

Web Scraping with Python (Indiameters)

What is web scraping?

Web scraping, **web harvesting**, or **web data extraction** is [data scraping](#) used for [extracting data](#) from [websites](#). Web scraping software may access the [World Wide Web](#) directly using the [Hypertext Transfer Protocol](#), or through a web browser. While web scraping can be done manually by a software user, the term typically refers to automated processes implemented using a [bot](#) or [web crawler](#). It is a form of copying, in which specific data is gathered and copied from the web, typically into a central local [database](#) or spreadsheet, for later [retrieval](#) or [analysis](#).

The crawler:

A web crawler, which we generally call a “spider,” is an artificial intelligence that browses the internet to index and search for content by following links and exploring, like a person with too much time on their hands. In many projects you first “crawl” the web or one specific website to discover URLs which then you pass on to your scraper.

The Spider:

A web scraper is a specialized tool designed to accurately and quickly extract data from a web page. Web scrapers vary widely in design and complexity, depending on the project. An important part of every scraper is the data locators (or selectors) that are used to find the data that you want to extract from the HTML file - usually xpath, css selectors, regex or a combination of them is applied.

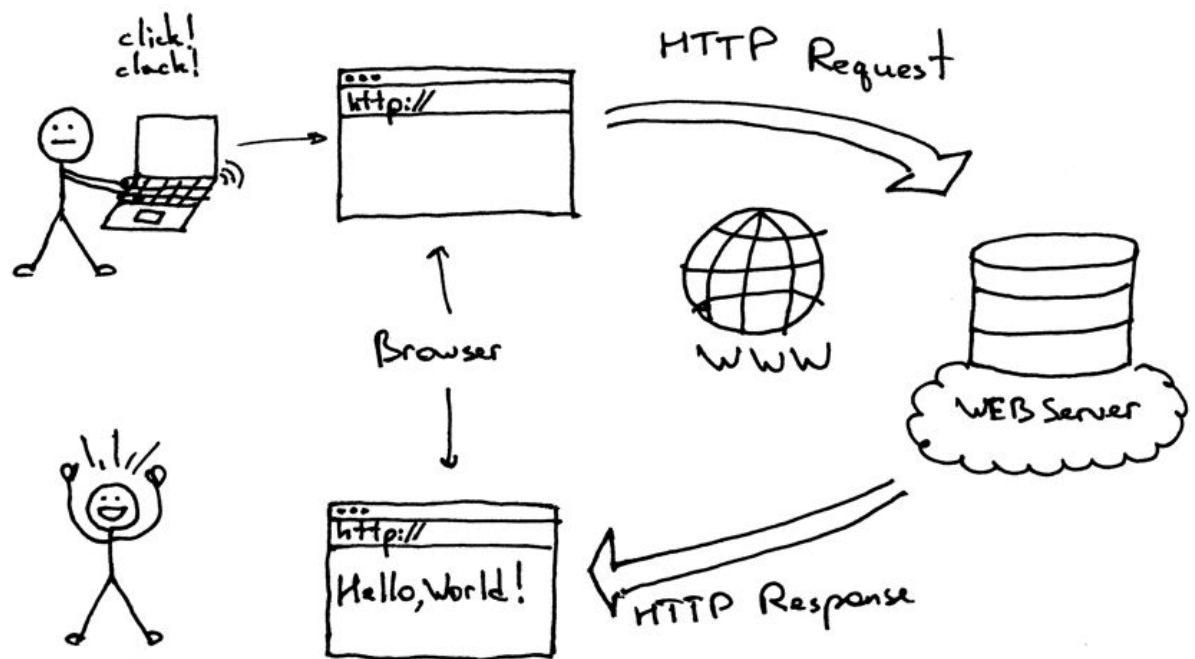
Web Scraping process:

1. Identify target website
2. Collect URLs of the pages where you want to extract data from
3. Make a request to these URLs to get the HTML of the page
4. Use locators to find the data in the HTML
5. Save the data in a JSON or CSV file or Excel or some other structured format

Libraries used in Web Scraping:

Requests: It is one of the most basic yet essential library for web scraping. ‘Requests’ lets us make HTML requests to the website’s server for retrieving the data on its page. Getting the HTML content of a web page is the first and

foremost step of web scraping.



lxml: We know the requests library cannot parse the HTML retrieved from a web page. Therefore, we require lxml, a high performance, blazingly fast, production-quality HTML, and XML parsing Python library.

