



# Extracting Stock Data Using a Web Scraping

Not all stock data is available via the API in this assignment; you will use web-scraping to obtain financial data. You will be quizzed on your results.

You will extract and share historical data from a web page using the BeautifulSoup library.

## Table of Contents

### 1. Extracting data using BeautifulSoup

- Download the web page Using Requests Library
- Parse HTML on a web page using BeautifulSoup
- Extract data and build a data frame

### 2. Extracting data using pandas

### 3. Exercise

Estimated Time Needed: **30 min**

```
In [ ]: !pip install pandas
        !pip install requests
```

```
!pip install bs4
!pip install html5lib
!pip install lxml
!pip install plotly
```

```
In [54]: import pandas as pd
import requests
from bs4 import BeautifulSoup
```

In Python, you can ignore warnings using the warnings module. You can use the filterwarnings function to filter or ignore specific warning messages or categories.

```
In [55]: import warnings
# Ignore all warnings
warnings.filterwarnings("ignore", category=FutureWarning)
```

## Using Webscraping to Extract Stock Data Example

We will extract Netflix stock data [https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMDeveloperSkillsNetwork-PY0220EN-SkillsNetwork/labs/project/netflix\\_data\\_webpage.html](https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMDeveloperSkillsNetwork-PY0220EN-SkillsNetwork/labs/project/netflix_data_webpage.html).

In this example, we are using yahoo finance website and looking to extract Netflix data.

Date	Open	High	Low	Close*	Adj Close**	Volume
Jun 01, 2021	504.01	536.13	482.14	528.21	528.21	78,560,600
May 01, 2021	512.65	518.95	478.54	502.81	502.81	66,927,600
Apr 01, 2021	529.93	563.56	499.00	513.47	513.47	111,573,300
Mar 01, 2021	545.57	556.99	492.85	521.66	521.66	90,183,900
Feb 01, 2021	536.79	566.65	518.28	538.85	538.85	61,902,300
Jan 01, 2021	539.00	593.29	485.67	532.39	532.39	139,988,600
Dec 01, 2020	492.34	545.50	491.29	540.73	540.73	77,564,100
Nov 01, 2020	478.87	518.73	463.41	490.70	490.70	91,788,900
Oct 01, 2020	506.03	572.49	472.21	475.74	475.74	154,302,400
Sep 01, 2020	532.60	557.39	458.60	500.03	500.03	118,796,900

Fig:- Table that we need to extract

On the following web page we have a table with columns name (Date, Open, High, Low, close, adj close volume) out of which we must extract following columns

- Date
- Open
- High
- Low
- Close

- Volume

## Steps for extracting the data

1. Send an HTTP request to the web page using the requests library.
2. Parse the HTML content of the web page using BeautifulSoup.
3. Identify the HTML tags that contain the data you want to extract.
4. Use BeautifulSoup methods to extract the data from the HTML tags.
5. Print the extracted data

### Step 1: Send an HTTP request to the web page

You will use the request library for sending an HTTP request to the web page.

```
In [56]: url = "https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMDeveloperSkillsNetwork-PY0220EN-SkillsNetwork/lab
```

The `requests.get()` method takes a URL as its first argument, which specifies the location of the resource to be retrieved. In this case, the value of the `url` variable is passed as the argument to the `requests.get()` method, because you will store a web page URL in a `url` variable.

You use the `.text` method for extracting the HTML content as a string in order to make it readable.

```
In [ ]: data = requests.get(url).text  
print(data)
```

### Step 2: Parse the HTML content

---

---

## What is parsing?

In simple words, parsing refers to the process of analyzing a string of text or a data structure, usually following a set of rules or grammar, to understand its structure and meaning. Parsing involves breaking down a piece of text or data into its individual components or elements, and then analyzing those components to extract the desired information or to understand their relationships and meanings.

---

Next you will take the raw HTML content of a web page or a string of HTML code which needs to be parsed and transformed into a structured, hierarchical format that can be more easily analyzed and manipulated in Python. This can be done using a Python library called **Beautiful Soup**.

## Parsing the data using the BeautifulSoup library

- Create a new BeautifulSoup object.

**Note:** To create a BeautifulSoup object in Python, you need to pass two arguments to its constructor:

1. The HTML or XML content that you want to parse as a string.
2. The name of the parser that you want to use to parse the HTML or XML content. This argument is optional, and if you don't specify a parser, BeautifulSoup will use the default HTML parser included with the library. here in this lab we are using "html5lib" parser.

```
In [58]: soup = BeautifulSoup(data, 'html.parser')
```

### Step 3: Identify the HTML tags

As stated above, the web page consists of a table so, we will scrape the content of the HTML web page and convert the table into a data frame.

