

# D26 : Scrapy 爬蟲編寫 (1)



簡報閱讀



範例與作業



問題討論

**xpath & 基本爬蟲介紹** >

本日知識點目標 >

以scrapy指令自動建立 >

常見的定位器 >

完成一隻簡單的爬蟲 >

## xpath & 基本爬蟲介紹



重要知識點複習 &gt;

解題時間 &gt;

Day 26 SCRAPY 爬蟲編寫-1

xpath &amp; 基本爬蟲介紹



出題教練：黃國展



## 本日知識點目標



## 本日知識點目標

- 上次我們成功地建立了爬蟲專案以及第一隻爬蟲程式
- 這次我們需要使用scrapy將我們要的資料打印出來

在spiders資料夾中你可以看到你的爬蟲程式以scrapy指令自動建立  
在程式中你可以看到：

```
name = 'ettoday'
```

這是你爬蟲的名字，之後你要執行這隻爬蟲時，你會需要告訴scrapy你要執行的是哪隻爬蟲

```
allowed_domains = ['www.ettoday.net']
```

這是你這隻爬蟲被允許請求的domain，他是一個list，也就是說當你需要請求的domain不只一個時，你可以將該domain加入allowed\_domains裡

```
start_urls = ['http://www.ettoday.net/']
```

你的爬蟲起始網站，他也是一個list，你爬蟲的第一步會先對這個list內的全部網址依序進行請求，你可以依照你的需求添加url進start\_urls

```
function parse
```

當你對一個url進行請求後，你需要對返回的response進行操作，而parse這個function就是這隻爬蟲默認進行處理response的function，若你沒有特別指定處理response的function，則由parse這個function進行處理

```
1 # -*- coding: utf-8 -*-
2 import scrapy
3
4
5 class EttodaySpider(scrapy.Spider):
6     name = 'ettoday'
7     allowed_domains = ['www.ettoday.net']
8     start_urls = ['http://www.ettoday.net/']
9
10    def parse(self, response):
11        pass
12
```

## 常見的定位器

- 當我們要進行爬蟲時，常見的定位器有兩種，分別是css selector 和 xpath

- 假設我們要取得回傳的response中title這個標籤裡的文字，以 css selector 及 xpath 分別回傳第一個匹配以己全部匹配資料

	css selector	xpath
第一個匹配的資料	<code>response.css('title::text').get()</code>	<code>response.xpath('//title/text()).get()</code>
全部匹配的資料	<code>response.css('title::text').getall()</code>	<code>response.xpath('//title/text()).getall()</code>

## 完成一隻簡單的爬蟲

- 那麼，知道如何定位資料後，我們開始試著完成一隻簡單的爬蟲
- 首先將我們要爬取的url 加入start\_urls中，並將你要對返回的response要進行的操作放入parse function中，你的爬蟲第一步就算完成了
- 假設我們要爬取某的文章的標題以及內文，並且將爬到的資料print出來，那麼你爬蟲的程式應該會類似這樣

```

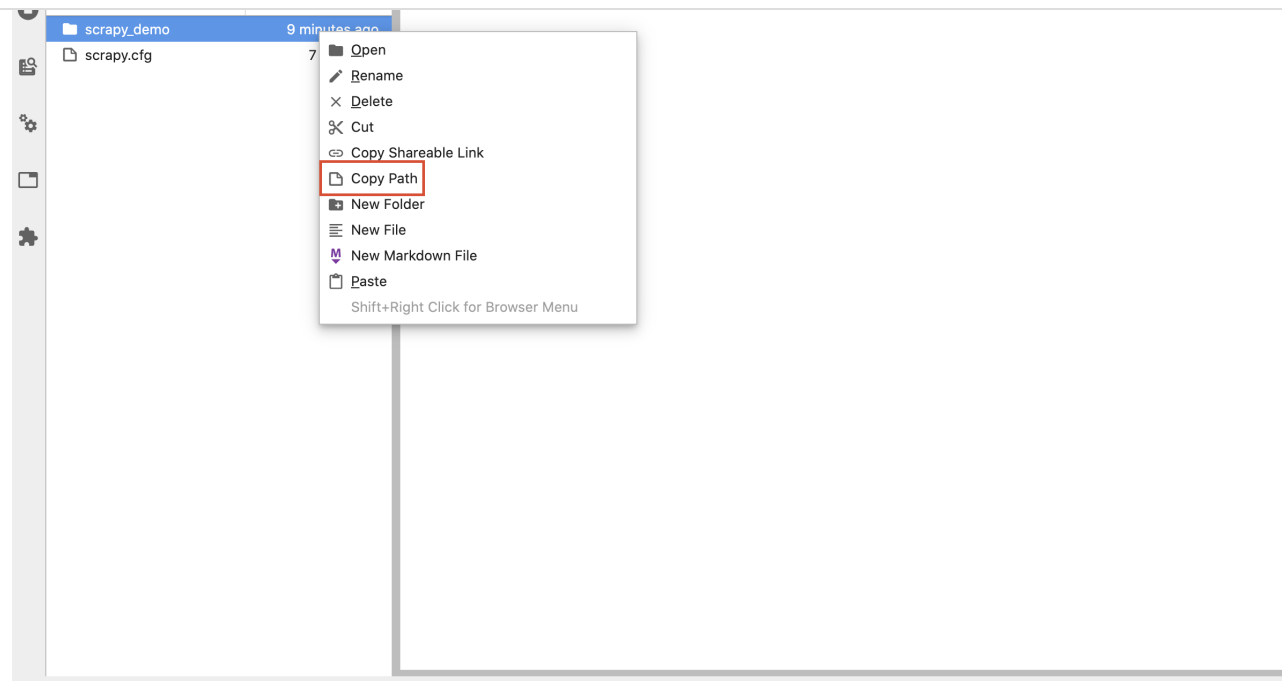
5 class EttodaySpider(scrapy.Spider):
6     name = 'ettoday'
7     allowed_domains = ['www.ettoday.net']
8     start_urls = ['目標新聞1', '目標新聞2', '目標新聞3']
9
10    def parse(self, response):
11        title = response.xpath('/xpath/of/title').get()
12        content = response.xpath('/xpath/of/content').getall()
13        print(title)
14        print(content)

