

CS104 PROJECT REPORT :

Python Web Crawler

Vishal Bysani
22B1061

14 June 2023

Introduction

I am Vishal Bysani. I have made a web crawler using Python for the CS104 project.

A web crawler, crawler or web spider, is a computer program that's used to search and automatically index website content and other information over the internet. These programs, or bots, are most commonly used to create entries for a search engine index.

Libraries

The libraries and dependencies used in this project are:

- **argparse** : This library helps in parsing the command line arguments[1]
- **BeautifulSoup** : BeautifulSoup is a Python library for getting data out of HTML, XML, and other markup languages.
- **requests** : Python requests is a library for making HTTP requests. It provides an easy-to-use interface that makes working with HTTP very simple, which means it simplifies the process of sending and receiving data from websites by providing a uniform interface for both GET and POST methods.
- **Numpy and Matplotlib** : I used matplotlib to generate the plots of the files distribution at various levels recursed by the web crawler[2]
- **warnings** : I used this library to not display unnecessary warnings.

Usage

The web crawler is run using the command:

```
python3 Web_Crawler.py -u <site-name> -t <threshold> -o <output-filename>
-s <user-input>
```

- -u: for the URL. This is a compulsory argument and must print an error if not provided
- -t: for the threshold of recursiveness. It must be greater than zero. If not provided, the web crawler will recurse till the end. In case of invalid (negative values), it must print an error
- -o: For an output file. If not provided then by default print on the command line.
- -s: If the user gives 'Y' argument, then the web crawler computes and prints the files sizes as well. It only prints the links if anything else is given as argument.

Code Structure

Function of base code

Using the base code from [3], we can scrape all urls from the given site and prints them all. However, it doesn't crawl into these urls again. Figure 1 was taken from [3]

```
http://example.webscraping.com///places/default/user/register?_next=/places/default/index
http://example.webscraping.com///places/default/user/register?_next=/places/default/index/places/default/user/regi
ster
http://example.webscraping.com///places/default/user/register?_next=/places/default/index/places/default/user/regi
ster/places/default/user/register
http://example.webscraping.com///places/default/user/register?_next=/places/default/index/places/default/user/regi
ster/places/default/user/register/places/default/user/register
http://example.webscraping.com///places/default/user/register?_next=/places/default/index/places/default/user/regi
ster/places/default/user/register/places/default/user/register/places/default/user/register
http://example.webscraping.com///places/default/user/register?_next=/places/default/index/places/default/user/regi
ster/places/default/user/register/places/default/user/register/places/default/user/register/places/default/user/re
gister
http://example.webscraping.com///places/default/user/register?_next=/places/default/index/places/default/user/regi
ster/places/default/user/register/places/default/user/register/places/default/user/register/places/default/user/re
gister/places/default/user/register
http://example.webscraping.com///places/default/user/register?_next=/places/default/index/places/default/user/regi
ster/places/default/user/register/places/default/user/register/places/default/user/register/places/default/user/re
gister/places/default/user/register/places/default/user/register
http://example.webscraping.com///places/default/user/register?_next=/places/default/index/places/default/user/regi
ster/places/default/user/register/places/default/user/register/places/default/user/register/places/default/user/re
gister/places/default/user/register/places/default/user/register/places/default/user/register
```

Figure 1: Output of a basic Web crawler

Implementation of modified code

I have modified the code from [3] in the following manner to meet the requirements of the problem statement:

- The code first extracts all the links in the *href* and *src* tags using BeautifulSoup and then modifies them to make valid urls.
- The code recursively crawls through the url provided by the user upto a specified threshold or until the end if not specified.
- The code also segregates all the files at a particular recursion level into various categories like HTML, CSS, JS, JPG, PNG and Others.

