py4kids (https://github.com/wgong/py4kids)

# Scrapy- Crawl & Scrape the Web

In this lesson we learn how to extract data from the web using powerful scrapy

In [2]:
```python
from jyquickhelper import add_notebook_menu
add_notebook_menu()
```

Out[2]:
- Scrapy
  - Install
  - Create a project
  - Customize spider
  - Launch spider
  - Extract data
  - Follow next link
- References

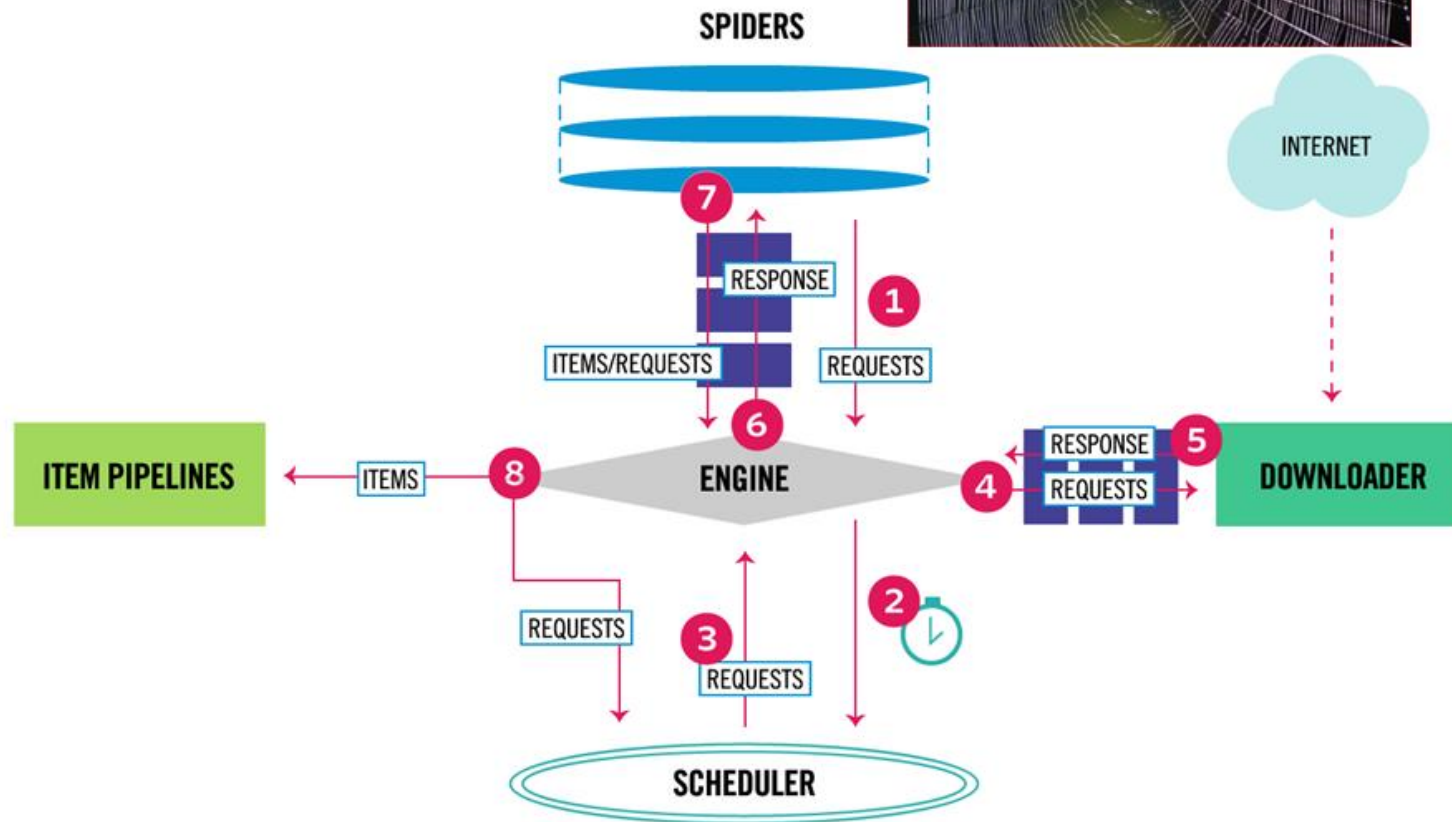## Scrapy (https://scrapy.org/)



