



AIMS - deep learning

網路爬蟲的介紹與實作

TA course 02

TAs: 蘇時頤, 廖柄淦

➤ Course website:

https://github.com/matteosoo/aimsfellows_DL





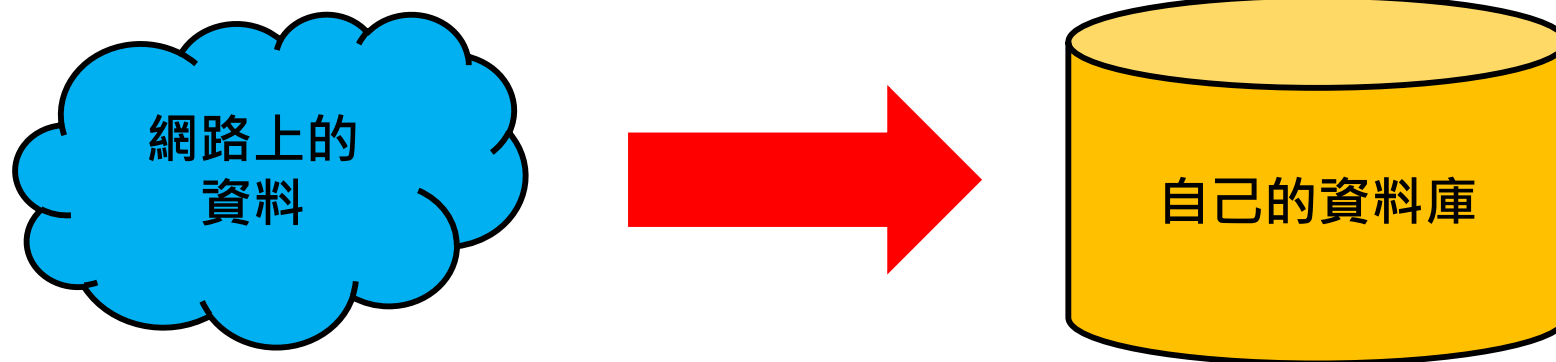
Catalog

- Python 網路爬蟲
 - 什麼是網路爬蟲
 - 網路爬蟲的應用
 - 瀏覽器送出的請求
 - Method Get與Post
 - 爬蟲前的注意事項
 - 撰寫網路爬蟲的步驟
 - HTTP狀態碼
 - 標籤與CSS選擇器
 - 爬蟲實作-PTT表特版當日圖片下載



