee388bb-4_QymOFDsJlMQ5nqptaxiab5gPueh8mb8441s08FLpCMF9cBh0btyhw58vmlz3w7huarRh_hlMmeT0BQ-zk-zTPwArSEllyyVnGHP64u8Jn64N72d3K3Dkukcb-SuB67W3ABqtdMap3N0bzKdCdpP8fccId=bc3b1aa573faee78rJj=3-/link>
</jobs>
<?xml>
<title>Data Scientist</title>
<location>JP Morgan Chase & CoNew York, NY 10017 (Midtown area)</location>
<link>https://www.indeed.com/rc/clk?jk=9c2522286993e9e58bb-4_QymOFDsJlMQ5nqptaw5TQ7d69yQpooP-PkZv8uH16KxrmPMPW-DC4xT8BnncVyyQpIgy-8BQR1TGEcy2065z21dEVDHf--hnt_Ayyg72W9XQ3K3Dkukcb-SuB67W3ABqtdMap3N0bzKdCdpP8fccId=46d8116f6e9eae8rJj=3-/link>
</jobs>
<?xml>
<title>Data Scientist - All Levels</title>
<location>Interclipse, Inc.San Antonio, TX</location>
<link>https://www.indeed.com/rc/clk?jk=ffdl7be34695286d8bb-4_QymOFDsJlMQ5nqpta1VhlymsUc1P037mu53bV1s5ReFz5TwAqP-N0Y9QxTYK223d68ADIpabL8uPVIr8TdxQ2Cshrp_TpZP3amP5Q3d6Ue1DVP4BPz1Kx61dKukcb-SuB67W3ABqtdMap3N0bzKdCdpP8fccId=46d8116f6e9eae8rJj=3-/link>
</jobs>
<?xml>
<title>Data Scientist</title>
<location>SpectrumCharlotte, NC 28217</location>
<link>https://www.indeed.com/rc/clk?jk=6b7b7985803d918618bb-4_QymOFDsJlMQ5nqpta65F7jcdR6-V6uk28P32m3H0DnTg75H0uffenT3aIreP0amgDATUCXtqE-Jy3b6v8u7L1kxLk4hV019AP2_Fgg9h4S0K3Dkukcb-SuB67W3ABqtdMap3N0bzKdCdpP8fccId=2d1d876e534b4d8rJj=3-/link>
</jobs>
<?xml>
<title>Senior Data Scientist (Remote)</title>
<location>Stryker CorporationRemote in Portage, MI 49802</location>
<link>https://www.indeed.com/rc/clk?jk=9283965db2e1c9288bb-4_QymOFDsJlMQ5nqpta_b4tGcmfUF_jpg6Qp9YEP-98Pygah7i7Dq5ba8ducAAJAdnd62971G1aUc0P9X0b1mP2_SkN8U1srN8YVnL7Md6v8uEQR3K3Dkukcb-SuB67W3ABqtdMap3N0bzKdCdpP8fccId=9d4d8ca3e88b9c8ffrJj=3-/link>
</jobs>
<?xml>
<title>Data Scientist</title>
<location>BECWashington State</location>
<link>https://www.indeed.com/rc/clk?jk=58ea5b849a79328bb-4_QymOFDsJlMQ5nqpta6v1d5JxV1CCq5gt48LQcLWFED1j1P1Jm88Rhm2-9X0FPe1BzKb07FT1P1VhNhuwPQg_2gesPM2d8uR1cn80JeggY961c5QR3K3Dkukcb-SuB67W3ABqtdMap3N0bzKdCdpP8fccId=8ab8e1618cb3cfc8rJj=3-/link>
</jobs>
<?xml>
<title>Data Scientist</title>
<location>Tri-Force Consulting ServicesChesterfield, VA</location>
<link>https://www.indeed.com/rc/clk?jk=38819119cc8d8d8bb-4_QymOFDsJlMQ5nqpta8tz-R1-Tmu2RFBvt-371Dn_M_VhmpwZDATawTCN71x5CD6--l9ys7PHx_gbuU1ZMlbfGndfL23w123ZQh4PQdW3DIPFwEfwXqQXT&kukcb-SuB67W3ABqtdMap3N0bzKdCdpP8fccId=e58F89388b8f5da8rJj=3-/link>
</jobs>
<?xml>
<title>Data Scientist - Generative AI (GenAI, Applied ML, LLMs)</title>
<location>TargetMinneapolis, MN 55449</location>
<link>https://www.indeed.com/rc/clk?jk=5ae457bee98b74f88bb-4_QymOFDsJlMQ5nqpta2a8_tgtw7FUE-UBQm7p3m89B72f6_m8XEUB8eupC83r-sr-Vh3W8hC1he1Q8u7aunwP18K1qf7g4tDVEfh2uW3K3Dkukcb-SuB67W3ABqtdMap3N0bzKdCdpP8fccId=15f43d82d6301ff28rJj=3-/link>
</jobs>
<?xml>
<title>Business Data Scientist, Devices and Services Systems</title>
<location>GoogleSan Francisco, CA</location>
<link>https://www.indeed.com/rc/clk?jk=9646185c818823de8bb-4_QymOFDsJlMQ5nqptaw02-qilp5fr-6J9tIClYgnk11P6Idqez8PqP_U8U81H-38Fw89In1hJy_Rw8WV_61s18rbuOVQd8c1ey6f865QXdeQR3K3Dkukcb-SuB67W3ABqtdMap3N0bzKdCdpP8fccId=a58449d8e91a5c8rJj=3-/link>
</jobs>
<?xml>
<title>Data Scientist</title>
<location>Baylor College of MedicineHybrid work in United States</location>
<link>https://www.indeed.com/rc/clk?jk=5684bcf83c8e9af8bb-4_QymOFDsJlMQ5nqptaVLDvxd8x50Rg18xJcsm48D8k0JfchevAU07-5_F3xJh88d8d8c8t8k18z8m8d8y8r7y8y8PvU7F8W8P1C8WQJ88fID_4as8xIE88kukcb-SuB67W3ABqtdMap3N0bzKdCdpP8fccId=5f589639cd588898rJj=3-/link>
</jobs>
<?xml>
<title>Data Scientist</title>
<location>Brigham Young UniversityProvo, UT (University area)</location>
<link>https://www.indeed.com/rc/clk?jk=5684bcf83c8e9af8bb-4_QymOFDsJlMQ5nqptaVLDvxd8x50Rg18xJcsm48D8k0JfchevAU07-5_F3xJh88d8d8c8t8k18z8m8d8y8r7y8y8PvU7F8W8P1C8WQJ88fID_4as8xIE88kukcb-SuB67W3ABqtdMap3N0bzKdCdpP8fccId=5f589639cd588898rJj=3-/link>
</jobs>
```