### Overview

An open source and collaborative framework for extracting the data you need from websites. In a fast, simple, yet extensible way.
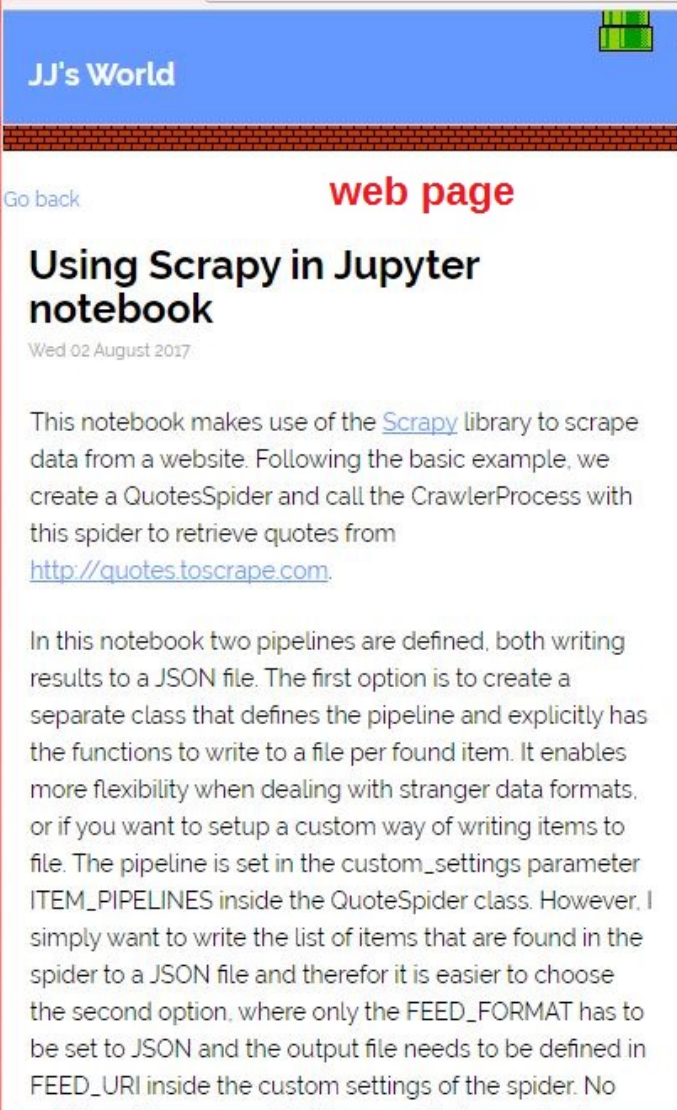
# Scrapy advantages

▶ **Faster than mechanize** because it uses asynchronous operations (Twisted).

▪ Scrapy has **better support for html parsing**.

▪ Scrapy has **better support for unicode characters, redirections, gzipped responses, encodings.**

▪ You can export the **extracted data directly to JSON, XML and CSV.**
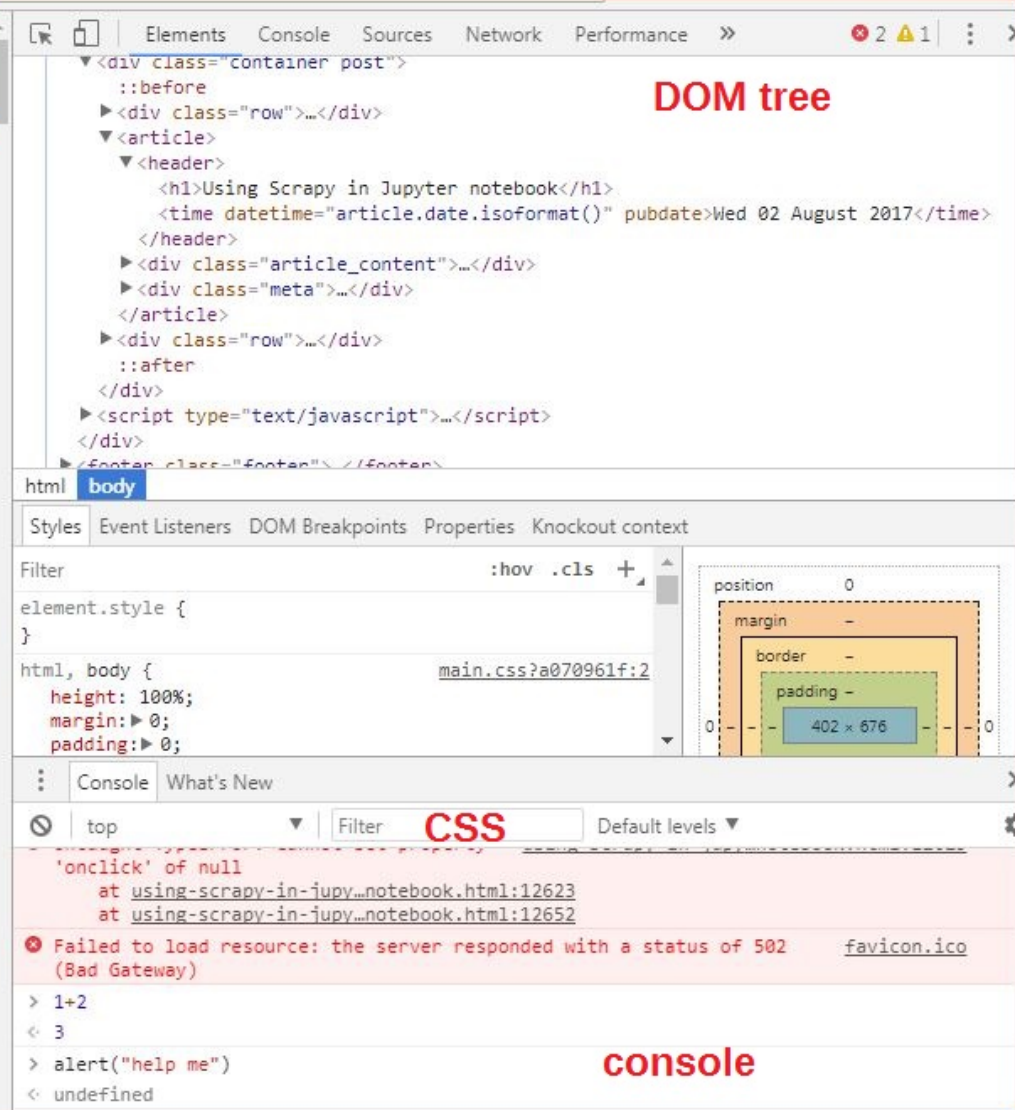
**How does it work**

