





Web Scraping Guide

What is Web Scraping?

Web scraping is a technique for extracting data from websites in an automated, efficient manner. This process transforms unstructured web data into a structured format, such as databases or spreadsheets, making it easier to analyze and use for various purposes, whether in B2B or B2C contexts.

Uses of Web Scraping

- Industry Statistics and Insights: Gathering data for market analysis.
- Price Comparison: Comparing prices across different online stores.
- Financial Analysis: Analyzing stock prices or cash flows.
- Academic Research: Collecting data for research purposes.

→ In our project, web scraping is used for research purposes to retrieve data from the Food and Drug Administration (FDA) website about medical devices' approval (accessdata.fda.gov/scripts/cdrh/cfdocs/cfPMN/pmn.cfm) and product recalls (accessdata.fda.gov/scripts/cdrh/cfdocs/cfRES/res.cfm).

Limitations of Web Scraping

- CAPTCHA: Used to prevent automated access.
- IP Blocking: Websites may block IP addresses suspected of being bots.
- Login Areas: Restricted access to certain sections.
- robots.txt: Files that instruct crawlers on which parts of the site can or cannot be accessed.

 \rightarrow In our project, by applying the web scraping technique we did not encounter these limits.

Legal Considerations

Web scraping is generally legal when used to obtain public information. It becomes illegal if it violates copyright laws, breaks terms of service, or is used for malicious purposes like phishing or identity theft.

 \rightarrow In our project, by applying the web scraping technique to public information our approach is legal and not subject to legal considerations.

Types of Web Scrapers

- 1. Browser Extensions: Tools that run within your web browser.
- 2. Downloadable Software: Programs installed on your computer.

3. <u>Self-Built Scrapers:</u> Custom-built tools using programming languages like Python, R, or PHP.









 \rightarrow In our project, we have adapted a browser extension called "Webscraper.io" (<u>https://webscraper.io</u>). This is a no-code tool that is easy to use because it does not require programming experience or knowledge and the site-mapping implementation is user-friendly.

How Web Scraping Works

Web scraping tools work by iterating through every URL on a page coded in HTML. These tools extract the required data assuming the structure of the website remains consistent.

Making the Web Scraping using Webscraper.io

1. Setting Up:

- Create a sitemap pointing to the starting domain (either a local host or a web existing domain)

 \rightarrow This is the case with a virtualization of the localhost:

Sitemaps	Sitemap 🔻	Create new sitemap 🔻
Sitemap name	product	_recalled
Start URL 1	http://dr	ugrecall.us/
	Create	Sitemap

 \rightarrow This is the case with an existing website:

Sitemap name	product_notification2
Start URL 1	https://www.accessdata.fda.gov/scripts/cdrh/cfdocs/cfpmn/pmn.cfm?start_search
	Save Sitemap

- Define selectors to target specific data fields in the HTML code

 \rightarrow This is an example of how to build a selector. Selectors can be of several types, we used principally 'text' and 'link' data type:

Finanziato dall'Unione europea NextGenerationEU	a Ministero dell'Università e della Ricerca			INSTITUTE OF ECONOMICS Sant'Anna Date of Advenuel Thebas - Pha
Sitemaps Site	emap product_recall - Crea	ate new sitemap 👻		* LIUC Università Cattaneo
Id	recall_link			-
Туре	Link	Data and iour	Machierral	
Selector	Select Element preview	Data preview	.MsoNormal a	
Link type	Link (read from href attribute	e)		-
Parent Selectors	_root recall_link			-
	Save selector Cancel			

\rightarrow This is a sample of selectors for a sitemap:

ID	Selector	Actions			
device_classification_name	td td td tr:contains('Device Classification Name') a	Element preview	Data preview	Edit	Delete
510k_number	td td td tr:contains('510(k) Number') td	Element preview	Data preview	Edit	Delete
device_name	td td td tr:contains('Device Name') td	Element preview	Data preview	Edit	Delete
applicant	td td tr[valign]:contains('Applicant') > td	Element preview	Data preview	Edit	Delete
applicant_contact	td td td tr:contains('Applicant Contact') td	Element preview	Data preview	Edit	Delete
correspondent	td td tr[valign]:contains('Correspondent') > td	Element preview	Data preview	Edit	Delete
correspondent_contact	td td td tr:contains('Correspondent Contact') td	Element preview	Data preview	Edit	Delete
regulation_number	td td td tr:contains('Regulation Number') a	Element preview	Data preview	Edit	Delete
classification_product_code	td td td td td a	Element preview	Data preview	Edit	Delete

2. In the case of a local webpage, you have to change 'localhost' to a Domain Name:

- Redirect the website name on your computer (hosts file in C:\Windows\System32\drivers\etc).

127.0.0.1	drugrecall.us
127.0.0.1	www.drugrecall.us

- Redirect Apache to the website address in the installation directory (httpd-vhosts.conf).



3. Start the web scraper:

- After creating Sitemap and Selectors, the scraping process can be initialized

	Sitemaps	Sitemap product_notification	Create n	ew sitemap 🔻	
_root		Selectors Selector graph			
	ID	Edit metadata	type	Multiple	Parent selectors
notify_link	Scrape	SelectorL	ink yes	_root	
Add new selector		Browse			
		Export Sitemap			
		Export data			

Once scraping is complete, data can be exported from the tool to an Excel spreadsheet or a .csv file.

Final Tips for Effective Web Scraping

- Know What to Search For: Identify the specific data fields you need.

 \rightarrow Focusing on the data of interest reduces the time needed for web scraping.

- Analyze Page Source: Understand the HTML structure (elements like tr/td/a).

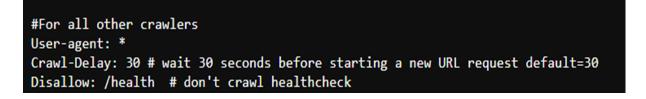
 \rightarrow In our process of web scraping, irrelevant elements have been added to the useful data because these are part of the HTML structure of the page (see example displayed in the figure below). To clean the data, in a second step, these elements were deleted.

Finanziato dall'Unione europea NextGenerationEU	Ministero dell'Università e della Ricerca	Italiadomani	SITATIS SECXXX ⁺	STILLE ECONOMICS Sant'Anna Event of Manual Bade Par
				* LIUC Università Cattaneo
Recall Status ¹	Open ³ , Classifie	d		
Recall Number	Z-2137-2024			
Recall Event ID	94622			
510(K)Number	<u>K233448</u>			
s match Chrome's language	Switch DevTools to Italian	Don't show again		
Network Performance	Memory Application S	ecurity Lighthouse	Recorder	Performance i
Create new sitemap 👻				

posted	recall_status	recall_number	recall_event_id	510k_number	product_classification	product
12,	Open <mark>3.</mark> Classified	Z-2137-2024	94622	K233448 <mark>24</mark>	Shunt, central nervous system and components	JXG

- Respect Crawl Rate: Follow guidelines in robots.txt to avoid being blocked.

 \rightarrow The FDA website does not allow the switch to another web page within 30 seconds. FDA's robots.txt is placed at the root domain (<u>https://www.fda.gov/robots.txt</u>):



- Sitemap Portability: Export and import sitemaps between different computers using JSON.

 \rightarrow Direct export is only possible for existing and available websites. In the case of local sites, changes in the starting URL are required.