

# 網頁資料擷取

擷取網頁資料的前題是不能觸犯著作權：

- <https://www.tipo.gov.tw/ct.asp?xItem=219598&ctNode=7561&mp=1> (<https://www.tipo.gov.tw/ct.asp?xItem=219598&ctNode=7561&mp=1>)
- <https://udn.com/news/story/6871/3221682> (<https://udn.com/news/story/6871/3221682>)

這個單元例子僅用來說明由網頁擷取資料的入門技巧。

- 使用re
- 使用BeautifulSoup
- 使用Selenium

## 解析網址

單純網址

In [ ]:

```
1 # coding=utf-8
2 from urllib.parse import urlparse
3 url = 'https://tw.stock.yahoo.com/news_list/url/d/e/'
4 up = urlparse(url)
5 print(up)
```

In [ ]:

```
1 # coding=utf-8
2 from urllib.parse import urlparse
3 url = 'https://www.cwb.gov.tw/V7/forecast/index.htm'
4 up = urlparse(url)
5 print(up)
```

有get參數

In [ ]:

```
1 # coding=utf-8
2 from urllib.parse import urlparse
3 url = 'https://tw.stock.yahoo.com/q/q?s=2330'
4 up = urlparse(url)
5 print(up)
```

In [ ]:

```
1 # coding=utf-8
2 from urllib.parse import urlparse
3 url = 'https://ecshweb.pchome.com.tw/search/v3.3/?q=pc&scope=all'
4 up = urlparse(url)
5 print(up.query.split('&'))
```

- **scheme**: 通訊協定
- **netloc**: 網域名稱
- **path**: 網頁所在路徑與檔名
- **query**: GET參數

---

## 透過requests.get(url)擷取網頁的內容

```
import requests  
r = requests.get(url)
```

註: 萬一被阻擋, 可以嘗試設定user agent

```
headers = {'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; WOW64) AppleWebKit/537.36  
(KHTML, like Gecko) Chrome/71.0.3578.80 Safari/537.36',  
           'Content-Type': 'application/x-www-form-urlencoded',  
           'Connection': 'Keep-Alive',  
           'Accept': 'text/html,application/xhtml+xml,application/xml;q=0.9,*/*;q=0.8',  
           'Accept-Encoding': 'gzip, deflate, sdch',  
           'Accept-Language': 'en-US,en;q=0.8'}  
  
r = requests.get(url, headers=headers)
```

你可以從chrome瀏覽器找到各式user agent字串: 按Ctrl-Shift-I啟動開發者工具, 在Network選單, 按右上選項圖示後選More tools>Network conditions。

## 在URL裡傳遞參數

In [ ]:

```
1 import requests  
2 keys = {'search_query': 'python'}  
3 r=requests.get('https://www.youtube.com/results',params=keys)  
4 print(r.url)
```

## 擷取網頁的內容

方式一: **r.content** (它是bytes型態)

可透過下面方式將bytes轉str

- 轉換方法一: `r.content.decode('utf-8')`
- 轉換方法二: `str(r.content,encoding='utf-8')`

str轉bytes

- 轉換方法一: `s.encode('utf-8')`
- 轉換方法二: `bytes(s, encoding='utf-8')`

In [ ]:

```
1 import requests
2 url = 'https://www.cwb.gov.tw/V7/forecast/index.htm'
3 r = requests.get(url)
4 for i in r.content.decode('utf-8').splitlines()[1:10]:
5     print(i)
```

方式二: **r.text** (它是**str**型態)

如下面所示，**requests**編碼方式為ISO-8859-1。

In [ ]:

```
1 import requests
2 url = 'https://www.cwb.gov.tw/V7/forecast/index.htm'
3 #html = requests.get(url).content.decode('utf-8', 'ignore').splitlines()
4 r = requests.get(url)
5 print(r.encoding)
```

用**r.encoding='utf-8'**將編碼改為'utf-8'編碼後，處理網頁內容。

In [ ]:

```
1 import requests
2 url = 'https://www.cwb.gov.tw/V7/forecast/index.htm'
3 r = requests.get(url)
4 r.encoding = 'utf-8'
5 for i in r.text.splitlines()[1:10]:
6     print(i)
```

