



Web Scraping in Python

What is Web Scraping?

Web scraping is the process of extracting data from websites and converting it into a structured format for analysis or storage. Python provides several libraries like **BeautifulSoup**, **Scrapy**, and **Selenium** to automate this process. Web scraping allows users to collect publicly available data efficiently without manual copying.

Use Cases of Web Scraping in Python

1. **Price Monitoring and Comparison** – Businesses track competitor pricing for dynamic pricing strategies.
2. **Market Research and Data Collection** – Researchers collect large amounts of online data for analysis.
3. **Lead Generation** – Companies scrape business directories and social media for contact information.
4. **News Aggregation** – Automating the collection of news articles from multiple sources.
5. **SEO and Web Analytics** – Extracting search engine rankings, keyword trends, and competitor data.
6. **Stock Market and Financial Data** – Fetching live or historical financial data for investment analysis.
7. **E-commerce Product Scraping** – Gathering product details, reviews, and ratings for trend analysis.

Benefits of Web Scraping in Python

- **Automation** – Eliminates the need for manual data extraction.
- **Efficiency** – Extracts large volumes of data quickly.
- **Cost-Effective** – Reduces the cost of data collection compared to hiring manual workers.
- **Real-time Updates** – Enables tracking of real-time data for informed decision-making.
- **Scalability** – Easily handles scraping of thousands or millions of web pages.

Python's web scraping capabilities make it a powerful tool for businesses, researchers, and analysts looking to extract valuable insights from online sources.

To begin with the Lab

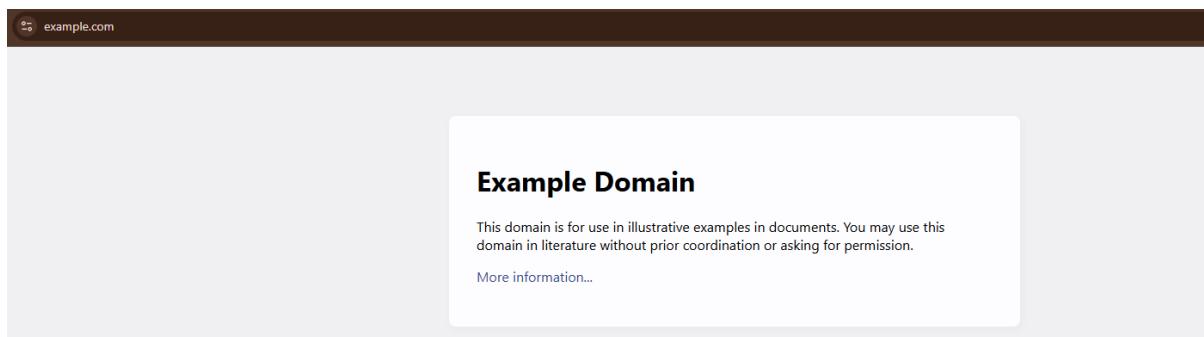
1. Before starting this, you should have a prior knowledge of basic HTML and CSS.
2. Now you need to install some libraries to work in this lab.

```
conda install requests
conda install lxml
conda install bs4
```

if you are not using the Anaconda Installation, you can use **pip install** instead of **conda install**, for example:

```
pip install requests
pip install lxml
pip install bs4
```

3. In a new tab, open example.com and right click anywhere on the page to open the menu, and from that choose view page source.



4. You will be able to see the basic HTML code for the page. Now we have done this to have a basic understanding of how things are going to work in this lab.

```
Line wrap □
1 <!DOCTYPE html>
2 <html>
3 <head>
4   <title>Example Domain</title>
5
6   <meta charset="utf-8" />
7   <meta http-equiv="Content-type" content="text/html; charset=utf-8" />
8   <meta name="viewport" content="width=device-width, initial-scale=1" />
9   <style type="text/css">
10  body {
11    background-color: #f0f0f2;
12    margin: 0;
13    padding: 0;
14    font-family: -apple-system, system-ui, BlinkMacSystemFont, "Segoe UI", "Open Sans", "Helvetica Neue", Helvetica, Arial, sans-serif;
15  }
16  div {
17    width: 600px;
18    margin: 5em auto;
19    padding: 2em;
20    background-color: #fdfdff;
21    border-radius: 0.5em;
22    box-shadow: 2px 3px 7px 2px rgba(0,0,0,0.02);
23  }
24  a:link, a:visited {
25    color: #38488f;
26    text-decoration: none;
27  }
28  @media (max-width: 700px) {
29    div {
30      margin: 0 auto;
31      width: auto;
32    }
33  }
34  </style>
35 </head>
36
37 <body>
38 <div>
39   <h1>Example Domain</h1>
40   <p>This domain is for use in illustrative examples in documents. You may use this
41   domain in literature without prior coordination or asking for permission.</p>
42   <p><a href="https://www.iana.org/domains/example">More information...</a></p>
43 </div>
44 </body>
45 </html>
46
47
```

