☰  ⬤  Wizard-Fingerz  /  **AppClickSeptemberCohort**              🔍  🗔  {⟨⟩}

`<>` **Code**    ⊙ Issues    �units Pull requests    ⏵ Actions    ▦ Projects    📖 Wiki    ⚠ Security    ⬚ In

**AppClickSeptemberCohort** / **Week4Class3.md**    ⧉                                                         ⋯

> ⬤ **Wizard-Fingerz** Add comprehensive guide for building a web scraper using Beautiful Soup
>
> fd4a183 · 2 days ago   🕓

251 lines (170 loc) · 5.91 KB

| Preview   Code   Blame | | 🎱  Raw ⧉ ⬇  ✏ ▾  ☰ |
| --- | --- |

# Building a Web Scraper Application using Beautiful Soup

## 🧠 1. What is Web Scraping?

**Web scraping** is the process of automatically extracting data from websites. Python is one of the most popular languages for this task — thanks to libraries like:

- `requests` → for fetching web pages
- `beautifulsoup4` → for parsing HTML and extracting data
- `lxml` → for fast XML/HTML parsing (optional)
- `pandas` → for storing and exporting scraped data (optional)

## ⚙️ 2. Installing Required Libraries

Before starting, you'll need to install the required Python packages.

```
pip install requests beautifulsoup4 lxml pandas
```

✅ **Explanation:**

- `requests` : sends HTTP requests to fetch webpage content
- `beautifulsoup4` : helps parse and search the HTML structure
- `lxml` : speeds up parsing
- `pandas` : organizes the data into tables or exports it (e.g., to CSV)

## 🌐 3. Importing Required Libraries

```python
import requests
from bs4 import BeautifulSoup
import pandas as pd
```

## 🧩 4. Fetching a Web Page

We'll start by fetching a webpage. For example, let's scrape **quotes** from: 👉
`https://quotes.toscrape.com`

```python
url = "https://quotes.toscrape.com"
response = requests.get(url)

# Check if request was successful
print(response.status_code)
```

✅ **Expected Output:** `200` means the page loaded successfully.

## 📑 5. Parsing HTML Content

Now, let's load the page into **BeautifulSoup** for parsing.

```python
soup = BeautifulSoup(response.text, 'lxml')

# Print the first 500 characters of the page to inspect
print(soup.prettify()[:500])
```

This will show you the HTML structure — which we'll use to find elements.

## 🔍 6. Finding and Extracting Data

Let's extract all **quotes** and **authors** from the page.

Inspect the page in your browser → you'll find that:

- Each quote is inside `<div class="quote">`
- The quote text is in `<span class="text">`
- The author is in `<small class="author">`

```python
quotes = soup.find_all('div', class_='quote')

for quote in quotes:
    text = quote.find('span', class_='text').text
    author = quote.find('small', class_='author').text
    print(f'"{text}" — {author}')
```

### ✅ Output Example:

```
"The world as we have created it is a process of our thinking." — Albert
Einstein
"It is our choices, Harry, that show what we truly are." — J.K. Rowling
```

## 📄 7. Saving Extracted Data

We can store all data in a **list of dictionaries** and save it to a CSV file.

```python
data = []

for quote in quotes:
    text = quote.find('span', class_='text').text
    author = quote.find('small', class_='author').text
    data.append({'Quote': text, 'Author': author})

# Convert to DataFrame and save to CSV
df = pd.DataFrame(data)
df.to_csv('quotes.csv', index=False)
```

```
print("✅ Data saved successfully to quotes.csv")
```

Now you'll have a file like:

| Quote | Author |
|---|---|
| "Life is what happens…" | John Lennon |
| "The world as we…" | Albert Einstein |

## 🔁 8. Scraping Multiple Pages

Most websites with paginated content have a **"Next"** button or URL pattern like:

```
https://quotes.toscrape.com/page/2/
```

Let's loop through multiple pages.

```python
page = 1
data = []

while True:
    url = f"https://quotes.toscrape.com/page/{page}/"
    response = requests.get(url)
    soup = BeautifulSoup(response.text, 'lxml')
    quotes = soup.find_all('div', class_='quote')

    if not quotes:
        break  # Stop when no more quotes

    for quote in quotes:
        text = quote.find('span', class_='text').text
        author = quote.find('small', class_='author').text
        data.append({'Quote': text, 'Author': author})

    page += 1

print(f"Scraped {len(data)} quotes!")
pd.DataFrame(data).to_csv('all_quotes.csv', index=False)
```

✅ **Result:** All quotes across pages will be saved to `all_quotes.csv`.

## 🗨️ 9. Real-World Tips for Web Scraping

1. **Always check robots.txt** Example: `https://example.com/robots.txt` → tells what's allowed to scrape.

2. **Add headers** to mimic a real browser:

```python
headers = {'User-Agent': 'Mozilla/5.0'}
response = requests.get(url, headers=headers)
```

3. **Avoid scraping too fast** — add short delays:

```python
import time
time.sleep(1)
```

4. **Use exception handling** to avoid crashes if a page fails to load.

5. **Don't scrape private or copyrighted content.**

## 💼 10. Mini Project Idea

**Project:** Build a scraper that:

- Extracts the latest job listings from a site like `https://realpython.github.io/fake-jobs/`
- Saves **Job Title**, **Company**, **Location**, and **Link** to a CSV file

**Bonus:** Schedule it to run every day and email you new results.

## 🧩 11. Example: Job Scraper Code

```python
import requests
from bs4 import BeautifulSoup
import pandas as pd

url = "https://realpython.github.io/fake-jobs/"
response = requests.get(url)
soup = BeautifulSoup(response.text, "lxml")
```

```python
jobs = soup.find_all('div', class_='card-content')

data = []
for job in jobs:
    title = job.find('h2', class_='title').text.strip()
    company = job.find('h3', class_='company').text.strip()
    location = job.find('p', class_='location').text.strip()
    link = job.find('a')['href']
    data.append({
        'Title': title,
        'Company': company,
        'Location': location,
        'Link': link
    })

df = pd.DataFrame(data)
df.to_csv("job_listings.csv", index=False)
print("✅ Jobs scraped and saved to job_listings.csv!")
```

# 🧩 12. Next Steps (Advanced Topics)

- Scrape with **Selenium** for dynamic sites (JavaScript-rendered)
- Use **APIs** instead of scraping where possible
- Store scraped data in **databases** (SQLite, PostgreSQL)
- Create a **Flask/Django dashboard** to display scraped data