

Introduction to Data Analytics

AIGC5000

Lecture10-Web Scraping

Instructor: Parisa Pouladzadeh

Email: parisa.pouladzadeh@humber.ca

What is Web Scraping ?

- Web scraping is the process of using bots to extract content and data from a website.
- web scraping extracts underlying HTML code and, with it, data stored in a database.
- It is a technique to fetch data and information from websites.
- Everything you see on a webpage can be scraped.

Use of Web Scraping

Companies can use this data to fix the optimal pricing for their products so that they can obtain maximum revenue.

Preparing dataset for your ML model.

Web scraping news sites can provide detailed reports on the current news to a company.

Companies can also use Web scraping for email marketing.

Use of Web Scraping

1. Marketing: Web scraping is used by many companies to collect information about their products or services from various social media websites to get a general public sentiment. Also, they extract email ids from various websites and then send bulk promotional emails to the owners of these email ids.

Use of Web Scraping

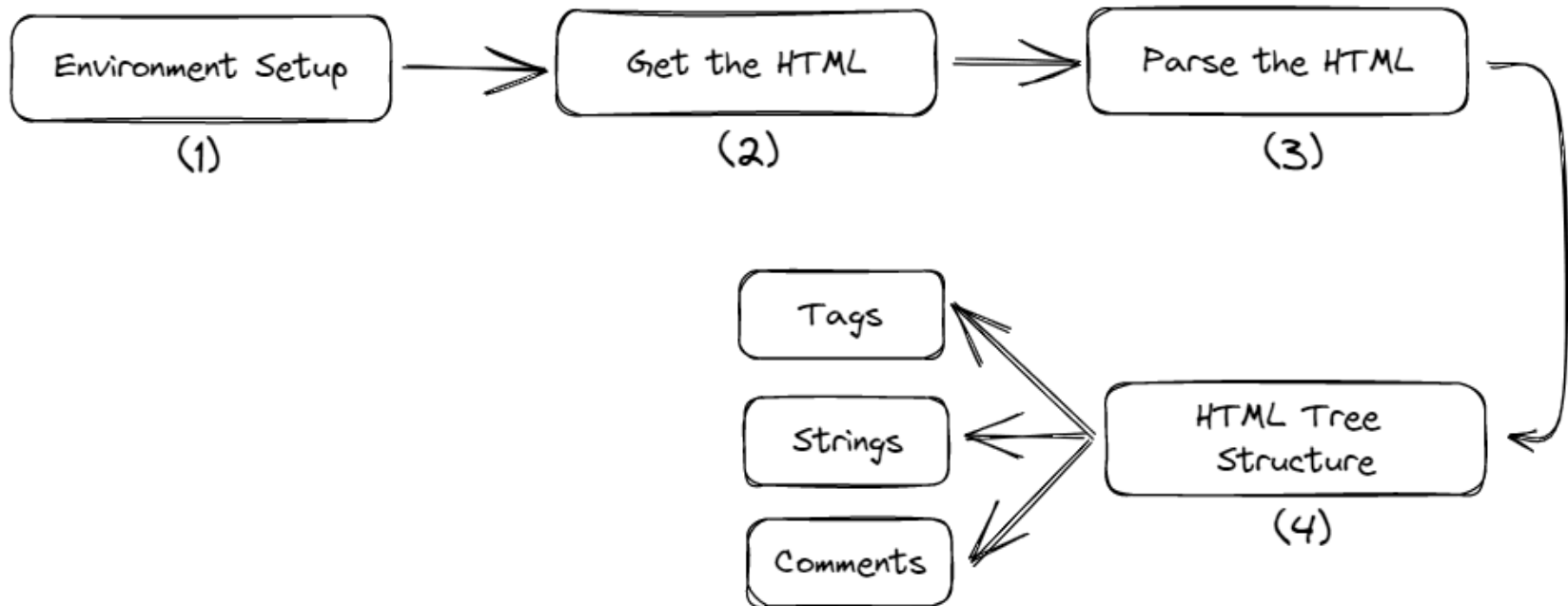
2. Content Creation: Web scraping can gather information from multiple sources like news articles, research reports, and blog posts. It helps the creator to create quality and trending content.

3. Price Comparison: Web scraping can be used to extract the prices of a particular product across multiple e-commerce websites to give a fair price comparison for the user. It also helps companies fix the optimal pricing of their products to compete with their competitors.

Use of Web Scraping

4. Job Postings: Web Scraping can also be used to collect data on various job openings across multiple job portals so that this information can help many job seekers and recruiters.

Code Implementation



Workflow and Libraries

Libraries: Requests and BeautifulSoup

The Requests library :

- is used to make HTTP requests in Python. It is an easy-to-use library with a lot of features. We will be using it to get the content from the URL of the webpage that we wish to scrap
- Send Request and Load the webpage.
(Requests, urllib, httpplib)

BeautifulSoup:

- is a library that is used to parse HTML documents. We create a new BeautifulSoup object by passing the HTML content and the type of parser we want to use

Setting Up the Environment

1. Installing the Libraries

\$ pip install requests

\$ pip install bs4



```
from bs4 import BeautifulSoup
import requests
import pandas as pd
```

Setting Up the Environment

The requests module extracts the HTML content from a URL. It extracts all the data in a raw format as a string that needs further processing.

