



Beautiful Soup



Beautiful Soup

A python Package Turning the Web into Data
for Data Scientists

supriya Sinha
MBA'26



Beautiful Soup

Introduction

- “Have you ever needed to extract data from a webpage when no API existed?”
- **Beautiful Soup** as a Python library for web scraping and parsing HTML/XML.
- Created by **Leonard Richardson**, with its initial release in 2004
- objective: Show how **Beautiful Soup** helps data scientists collect and clean unstructured web data.
- Document-
<https://beautiful-soup-.readthedocs.io/en/latest/>



About

- Collect data from websites (news, products, jobs, etc.)
- Build custom datasets for ML or NLP
- Track competitors / sentiment / pricing
- Clean messy web text into structured data
- Example Use Cases:
 - Scraping Amazon reviews for sentiment
 - Collecting article headlines for trend analysis



Integration in Data Science Workflow

WEB

**Beautiful
Soup**

Pandas

Analysis/ML

- Beautiful Soup = web → data bridge
- Great for quick, lightweight scraping
- Works seamlessly with Pandas / NLP / visualization
- Essential for data scientists collecting unstructured data

Core Features and Syntax

```
from bs4 import BeautifulSoup
import requests

# Fetch the HTML content of a webpage
url = "http://example.com"
response = requests.get(url)
html_content = response.text

# Create a Beautiful Soup object
soup = BeautifulSoup(html_content, 'html.parser')

# Find all paragraph tags
paragraphs = soup.find_all('p')

# Print the text content of each paragraph
for p in paragraphs:
    print(p.get_text())

# Find an element by its ID
element_by_id = soup.find(id='some_id')

# Find elements using a CSS selector
elements_by_css = soup.select('div.container a')
```



```
left: 10px;
right: 10px;
ng-left: 10px;
ng-right: 10px;
weight: bold;
low: auto;

<login id="login_form" >
    <div color='red'><b>Authent
    <div id="dError1" class="dError1
        <saml-auth-status>-1</s
        <script>window.top.lo
```



Beautiful Soup

Best Practice



- Use requests headers to mimic real browsers.
Example-
 - `headers = {"User-Agent": "Mozilla/5.0"}`
 - `requests.get(url, headers=headers)`
- Respect website rules — scrape responsibly and slowly (`time.sleep()` between requests).
- Validate and clean your HTML data (handle missing tags, broken links).
- Store results in structured formats (CSV, JSON, SQL).
- Log your scraping pipeline for reproducibility.
- Prefer APIs when available — cleaner, legal, and more stable.



Beautiful Soup

Limitations



- **Dynamic Content:** BeautifulSoup cannot render JavaScript; only parses static HTML.
- **Performance:** Slower for large-scale scraping; not ideal for big data extraction.
- **Fragile Structure:** Breaks if the website layout (HTML tags/classes) changes.
- **Legal & Ethical Concerns:** Some websites prohibit scraping — always check robots.txt or Terms of Service.
- **Rate Limits / Blocking:** Too many rapid requests may get your IP blocked.