# Extracting Structured Data from the Web Using Scrapy

GETTING STARTED SCRAPING WEBSITES USING SCRAPY



**Janani Ravi**
CO-FOUNDER, LOONYCORN

www.loonycorn.com

# Overview

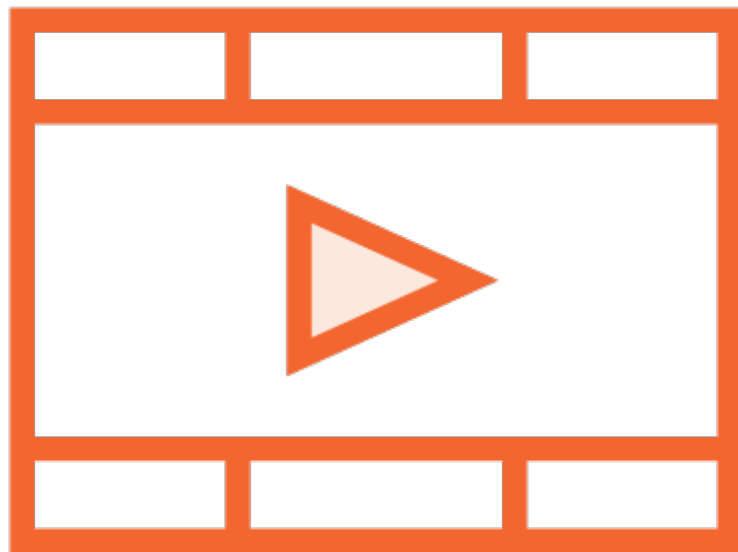Scrapy is an application framework for crawling websites

Allows data extraction in a structured format

The Scrapy shell is an interactive shell to quickly test data extraction

Selectors allow you to specify XPath and CSS classes to scrape information

# Prerequisites and Course Outline

# Prerequisite Courses

**Python: Getting Started**

**Python Fundamentals**

**Advanced Python**

# Software and Skills

**Be very comfortable programming in Python (Python 3)**

**Understand some basics of HTML and CSS**

# Course Outline

## Scraping websites
- Scrapy shell, XPath and CSS selectors

## Spiders
- Spiders, Items, Item Loaders, Item Pipelines

## Built-in services
- Logging, email notifications
- Debugging using the telnet console
- Broad crawls for parallel scraping
- Auto throttling crawls

## Crawlers on the Scrapy Cloud
- Deploying a Scrapy project on scrapinghub.com
- Scraping on the cloud using Portia

# Introducing Scrapy

# Scrapy

Scrapy is an application framework for crawling web sites and extracting structured data

*https://doc.scrapy.org/en/latest/intro/overview.html*

Originally built for web scraping
but now used for web crawling

# Scrapy

Scrapy is an application framework for crawling web
sites and extracting structured data

*https://doc.scrapy.org/en/latest/intro/overview.html*

# Scraping vs. Crawling

## Web Scraping

Extract data directly from web sites

Data analysis and somewhat unsavory reputation

Specific - "scrape prices from Amazon"

Small scale, results in specialized dataset

## Web Crawling

Download and index web sites

Performed by all search engines and associated with legitimate use

General - "crawl sites linked off Amazon"

Large scale, results in document corpus

Originally built for web scraping but now used for web crawling

# Scrapy

Scrapy is an application framework for crawling web sites and extracting structured data

*https://doc.scrapy.org/en/latest/intro/overview.html*

Framework vs. library: inversion of control

# Scrapy

Scrapy is an application framework for crawling web sites and extracting structured data

*https://doc.scrapy.org/en/latest/intro/overview.html*

# Library vs. Framework

## Library

You call library functions

You write the application and invoke library for specific portions

## Framework

Framework calls you

Framework defines the application and invokes your code for specific portions

Hollywood Principle: Don't call us, we'll call you

This is a defining characteristic of frameworks

Framework vs. library: inversion of control

# Scrapy

Scrapy is an application framework for crawling web sites and extracting structured data

*https://doc.scrapy.org/en/latest/intro/overview.html*

You must know what you are
looking for - tied to HTML format

# Scrapy

Scrapy is an application framework for crawling web
sites and extracting structured data

*https://doc.scrapy.org/en/latest/intro/overview.html*

Inherently somewhat fragile, like regular expressions and other related tools

# Scrapy

Scrapy is an application framework for crawling web sites and extracting structured data

*https://doc.scrapy.org/en/latest/intro/overview.html*

Specific HTML elements are selected
for processing using Selectors

# Scrapy

Scrapy is an application framework for crawling web
sites and extracting structured data