The bs4 is the Beautiful Soup module. It will parse the raw HTML content obtained from the `request` module in a well-structured format.

Implementation

2. Understanding the Webpage structure

The most important step in data scraping is to understand the structure of the webpage and the structure of the data we want to scrap

In order to see the HTML code for a website, do a “Right-Click” and click on “Inspect Source”

Example

https://www.imdb.com/chart/moviemeter/

IMDb Menu All Search IMDb

























IMDb Charts

Most Popular Movies

As determined by IMDb Users

Showing 100 Titles

Sort by: Ranking

Rank & Title	IMDb Rating	Your Rating
 Venom: Let There Be Carnage (2021) 1 ( 191)		 
 The Mitchells vs the Machines (2021) 2 ( 1)		 
 Those Who Wish Me Dead (2021) 3 ( 23)		 
 Mortal Kombat (2021) 4 ( 1)		 
 The Woman in the Window (2021) 5 ( 181)		 
 Wrath of Man (2021) 6 ( 3)		 

Back Forward Reload Save as... Print... Cast... Send to Samsung Phone Create QR code for this page Translate to English View page source Inspect

Alt+Left Arrow Alt+Right Arrow Ctrl+R Ctrl+S Ctrl+P Ctrl+U Ctrl+Shift+I

You Have Seen

0/100 (0%)

☐ Hide titles I've seen

IMDb Charts

- Box Office
- Most Popular Movies
- Top Rated Movies
- Top Rated English Movies
- Most Popular TV
- Top Rated TV
- Lowest Rated Movies

Top India Charts

- Top Rated Indian Movies
- Top Rated Malayalam Movies
- Top Rated Tamil Movies
- Top Rated Telugu Movies

Popular Movies by Genre

- Action
- Adventure
- Animation
- Biography

Example

- We are going to scrap the most popular movies from IMDb.
- This shows us the HTML code for the web page. Now we wish to scrap the following.
 - 1.The rank of the movie
 - 2.Name of the movie
 - 3.IMDb rating of the movie
- Ongoing through the HTML code, we see that this information is stored in the form of a table using the <table> tag

HTML Example






← → ↻ https://www.imdb.com/chart/moviemeter/ ☆ C T ⓘ ⚙ S ⋮

IMDb Charts

Most Popular Movies

As determined by IMDb Users

table.chart.full-width 625.6 × 6991.2

Rank & Title
 Venom: Let There Be Carnage (2021) 1 (📈 191)
 The Mitchells vs the Machines (2021) 2 (📉 1)
 Those Who Wish Me Dead (2021) 3 (📈 23)
 Mortal Kombat (2021) 4 (📉 1)
 The Woman in the Window (2021) 5 (📈 181)

```
<script type="text/javascript">it(typeof uet === 'function'){uet('bb','ChartWidget',{wb:1});}</script>
<span class="ab_widget">
  <input id="seen-config" type="hidden" data-caller="chtmvm" data-pagetype="chart" data-subpagetype="moviemeter" data-baseref="chtmvm">
  <div class="seen-collection" data-collectionid="moviemeter">
    <div class="article">
      <h3>IMDb Charts</h3>
      <div class="chart-social-sharing-widget" id="social-sharing-widget">...</div>
      <h1 class="header">Most Popular Movies</h1>
      <div class="byline">As determined by IMDb Users</div>
      <hr>
      <div class="listner">
        <div>...</div>
        <br class="clear">
        <table class="chart full-width" data-caller-name="chart-moviemeter"> == $0
          <colgroup>...</colgroup>
          <thead>...</thead>
          <tbody class="listner-list">
            <tr>
              <td class="posterColumn">...</td>
              <td class="titleColumn">...</td>
              <td class="ratingColumn imdbRating">
                </td>
              <td class="ratingColumn">...</td>
              <td class="watchlistColumn">...</td>
            </tr>
            <tr>
              <td class="posterColumn">...</td>
              <td class="titleColumn">...</td>
              <td class="ratingColumn imdbRating">...</td>
              <td class="ratingColumn">...</td>
              <td class="watchlistColumn">...</td>
            </tr>
            <tr>...</tr>
            <tr>...</tr>
            <tr>...</tr>
            <tr>...</tr>
            <tr>...</tr>
          </tbody>
        </table>
      </div>
    </div>
  </div>
</span>
```

div.pagecontent div#content-2-wide div#main div.article span.ab_widget div.seen-collection div.article div.listner table.chart.full-width ...