## 使用**re**模組擷取資訊

### HTML網頁格式

HTML標記	用途說明
<html>...</html>	標記...為HTML文件
<head>...</head>	標記...為HTML文件標頭
<title>...</title>	標記...為HTML文件標題，通常會顯示在瀏覽器標題列
<body>...</body>	標記...為HTML文件內容
<script>...</script>	標記...為描述語言
<h1>...</h1>	標記...為標題(等級為h1,...,h6)
<p>...</p>	標記...為文字段落
<div>...</div>	排版用格式標記，...通常為內文大段落或顯示分塊
<span>...</span>	類似<div>，通常用在小段落

<code>&lt;table&gt;...&lt;/table&gt;</code>	標記...為表格呈現內容
<code>&lt;img src='...'&gt;</code>	顯示圖形檔設定
<code>&lt;a href='...'&gt;</code>	外部連結設定

例子:使用`re`擷取網頁裡的新聞標題

<https://udn.com/news/cate/2/7226> (<https://udn.com/news/cate/2/7226>) 網頁裡新聞標題寫在`<h2> ... </h2>`段落，如下範例

```
<h2 style="width:100%">專利戰高通告贏蘋果 陸將禁售iPhone X以前機種 <time>22:29</time></h2>
<h2>先從海外版施行 日本漫畫週刊少年Jump也走向數位訂閱制<span class="i-video1"></span></h2>
```

In [ ]:

```
1 import requests
2 import re
3 url = 'https://udn.com/news/cate/2/7226'
4 html = requests.get(url).content.decode('utf-8')
5 for idx,title in enumerate(re.finditer(r'<h2[^>]*?>([<]*?)(<time.+?/time>)?(<span.+?/span>)?</h2>')):
6     print('{:4d}. {}'.format(idx+1,title.group(1)))
```

## 使用BeautifulSoup模組擷取資訊

<https://www.crummy.com/software/BeautifulSoup/bs4/doc/>  
(<https://www.crummy.com/software/BeautifulSoup/bs4/doc/>)

安裝

```
conda install -c anaconda beautifulsoup4
```

找到所有`<h2> ... </h2>`段落程式片段

```
from bs4 import BeautifulSoup
import requests
url = 'https://udn.com/news/cate/2/7226'
html = requests.get(url).content.decode('utf-8')
sp = BeautifulSoup(html,'html.parser')
for link in sp.find_all('h2'):
    print('{}'.format(link.text))
```

指令範例	說明
<code>sp.find('a',key))</code>	傳回第一個符合的內容
<code>sp.find_all('a',key))</code>	傳回所有符合的內容
<code>sp.title/sp.title.text</code>	傳回<title>網頁標題</title>

---

`sp.text` 傳回去掉HTML標籤的內容

例子:使用**BeautifulSoup**擷取網頁裡的新聞標題

In [ ]:

```
1 from bs4 import BeautifulSoup
2 import requests
3 url = 'https://udn.com/news/cate/2/7226'
4 html = requests.get(url).content.decode('utf-8')
5 sp = BeautifulSoup(html, 'html.parser')
6
7 for idx, link in enumerate(sp.find_all('h2')):
8     print('{:4d}. {}'.format(idx+1, link.text))
```

例子：擷取表格資料

中央氣象局月平均氣溫表格HTML原始碼:

[https://www.cwb.gov.tw/V7/climate/monthlyMean/Taiwan\\_tx.htm](https://www.cwb.gov.tw/V7/climate/monthlyMean/Taiwan_tx.htm)  
([https://www.cwb.gov.tw/V7/climate/monthlyMean/Taiwan\\_tx.htm](https://www.cwb.gov.tw/V7/climate/monthlyMean/Taiwan_tx.htm))

```
<table width="780" cellpadding="2" cellspacing="1" class="Form00" summary="排版用表格">
```

```
<tbody><tr height="44">
  <th width="150" height="44" class="tab01">地名</th>
  <th width="50" height="44" class="tab01" axis="month">一月</th>
  <th width="50" height="44" class="tab01" axis="month">二月</th>
  <th width="50" height="44" class="tab01" axis="month">三月</th>
  <th width="50" height="44" class="tab01" axis="month">四月</th>
  <th width="50" height="44" class="tab01" axis="month">五月</th>
  <th width="50" height="44" class="tab01" axis="month">六月</th>
  <th width="50" height="44" class="tab01" axis="month">七月</th>
  <th width="50" height="44" class="tab01" axis="month">八月</th>
  <th width="50" height="44" class="tab01" axis="month">九月</th>
  <th width="50" height="44" class="tab01" axis="month">十月</th>
  <th width="50" height="44" class="tab01" axis="month">十一月</th>
  <th width="50" height="44" class="tab01" axis="month">十二月</th>
  <th width="50" height="44" class="tab01">平均</th>
  <th width="100" height="44" class="tab01">統計期間</th>
</tr>
<tr height="44">
  <td height="44" class="active" axis="item">淡水</td>
  <td class="whitetd" width="36" height="44">15.2</td>
  <td class="whitetd" width="36" height="44">15.6</td>
  <td class="whitetd" width="36" height="44">17.4</td>
  <td class="whitetd" width="36" height="44">21.1</td>
  <td class="whitetd" width="36" height="44">24.5</td>
  <td class="whitetd" width="36" height="44">26.9</td>
  <td class="whitetd" width="36" height="44">28.8</td>
  <td class="whitetd" width="36" height="44">28.6</td>
  <td class="whitetd" width="36" height="44">26.7</td>
  <td class="whitetd" width="36" height="44">23.7</td>
  <td class="whitetd" width="36" height="44">20.6</td>
  <td class="whitetd" width="36" height="44">16.9</td>
  <td class="whitetd" width="36" height="44">22.2</td>
  <td class="whitetd" width="90" height="25">1981-2010</td>
</tr>
...
</table>
```

平均氣溫表格放在<table></table>標記內。可是那個網頁有許多<table>標記，不過其中class為Form00為平均氣溫表格。因此table = sp.find\_all('table',{ 'class': 'Form00'})可以鎖定以'class'為key, value為'Form00'的那個<table>標記。

```
url = 'https://www.cwb.gov.tw/V7/climate/monthlyMean/Taiwan_tx.htm'
html = requests.get(url).content.decode('utf-8')
sp = BeautifulSoup(html, 'html.parser')
table = sp.find('table', { 'class': 'Form00' })
```

表格裡每一列以<tr>，</tr>標記標示。下面指令找出<table></table>標記內每一列。

```
rows = table.find_all('tr')
```

第一列，為標題列，每一欄以<th>，</th>標記標示。下面指令找出<tr>，</tr>標記內每一欄。

```
title = [c.text for c in rows[0].find_all('th')]
```

其他列每一欄以<td>，</td>標記標示。下面指令找出<tr>，</tr>標記內每一欄，並存放在各自list。

```
data = [list() for _ in range(len(title))]  
  
for r in rows[1:]:  
    for col,cell_data in zip(data,r.find_all('td')):  
        try:  
            col.append(float(cell_data.text))  
        except ValueError:  
            col.append(cell_data.text)
```

In [ ]:

```
1 from bs4 import BeautifulSoup  
2 import requests  
3 import numpy as np  
4 import matplotlib.pyplot as plt  
5  
6 url = 'https://www.cwb.gov.tw/V7/climate/monthlyMean/Taiwan_tx.htm'  
7 html = requests.get(url).content.decode('utf-8')  
8 sp = BeautifulSoup(html,'html.parser')  
9 table = sp.find('table',{'class':'Form00'})  
10 rows = table.find_all('tr')  
11  
12 title = [c.text for c in rows[0].find_all('th')]  
13 data = [list() for _ in range(len(title))]  
14  
15 for r in rows[1:]:  
16     for col,cell_data in zip(data,r.find_all('td')):  
17         try:  
18             col.append(float(cell_data.text))  
19         except ValueError:  
20             col.append(cell_data.text)  
21  
22 #放入 numpy.ndarray  
23  
24 data_table= np.core.records.fromarrays(data)  
25 data_table.dtype.names = title  
26  
27 #資料標題  
28 print(data_table.dtype.names)  
29  
30 #取得第0列資料  
31 print(data_table[0])  
32  
33 #取得各觀測站五月均溫  
34 print(data_table['五月'])  
35
```

使用matplotlib繪出平均氣溫圖。

In [ ]:

```
1 import matplotlib.pyplot as plt
2 from matplotlib.font_manager import FontProperties
3 import os
4
5 font = FontProperties(fname=os.environ['WINDIR']+'\\Fonts\\kaiu.ttf', size=12)
6
7 plt.figure(figsize=(8,4))
8 for i in range(5):
9     r = list(data_table[i])
10    plt.plot(np.arange(len(data_table.dtype.names)-3),r[1:-2],label=r[0])
11 plt.legend(prop=font)
12 plt.title('平均氣溫',fontproperties=font)
13 plt.xticks(np.arange(len(title)-3),title[1:-2],fontproperties=font)
14 plt.ylabel('攝氏',fontproperties=font)
15 plt.show()
```

例子：擷取所有連結

HTML連結格式為：

```
<a href='https://tw.yahoo.com/?p=us'>本文</a>
```

所以

```
all_links = sp.find_all('a')
```

得到所有<a ....>....</a>段落。假設link為以上連結為例子，

- link.get('href')得到'<https://tw.yahoo.com/?p=us>' (<https://tw.yahoo.com/?p=us>)
- link.text得到'本文'

In [ ]:

```
1 from bs4 import BeautifulSoup
2 import requests
3 url = 'https://udn.com/news/index'
4 html = requests.get(url).content.decode('utf-8')
5 sp = BeautifulSoup(html,'html.parser')
6
7 for idx,link in enumerate(sp.find_all('a')):
8     href = link.get('href')
9     if href is not None and href.startswith('http'):
10         print('{:4d} text:{:<s}, link:{:>s}'.format(idx+1,link.text,href))
```

例子：擷取所有圖形檔

下面例子需要用到pillow模組

```
conda install -c anaconda pillow
```



In [ ]:

```
1 from bs4 import BeautifulSoup
2 import requests
3 from urllib.parse import urlparse
4 from urllib.request import urlopen
5 import matplotlib.pyplot as plt
6 from PIL import Image, ImageDraw, ImageFont
7 import numpy as np
8 import io
9 import re
10
11 url = 'https://udn.com/news/story/7934/3526132'
12
13 headers = {'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; WOW64) AppleWebKit/537.36 (KHTML
14             'Content-Type': 'application/x-www-form-urlencoded',
15             'Connection' : 'Keep-Alive',
16             'Accept': 'text/html,application/xhtml+xml,application/xml;q=0.9,*/*;q=0.8',
17             'Accept-Encoding': 'gzip, deflate, sdch',
18             'Accept-Language': 'en-US,en;q=0.8'}
19
20 uc = urlparse(url)
21 print(uc.scheme,uc.hostname)
22 domain = '{}://{}'.format(uc.scheme,uc.hostname)
23 print(domain)
24
25 html = requests.get(url,headers=headers).content.decode('utf-8')
26
27 sp = BeautifulSoup(html,'html.parser')
28 for idx,link in enumerate(sp.find_all(['a','img'])):
29     href = link.get('href')
30     src = link.get('src')
31     for t in [href, src]:
32         if t is not None and ('.jpg' in t or '.png' in t):
33             if t.startswith('http'):
34                 img_path = t
35             elif t.startswith('///'):
36                 img_path = 'https:'+t
37             else:
38                 domain+t
39             print(img_path)
40             print('filename:{}'.format(re.search('[^/]+((.jpg)|(.png))',img_path).group(1)))
41             image = urlopen(img_path)
42             img = Image.open(image)
43             plt.imshow(img)
44             plt.axis('off')
45             plt.show()
46
```

## 透過Selenium擷取網頁資料

安裝

```
conda install -c conda-forge selenium
```

安裝各瀏覽器的webdriver載點<https://www.seleniumhq.org/about/platforms.jsp>  
(<https://www.seleniumhq.org/about/platforms.jsp>)

- 如Chrome web driver:<https://sites.google.com/a/chromium.org/chromedriver/home>  
(<https://sites.google.com/a/chromium.org/chromedriver/home>)

並將webdriver(如chromedriver.exe)放在Python執行的目錄內。測試下面範例：

## 操作瀏覽器函式

webdriver方法	說明
refresh()	重新整理頁面
back()	回上一頁
forward()	到下一頁
close()	關視窗
quit()	結束瀏覽器
get(url)	瀏覽url這網址
current_url	目前網址
title	網頁標題
page_source	網頁原始碼
save_screenshot(pngfile)	存目前網頁畫面於png檔
get_window_position()	取得視窗左上角位置
set_window_position(x,y)	設定視窗左上角位置
maximize_window()	最大化視窗
get_window_size()	取得視窗大小
set_window_size(x,y)	設定視窗大小

In [ ]:

```
1 from selenium import webdriver
2 urls = ['https://www.cwb.gov.tw/V7/', 'https://tw.yahoo.com/?p=us']
3 web = webdriver.Chrome()
4 for idx,url in enumerate(urls):
5     web.get(url)
6     web.save_screenshot('screenshot_{}.png'.format(idx))
7 web.quit()
```

In [ ]:

```
1 import matplotlib.pyplot as plt
2 import matplotlib.image as mpimg
3 fig = plt.figure(figsize=(15,30))
4 for idx in range(len(urls)):
5     img = mpimg.imread('screenshot_{}.png'.format(idx))
6     plt.subplot(1,2,idx+1)
7     plt.imshow(img)
8     plt.axis('off')
9 plt.show()
```

## 網頁元素檢索功能

## 基本步驟範例

```
web = webdriver.Chrome()
web.get(url)

# 檢索網頁元素，通常元素id為唯一，比較好找
element = web.find_element_by_id(id)

# 操作網頁元素
element.send_keys(value)
element.submit() # 提交
```

webdriver元素檢索方法	說明
<b>find_element_by_X(value)</b>	使用X檢索，取得第一個符合的元素
-----	-----
find_element_by_class_name(name)	使用類別名稱檢索
find_element_by_css_selector(selector)	使用CSS選擇器檢索
find_element_by_id(id)	使用id檢索
find_element_by_link_text(text)	使用連結文字檢索
find_element_by_name(name)	使用名稱檢索
find_element_by_tag_name(name)	使用HTML標籤檢索
-----	-----
<b>find_elements_by_X</b>	使用X檢索，取得所有符合的元素
<hr/>	
webdriver元素操作方法	說明
clear()	清除內容
click()	點擊，通常用於按鈕、連結、選單
send_keys(value)	對此元素送出字串
submit()	提交
is_displayed()	此元素是否可見
is_enabled()	此元素是否可用
is_selected()	此元素是否被選定

In [ ]:

```
1 from selenium import webdriver
2
3 web = webdriver.Chrome()
4 web.maximize_window()
5 web.get("https://www.google.com")
6
7 #找到輸入框
8 element = web.find_element_by_name("q")
9
10 #輸入
11 element.send_keys("中央氣象局")
12
13 #提交
14 element.submit()
15
16 #web.close()
```

In [ ]:

```
1 from selenium import webdriver
2 import time
3
4 web = webdriver.Chrome()
5 web.maximize_window()
6 web.get("https://www.youtube.com")
7
8 #找到輸入框
9 element = web.find_element_by_id("search")
10
11 #輸入
12 element.send_keys("selenium Python")
13
14 #按搜尋
15 search_btn = web.find_element_by_id("search-icon-legacy")
16 search_btn.click()
17 # Get scroll height
18 last_height = -1
19 for idx in range(500):
20     # Scroll down to bottom
21     web.execute_script("window.scrollTo(0, window.scrollY + 800);")
22     # Wait to load page
23     time.sleep(.5)
24     current_height = web.execute_script("return window.scrollY")
25     if last_height == current_height:
26         print('stop')
27         break
28     last_height = current_height
29
```

In [ ]:

```
1 print(web.page_source.splitlines()[:5])
```

In [ ]:

```
1
```

