# 安裝 Selenium 和 WebDriver 來操作瀏覽器

# 網路爬蟲要應付的三種類型網站

- API ( 應用程式介面， Application Programming Interface)
  - Public：中央氣象局
    - https://opendata.cwb.gov.tw/fileapi/v1/opendataapi/{}?Authorization={}&format={}
  - Private：喜馬拉雅音頻分享平台
- 靜態網頁
  - 一般的小說網站：例如「愛下電子書」 https://tw.aixdzs.com/
- 動態網頁 ( 一般會安裝 python 的 Selenium 模組和 Chrome WebDriver)
  - 一般會牽扯到 java script 程式語言：例如 windy.com

# 安裝和更新 Selenium

# 下載 WebDriver (Google Chrome)

# 操控威騰 104 網頁的程式碼

挑戰：
註解的地方有 3 行是
Mozilla Firefox 瀏覽器的。
試著還原並做些修改，
再搭配 firefox 版本的
WebDriver 來操作。

程式開啟網頁後會模擬鍵盤
的 PgDn(PAGE_DOWN)，
操作 3 次，間隔 1 秒，之後
PgUp 也是 3 次，間隔 0.5
秒。最後，程式會尋找網頁
的「顯示全部」選項，找到
後模擬滑鼠左鍵的執行
(click)

```
wt_104_selenium_chrome_0.py - D:\Slides_2020\wt_demo\2020_12\code\wt_1...
File  Edit  Format  Run  Options  Window  Help
1  from selenium import webdriver
2  from selenium.webdriver.common.keys import Keys
3  #from selenium.webdriver.firefox.firefox_binary import FirefoxBinary #firefox
4  from selenium.webdriver.chrome.options import Options
5  import time
6
7  URL_HOME = "https://www.104.com.tw/company/1a2x6bimet"
8
9  # for firefox (mozilla)
10 #binary = FirefoxBinary(r'C:\InstalledApps\Mozilla Firefox\firefox.exe')
11 #browser = webdriver.Firefox(firefox_binary=binary)
12
13 # for chrome (google)
14 options = Options()
15 WEB_DRIVER_PATH = r'd:\PortableApps\chromedriver_win32\chromedriver.exe'
16 browser = webdriver.Chrome(WEB_DRIVER_PATH, options=options)
17 browser.get(URL_HOME)
18
19 #for i in range(3):
20 #    browser.find_element_by_tag_name('body').send_keys(Keys.PAGE_DOWN)
21 #    time.sleep(1)
22 browser.find_element_by_tag_name('body').send_keys(Keys.PAGE_DOWN)
23 time.sleep(1)
24 browser.find_element_by_tag_name('body').send_keys(Keys.PAGE_DOWN)
25 time.sleep(1)
26 browser.find_element_by_tag_name('body').send_keys(Keys.PAGE_DOWN)
27 time.sleep(1)
28
29 for i in range(3):
30     browser.find_element_by_tag_name('body').send_keys(Keys.PAGE_UP)
31     time.sleep(0.5)
32
33 MORE_BUTTON = 'div[class="dialog__show-all btn btn-sm btn-text pt-3"]'
34
35 while True:
36     show_more = browser.find_elements_by_css_selector(MORE_BUTTON)
37     if not show_more:
38         break
39     show_more[0].click()
40
```

載入 4 個「模組」(?)，其中
# 開頭那行代表 python 不會執行

威騰 104 網址

測試 firefox 瀏覽器用的，暫時被「註解」掉

產生一個 Options 物件 ( 剛好目前這個例子沒更改到預設值，不過未來會設定其他值 )

chrome webdriver 在我電腦內的路徑，大家的位置應該都不一樣

產生一個 google chrome 瀏覽器物件

使用剛剛產生的瀏覽器物件內的 get 功能去開啟威騰的 104 網頁

這 3 行被註解掉，用迴圈的方式做 PAGE_DOWN 三次

這 6 行和上面被拿掉的 3 行做相同的事，
每向下翻頁一次，就等 1 秒，如此重複 3 次。
這是比較「笨」的寫法。

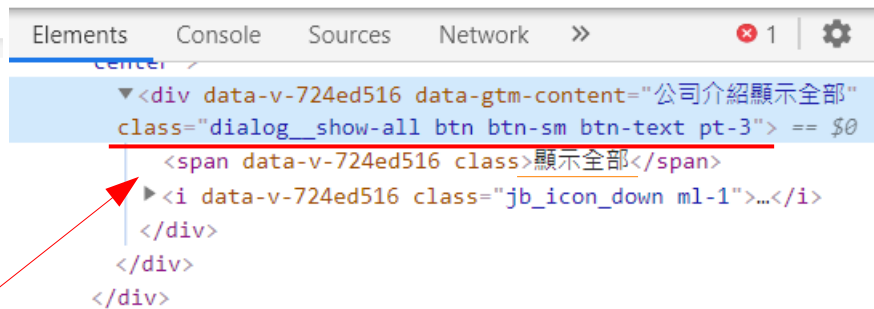向上翻頁 3 次，每次等 0.5 秒 ( 用到縮排 )

用來鎖定網頁「顯示全部」選項的資訊 ( 下一頁詳解 )

無限迴圈

使用 find_elements_by_css_selector 功能來尋找「顯示全部」

當找不到殘留的「顯示全部」選項的話，才會跳出永無止盡的迴圈

找到的話就模擬點擊滑鼠左鍵，這裡出現的 [0] 和我前面用 find_elements 有關，改成 find_element 單數才合理，有興趣人把它改掉吧～

# 進入瀏覽器的開發者模式



將滑鼠指標移到「顯示全部」上方，之後按右鍵