# Web Scraping
# Gathering Data from Websites Using Scrapy Framework
# (http://172.16.211.189:8088)

*By,*

**Mr. Raghudathesh GP**

Assistant Professor – Sr. Scale

Manipal School of Information Sciences

MAHE, Manipal

- Status codes are issued by a server in response to a client's request made to the server

  – **110 - Connection timed out**

  – **200 - Success**

  – **404 - Not Found**

https://en.wikipedia.org/wiki/List_of_HTTP_status_codes

- .json - is a minimal, readable format for structuring data. It is used primarily to transmit data between a server and web application.

```
{
    "date" : "29",
    "year" : "2019"

}
```

- JL - json lines: Every line is a json. Great for streaming data and easy for appending new jsons

```
{"day:"29","year":"2019"}
{"day:"12","year":"2012"}
{"day:"2","year":"2010"}
```

http://jsonlines.org/

# Why Scrapy ?

- Scrapy is a open source and collaborative framework for crawling/scraping the web
  - Scrapy is an excellent choice for focused crawls
  - Scrapy is written in Python
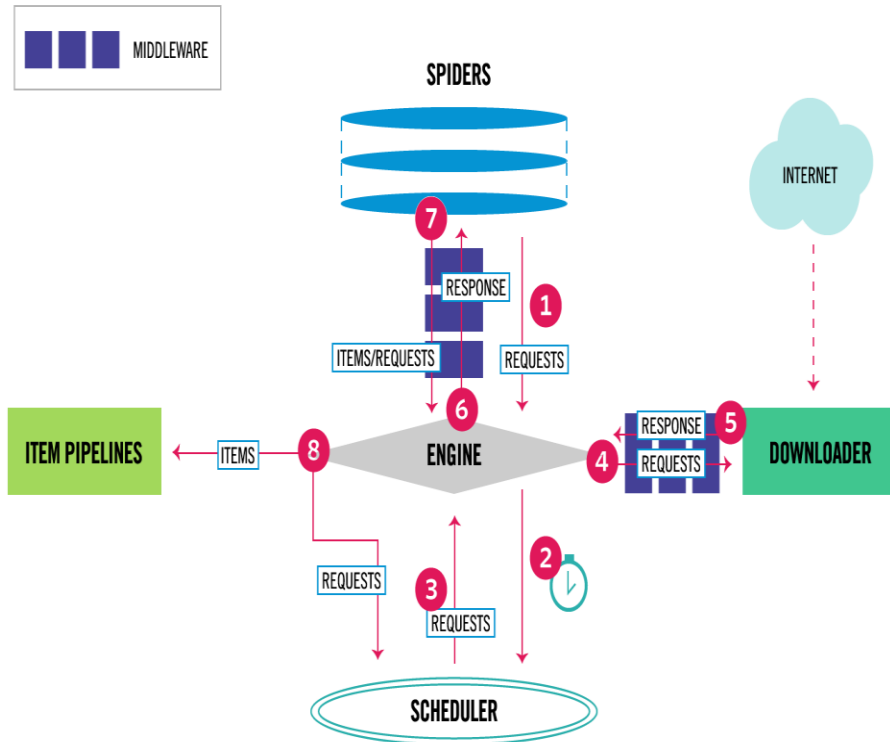  - Scrapy is faster than Heritr

Yadav, M., & Goyal, N. (2015). Comparison of Open Source Crawlers-A Review. *International Journal of Scientific & Engineering Research*, *6*(9), 1544-1551.

TABLE 3
COMPARISON OF OPEN SOURCE CRAWLERS IN TERMS OF VARIOUS PARAMETERS

| Open source Crawlers | Language | Operating System | License | Parallel |
|---|---|---|---|---|
| Scrapy | Python | Linux/Mac OS X/Windows | BSD License | Yes(During broad crawls) |
| Apache Nutch | Java | Cross-platform | Apache License 2.0 | Yes (Using Hadoop) |
| Heritrix | Java | Linux/Unix-like/Windows Unsupported | Apache License | Yes |
| WebSphinix | Java | Windows, Mac, Linux, Android, IOS | Apache Software License | Yes |
| JSpider | Java | Windows, Windows7, Window vista | GNU Library or Lesser General Public License version 2.0 (LGPLv2) | Yes |
| Gnu Wget | C | Cross-platform | GNU General Public License version 3 and later | |
| WIRE | C/C++ | | GPL LIcense | |
| Pavuk | C | Linux | GNU General Public License Version 2.0(GPLV2) | Yes |
| Teleport | - | Windows | Apache License | Yes |
| Web2disk | - | Windows | - | Yes |
| WebCopierPro | - | Windows/Mac OS X | - | No |
| WebHTTrack | C/C++ | Cross-Platform | GPL | Yes |

