

Web Crawler

Code and Documentation Written By
Steven Fan

Overview of System

The web crawler is written in Python 3 and mainly uses the Scrapy architecture for web crawling. The Windows operating system is utilized in building and testing this program. The following documentation will document the source code itself.

Import Code Blocks

```
import sys
```

Obtain user parameters passed in through the command line

```
import os
```

Creation and manipulation of files and directories

```
import scrapy
```

Web crawling

```
from scrapy.spiders import CrawlSpider, Rule
```

Utilization of a specific spider type: *CrawlSpider* and allows specifications of rules in the behavior of the spider

```
from scrapy.linkextractors import LinkExtractor
```

Extraction of web links from HTML to traverse the web

```
from scrapy.crawler import CrawlerProcess
```

Sets up the environment and process for the spider to function in

```
from scrapy.utils.project import get_project_settings
```

Extraction of spider settings to manipulate functionality of the spider

```
from scrapy.exceptions import CloseSpider
```

Shutting down of the spider, especially due to failed file writes or page limit exceeded

crawler Class Code Blocks

```
rules = (Rule\  
    (  
        LinkExtractor(allow = (r'^https?:\\/\\/ (www\\.)?[^.]*\\.gov[\\.]*$',)),  
        callback = 'parse_page',  
        follow = True  
    ),  
)
```

Within the *crawler* class, the rules are initialized for the spider. The *LinkExtractor* contains a regular expression that parses only web links that are government websites. The *callback* and *follow* rules allow the spider to crawl pages from these obtained links. By default, the Scrapy spider will also not fetch duplicate web links. Note that the current iteration of the web crawler only crawls through government websites. This is denoted by the regular expression that represents the respective syntax of these government websites.

```
def __init__(self, url = None):  
    super().__init__()   
    self.start_urls = url
```

In the `__init__` function in the *crawler* class, the function itself is imposed with *super()* in order to prevent inheritance issues. The `__init__` function also initializes the spider with the seed URLs provided by the user.

```
def parse_page(self, response):
```

The *parse_page* function deals with handling the HTML page obtained. The function writes the HTML page to its respective data file in the data folder. The function also counts how many pages (if applicable) the spider should generate depending on the user input. If a file fails to write or the limit of pages crawled is reached, the spider raises a flag through *CloseSpider()* and exits. The function obtains user input information through global variables.

Main Code Blocks

Lines 40 through 69 of the code handles user input and feedback. The code block simply parses the user input and determines the necessary parameters: seed file name and starting URLs, the number of pages the spider should crawl, the number of levels or depth the spider should be limited to crawling, and the folder and file name where the data should be stored in. The code block consists of *try* and *except* data structures in order to effectively provide feedback to the user on incorrect or useful input.

```
try:  
    os.makedirs(directory, exist_ok = True)  
except Exception:  
    print('[ERROR] Failed to establish output directory')  
    exit()  
  
i = 0  
filename = directory + '/' + directory + str(i) + '.html'  
while os.path.isfile(filename):  
    i += 1  
    filename = filename[0:len(directory + '/' + directory)] + str(i) + '.html'
```

This code block handles how HTML data is stored. The *try* block creates (if it does not exist) the output directory that the user specifies. The program stores data as `<template>/<template><int>`. For example, if the user specifies *data* as the output directory, then the first HTML page will be stored in *data/data0*. The integer value will increment on each HTML page to be saved. The *while* loop allows the program to append to an already existing output directory. Therefore, if the output directory contains existing data, the program will create files where it last left off. For example, if files *data/data0* to *data/data50* exist, the program starts writing files at *data/data51*.

```
settings = get_project_settings()
settings.update({'DEPTH_LIMIT' : levels})
process = CrawlerProcess(settings)
spider = crawler()
process.crawl(spider, seed)
process.start()
```

The last code block deals with initialization of the spider. First, the settings for the Scrapy spider are obtained. The settings are updated to allow restrictions in the levels or depth of fetched URLs from the seed URLs. The spider process is then generated with these settings and the appropriate spider type: *CrawlSpider* and class *crawler* are initialized. The spider process is then initialized with the respective spider along with another parameter, the seed URLs, which will be passed into the `__init__` function in the *crawler* class. The spider process is then executed, running the spider.

Limitations

The web crawler program contains a few limitations. Scrapy is built on a single-threaded system, and thus can not handle multi-threading of a spider class. However, Scrapy does allow parallelism where the next *GET* request for a web link can be initiated before the previous *GET* request has been returned. Another limitation is that the spider may fail to obtain information from particular websites. Generally, these issues include lack of authentication or limitations of *GET* requests imposed by the hosts themselves. The other limitation is that seed files must be prepended with *http://* or *https://* in order for them to be evaluated correctly.

Instructions

The web crawler is built on Scrapy. Thus, installation of Scrapy is required. It is recommended to install Scrapy on Anaconda when installing for Windows. The program files contain a *crawler.bat* file used for execution in Windows, the source code *crawler.py*, and the seed URLs file *seed.txt*. The *seed.txt* is an optional file and may be manipulated in name, file type, and data. The user can also create their own seed file. To execute the program, the user enters the following command:

crawler.bat <seed file> [# pages] [# levels] <output directory>

where <> are required and [] are optional

However, if the user is not using Windows, they can bypass the batch file and execute in Python directly.

python crawler.py <seed file> [# pages] [# levels] <output directory>

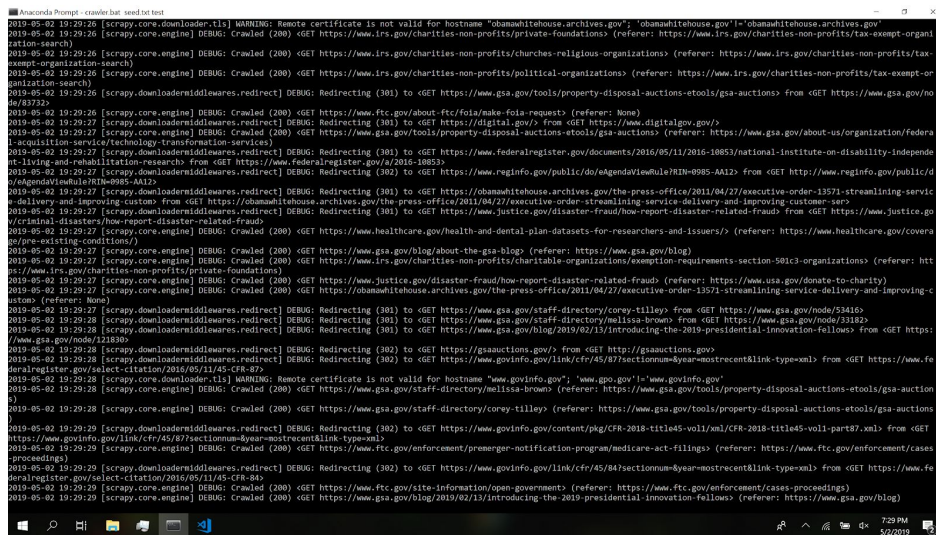
where <> are required and [] are optional

Note that inputting negative integers for pages and levels will be interpreted as infinite. The output directory will be the name of the folder and files the HTML code will be stored in.

Note: The web crawler program is compatible with the Web Document Index and Search program:

<https://github.com/steventfan/Web-Document-Index-and-Search>

Screenshots



```
Microsoft Prompt - crawler.bat - seed.txt
2019-05-02 19:29:26 [scrapy.core.downloader.cl] WARNING: Remote certificate is not valid for hostname 'obamawhitehouse.archives.gov'; 'obamawhitehouse.gov' is 'obamawhitehouse.archives.gov'
2019-05-02 19:29:26 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://www.irs.gov/charities-non-profits/private-foundations> (referer: https://www.irs.gov/charities-non-profits/tax-exempt-organization-search)
2019-05-02 19:29:26 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://www.irs.gov/charities-non-profits/churches-religious-organizations> (referer: https://www.irs.gov/charities-non-profits/tax-exempt-organization-search)
2019-05-02 19:29:26 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://www.irs.gov/charities-non-profits/political-organizations> (referer: https://www.irs.gov/charities-non-profits/tax-exempt-organization-search)
2019-05-02 19:29:26 [scrapy.downloadermiddlewares.redirect] DEBUG: Redirecting (301) to <GET https://www.gsa.gov/tools/property-disposal-auctions-etools/gsa-auctions> from <GET https://www.gsa.gov/no/469322>
2019-05-02 19:29:26 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://www.ftc.gov/about-ftc/foia/make-foia-request> (referer: None)
2019-05-02 19:29:27 [scrapy.downloadermiddlewares.redirect] DEBUG: Redirecting (301) to <GET https://digital.gov/> from <GET https://www.digital.gov/>
2019-05-02 19:29:27 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://www.gsa.gov/tools/property-disposal-auctions-etools/gsa-auctions> (referer: https://www.gsa.gov/about-us/organization/federal-acquisition-service/technology-transformation-services)
2019-05-02 19:29:27 [scrapy.downloadermiddlewares.redirect] DEBUG: Redirecting (301) to <GET https://www.federalregister.gov/documents/2016/05/11/2016-10853/national-institute-on-disability-independence-living-and-rehabilitation-research> from <GET https://www.federalregister.gov/4/2016-10853>
2019-05-02 19:29:27 [scrapy.downloadermiddlewares.redirect] DEBUG: Redirecting (302) to <GET https://www.reginfo.gov/public/do/AgendaViewRule?RIN=0985-AA12> from <GET http://www.reginfo.gov/public/d/AgendaViewRule?RIN=0985-AA12>
2019-05-02 19:29:27 [scrapy.downloadermiddlewares.redirect] DEBUG: Redirecting (301) to <GET https://obamawhitehouse.archives.gov/the-press-office/2011/04/22/executive-order-13571-streamlining-service-delivery-and-improving-customer-s> from <GET https://obamawhitehouse.archives.gov/the-press-office/2011/04/22/executive-order-streamlining-service-delivery-and-improving-customer-s>
2019-05-02 19:29:27 [scrapy.downloadermiddlewares.redirect] DEBUG: Redirecting (301) to <GET https://www.justice.gov/disaster-fraud/how-report-disaster-related-fraud> from <GET https://www.justice.gov/criminal-disaster/how-report-disaster-related-fraud>
2019-05-02 19:29:27 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://www.healthcare.gov/health-and-dental-plan-datasets-for-researchers-and-issuers/> (referer: https://www.healthcare.gov/covers-a/pre-existing-conditions/)
2019-05-02 19:29:27 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://www.gsa.gov/blog/about-the-gsa-blog> (referer: https://www.gsa.gov/blog)
2019-05-02 19:29:27 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://www.irs.gov/charities-non-profits/charitable-organizations/exemption-requirements-section-501(c)(3)-organizations> (referer: https://www.irs.gov/charities-non-profits/private-foundations)
2019-05-02 19:29:27 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://www.justice.gov/disaster-fraud/how-report-disaster-related-fraud> (referer: https://www.usa.gov/donate-to-charity)
2019-05-02 19:29:27 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://obamawhitehouse.archives.gov/the-press-office/2011/04/22/executive-order-13571-streamlining-service-delivery-and-improving-customer-s> (referer: None)
2019-05-02 19:29:27 [scrapy.downloadermiddlewares.redirect] DEBUG: Redirecting (301) to <GET https://www.gsa.gov/staff-directory/correy-tilley> from <GET https://www.gsa.gov/node/53416>
2019-05-02 19:29:28 [scrapy.downloadermiddlewares.redirect] DEBUG: Redirecting (301) to <GET https://www.gsa.gov/staff-directory/melissa-brown> from <GET https://www.gsa.gov/node/53182>
2019-05-02 19:29:28 [scrapy.downloadermiddlewares.redirect] DEBUG: Redirecting (301) to <GET https://www.gsa.gov/blog/2019/02/13/introducing-the-2019-presidential-innovation-fellows> from <GET https://www.gsa.gov/node/72180>
2019-05-02 19:29:28 [scrapy.downloadermiddlewares.redirect] DEBUG: Redirecting (302) to <GET https://gsaauctions.gov/> from <GET http://gsaauctions.gov>
2019-05-02 19:29:28 [scrapy.downloadermiddlewares.redirect] DEBUG: Redirecting (302) to <GET https://www.govinfo.gov/link/cfr/45/87/sectionnum-by-year-mostrecent&link-type=xm> from <GET https://www.federalregister.gov/select-citation/2016/05/11/45-CFR-87>
2019-05-02 19:29:28 [scrapy.core.downloader.cl] WARNING: Remote certificate is not valid for hostname 'www.govinfo.gov'; 'www.gpo.gov' is 'www.govinfo.gov'
2019-05-02 19:29:28 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://www.gsa.gov/staff-directory/melissa-brown> (referer: https://www.gsa.gov/tools/property-disposal-auctions-etools/gsa-auctions)
2019-05-02 19:29:28 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://www.gsa.gov/staff-directory/correy-tilley> (referer: https://www.gsa.gov/tools/property-disposal-auctions-etools/gsa-auctions)
2019-05-02 19:29:29 [scrapy.downloadermiddlewares.redirect] DEBUG: Redirecting (302) to <GET https://www.govinfo.gov/content/pkg/CFR-2018-title45-vol1/ml/CFR-2018-title45-vol1-part87.xml> from <GET https://www.govinfo.gov/link/cfr/45/87/sectionnum-by-year-mostrecent&link-type=xm>
2019-05-02 19:29:29 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://www.ftc.gov/enforcement/premerger-notification-program/medicare-act-filings> (referer: https://www.ftc.gov/enforcement/cases/proceedings)
2019-05-02 19:29:29 [scrapy.downloadermiddlewares.redirect] DEBUG: Redirecting (302) to <GET https://www.govinfo.gov/link/cfr/45/84/sectionnum-by-year-mostrecent&link-type=xm> from <GET https://www.federalregister.gov/select-citation/2016/05/11/45-CFR-84>
2019-05-02 19:29:29 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://www.ftc.gov/site-information/open-governments> (referer: https://www.ftc.gov/enforcement/cases/proceedings)
2019-05-02 19:29:29 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://www.gsa.gov/blog/2019/02/13/introducing-the-2019-presidential-innovation-fellows> (referer: https://www.gsa.gov/blog)
```


100% 90% 80% 70% 60% 50% 40% 30% 20% 10% 0%
