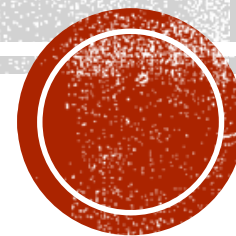


網路爬蟲實戰心法

國立臺灣大學共同教育中心

蔡芸琇



大綱

- 爬蟲基本概念介紹
- URL 的格式
- 網頁基本架構介紹
- 檢視原始碼
- 開發人員工具
- POSTMAN
- 使用 R 語言 **GET** Function
- 使用 R 語言 **POST** Function



爬蟲基本概念介紹

- 網頁的內容是由 HTML+CSS+JavaScript 組合而成。
- HTML+CSS+JavaScript 透過瀏覽器編譯後呈現給使用者。
- 爬蟲是一種自動抓取網頁內容 (HTML+CSS+JavaScript) 的程式。
- 爬蟲透過網址 (URL：Uniform Resource Locator) 提取網頁內容。

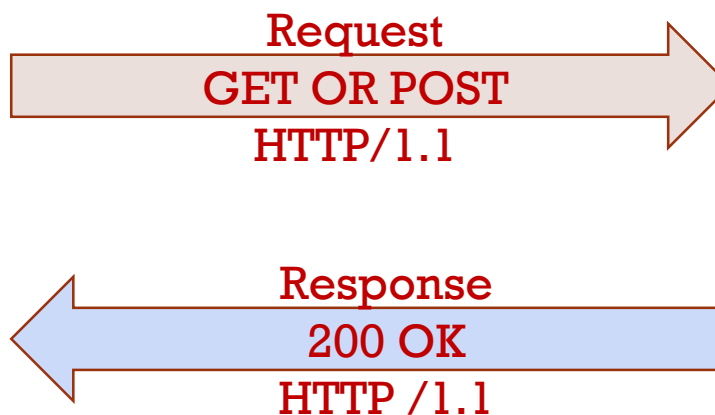


URL 的格式

- URL 分成三部分：<http://www.ntu.edu.tw/about/about.html>
 - 第一部分是協議 (http、https)。
 - 第二部分是存有該資源的主機位址 (140.112.8.116、www.ntu.edu.tw)。
 - 第三部分是主機上的具體目錄或文件名等 (about/about.html)。
 - 第一部分和第二部分用「://」符號隔開。
 - 第二部分和第三部分用「/」符號隔開。
 - 第一部分和第二部分不可缺少，第三部分看實際狀況決定是否加上。



網頁基本架構介紹



<http://httpbin.org/>
<https://www.w3.org/Protocols/rfc2616/rfc2616.html>



檢視原始碼

網路爬蟲並不默許爬蟲工作。因此在存取大量頁面時，爬蟲需要考慮到規劃、負載，還需要講「禮貌」。不願意被爬蟲存取、被爬蟲主人知曉的公開站點可以使用robots.txt檔案之類的方法避免存取。這個檔案可以要求機器人只對網站的一部分進行索引，或完全不作處理。

網路網路上的頁面極多，即使是最大的爬蟲系統也無法做出完整的索引。因此在公元2000年之前的全球資訊網出現初期，搜尋引擎經常找不到多少相關結果。現在的搜尋引擎在這方面已經進步很多，能夠即刻給出高品質結果。

爬蟲還可以驗證超連結和HTML代碼，用於網路抓取（參見資料驅動編程）。

目錄

- 命名
- 概述
- 設計者所面臨的挑戰
- 爬蟲策略
 - 選擇策略
 - 連結限制
 - URL 標識化
 - 路徑遞增爬取
 - 主題爬取
 - 重新存取策略
 - 平衡權限策略
 - 並列策略
- 另見
- 參考文獻

命名

網路爬蟲也可稱作網路蜘蛛^[1]、蜘蛛、自動索引程式（automatic indexer）^[2]，或（在FOAF軟體中）稱為網路爬虫（web scutter）。^[3]

