

Topic	Web Scraping 2	
Class Description	Students would be reworking the previously written code to scrape more data.	
Class	C128	
Class time	45 mins	
Goal	<ul style="list-style-type: none"> Scrape more data about all the exoplanets 	
Resources Required	<ul style="list-style-type: none"> Teacher Resources <ul style="list-style-type: none"> Laptop with internet connectivity Earphones with mic Notebook and pen Student Resources <ul style="list-style-type: none"> Laptop with internet connectivity Earphones with mic Notebook and pen 	
Class structure	Warm Up Teacher-led Activity Student-led Activity Wrap up	5 mins 15 min 15 min 5 min
<div> <div></div> <div>CONTEXT</div> <ul style="list-style-type: none"> Review the concepts learned in the earlier classes </div>		
Class Steps	Teacher Action	Student Action
Step 1: Warm Up (5 mins)	Hi <Student Name>! In the last class, we scraped exoplanet's data from NASA's website. Can you recall all the tools that we used in the last class?	ESR: - Selenium - BeautifulSoup

	<p>Great! Now, in today's class, we will scrape some more data from the same website. We got some data like distance from earth, planet size, etc. but today we will scrape more data so that when we perform analysis later, we can better predict the planets, for instance, to see if they are likely habitable, etc.</p> <p>Are you excited?</p>	<p>ESR: "Yes!"</p>
	<p>Before we start I have an exciting quiz question for you! Are you ready to answer this question?</p> <div data-bbox="730 945 954 1033" data-label="Image"> </div> <p>Teacher click on the button on the bottom right corner of your screen to start the In-Class Quiz.</p> <p>A quiz will be visible to both you and the student.</p> <p>Encourage the student to answer the quiz question.</p> <p>The student may choose the wrong option, help the student to think correctly about the question and then answer again.</p> <p>After the student selects the correct option, the <div data-bbox="599 1652 836 1715" data-label="Image"></div> button will start appearing on your screen.</p> <p>Click the End quiz to close the quiz pop-up and continue the class.</p>	<p>ESR: Yes!</p>

	Let's get started!	
Teacher Initiates Screen Share		
<p style="text-align: center;"><u>CHALLENGE</u></p> <ul style="list-style-type: none"> Scraping more data from the website and letting students lead the development this time. 		
Step 2: Teacher-led Activity (15 min)	<p>Teacher opens the same website that we scraped in the last class.</p> <p><Teacher opens the link from Teacher activity 1></p> <p>https://exoplanets.nasa.gov/exoplanet-catalog/</p>	
	<p>Let's look at this page again.</p> <p>Here, if we look closely, we can see that the name of these exo-planets is a hyperlink.</p> <p>Note: NASA's exoplanet catalog web page keeps updating as per the new planet discoveries. At the time of writing this document, the web page had 428 with 10 Planets per page showing a total of planets data 4280.</p> <p>By default it can have 25 planets per page.</p>	

< 1 of 428 >

Per page 10 ▾

NAME ↑	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
16 Cygni B b	69	1.78 Jupiters	6.25	1996
18 Delphini b	249	10.3 Jupiters	5.506	2008
1RXS J160929.1-210524 b	473	8 Jupiters	12.057	2008
24 Bootis b	314	0.91 Jupiters	5.58	2018
24 Sextantis b	236	1.99 Jupiters	6.441	2010
24 Sextantis c	236	0.86 Jupiters	6.441	2010

< 1 of 428 >

[Back to top](#)

Let's click on the link and see what kind of data we can find?