| Framework | Library |
|---|---|
| Provides ready to use tools, standards, templates, and policies for fast application development | Provides reusable function for our code |
| The framework controls calling of libraries for our code | Our code controls when and where to call a library |
| To leverage the benefit of a framework, a fresh application can be developed following the framework's guideline | Library can be added to augment the features of an existing application |
| Easy to create and deploy an application | Facilitates program binding |
| Helps us to develop a software application quickly | Helps us to reuse a software function |
| Intent of a framework is to reduce the complexity of the software development process | Intent of a library is to provide reusable software functionality |

https://www.baeldung.com/cs/framework-vs-library/

# Scrapy Architecture



Source - Architecture overview — Scrapy 2.4.1 documentation

1. The **Engine** gets the initial Requests to crawl from the **Spider**.
2. The **Engine** schedules the Requests in the **Scheduler** and asks for the next Requests to crawl.
3. The **Scheduler** returns the next Requests to the **Engine**.
4. The **Engine** sends the Requests to the **Downloader**, passing through the Downloader Middlewares (see process_request()).
5. Once the page finishes downloading the **Downloader** generates a Response (with that page) and sends it to the Engine, passing through the Downloader Middlewares (see process_response()).
6. The **Engine** receives the Response from the **Downloader** and sends it to the **Spider** for processing, passing through the Spider Middleware (see process_spider_input()).
7. The **Spider** processes the Response and returns scraped items and new Requests (to follow) to the **Engine** , passing through the Spider Middleware (see process_spider_output()).
8. The **Engine** sends processed items to Item Pipelines, then send processed Requests to the **Scheduler** and asks for possible next Requests to crawl.
9. The process repeats (from step 1) until there are no more requests from the **Scheduler** .

# Roadmap

# Understand the webpage we are crawling

> **scrapy startproject tutorial**

New Scrapy project 'tutorial', using template directory '/usr/local/lib/python3.5/site-packages/scrapy/templates/project', created in:
    /Users/mtodor/Projects/meetups/tutorial

You can start your first spider with:
    cd tutorial
    scrapy genspider example example.com

> **cd tutorial**

> **scrapy genspider example example.com**

> **scrapy crawl example -t json -o output.json**

# Scrapy basics

```
> cd tutorial && ls *
__init__.py
items.py
pipelines.py
settings.py
...
spiders:
__init__.py
__pycache__
fivethirtyeight.py
```
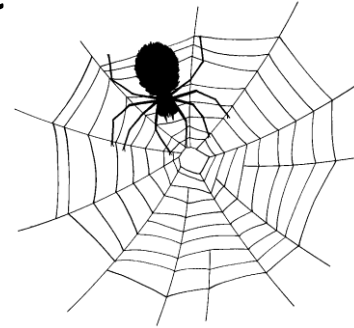
- Spiders are the place where you define the custom behaviour for crawling and parsing pages for a particular site (or, in some cases, a group of sites).

http://doc.scrapy.org/en/latest/topics/spiders.html

```
class ExampleSpider(scrapy.Spider):
    name = "example"
    allowed_domains = ["example.com"]
    start_urls = ['http://www.example.com/']
    def parse(self, response):

        ...
        process response
        ...
```

- Process **response** using selectors
  - xpath()
  - css()
  - extract()
  - re()

- Follow links
  - yield scrapy**.**Request(url, callback**=**self**.**parse_link)

```
1     """
2     This spider will crawl all the questions on datascience stack exchange
3
4     to run `scrapy crawl getquestions`
5
6     This is a follow spider and it will keep adding urls to the queue which have subdomain as /questions/
7     Question heading, question body and answers are extracted from the html
8     BeautifulSoup is used to extract the accepted answer.
9     """
10
11
12    from scrapy.spiders import CrawlSpider, Rule
13    from scrapy.linkextractors import LinkExtractor
14    from bs4 import BeautifulSoup
15    import jsonlines
16
17
18    class someSpider(CrawlSpider):
19        # name of spider
20        name = "getquestions"
21
22        # make the crawler stick to this domain
23        allowed_domains = ["datascience.stackexchange.com"]
24
25        # picking the most aswered question so that there are a lot of next URLs to crawl
26        start_urls = [
27            "https://datascience.stackexchange.com/questions/6107/what-are-deconvolutional-layers"
28        ]
29
30        # Rule — to crawl only subdomains which have "questions" in them and follow them
31        rules = (
32            Rule(LinkExtractor(allow="/questions/*"), callback="parse_obj", follow=Trues),
33        )
```

Description

Library imports

Defining class variables and rules

scrapy crawl quotes-1

**scrapy crawl <name_of_spider>**

# Scrapy item exporters

- Default exporters:
  - json': 'scrapy.exporters.JsonItemExporter'
  - 'jsonlines': 'scrapy.exporters.JsonLinesItemExporter'
  - 'jl': 'scrapy.exporters.JsonLinesItemExporter'
  - 'csv': 'scrapy.exporters.CsvItemExporter'
  - 'xml': 'scrapy.exporters.XmlItemExporter'
  - 'marshal': 'scrapy.exporters.MarshalItemExporter'
  - 'pickle': 'scrapy.exporters.PickleItemExporter'

http://doc.scrapy.org/en/latest/topics/exporters.html

- After an item has been scraped by a spider, it is sent to the Item Pipeline which processes it through several components that are executed sequentially.

