

week 3

agenda

RegExp 介紹

Facebook API 介紹

抓取 FB 蘋果即時新聞 內容

資料庫存取

爬蟲自動化

正規表達式

為什麼寫爬蟲會遇到 RegExp ?



瑞銀堅持以客為優先
做客戶的可靠後援

UBS 瑞銀集團

黃色：施力方向
紅色：痠痛處
綠色：坐骨神經圖示走向

《元氣網粉絲團》健康資訊不漏接

「壁咚」別為浪漫傷肩、肘、頸

·穿緊身牛仔褲久蹲 肌肉、神經受損 ·6大錯誤習慣 害你背超痛

< 線上影音直播 哨人堂 快艇夏季補強 美女野獸情侶檔 壁咚過度傷筋骨 >

2015/7/25 星期六

新北市 32 °C ~ 36 °C



即時



分享



熱門



影音

- 14:38 5獨木舟翻覆11人落海 1女斷指
- 14:38 藍挺吳思華提告 緣駁抹黑學運老梗
- 14:29 教育部告不告記者？一周內有答案
- 14:24 50嵐蜂蜜出包 業者：落實送驗SOP
- 14:19 網友王子變青蛙 陸女吞玻璃輕生
- 14:06 50嵐蜂蜜出包 3飲品暫停售
- 14:01 神秘薄霧環繞 冥王星再給地球人驚喜
- 13:58 蜂蜜出包 50嵐：加強原料把關
- 13:58 基隆獨木舟翻船 四人平安獲救

兆豐國際商銀

網路結匯選兆豐 出國旅遊好輕鬆 甜蜜好禮送給您

商品圖片僅供參考，實際商品請以店內為主

良友旅行

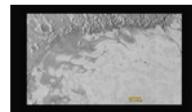
夏季旅遊特輯 www.ftstour.com.tw

嬉遊一夏 High 翻暑假

搜尋結果

搜尋 **冥王星** 的結果：共找到 58 筆 約 0.009秒**神秘薄霧環繞 冥王星再給地球人驚喜**

即時：2015/07/25



科學家今天表示，美國國家航空暨太空總署「新視野號」太空船上週飛越冰封



關注 udn



udn facebook

Elements Network Sources Timeline Profiles Resources Audits Console

▲ 2 >_

●

○

▼

View: Options: Preserve log Disable cache

Filter

Name

Path

x Headers Preview Response Cookies Timing

```

17 var pageList="";
18 var arr = new Array(rc);
19 var arrList = '';
20 var searchInfo='';
21 //searchInfo='搜尋 <u>' + q + '</u> 的結果：共找到 ' + numFound + '筆 約 ' + runTime ;
22 //searchInfo='搜尋 <u>' + q + '</u> 的結果：共找到 ' + numFound + '筆';
23
24 arr[0] = new Array(7);
25 arr[0][0] = 'http://udn.com/news/story/5/1078576';
26 arr[0][1] = 'http://uc.udn.com.tw/photo/2015/07/25/98/1144486.jpg';
27 arr[0][2] = '神秘薄霧環繞 <u>冥王星</u>再給地球人驚喜';
28 arr[0][3] = '即時';
29 arr[0][4] = '2015/07/25';
30 arr[0][5] = '科學家今天表示，美國國家航空暨太空總署「新視野號」太空船上週飛越冰封星球<u>冥王星</u>時，拍到的清晰輪廓圖可看到一大片大氣薄霧，地表特寫照顯示有氮冰流動。美國國家航空暨';
31 arr[0][6] = '';
32
33 arr[1] = new Array(7);
34 arr[1][0] = 'http://udn.com/news/story/5/1078223';
35 arr[1][1] = 'http://uc.udn.com.tw/photo/2015/07/25/98/1143868.jpg';
36 arr[1][2] = '新視野號揭祕 <u>冥王星</u>可能有氮冰河';
37 arr[1][3] = '即時';
38 arr[1][4] = '2015/07/25';
39 arr[1][5] = '「新視野號」(NEW HORIZONS) 太空船回傳的最新影像顯示，<u>冥王星</u>上可能有由氮冰組成的冰河。英國廣播公司(BBC)網站報導，科學家相信他們看到跡象顯示，冥';
40 arr[1][6] = '';
41
42 arr[2] = new Array(7);

```

Data 直接儲存在 javascript code 回傳
需要自己用正規表達截取出來

2015-07-24 台泥(1101) - 上市水泥工業

Stock Dog 是一個提供股票資訊的網站

總股數：3,692,176張

董監事及大股東持股比：29.07%

在外流通張數：2,618,852張

2015-07-24 台泥(1101) 走勢圖

Zoom all

From Jul 24, 2015 To Jul 24, 2015

昨收:37.35

均價:36.92

開盤	37.3	漲幅	-1.47%
最高	37.35	振福	1.47%
最低	36.8	漲跌	-0.55
收盤	36.8	昨收	37.35
總量	6719		

委賣價	委買量	委賣價	委賣量
36.80	354	36.85	4
36.75	122	36.90	183
36.70	202	36.95	103
36.65	337	37.00	165
36.60	852	37.05	199
委買賣差			1213
委買賣比			2.85

<https://www.stockdog.com.tw/stockdog/index.php?m=overview&sid=1101%20%E5%8F%B0%E6%B3%A5>



查詢總覽

股狗部落格

建議與留言

社群

付費方案

功能特色

行動版

登入

註冊

在外流通張數：2,618,852張

2015-07-24 台泥(1101) 走勢圖

Zoom

all

From Jul 24, 2015 To Jul 24, 2015

昨收:37.35

開盤:36.92

開盤	37.3	漲幅	-1.47%
最高	37.35	漲幅	1.47%
最低	36.8	漲跌	-0.55
收盤	36.8	昨收	37.35
總量	6719		
委買價	36.80	委賣價	36.85
	36.75		122
	36.70		202
	36.65		337
	36.60		852
委賣價差			1213
委賣賣比			2.85

Elements Network Sources Timeline Profiles Resources Audits Console

View: Options: Preserve log Disable cache

Filter All XHR Script Style Images Media Fonts Documents WebSockets Other Hide data URLs

Name Path

- index.php?m=overview&sid=1101 /stockdog
- guest_show.php?sid=1101>ype=17145197a2935cf104dd569a66e7ce4423 /stockdog ↑
- guest_show.php?click_btn3=1&click_btn6=1&click_btn7=1&sid=1101>ype... /stockdog

x Headers Preview Response Cookies Timing

```
<!DOCTYPE HTML>
<html>
<head>
<meta http-equiv=Content-Type content="text/html; charset=utf-8">
<title>2015-07-24 台泥(1101) 走勢圖</title>
<script type=text/javascript src=js/jquery/1.8.2/jquery.min.js></script>
<script type=text/javascript>$(function(){var data=[[Date.UTC(2015,6,24,09,00,00), 37.35, 0, 37.35,37.35], [Date.UTC(2015,6,24,09,00,03), 37.30, 133, null, null],[Date.UTC(2015,6,24,09,00,08), 37.30, 5, 37.30, 37.35], [Date.UTC(2015,6,24,09,00,13), 37.30, 1, 37.30, 37.35],[Date.UTC(2015,6,24,09,00,29), 37.35, 1, 37.30, 37.35], [Date.UTC(2015,6,24,09,00,54), 37.30, 25, 37.30, 37.35],[Date.UTC(2015,6,24,09,00,59), 37.30, 10, 37.30, 37.35], [Date.UTC(2015,6,24,09,01,04), 37.30, 6, 37.25, 37.30],[Date.UTC(2015,6,24,09,01,10), 37.30, 1, 37.25, 37.30], [Date.UTC(2015,6,24,09,01,15), 37.30, 1, 37.25, 37.30],[Date.UTC(2015,6,24,09,01,30), 37.20, 7, 37.25, 37.30],[Date.UTC(2015,6,24,09,01,35), 37.25, 9, 37.20, 37.25], [Date.UTC(2015,6,24,09,01,40), 37.25, 2, 37.20, 37.25],[Date.UTC(2015,6,24,09,01,45), 37.20, 4, 37.20, 37.25], [Date.UTC(2015,6,24,09,01,56), 37.20, 1, 37.20, 37.25],[Date.UTC(2015,6,24,09,02,01), 37.20, 2, 37.20, 37.25], [Date.UTC(2015,6,24,09,02,11), 37.25, 5, 37.20, 37.25],[Date.UTC(2015,6,24,09,02,16), 37.20, 6, 37.20, 37.25], [Date.UTC(2015,6,24,09,02,21), 37.25, 2, 37.20, 37.25],[Date.UTC(2015,6,24,09,02,37), 37.20, 5, 37.20, 37.25], [Date.UTC(2015,6,24,09,02,42), 37.25, 12, 37.20, 37.25],[Date.UTC(2015,6,24,09,02,47), 37.20, 12, 37.20, 37.25], [Date.UTC(2015,6,24,09,02,52), 37.15, 33, 37.15, 37.20],[Date.UTC(2015,6,24,09,02,57), 37.10, 30, 37.10, 37.15], [Date.UTC(2015,6,24,09,03,02), 37.00, 57, 37.00, 37.10],[Date.UTC(2015,6,24,09,03,08), 36.80, 24, 36.95, 37.00], [Date.UTC(2015,6,24,09,03,13), 36.80, 4, 36.80, 37.00], [Date.UTC(2015,6,24,09,03,18), 37.00, 31, 36.70, 37.00], [Date.UTC(2015,6,24,09,03,23), 37.00, 5, 36.90, 37.00], [Date.UTC(2015,6,24,09,03,28), 37.00, 3, 36.90, 37.00], [Date.UTC(2015,6,24,09,03,33), 37.00, 1, 36.95, 37.00], [Date.UTC(2015,6,24,09,03,38), 37.00, 1, 36.95, 37.00], [Date.UTC(2015,6,24,09,03,49), 37.00, 2, 36.90, 37.00]]
```

3 / 76 requests | 50.7 KB / 51.9 KB transferred | Finish: 3.93 s | DOMContentLoaded: 9.17145197a2935cf104dd569a66e7ce4423 1 of 1 Cancel 7



查詢總覽



段狗部落格

建議與留言

社群

付費方案

功能特色

行動版

登入

註冊

在外流通張數：2,618,852張

2015-07-24 台泥(1101) 走勢圖

[Zoom all](#)

From Jul 24, 2015 To Jul 24, 2015

開盤	37.3	漲幅	-1.47%
最高	37.35	振福	1.47%
最低	36.8	漲跌	-0.55
收盤	36.8	昨收	37.35
總量	6719		
委賣價	委買量	委賣價	委賣量
36.80	354	36.85	4
36.75	122	36.90	183
36.70	202	36.95	103
36.65	337	37.00	165
36.60	852	37.05	199
委賣賣差			1213
委賣賣比			2.85

昨收:37.35

均價:36.9

Elements Network Sources Timeline Profiles Resources Audits Console

View: Options: Preserve log Disable cache

Filter
Name
Path

[index.php?m=overview&sid=1101](#)
[/stockdog](#)

1

```
x Headers Preview Response Cookies Timing
133 $.ajax({url:'guest_show.php',cache:false,dataType:'html',type:'GET',data:
134 {sid:v1,Gstype:'27145197a2935cf104dd569a66e7ce4427',type:1,overview:'1'},
135 error:function(xhr{}),success:function(response){$('#ajax_pledger').html(response);$('#ajax_pledger').fadeIn();}}});
136
137
138 var k_i=0;var bs_i=0;
139 function change10bs(){if(bs_i++%2==1){document.getElementById('10bs_id').innerHTML='<a class="btn btn-success" onclick="change10bs()"切換至大圖'+'</a>'}
140 function setgraph(v1)
141 {if(v1==0){document.getElementById('g0').innerHTML='<iframe src="guest_show.php?sid=1101&Gstype=17145197a2935cf104dd569a66e7ce4423" width="100%" height="100%">'+'</iframe>'}
142 else if(v1==9){document.getElementById('g9').innerHTML='<iframe src="guest_show.php?date=2015-07-24&sid=1101&Ym=201507&GStvne=27145197a2935cf104dd569a66e7ce4423" width="100%" height="100%">'+'</iframe>'}
143 else if(v1==10){document.getElementById('g10').innerHTML='<iframe src="guest_show.php?date=2015-07-24&sid=1101&Ym=201507&Gstype=27145197a2935cf104dd569a66e7ce4423" width="100%" height="100%">'+'</iframe>'}
144 <script>document.getElementById('main_container').className='row-fluid';</script>
145 </div>
146 <meta property="og:title" content="籌碼分析_股狗網"/>
147 <title>台泥1101_籌碼分析_股狗網</title>
148 <div class="navbar navbar-fixed-bottom"><div class="navbar-inner"><div class="container-fluid" align="left" style="color:#fff; font-size:12px; padding-left:10px; padding-right:10px;">
149 <br/>
150 <ul>目前使用：電腦版，若為行動裝置可前往 <a href="https://www.stockdog.com.tw/m/index.php?cpc=0">行動版</a> </ul>
151 <ul>版權所有，網站由“股狗網資訊股份有限公司”設計，引用網頁中圖表，請註明出處來自股狗網</ul><ul>法律顧問：安成法律事務所。</ul>
152 <ul>資料來源：台灣證券交易所 TWSE、中華民國證券櫃檯買賣中心 OTC、公開資訊觀測站</ul><ul>使用本網站提供的服務，代表您同意我們的<a href="index.php?cpc=0">服務規範</a> </ul>
153
154 </div>
155 <script type="text/javascript" src="//use.typekit.net/yxt7dpk.js"></script>
156 <script type="text/javascript">try{Typekit.load();}catch(e{})</script>
```

3 / 76 requests | 50.7 KB / 51.9 KB transferred | Finish: 17145197a2935cf104dd56c 1 of 1 ▲ ▼ Cancel

正規表達練習

agilearning / RCrawlers

Branch: master ▾ RCrawlers / DataParsersInR / +

update codes c3h3 authored on 23 May latest commit f2e95faf35

..

RegEx_example1.R	create RPTT package	2 months ago
RegEx_example2.R	update codes	2 months ago
RegEx_example3.R	update codes	2 months ago
RegEx_jsonp2json2dataframe.R	create RPTT package	2 months ago
cssDemo1.R	update codes	2 months ago
cssDemo2.R	update codes	2 months ago
xpathDemo1.R	update codes	2 months ago
xpathDemo2.R	update codes	2 months ago
xpathDemo3.R	update codes	2 months ago

<https://github.com/agilearning/RCrawlers/tree/master/DataParsersInR>

練習 5分鐘

Facebook API

觀察 API 怎麼用

如何透過 API 拿到資料

如何使用 Search API

 Developers

My Apps

Products

Docs

Tools & Support

News

 Search in docs



Graph API Explorer

Application:

Access Token: CAACEdEose0cBANT4udFpveYy6uSIN3Fl0gV84

Graph API **FQL Query**

GET ▾

→ /v2.4/me?fields=id,name

Tools

[Graph API Explorer](#)

[App Insights](#)

[Object Browser](#)

[Ads Manager](#)

[URL Debugger](#)

[Access Token Tool](#)

[JS SDK Console](#)

[App Ads Helper](#)

Support

[Bugs](#)

[Platform Status](#)

Community

[Developer Group](#)

[Marketing Developers Group](#)

[On StackOverflow](#)



[On YouTube](#)

JS ▾

API Version: [?]

v2.4 ▾

BQU7S:

Debug

Get Token ▾

Debug Enabled ▾

Submit

Learn more about the Graph API syntax



Graph API Explorer

Application: [?]

Graph API Explorer ▾

Locale: [?]