It combines the speed and power of Element trees with the simplicity of Python. It works well when we're aiming to scrape large datasets. The combination of requests and lxml is very common in web scraping. It also allows you to extract data from HTML using XPath and CSS selectors.

BeautifulSoup is perhaps the most widely used Python library for web scraping. It creates a parse tree for parsing HTML and XML documents. BeautifulSoup automatically converts incoming documents to Unicode and outgoing documents to UTF-8. One of the primary reasons the BeautifulSoup library is so popular is that it is easier to work with and well suited for beginners. We can also combine BeautifulSoup with other parsers like **lxml**. But all this ease of use comes with a cost – it is slower than **lxml**. Even while using **lxml** as a parser, it is slower than pure **lxml**.

One major advantage of the BeautifulSoup library is that it works very well with poorly designed HTML and has a lot of functions. The combination of **Beautiful Soup** and **Requests** is quite common in the industry.

Selenium: There is a limitation to all the Python libraries we have discussed so far – we cannot easily scrape data from dynamically populated websites. It happens because sometimes the data present on the page is loaded through JavaScript. In simple words, if the page is not static, then the Python libraries mentioned earlier struggle to scrape the data from it.

That's where Selenium comes into play

Scrapy : Scrapy is not just a library; it is an entire web scraping framework. Scrapy provides spider bots that can crawl multiple websites and extract the data. With Scrapy, you can create your spider bots, host them on Scrapy Hub, or as an API. It allows you to create fully-functional spiders in a matter of a few minutes. You can also create pipelines using Scrapy.

The best thing about Scrapy is that it's asynchronous. It can make multiple HTTP requests simultaneously. This saves us a lot of time and increases our efficiency

Scraping Indian Cities population from the websites

Here we will be using **BeautifulSoup**, **requests** and **pandas** libraries to scrape the data. And we will be using "html.parser" as parser to parse the HTML content. Note: always check the legality of website whether the website you are going to scrape allows you to scrape data from it. You can check it by adding "/robots.txt" at the end of the URL of the website.

So let's scrape the cities population from the website:

URI of website: ("https://www.worldometers.info/world-population/india-population/")
And we will save our scrape data in excel format

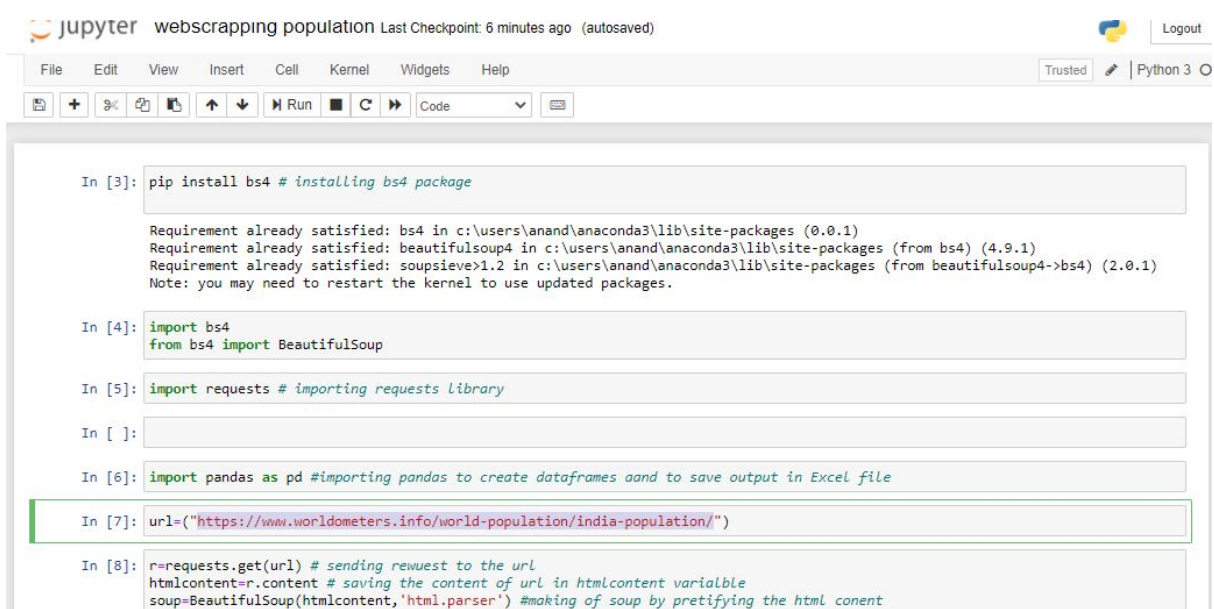
We will scrape below table:

Main Cities by Population in India

(includes boroughs, districts, urban agglomerations, etc.)

