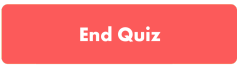


<b>Topic</b>	<b>Web Scraping 2</b>	
<b>Class Description</b>	Students would be reworking the previously written code to scrape more data.	
<b>Class</b>	<b>C128</b>	
<b>Class time</b>	<b>45 mins</b>	
<b>Goal</b>	<ul style="list-style-type: none"> <li>Scrape more data about all the exoplanets</li> </ul>	
<b>Resources Required</b>	<ul style="list-style-type: none"> <li>Teacher Resources               <ul style="list-style-type: none"> <li>Laptop with internet connectivity</li> <li>Earphones with mic</li> <li>Notebook and pen</li> </ul> </li> <li>Student Resources               <ul style="list-style-type: none"> <li>Laptop with internet connectivity</li> <li>Earphones with mic</li> <li>Notebook and pen</li> </ul> </li> </ul>	
<b>Class structure</b>	<b>Warm Up</b> <b>Teacher-led Activity</b> <b>Student-led Activity</b> <b>Wrap up</b>	<b>5 mins</b> <b>15 min</b> <b>15 min</b> <b>5 min</b>
<div> <div></div> <div>CONTEXT</div> <ul style="list-style-type: none"> <li>Review the concepts learned in the earlier classes</li> </ul> </div>		
<b>Class Steps</b>	<b>Teacher Action</b>	<b>Student Action</b>
<b>Step 1: Warm Up (5 mins)</b>	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?	<b>ESR:</b> - 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><b>ESR:</b> "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="735 951 958 1035" 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  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><b>ESR:</b> Yes!</p>

	Let's get started!	
<b>Teacher Initiates Screen Share</b>		
<p align="center"><b><u>CHALLENGE</u></b></p> <ul style="list-style-type: none"> <li>Scraping more data from the website and letting students lead the development this time.</li> </ul>		
<b>Step 2: Teacher-led Activity (15 min)</b>	<p><i>Teacher opens the same website that we scraped in the last class.</i></p> <p><i>&lt;Teacher opens the link from Teacher activity 1&gt;</i></p> <p><a href="https://exoplanets.nasa.gov/exoplanet-catalog/">https://exoplanets.nasa.gov/exoplanet-catalog/</a></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>	

&lt; 1 of 428 &gt;

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

&lt; 1 of 428 &gt;

[Back to top](#)

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

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

	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?	<b>ESR:</b> We need to save the hyperlink's href in our CSV.
	<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> <p>We have added an extra hyperlink into our header list. Now, we also</p>	

need to add this into the temp\_list variable list, before we append into the planet\_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>
    </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:

```

scraper_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", "hyperlink"]
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             hyperlink_li_tag = li_tags[0]
28             temp_list.append("https://exoplanets.nasa.gov"+hyperlink_li_tag.find_all("a", href=True)[0]["href"])
29             planet_data.append(temp_list)
30             browser.find_element_by_xpath('//*[@id="primary_column"]/footer/div/div/div/nav/span[2]/a').click()
31         with open("scraper_2.csv", "w") as f:
32             csvwriter = csv.writer(f)
33             csvwriter.writerow(headers)
34             csvwriter.writerows(planet_data)
35
36     scrape()
37

```

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.

**ESR:**

- We'll scrape data by using these links!



	Now that we have the links in planet_data, can you tell me what should be our next steps?	
	<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>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>	
<b>Teacher Stops Screen Share</b>		
	Now it's your turn. Please share your screen with me.	
<ul style="list-style-type: none"> <li>• <b>Ask Student to press ESC key to come back to panel</b></li> <li>• <b>Guide Student to start Screen Share</b></li> <li>• <b>Teacher gets into Fullscreen</b></li> </ul>		
<p style="text-align: center;"><b>ACTIVITY</b></p> <ul style="list-style-type: none"> <li>• <b>Student creates a new function to use all the hyperlinks one by one and scrape data from there</b></li> </ul>		
<b>Step 3: Student-Led Activity (15 min)</b>	<p><i>Ask the student to move the variables <b>headers</b> and <b>planet_data</b> to the global scope, i.e, below <code>time.sleep(10)</code> line.</i></p> <p>This is because we now would want to access these variables in multiple functions.</p>	<i>The student moves the variables.</i>



```

scraper_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 headers = ["name", "light_years_from_earth", "planet_mass", "stellar_magnitude", "discovery_date", "hyperlink"]
11 planet_data = []

```

Great! Now, let's add the new headers, that is, the new data that is available on the new page we just discovered.

```

headers = ["name",
"light_years_from_earth",
"planet_mass", "stellar_magnitude",
"discovery_date", "hyperlink",
"planet_type", "planet_radius",
"orbital_radius", "orbital_period",
"eccentricity"]

```

*The student adds more headers.*

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.

Also, let's comment out the CSV saving code. We want to save a csv with half the data, right?

*The student creates a new function.*

```

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])

```

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?

**ESR:**  
We can get the page's HTML by making a GET request.

That's right! For that, we will import requests module

```
import requests
```

And we will write the following code inside the new function we created:

```

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

```

Here, we are first getting the page, and then we are parsing the contents of the page as HTML.

