# Scraping the Web with Scrapy

# Hello, there!

- I'm **Matt Lebrun**
- PyClub Ring Leader at **Save22, Inc.**
- Trustee of **Python.PH**
- Been in love with **Python** since 2012
- *#peanutbutterlife*

IF SOMEONE EVER TELLS YOU THAT
YOU'RE PUTTING TOO MUCH PEANUT
BUTTER ON YOUR BREAD,

STOP TALKING TO THEM. YOU DON'T
NEED THAT KIND OF NEGATIVITY IN
YOUR LIFE.

# So you want to get data from the web?

- Content extraction
- Page navigation

# Content extraction (Scraping)

# Content extraction (Scraping)

- request
- beautifulsoup
- lxml

# Page navigation (Crawling)

# Page navigation (Crawling)

- custom script

# What is Scrapy?

- Web scraping framework

- Handles scraping *and* crawling

- Has an interactive shell for testing **CSS** and **XPath** expressions

- Support for generating multiple feed formats (CSV, JSON, XML) and storing them in multiple backends (FTP, S3, local filesystem)

- Encoding support and auto detection

- Highly extensible via different signal hooks and middlewares

- Python3 support!

- (more)

# Install Scrapy

```
$ pip install scrapy
```

*Note: There might be caveats depending on you OS*

# Scrapy script

```python
 1  import scrapy¬
 2  from scrapy.loader import ItemLoader¬
 3  ¬
 4  ¬
 5  class Product(scrapy.Item):¬
 6      title = scrapy.Field()¬
 7      current_price = scrapy.Field()¬
 8      url = scrapy.Field()¬
 9      sku = scrapy.Field()¬
10      primary_image =scrapy.Field()¬
11  ¬
12  ¬
13  class LazadaScraper(scrapy.Spider):¬
14      name = 'lazada_scraper'¬
15      start_urls = [¬
16          'http://www.lazada.com.ph/catalog/?q=peanut+butter',¬
17      ]¬
18  ¬
19      def parse(self, response):¬
20          product_grid = response.xpath(¬
21              '//div[@data-component="product_list"]'¬
22              '/div[contains(@class, "product-card")]'¬
23          )¬
24  ¬
25          for selector in product_grid:¬
26              loader = ItemLoader(Product(), selector=selector)¬
27              loader.add_xpath('title', 'a/div/div/span/@title')¬
28              loader.add_xpath('current_price', '@data-price')¬
29              loader.add_xpath('url', 'a/@href')¬
30              loader.add_xpath('sku', '@data-sku')¬
31              loader.add_xpath('primary_image', 'a/div/img/@data-original')¬
32              yield loader.load_item()¬
```

lazada_scraper.py

# Scrapy script

```
 1  import scrapy¬
 2  from scrapy.loader import ItemLoader¬
 3  ¬
 4  ¬
 5  class Product(scrapy.Item):¬
 6      title = scrapy.Field()¬
 7      current_price = scrapy.Field()¬
 8      url = scrapy.Field()¬
 9      sku = scrapy.Field()¬
10      primary_image =scrapy.Field()¬
11  ¬
12  ¬
13  class LazadaScraper(scrapy.Spider):¬
14      name = 'lazada_scraper'¬
15      start_urls = [¬
16          'http://www.lazada.com.ph/catalog/?q=peanut+butter',¬
17      ]¬
18  ¬
19      def parse(self, response):¬
20          product_grid = response.xpath(¬
21              '//div[@data-component="product_list"]'¬
22              '/div[contains(@class, "product-card")]'¬
23          )¬
24  ¬
25          for selector in product_grid:¬
26              loader = ItemLoader(Product(), selector=selector)¬
27              loader.add_xpath('title', 'a/div/div/span/@title')¬
28              loader.add_xpath('current_price', '@data-price')¬
29              loader.add_xpath('url', 'a/@href')¬
30              loader.add_xpath('sku', '@data-sku')¬
31              loader.add_xpath('primary_image', 'a/div/img/@data-original')¬
32              yield loader.load_item()¬
```