概述

網路爬蟲始於一張被稱作種子的統一資源位址（URLs）列表。當網路爬蟲存取這些統一資源定位器時，它們會甄別出頁面上所有的超連結，並將它們寫入一張「待訪列表」，即所謂「爬行疆域」（crawl frontier）。此疆域上的統一資源位址將被按照一套策略遞迴存取。如果爬蟲在它執行的過程中複製權限和儲存網站上的資訊，這些檔案通常儲存，使他們可以被檢視、閱讀和瀏覽他們的網站上實時更新的資訊，並儲存為網站的「快照」。大容量的體積意味著網路爬蟲只能在給定時間內下載有限數量的網頁，所以要優先考慮其下載。高變化率意味著網頁可能已經被更新或者刪除。一些被何跟蹤端軟體生成的URLs（統一資源定位符）也使得網路爬蟲很難避免檢索到重複內容。

設計者所面臨的挑戰

網路網路資源極其浩繁，這意味著網路爬蟲在一定時間內只能下載有限數量的網頁，因此



開發人員工具

The screenshot displays the PChome 24h購物 website's product page for the Acer Predator Z271T 27-inch VA monitor. The page is in Chinese and features a prominent banner for the 'EASY SHOP Audrey' sale, which is running from 3/14 to 3/21. The banner includes a large '開館慶' (Grand Opening) text and a '卡碟商品 指定 任2入94折' (Card/Disc items, specified 2 for 94% off) promotion. The main product, the Acer Predator Z271T, is shown in a black and red design, with a price of 24880. The right sidebar displays the 'Predator Z35' monitor, which is a 35-inch VA curved monitor, with a price of 34900. The website includes various navigation links, such as '線上購物', '24h購物', '購物中心', and '商店'. The bottom of the page shows the price of 34900. The website is in Chinese and includes various promotional banners and navigation links.