You will create an empty data frame using the **pd.DataFrame()** function with the following columns:

- "Date"
- "Open"
- "High"

- "Low"
- "Close"
- "Volume"

```
In [59]: netflix_data = pd.DataFrame(columns=["Date", "Open", "High", "Low", "Close", "Adj Close", "Volume"])
netflix_data
```

```
Out[59]:
```

Date	Open	High	Low	Close	Adj Close	Volume
------	------	------	-----	-------	-----------	--------

## Working on HTML table

These are the following tags which are used while creating HTML tables.

- `<table>`: This tag is a root tag used to define the start and end of the table. All the content of the table is enclosed within these tags.
- `<tr>`: This tag is used to define a table row. Each row of the table is defined within this tag.
- `<td>`: This tag is used to define a table cell. Each cell of the table is defined within this tag. You can specify the content of the cell between the opening and closing tags.
- `<th>`: This tag is used to define a header cell in the table. The header cell is used to describe the contents of a column or row. By default, the text inside a tag is bold and centered.
- `<tbody>`: This is the main content of the table, which is defined using the tag. It contains one or more rows of elements.

## Step 4: Use a BeautifulSoup method for extracting data

We will use **find()** and **find\_all()** methods of the BeautifulSoup object to locate the table body and table row respectively in the HTML.

- The *find()* method will return particular tag content.

- The `find_all()` method returns a list of all matching tags in the HTML.

```
In [60]: # First we isolate the body of the table which contains all the information
# Then we loop through each row and find all the column values for each row
for row in soup.find("tbody").find_all('tr'):
    col = row.find_all("td")
    date = col[0].text
    Open = col[1].text
    high = col[2].text
    low = col[3].text
    close = col[4].text
    adj_close = col[5].text
    volume = col[6].text

    # Finally we append the data of each row to the table
    netflix_data = pd.concat([netflix_data, pd.DataFrame({"Date": [date], "Open": [Open], "High": [high],
                                                         "Low": [low], "Close": [close], "Adj Close": [adj_close],
                                                         "Volume": [volume]})], ignore_index=True)
```

## Step 5: Print the extracted data

We can now print out the data frame using the `head()` or `tail()` function.

```
In [61]: netflix_data.head()
```

```
Out[61]:
```

	Date	Open	High	Low	Close	Adj Close	Volume
0	Jun 01, 2021	504.01	536.13	482.14	528.21	528.21	78,560,600
1	May 01, 2021	512.65	518.95	478.54	502.81	502.81	66,927,600
2	Apr 01, 2021	529.93	563.56	499.00	513.47	513.47	111,573,300
3	Mar 01, 2021	545.57	556.99	492.85	521.66	521.66	90,183,900
4	Feb 01, 2021	536.79	566.65	518.28	538.85	538.85	61,902,300

# Extracting data using pandas library

We can also use the pandas `read_html` function from the pandas library and use the URL for extracting data.

## What is read\_html in pandas library?

`pd.read_html(url)` is a function provided by the pandas library in Python that is used to extract tables from HTML web pages. It takes in a URL as input and returns a list of all the tables found on the web page.

```
In [62]: read_html_pandas_data = pd.read_html(url)
```

Or you can convert the BeautifulSoup object to a string.

```
In [ ]: read_html_pandas_data = pd.read_html(str(soup))
```

Because there is only one table on the page, just take the first table in the returned list.

```
In [63]: netflix_dataframe = read_html_pandas_data[0]  
  
netflix_dataframe.head()
```

```
Out[63]:
```

	Date	Open	High	Low	Close*	Adj Close**	Volume
0	Jun 01, 2021	504.01	536.13	482.14	528.21	528.21	78560600
1	May 01, 2021	512.65	518.95	478.54	502.81	502.81	66927600
2	Apr 01, 2021	529.93	563.56	499.00	513.47	513.47	111573300
3	Mar 01, 2021	545.57	556.99	492.85	521.66	521.66	90183900
4	Feb 01, 2021	536.79	566.65	518.28	538.85	538.85	61902300