Scraped item field definition

lazada_scraper.py

# Scrapy script

```python
 1  import scrapy¬
 2  from scrapy.loader import ItemLoader¬
 3  ¬
 4  ¬
 5  class Product(scrapy.Item):¬
 6      title = scrapy.Field()¬
 7      current_price = scrapy.Field()¬
 8      url = scrapy.Field()¬
 9      sku = scrapy.Field()¬
10      primary_image =scrapy.Field()¬
11  ¬
12  ¬
13  class LazadaScraper(scrapy.Spider):¬
14      name = 'lazada_scraper'¬
15      start_urls = [¬
16          'http://www.lazada.com.ph/catalog/?q=peanut+butter',¬
17      ]¬
18  ¬
19      def parse(self, response):¬
20          product_grid = response.xpath(¬
21              '//div[@data-component="product_list"]'¬
22              '/div[contains(@class, "product-card")]'¬
23          )¬
24  ¬
25          for selector in product_grid:¬
26              loader = ItemLoader(Product(), selector=selector)¬
27              loader.add_xpath('title', 'a/div/div/span/@title')¬
28              loader.add_xpath('current_price', '@data-price')¬
29              loader.add_xpath('url', 'a/@href')¬
30              loader.add_xpath('sku', '@data-sku')¬
31              loader.add_xpath('primary_image', 'a/div/img/@data-original')¬
32              yield loader.load_item()¬
```

Main spider code

lazada_scraper.py

# Scrapy script

```python
1  import scrapy¬
2  from scrapy.loader import ItemLoader¬
3  ¬
4  ¬
5  class Product(scrapy.Item):¬
6      title = scrapy.Field()¬
7      current_price = scrapy.Field()¬
8      url = scrapy.Field()¬
9      sku = scrapy.Field()¬
10     primary_image =scrapy.Field()¬
11 ¬
12 ¬
13 class LazadaScraper(scrapy.Spider):¬
14     name = 'lazada_scraper'¬
15     start_urls = [¬
16         'http://www.lazada.com.ph/catalog/?q=peanut+butter',¬
17     ]¬
18 ¬
19     def parse(self, response):¬
20         product_grid = response.xpath(¬
21             '//div[@data-component="product_list"]'¬
22             '/div[contains(@class, "product-card")]'¬
23         )¬
24 ¬
25         for selector in product_grid:¬
26             loader = ItemLoader(Product(), selector=selector)¬
27             loader.add_xpath('title',  'a/div/div/span/@title')¬
28             loader.add_xpath('current_price',  '@data-price')¬
29             loader.add_xpath('url',  'a/@href')¬
30             loader.add_xpath('sku',  '@data-sku')¬
31             loader.add_xpath('primary_image',  'a/div/img/@data-original')¬
32             yield loader.load_item()¬
```

Define what URLs to scrape

lazada_scraper.py

# Scrapy script

```
 1  import scrapy¬
 2  from scrapy.loader import ItemLoader¬
 3  ¬
 4  ¬
 5  class Product(scrapy.Item):¬
 6      title = scrapy.Field()¬
 7      current_price = scrapy.Field()¬
 8      url = scrapy.Field()¬
 9      sku = scrapy.Field()¬
10      primary_image =scrapy.Field()¬
11  ¬
12  ¬
13  class LazadaScraper(scrapy.Spider):¬
14      name = 'lazada_scraper'¬
15      start_urls = [¬
16          'http://www.lazada.com.ph/catalog/?q=peanut+butter',¬
17      ]¬
18  ¬
19      def parse(self, response):¬
20          product_grid = response.xpath(¬
21              '//div[@data-component="product_list"]'¬
22              '/div[contains(@class, "product-card")]'¬
23          )¬
24  ¬
25          for selector in product_grid:¬
26              loader = ItemLoader(Product(), selector=selector)¬
27              loader.add_xpath('title',  'a/div/div/span/@title')¬
28              loader.add_xpath('current_price',  '@data-price')¬
29              loader.add_xpath('url',  'a/@href')¬
30              loader.add_xpath('sku',  '@data-sku')¬
31              loader.add_xpath('primary_image', 'a/div/img/@data-original')¬
32              yield loader.load_item()¬
```

