

Web Scraping 101

Module 1 - No-code Tools

Mar 2025



Our Goals Today

After this workshop, you will be able to

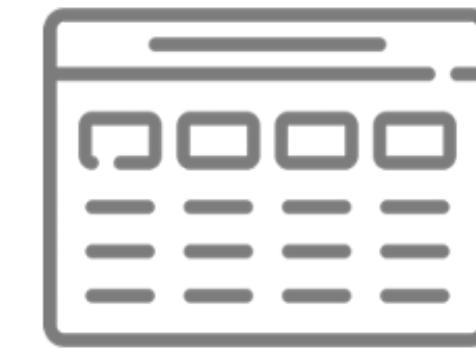
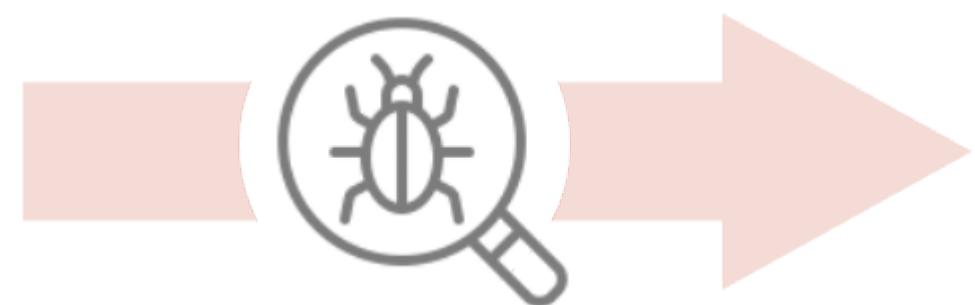
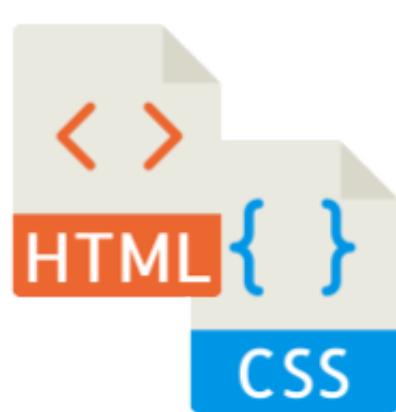
- Understand the basic concepts and principles of web scraping
- Identify available web scraping tools and their appropriate uses
- Use **Power Query** and **Web Scraper** to extract content of interest from the web
- Make wise decisions when collecting data from the web

Web Scraping Basics



What is Web Scraping

- > the process of **extracting information** from the web
- > **automate, rerun** the process



Websites

Identify your target website

Scraping

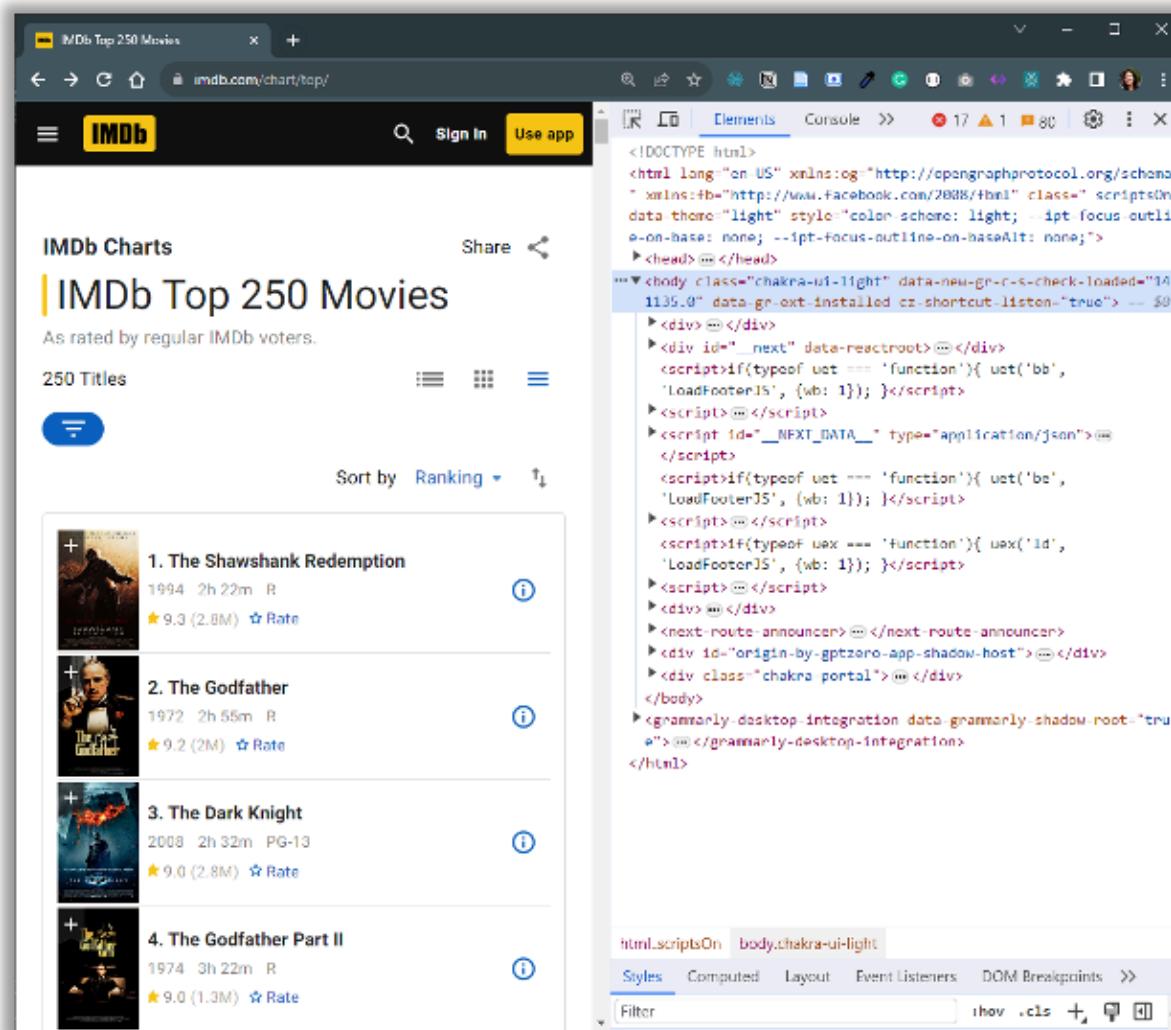
Extract info from the web

Structured Data

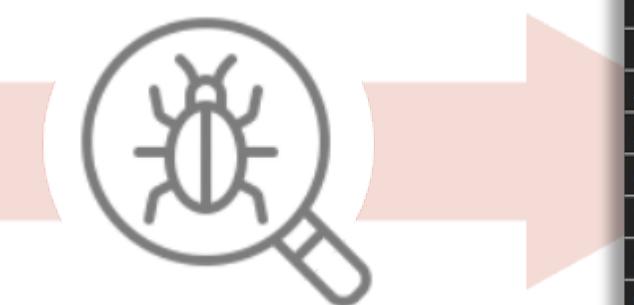
Store it in a structured format / database

What is Web Scraping

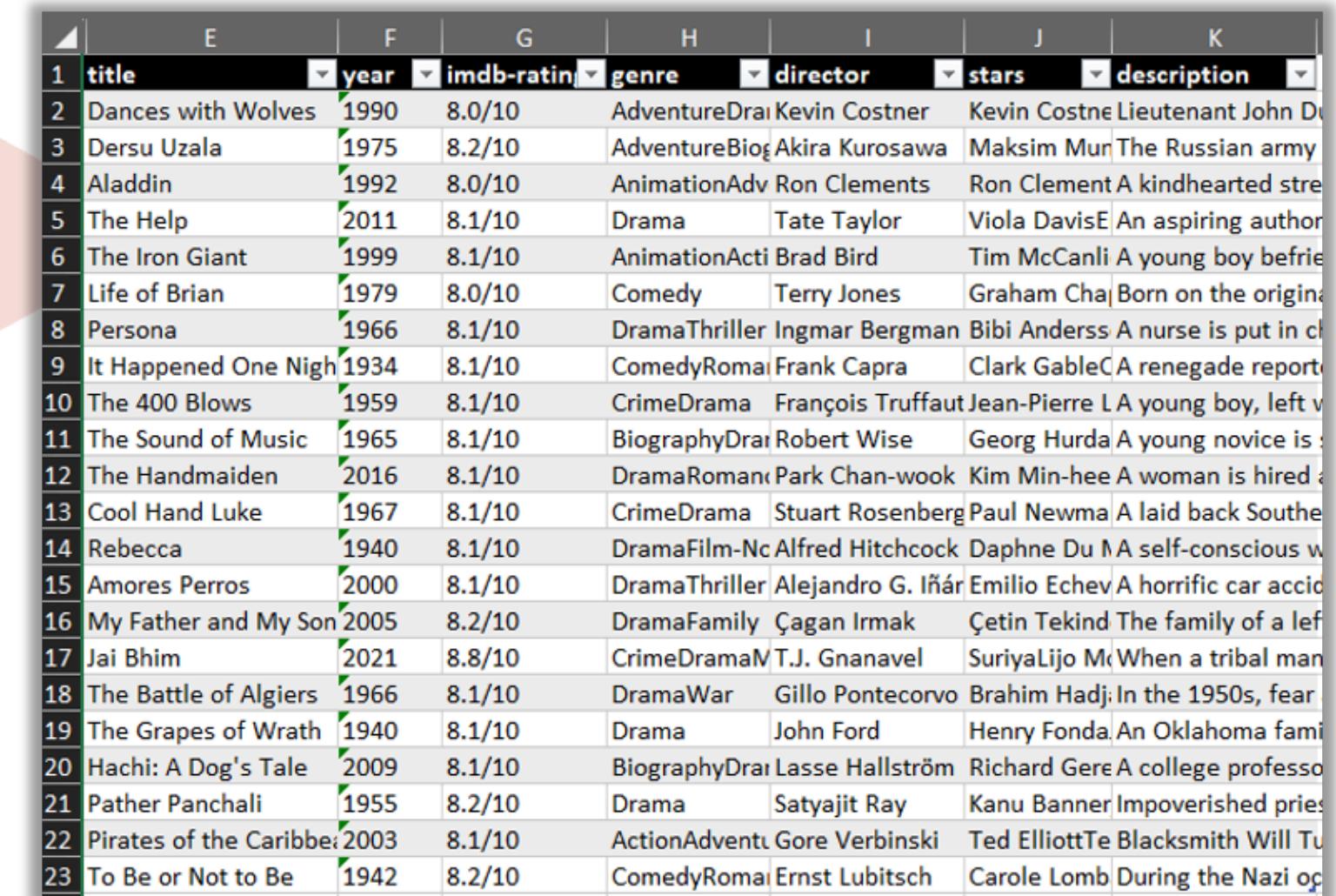
- > the process of **extracting information** from the web
- > **automate, rerun** the process



The screenshot shows the IMDb Top 250 Movies page. It displays a list of movies with their titles, years, and ratings. A magnifying glass icon is overlaid on the page, indicating the process of extracting data.



Scraping
Extract info
from the web



The screenshot shows an Excel spreadsheet with a table of movie data. The columns are labeled: title, year, imdb-rating, genre, director, stars, and description. The data includes various movies with their details such as title, year of release, rating on IMDb, genre, director, lead actors, and a brief description.