The student follows instructions.

	<p>Ask the student to create a new list <b>new_planet_data</b> to save data from these new pages, and ask them to scrape the data like before.</p> <p>Help the student if required. The code should look something like this:</p>	
	<pre>def scrape_more_data(hyperlink):     page = requests.get(hyperlink)     soup = BeautifulSoup(page.content, "html.parser")     for tr_tag in soup.find_all("tr", attrs={"class": "fact_row"}):         td_tags = tr_tag.find_all("td")         temp_list = []         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_planet_data.append(temp_list)</pre>	
	<p>Great job! Now we have 2 lists, <b>planet_data</b> and <b>new_planet_data</b>.</p> <p>What we want to do is, we want to merge this data. Adding 2 lists creates 1 final list with elements from both the lists in the same order.</p>	<p>The student merges the data.</p>
	<pre>final_planet_data = []  for index, data in enumerate(planet_data):     final_planet_data.append(data + final_planet_data[index])</pre>	

Finally, we will create a csv with our **headers** and **final\_planet\_data**.

Our final code looks something like this:

*The student creates a CSV.*

```
scraper_2.py
1  from selenium import webdriver
2  from bs4 import BeautifulSoup
3  import requests
4  import time
5  import csv
6
7  START_URL = "https://exoplanets.nasa.gov/exoplanet-catalog/"
8  browser = webdriver.Chrome("/Users/apoorvelous/Downloads/chromedriver")
9  browser.get(START_URL)
10 time.sleep(10)
11 headers = ["name", "light_years_from_earth", "planet_mass", "stellar_magnitude", "discovery_date", "hyperlink", "planet_type", "planet_r"]
12 planet_data = []
13 new_planet_data = []
14
15 def scrape():
16     for i in range(0, 428):
17         soup = BeautifulSoup(browser.page_source, "html.parser")
18         for ul_tag in soup.find_all("ul", attrs={"class", "exoplanet"}):
19             li_tags = ul_tag.find_all("li")
20             temp_list = []
21             for index, li_tag in enumerate(li_tags):
22                 if index == 0:
23                     temp_list.append(li_tag.find_all("a")[0].contents[0])
24                 else:
25                     try:
26                         temp_list.append(li_tag.contents[0])
27                     except:
28                         temp_list.append("")
29             hyperlink_li_tag = li_tags[0]
30             temp_list.append("https://exoplanets.nasa.gov"+hyperlink_li_tag.find_all("a", href=True)[0]["href"])
31             planet_data.append(temp_list)
32             browser.find_element_by_xpath('//*[id="primary_column"]/footer/div/div/div/nav/span[2]/a').click()
```

```

34 def scrape_more_data(hyperlink):
35     page = requests.get(hyperlink)
36     soup = BeautifulSoup(page.content, "html.parser")
37     for tr_tag in soup.find_all("tr", attrs={"class": "fact_row"}):
38         td_tags = tr_tag.find_all("td")
39         temp_list = []
40         for td_tag in td_tags:
41             try:
42                 temp_list.append(td_tag.find_all("div", attrs={"class": "value"})[0].contents[0])
43             except:
44                 temp_list.append("")
45         new_planet_data.append(temp_list)
46
47 scrape()
48 for data in planet_data:
49     scrape_more_data(data[5])
50
51 final_planet_data = []
52
53 for index, data in enumerate(planet_data):
54     final_planet_data.append(data + final_planet_data[index])
55
56 with open("final.csv", "w") as f:
57     csvwriter = csv.writer(f)
58     csvwriter.writerow(headers)
59     csvwriter.writerows(final_planet_data)

```

Let's run this code to see if it works fine and generates the desired result.

*Student runs the code.*

Although it is the running version of the code, scraping can take a lot of time sometimes (like for scraping 4,277 pages in this case) therefore we'll provide you the final csv.  
<Student Activity 1>

If you want you can also try running your code after the class to check the output

Student runs the code after class to get the output or downloads the csv from Student Activity 1

Teacher Guides Student to Stop Screen Share		
<b>FEEDBACK</b> <ul style="list-style-type: none"> <li>• Appreciate the student for their efforts</li> <li>• Identify 2 strengths and 1 area of progress for the student</li> </ul>		
<b>Step 4:</b> <b>Wrap-Up</b> <b>(5 min)</b>	So, in this project class we revisited the concepts from the previous class and you did the majority of the scraping yourself! Congratulations!	<b>ESR:</b> Thanks!
	Next class, we will be learning new concepts and building new projects.	-
<div> <div>Teacher Clicks</div> <div>✕ End Class</div> </div>		

Activity	Activity Name	Links
Teacher activity 1	solution	<a href="https://github.com/whitehatjr/web-scraping-2">https://github.com/whitehatjr/web-scraping-2</a>
Student Activity 1	final csv	<a href="https://raw.githubusercontent.com/whitehatjr/web-scraping-2/master/final.csv">https://raw.githubusercontent.com/whitehatjr/web-scraping-2/master/final.csv</a>