什麼是 網路爬蟲

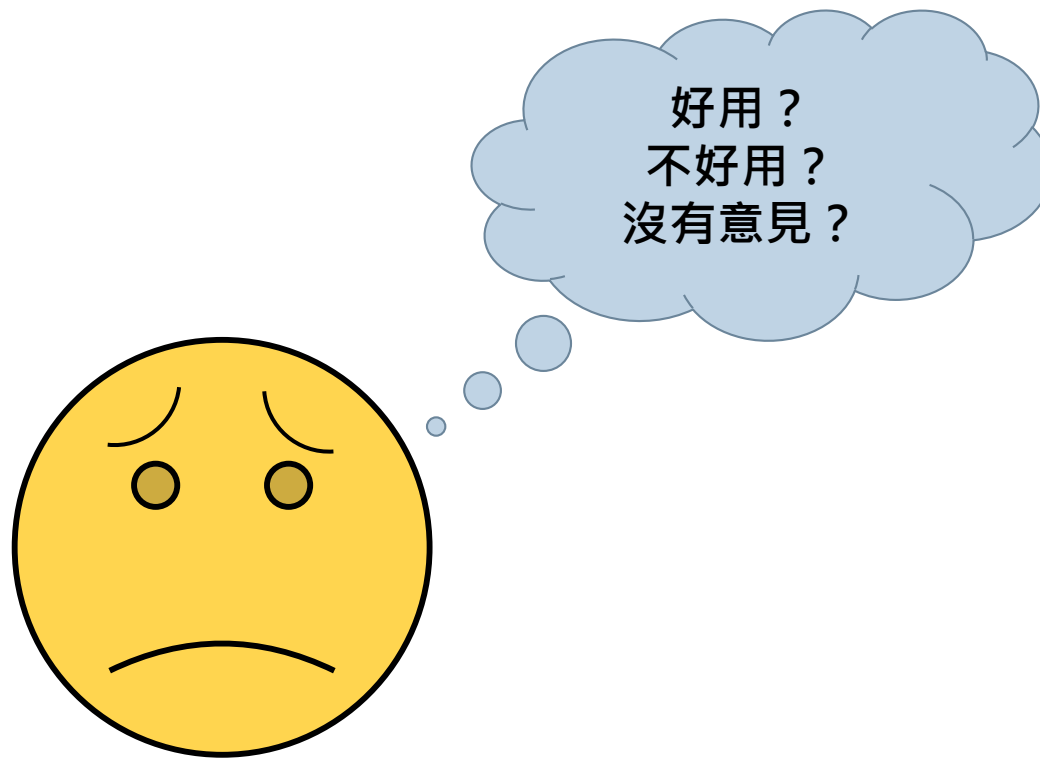
- 是一種利用http request抓取網路資料的技術





網路爬蟲的 應用

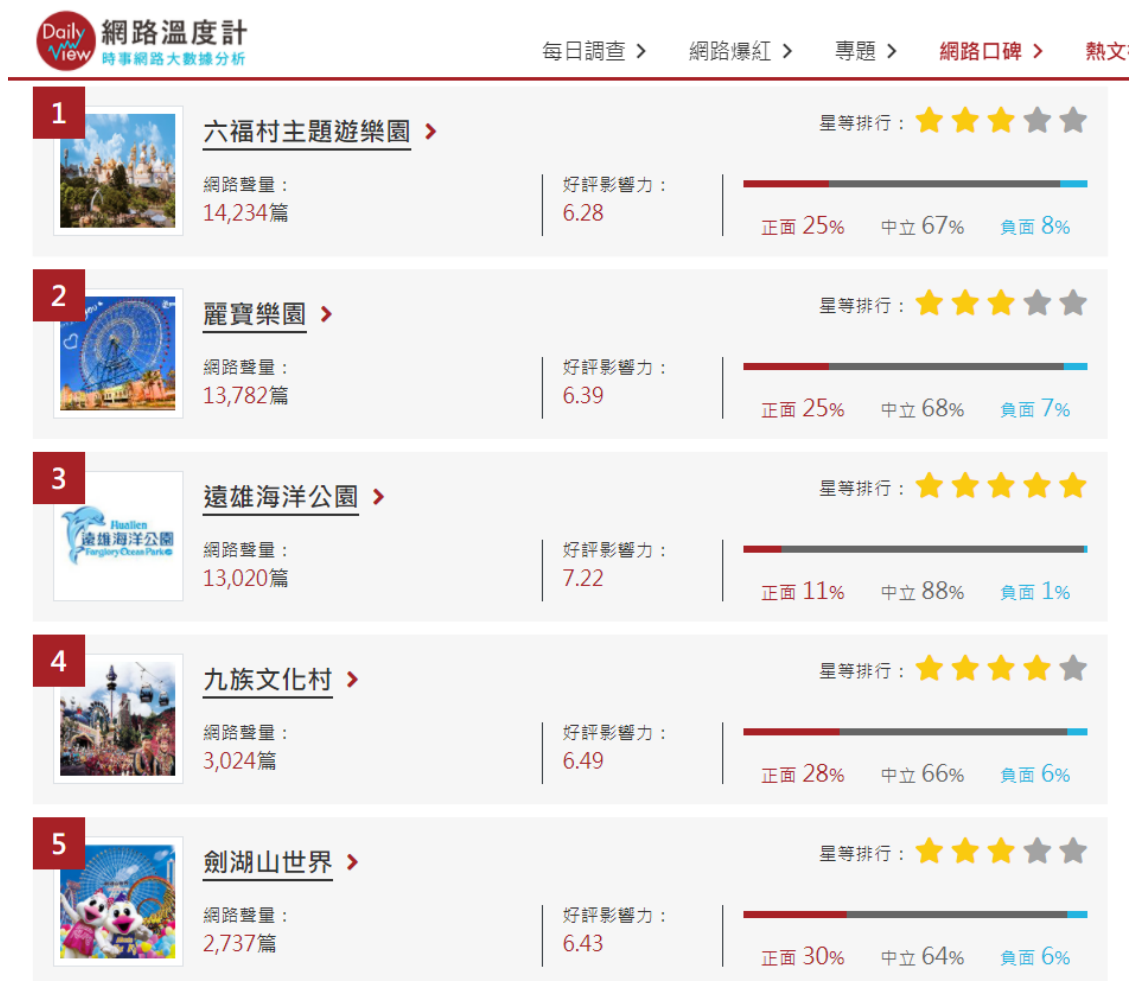
- 網路輿情分析
 - 分析網路上消費者的意見
 - 了解目標客戶對於品牌、產品或是服務的想法





