

Chapter 2: Web Scraping, Applications and Tools

DR. LINDA MAHMOUDI



”يرفع الله الذين آمنوا منكم والذين أوتوا
العلم درجات“



Tools for Web Scraping

Web Scraping tools are specifically developed for extracting data from the internet. Also, known as web harvesting tools or data extraction tools, they are useful for anyone trying to collect specific data from websites as they provide the user with structured data extracting data from a number of websites. Some of the most popular Web Scraping tools are:

- [Import.io](#)
- [Webhose.io](#)
- [Dexi.io](#)
- [Scrapinghub](#)
- [Parsehub](#)

What tools can you use to scrape the web?

We've covered the basics of how to scrape the web for data, but how does this work from a technical standpoint? Often, web scraping requires some knowledge of programming languages, the most popular for the task being [Python](#). Luckily, Python comes with a huge number of [open-source libraries](#) that make web scraping much easier.

These include:

Python Libraries

[BeautifulSoup](#): is Python library, commonly used to parse data from XML and HTML documents. Organizing this parsed content into more accessible trees, BeautifulSoup makes navigating and searching through large swathes of data much easier. It's the go-to tool for many data analysts.

[Scrapy](#): is a Python-based application framework that crawls and extracts structured data from the web. It's commonly used for data mining, information processing, and for archiving historical content. As well as web scraping (which it was specifically designed for) it can be used as a general-purpose web crawler, or to extract data through APIs.

Python Libraries

[Pandas](#): is another multi-purpose Python library used for data manipulation and indexing. It can be used to scrape the web in conjunction with BeautifulSoup. The main benefit of using pandas is that analysts can carry out the entire data analytics process using one language (avoiding the need to switch to other languages, such as R).

Parsehub

A bonus tool, in case you're not an experienced programmer!

[Parsehub](#) is a free online tool (to be clear, this one's not a Python library) that makes it easy to scrape online data. The only catch is that for full functionality you'll need to pay. But the free tool is worth playing around with, and the company offers excellent customer support.

How does a web scraper function?

So, we now know what web scraping is, and why different organizations use it. **But how does a web scraper work?** While the exact method differs depending on the software or tools you're using, all web scraping bots follow three basic principles:

Step 1: Making an HTTP request to a server

Step 2: Extracting and parsing (or breaking down) the website's code

Step 3: Saving the relevant data locally

Step 1: Making an HTTP request to a server

- As an individual, when you visit a website via your browser, you send what's called an HTTP request.
- This is basically the digital equivalent of knocking on the door, asking to come in. Once your request is approved, you can then access that site and all the information on it.
- Just like a person, a web scraper needs permission to access a site.
- Therefore, the first thing a web scraper does is send an HTTP request to the site they're targeting.

Step 2: Extracting and parsing the website's code

- ❖ Once a website gives a scraper access, the bot can read and extract the site's HTML or XML code.
- ❖ This code determines the website's content structure.
- ❖ The scraper will then parse the code (which basically means breaking it down into its constituent parts) so that it can identify and extract elements or objects that have been predefined by whoever set the bot loose! These might include specific text, ratings, classes, tags, IDs, or other information

Step 3: Saving the relevant data locally

- ❖ Once the HTML or XML has been accessed, scraped, and parsed, the web scraper will then store the relevant data locally.
- ❖ As mentioned, the data extracted is predefined by you (having told the bot what you want it to collect).
- ❖ Data is usually stored as structured data, often in an Excel file, such as a .csv or .xls format, mysql database.

- With these steps complete, you're ready to start using the data for your intended purposes. Easy, eh? And it's true...these three steps *do* make data scraping seem easy. In reality, though, the process isn't carried out just once, but countless times. This comes with its own swathe of problems that need solving.
- For instance, badly coded scrapers may send too many HTTP requests, which can crash a site. Every website also has different rules for what bots can and can't do. Executing web scraping code is just one part of a more involved process. Let's look at that now.