English (US) ▾

API Version: [?]

v2.4 ▾

Access Token: CAACEdEose0cBANT4udFpveYy6uSIN3Ft0gV84cl0V8qvoKq2g4suAQ3lFHumF5K0Mqx5KV9vhCEz2VDABQU7S:

Debug

Get Token ▾

Graph API FQL Query

GET ▾

→ /v2.4/me?fields=id,name

Debug Enabled

Learn more about th

Get Access Token

Get App Token

Learning By Hacking

Taiwan R User Group

NUK useR Meetup

Graph API Explorer Application: [?] Graph API Explorer Locale: [?] English (US) API Version: [?] v2.4

Access Token: CAACEdEose0cBANT4udFpveYy6uSIN3Ft0gV84cl0V8qvoKq2g4suAQ3lFHUmF5K0Mqx5KV9vhCEz2VDABQU7S Debug Get Token

Graph API FQL Query GET → /v2.4/me?f User Data Permissions Extended Permissions

Select Permissions

User Data Permissions

user_about_me user_actions.books user_actions.fitness
user_actions.music user_actions.news user_actions.video
user_birthday user_education_history user_events
user_friends user_games_activity user_hometown
user_likes user_location user_managed_groups
user_photos user_posts user_relationship_details
user_relationships user_religion_politics user_status
user_tagged_places user_videos user_website
user_work_history

Public profile included by default.

Get Access Token Clear Cancel

先全選就對了
2小時後會自動 expired



Graph API Explorer

Application: [?]

Graph API Explorer ▾

Locale: [?]

English (US) ▾

API Version: [?]

v2.4 ▾

Access Token:

CAACEdEose0cBANT4udFpveYy6uSIN3Ft0gV84cl0V8qvoKq2g4suAQ3lFHumF5K0Mqx5KV9vhCEz2VDABQU7S:

Debug

Get Token ▾

Graph API

FQL Query

GET ▾

→ /v2.4/me?fi

Select Permissions

User Data Permissions

Extended Permissions

 ads_management ads_read email manage_pages publish_actions publish_pages read_custom_friendlists read_insights read_page_mailboxes rsvp_event

Public profile included by default.

Get Access Token

Clear

Cancel

Graph API Explorer

Application: [?]

Graph API Explorer [?]

Locale: [?]

English (US) [?]

API Version: [?]

v2.4 [?]

Access Token: CAACEdEose0cBABHFrzTmnYNRKzKuT7pMqzQGmKwGlGmSOY3wZC2anVjvikhLIUFLhv3ZCcE65tsXQDtpZ

Debug

Get Token [?]

Graph API

Facebook

GET [?]

Token Expired

Refresh

You can also use ALT T or click "Get Token"

Debug Enabled [?]

Submit

Learn more about the Graph API syntax

Node: 352962731493606

 likes

+ Search for a field

```
{  
  "error": {  
    "message": "Error validating access token: Session has expired on Saturday, 25-Jul-15 02:00:00 PDT. The current time is Saturday, 25-Jul-15 02:04:39 PDT.",  
    "type": "OAuthException",  
    "code": 190,  
    "error_subcode": 463  
  }  
}
```



兩個小時就會 expired !!



Graph API Explorer

Application: [?]

Graph API Explorer ▾

Locale: [?]

English (US) ▾

API Version: [?]

v2.4 ▾

Access Token: CAACEdEose0cBABHFrzTmnYNRKzKuT7pMqzQGmKwGIJmSOY3wZC2anVjvikhLIUFLhv3ZCcE65tsXQDtpZB5j

Debug

Get Token ▾

Graph API FQL Query

GET ▾

→ /v2.4/me?fields=id,name

Debug Enabled ▾

▶ Submit

Learn more about the Graph API syntax

Node: me

- id
- name

+ Search for a field

{
 "id": "1803156494",
 "name": "陳嘉葳"
}

大家應該會來到這個畫面

GET ▾

→ /v2.4/me?fields=id,name

Debug Enabled ▾

▶ Submit

Learn more about the Graph API syntax

Node: me



id



name



+

fields

about

address

age_range

bio

birthday

can_send_to_mobile

context

cover

created_time

currency

devices

education

email

favorite_athletes

favorite_teams

first_name

gender

hometown

inspirational_people

install_type

installed

interested_in

is_eligible_promo

is_verified

{
 "id": "1803156494",
 "name": "陳嘉葳"
}

Response received in 293 ms

Save Session

Graph API Explorer - Facebook

Search in docs

work
connections
accounts
activities
adaccounts
albums
apprequests
books
checkins
events
family
feed
friendlists
friendrequests
friends
games
groups
home
inbox
interests
likes
links
locations
movies
music
mutualfriends
notes
notifications
outbox
payments
permissions
photos
picture
posts
questions
scores
statuses
subscribedto
subscribers
tagged
television
updates
videos

Application: [?] Graph API Explorer [?] Locale: [?] English (US) [?] API Version: [?] v2.4 [?]

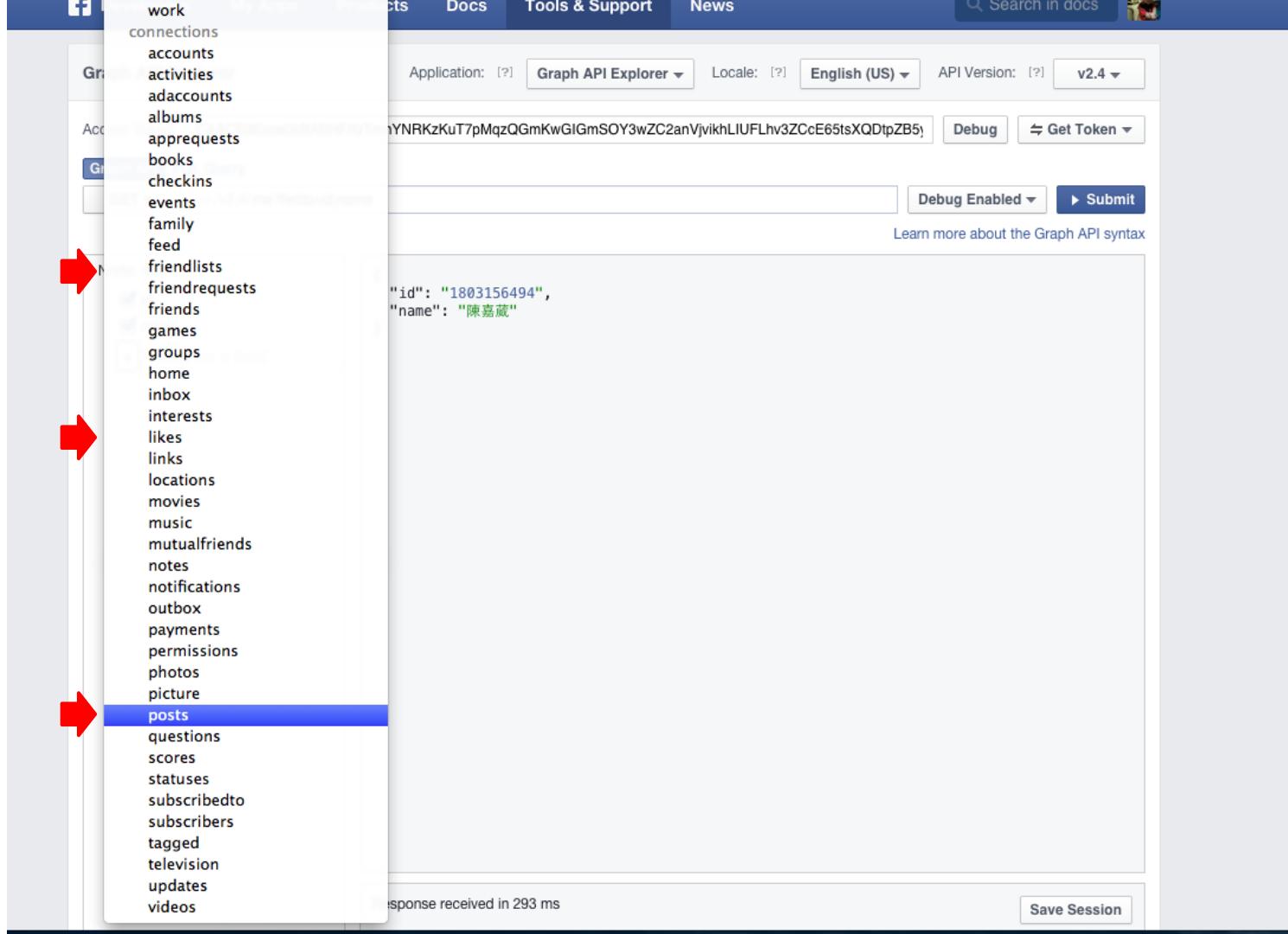
Debug Get Token

Debug Enabled Submit

Learn more about the Graph API syntax

```
{"id": "1803156494", "name": "陳嘉葳"}
```