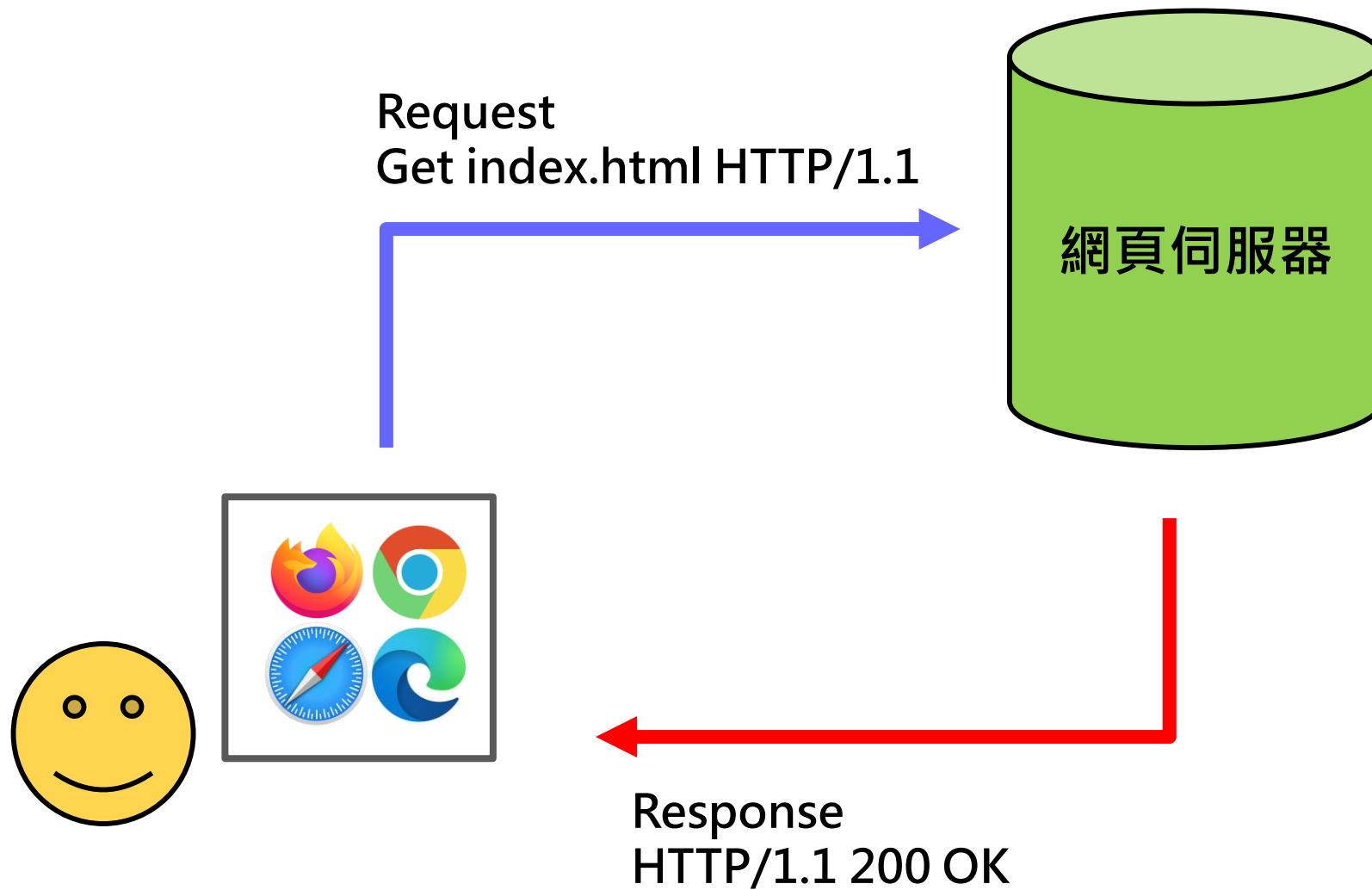
網路爬蟲的應用

- 網路輿情分析
 - 針對某一題進行群眾意見的分析





平常是如何取得網頁的內容





瀏覽器送出的請求

- 模擬瀏覽器送出請求的方法
- 但在這之前要先了解瀏覽器送了什麼請求

The screenshot shows the National Tsing Hua University (NTHU) website. The header includes the NTHU logo and name in Chinese and English. A navigation menu lists various links like '網站導覽', '通訊錄', '捐款', 'English', '在校生', 'International Students', '僑生', '陸生', '教職員', '校友', '未來學生', '訪客', '首頁故事', '清華簡訊', 'Newsletter', and social media icons. A purple banner below the header contains text about admission and a video premiere. The Chrome DevTools Network tab is open, showing a list of requests. The first request is 'www.nthu.edu.tw' (GET, 200, 40.7 kB) and the second is '4H6mthdcO_8' (GET, 200, 11.1 kB). The bottom status bar shows '2 / 90 requests', '51.8 kB / 142 kB transferred', '78.0 kB / 7.1 MB resources', 'Finish: 12.92 s', 'DOMContentLoaded: 463 ms', and 'Load: 1.33 s'.