Customization

- The code keeps printing the number of files and the url crawled by it to indicate the progress of the program.²

```

http://www.w3schools.com/lib/w3schools32.css
247]
http://www.w3schools.com/js/tryit.asp?filename=tryjs_default
248]
http://www.w3schools.com/images/w3lynx_200.png
249]
250]
http://www.w3schools.com/lib/fonts/freckle-face-v9-latin-regular.woff2
251]
http://www.w3schools.com/images/colorpicker.png
252]
http://www.w3schools.com/r/r_exercises.asp
253]
http://www.w3schools.com/python/scipy/scipy_quiz.php
254]
http://www.w3schools.com/cssref/default.asp
255]
http://www.w3schools.com/aws/index.php
256]
http://www.w3schools.com/bootstrap/bootstrap_exercises.asp
257]
http://www.w3schools.com/cpp/cpp_exercises.asp
258]
http://www.w3schools.com/bootcamp/bootcamp_sql.php
259]
http://www.w3schools.com/w3css/w3css_references.asp
260]
http://www.w3schools.com/python/pandas/pandas_exercises.asp
261]
http://www.w3schools.com/bootcamp/index.php
262]
263]
http://www.w3schools.com/colors/colors_picker.asp
264]
265]
http://www.w3schools.com/php/php_examples.asp
266]
http://www.w3schools.com/nodejs/default.asp
267]

```

Figure 2: Indication of running of the program

- The modified code also identifies if a link is of internal or external domain, and selectively crawls only the internal ones. It separately counts the number of internal and external domain files at each level and prints them in the end.
- It also plots bar graphs showing the various categories of files in internal and external domains at each recursion level and saves this plot in Files.Distribution.png.