Save Session



GET ▾

→ /v2.4/ me?fields=id,name,posts{message}

Debug Enabled  Submit

Learn more about the Graph API syntax

Node: me

- id
 - name
 - posts
 - message
- + Search for a field

+ Search for a field

```
{  
  "id": "1803156494",  
  "name": "陳嘉歲",  
  "posts": {  
    "data": [  
      {  
        "id": "1803156494_10203280881301813"  
      },  
      {  
        "id": "1803156494_10203270934573151"  
      },  
      {  
        "message": "AI 自己學習互相PK接球遊戲  
[code] https://github.com/ugo-nama-kun/DQN-chainer ○  
",  
        "id": "1803156494_10203269291052064"  
      },  
      {  
        "message": "今天小育來分享資安方面的議題，好多人一直發問撐超久 XD  
林育漢",  
        "id": "1803156494_10203266504142393"  
      },  
      {  
        "id": "1803156494_10203266151093567"  
      },  
      {  
        "message": "[CNN] 回來看錄影還是覺得很有意思 XD  
感謝 Liang Bo Wang  
",  
        "id": "1803156494_10203262121032818"  
      },  
      {  
        "message": "scikit-learn 官網從前天掛點到現在，我好怕明天上班怎辦 ...",  
        "id": "1803156494_10203261870826563"  
      }  
    ]  
  }  
}
```

Response received in 1407 ms

Save Session

練習 5分鐘

除了 post 還有 message 以外
找找看有沒有其它 你想要的 data :)

要如何用程式拿到
在Graph API Explorer 看到的 data ??

Access Token: CAACEdEose0cBABHFrzTmnYNRKzKuT7pMqzQGmKwGlGmSOY3wZC2anVjvikhLlIuFLhv3ZCcE65tsXQDtpZB5y

Debu

Get Token

Graph API FQL Query

GET ▾

→ /v2.4/me?fields=id,name,posts{message}

Debug Enabled

 **Submit**

[Learn more about the Graph API syntax](#)

Node: me

 id

name

posts

message

Search for a field

```
{  
  "id": "1803156494",  
  "name": "陳嘉葳",  
  "posts": {  
    "data": [  
      {  
        "id": "1803156494_10203280881301813"  
      },  
      {  
        "id": "1803156494_10203270934573151"  
      },  
      {  
        "message": "AT 自己 跟朋友互相PK接棒遊戲"
```

Elements Network Sources Timeline Profiles Resources Audits Consol

View: Options: Preserve log Disable cache

filter
ame
ath

me?access_token=CAACEdEose0cBABHFrzTmnYNRKzKuT...
graph.facebook.com/v2.4

bz
/ajəks/

Script Style Images Media Fonts Documents WebSockets Other Hide data URLs

x Headers Preview Response Timing

▼ General

Remote Address: 31.13.70.1:443
Request URL: https://graph.facebook.com/v2.4/me?access_token=CAACEdEose0cBABHFrzTmnYNRKzKuT7pMqzQGmKwGIGmSOY3wZC2anVjvikhLIUFLhv3ZCcE65QDtpZB5yAcp8MeA5bZCioKYv0Ph7gFcugSF9LQL3ufsvnAWshnw2VHigBr15V9ciZCBS3Ea9d2DvfmlLekDqnbbLixolh5fw21tykFlT0yySi0R4MdNWM1kIMt72xJ6K5ZAaEr4EcHDKaYZD&debug=all&fields=id%2Cname%2Cposts%7Bmessage%7D&format=json&method=get&pretty=0&suppress_http_code=1
Request Method: GET
Status Code: 200 OK

▼ Response Headers

access-control-allow-origin: *
cache-control: private, no-cache, no-store, must-revalidate
content-encoding: gzip
content-length: 1532
content-type: application/json; charset=UTF-8
date: Sat, 25 Jul 2015 07:56:02 GMT
etag: "4f9ee61bf3ac2cc72f84d9f2c66f8006260cefc2"
expires: Sat, 01 Jan 2000 00:00:00 GMT
facebook-api-version: v2.4
pragma: no-cache
status: 200 OK
vary: Accept-Encoding
via: HTTP/1.1