```

Saving completed

Ln 14, Col 22 Spaces: 4 ettoday.py

- 那麼爬蟲編輯完之後要怎麼執行呢？
- scrapy要執行爬蟲的時候需要使用terminal，先開啟terminal 並移動到爬蟲專案的資料夾中
- 如果不知道路徑是什麼，可以對該資料夾點擊右鍵並選擇copy path



- 到該資料夾底下之後，執行指令 `scrapy crawl spider_name`
- 這個spider\_name是你要執行的爬蟲的名稱，範例爬蟲名稱是ettoday，因此執行指令`scrapy crawl ettoday`

## 執行指令後

執行指令後，你會看到scrapy 開始執行該爬蟲，並印出log

此時你會看到你要print出來的資料

我這邊是以<https://www.ettoday.net/news/20201004/1824032.htm>做示範



```
2020-10-04 10:29:59 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://www.ettoday.net/robots.txt> (referer: None)
2020-10-04 18:29:59 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://www.ettoday.net/news/20201004/1824032.htm> (referer: None)
2020-10-04 18:29:59 [scrapy.core.engine] INFO: Closing spider (finished)
2020-10-04 18:29:59 [scrapy.statscollectors] INFO: Dumping Scrapy stats:
{'downloader/request_bytes': 463,
 'downloader/request_count': 2,
 'downloader/request_method_count/GET': 2,
 'downloader/response_bytes': 42297,
 'downloader/response_count': 2,
 'downloader/response_status_count/200': 2,
 'finish_reason': 'finished',
 'finish_time': datetime.datetime(2020, 10, 4, 10, 29, 59, 280416),
 'log_count/DEBUG': 2,
 'log_count/INFO': 9,
 'memusage/max': 48828416,
 'memusage/startup': 48828416,
 'response_received_count': 2,
 'robotstxt/request_count': 1,
 'robotstxt/response_count': 1,
 'robotstxt/response_status_count/200': 1,
 'scheduler/dequeued': 1,
 'scheduler/dequeued/memory': 1,
 'scheduler/enqueued': 1,
 'scheduler/enqueued/memory': 1,
 'start_time': datetime.datetime(2020, 10, 4, 10, 29, 59, 19128)}
2020-10-04 18:29:59 [scrapy.core.engine] INFO: Spider closed (finished)
```

## 重要知識點複習



1. 瞭解如何在scrapy中使用 xpath 以及 css selector 對資料進行定位
2. 了解默認產生的爬蟲程式架構以及簡易編寫流程
3. 如何執行scrapy的爬蟲程式

**解題時間**

# 解題時間 LET'S CRACK IT

Sample Code & 作業

開始解題



下一步：閱讀範例與完成作業