<div><div>PLANET TYPE</div><div>Gas Giant</div></div>		<div><div>DISCOVERY DATE</div><div>2007</div></div>
<div><div>MASS</div><div>19.4 Jupiters</div></div>		<div><div>PLANET RADIUS</div><div>1.08 x Jupiter (estimate)</div></div>
<div><div>ORBITAL RADIUS</div><div>1.29 AU</div></div>		<div><div>ORBITAL PERIOD</div><div>326 days</div></div>
<div><div>ECCENTRICITY</div><div>0.23</div></div>		<div><div>DETECTION METHOD</div><div>Radial Velocity</div></div>

	<p>Great! Now, let's say we want to scrape this data as well. Can you tell me what's the first change that we'll have to make in our previous code?</p>	<p>ESR:</p> <p>We need to save the hyperlink's href in our CSV.</p>
	<p>That's great! Let's get started.</p> <p>We will add a new column in our header. Our header variable would now look like this:</p> <pre>headers = ["name", "light_years_from_earth", "planet_mass", "stellar_magnitude", "discovery_date", "hyperlink"]</pre>	

```

updated_scraper.py > scrape
1 ~ from selenium import webdriver
2 from selenium.webdriver.common.by import By
3 from bs4 import BeautifulSoup
4 import time
5 import csv
6
7 # NASA Exoplanet URL
8 START_URL = "https://exoplanets.nasa.gov/exoplanet-catalog/"
9
10 # Webdriver
11 browser = webdriver.Edge("C:/Whitehat_jr/PRO-127-130/msedgedriver.exe")
12 browser.get(START_URL)
13
14 time.sleep(10)
15
16 planets_data = []
17
18 headers = ["name", "light_years_from_earth", "planet_mass", "stellar_magnitude", "discovery_date", "hyperlink"]
19

```

We have added an extra hyperlink into our header list. Now, we also need to add this into the temp_list variable list, before we append into the planets_data.

Before we do that, let's investigate the href url in these hyperlinks:

	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
16 Cygni B b	69	1.78 Jupiters	6.25	1996

```

Elements Console Sources Network Performance Memory Application Security Lighthouse
<div>
  <div class="datasearch extra_wide_content tbl" id="results">
    <ul class="header"></ul>
    <ul class="exoplanet">
      <li>
        <a href="/exoplanet-catalog/6988/11-comae-berenices-b/">11 Comae Berenices b</a>
      </li>
    </ul>
  </div>

```

Here, we can see that these links do not have <https://exoplanets.nasa.gov> before them. We will have to add them.

Now to achieve this, we will do the following:

Note: Scrape 4 pages only. Thus in the for loop range is given from 1 to 5. If each page has data of 25 planets, total 100 hyperlinks will be scraped in the later section.

```
def scrape():
    for i in range(1,5):
        while True:
            time.sleep(2)

            soup = BeautifulSoup(browser.page_source, "html.parser")

            for ul_tag in soup.find_all("ul", attrs={"class", "exoplanet"}):
                li_tags = ul_tag.find_all("li")
                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("")

                # Get Hyperlink Tag
                hyperlink_li_tag = li_tags[0]

                temp_list.append("https://exoplanets.nasa.gov"+ hyperlink_li_tag.find_all("a", href=True)[0]["href"])

                planets_data.append(temp_list)

            browser.find_element(By.XPATH, value='//*[@id="primary_column"]/footer/div/div/div/nav/span[2]/a').click()

            print(f"Page {i} scraping completed")

# Calling Method
scrape()
```

We have added:

```
hyperlink_li_tag = li_tags[0]

temp_list.append("https://exoplanets.
nasa.gov"+hyperlink_li_tag.find_all("
a", href=True)[0]["href"])
```

Here, first we are creating a variable `hyperlink_li_tag` and then we are using this variable to find all the anchor tag with href, take the first anchor tag (since we know there's only one anchor tag in all li tags) and then we are taking out the href from it.