Page scraping code

lazada_scraper.py

# Scrapy script

```
 1  import scrapy¬
 2  from scrapy.loader import ItemLoader¬
 3  ¬
 4  ¬
 5  class Product(scrapy.Item):¬
 6      title = scrapy.Field()¬
 7      current_price = scrapy.Field()¬
 8      url = scrapy.Field()¬
 9      sku = scrapy.Field()¬
10      primary_image =scrapy.Field()¬
11  ¬
12  ¬
13  class LazadaScraper(scrapy.Spider):¬
14      name = 'lazada_scraper'¬
15      start_urls = [¬
16          'http://www.lazada.com.ph/catalog/?q=peanut+butter',¬
17      ]¬
18  ¬
19      def parse(self, response):¬
20          product_grid = response.xpath(¬
21              '//div[@data-component="product_list"]'¬
22              '/div[contains(@class, "product-card")]'¬
23          )¬
24  ¬
25          for selector in product_grid:¬
26              loader = ItemLoader(Product(), selector=selector)¬
27              loader.add_xpath('title', 'a/div/div/span/@title')¬
28              loader.add_xpath('current_price', '@data-price')¬
29              loader.add_xpath('url', 'a/@href')¬
30              loader.add_xpath('sku', '@data-sku')¬
31              loader.add_xpath('primary_image', 'a/div/img/@data-original')¬
32              yield loader.load_item()¬
```

Find the product grid in the page response

lazada_scraper.py

# Scrapy script

```python
1  import scrapy¬
2  from scrapy.loader import ItemLoader¬
3  ¬
4  ¬
5  class Product(scrapy.Item):¬
6      title = scrapy.Field()¬
7      current_price = scrapy.Field()¬
8      url = scrapy.Field()¬
9      sku = scrapy.Field()¬
10     primary_image =scrapy.Field()¬
11 ¬
12 ¬
13 class LazadaScraper(scrapy.Spider):¬
14     name = 'lazada_scraper'¬
15     start_urls = [¬
16         'http://www.lazada.com.ph/catalog/?q=peanut+butter',¬
17     ]¬
18 ¬
19     def parse(self, response):¬
20         product_grid = response.xpath(¬
21             '//div[@data-component="product_list"]'¬
22             '/div[contains(@class, "product-card")]'¬
23         )¬
24 ¬
25         for selector in product_grid:¬
26             loader = ItemLoader(Product(), selector=selector)¬
27             loader.add_xpath('title', 'a/div/div/span/@title')¬
28             loader.add_xpath('current_price', '@data-price')¬
29             loader.add_xpath('url', 'a/@href')¬
30             loader.add_xpath('sku', '@data-sku')¬
31             loader.add_xpath('primary_image', 'a/div/img/@data-original')¬
32             yield loader.load_item()¬
```

Loop through each grid item and parse data

lazada_scraper.py

# Run the Scrapy spider

```
$ scrapy runspider lazada_scraper.py -t csv -o
lazada_scraper.csv
```

# Item data extraction tips

- Use scrapy shell
- Test and study how to extract each field (CSS, XPath, Beautifulsoup, RegEx, etc)
- Use the python debugger (pdb)

# Scrapy ecosystem

- Scrapyd
- Splash
- ScrapingHub

# Further reading

- **How to Install Scrapy in Windows**
  `https://scraper24x7.wordpress.com/2016/03/19/how-to-install-scrapy-in-windows/`

- **Official Scrapy Documentation**
  `https://doc.scrapy.org/en/1.2/`

# Q & A

# You guys are awesome!

Let's keep in touch ;)

email: **matt@lebrun.org**
social: **cr8ivecodesmith**

Looking for internship?
`http://goo.gl/forms/mOQ39cm3Gy`