5. Come back to Jupyter Notebook and start writing your code. First, start by importing requests, and we are defining the type(res) for the requests.

- The code uses the **requests** library in Python to send an HTTP GET request to the website "<http://www.example.com>" and stores the response in the variable **res**. The **type(res)** command checks the type of the response object, which is typically **requests.models.Response**. Finally, **res.text** extracts and returns the content of the webpage as a string, usually containing the HTML structure of the page. If there are network restrictions or firewalls, the request may fail, requiring multiple attempts.
- In cell 8, you can see that we got the source in the form of a string and we can change that using the beautiful soup from **bs4** library.

```
[2]: import requests
[3]: # Step 1: Use the requests Library to grab the page
# Note, this may fail if you have a firewall blocking Python/Jupyter
# Note sometimes you need to run this twice if it fails the first time
res = requests.get("http://www.example.com")
[4]: type(res)
[4]: requests.models.Response
[8]: res.text
[8]: '<!DOCTYPE html><html><head><title>Example Domain</title></head><body><meta charset="utf-8" /><meta http-equiv="Content-type" content="text/html; charset=utf-8" /><meta name="viewport" content="width=device-width, initial-scale=1" /><style type="text/css">n body {n background-color: #f0f0f2;n margin: 0;n padding: 0;n font-family: -apple-system, system-ui, BlinkMacSystemFont, "Segoe UI", "Open Sans", "Helvetica Neue", Helvetica, Arial, sans-serif;n }n div {n width: 600px;n margin: 5em auto;n padding: 2em;n background-color: #fdfdff;n border-radius: 0.5em;n box-shadow: 2px 3px 7px 2px rgba(0,0,0,0.02);n }n a:link, a:visited {n color: #38488E;n text-decoration: none;n }n @media (max-width: 700px) {n div {n margin: 0 auto;n width: auto;n }n }</style></head><body><div><h1>Example Domain</h1><p>This domain is for use in illustrative examples in documents. You may use this domain in literature without prior coordination or asking for permission.</p><p><a href="https://www.iana.org/domains/example">More information...</a></p></div></body></html>'
```

- The code uses the **BeautifulSoup** library from **bs4** to parse the HTML content retrieved from the previous request (**res.text**). The "**lxml**" parser is specified for efficient and fast parsing. The resulting **soup** object represents the structured HTML document, allowing easy navigation, searching, and manipulation of elements within the webpage.

```
[10]: import bs4
[12]: soup = bs4.BeautifulSoup(res.text,"lxml")
[14]: soup
[14]: <!DOCTYPE html>
<html>
<head>
<title>Example Domain</title>
<meta charset="utf-8"/>
<meta content="text/html; charset=utf-8" http-equiv="Content-type"/>
<meta content="width=device-width, initial-scale=1" name="viewport"/>
<style type="text/css">
  body {
    background-color: #f0f0f2;
    margin: 0;
    padding: 0;
    font-family: -apple-system, system-ui, BlinkMacSystemFont, "Segoe UI", "Open Sans", "Helvetica Neue", Helvetica, Arial, sans-serif;
  }
  div {
    width: 600px;
    margin: 5em auto;
```

- The code extracts the **<title>** tag from the parsed HTML using **soup.select('title')**, which returns a list of matching elements. The first element (**title_tag[0]**) is accessed, confirming its type as a **Tag** object. Finally, **title_tag[0].get_text()** retrieves the text content inside the **<title>** tag, effectively extracting the page's title.

```
[16]: soup.select('title')

[16]: [<title>Example Domain</title>]

[18]: title_tag = soup.select('title')

[20]: title_tag[0]

[20]: <title>Example Domain</title>

[22]: type(title_tag[0])

[22]: bs4.element.Tag

[24]: title_tag[0].getText()

[24]: 'Example Domain'
```

10. Now we are going to **Grab all the elements of a class**, for that we are going to use a Wikipedia page link https://en.wikipedia.org/wiki/Grace_Hopper.