Console What's New

HTML structure

3. Understanding HTML structure and tags

All the data on a webpage is stored inside different tags based on the content and presentation. The various important tags are :

1. <div> tag: It is used to define a section of the webpage. For example, you can observe that the table which we want to scrap is present inside a <div> tag with class =” lister” and this is further inside a <div> tag with class = “article”

HTML structure

2. <a> tag: This is used to create a hyperlink to another webpage. It consists of the “href” attribute which specifies the URL of the webpage it links to. For example, when you click on a word and it directs you to another webpage on Wikipedia, it is because the word is present inside a tag along with the href attribute containing the URL of the web page you are directed to

3. tag: This is used to create unordered lists on the webpage. Each entry of the list is marked with a tag

4. tag: This is used to create ordered lists on the webpage. Similar to the tag, each entry is marked with a tag

HTML structure

5.<table> tag: This is used to create a table on the webpage. For example, we can see on the IMDb page, the information we require is present in a tabular form.

6.<tr> tag: This is used to specify the different rows under the table. For example, the different movies on the IMDb page are present in different rows and hence, in different <tr> tags

7.<td> tag: This is used to specify the content for each column for a given row in a table. For example, the content for the different columns (“Rank and Title”, “IMDb Rating”, “My Rating”) for each row is specified under the <td> tag

Implementation

4. Creating a BeautifulSoup object



```
url="https://www.imdb.com/chart/moviemeter/"  
res=requests.get(url)  
soup=BeautifulSoup(res.text,"html.parser")  
print(soup)
```

Example

This creates an object named “soup” which has the HTML code for the URL given and can be used to select certain sections of the data

[illegible]

Implementation

5. Extracting the data using BeautifulSoup

```
▼ <table class="chart full-width" data-caller-name="chart-moviemeter"> == $0
  ▶ <colgroup>...</colgroup>
  ▶ <thead>...</thead>
  ▼ <tbody class="lister-list">
    ▼ <tr>
      ▶ <td class="posterColumn">...</td>
      ▶ <td class="titleColumn">...</td>
      <td class="ratingColumn imdbRating">
        </td>
      ▶ <td class="ratingColumn">...</td>
      ▶ <td class="watchlistColumn">...</td>
    </tr>
    ▼ <tr>
      ▶ <td class="posterColumn">...</td>
      ▶ <td class="titleColumn">...</td>
      ▶ <td class="ratingColumn imdbRating">...</td>
      ▶ <td class="ratingColumn">...</td>
      ▶ <td class="watchlistColumn">...</td>
    </tr>
    ▶ <tr>...</tr>
    ▶ <tr>...</tr>
    ▶ <tr>...</tr>
    ▶ <tr>...</tr>
    ▶ <tr>...</tr>
    ▶ <tr>...</tr>
    ▶ <tr>...</tr>
    ▶ <tr>...</tr>
    ▶ <tr>...</tr>
    ▶ <tr>...</tr>
    ▶ <tr>...</tr>
    ▶ <tr>...</tr>
```

The information regarding the movies is stored inside a <table> tag with the class attribute as “chart full-width”

So we will use the “find_all” function to find all such tables in the soup object



```
for i in soup.find_all("table",class_="chart full-width"):
    print(i)
```

Within the table, the data for each movie is stored in a `<tr>` tag so we will loop through these movies and store their data in a dataframe

The `<td>` tag containing the name and rank of the title has the class = “titleColumn”. The name of the movie is stored in a tag since it connects it to the webpage of that particular movie while the rank is stored inside a `<div>` tag with the class attribute as “velocity”

Also, the <td> tag containing the IMDb rating of the movie has the class attribute as “ratingColumn imdbRating”

```
▼<tr>
  ▶<td class="posterColumn">...</td>
  ▼<td class="titleColumn">
    <a href="/title/tt7097896/?pf_rd_m=A2FGELUUNOQJNL&pf_rd_p=ea4e08e1-c8a3-47b5-ac3a...QCPHR5C&p:
    Carnage</a>
    <span class="secondaryInfo">(2021)</span>
    ▼<div class="velocity">
      "1
      "
      ▶<span class="secondaryInfo">...</span>
    </div>
  </td>
  <td class="ratingColumn imdbRating">
    </td>
  ▶<td class="ratingColumn">...</td>
  ▶<td class="watchlistColumn">...</td>
```

We can extract the text within a tag using the following way:

```
▶ for i in soup.find_all("table",class_="chart full-width"):  
    for j in i.find_all("tr"):  
        for k in j.find_all("td",class_="titleColumn"):  
            for l in k.find_all("a"):  
                print("Name of the movie is ",l.text)
```

```
Name of the movie is  Venom: Let There Be Carnage  
Name of the movie is  The Mitchells vs the Machines  
Name of the movie is  Those Who Wish Me Dead  
Name of the movie is  Mortal Kombat  
Name of the movie is  The Woman in the Window  
Name of the movie is  Wrath of Man  
Name of the movie is  Spiral  
Name of the movie is  Oxygen  
Name of the movie is  Tom Clancy's Without Remorse  
Name of the movie is  Tenet  
Name of the movie is  Things Heard & Seen  
Name of the movie is  Army of the Dead  
Name of the movie is  Nomadland  
Name of the movie is  Nobody  
Name of the movie is  The Green Knight  
Name of the movie is  Radhe  
Name of the movie is  A Quiet Place Part II  
Name of the movie is  Promising Young Woman
```