url 複製下來貼到瀏覽器

url 複製下來貼到瀏覽器

https://graph.facebook.com/v2.4/me?access_token=CAACEdEose0cBABHFrzTmnYNRKzKuT7pMqzQGmKwGIGmSOY3wZC2anVjvi

```
{"id": "1803156494", "name": "\u9673\u5609\u8473", "posts": {"data": [{"id": "1803156494_10203280881301813"}, {"id": "1803156494_10203270934573151"}, {"message": "AI \u81ea\u5df1\u5b78\u7fd2\u4e92\u76f8PK\u63a5\u7403\u904a\u6232[n]code] https://github.com/ugo-nama-kun/DQN-chainer\n", "id": "1803156494_10203269291052064"}, {"message": "\u4eca\u5929\u5c0f\u80b2\u4f86\u5206\u4eab\u8cc7\u5b89\u65b9\u9762\u7684\u8b70\u984c\uuff0c\u597d\u591a\u4eba\u4e00\u76f4\u767c\u554f\u6490\u8d85\u4e45 XD\n\u6797\u80b2\u6f22", "id": "1803156494_10203266504142393"}, {"id": "1803156494_10203266151093567"}, {"message": "[CNN] \u56de\u4f86\u770b\u9304\u5f71\u9084\u662f\u89ba\u5f97\u5f88\u6709\u8da3 XD\n\u611f\u8bld Liang Bo Wang", "id": "1803156494_10203262121032818"}, {"message": "scikit-learn \u5b98\u7db2\u5f9e\u524d\u5929\u639b\u9ede\u5230\u73fe\u5728\uuff0c\u6211\u597d\u6015\u660e\u5929\u4e0a\u73ed\u600e\u8fa6 ...", "id": "1803156494_10203261870826563"}, {"message": "Wilhelm Fumin \u6628\u5929\u624d\u525b\u804a\u5230\u4eca\u5929\u5c31\u770b\u5230\u6709\u5de5\u5177\u53ef\u4ee5\u7528 XD", "id": "1803156494_10203261531498080"}, {"id": "1803156494_10203249754323658"}, {"id": "1803156494_10203249441395835"}, {"message": "\u90a3\u500b\u767c\u73fe\u51a5\u738b\u661f\u7684\u5e74\u8f15\u4eba", "id": "1803156494_10203248470251557"}, {"id": "1803156494_10203244761478840"}, {"message": "\u4e09\u500b\u6642\u7684\u6700\u4f73\u5316\u6559\u5b78\uuff0c\u53ef\u4ee5\u4e45\u4e86~\n", "id": "1803156494_10203238922879"}, {"message": "SciPy 2015 \u7684\u5f71\u7247\u5df2\u7d93\u90fd\u653e\u51fa\u4f86\u4e86\uuff0c\u597d\u8fc5\u901f\u5594 \uff20\uuff20", "id": "1803156494_10203238907412492"}, {"message": "\u597d\u5f37\u5927\u5594\u00ff", "id": "1803156494_10203230737368246"}, {"message": "Excel power view \u597d\u5f37\u5594\u00ff01\n\u53ef\u4ee5\u76f4\u63a5\u9023\u7d50 bing map \u756b3d\u76f4\u65b9\u5716\u548c\u5bc6\u5ea6\u5716\n\u5fae\u8edf\u8d85\u5f37\u7121\u8aa4 \u738b\u4eae\u535a", "id": "1803156494_10203230644885934"}, {"id": "1803156494_1020322709087544"}, {"message": "http://icml.cc/2015/?page_id=97", "id": "1803156494_10203217897607260"}, {"id": "1803156494_10203211265041450"}, {"message": "\u4eca\u5e74 COSCUP \u6709\u4e00\u5834\u95dc\u65bc\u795e\u7d93\u5716\u9748\u6a5f \u7684 hands-on", "id": "1803156494_10203208560693843"}, {"id": "1803156494_10203208420650342"}, {"id": "1803156494_10203206612445138"}, {"paging": {"previous": "https://graph.facebook.com/v2.4/1803156494/posts?", "fields": "message&since=1437658261&access_token=CAACEdEose0cBABHFrzTmnYNRKzKuT7pMqzQGmKwGIGmSOY3wZC2anVjvikhLIUFLhv3ZCcE65tsXQDtpZB5yJAcP8MeAESbZCioKYv0Ph7gCfugSF9LIQLSL3ufsvnAWsHnw2VHi gBrI5V9ciZCBS3Ea9d2DvfLekDqnbBLixolH5fWz1tykFlT0yyso1R4MdNWm1kIMt72xJ6K5ZAAEMDKr4EcHDKaYZD&limit=25&__paging_token=enc_AdDiCQAR2TDDEcUwnFaOGGvGmgIIIdq4G3QPpvr17HiyJqAcYETxsBUG2pn1EZ FMxdlE2A82BQMFi3ULq33gF4ZCI&__previous=1", "next": "https://graph.facebook.com/v2.4/1803156494/posts?"}, {"fields": "message&access_token=CAACEdEose0cBABHFrzTmnYNRKzKuT7pMqzQGmKwGIGmSOY3wZC2anVjvikhLIUFLhv3ZCcE65tsXQDtpZB5yJAcP8MeAESbZCioKYv0Ph7gCfugSF9LIQLSL3ufsvnAWsHnw2VHi gBrI5V9ciZCBS3Ea9 d2DvfLekDqnbBLixolH5fWz1tykFlT0yyso1R4MdNWm1kIMt72xJ6K5ZAAEMDKr4EcHDKaYZD&limit=25&until=1436228395&__paging_token=enc_AdCXGllW9dZCaKm1TeF4vfMYkd0C6myf8K6t2Bt1ehGJC3jv47rf62qvDZCbe Y8v7uf1SnZALGuWoTBCPaPnbEVlrZB"}}}
```

```
In [1]: library(httr)
library(rjson)

In [2]: url = "https://graph.facebook.com/v2.4/me?access_token=CAACEdEose0cBABHFrfzTmnYNRKzKuT7pMqzQGmKwGIGmSOY3wZC2anVjvikhLIU"

In [3]: res = GET(url)

In [4]: content(res)

Out[4]: $id
'1803156494'
$name
'陳嘉歲'
$posts
$data
1. $id = '1803156494_10203280881301813'
2. $id = '1803156494_10203270934573151'
3. $message
'AI 自己學習互相PK接球遊戲 [code] https://github.com/ugo-nama-kun/DQN-chainer '
$id
'1803156494_10203269291052064'

4. $message
'今天小育來分享資安方面的議題，好多人一直發問撐超久 XD 林育漢'
$id
'1803156494_10203266504142393'

5. $id = '1803156494_10203266151093567'
6. $message
'[CNN] 回來看錄影還是覺得很有趣 XD 感謝 Liang Bo Wang '
$id
'1803156494_10203262121032818'
```

練習 5分鐘

試試看抓其他你想要的 data

如果我想抓 蘋果新聞
facebook 裡面的 data ...

https://www.facebook.com/apple.realtimenews?fref=ts

蘋果日報即時新聞

嘉威 Home 1

Create Page

Recent
2015
2014
2013

Sponsored

GQ Taiwan
S.S.22 夏日男POPUP Sh...
我們都認真工作，但絕不想因此而放棄食玩愛自由的Surfing精神。GQ和LEXUS聯手日本最強選物店BEAMS所打造的S.S.22夏日男POPUP Shop，將成為台灣首...
Thursday, July 23 at 1:00pm
Join Event · 1,260 people are going

668k people like this
Joseph Jiang and 6 other friends

Invite friends to like this Page

ABOUT

我十七歲的時候，雖然國文成績優異，不可能寫錯字，但是，我十七歲幹嘛呢？當時的我，有勇氣去質疑所讀的教科書嗎？」

「知識不是只有犀利與正確。知識也有它溫柔感性的一面。知識的溫柔，我們稱為智慧。」

只限於語言表面材料的細節是否正確，也必須通盤整體考量，這一段語文內容所指涉的意涵，從微觀到巨觀，學生經歷了什麼樣成長與補強的過程。

這包括：學生寫這一句話的動機是什麼？這一句話所傳達的意涵又是什麼？

要先找出蘋果即時新聞的 facebook ID

可以用 facebook search API 來找 :)



Marketing API

Messenger

Pages

Payments for Games

Sharing

Social Plugins

App Development

APIs and SDKs

Graph API**Using the Graph API**

Reference

Common Scenarios

Other APIs

Advanced

SDK for iOS

SDK for Android

SDK for JavaScript

SDK for PHP

SDK for Unity

Searching

You can search over many public objects in the social graph with the `/search` endpoint. The format for searches is:

```
GET graph.facebook.com  
/search?  
q={your-query}&  
[type={object-type}]({#searchtypes})
```



All Graph API search queries require an [access token](#) included in the request. The type of access token you need depends on the type of search you're executing.

- Searches across Page and Place objects requires an app access token.
- All other endpoints require a user access token.

Available Search Types

We support search for the following types:

Type	Description	`q` value
user	Search for a person (if they allow their name to be searched for).	Name.
page	Search for a Page.	Name.
event	Search for an event.	Name.
group	Search for a Group.	Name.
place	Search for a place. You can narrow your search to a specific location and distance by adding the <code>center</code> parameter (with <code>latitude</code> and <code>longitude</code>) and an optional <code>distance</code> parameter.	Name.

and distance by adding the `center` parameter (with `latitude` and `longitude`) and an optional `distance` parameter:

Basics
Reading
Publishing
Updating
Deleting
Searching
Available Search Types
Handling Errors
Debugging API Requests





Marketing API

Messenger

Pages

Payments for Games

Sharing

Social Plugins

App Development

APIs and SDKs

Graph API

Using the Graph API

Reference

Common Scenarios

Other APIs

Advanced

SDK for iOS

SDK for Android

SDK for JavaScript

SDK for PHP

SDK for Unity

We support search for the following types:

Type	Description	`q` value
user	Search for a person (if they allow their name to be searched for).	Name.
page	Search for a Page.	Name.
event	Search for an event.	Name.
group	Search for a Group.	Name.
place	Search for a place. You can narrow your search to a specific location and distance by adding the <code>center</code> parameter (with latitude and longitude) and an optional <code>distance</code> parameter:	Name.
placetopic	Returns a list of possible place Page topics and their IDs. Use with <code>topic_filter=all</code> parameter to get the full list.	None.
ad_*	A collection of different search options that can be used to find targeting options.	See Targeting Options docs

Example:

官方文件給的範例，
查詢在這個經緯度周圍一公里 內的咖啡廳

```
GET graph.facebook.com
/search?
q=coffee&
type=place&
center=37.76,-122.427&
distance=1000
```

English (US) ▾

Basics
Reading
Publishing
Updating
DeletingSearching
Available Search Types
Handling Errors
Debugging API Requests

Access Token: CAACEdEose0cBABHFrzTmnYNRKzKuT7pMqzQGmKwGiGmSOY3wZC2anVjvikhLIUFLhv3ZCcE65tsXQDtpZB5

Debug

Get Token ▾

Graph API FQL Query

GET ▾

→ /v2.4/ search?q=coffee&type=place¢er=37.76,-122.427&distance=1000

Debug Enabled ▾

Submit

Learn more about the Graph API syntax

Node: search

(No fields expansion available).