	<p>Now that we have the links in planet_data, can you tell me what should be our next steps?</p> <p>Perfect, we will create a new function that will take these hyperlinks one by one, get the HTML and then we will scrape the data.</p>	<p>ESR:</p> <ul style="list-style-type: none"> - We'll scrape data by using these links!
	<p>To make sure that we are scraping pages one by one, we'll add a code in the scrape() function to check the current page number.</p>	

```
def scrape():
    for i in range(1,5):
        while True:
            time.sleep(2)

            soup = BeautifulSoup(browser.page_source, "html.parser")

            # Check page number
            current_page_num = int(soup.find_all("input", attrs={"class", "page_num"})[0].get("value"))

            if current_page_num < i:
                browser.find_element(By.XPATH, value='//*[@id="primary_column"]/footer/div/div/div/nav/span[2]/a').click()
            elif current_page_num > i:
                browser.find_element(By.XPATH, value='//*[@id="primary_column"]/footer/div/div/div/nav/span[1]/a').click()
            else:
                break
```

	<p>Earlier, we used selenium because we wanted to click a button on the page (next button) but this time, we do not want to interact with the browser, therefore we can do this without selenium.</p> <p>Let's get started!</p>	
--	---	--

Teacher Stops Screen Share

	<p>Now it's your turn. Please share your screen with me.</p>	
--	--	--

- Ask Student to press ESC key to come back to panel
- Guide Student to start Screen Share
- Teacher gets into Fullscreen

ACTIVITY

- Student creates a new function to use all the hyperlinks one by one and scrape data from there

Step 3: Student-Led Activity (15 min)

Ask the student to move the variables **headers** and **planets_data** to the global scope, i.e, below `time.sleep(10)` line.

This is because we now would want to access these variables in multiple functions.

Let's add the new headers, that is, the new data that is available on the new page we just discovered.

The student moves the variables.

```

1  from selenium import webdriver
2  from selenium.webdriver.common.by import By
3  from bs4 import BeautifulSoup
4  import time
5  import pandas as pd
6  import requests
7  import csv
8
9  # NASA Exoplanet URL
10 START_URL = "https://exoplanets.nasa.gov/exoplanet-catalog/"
11
12 # Webdriver
13 browser = webdriver.Edge("C:/Whitehat_jr/PRO-127-130/msedgedriver.exe")
14 browser.get(START_URL)
15
16 time.sleep(10)
17
18 planets_data = []
19
20 headers = ["name", "light_years_from_earth", "planet_mass", "stellar_magnitude", "discovery_date", "hyperlink",
21            "planet_type", "planet_radius", "orbital_radius", "orbital_period", "eccentricity"]
22

```

	<pre>headers = ["name", "light_years_from_earth", "planet_mass", "stellar_magnitude", "discovery_date", "hyperlink", "planet_type", "planet_radius", "orbital_radius", "orbital_period", "eccentricity"]</pre>	<i>The student adds more headers.</i>
	<p>Great, now let's create a new function and call that function. We will call the function in loop and pass the hyperlink we saved with the earlier function into this function.</p> <p>Also, let's comment out the CSV saving code. We want to save a csv with half the data, right?</p>	<i>The student creates a new function.</i>
<pre>31 # with open("scrapper_2.csv", "w") as f: 32 # csvwriter = csv.writer(f) 33 # csvwriter.writerow(headers) 34 # csvwriter.writerows(planet_data) 35 36 def scrape_more_data(hyperlink): 37 pass 38 39 scrape() 40 for data in planet_data: 41 scrape_more_data(data[5])</pre>		

	<p>Okay, now earlier, we created a soup object where we passed the browser's page source and parsed it as html. This time, since we are not going to use selenium, how can we do it?</p>	<p>ESR:</p> <p>We can get the page's HTML by making a GET request.</p>
	<p>That's right! For that, we will import requests module</p> <pre>import requests</pre> <p>And we will write the following code inside the new function we created:</p> <pre>page = requests.get(hyperlink) soup = BeautifulSoup(page.content, "html.parser")</pre> <p>Here, we are first getting the page, and then we are parsing the contents of the page as HTML.</p>	<p>The student follows instructions.</p>
	<p><i>Ask the student to create a new list new_planets_data to save data from these new pages, and ask them to scrape the data like before.</i></p> <p><i>Help the student if required.</i></p> <p><i>The code should look something like this:</i></p>	

```
new_planets_data = []

def scrape_more_data(hyperlink):
    try:
        page = requests.get(hyperlink)

        soup = BeautifulSoup(page.content, "html.parser")

        temp_list = []

        for tr_tag in soup.find_all("tr", attrs={"class": "fact_row"}):
            td_tags = tr_tag.find_all("td")

            for td_tag in td_tags:
                try:
                    temp_list.append(td_tag.find_all("div", attrs={"class": "value"})[0].contents[0])
                except:
                    temp_list.append("")

            new_planets_data.append(temp_list)

    except:
        time.sleep(1)
        scrape_more_data(hyperlink)
```

Let's call the method and check the list new_planets_data.

```
for index, data in enumerate(planets_data):
    scrape_more_data(data[5])
    print(f"scraping at {index+1} is completed.")

print(new_planets_data[0:10])
```

Save the code and run using virtual environment.

	<pre> Page 1 scraping completed Page 2 scraping completed Page 3 scraping completed Page 4 scraping completed scraping at hyperlink 1 is completed. scraping at hyperlink 2 is completed. scraping at hyperlink 3 is completed. scraping at hyperlink 4 is completed. scraping at hyperlink 5 is completed. scraping at hyperlink 6 is completed. scraping at hyperlink 7 is completed. scraping at hyperlink 8 is completed. </pre>	
	<p>Since we are running it for only four pages, it will start scraping the data from respective hyperlink.</p>	
<pre> [['\nGas Giant\n', '\n2007\n', '\n19.4 Jupiters\n', '\n1.08 x Jupiter', '\n1.29 AU\n', '\n326 days\n', '\n', '\nGas Giant\n', '\n2009\n', '\n14.74 Jupiters\n', '\n1.09 x Jupiter', '\n1.53 AU\n', '\n1.4 years\n', '\n'], ['\nGas Giant\n', '\n2008\n', '\n4.8 Jupiters\n', '\n1.15 x Jupiter', '\n0.83 AU\n', '\n185.8 days\n', '\n'], ['\nGas Giant\n', '\n2002\n', '\n4.66 Jupiters\n', '\n1.15 x Jupiter', '\n2.93 AU\n', '\n4.9 years\n', '\n'], ['\nGas Giant\n', '\n1996\n', '\n1.78 Jupiters\n', '\n1.2 x Jupiter', '\n1.66 AU\n', '\n2.2 years\n', '\n'], ['\nGas Giant\n', '\n2020\n', '\n4.32 Jupiters\n', '\n1.15 x Jupiter', '\n1.45 AU\n', '\n1.6 years\n', '\n'], ['\nGas Giant\n', '\n2008\n', '\n10.3 Jupiters\n', '\n1.11 x Jupiter', '\n2.6 AU\n', '\n2.7 years\n', '\n'], ['\nGas Giant\n', '\n2008\n', '\n8 Jupiters\n', '\n1.664 x Jupiter\n', '\n330.0 AU\n', '\n650 days\n', '\n'], ['\nGas Giant\n', '\n2018\n', '\n0.91 Jupiters\n', '\n1.24 x Jupiter', '\n0.19 AU\n', '\n0.04\n', '\n'], ['\nGas Giant\n', '\n2010\n', '\n1.99 Jupiters\n', '\n1.19 x Jupiter', '\n1.333 AU\n', '\n', '\n0.09\n', '\n']] </pre>		
	<p>Great job! Now we have 2 lists, planets_data and new_planets_data.</p> <p>In new_planets_data, a special character '\n' is present. we need to remove it before saving it to csv file.</p> <p>Also, we want to merge the two lists. Adding 2 lists creates 1 final list with elements from both the lists in the same order.</p>	<p><i>The student merges the data.</i></p>

```
final_planet_data = []

for index, data in enumerate(planets_data):
    new_planet_data_element = new_planets_data[index]
    new_planet_data_element = [elem.replace("\n", "") for elem in new_planet_data_element]
    new_planet_data_element = new_planet_data_element[:7]
    final_planet_data.append(data + new_planet_data_element)
```

Finally, we will create a csv with our headers and **final_planet_data**.

The student creates a CSV.

```
with open("final.csv", "w") as f:
    csvwriter = csv.writer(f)
    csvwriter.writerow(headers)
    csvwriter.writerows(final_planet_data)
```

Let's run this code to see if it works fine and generates the **final.csv**.

Student runs the code.

```
1  name,light_years_from_earth,planet_mass,stellar_magnitude,discovery_date,hyperlink,planet_type,planet_radius,orbital_radius,orbital
2
3  11 Comae Berenices b,304,19.4 Jupiters,4.72307,2007,https://exoplanets.nasa.gov/exoplanet-catalog/6988/11-comae-berenices-b/,Gas Gi
4
5  11 Ursae Minoris b,409,14.74 Jupiters,5.013,2009,https://exoplanets.nasa.gov/exoplanet-catalog/6989/11-ursae-minoris-b/,Gas Giant,2
6
7  14 Andromedae b,246,4.8 Jupiters,5.23133,2008,https://exoplanets.nasa.gov/exoplanet-catalog/6990/14-andromedae-b/,Gas Giant,2008,4.
8
9  14 Herculis b,58,4.66 Jupiters,6.61935,2002,https://exoplanets.nasa.gov/exoplanet-catalog/6991/14-herculis-b/,Gas Giant,2002,4.66
10
11 16 Cygni B b,69,1.78 Jupiters,6.215,1996,https://exoplanets.nasa.gov/exoplanet-catalog/6992/16-cygni-b-b/,Gas Giant,1996,1.78 Jupit
12
13 17 Scorpii b,408,4.32 Jupiters,5.22606,2020,https://exoplanets.nasa.gov/exoplanet-catalog/8016/17-scorpii-b/,Gas Giant,2020,4.32 Ju
14
15 18 Delphini b,249,10.3 Jupiters,5.51048,2008,https://exoplanets.nasa.gov/exoplanet-catalog/6993/18-delphini-b/,Gas Giant,2008,10.3
16
17 1RXS J160929.1-210524 b,454,8 Jupiters,12.618,2008,https://exoplanets.nasa.gov/exoplanet-catalog/7061/1rxs-j1609291-210524-b/,Gas G
18
19 24 Bootis b,313,0.91 Jupiters,5.59,2018,https://exoplanets.nasa.gov/exoplanet-catalog/7274/24-bootis-b/,Gas Giant,2018,0.91 Jupiter
20
21 24 Sextantis b,235,1.99 Jupiters,6.4535,2010,https://exoplanets.nasa.gov/exoplanet-catalog/6994/24-sextantis-b/,Gas Giant,2010,1.99
22
23 24 Sextantis c,235,0.86 Jupiters,6.4535,2010,https://exoplanets.nasa.gov/exoplanet-catalog/6995/24-sextantis-c/,Gas Giant,2010,0.86
```

	<p>Although we have only scraped 4 pages, there are approximately 200 pages with 25 planets on each page. The scraping can take a lot of time sometimes (like for scraping 5000 planets data in this case) therefore we'll provide you the final.csv with the data of all the exoplanets.</p> <p><Student Activity 1></p> <p>https://github.com/procodingclass/PRO-129-Datasets</p> <p>If you want you can also try running your code after the class to check the output</p>	<p>Student runs the code after class to get the output or downloads the csv from Student Activity 1</p>
Teacher Guides Student to Stop Screen Share		
<p><u>FEEDBACK</u></p> <ul style="list-style-type: none"> • Appreciate the student for their efforts • Identify 2 strengths and 1 area of progress for the student 		
Step 4: Wrap-Up (5 min)	<p>So, in this project class we revisited the concepts from the previous class and you did the majority of the scraping yourself! Congratulations!</p>	ESR: Thanks!
	<p>Next class, we will be learning new concepts and building new projects.</p>	-
<p>Teacher Clicks ✕ End Class</p>		

Activity	Activity Name	Links
Teacher activity 1	solution	https://github.com/procodingclass/P-RO-C128-RefCode
Student Activity 1	final csv	https://github.com/procodingclass/P-RO-129-Datasets