Name	Method	Status	Protocol	Type	Initiator	Size	Time
www.nthu.edu.tw	GET	200	http/1.1	document	Other	40.7 kB	
4H6mthdcO_8	GET	200	h2	document	(index)	11.1 kB	



Method Get 與 Post

- Get像是明信片
 - 網頁傳遞的參數可以透過網址發現

原始網址：

https://www.ptt.cc/bbs/C_Chat/index.html

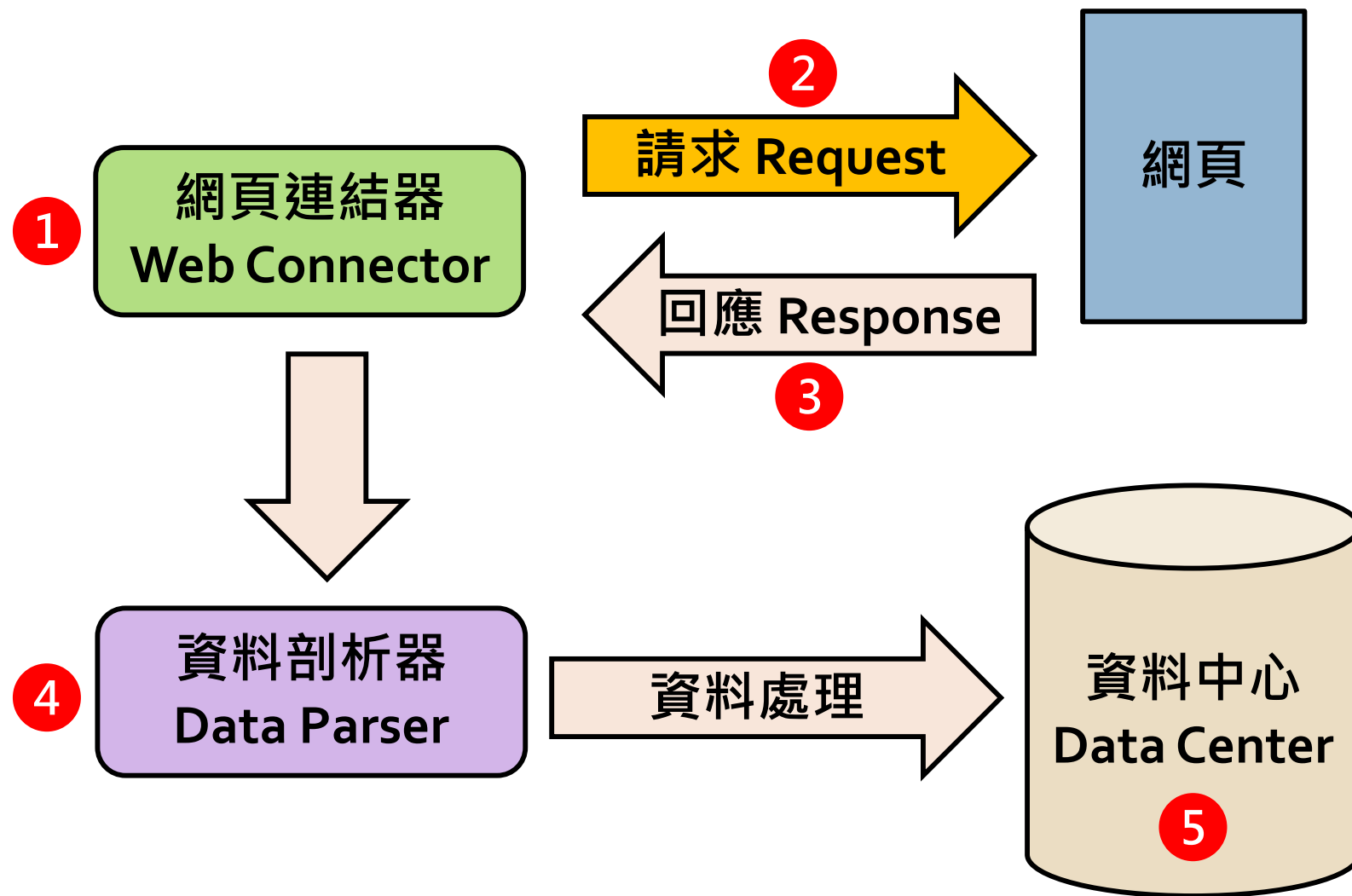
經過查詢：

https://www.ptt.cc/bbs/C_Chat/search?q=%E7%B6%B4%E6%AD%8C

- Post像是信封
 - 網頁傳遞的參數看不到
- 本次實作範例要爬取的網站是PTT表特版
 - Method屬於Get
 - 只介紹Python抓取Method為Get的網站



使用Python撰寫網路爬蟲的步驟





為什麼要用 Python實作網 路爬蟲

- 有各種爬蟲框架
 - Scrapy
 - urllib2
 - BeautifulSoup
 - lxml
- 多線程、進程模型成熟穩定
 - 爬蟲是典型的多任務處理 → 需要等待
 - 多線程、進程可以優化效率
- 抓取下來的資料處理
 - 分詞處理
 - NLTK
 - Jieba
 - CKIP



在爬蟲前的 注意事項

- 要檢視目標網站的規範「robots.txt」
- Robots協議
 - Robots Exclusion Protocol
 - 爬蟲協議、機器人協議
 - 網站通過此協議告訴搜尋引擎只有哪些頁面可以抓取
 - 這只是一種互相尊重的協議



使用Python實 作網路爬蟲

PTT表特版 圖片下載

▼ 爬蟲的前置準備

```
[ ] 1 #匯入常用的套件  
    2 import requests  
    3 import os  
    4 import re  
    5 import time  
    6 import urllib.request as ur  
    7 from bs4 import BeautifulSoup
```

```
[ ] 1 #定義全域變數  
    2 file_path = '/content/gdrive/My Drive/TA/1091/crawler/'  
    3 ptt_url_head = 'https://www.ptt.cc'
```



取得網頁

▼ 定義爬蟲所需要用到的函式

```
[ ] 1 #取得網頁
    2 def get_web_page(url):
    3     #取得網頁的回應
    4     response = requests.get(url=url, cookies={'over18':'1'})
    5
    6     #如果回應的狀態碼是200以外的
    7     #代表網頁有問題
    8     if response.status_code != 200:
    9         print('Invalid url' + url)
    10        return None
    11    else:
    12        #如果回應的狀態碼是200代表網頁正常
    13        #回傳網頁的所有文字。
    14        return response.text
```



關於HTTP Status Code

- HTTP狀態碼