第一間

第二間

...

```
{  
  "data": [  
    {  
      "category": "Local business",  
      "category_list": [  
        {  
          "id": "162264673824073",  
          "name": "Specialty Grocery Store"  
        },  
        {  
          "id": "128673187201735",  
          "name": "Coffee Shop"  
        }  
      ],  
      "location": {  
        "street": "4023 18th St",  
        "city": "San Francisco",  
        "state": "CA",  
        "country": "United States",  
        "zip": "94114-2501",  
        "latitude": 37.7607613,  
        "longitude": -122.4332886  
      },  
      "name": "Philz Coffee - Castro",  
      "id": "151116474914629"  
    },  
    {  
      "category": "Local business",  
      "category_list": [  
        {  
          "id": "197871390225897",  
          "name": "Cafe"  
        }  
      ],  
      "location": {  
        "street": "18th & Valencia St",  
        "city": "San Francisco",  
        "state": "CA",  
        "country": "United States",  
        "zip": "94114-2501",  
        "latitude": 37.7749235,  
        "longitude": -122.4332886  
      },  
      "name": "Philz Coffee - Castro",  
      "id": "151116474914629"  
    }  
  ]  
}
```



上一頁的範例
在 Graph API Explorer 的用法

第一間咖啡廳的類別名稱
(屬於兩種類別)

第一間咖啡廳的所在國家，街道，經緯度等資訊

第一間咖啡廳的名稱 與 id

Response received in 988 ms

Save Session



Graph API Explorer

Application: [?]

Graph API Explorer [?]

Locale: [?]

English (US) [?]

API Version: [?]

v2.4 [?]

Access Token: CAACEdEose0cBABHFrzTmnYNRKzKuT7pMqzQGmKwGIGmSOY3wZC2anVjvikhLIUFLhv3ZCcE65tsXQDtpZB5|

Debug

Get Token [?]

Graph API FQL Query

GET [?]

→ /v2.4/search?q=coffee&type=place&center=37.76,-122.427&distance=1000

Debug Enabled [?]

Submit

Learn more about the Graph API syntax

Node: search

(No fields expansion available).

```
{
  "data": [
    {
      "category": "Local business",
      "category_list": [
        {
          "id": "162264673824073",
          "name": "Starbucks"
        }
      ]
    }
  ]
}
```

Elements Network Sources Timeline Profiles Resources Audits Console



View: Options: Preserve log Disable cache

Filter

Name Path

search?access_token=CAACEdEose0cBABHFrzTmnYNRKz...
graph.facebook.com/v2.4

All | XHR

Script Style Images Media Fonts Documents WebSockets Other □ Hide data URLs

Headers Preview Response Timing

General

Remote Address: 31.13.70.1:443

Request URL: https://graph.facebook.com/v2.4/search?access_token=CAACEdEose0cBABHFrzTmnYNRKzKuT7pMqzQGmKwGIGmSOY3wZC2anVjvikhLIUFLhv3ZCcE65tsXQDtpZB5|JAcP8MeAESbZCioKYv0Ph7gCfugSF9LIQL3ufsvnAwshnw2VHigBrI5V9ciZCBS3Ea9d20vfmLekDqnbLixolH5fw21tykFlT0yySi0R4MdNWM1kIMt72xJ6K5ZAaEMDKr4EcHdKYZD&center=37.76%2C-122.427&debug=all&distance=1000&format=json&method=get&pretty=0&q=coffee&suppress_http_code=1&type=place

Request Method: GET

Status Code: 200 OK

Response Headers

access-control-allow-origin: *
cache-control: private, no-cache, no-store, must-revalidate
content-encoding: gzip
content-length: 1686
content-type: application/json; charset=UTF-8
date: Sat, 25 Jul 2015 08:23:20 GMT
etag: "58a06d7936cb3457b37bee144d40750e37dfd51e"
expires: Sat, 01 Jan 2000 00:00:00 GMT
facebook-api-version: v2.4
pragma: no-cache
status: 200 OK
vary: Accept-Encoding
version: HTTP/1.1
x-fb-debug: JLE66e0vIwMDGw2mH2oqCohd2ukFCMY7NCTmXrr/VB0UqzGpwpG0f63YJEIK0UPHIjVKenEZKk7xD6hcbcRa4g==

複製到瀏覽器就可以拿到附近咖啡廳資料
也可以直接用 R 送出 GET 去拿資料



Graph API Explorer

Application: [?]

Graph API Explorer ▾

Locale: [?]

English (US) ▾

API Version: [?]

v2.4 ▾

Access Token: CAACEdEose0cBABHFrfzTmnYNRKzKuT7pMqzQGmKwGlGmSOY3wZC2anVjvikhLIUFLhv3ZCcE65tsXQDtpZB5j

Debug

Get Token ▾

Graph API FQL Query

GET ▾

→ /v2.4/search?q=蘋果日報&type=group

Debug Enabled ▾

Submit

type 要換成 group

Learn more about the Graph API syntax

Node: search

```
{  
  "data": [  
    {  
      "name": "蘋果日報 表格化",  
      "privacy": "CLOSED",  
      "id": "456650471114601"  
    },  
    {  
      "name": "蘋果日報壹傳媒的賣賣秘辛",  
      "privacy": "CLOSED",  
      "id": "349973165098837"  
    },  
    {  
      "name": "蘋果日報動新聞NMNews",  
      "privacy": "OPEN",  
      "id": "1468132446748773"  
    },  
    {  
      "name": "蘋果日報動新聞",  
      "privacy": "CLOSED",  
      "id": "101237726708440"  
    },  
    {  
      "name": "蘋果日報每日新聞俱樂部",  
      "privacy": "OPEN",  
      "id": "142228275838151"  
    },  
    {  
      "name": "反對蘋果日報大聯盟",  
      "privacy": "OPEN",  
      "id": "47745043598"  
    },  
    {  
      "name": "蘋果日報八卦週刊-",  
      "privacy": "OPEN",  
      "id": "142228275838151"  
    }  
  ]  
}
```

Response received in 438 ms

Save Session



Graph API Explorer

Application: [?]

Graph API Explorer ▾

Locale: [?]

English (US) ▾

API Version: [?]

v2.4 ▾

Access Token: CAACEdEose0cBABHFrzTmnYNRKzKuT7pMqzQGmKwGlGmSOY3wZC2anVjvikhLIUFLhv3ZCcE65tsXQDtpZB5

Debug

Get Token ▾

Graph API FQL Query

GET ▾

→ /v2.4/search?q=蘋果日報即時新聞&type=page

Debug Enabled ▾

Submit

type 要換成 page

Learn more about the Graph API syntax

Node: search

```
{  
  "data": [  
    {  
      "name": "蘋果日報即時新聞",  
      "id": "352962731493606"  
    },  
    {  
      "name": "蘋果日報即時新聞",  
      "id": "1497703573830172"  
    }  
  ],  
  "paging": {  
    "cursors": {  
      "before": "MAZDZD",  
      "after": "MQZDZD"  
    }  
  }  
}
```

id 是這一組

Graph API Explorer Application: [?] Graph API Explorer Locale: [?] English (US) API Version: [?] v2.4

Access Token: CAACEdEose0cBABHFrzTmnYNRKzKuT7pMqzQGmKwGlGmSOY3wZC2anVjvikhLlUFLhv3ZCcE65tsXQDtpZB5] Debug Get Token

Graph API FQL Query

GET → /v2.4/352962731493606?fields=posts{message}

Debug Enabled Submit

Learn more about the Graph API syntax

Node: 352962731493606

posts

message

+ Search for a field

+ Search for a field

得到蘋果即時新聞的 po 文

```
{ "posts": { "data": [ { "message": "「我十七歲的時候，雖然國文成績優異，不可能寫錯字，但是，我十七歲幹嘛呢？當時的我，有勇氣去質疑所讀的教科書嗎？」", "id": "352962731493606_729010403888835" }, { "message": "「知識不是只有犀利與正確。知識也有它溫柔感性的一面。知識的溫柔，我們稱為智慧。」", "id": "352962731493606_729011380555404" }, { "message": "「這位妹妹可以表示：我爸爸走丟了...」", "id": "352962731493606_729009097222299" }, { "message": "柯文哲表示，他過去在台大服務，跟前中正一分局長方仰寧熟識，清楚知道他因為太陽花一事，身揹46條案件，他認為這樣的狀況是不正常的，他的態度一致，就是不希望讓基層扛責任，未來會跟警政署講好規則，「每次都叫下面來扛，我不喜歡這樣」。", "id": "352962731493606_729008807222328" }, { "message": "#黃安 這次也不缺席", "id": "352962731493606_728970127226196" }, { "message": "有些影像淡淡的卻很深刻，我想就是這種影片吧。", "id": "352962731493606_728988660557676" } ] } }
```

Response received in 524 ms Save Session



Graph API Explorer

Application: [?]