http://doc.scrapy.org/en/latest/topics/item-pipeline.html

- settings.py
  - LOG_LEVEL = 'INFO'
  - FEED_EXPORTERS = {'json': 'wiki_logs.exporters.UnicodeJsonItemExporter'}
    - Create a custom JSON exporter because the builtin one is brain damaged and forces ASCII output
  - ROBOTSTXT_OBEY = False
    - or, be polite and respect robots.txt

```
.
└── stackcrawl
    ├── scrapy.cfg
    └── stackcrawl
        ├── __init__.py
        ├── items.py
        ├── middlewares.py
        ├── pipelines.py
        ├── __pycache__
        │   ├── __init__.cpython-38.pyc
        │   └── settings.cpython-38.pyc
        ├── settings.py
        └── spiders
            ├── body_scrapy.py
            ├── data.jsonl
            ├── __init__.py
            ├── links.txt
            └── __pycache__
                ├── body_scrapy.cpython-38.pyc
                └── __init__.cpython-38.pyc

5 directories, 14 files
```

```
# Scrapy settings for stackcrawl project
#
# For simplicity, this file contains only settings considered important or
# commonly used. You can find more settings consulting the documentation:
#
#     https://docs.scrapy.org/en/latest/topics/settings.html
#     https://docs.scrapy.org/en/latest/topics/downloader-middleware.html
#     https://docs.scrapy.org/en/latest/topics/spider-middleware.html

BOT_NAME = 'stackcrawl'

SPIDER_MODULES = ['stackcrawl.spiders']
NEWSPIDER_MODULE = 'stackcrawl.spiders'

LOG_LEVEL="INFO"
# Crawl responsibly by identifying yourself (and your website) on the user-agent
#USER_AGENT = 'stackcrawl (+http://www.yourdomain.com)'

# Obey robots.txt rules
ROBOTSTXT_OBEY = True

# Configure maximum concurrent requests performed by Scrapy (default: 16)
#CONCURRENT_REQUESTS = 32

# Configure a delay for requests for the same website (default: 0)
# See https://docs.scrapy.org/en/latest/topics/settings.html#download-delay
# See also autothrottle settings and docs
DOWNLOAD_DELAY = 5
```
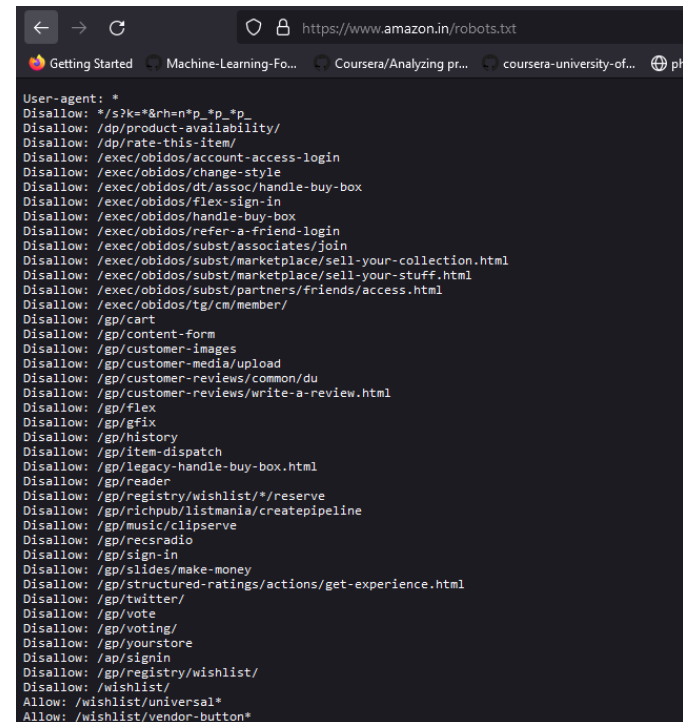
**Avoid getting banned**

- rotate your user agent from a pool of well-known ones from browsers (google around to get a list of them)

- disable cookies (see COOKIES_ENABLED) as some sites may use cookies to spot bot behaviour

- use download delays (5 or higher). See DOWNLOAD_DELAY setting.

- if possible, use Google cache to fetch pages, instead of hitting the sites directly

- What Makes a Crawler Polite?
  - A polite crawler respects **robots.txt**
  - A polite crawler never degrades a website's performance
  - A polite crawler identifies its creator with contact information
  - A polite crawler is not a pain a for system administrators

https://blog.scrapinghub.com/2016/08/25/how-to-crawl-the-web-politely-with-scrapy/

# Queries

# THANK YOU...