1	E	F	G	H	I	J	K	
1	title	year	imdb-rating	genre	director	stars	description	
2	Dances with Wolves	1990	8.0/10	Adventure	Draí Kevin Costner	Kevin Costner	Lieutenant John D	
3	Dersu Uzala	1975	8.2/10	Adventure	Bio	Akira Kurosawa	Maksim Mur	The Russian army
4	Aladdin	1992	8.0/10	Animation	Adv	Ron Clements	Ron Clements	A kindhearted stre
5	The Help	2011	8.1/10	Drama	Tate Taylor	Viola Davis	E An aspiring autho	
6	The Iron Giant	1999	8.1/10	Animation	Acti	Brad Bird	Tim McCanli	A young boy befrie
7	Life of Brian	1979	8.0/10	Comedy	Terry Jones	Graham Cha	Born on the origina	
8	Persona	1966	8.1/10	Drama	Thriller	Ingmar Bergman	Bibi Andersss	A nurse is put in ch
9	It Happened One Nigh	1934	8.1/10	Comedy	Roma	Frank Capra	Clark Gable	CA renegade report
10	The 400 Blows	1959	8.1/10	Crime	Drama	François Truffaut	Jean-Pierre L	A young boy, left w
11	The Sound of Music	1965	8.1/10	Biography	Dra	Robert Wise	Georg Hurda	A young novice is s
12	The Handmaiden	2016	8.1/10	Drama	Roman	Park Chan-wook	Kim Min-hee	A woman is hired a
13	Cool Hand Luke	1967	8.1/10	Crime	Drama	Stuart Rosenberg	Paul Newma	A laid back Souther
14	Rebecca	1940	8.1/10	Drama	Film-Nc	Alfred Hitchcock	Daphne Du M	A self-conscious w
15	Amores Perros	2000	8.1/10	Drama	Thriller	Alejandro G. Iñárr	Emilio Echev	A horrific car accide
16	My Father and My Son	2005	8.2/10	Drama	Family	Çagan Irmak	Çetin Tekind	The family of a lef
17	Jai Bhim	2021	8.8/10	Crime	Drama	M T.J. Gnanavel	Suriya	When a tribal man
18	The Battle of Algiers	1966	8.1/10	Drama	War	Gillo Pontecorvo	Brahim Hadj	In the 1950s, fear
19	The Grapes of Wrath	1940	8.1/10	Drama		John Ford	Henry Fonda	An Oklahoma fami
20	Hachi: A Dog's Tale	2009	8.1/10	Biography	Dra	Lasse Hallström	Richard Gere	A college profess
21	Pather Panchali	1955	8.2/10	Drama		Satyajit Ray	Kanu Banner	Impoverished pries
22	Pirates of the Caribbean: Dead Man's Chest	2003	8.1/10	Action	Adventu	Gore Verbinski	Ted Elliott	Te Blacksmith Will Tu
23	To Be or Not to Be	1942	8.2/10	Comedy	Roma	Ernst Lubitsch	Carole Lomb	During the Nazi op

Things to consider before you scrape

Adapted from: DeVito, N. J., Richards, G. C., & Inglesby, P. (2020). How we learnt to stop worrying and love web scraping. *Nature*, 585(7826), 621–622.
<https://doi.org/10.1038/d41586-020-02558-0>

1. Can I get the data in an easier way?

Is the data **already available** somewhere?

Things to consider before you scrape

Adapted from: DeVito, N. J., Richards, G. C., & Inglesby, P. (2020). How we learnt to stop worrying and love web scraping. *Nature*, 585(7826), 621–622. <https://doi.org/10.1038/d41586-020-02558-0>

1. Can I get the data in an easier way?

Is the data **already available** somewhere?

Find raw data

Google
Dataset Search



kaggle
(data + code)



Hugging Face
(data + code)

Library databases, e.g.

Passport
(consumer data)

CEIC
(economy data)

Find research data

zenodo

figshare

DataSpace@HKUST

(HKUST data repository)

re3data.org
REGISTRY OF RESEARCH DATA REPOSITORIES

(find data repository by subject)

Things to consider before you scrape

Adapted from: DeVito, N. J., Richards, G. C., & Inglesby, P. (2020). How we learnt to stop worrying and love web scraping. *Nature*, 585(7826), 621–622.
<https://doi.org/10.1038/d41586-020-02558-0>

2. Can this website be scraped, technically?

- Any restrictions, e.g. CAPTCHAS (“I’m not a robot” test)?
- Any rules to follow, e.g. no. of requests per min?

Things to consider before you scrape

Adapted from: DeVito, N. J., Richards, G. C., & Inglesby, P. (2020). How we learnt to stop worrying and love web scraping. *Nature*, 585(7826), 621–622.
<https://doi.org/10.1038/d41586-020-02558-0>

2. Can this website be scraped, technically?

- Any restrictions, e.g. CAPTCHAS (“I’m not a robot” test)?
- Any rules to follow, e.g. no. of requests per min?

3. Can this website be scraped, legally?

- Any restrictions on copyright, or privacy issues in scraping or sharing?

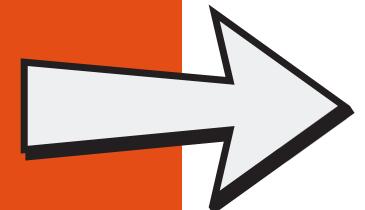
Things to consider before you scrape

2. Can this website be scraped, technically?

- Any restrictions, e.g. CAPTCHAS (“I’m not a robot” test)?
- Any rules to follow, e.g. no. of requests per min?

3. Can this website be scraped, legally?

- Any restrictions on copyright, or privacy issues in scraping or sharing?



Check “Terms of Use”, “robots.txt”, any official API available, copyright / data sharing policies.

Adapted from: DeVito, N. J., Richards, G. C., & Inglesby, P. (2020). How we learnt to stop worrying and love web scraping. *Nature*, 585(7826), 621–622.
<https://doi.org/10.1038/d41586-020-02558-0>

Things to consider before you scrape

Adapted from: DeVito, N. J., Richards, G. C., & Inglesby, P. (2020). How we learnt to stop worrying and love web scraping. *Nature*, 585(7826), 621–622. <https://doi.org/10.1038/d41586-020-02558-0>

Example

<https://x.com/en/tos>

X Terms of Service

Download PDF

- You must abide by the Services' acceptable use terms:**
You may not access the Services in any way other than through the currently available, published interfaces that we provide. For example, this means that you cannot scrape the Services without X's express written permission, try to work around any technical limitations we impose, or otherwise attempt to disrupt the operation of the Services.

to give instructions to web crawlers or bots about which parts of a website they are allowed or disallowed to access

<https://x.com/robots.txt>

User-Agent: FacebookBot
Disallow: *

User-agent: facebookexternalhit
Disallow: *

User-agent: Discordbot
Disallow: *

User-agent: Bingbot
Disallow: *

[How to read this](#)

More examples:

<https://arxiv.org/robots.txt> (allow, with crawl delay)

<https://www.amazon.com/robots.txt> (block GPTbot)

<https://stackoverflow.com/robots.txt> (block all bots)

Web Scraping Tools



Not all scraping tasks need programming!

Pick the tool that fits your needs.



No-code tools

(Browser plug-in, e.g. WebScraper)

- “Point-and-click”
- Static and simple websites
- Small scale (a few pages)



Programming

(e.g. Python selenium, scrapy)

- Dynamic websites
- Large scale projects
- Integrate with databases, build complex workflows

Not all scraping tasks need programming!

Pick the tool that fits your needs.



No-code tools



Power Query



Web
Scraper



Programming

(e.g. Python selenium, scrapy)

- Dynamic websites
- Large scale projects
- Integrate with databases, build complex workflows



Power Query

A built-in Excel tool, for connecting and manipulating data from various sources, including web sources.

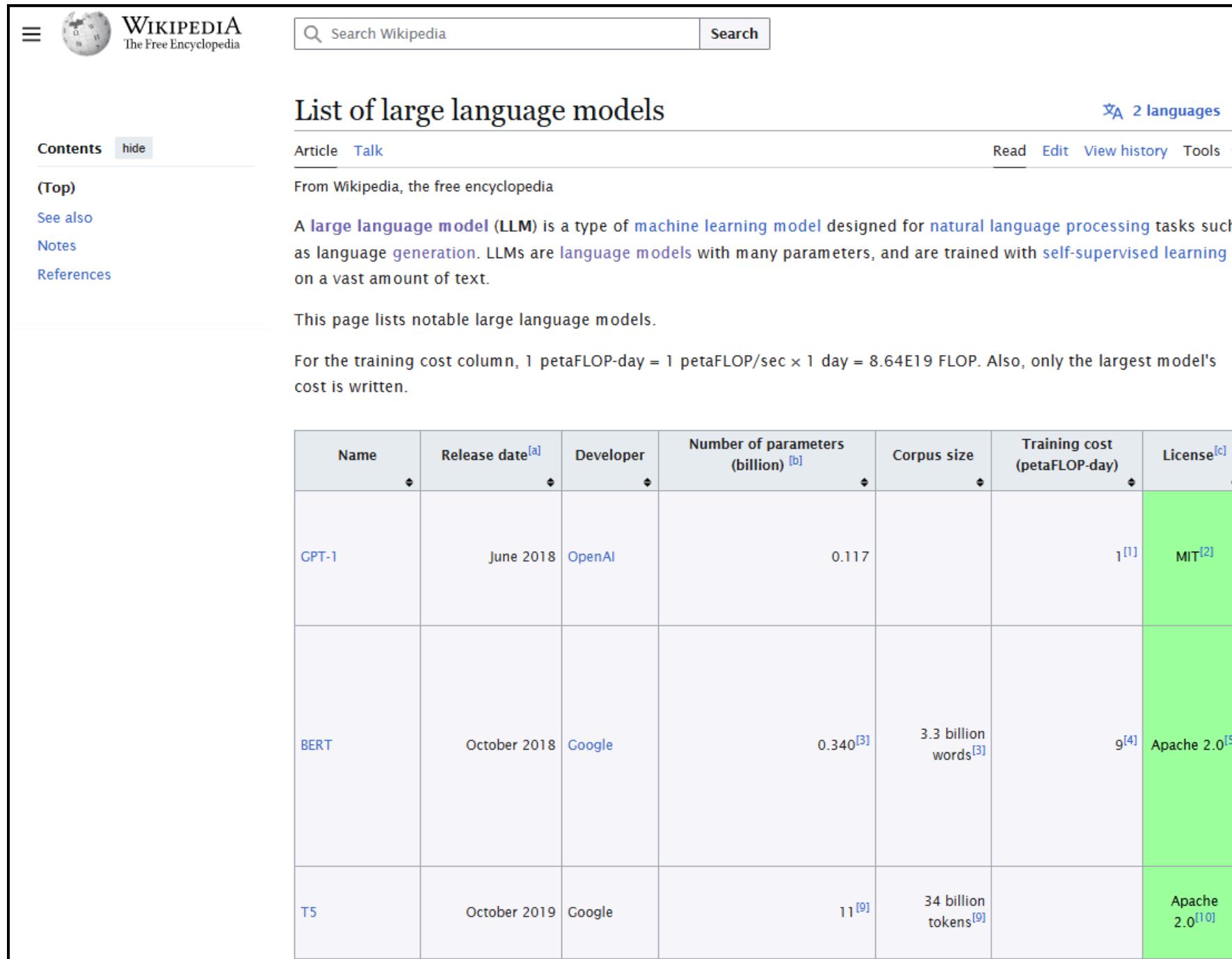
1. Get the page URL
2. Select the info you want
3. Extract the data
4. Save in Excel (for further data manipulation)

Note: Some features (e.g. get data from Web) **are NOT supported in Mac** ([learn more](#)).
So, use Windows PCs to run Power Query.