 - 用來表示網頁伺服器HTTP回應狀態的3位數代碼
 - 常見的代碼與其意義如下：
 - 200 OK
 - 403 Forbidden
 - 404 Not Found
 - 502 Bad Gateway
- 其他的狀態碼可以參考以下網站：

[HTTP狀態碼 - 維基百科，自由的百科全書](#)



response.txt

```
<!DOCTYPE html>
<html>
  <head>
    <meta charset="utf-8">

    <meta name="viewport" content="width=device-width, initial-scale=1">

    <title>看板 Beauty 文章列表 - 批踢踢實業坊</title>

    <link rel="stylesheet" type="text/css" href="//images.ptt.cc/bbs/v2.27/bbs-common.css">
    <link rel="stylesheet" type="text/css" href="//images.ptt.cc/bbs/v2.27/bbs-base.css" media="screen">
    <link rel="stylesheet" type="text/css" href="//images.ptt.cc/bbs/v2.27/bbs-custom.css">
    <link rel="stylesheet" type="text/css" href="//images.ptt.cc/bbs/v2.27/pushstream.css" media="screen">
    <link rel="stylesheet" type="text/css" href="//images.ptt.cc/bbs/v2.27/bbs-print.css" media="print">

    </head>
    <body>

      <div id="topbar-container">
        <div id="topbar" class="bbs-content">
          <a id="logo" href="/bbs/">批踢踢實業坊</a>
          <span>&rsquo;</span>
          <a class="board" href="/bbs/Beauty/index.html"><span class="board-label">看板 </span>Beauty</a>
          <a class="right small" href="/about.html">關於我們</a>
          <a class="right small" href="/contact.html">聯絡資訊</a>
        </div>
      </div>

      <div id="main-container">
        <div id="action-bar-container">
          <div class="action-bar">
            <div class="btn-group btn-group-dir">
              <a class="btn selected" href="/bbs/Beauty/index.html">看板</a>
            </div>
          </div>
        </div>
      </div>
    </body>
  </html>
```