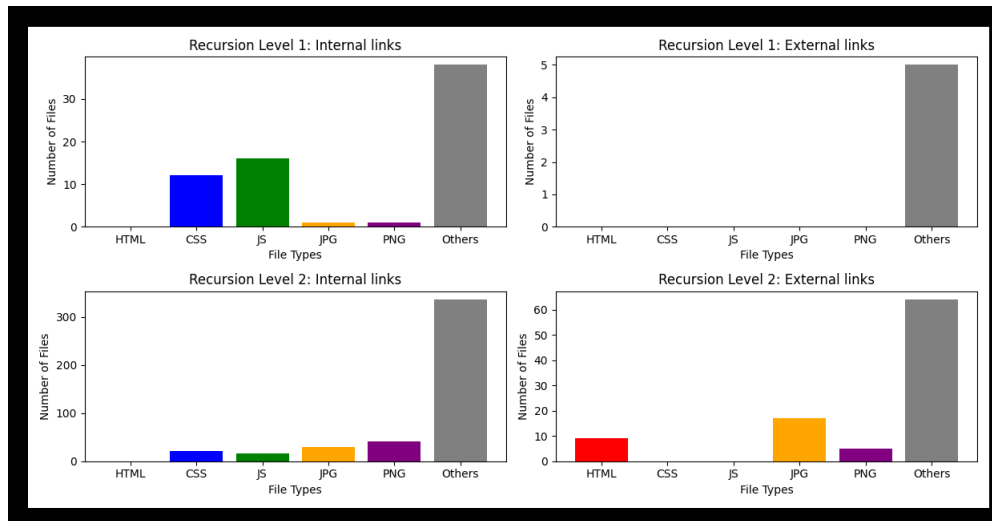


Figure 3: Bar Graph obtained from modified code

- The code also finds file sizes based upon the argument given by the user to the `s` tag.

```

At recursion level: 1
Total number of files in this level: 73
Of these files, the various types are:
INTERNAL LINKS: 68
HTML: 0

CSS: 12
('http://mahitgadhiwala.com/wp-content/plugins/header-footer-elementor/inc/widgets-css/frontend.css?ver=1.6.4', 7645)
('http://mahitgadhiwala.com/wp-content/plugins/elementor/assets/lib/animations/animations.min.css?ver=3.4.6', 2572)
('http://mahitgadhiwala.com/wp-content/plugins/elementor/assets/lib/font-awesome/css/solid.min.css?ver=5.15.3', 311)
('http://mahitgadhiwala.com/wp-content/themes/astra/assets/css/minified/main.min.css?ver=3.7.2', 8241)
('http://mahitgadhiwala.com/wp-content/plugins/elementor/assets/lib/font-awesome/css/brands.min.css?ver=5.15.3', 312)
('http://mahitgadhiwala.com/wp-content/plugins/elementor/assets/css/frontend.min.css?ver=3.4.6', 17535)
('http://mahitgadhiwala.com/wp-content/plugins/elementor/assets/lib/font-awesome/css/fontawesome.min.css?ver=5.15.3', 12392)
('http://mahitgadhiwala.com/wp-content/uploads/elementor/css/post-56.css?ver=1667884167', 3363)
('http://mahitgadhiwala.com/wp-content/plugins/elementor/assets/lib/elcons/css/elementor-icons.min.css?ver=5.13.0', 3812)
('http://mahitgadhiwala.com/wp-content/themes/astra/assets/css/minified/block-library/style.min.css?ver=5.8.7', 10421)
('http://mahitgadhiwala.com/wp-content/uploads/elementor/css/post-55.css?ver=1635158278', 330)
('http://mahitgadhiwala.com/wp-content/plugins/header-footer-elementor/assets/css/header-footer-elementor.css?ver=1.6.4', 322)

JavaScript: 16
('http://mahitgadhiwala.com/wp-includes/js/wp-embed.min.js?ver=5.8.7', 805)
('http://mahitgadhiwala.com/wp-includes/js/wp-util.min.js?ver=5.8.7', 710)
('http://mahitgadhiwala.com/wp-content/plugins/elementor/assets/lib/dialog/dialog.min.js?ver=4.8.1', 4479)
('http://mahitgadhiwala.com/wp-content/plugins/elementor/assets/lib/jquery/jquery-migrate.min.js?ver=3.3.2', 4166)
('http://mahitgadhiwala.com/wp-content/plugins/elementor/assets/js/preloaded-modules.min.js?ver=3.4.6', 10016)
('http://mahitgadhiwala.com/wp-includes/js/underscore.min.js?ver=1.13.1', 7325)
('http://mahitgadhiwala.com/wp-content/plugins/wpforms-lite/assets/js/integrations/elementor/frontend.min.js?ver=1.7.0', 375)
('http://mahitgadhiwala.com/wp-includes/js/jquery/ui/core.min.js?ver=1.12.1', 6871)
('http://mahitgadhiwala.com/wp-content/plugins/elementor/assets/lib/waypoints/waypoints.min.js?ver=4.0.2', 3954)
('http://mahitgadhiwala.com/wp-content/plugins/elementor/assets/js/webpack.runtime.min.js?ver=3.4.6', 3089)
('http://mahitgadhiwala.com/wp-content/plugins/elementor/assets/lib/swiper/swiper.min.js?ver=5.3.6', 36405)
('http://mahitgadhiwala.com/wp-content/plugins/elementor/assets/lib/share-link/share-link.min.js?ver=3.4.6', 2052)
('http://mahitgadhiwala.com/wp-content/themes/astra/assets/js/minified/frontend.min.js?ver=3.7.2', 3818)
('http://mahitgadhiwala.com/wp-includes/js/jquery/jquery.min.js?ver=3.6.0', 30915)
('http://mahitgadhiwala.com/wp-content/plugins/elementor/assets/js/frontend.min.js?ver=3.4.6', 11544)
('http://mahitgadhiwala.com/wp-content/plugins/elementor/assets/js/frontend-modules.min.js?ver=3.4.6', 5476)

JPG: 1
('http://mahitgadhiwala.com/wp-content/uploads/elementor/thumbs/Screenshot-88-e1632570445250-pdmsp6sol1f8rx7ccpdlrkzlpixltmg0xr6o126o1

PNG: 1
('http://mahitgadhiwala.com/wp-content/uploads/2021/09/Screenshot__92_-removebg-preview-e1632571908657.png', 154389)
Others: 38

```

Figure 4: Output redirected to the output file

As it can be seen from figures 3 and 4, the modified web crawler categorizes the files at various recursion levels. Further the web-crawler also finds the file sizes of all the files and prints it along with the links.

Description of the Submitted Folder

I have submitted all the files and sub-folder under the folder named **22B1061_project**. It includes the following:

- Web.Crawler.py: It contains the python code.
- report.pdf: It contains the report in the pdf format.
- report.tex: It contains the latex code for the pdf.
- report.bib: It contains all the references I used for this project
- BarGraphs.png: It shows the bar graphs it prints for various recursion levels.
- OutputFilePrintFormat.png: It shows the format in which it prints into the output file given by the user

- Screenshot.jpeg: It shows the output of code which I took as a base for my web crawler.

References

- [1] URL: <https://docs.python.org/3/library/argparse.html>.
- [2] URL: https://matplotlib.org/stable/gallery/subplots_axes_and_figures/subplots_demo.html.
- [3] URL: <https://www.geeksforgeeks.org/python-program-to-recursively-scrape-all-the-urls-of-the-website/>.