Demo

Extract info from a Wikipedia page

https://en.wikipedia.org/wiki/List_of_large_language_models

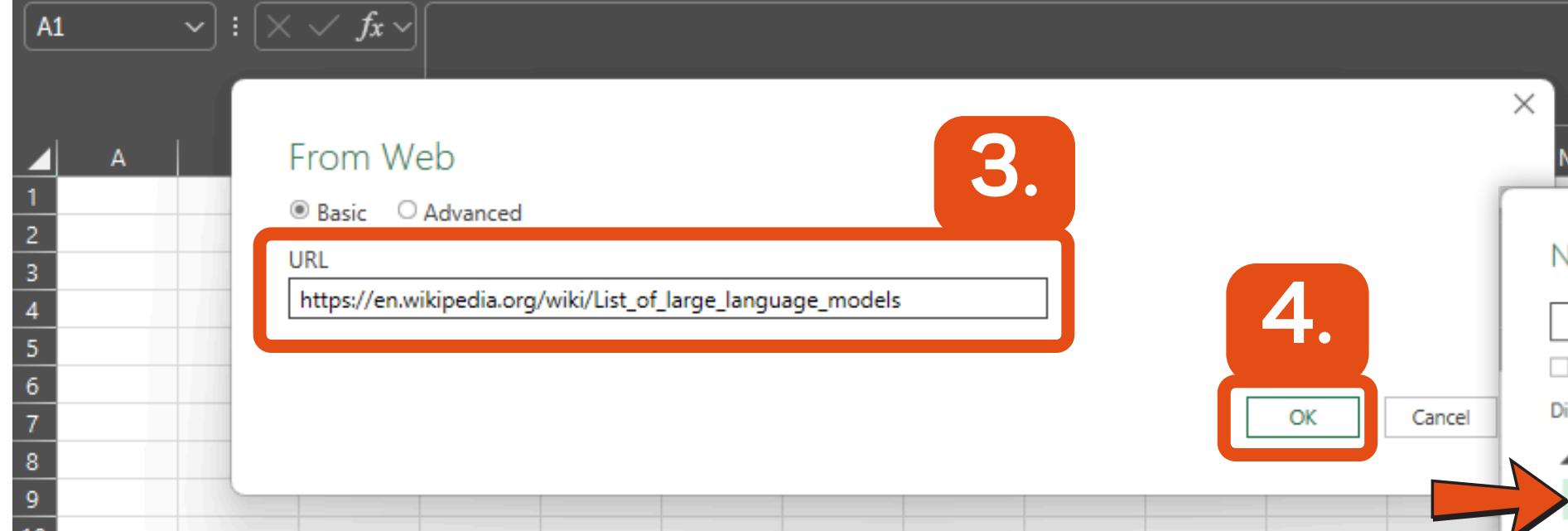
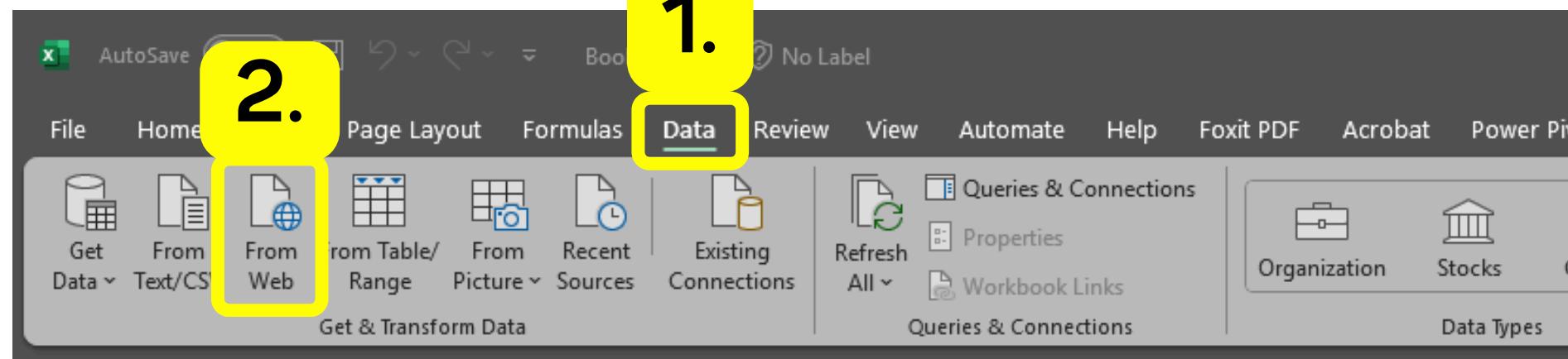


The screenshot shows a Wikipedia article titled "List of large language models". The page includes a sidebar with navigation links like Contents, Article, Talk, Read, Edit, View history, Tools, and a search bar at the top. The main content area contains a brief introduction about Large Language Models (LLMs) and their training. Below this is a table listing three LLMs: GPT-1, BERT, and T5. The table has columns for Name, Release date, Developer, Number of parameters (billion), Corpus size, Training cost (petaFLOP-day), and License.

Name	Release date ^[a]	Developer	Number of parameters (billion) ^[b]	Corpus size	Training cost (petaFLOP-day)	License ^[c]
GPT-1	June 2018	OpenAI	0.117		1 ^[1]	MIT ^[2]
BERT	October 2018	Google	0.340 ^[3]	3.3 billion words ^[3]	9 ^[4]	Apache 2.0 ^[5]
T5	October 2019	Google	11 ^[6]	34 billion tokens ^[9]		Apache 2.0 ^[10]

What we want:

- A cleaned list of Large Language Models
- Saved in Excel



- Data > From Web > Type in URL > OK
- Select target table > Transform Data

Column1	Column2	Column3
Name	Release date[a]	Developer
GPT-1	June 2018	OpenAI
BERT	October 2018	Google
T5	October 2019	Google
XLNet	June 2019	Google
GPT-2	February 2019	OpenAI
GPT-3	May 2020	OpenAI
GPT-Neo	March 2021	EleutherAI
GPT-J	June 2021	EleutherAI
Megatron-Turing NLG	October 2021 [28]	Microsoft and Nvidia
Ernie 3.0 Titan	December 2021	Baidu
Claude[32]	December 2021	Anthropic
GLaM (Generalist Language Model)	December 2021	Google
Gopher	December 2021	DeepMind
LaMDA (Language Models for Dialog Applications)	January 2022	

Target webpage:

https://en.wikipedia.org/wiki/List_of_large_language_models

Transform Data



Hands-on (1)

Get the list of IMDb Top 250 Movies

https://www.imdb.com/chart/top/?ref_=nv_mv_250

The screenshot shows the IMDb Top 250 Movies chart. At the top, it says "IMDb Charts" and "IMDb Top 250 Movies". Below that, it says "As rated by regular IMDb voters." and "250 Titles". There are three movie cards displayed: 1. The Shawshank Redemption (1994, 2h 22m, R, 9.3 (3M) Rate), 2. The Godfather (1972, 2h 55m, R, 9.2 (2.1M) Rate), and 3. The Dark Knight (2008, 2h 32m, PG-13, 9.0 (2.9M) Rate). An orange arrow points from this screenshot to the Power Query table below.

A	B	C
Rank	Title	Rating
1	1 The Shawshank Redemption	9.3
2	2 The Godfather	9.2
3	3 The Dark Knight	9
4	4 The Godfather Part II	9
5	5 12 Angry Men	9
6	6 The Lord of the Rings: The Return of the King	9
7	7 Schindler's List	9
8	8 Pulp Fiction	8.9
9	9 The Lord of the Rings: The Fellowship of the Ring	8.9
10	10 The Good, the Bad and the Ugly	8.8
11	11 Forrest Gump	8.8
12	12 The Lord of the Rings: The Two Towers	8.8
13	13 Fight Club	8.8
14	14 Inception	8.8
15	15 Star Wars: Episode V - The Empire Strikes Back	8.7
16	16 The Matrix	8.7
17	17 Goodfellas	8.7
18	18 One Flew Over the Cuckoo's Nest	8.7
19	19 Interstellar	8.7
20	20 Se7en	8.6
21	21 It's a Wonderful Life	8.6
22	22 Seven Samurai	8.6
23	23 The Silence of the Lambs	8.6
24	24 Saving Private Ryan	8.6
25	25 City of God	8.6

What we want:

List of movies

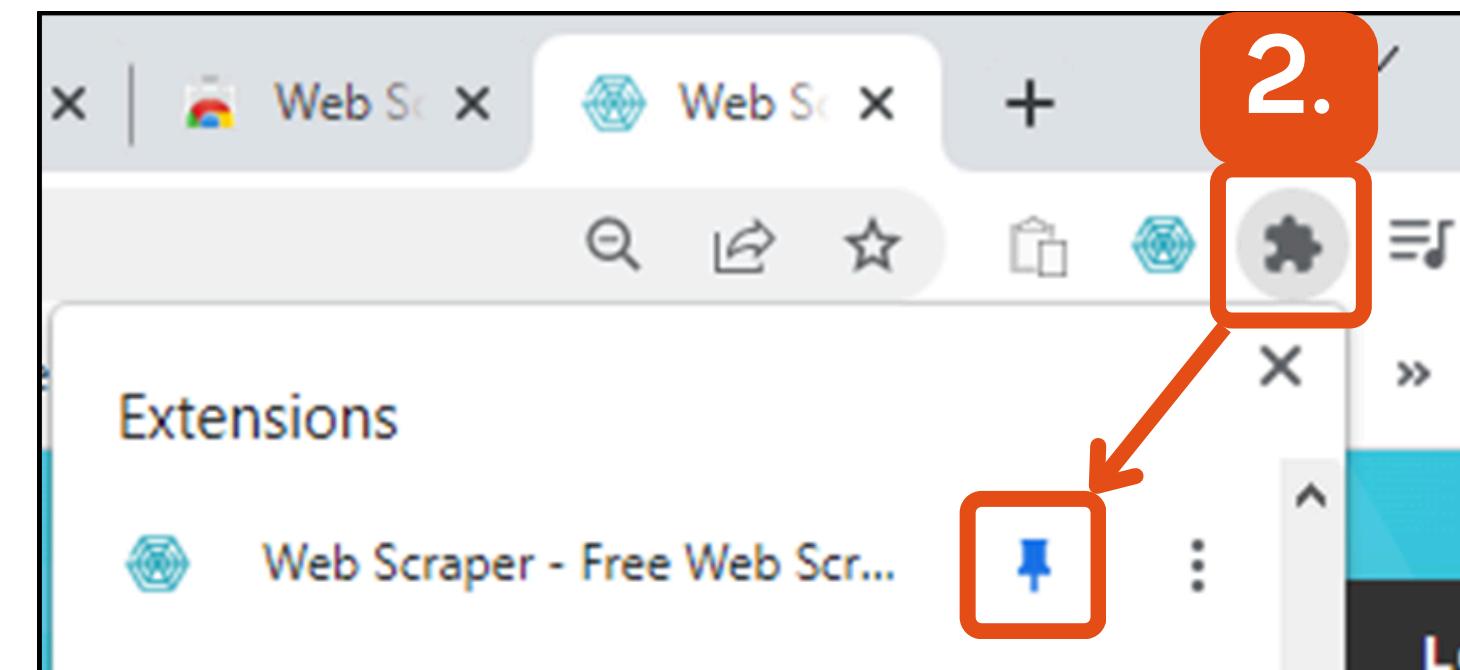
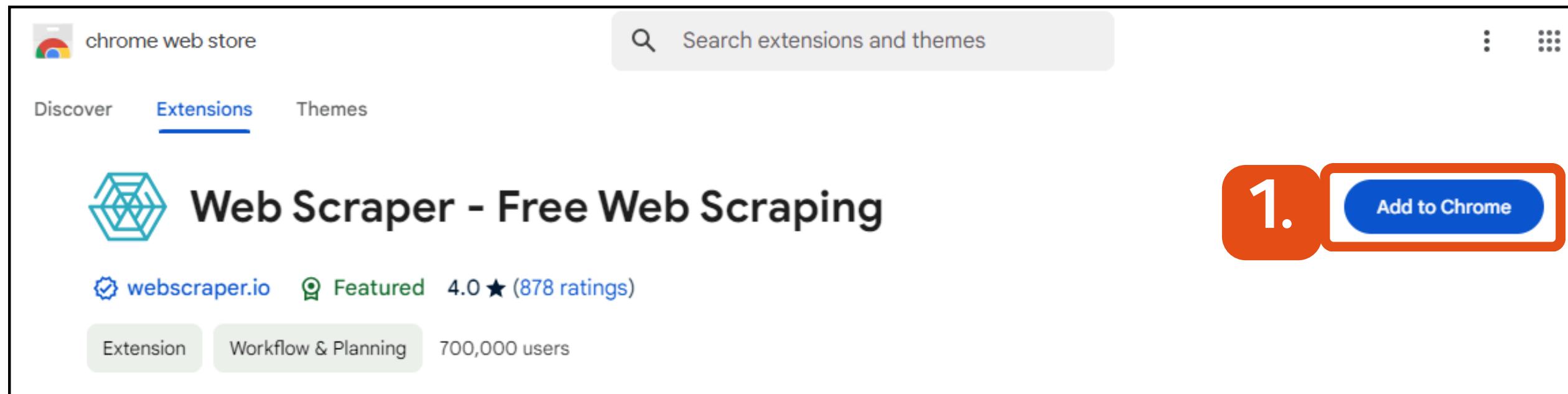
- Rank
- Title
- Rating