*https://doc.scrapy.org/en/latest/intro/overview.html*

Scrapy supports selectors specified in CSS and XPath

# Scrapy

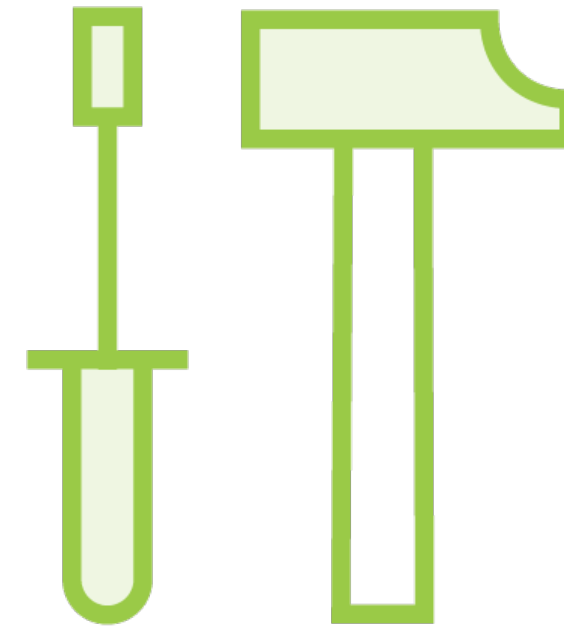Scrapy is an application framework for crawling web sites and extracting structured data

*https://doc.scrapy.org/en/latest/intro/overview.html*

# Benefits of Scrapy

**Asynchronous Callbacks**

Requests and callbacks are scheduled and processed asynchronously

**Granular Control**

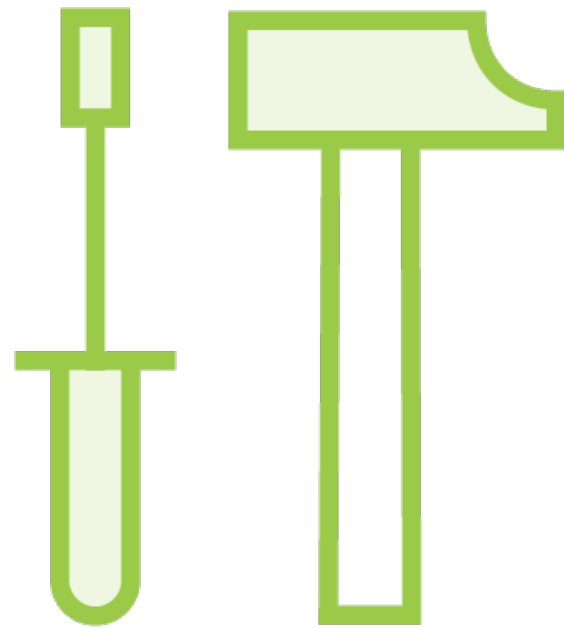Settings to govern politeness of crawl, error handling etc.

