

Documentation

Aabishkar Karki
karkiaabishkar@gmail.com

April 5 2021

To create the crawler project I have used scrapy framework from python which helps in crawling data from websites easily. Explanations of the functions and classes can be found in the comments of the code. In this documentation I have just highlighted the important parts of the project.

When the project is run it goes to "packages_spider.py" inside the spiders folder. This file contains the definition of the crawler (spider) which has been implemented as a class. It is going to be used to crawl data from those two websites. To store the crawled/scraped data a container class [Figure 1] has been created in "items.py" file that is inside Crawler/Crawler.

```
import scrapy

'''
    This class is used to store the data that has been scraped/crawled from the given websites.
'''
class Packages(scrapy.Item):
    # define the fields for your item here like:
    # name = scrapy.Field() -> dictionary key=name value=value

    package_url=scrapy.Field()
    package_description=scrapy.Field()
    package_name=scrapy.Field()
```

Figure 1: Class to store crawled data

In the "main.py" file I have added the details for the website that we need to crawl the data from. For this project I have added two process functions which calls the same crawler to crawl that website but gives their own respective websites output. The crawler has been given a name so once the Crawler is activated using its name it opens the links that are present in its 'start_urls' property. Both of the crawler have been given only one start url that opens the first page of each websites list of packages. After that the parse function is automatically called by scrapy and using "XPath" the relevant details are extracted/crawled from the webpage like package name, url and descriptions and stored in the Class that was created earlier. Since both the websites need to extract same data the same class has been used in both of the crawler to store data. The parse function also checks whether we have extracted 100 package datas from the website. Once 100 packages are extracted the function ends. If the initial page contains < 100 packages then the url of the next page is extracted using XPath and then the 'parse' function is called on that url and again the same process is repeated.

The output of both the crawler can be seen in Figure 2 and Figure 3. For the project description I have extracted the project description that was written in the main web page that contained the list of packages. I tried extracting the project description by opening each of the package's url but the description was too large hence could not be easily stored in tables like this. Output of unittests can be seen in Figure 4. In this unittests I have only tested the parse function of the crawler as it is the most important part of the program.

New websites can be easily scraped using this project if similar data needs to be extracted. We just need to pass the details from the "main.py" file to the crawler/spider class that does the scraping. However for websites where we need to extract a variety of other data new Classes have to be created to store and crawl the data but the other processes are still the same.

Resources that have been used in the project:

- Scrapy

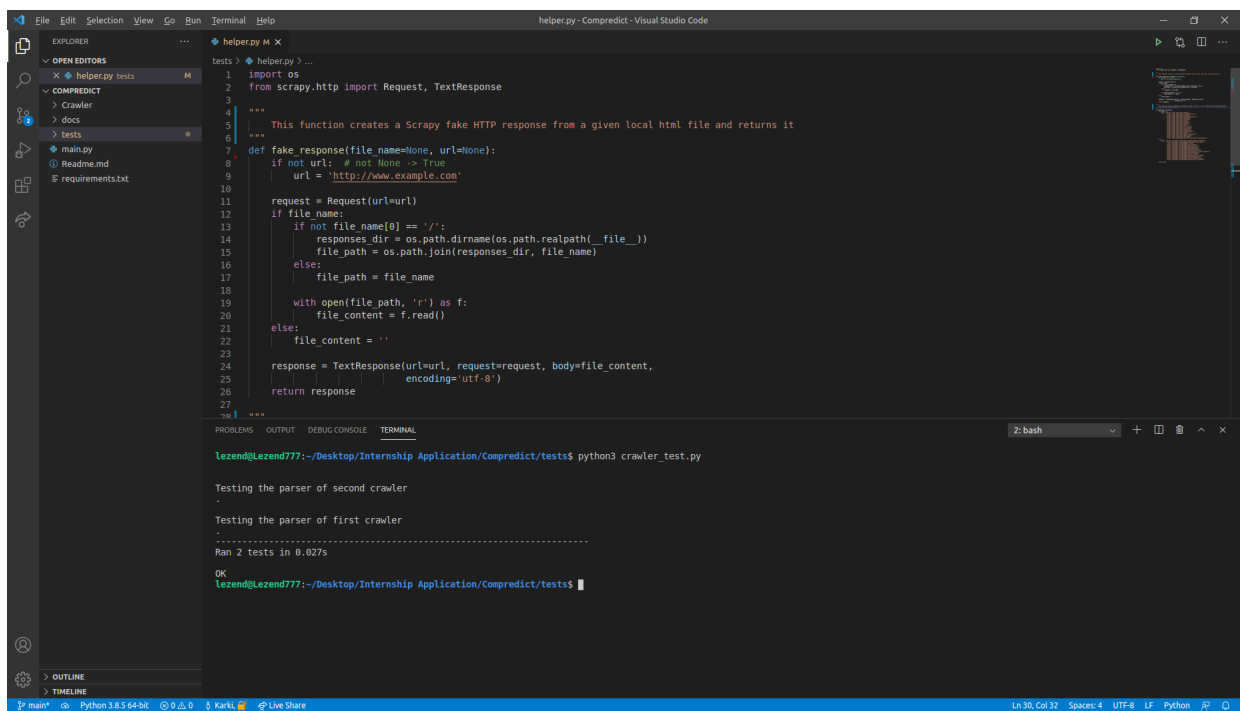
- Xpath
- Unittests

package_name	package_description	package_url
fnlib1.18.0	Neural Network Libraries	https://pypi.org/project/fnlib/
reaction-generator1.4.164	Reaction Network Generator	https://pypi.org/project/reaction-generator/
ontobio2.6.4	Library for working with GBO Library Ontologies and associations	https://pypi.org/project/ontobio/
large-image-tasks1.4.3	Order Worker tasks for Large Image.	https://pypi.org/project/large-image-tasks/
modemcli.6.4	ModemCLI: high performance rendering for Python 3	https://pypi.org/project/modemcli/
ocrmypdf1.7.1.3	OCRmyPDF adds an OCR text layer to scanned PDF files, allowing them to be searched	https://pypi.org/project/ocrmypdf/
obspy1.2.2	ObsPy - a Python framework for seismological observations.	https://pypi.org/project/obspy/
gpcoder0.38.1	Gpcoder is a simple and consistent gpcoding library.	https://pypi.org/project/gpcoder/
klito1.2.0	Conventions for Python + Apache Beam	https://pypi.org/project/klito/
klito-core21.2.0	Core klito library for common functionality	https://pypi.org/project/klito-core/
onnruntime1.7.0	ONNX Runtime is a runtime accelerator for Machine Learning models	https://pypi.org/project/onnxruntime/
simpy-py1.6.7	Fast and light weight simulator of rigid poly-articulated systems.	https://pypi.org/project/simpy-py/
klito-exec21.2.0	Component of klito project that reduces boilerplate while writing YAML-config-driven, Dockerized python Beam pipelines.	https://pypi.org/project/klito-exec/
pymankendall1.4.1	A python package for non-parametric Mann-Kendall family of trend tests.	https://pypi.org/project/pymankendall/
jnrl.1.2	Jnrl is the cloud-native neural search solution powered by the state-of-the-art AI and deep learning	https://pypi.org/project/jnrl/
python-opensesame3.3.8	A graphical experiment builder for the social sciences	https://pypi.org/project/python-opensesame/
bohrium-api11.0.post60	Bohrium Python API	https://pypi.org/project/bohrium-api/
pytripe93.3.5	PyTrie	https://pypi.org/project/pytripe/
pytorch-transformers-pvt-nightly1.2.0.dev201909281000	PyTorch: Python-based Simulations of Chemistry Framework	https://pypi.org/project/pytorch-transformers-pvt-nightly/
diatlet41.1	Repository of pre-trained NLP Transformer models: BERT & RoBERTa, GPT & GPT-2, Transformer-XL, XLNet and XLM	https://pypi.org/project/diatlet/
allennlp-py-nightly0.9.1.dev201910011800	Web Client for Visualizing Pandas Objects	https://pypi.org/project/allennlp-py-nightly/
napienotbook2.89	An open-source ML research library, built on PyTorch.	https://pypi.org/project/napienotbook/
mnem22.1	Dashboard for financial market data	https://pypi.org/project/mnem/
bioinfo2.7.0	MNE python project for MEG and EEG data analysis.	https://pypi.org/project/bioinfo/
tsufile0.92	A comprehensive library for computational molecular biology	https://pypi.org/project/tsufile/
Ktrain0.26.2	USO Scientific Collaboration Algorithm Library - minimal Python package	https://pypi.org/project/ktrain/
nippyee1.6.0	ktapi is a wrapper for TensorFlow Keras that makes deep learning and AI more accessible and easier to apply	https://pypi.org/project/nippyee/
dm-reverb-nightly0.3.0.dev20210405	Neurologist in Python: Pipelines and Interfaces	https://pypi.org/project/dm-reverb-nightly/
dtwl4.0	Reverb is an efficient and easy-to-use data storage and transport system designed for machine learning research.	https://pypi.org/project/dtwl/
biosteam2.25.1	Python DTW Module	https://pypi.org/project/biosteam/
boost-histogram1.0.1	The Boost: Histogram Python wrapper.	https://pypi.org/project/boost-histogram/
PyGeddy021.3.3	Pure Python geddy tools	https://pypi.org/project/pygeddy/
mlmln0.3.3	Little for Machine Learning projects	https://pypi.org/project/mlmln/
MCEq1.2.1	Numerical cascade equation solver	https://pypi.org/project/mceq/
PyOLM1.99.3	OpenGL Mathematics library for Python	https://pypi.org/project/pyolm/
torchmetrics0.2.0	PyTorch native Metrics	https://pypi.org/project/torchmetrics/
xrtgen0.14.2	XrtGen is a Python library for 3D grids, surfaces, wells, etc.	https://pypi.org/project/xrtgen/
deep-daze0.9.0	Deep Daze	https://pypi.org/project/deep-daze/
MCCPy0.8.5	MCC: parsing and manipulation in Python	https://pypi.org/project/mccpy/
minia2.4.0	Supports Speedy Python for MINUT2 in C++	https://pypi.org/project/minia/
pandas-summary0.0.7	An extension to pandas describe function.	https://pypi.org/project/pandas-summary/
news0.7.1	Statistical learning for newsprocessing in Python	https://pypi.org/project/news/
word2number1.1	Convert number words eg. three hundred and forty two to numbers (342).	https://pypi.org/project/word2number/