#	CITY NAME	POPULATION
1	Mumbai	12,691,836
2	Delhi	10,927,986
3	Bengaluru	5,104,047
4	Kolkata	4,631,392
5	Chennai	4,328,063
6	Ahmedabad	3,719,710
7	Hyderabad	3,597,816
8	Pune	2,935,744
9	Surat	2,894,504
10	Kanpur	2 823 249

(Below is the snap of the code with its proper explanation and output):



The screenshot shows a Jupyter Notebook titled "webscrapping population" with a last checkpoint 6 minutes ago. The interface includes a menu bar (File, Edit, View, Insert, Cell, Kernel, Widgets, Help) and a toolbar with icons for file operations, running, and code execution. The notebook contains eight code cells:

```
In [3]: pip install bs4 # installing bs4 package

Requirement already satisfied: bs4 in c:\users\anand\anaconda3\lib\site-packages (0.0.1)
Requirement already satisfied: beautifulsoup4 in c:\users\anand\anaconda3\lib\site-packages (from bs4) (4.9.1)
Requirement already satisfied: soupsieve>1.2 in c:\users\anand\anaconda3\lib\site-packages (from beautifulsoup4->bs4) (2.0.1)
Note: you may need to restart the kernel to use updated packages.

In [4]: import bs4
        from bs4 import BeautifulSoup

In [5]: import requests # importing requests library

In [ ]: 

In [6]: import pandas as pd #importing pandas to create dataframes aand to save output in Excel file

In [7]: url=("https://www.worldometers.info/world-population/india-population/")

In [8]: r=requests.get(url) # sending request to the url
        htmlcontent=r.content # saving the content of url in htmlcontent variable
        soup=BeautifulSoup(htmlcontent,'html.parser') #making of soup by pretifying the html content
```



```

ut[61]: [ <tr style="background-color:#F8F8F8;"> <th style="width:20px;">#</th> <th style="font-weight: 100; font-size:16px; text-align:left; padding-top:5px; padding-bottom:5px">CITY NAME</th> <th style="font-weight: 100; font-size:16px; text-align:left; padding-top:5px; padding-bottom:5px">POPULATION</th> </tr>
<tr> <td style="color:#999; text-align:center; vertical-align:middle">1</td> <td style="font-weight: bold; font-size:17px; text-align:left; padding-left:10px; padding-top:5px; padding-bottom:5px">Mumbai</td> <td style="font-weight: bold; font-size:17px; text-align:left; padding-left:10px; padding-top:5px; padding-bottom:5px">12,691,836</td> </tr>
<tr> <td style="color:#999; text-align:center; vertical-align:middle">2</td> <td style="font-weight: bold; font-size:17px; text-align:left; padding-left:10px; padding-top:5px; padding-bottom:5px">Delhi</td> <td style="font-weight: bold; font-size:17px; text-align:left; padding-left:10px; padding-top:5px; padding-bottom:5px">10,927,986</td> </tr>
<tr> <td style="color:#999; text-align:center; vertical-align:middle">3</td> <td style="font-weight: bold; font-size:17px; text-align:left; padding-left:10px; padding-top:5px; padding-bottom:5px">Bengaluru</td> <td style="font-weight: bold; font-size:17px; text-align:left; padding-left:10px; padding-top:5px; padding-bottom:5px">5,104,047</td> </tr>
<tr> <td style="color:#999; text-align:center; vertical-align:middle">4</td> <td style="font-weight: bold; font-size:17px; text-align:left; padding-left:10px; padding-top:5px; padding-bottom:5px">Kolkata</td> <td style="font-weight: bold; font-size:17px; text-align:left; padding-left:10px; padding-top:5px; padding-bottom:5px">4,631,392</td> </tr>
<tr> <td style="color:#999; text-align:center; vertical-align:middle">5</td> <td style="font-weight: bold; font-size:17px; text-align:left; padding-left:10px; padding-top:5px; padding-bottom:5px">Chennai</td> <td style="font-weight: bold; font-size:17px; text-align:left; padding-left:10px; padding-top:5px; padding-bottom:5px">4,328,063</td> </tr>
<tr> <td style="color:#999; text-align:center; vertical-align:middle">6</td> <td style="font-weight: bold; font-size:17px; text-align:left; padding-left:10px; padding-top:5px; padding-bottom:5px">Ahmedabad</td> <td style="font-weight: bold; font-size:17px; text-align:left; padding-left:10px; padding-top:5px; padding-bottom:5px">4,328,063</td> </tr>