# Exercise: use webscraping to extract stock data

Use the `requests` library to download the webpage [https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMDeveloperSkillsNetwork-PY0220EN-SkillsNetwork/labs/project/amazon\\_data\\_webpage.html](https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMDeveloperSkillsNetwork-PY0220EN-SkillsNetwork/labs/project/amazon_data_webpage.html). Save the text of the response as a variable named `html_data`.

```
In [ ]: amazon_url = "https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMDeveloperSkillsNetwork-PY0220EN-SkillsNetw
html_data = requests.get(amazon_url).text
print(html_data)
```

Parse the html data using `beautiful_soup` using `html.parser`.

```
In [68]: amazon_soup = BeautifulSoup(html_data, 'html.parser')
```

**Question 1:** What is the content of the title attribute?

```
In [69]: print(amazon_soup.title)
print(amazon_soup.title.text)
```

```
<title>Amazon.com, Inc. (AMZN) Stock Historical Prices & Data - Yahoo Finance</title>
Amazon.com, Inc. (AMZN) Stock Historical Prices & Data - Yahoo Finance
```

Using BeautifulSoup, extract the table with historical share prices and store it into a data frame named `amazon_data`. The data frame should have columns Date, Open, High, Low, Close, Adj Close, and Volume. Fill in each variable with the correct data from the list `col`.

```
In [74]: amazon_data = pd.DataFrame(columns=["Date", "Open", "High", "Low", "Close", "Adj Close", "Volume"])

for row in amazon_soup.find("tbody").find_all("tr"):
    col = row.find_all("td")
    date = col[0].text
    Open = col[1].text
    high = col[2].text
    low = col[3].text
    close = col[4].text
    adj_close = col[5].text
    volume = col[6].text
```

```
amazon_data = pd.concat([amazon_data, pd.DataFrame({"Date": [date], "Open": [Open], "High": [high], "Low": [low],
                                                    "Close": [close], "Adj Close": [adj_close], "Volume": [volume]})],
                        ignore_index=True)
```

Print out the first five rows of the `amazon_data` data frame you created.

```
In [75]: amazon_data.head()
```

```
Out[75]:
```

	Date	Open	High	Low	Close	Adj Close	Volume
0	Jan 01, 2021	3,270.00	3,363.89	3,086.00	3,206.20	3,206.20	71,528,900
1	Dec 01, 2020	3,188.50	3,350.65	3,072.82	3,256.93	3,256.93	77,556,200
2	Nov 01, 2020	3,061.74	3,366.80	2,950.12	3,168.04	3,168.04	90,810,500
3	Oct 01, 2020	3,208.00	3,496.24	3,019.00	3,036.15	3,036.15	116,226,100
4	Sep 01, 2020	3,489.58	3,552.25	2,871.00	3,148.73	3,148.73	115,899,300

**Question 2:** What are the names of the columns in the data frame?

```
In [76]: print(list(amazon_data.columns.values))
print(list(amazon_data))
```

```
['Date', 'Open', 'High', 'Low', 'Close', 'Adj Close', 'Volume']
['Date', 'Open', 'High', 'Low', 'Close', 'Adj Close', 'Volume']
```

**Question 3:** What is the `Open` of the last row of the `amazon_data` data frame?

```
In [77]: print(amazon_data.tail())
print(amazon_data.iloc[-1,1])
```

	Date	Open	High	Low	Close	Adj Close	Volume
56	May 01, 2016	663.92	724.23	656.00	722.79	722.79	90,614,500
57	Apr 01, 2016	590.49	669.98	585.25	659.59	659.59	78,464,200
58	Mar 01, 2016	556.29	603.24	538.58	593.64	593.64	94,009,500
59	Feb 01, 2016	578.15	581.80	474.00	552.52	552.52	124,144,800
60	Jan 01, 2016	656.29	657.72	547.18	587.00	587.00	130,200,900
		656.29					

## About the Authors:

[Joseph Santarcangelo](#) has a PhD in Electrical Engineering, his research focused on using machine learning, signal processing, and computer vision to determine how videos impact human cognition. Joseph has been working for IBM since he completed his PhD.

Azim Hirjani  
Akansha yadav

## Change Log

Date (YYYY-MM-DD)	Version	Changed By	Change Description
02-05-2023	1.3	Akansha yadav	Updated Lab content under maintenance
2021-06-09	1.2	Lakshmi Holla	Added URL in question 3
2020-11-10	1.1	Malika Singla	Deleted the Optional part
2020-08-27	1.0	Malika Singla	Added lab to GitLab

© IBM Corporation 2020. All rights reserved.