

Web Scrapping



What is our GOAL for this MODULE?

The goal of this module is to learn scraping to get the useful data from different sites.

What did we ACHIEVE in the class TODAY?

- We scraped the stars data from Nasa's site

Which CONCEPTS/CODING BLOCKS did we cover today?

- Usage of selenium
- Usage of BeautifulSoup
- Getting data from HTML of a page

How did we DO the activities?

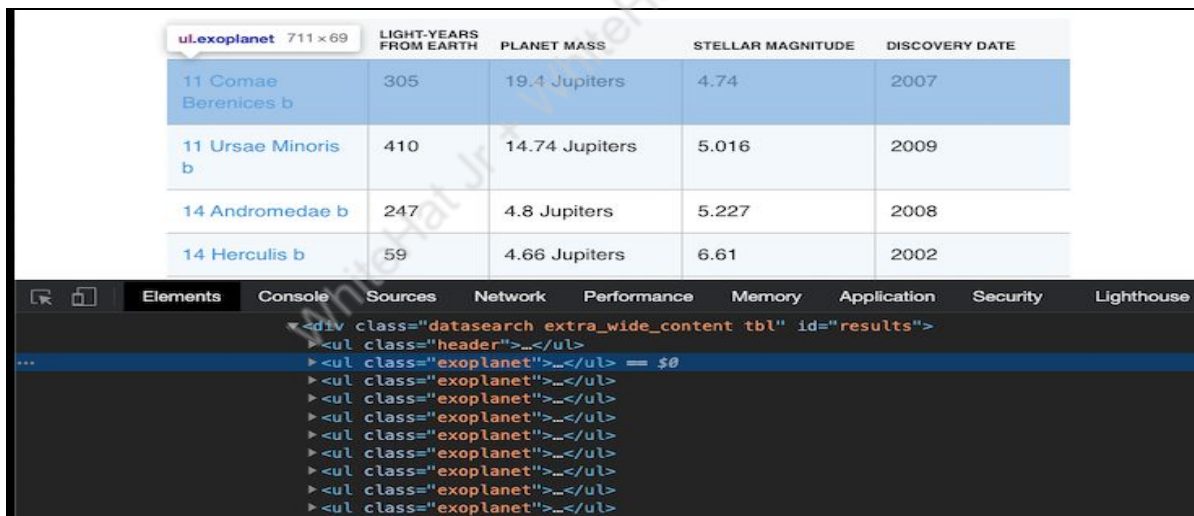
1. We created a virtual environment.
2. We installed the necessary libraries.
3. We imported the selenium and BeautifulSoup in our code.

```
from selenium import webdriver
from bs4 import BeautifulSoup
import time
import csv
```

4. We also installed the chrome driver.
5. Then we opened the link we want to scrape using selenium.

```
START_URL = "https://exoplanets.nasa.gov/exoplanet-catalog/"
browser = webdriver.Chrome("/path/to/chromedriver")
browser.get(START_URL)
time.sleep(10)
```

6. We got all the tags from the table to get the data from it.



ul.exoplanet 711 x 69	LIGHT-YEARS FROM EARTH	PLANET MASS	STELLAR MAGNITUDE	DISCOVERY DATE
11 Comae Berenices b	305	19.4 Jupiters	4.74	2007
11 Ursae Minoris b	410	14.74 Jupiters	5.016	2009
14 Andromedae b	247	4.8 Jupiters	5.227	2008
14 Herculis b	59	4.66 Jupiters	6.61	2002


```
<div class="datasearch extra_wide_content tbl" id="results">
  <ul class="header"></ul>
  <ul class="exoplanet"></ul> = $0
  <ul class="exoplanet"></ul>
  <ul class="exoplanet"></ul>
  <ul class="exoplanet"></ul>
  <ul class="exoplanet"></ul>
  <ul class="exoplanet"></ul>
  <ul class="exoplanet"></ul>
  <ul class="exoplanet"></ul>
  <ul class="exoplanet"></ul>
```

7. We found all the ul tags to get the data.

```
soup = BeautifulSoup(browser.page_source, "html.parser")
for ul_tag in soup.find_all("ul", attrs={"class", "exoplanet"}):
```

8. Then we found all the li tags.

```
li_tags = ul_tag.find_all("li")
```

9. We found that all the data is inside the anchor tags.

```
temp_list = []
for index, li_tag in enumerate(li_tags):
    if index == 0:
        temp_list.append(li_tag.find_all("a")[0].contents[0])
    else:
        try:
            temp_list.append(li_tag.contents[0])
        except:
            temp_list.append("")
planet_data.append(temp_list)
```

10. We code to loop the code 428 times.

```

scrapper_2.py
1  from selenium import webdriver
2  from bs4 import BeautifulSoup
3  import time
4  import csv
5
6  START_URL = "https://exoplanets.nasa.gov/exoplanet-catalog/"
7  browser = webdriver.Chrome("/Users/apoorvelous/Downloads/chromedriver")
8  browser.get(START_URL)
9  time.sleep(10)
10
11 def scrape():
12     headers = ["name", "light_years_from_earth", "planet_mass", "stellar_magnitude", "discovery_date"]
13     planet_data = []
14     for i in range(0, 428):
15         soup = BeautifulSoup(browser.page_source, "html.parser")
16         for ul_tag in soup.find_all("ul", attrs={"class", "exoplanet"}):
17             li_tags = ul_tag.find_all("li")
18             temp_list = []
19             for index, li_tag in enumerate(li_tags):
20                 if index == 0:
21                     temp_list.append(li_tag.find_all("a")[0].contents[0])
22                 else:
23                     try:
24                         temp_list.append(li_tag.contents[0])
25                     except:
26                         temp_list.append("")
27             planet_data.append(temp_list)
28             browser.find_element_by_xpath('//*[@id="primary_column"]/footer/div/div/div/nav/span[2]/a').click()
29             with open("scrapper_2.csv", "w") as f:
30                 csvwriter = csv.writer(f)
31                 csvwriter.writerow(headers)
32                 csvwriter.writerows(planet_data)
33
34     scrape()
  
```

What's NEXT?

In the next class, we will explore more of scraping and data cleaning.

EXTEND YOUR KNOWLEDGE:

You can read the following blog on scraping with selenium to understand more:

<https://medium.com/ymedialabs-innovation/web-scraping-using-beautiful-soup-and-selenium-for-dynamic-page-2f8ad15efe25>