11. If you go to this page and open the inspect menu by right-clicking in any white space area of the page then you will be able to see the elements of the page.

The screenshot shows a Wikipedia article about Grace Hopper. The page content includes her biography, a photo, and a table of awards. To the right, the Chrome DevTools Elements tab is open, displaying the DOM tree. The tree shows various HTML elements like `<div>`, ``, and `` with classes such as `vector-toc-level-1`, `vector-toc-list-item`, and `vector-toc-list-item-expanded`. The `title` element is also visible in the tree.

12. The code retrieves the HTML content of the Wikipedia page for **Grace Hopper** using the requests library and stores the response. It then parses the HTML using BeautifulSoup with the **lxml** parser, creating a soup object. This object represents the entire webpage structure, allowing easy extraction and manipulation of its elements, such as text, links, and tables.

```
[26]: # First get the request
res = requests.get('https://en.wikipedia.org/wiki/Grace_Hopper')

[28]: # Create a soup from request
soup = bs4.BeautifulSoup(res.text,"lxml")

[30]: soup
```

[30]: <!DOCTYPE html>
<html class="client-nojs vector-feature-language-in-header-enabled vector-feature-language-in-main-page-header-disabled vector-feature-page-tools-pinned-disabled vector-feature-toc-pinned-clientpref-1 vector-feature-main-menu-pinned-disabled vector-feature-limited-width-clientpref-1 vector-feature-limited-width-content-enabled vector-feature-custom-font-size-clientpref-1 vector-feature-appearance-pinned-clientpref-1 vector-feature-night-mode-enabled skin-theme-clientpref-day vector-sticky-header-enabled vector-toc-available" dir="ltr" lang="en">
<head>
<meta charset="utf-8"/>
<title>Grace Hopper - Wikipedia</title>
<script>(function(){var className="client-js vector-feature-language-in-header-enabled vector-feature-language-in-main-page-header-disabled vector-feature-page-tools-pinned-disabled vector-feature-toc-pinned-clientpref-1 vector-feature-main-menu-pinned-disabled vector-feature-limited-width-clientpref-1 vector-feature-limited-width-content-enabled vector-feature-custom-font-size-clientpref-1 vector-feature-appearance-pinned-clientpref-1 vector-feature-night-mode-enabled skin-theme-clientpref-day vector-sticky-header-enabled vector-toc-available";var cookie=document.cookie.match(/(?:^| ;)enwikimwclientprefReferences=(\[\.\.\.\])+);if(cookie){cookie[1].split('%2C').forEach(function(pref){className.replace(new RegExp('(^|)'+pref.replace(/-/g,'_')),'\$1clientpref-\w+\|\w+|g,'')+clientpref-\w+(\ \$)', '\$1'+pref+\$2)});document.documentElement.className=className;}();RLCONF={"wgBreakFrames":false,"wgSeparatorTransformTable":["",""], "wgDigitTransformTable":["",""], "wgDefaultDateFormat":"dmy","wgMonthNames":["","January","February","March","April","May","June","July","August","September","October","November","December"], "wgRequestId":"97b025c1-91b9-4f1e-97a1-555401da1ab7", "wgCanonicalNamespace":"","wgCanonicalPageName":false, "wgNamespaceNumber":0, "wgPageName":"Grace_Hopper", "wgTitle":"Grace Hopper", "wgCurRevisionId":1283235419, "wgRevisionId":1283235419, "wgArticleId":12500, "wgIsAnon":true, "wgIsRedirect":false, "wgRevision":null, "wgUserName":null, "wgUserGroups":["All pages", "Administrators"]});</script>

