

Web Scraping :

- Technique to fetch data from websites
- automation of data extraction process from websites.
- Construction of agent to download / parse / organize data from the web in an automated manner.
- technique to collect content and data from the internet.
- done with the help of web scraping softwares (web scrapers) automatically load and extract based on user requirements.
- saved in local file analyzed as needed.

We can use : Java / Python / Ruby / Node

Some software companies have designed tools that enable other people to use scraping techniques by means of attractive and powerful user interface.

Video / image / text / ~~url~~ or sentiments.

Some definitions:

→ HTTP: hypertext transfer Protocol

Machine interchange information

Transported over the internet enable data exchange (via WWW.)

-declarative Document type definition (HTML, XML, JSON...)

→ HTML: hypertext Markup Language (Standard)

Parsing

analyzing strings
and symbols to
reveal the data
you need only

Crawling

Moving across / Through
a website to gather data
from more than 1 URL.

→ JSON: JavaScript Open Notation
Readable Text used to transmit data objects
consisting of attribute-value pairs.

→ API: Application programming interface
Set of rules / protocols to build a software application. In the context
web.
Used to gather clean data from a website (data that is not
wrapped in HTML, Javascript, bound in HTTP...)

Web Scraping Techniques:

Manual Extraction Techniques

Manually copy pasting

Automated Extraction Techniques:

Automatically based on user
requirement.

Parse
make something understandable by analyzing it part by part
Convert information to easier form.

HTML Parsing: Extracting information based on user requirements
Using ~~HTML~~ JavaScript

DOM Parsing: Document Object Model

defines an interface that enables a user to modify/update
style / Structure / Content of XML document

Challenges of web scraping:

Data Warehousing:

the data warehousing infrastructure should be properly built
before storing / Exporting data.

Website structure changes:

Every user interface is updated to improve experience.
⇒ requires web scrapers to modify their code.

Anti-Scraping Technologies:

Some websites use anti-scraping technologies,
they use IP blocking mechanism.

Quality of data extracted:

Records that do not meet the quality of information will affect
the integrity of data.

Scraping libraries in Python:

Beautiful Soup: parse data from XML and HTML.
makes navigating and searching through large swathes much easier.

Scrapy: Python-based application framework that crawls and extracts structured data from the web.
It can be used as a general-purpose web crawler, extract data through API.

Pandas: Used for data manipulation / indexing
" to scrape the web in conjunction with BeautifulSoup.
if we use Pandas we can continue to analyse using (language).

Parsehubs (not in Python)

Web Scrapers

Step 1: Making HTTP request to a server (Knocking the door)

Step 2: Extracting and parsing the website's code (The bot can read/extract HTML/XML code)

Step 3: Saving relevant data (stored in excel file
Such as .CSV, .XLS format, mysql database)

Steps of Scraping:

- 1) Find URL
- 2) Inspect the Page
- 3) Identify the data you want to extract
- 4) Write the necessary code
- 5) Execute the code (python web-scrap.py)
- 6) Store the data (.CSV)

PYTHON IDLE

python installation comes with an Integrated Development and Learning Environment. (Edit .py files with ease)

Best place shell to experiment python: shell (a basic Read-Eval-Print loop)

```
from urllib.request import urlopen
url = "http://olympus.realpython.org/profiles/aphrodite"
```

```
page = urlopen(url)
```

```
html_bytes = page.read
```

```
html = html_bytes.decode("utf-8")
```

```
print(html) → (see the content of the page)
```

Extract text from HTML:

```
title_index = html.find("<title>")
```

title_index

↑ index of the <title>

```
start_index = title_index + len("<title>")
```

start_index

↑ index of the title itself.

```
end_index = html.find("</title>")
```

end_index

```
title = html[start_index:end_index]
```

title

} extract the title by slicing HTML string

Regular Expressions

```
import re
```

Use special characters (metacharacters)

```
re.findall("ab*c", "ac")
```

↓

we can replace it with "Beautiful Soup"

```

from bs4 import BeautifulSoup
from urllib.request import urlopen
    ↖ requesting page content
url = "http://olympus.realpython.org/profiles/dionysus"
page = urlopen(url)
html = page.read().decode("utf-8")
soup = BeautifulSoup(html, "html.parser")

```

```
print(soup.get_text())
```

```
soup.find_all("img") ↖ it gives me image url
```

```
image1, image2 = soup.find_all("img")
```

```
image1.name
```

```
image1["src"]
```

```
soup.title (when using BeautifulSoup it corrects mistakes)
```

<TITLE> → <title>

```
soup.title.string ↖ string between title tags,
```

```
soup.find_all("img", src="/static/dionysus.jpg")
```

Interact with HTML forms:

MechanicalSoup: web browser with no graphical user interface

```
import mechanicalsoup
```

```
browser = mechanicalsoup.Browser
```

```
url = "http://olympus.realpython.org/login"
```

```
page = browser.get(url)
```

```
page ← <Response [200]>
```

```
type(page.soup) → mechanicalsoup uses BeautifulSoup to parse HTML.
```

```
page.soup ↖ view the HTML
```

```
import mechanicalsoup
```

Username:
Password:

#1

```
browser = mechanicalsoup.Browser()  
url = "http://olympus.real.org/login"  
login_page = browser.get(url)  
login_html = login_page.soup
```

#2

```
form = login_html.select("form")[0]  
form.select("input")[0]["value"] = "Zeus"  
" " " " [1] " " = "Thunder"
```

#3

```
profiles_page = browser.submit(form, login_page.url)
```

profiles_page.url → to confirm that the submission successfully redirected to the profiles_page.

login_html.select("form") returns a list of all <form> elements on the page
browser.submit() → to submit the form.

Obtain the URL for each link on the profiles page.

```
links = profiles_page.soup.select("a")
```

```
for link in links:
```

```
    address = link["href"]
```

```
    text = link.text
```

```
    print(f"{text} : {address}")
```

} iterate through links.

storing scraped data to database

storing data to CSV file:

CSV: comma-separated values
Supported by Microsoft Excel

to write to a CSV file in python:

1 - Open CSV file in the write mode:

```
import CSV
csvfile = open('test.csv', 'w+')
```

read and write
'r': read
'w': write

to support writing non-ASCII values to a CSV file:
specify the character encoding!

```
csvfile = open('test.csv', 'w', encoding='UTF8')
```

2 - Create a CSV writer object:

```
Writer = CSV.writer(csvfile)
```

pass the opened file as its argument.

3 - Write data to the CSV file:

```
Writer.writerow(('number', 'number plus 2', 'number times 2'))
```

for i in range(10):

```
Writer.writerow((i, i+2, i*2))
```

4 - Close the CSV file:

```
csvfile.close()
```