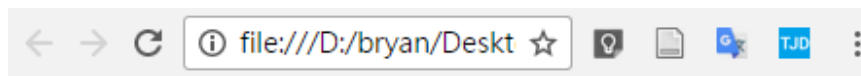


要如何取得需要的資料

- 網頁是由標籤 (Tag) 所組成的階層式文件
- 主要由3個部分構成：
 - HTML (骨架)
 - CSS (階層樣式表，樣式)
 - JavaScript(負責與使用者互動的程式功能)
- 可以透過標籤與相關屬性去定位資料的位置



階層式文件



我是變色且置中的抬頭

我是段落一

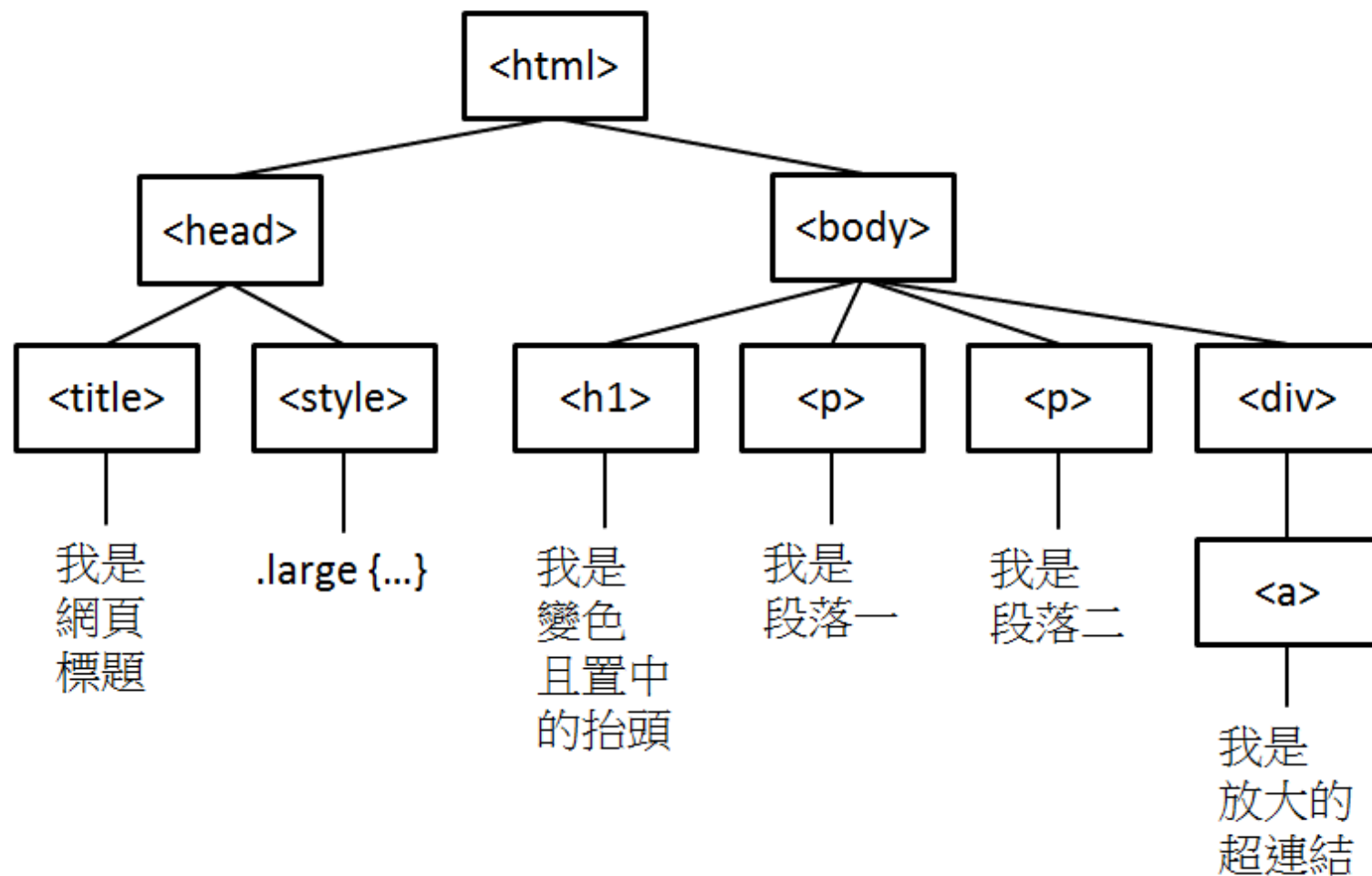
我是段落二

我是放大的超連結

```
test.html x
1 <html>
2   <head>
3     <title>我是網頁標題</title>
4     <style>
5       .large {
6         color:blue;
7         text-align: center;
8       }
9     </style>
10  </head>
11  <body>
12    <h1 class="large">我是變色且置中的抬頭</h1>
13    <p id="p1">我是段落一</p>
14    <p id="p2" style="">我是段落二</p>
15    <div><a href='https://www.cycu.edu.tw/' style="font-size:200%;">我是放大的超連結</a></div>
16  </body>
17 </html>
```



網頁架構





CSS Selector

- CSS
 - 階層式樣式表，Cascading Style Sheets。
 - 用來位網頁添加樣式的電腦語言
- CSS選擇器
 - 可以根據class或是id，將符合的標籤套上樣式。
 - class以「.»作為開頭
 - id以「#」作為開頭

選擇器	範例	意思
.class	.large	選擇所有class=".large"的標籤
#id	#name	選擇id="name"的標籤



BeautifulSoup

- 使用BeautifulSoup進行解析
- 解析完畢後可以使用以下的函式去定位資料區塊：
 - **find()**
 - **find_all()**
 - **select()**

```
1 #定義最一開始的網頁超連結
2 ptt_beauty_url = 'https://www.ptt.cc/bbs/Beauty/index.html'
3 page = get_web_page(ptt_beauty_url)
4
5 #解析
6 soup = BeautifulSoup(page, 'html.parser')
7 #使用 find() 尋找 body 標籤的文字
8 print(soup.find('body').text)
```

批踢踢實業坊
>
看板 Beauty
關於我們
聯絡資訊



取得想要的 資料

```
[ ] 1 #取得網頁的內容
    2 def get_articles(web_page, date):
    3     #解析網頁
    4     soup = BeautifulSoup(web_page, 'html.parser')
    5
    6     #取得上一頁按鈕的超連結
    7     a_tags = soup.find_all('a', 'btn wide')
    8     previous_page_url = a_tags[1]['href']
    9
    10    #用來儲存取得的文章資料
    11    articles = []
    12
    13    divs = soup.find_all('div', 'r-ent')
    14    for d in divs:
    15        #符合條件而且超連結存在，代表文章沒有被刪除。
    16        if d.find('div', 'date').text == date and d.find('a'):
    17
    18            #取得文章的超連結
    19            href = d.find('a')['href']
    20
    21            #取得文章的網頁
    22            content = get_web_page(ptt_url_head + href)