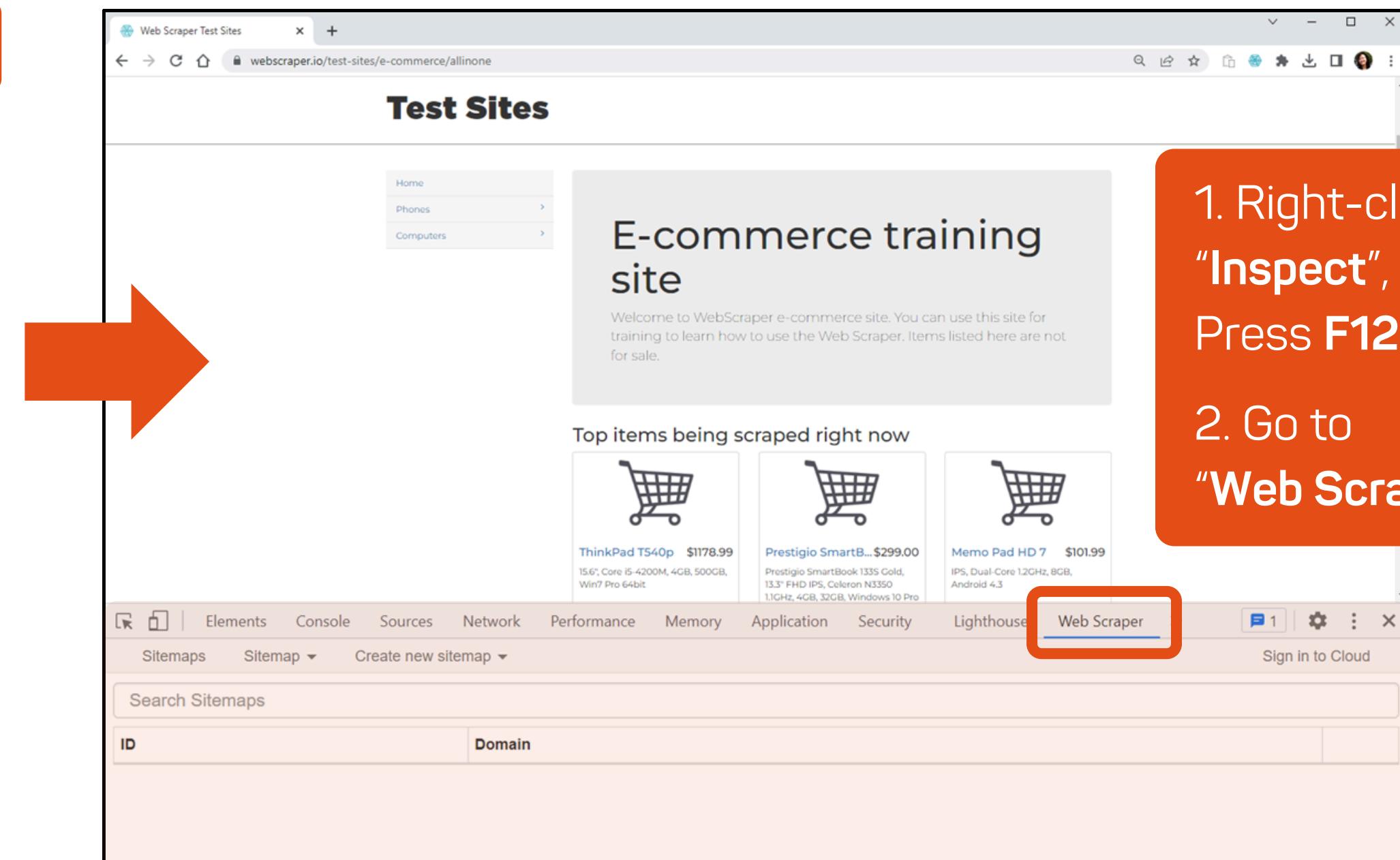
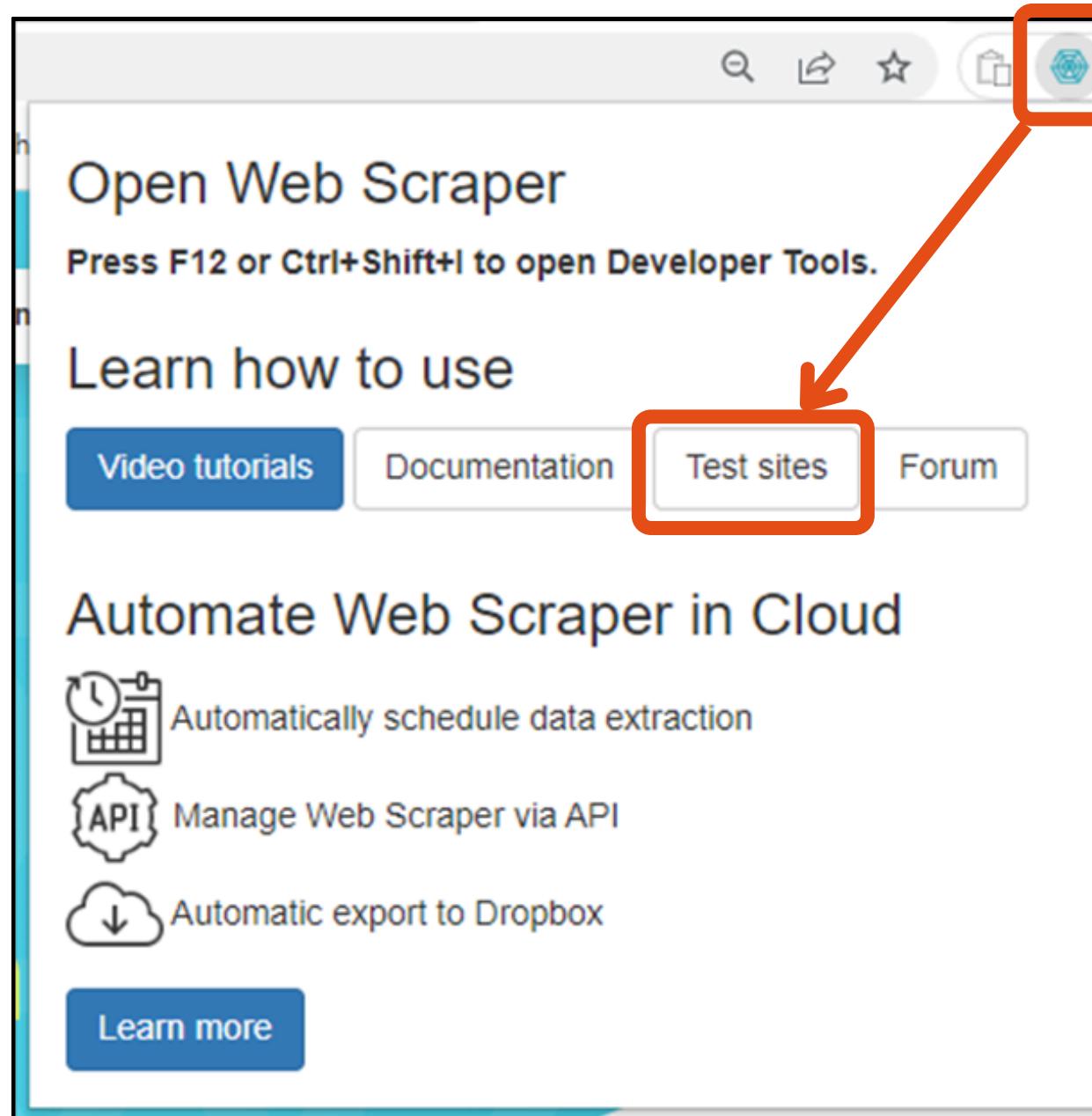
A browser plug-in for web scraping.

- 1.Understand basic structure of the page
- 2.Point and click the “field”
- 3.Extract the data
- 4.Download in .csv or .xlsx files

Install Web Scraper plug-in



Launch Web Scraper



1. Right-click and "Inspect", or Press **F12**
2. Go to "Web Scraper"

Some HTML Basics

HTML - Structure of a webpage

Webpage Preview

```
1 <html>
2
3 <head>
4   <title>Title of the page</title>
5 </head>
6
7 <body>
8   <div>
9     <h1>Header</h1>
10    <p>This is a paragraph.</p>
11    <p>This is another paragraph.</p>
12  </div>
13
14  <div>
15    <p>Here is a link:</p>
16    <a href="https://library.hkust.edu.hk/">Library homepage</a>
17  </div>
18 </body>
19
20 </html>
```

Tags

<body>...</body>
<div>...</div>
<h1>...</h1>
<p>...</p>

Header

This is a paragraph.

This is another paragraph.

Here is a link:

[Library homepage](https://library.hkust.edu.hk/)

Hierarchy

Like family tree, e.g.

<div> Parent

<p> Child

<https://codepen.io/asterzhao/pen/oNaXYKw?editors=1100>

Some HTML Basics

The screenshot shows a code editor with two tabs: 'HTML' and 'CSS'. The 'HTML' tab contains the following code:

```
7 <body>
8   <div class="some-class">
9     <h1>Header</h1>
10    <p>This is a paragraph.</p>
11    <p id="unique-id">This is another paragraph.</p>
12  </div>
13
14  <div>
15    <p>Here is a link:</p>
16    <a href="https://library.hkust.edu.hk/">Library homepage</a>
17  </div>
18 </body>
```

The 'CSS' tab contains the following code:

```
1 .some-class {
2   color: green
3 }
4
5 #unique-id {
6   color: red
7 }
```

Annotations highlight specific elements:

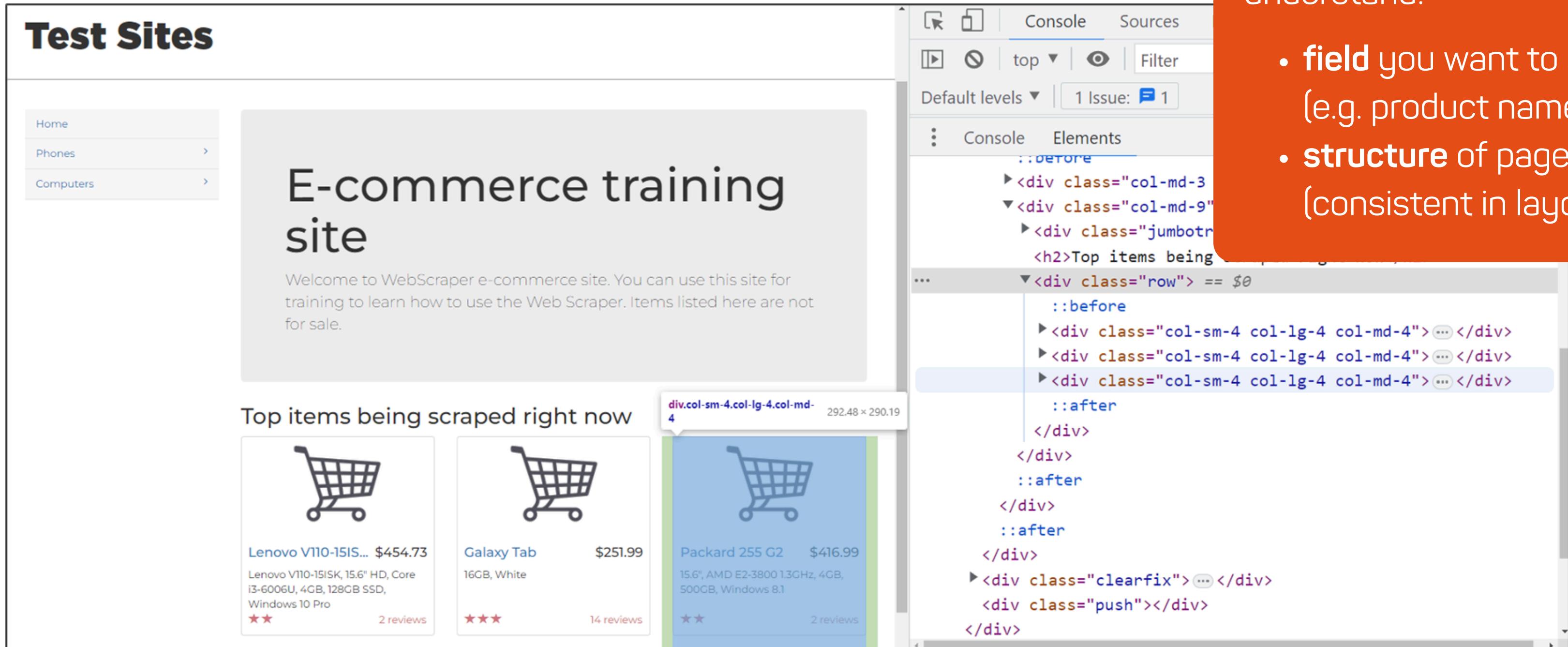
- A callout box labeled 'Attributes' lists: <class = "...>, <id = "...>, <href = "...>. It also includes the text 'Useful for scrapping!'.
- A callout box labeled 'HTML - Structure of a webpage' points to the 'Header' element (<h1>Header</h1>).
- A callout box labeled 'CSS - Style of a webpage' points to the class definition (.some-class { color: green; }) and the ID definition (#unique-id { color: red; }).
- A callout box labeled 'Webpage Preview' points to the 'Header' element in the HTML structure.

