

## Advanced Web Development and Web Scraping Spring 2022, Assignment #7 – Web Scraping Assignment

**Note:** For all assignments, *get* the web page using python's request library.

1. Scrape the sample Schedule page (<https://gdancik.github.io/CSC-301/data/notes/schedule.html>) to output the number of courses being taught and the *total* number of credits, in the following format:

```
Dr. Dancik is teaching 3 courses (9 credits)
```

Note: Your scraper should work for different data (e.g., a different instructor, or a different number of courses) and you must scrape/calculate the following:

- the instructor's name
- the number of courses
- the total number of credits

Note: In order to add the credits, you will need to convert each 'number' to an integer, using the *int* function, e.g., `int('3')` will return 3.

2. Scrape the title and rating for 5 movies from IMDB, whose links are given below, and construct a bar graph that shows the rating for each movie. Give your graph an informative title. The following links should be used:

- [https://www.imdb.com/title/tt0109830/?ref=fn\\_al\\_tt\\_1](https://www.imdb.com/title/tt0109830/?ref=fn_al_tt_1)
- [https://www.imdb.com/title/tt0076759/?ref=fn\\_tt\\_tt\\_1](https://www.imdb.com/title/tt0076759/?ref=fn_tt_tt_1)
- [https://www.imdb.com/title/tt0368226/?ref=mv\\_sr\\_2](https://www.imdb.com/title/tt0368226/?ref=mv_sr_2)
- Select another movie from IMDB and include the URL
- Select another movie from IMDB and include the URL

In order to do this, you should create a list containing the URLs and iterate through the list to scrape the relevant information for each movie. Note that all pages have the same format. **After submitting a request to a page, you must sleep for 1 second so that you do not abuse IMDB's servers.** This is done by importing the *time* module and calling `time.sleep(seconds)`:

```
import time
time.sleep(1)
```

Notes: (1) The rating will need to be converted to a float (decimal) using the float function (e.g., `float("3.1")` will return the number 3.1); (2) The title strings may contain a `'\xa0'`, which is code for a non-breaking space. These do not have to be removed, but if you want to remove them, you can use python's *strip* method.

- Scrape the job listing site *Indeed* (<https://www.indeed.com/>) to find jobs for a search of your choice. You will need to copy the URL for your search and use it in your Python script. Note that because *Indeed* uses the GET method to retrieve information, user input is visible in the URL. For example, a search for “computer programmer” in “Hartford, CT” has the following URL:  
<https://www.indeed.com/jobs?q=computer+programmer&l=Hartford%2C+CT>

Your web scraper should create a *pandas* data frame that contains the following information.

	Title	Company	Location	Salary
0	Entry Level Computer Programmer	Revature	Hartford, CT+2 locations	?
1	Systems Analyst/Programmer RPG/AS400	National Waste Associates	Glastonbury, CT 06033	\$70,000 - \$90,000 a year
2	Lead Web Programmer experienced with client-si...	Solution Innovators	South Windsor, CT 06074	\$80,000 - \$100,000 a year
3	Service/Help Desk Technician (Information Tech...	State of Connecticut - Department of Public He...	Hartford, CT 06106 (Frog Hollow area)	\$63,126 - \$80,821 a year
4	Information Technology Analyst 1 (40 Hour) (So...	State of Connecticut - Department of Revenue...	Hartford, CT 06103 (Downtown area)	\$72,145 - \$92,373 a year
5	Systems Analyst	Connecticut Children's Medical Center	Hartford, CT 06106	?

You therefore will need to extract the job title, company, location, and salary from each job listing. Not all job listings have salary information; if they do not, you should display a question mark (?).

Notes:

- in Jupyter Notebook, a dollar sign (\$) in a pandas data frame is interpreted as denoting a mathematical expression, which changes the formatting. In order to turn this formatting off, run the following statement after importing the pandas module: `pd.options.display.html.use_mathjax = False`
- Some job titles are preceded by the word ‘new’. However, the job titles themselves are found in *span* elements that have a *title* attribute. The recommended way to get these elements is to use a function that takes an element as an input and returns *True* if the element has a *title* attribute. An example can be found here:  
<https://www.crummy.com/software/BeautifulSoup/bs4/doc/#a-function>