    23            content_soup = BeautifulSoup(content, 'html.parser')
    24            content_divs = content_soup.find_all('div', id='main-container')
    25
    26            #用來儲存文章標題
    27            title = ''
    28
    29            #用來儲存要下載的圖片超連結清單
    30            image_url_list = []
```



取得想要的 資料 - 接續上頁

```
31
32     #開始進行取得
33     for cd in content_divs:
34         #取得文章標題
35         content_spans = cd.find_all('span', 'article-meta-value')
36         #去除取得的文章標題的多餘空白
37         title = content_spans[2].text.strip()
38
39         #取得要下載的圖片超連結清單
40         content_a_hrefs = cd.find_all('a', {'href': re.compile('https://(imgur|i\.imgur)\.com/.+')})
41         for cah in content_a_hrefs:
42             #cah的長度不是0代表有取得圖片的超連結
43             if len(cah) != 0:
44                 #這裡必須處理網址最後面沒有.jpg的情況
45                 if cah.text.find('.jpg') == -1:
46                     image_url_list.append(cah.text+'.jpg')
47                 else:
48                     image_url_list.append(cah.text)
49
50         #儲存取得的文章標題與下載圖片的超連結清單
51         articles.append({'title': title, 'image_url_list': image_url_list})
52
53     #回傳文章資料與上一頁的超連結
54     return articles, previous_page_url
```



儲存圖片

```
[ ] 1 #儲存圖片
    2 def save_images(folder_name, image_url_list):
    3     #將資料夾路徑與名稱串接起來
    4     folder_path = os.path.join(file_path, folder_name)
    5
    6     #創建資料夾
    7     os.makedirs(folder_path)
    8
    9     print('-----')
   10     print('Folder name:', folder_name)
   11
   12     #使用例外處理不可預知的例外發生
   13     try:
   14         for image_url in image_url_list:
   15             #取得圖片超連結的最後一段做為儲存圖片的檔名
   16             #例如kes83e.jpg
   17             file_name = image_url.split('/')[-1]
   18
   19             #透過超連結下載圖片
   20             ur.urlretrieve(image_url, os.path.join(folder_path, file_name))
   21             print('{} saved...'.format(file_name))
   22     except Exception as e:
   23         #發生例外時印出例外
   24         print(e)
```