<https://codepen.io/asterzhao/pen/oNaXYKw?editors=1100>

Demo

WebScraper - Test site

<https://webscraper.io/test-sites/e-commerce/allinone>



Test Sites

Home

Phones >

Computers >

E-commerce training site

Welcome to WebScraper e-commerce site. You can use this site for training to learn how to use the Web Scraper. Items listed here are not for sale.

Top items being scraped right now

Item	Price	Description	Reviews
Lenovo V110-15ISK	\$454.73	Lenovo V110-15ISK, 15.6" HD, Core i3-6006U, 4GB, 128GB SSD, Windows 10 Pro	2 reviews
Galaxy Tab	\$251.99	16GB, White	14 reviews
Packard 255 G2	\$416.99	15.6", AMD E2-3800 1.3GHz, 4GB, 500GB, Windows 8.1	2 reviews

"Inspect" the page helps you understand:

- **field** you want to scrape (e.g. product name, price)
- **structure** of pages (consistent in layout?)

```

<div>
  <div> Home </div>
  <div> Phones </div>
  <div> Computers </div>
</div>

<div>
  <h1>E-commerce training site</h1>
  <p>Welcome to WebScraper e-commerce site. You can use this site for training to learn how to use the Web Scraper. Items listed here are not for sale.</p>
  <h2>Top items being scraped right now</h2>
  <table border="1">
    <thead>
      <tr>
        <th>Item</th>
        <th>Price</th>
        <th>Description</th>
        <th>Reviews</th>
      </tr>
    </thead>
    <tbody>
      <tr>
        <td>Lenovo V110-15ISK</td>
        <td>$454.73</td>
        <td>Lenovo V110-15ISK, 15.6" HD, Core i3-6006U, 4GB, 128GB SSD, Windows 10 Pro</td>
        <td>2 reviews</td>
      </tr>
      <tr>
        <td>Galaxy Tab</td>
        <td>$251.99</td>
        <td>16GB, White</td>
        <td>14 reviews</td>
      </tr>
      <tr>
        <td>Packard 255 G2</td>
        <td>$416.99</td>
        <td>15.6", AMD E2-3800 1.3GHz, 4GB, 500GB, Windows 8.1</td>
        <td>2 reviews</td>
      </tr>
    </tbody>
  </table>
</div>

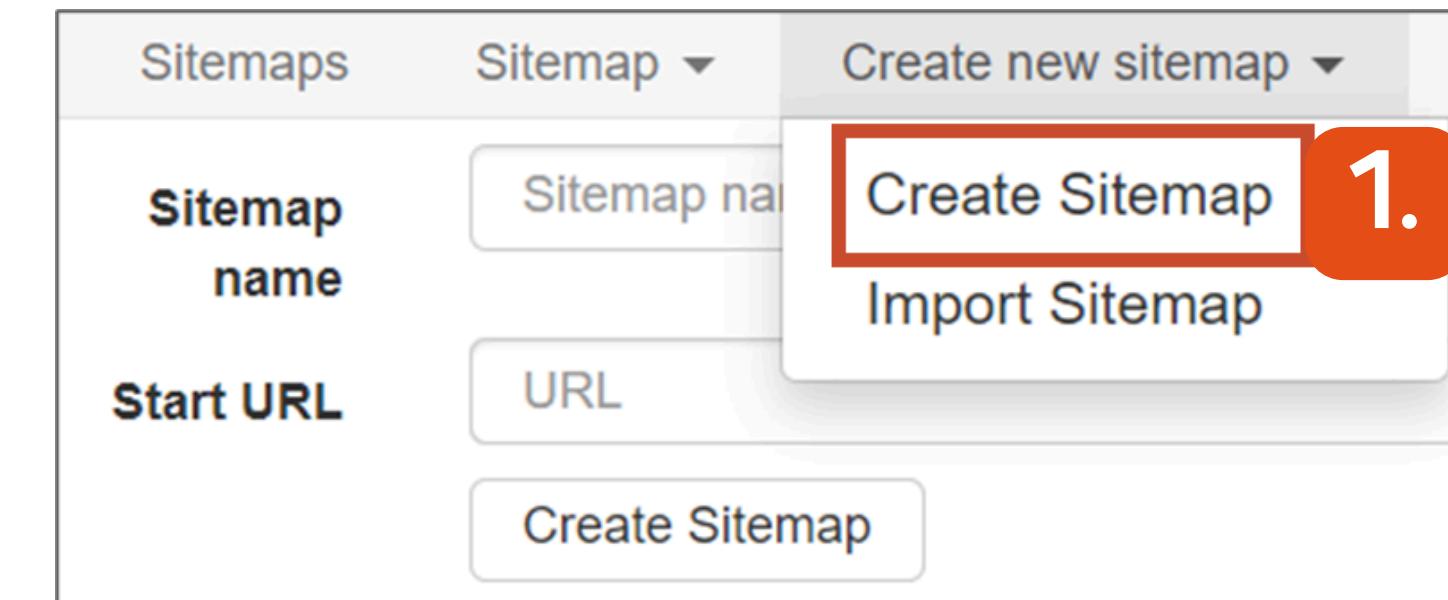
```

Demo

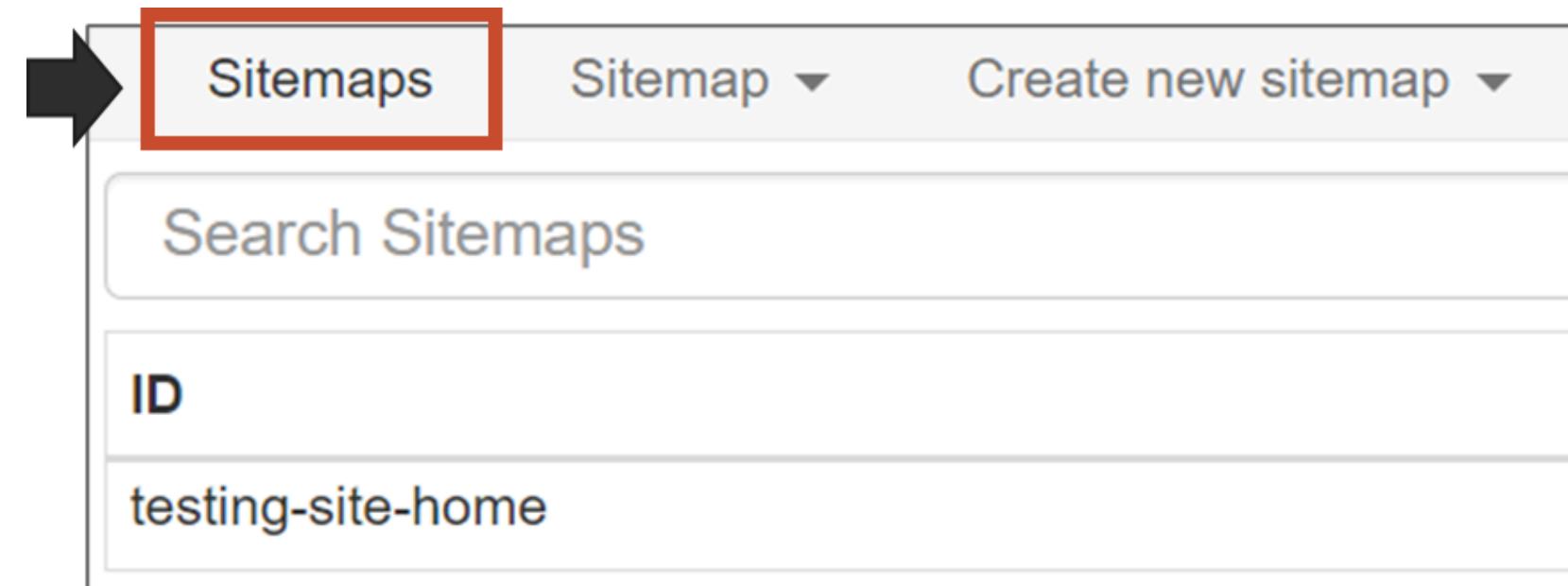
WebScraper - Test site

Steps:

1. Create a sitemap



View all
created
Sitemaps



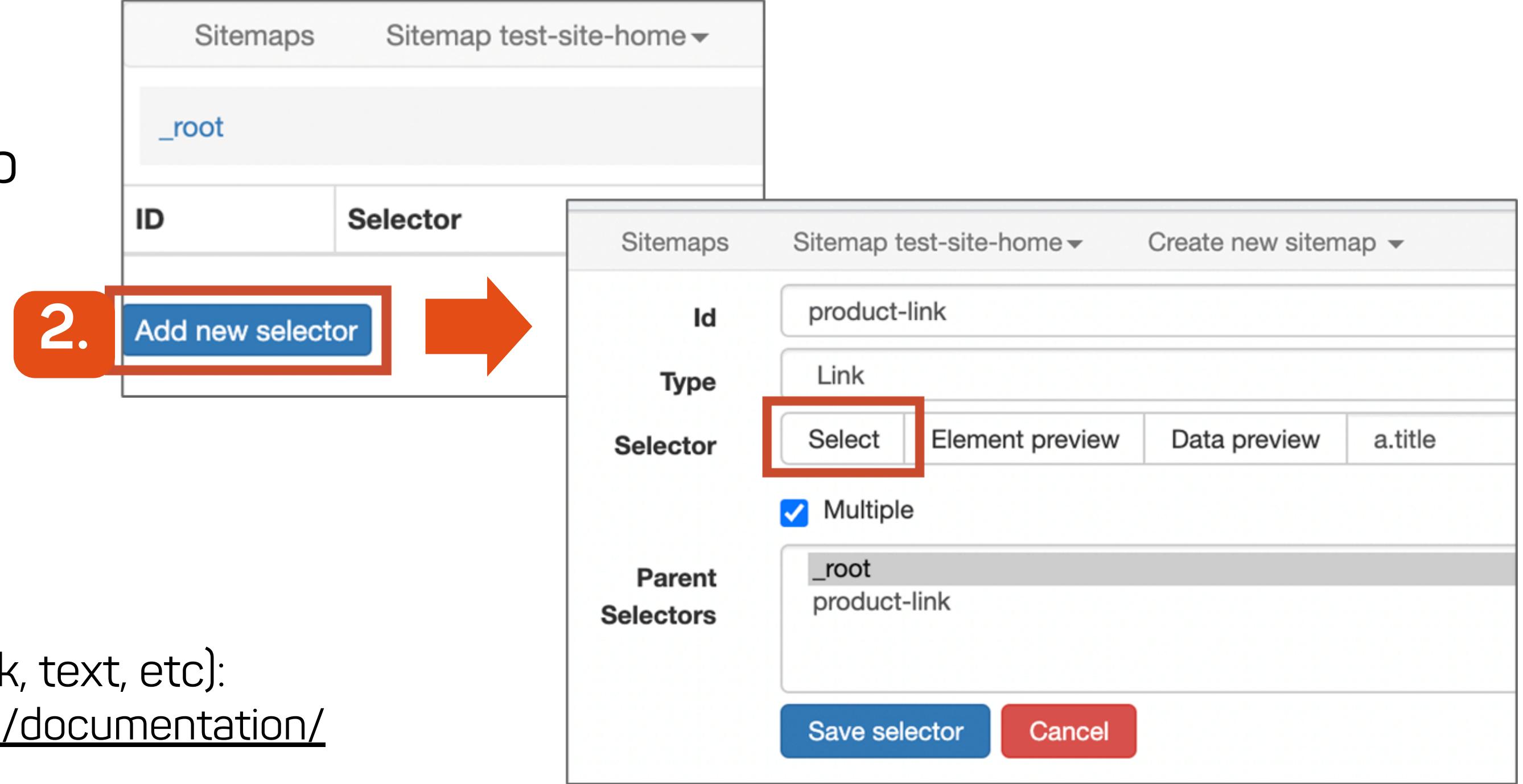
Demo

WebScraper - Test site

Steps:

1. Create a sitemap

2. Add selectors



The screenshot shows the WebScraper interface. On the left, there's a list of sitemaps under the heading "Sitemaps". One sitemap is selected, titled "Sitemap test-site-home". Inside this sitemap, there's a single item named "_root". Below this, there's a table with two columns: "ID" and "Selector". A large orange button labeled "2. Add new selector" is overlaid on the "Selector" column. An orange arrow points from this button to a detailed view on the right.

Screenshot of the Selector Configuration Dialog:

- Id:** product-link
- Type:** Link
- Selector:** Select (selected) Element preview Data preview a.title
- Multiple
- Parent Selectors:** _root product-link

Buttons at the bottom:

- Save selector
- Cancel

More about Selectors (link, text, etc):

<https://www.webscraper.io/documentation/selectors>

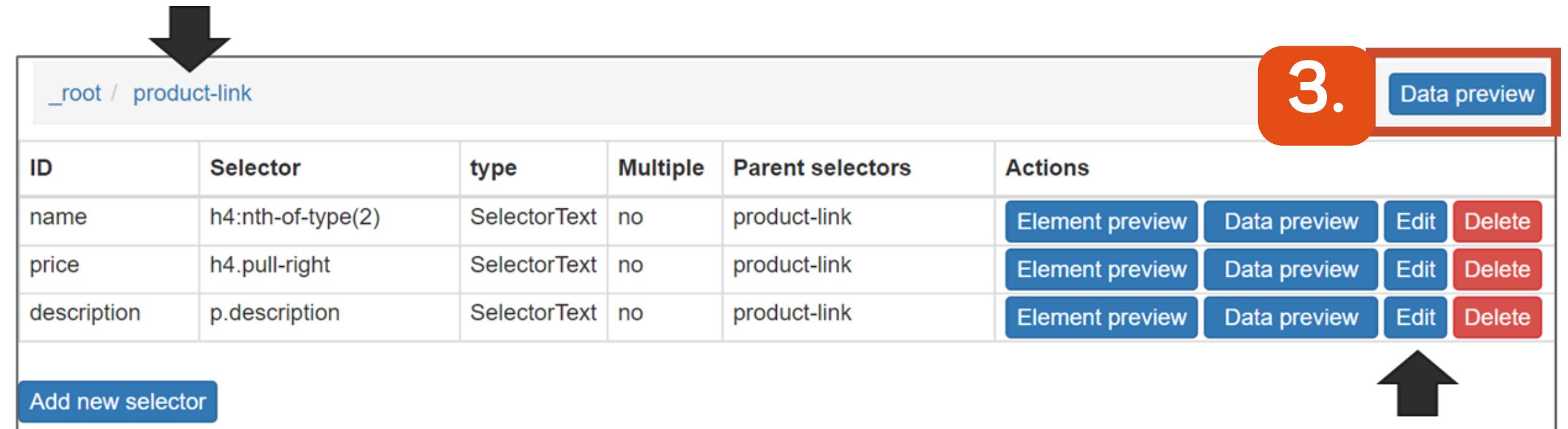
Demo

WebScraper - Test site

Steps:

1. Create a sitemap
2. Add selectors
3. Preview data

Navigate to upper level



ID	Selector	type	Multiple	Parent selectors	Actions			
name	h4:nth-of-type(2)	SelectorText	no	product-link	<button>Element preview</button>	<button>Data preview</button>	<button>Edit</button>	<button>Delete</button>
price	h4.pull-right	SelectorText	no	product-link	<button>Element preview</button>	<button>Data preview</button>	<button>Edit</button>	<button>Delete</button>
description	p.description	SelectorText	no	product-link	<button>Element preview</button>	<button>Data preview</button>	<button>Edit</button>	<button>Delete</button>
Add new selector								

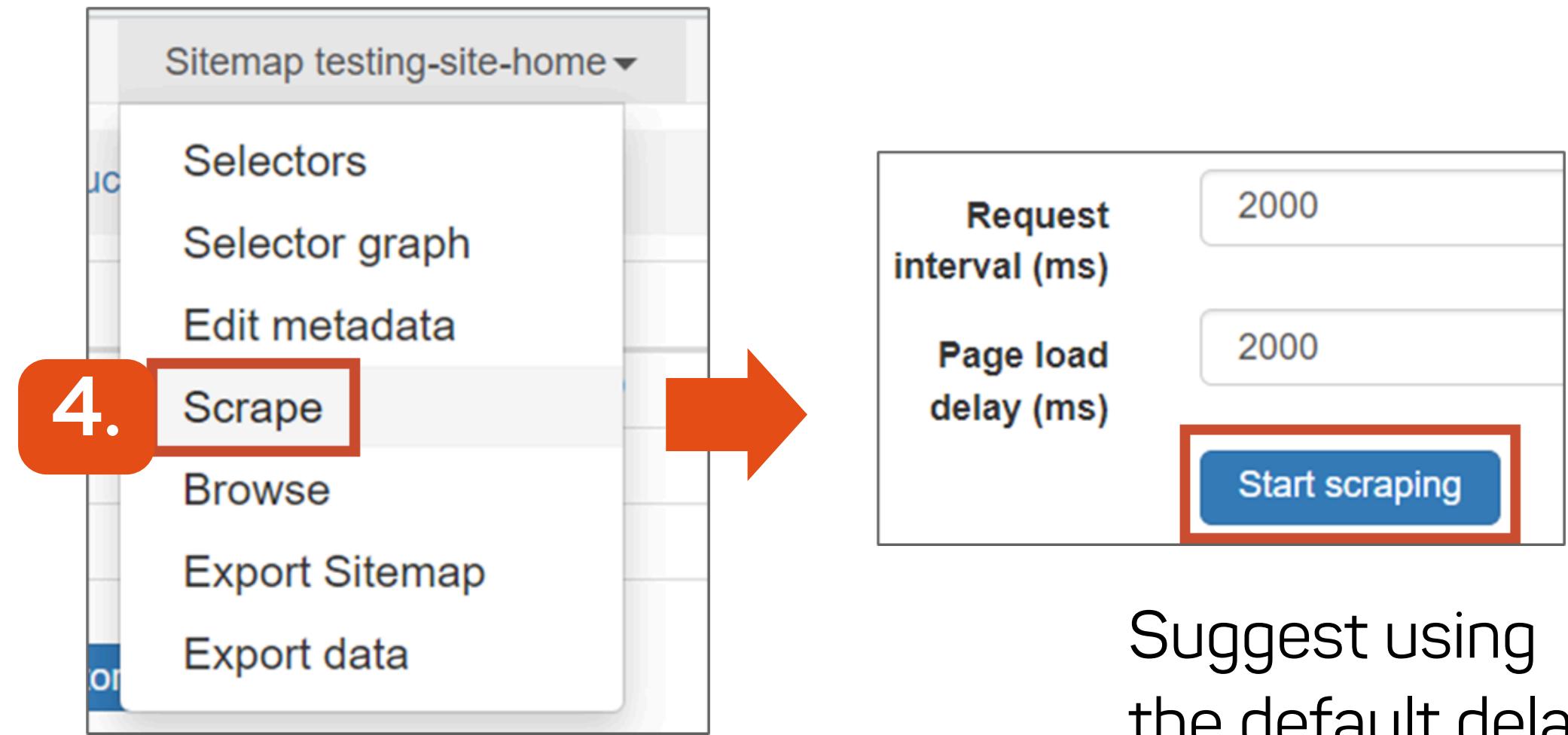
Edit selectors

Demo

WebScraper - Test site

Steps:

1. Create a sitemap
2. Add selectors
3. Preview data
4. **Scrape**



Suggest using
the default delay

(**2000 ms** - can be
longer but NOT shorter)

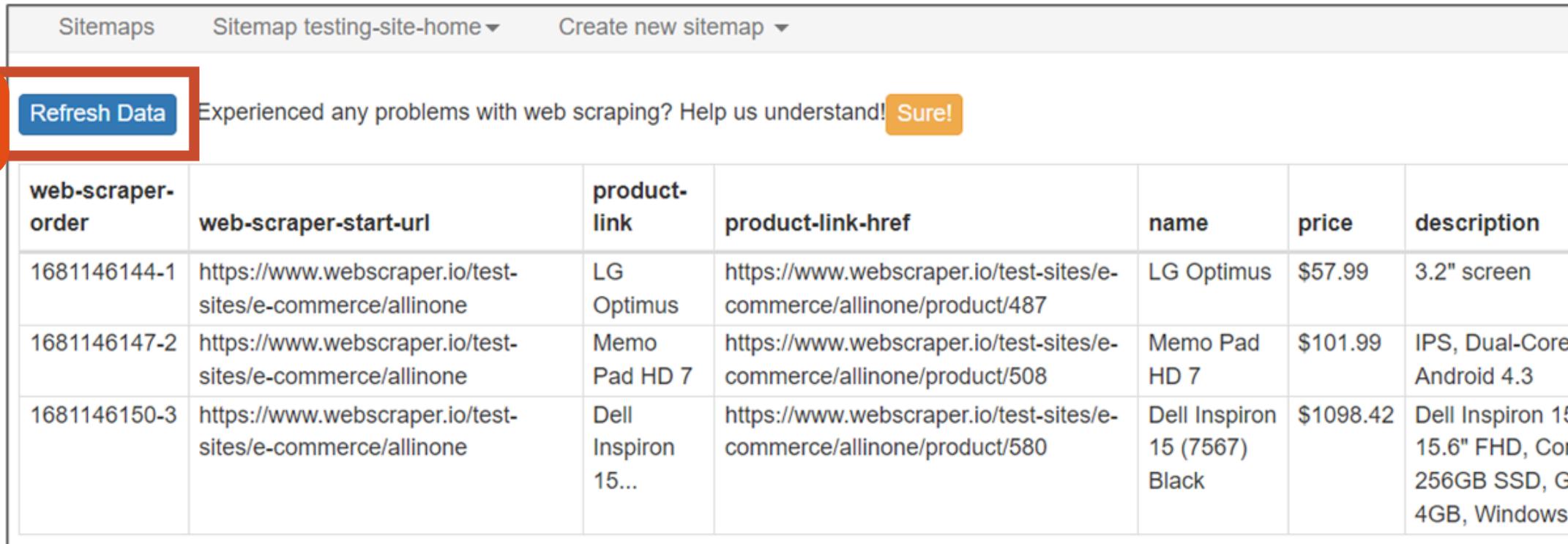
Demo

WebScraper - Test site

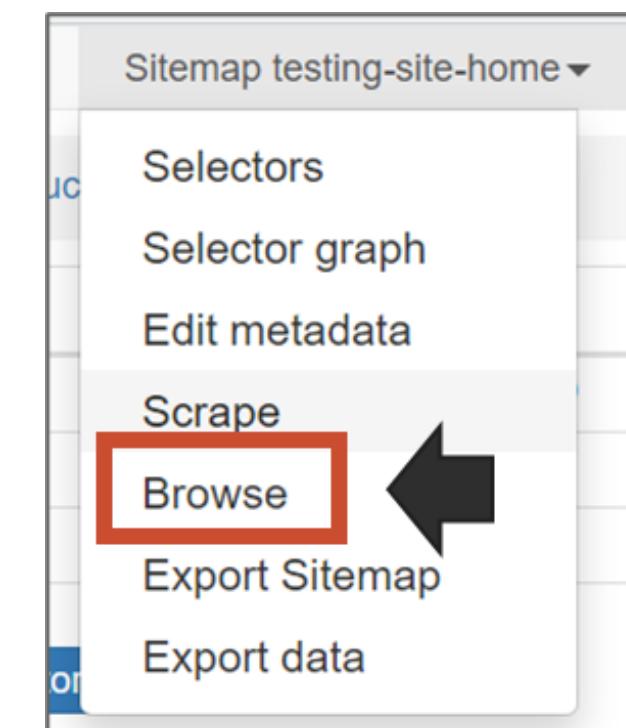
Steps:

1. Create a sitemap
2. Add selectors
3. Preview data
4. Scrape
5. Refresh to check results

5.



web-scraper-order	web-scraper-start-url	product-link	product-link-href	name	price	description
1681146144-1	https://www.webscraper.io/test-sites/e-commerce/allinone	LG Optimus	https://www.webscraper.io/test-sites/e-commerce/allinone/product/487	LG Optimus	\$57.99	3.2" screen
1681146147-2	https://www.webscraper.io/test-sites/e-commerce/allinone	Memo Pad HD 7	https://www.webscraper.io/test-sites/e-commerce/allinone/product/508	Memo Pad HD 7	\$101.99	IPS, Dual-Core Android 4.3
1681146150-3	https://www.webscraper.io/test-sites/e-commerce/allinone	Dell Inspiron 15...	https://www.webscraper.io/test-sites/e-commerce/allinone/product/580	Dell Inspiron 15 (7567) Black	\$1098.42	Dell Inspiron 15 15.6" FHD, Cor 256GB SSD, G 4GB, Windows



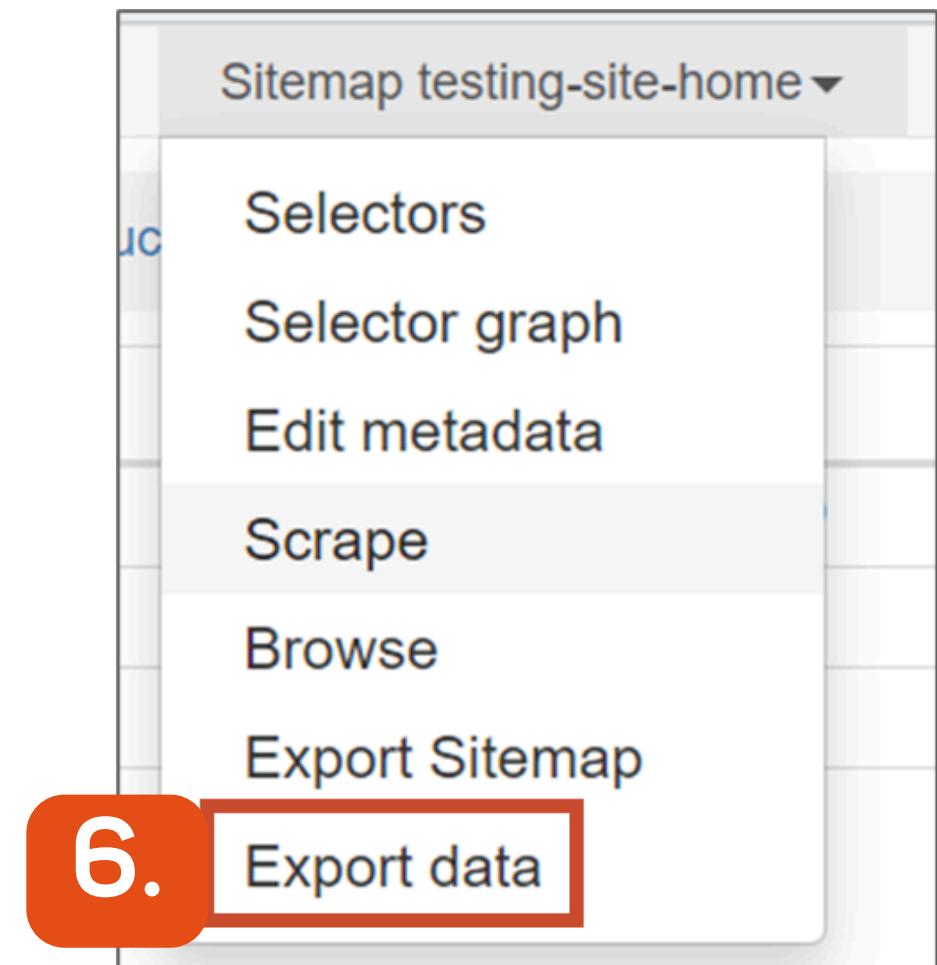
If you left the page and wish to go back to result list, click “**Browse**”.

Demo

WebScraper - Test site

Steps:

1. Create a sitemap
2. Add selectors
3. Preview data
4. Scrape
5. Refresh to check results
6. **Export data** (.csv or .xlsx)



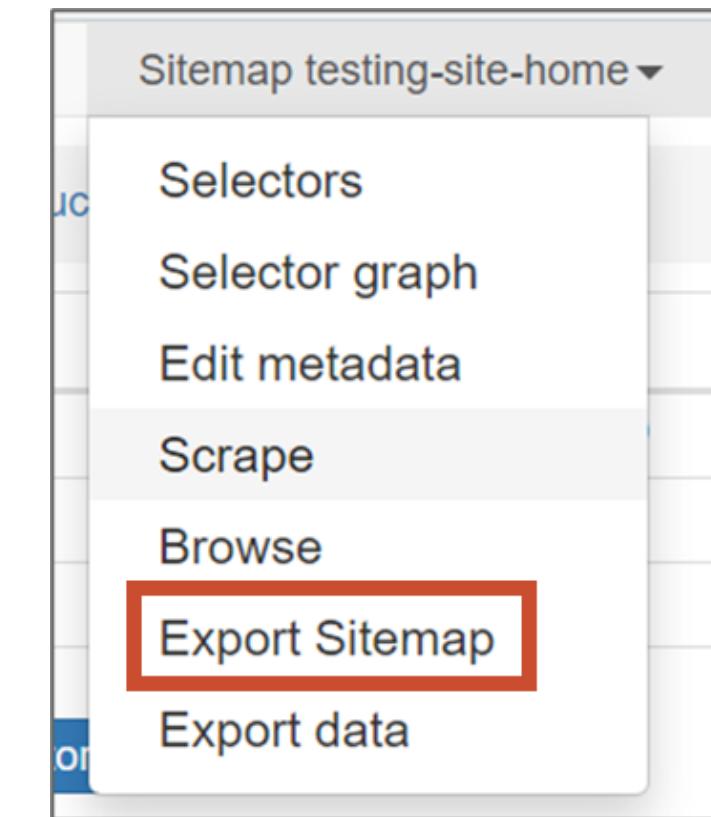
Demo

WebScraper - Test site

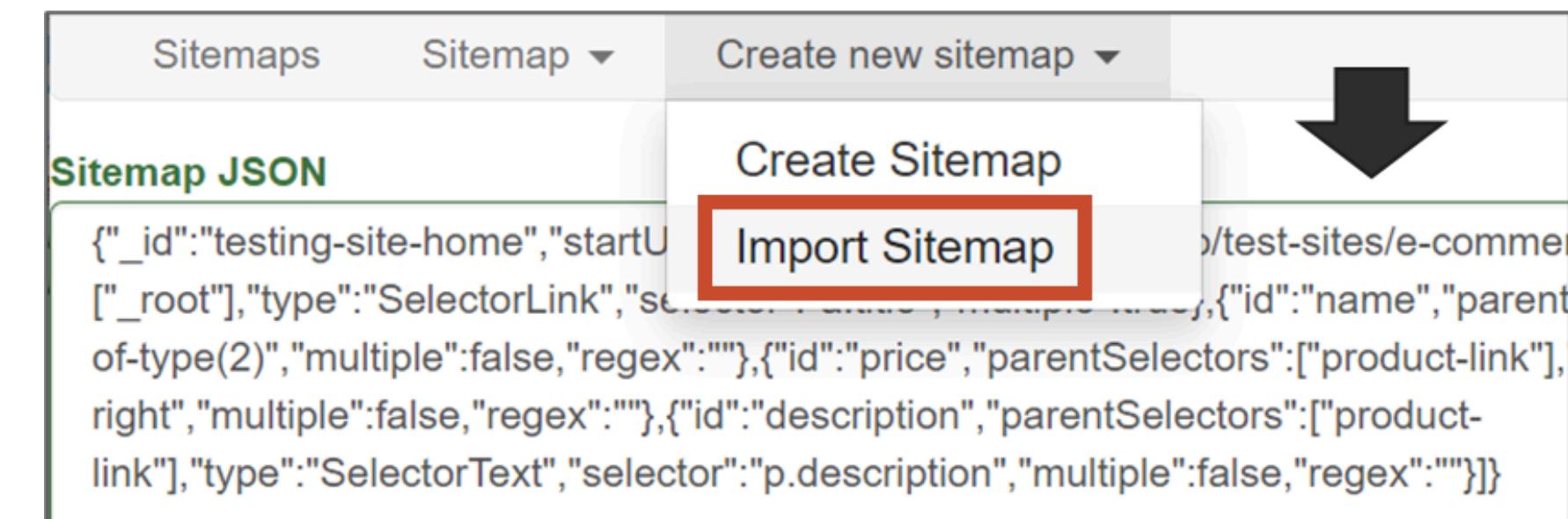
Optional:

Export Sitemap for backup

(to share or reuse in another device)



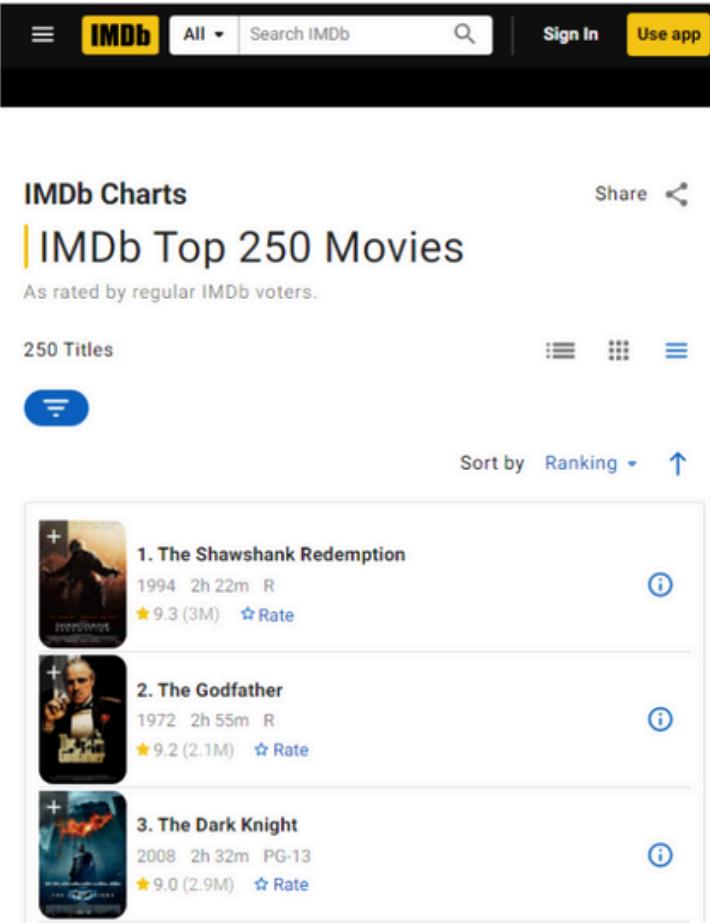
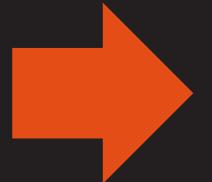
Sitemap is in
JSON format