How to scrape the web (step-by-step)

Step one: Find the URLs you want to scrape

Step two: Inspect the page

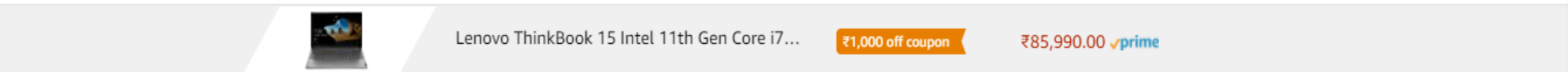
Step three: Identify the data you want to extract

Step four: Write the necessary code

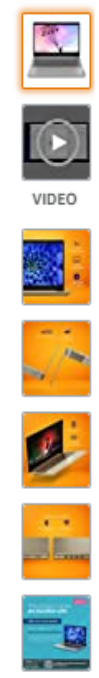
Step five: Execute the code

Step six: Storing the data

- Excel formats are the most common/ mysql database



Computers & Accessories > Laptops > Traditional Laptops Sponsored



Lenovo Ideapad Slim 3 10th Gen Intel Core i3 15.6" (39.62cm) FHD Thin & Light Laptop (8GB/256 GB SSD/UHD Graphics/Windows 10/MS Office/2 Year Warranty/Platinum Grey/1.7Kg), 81WB012DIN

Visit the Lenovo Store 4.5 stars 456 ratings | 127 answered questions

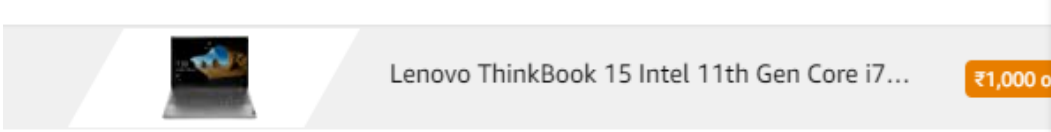
M.R.P.: ₹52,290.00 Deal of the Day: ₹35,990.00 Ends in 11h 19m 41s You Save: ₹16,300.00 (31%) Inclusive of all taxes

FREE delivery: Wednesday, Oct 13 Details EMI starts at ₹1,694. No Cost EMI available EMI options

Save Extra with 4 offers

Share [Email] [Facebook] [Twitter] [Pinterest]

- With Exchange Up to ₹ 18,050.00 off
Without Exchange ₹ 35,990.00 ₹ 52,290.00
Add a Protection Plan:
1 Year Extended Warranty for ₹2,099.00
2 Year Care Plan (2 Year Extended Warranty Plan Post Expiry of 1st Year of Manufacturer Warranty for ₹1,241.00
Ques Care Total Protection Plan for 1 Year Accidental Damage and Liquid Damage for Laptop Between 35001



Computers & Accessories > Laptops > Traditional Laptops



VIDEO



Click to open expanded view

- Back Alt+Left Arrow
- Forward Alt+Right Arrow
- Reload Ctrl+R
- Save as... Ctrl+S
- Print... Ctrl+P
- Cast...
- Create QR code for this page
- Translate to English
- View page source Ctrl+U
- Inspect Ctrl+Shift+I

Elements Console Sources Network Performance Memory

```

</div>
<div id="ppd">
  <div id="rightCol" class="rightCol rightCol-bbcxo">
    <div id="leftCenterCol" class="leftCenterCol">
      </div>
    <div id="leftCol" class="leftCol"></div>
    <div id="centerCol" class="centerColAlign centerCol">
      <div id="title_feature_div" class="celwidget" data-cel-widget="title_feature_div">
        <style type="text/css">
          .product-title-word-break {
            word-break: break-word;
          }
        </style>
        <div id="titleSection" class="a-section a-spacing">
          <h1 id="title" class="a-size-large a-spacing">
            <span id="productTitle" class="a-size-larg
              "
    
```

Styles Computed Layout Event Listeners

Filter :hov .cls +

```

element.style {
}
#productTitle {
  font-size: 19px!important;
}
.product-title-word-break {
  word-break: break-word;
}
.a-size-large {
  text-rendering: optimizeLegibility;
}
.a-size-large {
  font-size: 24px!important;
  line-height: 32px!important;
}
* {
  -moz-box-sizing: border-box;
  -webkit-box-sizing: border-box;
  box-sizing: border-box;
}
Inherited from h1#title.a-s...
h1 {
  font-weight: 400;
  font-size: 28px;
  line-height: 36px;
}
h1, h2, h3, h4 {
  text-rendering: optimizeLegibility;
}
h1 {
  font-size: 2em;
  font-weight: bold;
}

```

Lenovo Ideapad Slim 3 10th Gen Intel Cor Graphics/Windows 10/MS Office/2 Year War

... acing-none span#productTitle.a-size-large.product-title-word-break ...

Step 5: Run the code by using the below command:

```
python web-scrap.py
```

Step 6: A file called “products.csv” is created, which contains the extracted data.