開始執行 爬蟲

▾ 開始執行爬蟲

```
[ ] 1 #定義最一開始的網頁超連結
2 ptt_beauty_url = 'https://www.ptt.cc/bbs/Beauty/index.html'
3
4 #today是爬蟲的條件
5 #意思是只爬取指定日期的文章
6 today = time.strftime('%m', time.localtime()) + time.strftime('/%d', time.localtime())
7 print('-----')
8 print('Today is', today)
9
10 #開始進行爬蟲
11 print('-----')
12 print('Start')
13 while True:
14     #呼叫get_web_page()函式取得網頁的所有文字
15     page = get_web_page(ptt_beauty_url)
16
17     #如果page不是None就進行爬蟲
18     if page:
19         #呼叫get_articles()函式取得文章與上一頁的超連結
20         articles, previous_page_link = get_articles(page, today)
21
22         #如果articles不是None就下載圖片
23         if articles:
24             for article_content in articles:
25                 #文章標題作為資料夾名稱
26                 title = article_content['title']
27
28                 #取得要下載的圖片超連結清單
29                 image_url_list = article_content['image_url_list']
30
31                 #呼叫save_images()函式下載圖片
32                 save_images(title, image_url_list)
33
34                 #更新網址進行下一輪的爬蟲
35                 ptt_beauty_url = ptt_url_head + previous_page_link
36         else:
37             #如果articles是None代表沒有符合條件的文章了
38             #此時可以停止爬蟲
39             break
40 print('-----')
41 print('End')
```




執行結果

Today is 10/02

Start

Folder name: [正妹] 奶大還是臉正？抖幾？

WgFpFuu.jpg saved...

Ib7ImPt.jpg saved...

1yl49ib.jpg saved...

4z32Z7s.jpg saved...

Folder name: [正妹] 攝影師真的很厲害

n0abs1H.jpg saved...

45C9g6Y.jpg saved...

Sgupv6w.jpg saved...

HDyYZAB.jpg saved...

z8vktSk.jpg saved...

因為印出來的訊息太長了，
所以中間的訊息就不貼出來給
大家看。

oNunG1k.jpg saved...

Qc1lb9t.jpg saved...

Qv9DhK0.jpg saved...

jw9D7uz.jpg saved...

Folder name: [正妹] 畫漫畫的

02C0PTv.jpg saved...

rgUT31T.jpg saved...

l5zqNGz.jpg saved...

TJh0ZHv.jpg saved...

59ajKeC.jpg saved...

Pxg7Y9r.jpg saved...

oTdyvLG.jpg saved...

lOX8ulz.jpg saved...

End



執行結果

我的雲端硬碟 > TA > 1091 > crawler ▾ 👤



資料夾

名稱 ↑

👤 [正妹] 乃木坂46 早川聖...

👤 [正妹] 大尺碼 | 肉特(18...

👤 [正妹] 小臉鳥仔腳

👤 [正妹] 奶大還是臉正? ...

👤 [正妹] 射了

👤 [正妹] 烤肉之夜後來點...

👤 [正妹] 馬來西亞臭臉小...

👤 [正妹] 清新又氣質

👤 [正妹] 畫漫畫的

👤 [正妹] 媽咪的愛

👤 [正妹] 愛德琳

👤 [正妹] 葛佳慧 ジュリ

👤 [正妹] 漂亮臉孔修長身...

👤 [正妹] 瑪麗亞

👤 [正妹] 攝影師真的很厲害

👤 [正妹] 櫻

👤 [帥哥] 室剛

👤 [神人] 求神ig

👤 [神人] 清秀可人的妹

檔案



🔗 crawler_ptt_beauty.ipynb



執行結果

我的雲端硬碟 > ... > crawler > [正妹] 乃木坂46 早川聖來 ▾ 👤



檔案

名稱 ↑



0tzegOw.jpg



3rWYw7E.jpg



8wurrWQ.jpg



CgNv8gT.jpg



CwyXWus.jpg



DCSJ9Po.jpg



ExnJw41.jpg



GeGKPRX.jpg



mZyfSli.jpg



NNyp8RP.gif.jpg