Hands-on (2)

Get the IMDb Top 250 Movies - with more details

https://www.imdb.com/chart/top/?ref_=nv_mv_250

Rank	Title	Year	IMDb Rating	Genre	Director	Stars	Description
1	Dances with Wolves	1990	8.0/10	Adventure	Dra	Kevin Costner	Kevin Costne Lieutenant John D
2	Dersu Uzala	1975	8.2/10	Adventure	Bio	Akira Kurosawa	Maksim Mur The Russian army
3	Aladdin	1992	8.0/10	Animation	Adv	Ron Clements	Ron Clement A kindhearted stre
4	The Help	2011	8.1/10	Drama		Tate Taylor	Viola Davis E An aspiring autho
5	The Iron Giant	1999	8.1/10	Animation	Acti	Brad Bird	Tim McCanli A young boy befre
6	Life of Brian	1979	8.0/10	Comedy		Terry Jones	Graham Cha Born on the origina
7	Persona	1966	8.1/10	Drama	Thriller	Ingmar Bergman	Bibi Andersss A nurse is put in ch
8	It Happened One Nigh	1934	8.1/10	Comedy	Roma	Frank Capra	Clark Gable C A renegade report
9	The 400 Blows	1959	8.1/10	Crime	Drama	François Truffaut	Jean-Pierre L A young boy, left w
10	The Sound of Music	1965	8.1/10	Biography	Dra	Robert Wise	Georg Hurda A young novice is s
11	The Handmaiden	2016	8.1/10	Drama	Roman	Park Chan-wook	Kim Min-hee A woman is hired a
12	Cool Hand Luke	1967	8.1/10	Crime	Drama	Stuart Rosenberg	Paul Newma A laid back Souther
13	Rebecca	1940	8.1/10	Drama	Film	Alfred Hitchcock	Daphne Du M A self-conscious w
14	Amores Perros	2000	8.1/10	Drama	Thriller	Alejandro G. Iñárr	Emilio Echev A horrific car acci
15	My Father and My Son	2005	8.2/10	Drama	Family	Çagan Irmak	Çetin Tekind The family of a lef
16	Jai Bhim	2021	8.8/10	Crime	Drama	V.T.J. Gnanavel	SuriyaLijo M When a tribal man
17	The Battle of Algiers	1966	8.1/10	Drama	War	Gillo Pontecorvo	Brahim Hadj In the 1950s, fear
18	The Grapes of Wrath	1940	8.1/10	Drama		John Ford	Henry Fonda An Oklahoma fami
19	Hachi: A Dog's Tale	2009	8.1/10	Biography	Dra	Lasse Hallström	Richard Gere A college professo
20	Pather Panchali	1955	8.2/10	Drama		Satyajit Ray	Kanu Banner Impoverished pries
21	Pirates of the Caribbean	2003	8.1/10	Action	Adventu	Gore Verbinski	Ted Elliott Te Blacksmith Will Tu
22	To Be or Not to Be	1942	8.2/10	Comedy	Roma	Ernst Lubitsch	Carole Lomb During the Nazi op

What we want:

List of movies

- Rank
- Title
- Rating
- **Genre**
- **Director**
- **Stars**
- **Description**

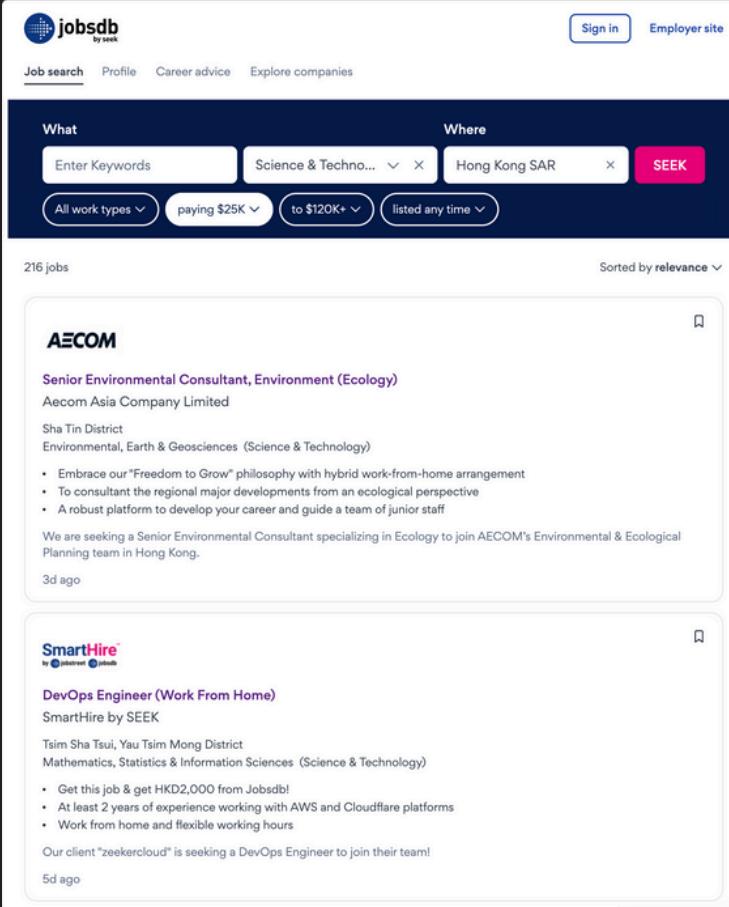
Check before scrape: <https://www.imdb.com/robots.txt>

Conditions of Use: https://www.imdb.com/conditions?ref_=ft_cou

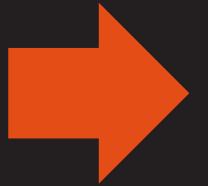
Hands-on (3)

Get the list of jobs from JobsDB in an interested industry

<https://hk.jobsdb.com/jobs/in-Hong-Kong-SAR>



The screenshot shows the JobsDB homepage with a search bar for 'Enter Keywords' (Science & Techno...) and 'Where' (Hong Kong SAR). Below the search bar are filters for 'All work types', 'paying \$25K+', '\$120K+', and 'listed any time'. The results section shows 216 jobs sorted by relevance. Two job listings are visible: one for AECOM (Senior Environmental Consultant) and another for SmartHire (DevOps Engineer). Both listings include company names, locations, and brief descriptions.



A	B	C	D	E
1	job-link	company	location	subject
	Field Application Specialist	Medikonia Limited	Mong Kok, Yau Tsim Mong District	Biological & Biomedical Sciences (Science & Technology)
				What You'll Do We are looking for enthusiastic people to join our scientific support team. We are passionate about bringing cutting-edge and affordable solutions to the scientists, helping Hong Kong move forward in the life sciences sectors. Your day-to-day life will be focusing on providing guidance on biological education, applications, troubleshooting and product selection to respond to our clients' scientific and technical needs. As a Field
2	Data Scientist - Recommendation Engine & AI Innovation	TVB New Media Group Limited	Tseung Kwan O Industrial Estate, Sai Kung District	Mathematics, Statistics & Information Sciences (Science & Technology)
				We are seeking a motivated and curious entry-level Data Scientist to join our dynamic data team. We prioritize growth mindset over years of experience and offer ample opportunities for learning, hands-on projects, and mentorship in cutting-edge machine learning and AI technologies. Responsibilities: Collaborate on designing, implementing, and improving Recommendation Engine systems in productionExplore and integrate Generative AI (GenAI) and Large
3	Chemist / Assistant Chemist		Sheung Wan, Central and Western District	Laboratory & Technical Services (Science & Technology)
				The position will be responsible for: Conduct laboratory testing and inspection according to prescribed procedures/standards. Prepare test sample/ standards/ working solution and operate analytical instruments accurately. Operate various instruments and hand-on experience on equipment maintenance & calibration. Assist to update quality documents and technical documents according to HOKLAS/CNAS & ISO/IEC 17025 requirements,
4	Senior Engineer	GRST Holdings Limited	Sha Tin District	Chemistry & Physics (Science & Technology)
				Ref. No.: GRST-(R&D Centre)-2025004 GRST Holdings Limited, a company within the GRST Group, is an energy technology company with its headquarters located at Hong Kong Science Park. Through our state-of-the-art research facilities, we develop green and sustainable technology for the manufacturing and recycling of batteries and other energy storage systems. Alongside its research and development capabilities, the company also possesses a
5	Data Scientist/AI Engineer	ST Partnership Limited	Central, Central and Western District	Mathematics, Statistics & Information Sciences (Science & Technology)
				Join a leading bank's strong data team as an AI Engineer, collaborating with cross-functional teams to build and integrate machine learning models and AI systems with a focus on practical application. Key Responsibilities: Design and implement AI models to address challenges and improve system capabilities. Collaborate with data scientists, stakeholders, and product managers to develop and deploy AI solutions, ensuring clear communication of technical
6	Senior Project Engineer - Testing and Certification	TUV SUD Hong Kong Limited	Yuen Long, Yuen Long District	Laboratory & Technical Services (Science & Technology)
				Are you ready to embark on an exciting professional journey where innovation, collaboration, and excellence are at the heart of everything we do? Look no further! We are seeking passionate individuals like you to join our dynamic team. This position is responsible for daily chemical project handling, report review and supporting sales for quoting price and technical inquiry. **What You'll Do** Manage project with reasonable manner and reports to project
7	Senior Engineer - Materials Scientist	SAE Magnetics (Hong Kong) Ltd	Science Park, Tai Po District	Materials Sciences (Science & Technology)
				Successful candidate will join the Materials Development Department to develop the latest magnetic recording heads in the R&D and prototype stages within a big cross-functional team. Responsibility: Work closely with material science lab on material characterization, suggest possible improvement to achieve in-depth knowledge of material properties; Work with the teams of material characterization, reliability, process and wafer design to understand the
8	Biomedical Officer I/II/III (Ref: JDB/20240126/EST_BOIII)	Union Medical Centre Ltd	Sha Tin District	Laboratory & Technical Services (Science & Technology)
				Job Duties: To oversee day to day management of all medical equipment. To perform emergency repair & quality assurance check for equipment. To monitor the maintenance schedule. To implement T&C works on new installation and equipment. To perform any ad hoc responsibilities as assigned by supervisors. Job Requirements: Higher Diploma or above with major in Electronics/Electrical/Biomedical Engineering. At least 1 year's relevant working experience and
9	Biomedical Engineer (Ref: MRC2025-007)	Multi-Scale Medical Robotics Center Limited	Sha Tin District	Biological & Biomedical Sciences (Science & Technology)
				Multi-Scale Medical Robotics Center. Multi-Scale Medical Robotics Center Limited (MRC) was established by the Faculty of Medicine and the Faculty of Engineering of the Chinese University of Hong Kong in collaboration with ETH Zürich, Imperial College London, Johns Hopkins University and the University of Hong Kong. MRC focusing on technological innovation with a strong emphasis on clinical translation and direct patient benefits, serves as a
10				

What we want:

List of jobs

- Position
- Company
- Location
- Industry
- Job description

Check before scrape: <https://hk.jobsdb.com/robots.txt>

Conditions of Use: <https://hk.jobsdb.com/terms>

Learn more

- **Common Scenarios:**
 - Scrape from **individual subpages** (covered in class)
 - Scrape **multiple items** from **one single page** ([how-to video](#))
 - **Pages with paginations**
 - “Next” button ([how-to video](#))
 - Page number buttons ([how-to video](#))
- **Documentation:**
 - <https://webscraper.io/documentation/scraping-a-site>

Take-away message

- “When to use what”
 - Simple static webpage -> **Power Query**
 - Static & a few webpages (w/ pagination) -> **Web Scraper plug-in**
 - Dynamic webpage, large scale projects -> **Programming, API**
- **Check before you scrape:**
 - Is the data available somewhere?
 - Technical limitations
 - Legal and ethical restrictions