# Asynchronous Callbacks

**Speed**

**Parallelism**

**Fault-tolerance**

# Granular Control

**Download delays between requests**

**Limit on concurrent connections**

- Per IP

- Per domain

**Auto-throttling extension**

# Demo

Install and set up Scrapy on your local machine

Basic introduction to Scrapy components

# Demo

**Introducing the Scrapy shell**

# How Scrapy Works

# How Scrapy Works

**Spiders**



**Internet**

**Middleware**

**7**

Response

**1**

Request

Items/Requests

**6**

**Item Pipelines** ← Items **8** **Scrapy Engine** **4** Request → **Downloader**

**5** ← Response

**Requests**

**3** Requests

**2**

**Scheduler**

# How Scrapy Works

**Spiders**

**Middleware**

Response

Items/Requests

Request

Controls
data flow
and triggers
all events

**Internet**

**Item Pipelines**

Items

**Scrapy Engine**

Response

Request

**Downloader**

Requests

Requests

**Scheduler**

# How Scrapy Works

**Spiders**

Spiders are classes written by the user

Middleware

Response

Items/Requests

Request

Internet

Item Pipelines

Items

Scrapy Engine

Response

Request

Downloader

Requests

Requests

Scheduler

# How Scrapy Works

**Spiders**

Define how to parse responses and extract items

**Middleware**

Response

Items/Requests

Request

Internet

**Item Pipelines**

Items

**Scrapy Engine**

Response

Request

**Downloader**

Requests

Requests

**Scheduler**

# How Scrapy Works

Spiders

Engine
forwards
Request to
the Scheduler

Internet

Middleware

Response

Items/Requests

Request

Item Pipelines          Items          Scrapy Engine          Response

Downloader

Request

Requests

2

Requests

Scheduler

# How Scrapy Works

Engine also asks for the next Request from the Scheduler

Spiders

Middleware

Response

Items/Requests

Request

Internet

Item Pipelines

Items

Scrapy Engine

Response

Request

Downloader

Requests

Requests

**2**

Scheduler

# How Scrapy Works

Spiders

Scheduler responds with next Request for Engine to process

Middleware

Response

Items/Requests

Request

Internet

Item Pipelines

Items

**Scrapy Engine**

Response

Request

Downloader

**3**

Requests

**Scheduler**

# How Scrapy Works

Engine requests the Downloader to get this from the internet

**Spiders**

**Middleware**

Response

Items/Requests

Request

**Item Pipelines**

Items

**Scrapy Engine**

**4** Request

**Middleware**

**Downloader**

**Internet**

Requests

Requests

**Scheduler**

# How Scrapy Works

**This request is passed via Downloader Middleware**

Spiders

Middleware

Response

Items/Requests

Request

Internet

**Item Pipelines**

Items

**Scrapy Engine**

Response

**4** Request

**Downloader**

**Middleware**

Requests

Requests

**Scheduler**

# How Scrapy Works

**Downloader Middleware** are hooks between Engine and Downloader

Spiders

**Middleware**

Response

Items/Requests

Request

**Item Pipelines**

Items

**Scrapy Engine**

**4** Request

**Middleware**

**Downloader**

Internet

Response

Requests

Requests

Scheduler

# How Scrapy Works

The downloader fetches the URL from the internet...

Spiders

Internet

Middleware

Response

Items/Requests

Request

Item Pipelines

Items

Scrapy Engine

Response

Request

Downloader

Middleware

Requests

Requests

Scheduler

# How Scrapy Works

...and sends it back to the Scrapy Engine

**Spiders**

Middleware

Response

Items/Requests

Request

**Internet**

**Item Pipelines**

Items

**Scrapy Engine**

Response

**5**

**Downloader**

**Middleware**

Request

Requests

Requests

**Scheduler**

# How Scrapy Works

**The Downloader passes back a Response object to the Engine**

Internet

Spiders

Middleware

Response

Items/Requests

Request

Item Pipelines

Items

Scrapy Engine

Response **5**

Downloader

Middleware

Requests

Requests

Scheduler

# How Scrapy Works

**Spiders**

This too is passed back on the Downloader Middleware

**Middleware**

Response

Items/Requests

Request

**Internet**

**Item Pipelines**

Items

**Scrapy Engine**

Response

**5**

**Downloader**

Request

**Middleware**

Requests

Requests

**Scheduler**

# How Scrapy Works

Spiders

**Middleware**

Response

**The engine forwards that Response object back to the Spider**

Internet

Items/Requests

Request

**6**

**Scrapy Engine**

Item Pipelines

Items

Response

Request

Downloader

**Middleware**

Requests

Requests

Scheduler

# How Scrapy Works

The Spider Middleware is a set of hooks between the Engine and the Spider class

**Spiders**

**Middleware**

Response

**6**

**Scrapy Engine**

**Item Pipelines**

Items

Items/Requests

Request

**Internet**

Response

Request

**Middleware**

**Downloader**

Requests

Requests

**Scheduler**

# How Scrapy Works

**Spiders**

The Spider processes the Response and returns scraped items

**Middleware**

**7**

Response

Items/Requests

Request

**Scrapy Engine**

Item Pipelines

Items

Response

Request

Downloader

Internet

**Middleware**

Requests

Requests

**Scheduler**

# How Scrapy Works

Spiders

The downloaded items
are sent to Item Pipelines

Middleware

Response

Items/Requests

Request

Internet

**Item Pipelines**

Items

**8**

**Scrapy Engine**

Response

Request

**Downloader**

Middleware

Requests

Requests

Scheduler

# How Scrapy Works

**Spiders**

Item Pipelines are used for cleaning, validation, persisting to database etc.

**Middleware**

Response

Items/Requests

Request

**Item Pipelines**

Items

**8**

**Scrapy Engine**

Response

Request

**Downloader**

**Internet**

**Middleware**

Requests

Requests

**Scheduler**

# How Scrapy Works

Any additional Requests are sent to the Scheduler to be added to the crawl queue

**Spiders**

**Middleware**

Response

Items/Requests

Request

**Internet**

**Item Pipelines**

Items

**8**

**Scrapy Engine**

Response

Request

**Downloader**

**Middleware**

Requests

Requests

**Scheduler**

# How Scrapy Works

**Spiders**

**Middleware**

Response

Items/Requests

Request

**Internet**

The process continues until the Scheduler has no more items to send to the Engine

**Item Pipelines**

Items

**Scrapy Engine**

Response

Request

**Downloader**

**Middleware**

**3**

Requests

Requests

**Scheduler**

# Demo

**Working with Selectors using XPath and CSS classes**

# Selector

Specification of what HTML elements ought to be selected for processing. Scrapy supports XPath and CSS selectors.

# Scrapy Selectors

## XPath

Select nodes in an XML (or HTML) document

## CSS

Select HTML elements (usually to associate styles with them)

Scrapy selectors are built atop the lxml library

# Demo

**Using regular expressions with Selectors**

# Summary

Scrapy is an application framework for crawling websites to extract structured data

The Scrapy shell is an interactive shell to quickly test data extraction

Selectors allow you to specify XPath and CSS classes to scrape information