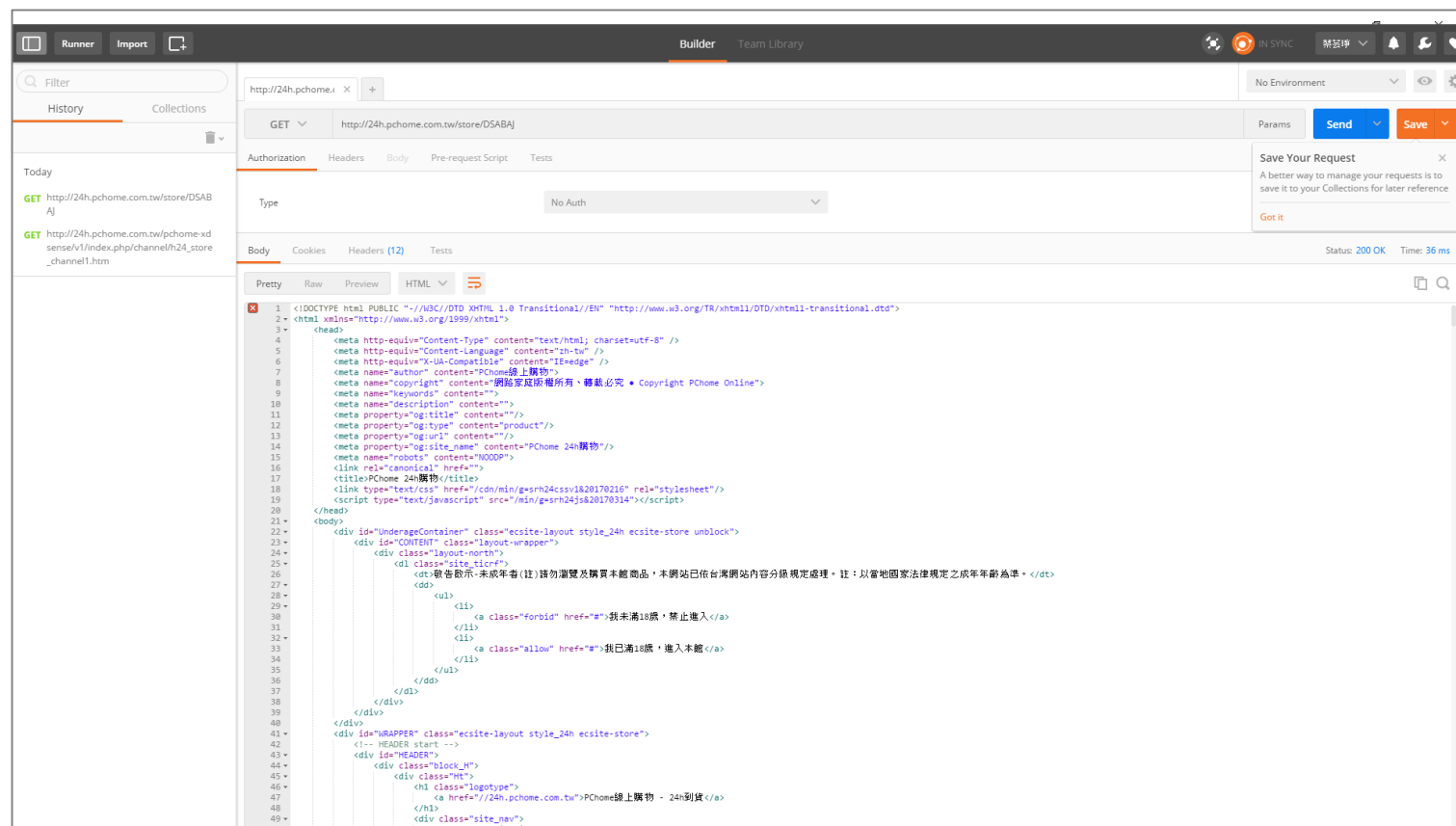


開發人員工具

- 在網頁任何位置按右鍵，選擇“檢查”，就可以看到相對應的原始碼。
- **Network** 頁面，可以看到網頁各項的執行細節。
- **Console** 頁面，可以檢查錯誤訊息。
- 直接對 **CSS** 樣式表更改參數，畫面就會直接可以預覽。
- 爬蟲觀察技巧：<http://tech-marsw.logdown.com/blog/2016/01/10/crawler-tips-mining-chrome>。



POSTMAN



R 語言 GET FUNCTION

```
GET("http://google.com/")
GET("http://google.com/", path = "search")
GET("http://google.com/", path = "search", query = list(q = "ham"))

# See what GET is doing with httpbin.org
url <- "http://httpbin.org/get"
GET(url)
GET(url, add_headers(a = 1, b = 2))
GET(url, set_cookies(a = 1, b = 2))
GET(url, add_headers(a = 1, b = 2), set_cookies(a = 1, b = 2))
GET(url, authenticate("username", "password"))
GET(url, verbose())

# You might want to manually specify the handle so you can have multiple
# independent logins to the same website.
google <- handle("http://google.com")
GET(handle = google, path = "/")
GET(handle = google, path = "search")
```

R Package httr : <https://cran.r-project.org/web/packages/httr/index.html>

User Guide : <https://cran.r-project.org/web/packages/httr/httr.pdf>



R 語言 **POST** FUNCTION

```
b2 <- "http://httpbin.org/post"
POST(b2, body = "A simple text string")
POST(b2, body = list(x = "A simple text string"))
POST(b2, body = list(y = upload_file(system.file("CITATION"))))
POST(b2, body = list(x = "A simple text string"), encode = "json")

# Various types of empty body:
POST(b2, body = NULL, verbose())
POST(b2, body = FALSE, verbose())
POST(b2, body = "", verbose())
```

R Package httr : <https://cran.r-project.org/web/packages/httr/index.html>

User Guide : <https://cran.r-project.org/web/packages/httr/httr.pdf>



牛刀小試

- 利用 GET 或 POST 存下 10 頁以上的網頁原始碼
- http://www.cookbook-r.com/Data_input_and_output/Writing_data_to_a_file/
- 提示 FUNCTIONS :
 - FOR
 - DUMP
 - GET
 - POST
- 回家練習參考：<https://github.com/pecu/RCrawler101-201504/tree/master/CaseStudies>