Figure 2: Output of First Crawler

package_name	package_description	package_url		
symfony-polyfill-redis0.9.0	Symfony polyfill for the Redis extension	https://packagist.org/packages/symfony/polyfill-redis		
symfony-polyfill-curl0.9.0	Symfony polyfill for cURL functions	https://packagist.org/packages/symfony/polyfill-curl		
psr/log	Common interface for logging libraries	https://packagist.org/packages/psr/log		
psr/container	Common Container Interface (PHP FIG PSR-11)	https://packagist.org/packages/psr/container		
symfony/console	Eases the creation of beautiful and testable command line interfaces	https://packagist.org/packages/symfony/console		
guzzlehttp/psr7	PSR-7 message implementation that also provides common utility methods	https://packagist.org/packages/guzzlehttp/psr7		
webmozart/assert	Assertions to validate method input/output with nice error messages.	https://packagist.org/packages/webmozart/assert		
guzzlehttp/promises	Guzzle promises library	https://packagist.org/packages/guzzlehttp/promises		
symfony-polyfill-intl-normalizer	Symfony polyfill for intl's Normalizer class and related functions	https://packagist.org/packages/symfony/polyfill-intl-normalizer		
psr/http-message	Common interface for HTTP messages	https://packagist.org/packages/psr/http-message		
symfony/finder	Finds files and directories via an intuitive fluent interface	https://packagist.org/packages/symfony/finder		
symfony/polyfill-php80	Symfony polyfill backporting some PHP 8.0+ features to lower PHP versions	https://packagist.org/packages/symfony/polyfill-php80		
guzzlehttp/guzzle	Guzzle is a PHP HTTP client library	https://packagist.org/packages/guzzlehttp/guzzle		
doctrine/instantiator	A small, lightweight utility to instantiate objects in PHP without invoking their constructors	https://packagist.org/packages/doctrine/instantiator		
symfony/polyfill-php72	Symfony polyfill backporting some PHP 7.2+ features to lower PHP versions	https://packagist.org/packages/symfony/polyfill-php72		
symfony/event-dispatcher	Provides tools that allow your application components to communicate with each other by dispatching events and listening to them	https://packagist.org/packages/symfony/event-dispatcher		
symfony/polyfill-intl-idn	Symfony polyfill for intl's idn_to_ascii and idn_to_utf8 functions	https://packagist.org/packages/symfony/polyfill-intl-idn		
symfony/process	Executes commands in sub-processes	https://packagist.org/packages/symfony/process		
symfony/polyfill-php73	Symfony polyfill backporting some PHP 7.3+ features to lower PHP versions	https://packagist.org/packages/symfony/polyfill-php73		
phpdocumentor/reflection-docblock	With this component, a library can provide support for annotations via DocBlocks or otherwise retrieve information that is embedded in a DocBlock.	https://packagist.org/packages/phpdocumentor/reflection-docblock		
ralouphie/getallheaders	A polyfill for getAllHeaders	https://packagist.org/packages/ralouphie/getallheaders		
symfony/service-contracts	Generic abstractions related to writing services	https://packagist.org/packages/symfony/service-contracts		
phpdocumentor/type-resolver	A PSR-5 based resolver of Class names, Types and Structural Element Names	https://packagist.org/packages/phpdocumentor/type-resolver		
phpdocumentor/reflection-common	Common reflection classes used by phpdocumentor to reflect the code structure	https://packagist.org/packages/phpdocumentor/reflection-common		
doctrine/lexer	PHP Doctrine Lexer parser library that can be used in Top-Down, Recursive Descent Parsers.	https://packagist.org/packages/doctrine/lexer		
monolog/monolog	Send your logs to files, sockets, inboxes, databases and various web services	https://packagist.org/packages/monolog/monolog		
sebastian/diff	Diff implementation	https://packagist.org/packages/sebastian/diff		
phpunit/phpunit	The PHP Unit Testing framework.	https://packagist.org/packages/phpunit/phpunit		
phpspec/phpspec	Highly opinionated mocking framework for PHP 5.3+	https://packagist.org/packages/phpspec/phpspec		
phpunit/php-code-coverage	Library that provides collection, processing, and rendering functionality for PHP code coverage information.	https://packagist.org/packages/phpunit/php-code-coverage		
sebastian/exporter	Provides the functionality to export PHP variables for visualization	https://packagist.org/packages/sebastian/exporter		
sebastian/recursion-context	Provides functionality to recursively process PHP variables	https://packagist.org/packages/sebastian/recursion-context		
sebastian/comparator	Provides the functionality to compare two values for equality	https://packagist.org/packages/sebastian/comparator		
phpspec/php-dsl	Utility class for timing	https://packagist.org/packages/phpspec/php-dsl		
phpspec/deep-copy	Create deep copies (clones) of your objects	https://packagist.org/packages/phpspec/deep-copy		
phpunit/php-file-iterator	FilteredIterator implementation that filters files based on a list of suffixes.	https://packagist.org/packages/phpunit/php-file-iterator		
sebastian/environment	Provides functionality to handle HTTP/PHP environments	https://packagist.org/packages/sebastian/environment		
symfony/var-dumper	Provides mechanisms for walking through any arbitrary PHP variable	https://packagist.org/packages/symfony/var-dumper		
sebastian/version	Library that helps with managing the version number of Git-hosted PHP projects	https://packagist.org/packages/sebastian/version		
sebastian/global-state	Snapshotting of global state	https://packagist.org/packages/sebastian/global-state		
phpunit/php-text-template	Simple template engine.	https://packagist.org/packages/phpunit/php-text-template		
doctrine/inflector	PHP Doctrine Inflector is a small library that can perform string manipulations with regard to upper/lowercase and singular/plural forms of words.	https://packagist.org/packages/doctrine/inflector		
symfony/http-foundation	Defines an object-oriented layer for the HTTP specification	https://packagist.org/packages/symfony/http-foundation		
sebastian/code-unit-reverse-lookup	Looks up which function or method a line of code belongs to	https://packagist.org/packages/sebastian/code-unit-reverse-lookup		

Figure 3: Output of Second Crawler



The image shows a Visual Studio Code editor window with a Python file named `helper.py` open. The file contains a function `fake_response` that simulates an HTTP response. The function takes `file_name` and `url` as arguments. It checks if `url` is not None and if `file_name` is not None. If both are provided, it reads the content of the file specified by `file_name` and returns a `TextResponse` object with the `url`, `request`, and `body` (file content). If `file_name` is None, it returns an empty `TextResponse` object.

```
1 import os
2 from scrapy.http import Request, TextResponse
3
4 """
5 This function creates a Scrapy fake HTTP response from a given local html file and returns it
6 """
7 def fake_response(file_name=None, url=None):
8     if not url: # not None -> True
9         url = 'http://www.example.com'
10
11     request = Request(url=url)
12     if file_name:
13         if not file_name[0] == '/':
14             responses_dir = os.path.dirname(os.path.realpath(__file__))
15             file_path = os.path.join(responses_dir, file_name)
16         else:
17             file_path = file_name
18
19         with open(file_path, 'r') as f:
20             file_content = f.read()
21     else:
22         file_content = ''
23
24     response = TextResponse(url=url, request=request, body=file_content,
25                             encoding='utf-8')
26     return response
27 """
```

The terminal output shows the execution of the test script `crawler_test.py`. The output indicates that the parser for the second crawler is being tested, followed by the parser for the first crawler. The tests pass, and the output shows "Ran 2 tests in 0.027s" and "OK".

```
lezend@lezend777:~/Desktop/Internship Application/Compredict/tests$ python3 crawler_test.py
Testing the parser of second crawler
.
Testing the parser of first crawler
.
Ran 2 tests in 0.027s
OK
lezend@lezend777:~/Desktop/Internship Application/Compredict/tests$
```

Figure 4: Output of Unittests