Graph API Explorer

Locale: [?]

English (US)

API Version: [?]

v2.4

Access Token:

CAACEdEose0cBABHFrzTmnYNRKzKuT7pMqzQGmKwGIGmSOY3wZC2anVjvikhLIUFhv3ZCcE65tsXQDtpZB5|

Debug

Get Token

Graph API FQL Query

GET

→ /v2.4/352962731493606?fields=posts{message}

Debug Enabled

Submit

Learn more about the Graph API syntax

Node: 352962731493606

 posts message

+ Search for a field

```
{
  "posts": {
    "data": [
      {
        "message": "「我十七歲的時候，雖然國文成績優異，不可能寫錯字，但是，我十七歲幹嘛呢？當時的我，有勇氣去質疑所讀的教科書嗎？」"
      }
    ]
  }
}
```

Elements Network Sources Timeline Profiles Resources Audits Console

View: Options: Preserve log Disable cacheFilter All | XHR Script Style Images Media Fonts Documents WebSockets Other Hide data URLsName
Path

352962731493606?access_token=CAACEdEose0cBABHFr...

graph.facebook.com/v2.4

bz
/ajax

General

Remote Address: 31.13.70.1:443
 Request URL: https://graph.facebook.com/v2.4/352962731493606?access_token=CAACEdEose0cBABHFrzTmnYNRKzKuT7pMqzQGmKwGIGmSOY3wZC2anVjvikhLIUFhv3ZCcE65tsXQDtpZB5|Lhv3ZCcE65tsXQDtpZB5yJAcP8MeAESbZCioKYv0Ph7gCfugSF9LIQSL3ufsvnAWsHnw2VHigBrI5V9ciZCBS3Ea9d2DvfmLekDqnbLixolH5fWz1tykFlT0yyxi0R4MdNWm1kIMt72xJ6K5ZAaEMDKr4EcHDKaYZD&debug=all&fields=posts%7Bmessage%7D&format=json&method=get&pretty=0&suppress_http_code=1
 Request Method: GET
 Status Code: 200 OK

Response Headers

access-control-allow-origin: *
 cache-control: private, no-cache, no-store, must-revalidate
 content-encoding: gzip
 content-length: 2473
 content-type: application/json; charset=UTF-8
 date: Sat, 25 Jul 2015 08:29:32 GMT
 etag: "8bb98f5dd9a382a94d294ba6b50c14d43e24de79"
 expires: Sat, 01 Jan 2000 00:00:00 GMT
 facebook-api-version: v2.4
 pragma: no-cache
 status: 200 OK
 vary: Accept-Encoding
 version: HTTP/1.1
 x-fb-debug: ucQw5BURhAmnhbNzMVqsGFhp/BNZBtz9TuVPj1RDYAr3K1Sknj6wXFb+jCpv5VGZy00/Wih0XK3qcJlZMAYaa0==

複製到瀏覽器就可以拿到蘋果新聞 po 文資料


```
In [1]: library(httr)
library(rjson)
```

```
In [14]: url = "https://graph.facebook.com/v2.4/352962731493606?access_token=CAACEdEose0cBABHFrfrzTmnYNRKzKuT7pMqzQGmKwGIGmSOY3wZt
```

```
In [15]: res = content(GET(url))
```

```
In [16]: res
```

```
Out[16]: $posts
```

```
$data
```

1. \$message

'敘利亞、利比亞及厄利垂亞政局動盪，製造出大量試圖偷渡往歐洲的難民...但卻不是每個人都有能活下來的好運氣...'

\$id

'352962731493606_729009783888897'

2. \$message

'「我十七歲的時候，雖然國文成績優異，不可能寫錯字，但是，我十七歲幹嘛呢？當時的我，有勇氣去質疑所讀的教科書嗎？」 「知識不是只有犀利與正確。知識也有它溫柔感性的一面。知識的溫柔，我們稱為智慧。」 '

\$id

'352962731493606_729010403888835'

3. \$message

'台灣首例'

\$id

'352962731493606_729011380555404'

4. \$message

'這位妹妹可以表示：我爸爸走丢了...'

\$id

'352962731493606_729009097222299'

5. \$message

透過 R 程式拿到蘋果新聞 po文資料

練習 5分鐘

res\$... 底下還有 post, paging ...

請觀察一下這些是在做什麼 :)

把抓下來的 data 進行整理
變成表格的樣子 !

```
In [18]: do.call(rbind, res$posts$data)
```

Out[18]:

message	id
敘利亞、利比亞及厄利垂亞政局動盪，製造出大量試圖偷渡往歐洲的難民...但卻不是每個人都有能活下來的好運氣...	352962731493606_729009783888897
「我十七歲的時候，雖然國文成績優異，不可能寫錯字，但是，我十七歲幹嘛呢？當時的我，有勇氣去質疑所讀的教科書嗎？」 「知識不是只有犀利與正確。知識也有它溫柔感性的一面。知識的溫柔，我們稱為智慧。」	352962731493606_729010403888835
台灣首例	352962731493606_729011380555404
這位妹妹可以表示：我爸爸走丢了...	352962731493606_729009097222299
柯文哲表示，他過去在台大服務，跟中正一分局長方仰寧熟識，清楚知道他因為太陽花一事，身揹46條案件，他認為這樣的狀況是不正常的，他的態度一致，就是不希望讓基層扛責任，未來會跟警政署講好規則，「每次都叫下面來扛，我不喜歡這樣」。	352962731493606_729008807222328
#黃安 這次也不缺席	352962731493606_728970127226196
有些影像淡淡的卻很深刻，我想就是這種影片吧。解婕翎 #一人一半才是伴 #永和	352962731493606_728988660557676
讚 #台鐵	352962731493606_728967820559760
#中肯文 #事情不是你想的那麼簡單	352962731493606_728950990561443
好離譜的大車 #大車就是任性	352962731493606_728950230561519
噴噴噴，你知道的太多了 #得罪方丈還想走	352962731493606_728906400565902
對發票的快樂日子又來囉，祝大家中獎 #統一發票	352962731493606_728947397228469
受傷最重的是騎士，最無辜的也是騎士 #路邊開車門	352962731493606_728905737232635
騎士一去不回頭呀 #阿斯拉 #風見隼人 #今天的我沒有極限	352962731493606_728905600565982
佔領部長室未成年學生聲明 #最新	352962731493606_728908663899009
好珍貴的畫面 #讓爸媽看一下 #記得當時年紀小	352962731493606_728904960566046
蠻可愛的呀 #心意最重要 #不收可能更恐怖	352962731493606_728865293903346
踹共啦！	352962731493606_728865147236694

上一步沒新訊息後，應該跟上一步一樣的問題，數字加上後面接續的字母即為數字的序號。

常聽說 ...

要抓 facebook 的按 

要自己一直往下捲很麻煩 :(

(要做社群分析, 或觀點挖掘)



Graph API Explorer

Application: [?]

Graph API Explorer [?]

Locale: [?]

English (US) [?]

API Version: [?]

v2.4 [?]

Access Token: CAACEdEose0cBANqVESvUWtZAq3hH8txuE27U3YmOXD5wHwUkIxcGuhyzC5cNYt7QZABmx3qWB6ZA7uyOkpl

Debug

Get Token [?]

Graph API FQL Query

GET [?]

→ /v2.4/352962731493606/?fields=posts{likes}

Debug Enabled [?]

