# Scrapy

Scrapy is the most powerful web scraping framework in python, but it's a bit complicated to set up, so check my guide or its documentation to set it up.

## Creating a project and spidr

To create a new project, run the following command in the terminal.

```
Scrapy startproject my_first_spider
```

To create a new spider, first change the directory.

```
cd my_first_spider
```

### Create an spider

```
Scrapy genspider example example.com
```

## The Basic Template

When you create a spider, you obtain a template with the following content.

```
import scrapy
class Example Spider(scrapy.spider):
    name = 'example'
    allowed_domains = ['example.com']          ⎫ Class
    start_urls = ['http://example.com/']       ⎭

    def parse(self, response):    ⎱ Parse method
        pass                      ⎰
```

The class is built with the data we introduced in the previous command, but the parse method needs to be built by us. To build it, use the functions below.

# Finding elements

To find elements in Scrapy, use the response argument from the parse method

```
response.xpath('//tag[@AttributeName="value"]')
```

# Getting the text

To obtain the text element we use text() and either .get() or .getall(). For example:

```
response.xpath('//h1/text()').get()
response.xpath('//tag[@Attribute="value"]/text()').getall()
```

# Return data extracted

To see the data extracted we have to use the yield keyword

```
def parse(self, response):
    title = response.xpath('//h1/text()').get()

    # Return data extracted
    yield {'titles': title}
```

# Run the spider and export data to CSV or JSON

```
Scrapy crawl example
scrapy crawl example -o name_of_file.csv
scrapy crawl example -o name_of_file.json
```