**DOM tree in browser**

## Install

```
cmd> conda install scrapy
```

```
cmd> pip install scrapy
```

## Create a project

- open a terminal window,
- create a working folder,
- type

```
cmd> scrapy startproject tutorial
```



## Customize spider

- specify download web addresses (URL)
- define parser logic

```
quotes_spider.py
1        import scrapy
2
3        class QuotesSpider(scrapy.Spider):
4            name = "quotes"
5
6            start_urls = [
7                'http://quotes.toscrape.com/page/1/',
8                'http://quotes.toscrape.com/page/2/',
9            ]
10
11           def parse(self, response):
12               page_number = response.url.split("/")[-2]
13               filename = 'quotes-%s.html' % page_number
14               with open(filename, 'wb') as f:
15                   f.write(response.body)
```

## Launch spider

```
CMD> scrapy crawl quotes
```

Two files will be created in the current folder:

- quotes-1.html:
    - web copy (http://quotes.toscrape.com/page/1/)
    - local copy (./scrapy-tutorial/quotes-1.html)
- quotes-2.html
    - web copy (http://quotes.toscrape.com/page/2/)
    - local copy (./scrapy-tutorial/quotes-2.html)

## Extract data

```
CMD> scrapy crawl quotes
```

Two files will be created in the current folder:

- quotes-1.html:
  - web copy (http://quotes.toscrape.com/page/1/)
  - local copy (./scrapy-tutorial/quotes-1.html)
- quotes-2.html
  - web copy (http://quotes.toscrape.com/page/2/)
  - local copy (./scrapy-tutorial/quotes-2.html)

Extract data and save them into CSV file:

```
11    def parse(self, response):
12        page_number = response.url.split("/")[-2]
13        filename = 'quotes-%s.html' % page_number
14        with open(filename, 'wb') as f:
15            f.write(response.body)
16
17        # extract quote, author, tags and store data in csv file
18        import csv
19        file_csv = 'quotes-%s.csv' % page_number
20        with open(file_csv, 'w') as f_csv:
21            csv_writer = csv.writer(f_csv, lineterminator='\n', quoting=csv.QUOTE_ALL)
22            for quote in response.css('div.quote'):
23                txt = quote.css('span.text::text').extract_first()
24                author = quote.css('small.author::text').extract_first()
25                tags = quote.css('div.tags a.tag::text').extract()
26                tags_str = " : ".join(tags)
27                csv_writer.writerow([author,txt,tags_str])
28
```

## Follow next link

To download all the quotes, follow "Next Page" line recursively:

```
28
29            # follow the link
30            next_page = response.css('li.next a::attr(href)').extract_first()
31            if next_page is not None:
32                yield response.follow(next_page, callback=self.parse)
```

# References

- Basic Tutorial (https://doc.scrapy.org/en/latest/intro/tutorial.html)
- Using Scrapy in Jupyter notebook (http://www.jitsejan.nl/using-scrapy-in-jupyter-notebook.html)
- Build your 1st web crawler (https://medium.com/python-pandemonium/develop-your-first-web-crawler-in-python-scrapy-6b2ee4baf954)
- How To Crawl A Web Page with Scrapy and Python (https://www.digitalocean.com/community/tutorials/how-to-crawl-a-web-page-with-scrapy-and-python-3)
- Scrapy Tutorial: Web Scraping Craigslist (http://python.gotrained.com/scrapy-tutorial-web-scraping-craigslist/)

In [ ]: