

程式設計 (二)

selenium fb 爬蟲實作

Ming-Hung Wang 王銘宏

tonymhwang@cs.ccu.edu.tw

Department of Computer Science and Information Engineering
National Chung Cheng University

Spring Semester, 2022

本章目錄

1. 介紹
2. 前置作業
3. 函式介紹
4. facebook 爬蟲
5. 總結

什麼是 selenium

- selenium 可以開啟瀏覽器，並根據使用者編寫的程式碼，模擬對瀏覽器的操作（輸入、點擊……等）。
- selenium 經常搭配 beautifulsoup 使用，以解析進行操作後的網頁資料。

前置作業

安裝 selenium、beautifulsoup。

pip install selenium

pip install bs4

前置作業

下載 webdriver。

ChromeDriver

前置作業

查看 Chrome 版本，並下載對應版本的 webdriver。

關於 Chrome

 Google Chrome

 Chrome 目前是最新版本
版本 99.0.4844.84 (正式版本) (64 位元)

前往 Chrome 說明頁面 

回報問題 

selenium 函式介紹

- **driver = webdriver.Chrome(path)**
開啟 webdriver。
- **driver.get(url)**
操作瀏覽器進入網頁。

selenium 函式介紹

- **`ele = driver.find_element_by_class_name(value)`**
尋找第一個符合 `class name = value` 的網頁元素並回傳。
- **`eles = driver.find_elements_by_class_name(value)`**
尋找所有符合 `class name = value` 的網頁元素並回傳。

selenium 尋找網頁元素函式

- `ele = driver.find_element_by_id(value)`
- `ele = driver.find_element_by_name(value)`
- `ele = driver.find_element_by_xpath(value)`

selenium 函式介紹

- **ele.send_keys(value)**
將 value 填入 ele 中。
- **ele.click()**
對 ele 進行點擊。

beautifulsoup 函式介紹

- `soup = BeautifulSoup(text, parser)`
解析網頁內容。

beautifulsoup 函式介紹

- **`ele = soup.find(type, {key : value})`**
回傳第一個符合 `key = value` 的 `type` 網頁元素。
ex. **`ele = soup.find('div', {'id' : 'myid'})`**
- **`ele = soup.find_all(type, {key : value})`**
回傳所有符合 `key = value` 的 `type` 網頁元素。

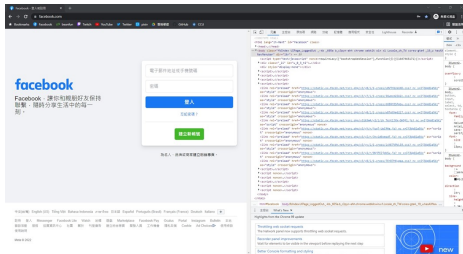
開啟 webdriver

首先設定瀏覽器的選項，將放大視窗加入選項，然後開啟 webdriver。

```
options = webdriver.ChromeOptions()  
options.add_argument("--start-maximized")  
  
driver = webdriver.Chrome('Crawler/chromedriver.exe', options = options)
```

登入

首先查看登入頁面的 html 碼。



登入

尋找輸入帳號密碼的網頁元素。

```
<form class="bvtf" data-testid="royal_login_form" action="/login/privacy_mutation_tokenkey70e0X11jow.C2jcnvhd1v198w011jowhJQ3MDv18QulC3jV0sc210ZV9a2C18tgeHjI5V0c5Mf.c10TQ3vQK30K30" method="post" onsubmit id="u_8_a_6a">
  <input type="hidden" name="jazoest" value="2940" autocomplete="off">
  <input type="hidden" name="lsd" value="AVqp64e148" autocomplete="off">
  <div>
    <div class="_810e">
      <input type="text" class="inputtext _55r1_810e" name="email" id="email" data-testid="royal_email" placeholder="電子郵件地址或手機號碼" autofocus="1" aria-label="電子郵件地址或手機號碼"/>
    </div>
    <div class="_810e">
      <div class="_55r1_810e" id="passContainer">
        <input type="password" class="inputtext _55r1_810e" name="pass" id="pass" data-testid="royal_pass" placeholder="密碼" aria-label="密碼"/>
        <div class="_5157 hidden_elem" id="u_8_a_8t"></div>
      </div>
    </div>
    <input type="hidden" autocomplete="off" name="login_source" value="csmet_headerless_login">
    <input type="hidden" autocomplete="off" name="next" value="">
    <div class="_810g"></div>
    <div class="_810j"></div>
    <div class="_810c"></div> </div>
    <div class="_810g"></div>
  </form>
```

登入

進入 facebook 首頁。

```
driver.get('https://www.facebook.com/')  
time.sleep(1)
```


登入

找到輸入帳號密碼的網頁元素並填入帳號密碼。

```
context = driver.find_element_by_id('email')
context.send_keys('peter123@gmail.com')
context = driver.find_element_by_id('pass')
context.send_keys('ptr870229')
```

登入

尋找登入按鈕的網頁元素。

```
*<form class="v4t4" data-testid="royal_login_form" action="/login/privacyMutation_token=eyJ0eXkiOiJmZCJmcmh0dGV1dD0wM1Jmcmh0dGV1dDQwLjC3fWxz22Gv9pC1I8Pqg0YjI3fDc3OTc1OTQ2fW50N3R" method="post" onsubmit id="u_0_a_0a">
  <input type="hidden" name="jsoect" value="2948" autocomplete="off">
  <input type="hidden" name="lnd" value="Wqpp44e548" autocomplete="off">
  <div>
    <div class="f1ua">
      <input type="text" class="inputtext_55r1_f1uy" name="email" id="email" data-testid="royal_email" placeholder="電子郵件地址或手機號碼" autofocus="1" aria-label="電子郵件地址或手機號碼">
    </div>
    <div class="f1ua">
      <div class="f1uy_55r1_508" id="passContainer">
        <input type="password" class="inputtext_55r1_f1uy_50pl" name="pass" id="pass" data-testid="royal_pass" placeholder="密碼" aria-label="密碼">
        <div class="51a7 hidden_elem" id="u_0_b_08"></div>
      </div>
    </div>
    <input type="hidden" autocomplete="off" name="login_source" value="comet_headerless_login">
    <input type="hidden" autocomplete="off" name="next" value">
    <div class="51tg">
      <button value="1" class="f0re_f1u0_f1u1_f1u1 selected_51u" name="login" data-testid="royal_login_button" type="submit" id="u_0_a_0a"></button>
    </div>
    <div class="51tj"></div>
    <div class="51ci"></div> <flex>
    <div class="51tg"></div>
  </form>
```

登入

找到登入按鈕的網頁元素並點擊。

```
commit = driver.find_element_by_name('login')  
commit.click()  
time.sleep(1)
```

登入

登入後看到以下通知，會阻止我們對瀏覽器進行操作。



登入

在瀏覽器的選項添加以下設定，可以封鎖通知。

```
prefs = {  
    'profile.default_content_setting_values' : {  
        'notifications': 2  
    }  
}  
options.add_experimental_option('prefs', prefs)
```

爬取粉絲專頁貼文

首先進入粉絲專頁頁面。

由於 facebook 是動態網頁，往下滾動才會顯示更多貼文，所以我們利用 javascript 指令來滾動瀏覽器。

```
driver.get('https://www.facebook.com/appledaily.tw')  
  
driver.execute_script('window.scrollTo(0, document.body.scrollHeight)')  
time.sleep(1)
```

爬取粉絲專頁貼文

找到目前所有貼文。

```
posts = driver.find_elements_by_class_name('du4w351b k4urcfbm l9j0dhe7 sjgh65i0'.replace(' ', '.'))
```

爬取粉絲專頁貼文

如果要讀取貼文全部內容的話，我們需要點擊顯示更多。



爬取粉絲專頁貼文

找到顯示更多的按鈕，這裡有多個網頁元素共用同個 class name，所以利用網頁元素的文字來找到我們想要點擊的按鈕。

```
<div class="oajrlxb2 g5ia77u1 qu0x051f esr5mh6w e9989ue4 r7d6kgcz rq0escxv nhd2j8a9  
zcic4wl gpro0wi8 oo9gr5id lrazzd5p" role="button" tabindex="0">顯示更多</div> == $0
```

```
for post in posts:  
    btms = post.find_elements_by_class_name('oajrlxb2 g5ia77u1 qu0x051f esr5mh6w e9989ue4  
    for btm in btms:  
        if btm.text == '顯示更多':  
            btm.click()  
            break
```

爬取粉絲專頁貼文

click() 會自動將瀏覽器滾動可以點擊的位置。

但是 facebook 上方的區塊會使點擊出錯，如果欲點擊的網頁元素在瀏覽器當前位置的上方，可能出現無法點擊的情況。



爬取粉絲專頁貼文

可以利用 javascript 指令來滾動瀏覽器到頂部，讓所有需要點擊的網頁元素都在當前位置的下方。

```
driver.execute_script('document.documentElement.scrollTop=0')  
time.sleep(1)
```

爬取粉絲專頁貼文

如果想要取得貼文的 ID ，需要從發布時間的連結取得。
但是連結裡面的 href 是空的，需要鼠標移動到該元素上面才會產生 href。

```
▼<span class="toJvnm2t a6sixzi8 abs2jz4q a8s20v7p tlp8iaqh k5wvi7nf q3lf05jv pk4s997a biomatt0 cebpdrjk qowsmv63 owhemhu dplhu0rb dhp6lc6y iyyx5f41">  
  ▶<a aria-label="4分鐘" class="osjflxb2 g5ia77u1 qu0x051f esr5mh6w e9989ue4 r7d6kgcz rq0escv nhd2j8a9 nc684n16 p7hj1n8o kvgm6g5 cxmmr5t8 oygrvhab hc  
    esuyzwur flsip0of lzcic4w1 gml0mx0 goro0wi8 blv8xokw" href="" role="link" tabindex="0">...</a>
```

爬取粉絲專頁貼文

首先找到需要連結的元素，然後利用 ActionChains 模擬鼠標移動到該元素上。

```
link = post.find_element_by_class_name('oajrlxb2 g5ia77u1')  
ActionChains(driver).move_to_element(link).perform()
```

爬取粉絲專頁貼文

回去查看網頁 html 碼，發現 href 的欄位出現了貼文連結。

```
<span class="toJvnm2t a6s1xzi8 abs2jz4q a8s20v7p tip8laqh K5wvi7nf q3lf05jv pk4s997a biomatto cebpdrjk qowdmv63 omwemhu dpihu0rb dhp6ic6y iyyx5f41">  
  <a aria-label="11分鐘" class="oaj1xb2 g5la77ul qu0x051f esr5mh6w e9989ue4 f70skgcr r00escxv nhd2j8e9 mc684n16 p7hj1n8o kvgn6g5 cxxmr5t8 oaygrvhab hcukyx3k j63vyjys  
    esuyzuar #isip00f l2c1c4ul gna10nx0 gpro0hi8 blv8xokw" href="https://www.facebook.com/appledaily.tw/posts/10161171306212069?__cft__vFnKgSK16zYLBuzYhJXo5o13VvGszVOC  
    #5jw8gtAowJL88Yh8bY8Ynth8__tn__=N2CON2CP>R" role="link" tabindex="0"></a>  
</span>
```

爬取粉絲專頁貼文

利用 BeautifulSoup() 來解析網頁資料。
找到目前所有貼文。

```
soup = BeautifulSoup(driver.page_source, 'lxml')  
  
posts = soup.find_all('div', {'class' : 'du4w351b k4urcfbm l9j0dhe7 sjgh65i0'})
```

爬取粉絲專頁貼文

由於每個貼文不一定都具備內容、按讚人數、留言人數、分享人數，
所以利用 try except 來確保不會有 bug。

```
try:
    message = post.find('div', {'class' : 'ecm0bbzt hv4rvrfc ihqw7lf3 datilw0a'})
    message = re.sub('\n', '', message.text)
except:
    message = ''

try:
    like_cnt = post.find('span', {'class' : 'pcp9lwgn'})
    like_cnt = re.findall('[0-9*]', like_cnt.text)[0]
except:
    like_cnt = 0

try:
    comment_cnt = post.find('span', {'class' : 'd2edcug0 hpfvrmgz qv66sw1b c1et5uq'})
    comment_cnt = re.findall('(.*?)則留言', comment_cnt.text)[0]
except:
    comment_cnt = 0

try:
    share_cnt = post.find('span', {'class' : 'd2edcug0 hpfvrmgz qv66sw1b c1et5uq'})
    share_cnt = re.findall('(.*?)次分享', share_cnt.text)[0]
except:
    share_cnt = 0
```


爬取粉絲專頁貼文

成功。

```
{
  "post": [
    {
      "post_id": "10161172324092069",
      "message": "怎麼可以隨意進入別人的住宅呢 (´皿`)>新聞追追追 逃出內幕追真相【◆◆獨家】發現「羊身吊屍」！網紅兩隻「女孩死屋內」空宅 這下慘了! https://bit.ly/3q8ueY3 #真花",
      "like": "3",
      "comment": "2",
      "share": 0
    },
    {
      "post_id": "10161171306212069",
      "message": "香李真花李哥又到了...看《蘋果》·解鎖解憂解無聊【台南香李#大賞花蜜點點點 走過木棉花·李調甲·紅花風鈴木棉花世界】https://bit.ly/3q8ueY3 #真花",
      "like": "77",
      "comment": 0,
      "share": 0
    },
    {
      "post_id": "10161172312067069",
      "message": "",
      "like": "4",
      "comment": "30",
      "share": 0
    },
    {
      "post_id": "10161171318732069",
      "message": "光看島新聞，就覺得這部的交通好複雜..... App在手，News我有 快下載看新聞更快捷【台中遠路口1個月竟出6663顆罰單 魔王網路網速網413萬】https://bit.ly/3q8ueY3 #真花",
      "like": "2",
      "comment": "4",
      "share": 0
    }
  ]
}
```

selenium requests 比較

- selenium
 - 模擬瀏覽器操作，較為方便理解。
 - 會得到所有網頁資料，速度較慢。
- requests
 - 送出請求要求資料，參數較多，不易理解。
 - 可以只要求目標資料，速度較快。

END