13. The code attempts to extract all elements from the Wikipedia page that have the class "**vector-toc-numb**" using `soup.select()`. This class is commonly associated with **table of contents numbering** on Wikipedia. However, the exact class name might vary based on the Wikipedia page's structure or depending on regional variations of the site. If the class is different, the selector may return an empty list.

```
[34]: # note depending on your IP Address,
# this class may be called something different
soup.select(".vector-toc-numb")
```

```
[34]: [<span class="vector-toc-numb">1</span>,
        <span class="vector-toc-numb">2</span>,
        <span class="vector-toc-numb">2.1</span>,
        <span class="vector-toc-numb">2.2</span>,
        <span class="vector-toc-numb">2.3</span>,
        <span class="vector-toc-numb">2.4</span>,
        <span class="vector-toc-numb">3</span>,
        <span class="vector-toc-numb">4</span>,
        <span class="vector-toc-numb">5</span>,
        <span class="vector-toc-numb">6</span>,
        <span class="vector-toc-numb">7</span>,
        <span class="vector-toc-numb">8</span>,
        <span class="vector-toc-numb">8.1</span>,
        <span class="vector-toc-numb">8.2</span>,
        <span class="vector-toc-numb">9</span>,
        <span class="vector-toc-numb">9.1</span>,
        <span class="vector-toc-numb">9.2</span>,
        <span class="vector-toc-numb">9.3</span>]
```

14. This code extracts and prints the text content of all HTML elements with the class "vector-toc-numb" from the Wikipedia page on Grace Hopper. These elements likely represent section numbers in the table of contents. The `soup.select()` method finds all matching elements, and the loop iterates over them, printing their text values.

```
[36]: for item in soup.select(".vector-toc-numb"):
    print(item.text)
```

```
1
2
2.1
2.2
2.3
2.4
3
4
5
6
7
8
8.1
8.2
9
9.1
9.2
9.3
```

15. Now we are going to get an image from a Wikipedia page, the initial steps are same as before.

```
[38]: res = requests.get("https://en.wikipedia.org/wiki/Deep_Blue_(chess_computer)")
[40]: soup = bs4.BeautifulSoup(res.text, 'lxml')
[50]: soup
```

```
<!DOCTYPE html>
<html class="client-nojs vector-feature-language-in-header-enabled vector-feature-language-in-main-page-header-disabled vector-feature-page-tools-pinned-disabled vector-feature-toc-pinned-clientpref-1 vector-feature-main-menu-pinned-disabled vector-feature-limited-width-clientpref-1 vector-feature-limited-width-content-enabled vector-feature-custom-font-size-clientpref-1 vector-feature-appearance-pinned-clientpref-1 vector-feature-night-mode-enabled skin-theme-clientpref-day vector-sticky-header-enabled vector-toc-available" dir="ltr" lang="en">
<head>
<meta charset="utf-8"/>
<title>Deep Blue (chess computer) - Wikipedia</title>
<script>(function(){var className="client-js vector-feature-language-in-header-enabled vector-feature-language-in-main-page-header-disabled vector-feature-page-tools-pinned-disabled vector-feature-toc-pinned-clientpref-1 vector-feature-main-menu-pinned-disabled vector-feature-limited-width-clientpref-1 vector-feature-limited-width-content-enabled vector-feature-custom-font-size-clientpref-1 vector-feature-appearance-pinned-clientpref-1 vector-feature-night-mode-enabled skin-theme-clientpref-day vector-sticky-header-enabled vector-toc-available";var cookie=document.cookie.match(/(?:^| )_jenwikimwclientprefReferences\[(^;+|)\);if(cookie){cookie[1].split(';%2C').forEach(function(pref){className.replace(new RegExp('(^| )'+pref.replace(/-/clientpref-\w+$|[^w-]+/g,'')+'-clientpref-\w+( |$)'),'$1'+pref+'$2'))};}document.documentElement.className=className;})();RLCONF={"wgBreakFrames":false,"wgSeparatorTransformTable":["",""], "wgDigitTransformTable":["",""], "wgDefaultDateFormat": "dmy", "wgMonthNames": ["", "January", "February", "March", "April", "May", "June", "July", "August", "September", "October", "November", "December"], "wgRequestId": "65f34b14-050d-4681-860f-539cdf461348", "wgCanonicalNamespace": "", "wgCanonicalSpecialPageName": false, "wgNamespaceNumber": 0, "wgPageName": "Deep_Blue_(chess_computer)", "wgTitle": "Deep Blue (chess computer)", "wgCurRevisionId": 128001401, "wgRevisionId": 128001401, "wgArticleId": 49387, "wgIsArticle": true, "wgIsRedirect": false, "wgIsMain": "view", "wgIsContent": null, "wgIsContentGroup": null}
```