```

```
In [64]: City_Name=[]
for row in anda: # with the help of nested loop extracting all the rows of the cityname columns and saving it in list
    for cell in row.find_all("td",attrs={'style':'font-weight: bold; font-size:17px; text-align:left; padding-left:10px; padding-top:5px; padding-bottom:5px'}):
        sa=cell.text
        City_Name.append(sa)
City_Name
```

```
Out[64]: ['Mumbai',
'Delhi',
'Bengaluru',
'Kolkata',
'Chennai',
'Ahmedabad',
'Hyderabad',
'Pune',
'Surat',
'Kanpur',
'Jaipur',
'Navi Mumbai',
'Lucknow',
'Nagpur',
'Indore',
'Patna',
'Bhopal',
'Ludhiana',
'Tirunelveli',
'Agra',
'Vadodara',
'Gorakhpur',
'Nashik',
'Raipur']
```

Activate Wi

```
In [66]: City_population=[]
for row in anda: # with the help of nested loop extracting all the rows of the Population columns and saving it in list
    for cell in row.find_all("td",attrs={"style":" font-size:17px; text-align:left; padding-left:10px; padding-top:5px; padding-bottom:5px"}):
        sd=cell.text
        City_population.append(sd)
City_population
```

```
Out[66]: ['12,691,836',
'10,927,986',
'5,104,047',
'4,631,392',
'4,328,063',
'3,719,710',
'3,597,816',
'2,935,744',
'2,894,504',
'2,823,249',
'2,711,758',
'2,600,000',
'2,472,011',
'2,228,018',
'1,837,041',
'1,599,920',
'1,599,914',
'1,545,368',
'1,435,844',
'1,430,055',
'1,409,476',
'1,324,570',
'1,289,497',
'1,284,600']
```

Activate Wi
Go to Settings

OUTPUT:

[illegible]

22	20	Vadodara	1,409,476				
23	21	Gorakhpur	1,324,570				
24	22	Nashik	1,289,497				
25	23	Pimpri	1,284,606				
26	24	Kalyan	1,262,255				
27	25	Thane	1,261,517				
28	26	Meerut	1,223,184				
29	27	Nowrangapur	1,220,946				
30	28	Faridabad	1,220,229				
31	29	Ghaziabad	1,199,191				
32	30	Dombivli	1,193,000				
33	31	Rajkot	1,177,362				
34	32	Varanasi	1,164,404				
35	33	Amritsar	1,092,450				
36	34	Allahabad	1,073,438				
37	35	Visakhapatnam	1,063,178				
38	36	Tenali	1,034,724				
39	37	Jabalpur	1,030,168				
40	38	Haora	1,027,672				
41	39	Aurangabad	1,016,441				
42	40	Shivajinagar	1,000,000				
47		Guwahati	899,094				
48		Gwalior	882,458				
49		Vijayawada	874,587				
50		Mysore	868,313				
51		Ranchi	846,454				
52		Hubli	840,214				
53		Jalandhar	785,178				
54		Thiruvananthapuram	784,153				
55		Salem	778,396				
56		Tiruchirappalli	775,484				
57		Kota	763,088				
58		Bhubaneswar	762,243				
59		Aligarh	753,207				
60		Bareilly	745,435				
61		Moradabad	721,139				
62		Bhiwandi	707,035				
63		Raipur	679,995				
64		Gorakhpur	674,246				
65		Bhilai	625,138				
66		Jamshedpur	616,338				
67		Borivli	609,617				

Conclusion: Hence we successfully scrape the table data from the target URL and we can further use this data for data visualization and in chatbot .