16. The code searches the parsed HTML (soup) for elements with the class .mw-file-element. In Wikipedia pages, this class is commonly used for images or media files embedded in the article. The output will be a list of elements (typically tags) that match this class, allowing further extraction of image URLs or related attributes.

```
[60]: soup.select('.mw-file-element')

[60]: [, , , , , ]
```

17. This code extracts the second image element from the Wikipedia page using the **.mw-file-element** class, retrieves its **src** attribute (the image URL), and then downloads the image using `requests.get()`. The `image_link.content` holds the raw binary data of the image, which can be saved to a file using binary write mode (`wb`).

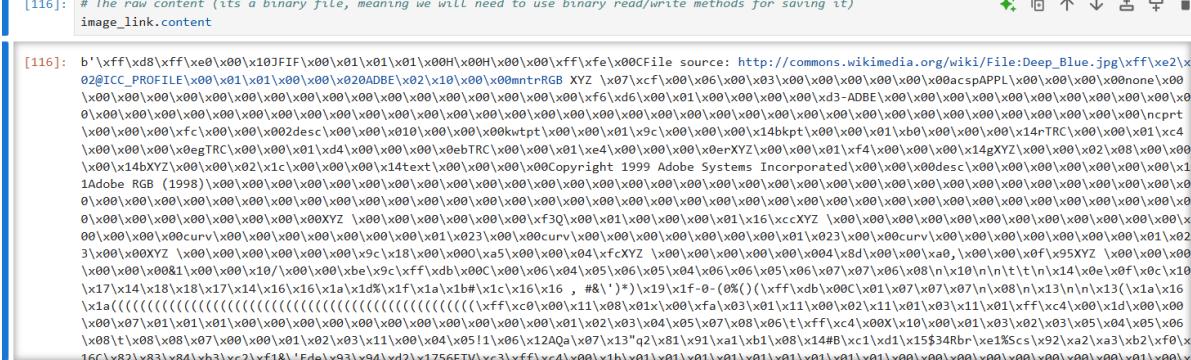
```
[86]: image_info = soup.select('.mw-file-element')[1]

[108]: image_info['src']

[108]: '//upload.wikimedia.org/wikipedia/commons/thumb/b/be/Deep_Blue.jpg/250px-Deep_Blue.jpg'

[114]: image_link = requests.get('https://upload.wikimedia.org/wikipedia/commons/thumb/b/be/Deep_Blue.jpg/250px-Deep_Blue.jpg')

[116]: # The raw content (its a binary file, meaning we will need to use binary read/write methods for saving it)
      image_link.content
```



18. Now, using the code below, we open a new file named `my_new_file_name.jpg` in **binary write mode (wb)**, write the downloaded image data (`image_link.content`) into the file, and then close the file to ensure proper saving. This effectively stores the image on your local system.

```
[118]: f = open('my_new_file_name.jpg', 'wb')

[120]: f.write(image_link.content)

[120]: 20190
```

19. So, in your current directory of Jupyter Notebook, you will see that an image has been created or downloaded. Now open this image.

Screenshot of a Jupyter Notebook interface showing a file list and a thumbnail image.

File List:

Name	Last Modified	File Size
Example_Top_Level	2 days ago	
extracted_content	yesterday	
final_unzip	yesterday	
target	10 months ago	
Untitled.ipynb	2 hours ago	1.6 KB
webScraping.ipynb	5 minutes ago	894 KB
zippingAndUnzipping.ipynb	yesterday	4.5 KB
comp_file.zip	yesterday	262 B
example.zip	yesterday	262 B
my_new_file_name.jpg	32 minutes ago	19.7 KB

Thumbnail Image:

20. Now we are moving to another example that is working with multiple pages and items.
21. The code below scrapes book listings from **Books to Scrape**, a test website for web scraping practice. It constructs a **URL template** (`base_url`) for different pages by using `.format()`. Then, it sends a **GET request** to fetch page 1, converts the response into a **BeautifulSoup object**, and selects all elements with the class `.product_pod`, which represent individual books.

```
[124]: base_url = 'http://books.toscrape.com/catalogue/page-{}.html'
[126]: res = requests.get(base_url.format('1'))
[128]: soup = bs4.BeautifulSoup(res.text,"lxml")
[130]: soup.select(".product_pod")

[130]: [<article class="product_pod">
    <div class="image_container">
        <a href="a-light-in-the-attic_1000/index.html"></a>
    </div>
    <p class="star-rating Three">
        <i class="icon-star"></i>
        <i class="icon-star"></i>
        <i class="icon-star"></i>
        <i class="icon-star"></i>
        <i class="icon-star"></i>
    </p>
    <h3><a href="a-light-in-the-attic_1000/index.html" title="A Light in the Attic">A Light in the ...</a></h3>
    <div class="product_price">
        <p class="price_color">£51.77</p>
        <p class="instock availability">
            <i class="icon-ok"></i>
        </p>
    </div>
</article>
```

22. This code extracts and analyzes book products from the **Books to Scrape** website.
23. It selects all elements with the class `.product_pod`, representing book listings. The first product is stored in `example`, and its attributes are examined using `.attrs`.
24. Finally, `list(example.children)` retrieves all direct child elements within the book's HTML structure, helping understand its internal layout for further data extraction.

```
[132]: products = soup.select(".product_pod")
[134]: example = products[0]
[136]: type(example)
[136]: bs4.element.Tag
[138]: example.attrs
[138]: {'class': ['product_pod']}
[140]: list(example.children)

[140]: ['\n',
    <div class="image_container">
        <a href="a-light-in-the-attic_1000/index.html"></a>
    </div>,
    '\n',
    <p class="star-rating Three">
        <i class="icon-star"></i>
        <i class="icon-star"></i>
        <i class="icon-star"></i>
        <i class="icon-star"></i>
        <i class="icon-star"></i>
    </p>,
    '\n',
    <h3><a href="a-light-in-the-attic_1000/index.html" title="A Light in the Attic">A Light in the ...</a></h3>,
    '\n',
    <div class="product_price">
        <p class="price_color">£51.77</p>
    </div>
]
```

25. The code analyzes the **Books to Scrape** website's HTML structure to extract information about book ratings and titles. It searches for books with a **three-star or two-star rating** using CSS class selectors. Then, it finds all `<a>` (anchor) tags within a product listing, selects the second anchor tag, and extracts the **book title** from its title attribute, providing insight into the book's details.
26. In the end, we extracted the book titles that have a two-star rating from the book to scrape website.

```
[142]: example.select('.star-rating.Three')
[142]: [<p class="star-rating Three">
  <i class="icon-star"></i>
  <i class="icon-star"></i>
  <i class="icon-star"></i>
  <i class="icon-star"></i>
  <i class="icon-star"></i>
</p>]

[144]: example.select('.star-rating.Two')
[144]: []

[146]: example.select('a')
[146]: [<a href="a-light-in-the-attic_1000/index.html"></a>,
 <a href="a-light-in-the-attic_1000/index.html" title="A Light in the Attic">A Light in the ...</a>]

[148]: example.select('a')[1]
[148]: <a href="a-light-in-the-attic_1000/index.html" title="A Light in the Attic">A Light in the ...</a>

[150]: example.select('a')[1]['title']
[150]: 'A Light in the Attic'

[154]: two_star_titles
```

```
[154]: ['Starving Hearts (Triangular Trade Trilogy, #1)',
 'Libertarianism for Beginners',
 'It's Only the Himalayas',
 'How Music Works',
 'Maude (1883-1993):She Grew Up with the country',
 "You can't bury them all: Poems",
 'Reasons to Stay Alive',
 'Without Borders (Wanderlove #1)',
 'Soul Reader',
 'Security',
 'Saga, Volume 5 (Saga (Collected Editions) #5)',
 'Reskilling America: Learning to Labor in the Twenty-First Century',
 'Political Suicide: Missteps, Peccadilloes, Bad Calls, Backroom Hijinx, Sordid Past, Rotten Breaks, and Just Plain Dumb Mistakes in the Annals of American Politics',
 'Obsidian (Lux #1)',
 'My Paris Kitchen: Recipes and Stories',
 'Masks and Shadows',
 'Lumberjanes Vol 2: Friendship to the Max (Lumberjanes #5-8)'
```