Submit

Learn more about the Graph API syntax

Edge: 352962731493606/

 posts likes

+ Search for a field

+ Search for a field

+ Search for a field

```
{
  "posts": {
    "data": [
      {
        "id": "352962731493606_729025537220655",
        "likes": {
          "data": [
            {
              "id": "535625486563366"
            },
            {
              "id": "733468196697926"
            },
            {
              "id": "854482524564749"
            },
            {
              "id": "1413973885556387"
            }
          ]
        }
      }
    ]
  }
}
```

Elements Network Sources Timeline Profiles Resources Audits Console

View: Preserve log Disable cache

Filter All | XHR Script Style Images Media Fonts Documents WebSockets Other Hide data URLs

Name Path

?access_token=CAACEdEose0cBANqVESvUWtZAq3hH8txuE27U3YmOXD5wHwUkIxcGuhyzC5cNYt7QZABmx3qWB6ZA7uyOkpl graph.facebook.com/v2.4/352962731493606

General

Remote Address: 31.13.70.1:443
 Request URL: https://graph.facebook.com/v2.4/352962731493606/?access_token=CAACEdEose0cBANqVESvUWtZAq3hH8txuE27U3YmOXD5wHwUkIxcGuhyzC5cNYt7QZABmx3qWB6ZA7uyOkpl&method=get&pretty=0&suppress_http_code=1
 Request Method: GET
 Status Code: 200 OK

Response Headers

access-control-allow-origin: *
 cache-control: private, no-cache, no-store, must-revalidate
 content-encoding: gzip

```
In [21]: url = "https://graph.facebook.com/v2.4/352962731493606/?access_token=CAACEdEose0cBANqVESvUWtZAq3hH8txuE27U3YmOXD5wHwUkI"

In [22]: res = content(GET(url))

In [23]: res$posts$data
```

Out[23]: 1. \$id
'352962731493606_729025537220655'
\$likes
\$data
A. \$id = '256515907884819'
B. \$id = '813815345305329'
C. \$id = '614920171928586'
D. \$id = '1488454388101539'
E. \$id = '649265348520420'
F. \$id = '667586220013890'
G. \$id = '737110892979155'
H. \$id = '253410108192365'
I. \$id = '1444147985909755'
J. \$id = '1519188568354590'
K. \$id = '1423990601202092'
L. \$id = '719911508103768'
M. \$id = '14277275812679579'

```
In [ ]:
```

```
In [ ]:
```

```
In [26]: res$posts$data[[1]]
```

```
Out[26]: $id  
'352962731493606_729025537220655'
```

```
$likes
```

```
$data
```

1. \$id = '256515907884819'
2. \$id = '813815345305329'
3. \$id = '614920171928586'
4. \$id = '1488454388101539'
5. \$id = '649265348520420'
6. \$id = '667586220013890'
7. \$id = '737110892979155'
8. \$id = '253410108192365'
9. \$id = '1444147985909755'
10. \$id = '1519188568354590'
11. \$id = '1423990601202092'
12. \$id = '719911508103768'
13. \$id = '1377375812578572'
14. \$id = '763108580380072'
15. \$id = '864220310261151'
16. \$id = '668151756588804'
17. \$id = '225784700956218'
18. \$id = '103922183276595'
19. \$id = '774755135879740'
20. \$id = '715041565258389'
21. \$id = '1008525422494853'
22. \$id = '850074498425143'
23. \$id = '1474845946130861'
24. \$id = '1587771148155712'
25. \$id = '1428255687495149'

```
$spaging
```

```
$scursors
```

```
In [27]: res$posts$data[[2]]
```

```
Out[27]: $id
```

```
'352962731493606_729009917222217'
```

```
$likes
```

```
$data
```

1. \$id = '642486195894961'
2. \$id = '806449262700338'
3. \$id = '865610850122726'
4. \$id = '238986379639869'
5. \$id = '1568492620042765'
6. \$id = '10202824355027728'
7. \$id = '710321349032346'
8. \$id = '861015497249010'
9. \$id = '1381075902193874'
10. \$id = '810318048989326'
11. \$id = '1554420924779466'
12. \$id = '110992655904806'
13. \$id = '918824734803259'
14. \$id = '234064736791495'
15. \$id = '809347762412102'
16. \$id = '532642473514928'
17. \$id = '386630278143094'
18. \$id = '577653822332246'
19. \$id = '695334383841746'
20. \$id = '835034753184482'
21. \$id = '843649072318968'
22. \$id = '1431661330428316'
23. \$id = '1519188568354590'
24. \$id = '853438874681914'
25. \$id = '253319124854701'

```
$spaging
```

```
$scursors
```

```
$after
```

```
'MjUzMzE5MTI0ODU0NzAx'
```

一次給 25筆，
需要自己不斷去 GET 下一頁

(蘋果新聞平均一篇文章超過 500 叢)



**不斷 for loop 就對了
太瑣碎就交給同學們自己嘗試**

順便出個 A Crawler A Day 題目

大家可以做一下 ”觀點分析“

抓好的 data 可以直接
存到資料庫嗎？

Graph API FQL Query

GET → /v2.4/352962731493606/?fields=posts Debug Enabled ▾ Submit

Learn more about the Graph API syntax

Edge: 352962731493606/
 posts
+ Search for a field
+ Search for a field

```
{
  "posts": [
    {
      "data": [
        {
          "message": "#中國 一次的英雄行為，卻為家人帶來巨變，但梁華並無後悔當年捉賊的行為，被問及如再重新選擇，還抓賊嗎？梁華肯定回應：「抓！」",
          "created_time": "2015-07-25T09:14:26+0000",
          "id": "352962731493606_729028317220377"
        },
        {
          "message": "十張美不勝收的照片，這張FB截圖可不是第一名喔",
          "created_time": "2015-07-25T09:00:44+0000",
          "id": "352962731493606_729025537220655"
        },
        {
          "message": "好久不見陳老師",
          "created_time": "2015-07-25T08:45:01+0000",
          "id": "352962731493606_72900991722217"
        },
        {
          "message": "敘利亞、利比亞及厄利垂亞政局動盪，製造出大量試圖偷渡往歐洲的難民...但卻不是每個人都有能活下來的好運氣...",
          "created_time": "2015-07-25T08:30:00+0000",
          "id": "352962731493606_729009783888897"
        },
        {
          "message": "[我十七歲的時候，雖然國文成績優異，不可能寫錯字，但是，我十七歲的時候呢？當時的我，"
        }
      ]
    }
  ]
}
```

Elements Network Sources Timeline Profiles Resources Audits Console

View: Options: Preserve log Disable cache

Filter Name Path

?access_token=CAACEdEose0cBANqVESvUWtZAq3hH8txuE...
graph.facebook.com/v2.4/352962731493606

All XHR Script Style Images Media Fonts Documents WebSockets Other Hide data URLs

Headers Preview Response Timing

General

Remote Address: 31.13.70.1:443
Request URL: https://graph.facebook.com/v2.4/352962731493606/?access_token=CAACEdEose0cBANqVESvUWtZAq3hH8txuE27U3Ym0XD5wHwUkIxGuhyzC5cNy...
QZABmx3qWB6ZA7u0kpUhrl56cnv03lZBkeAECjqZBLKpg5ZA7HCrlR6atu0EfRl09dw2dkI7ExyzFPC13Yt40Y1JzGvphCp6H2ZC6XfKUhxg9nXnwZC8oWPyZBdc...
ZB6GCKVuwZA6HzmNumR9NeIXwTtjcZD&debug=all&fields=posts&format=json&method=get&pretty=0&suppress_http_code=1
Request Method: GET
Status Code: 200 OK

Response Headers

access-control-allow-origin: *
cache-control: private, no-cache, no-store, must-revalidate
content-encoding: gzip

1 / 4 requests | 3.4 KB / 4.6 KB transferred

```
In [2]: library(httr)
library(stringr)
library(dplyr)
```

```
In [3]: url = "https://graph.facebook.com/v2.4/352962731493606/?access_token=CAACEdEose0cBAPtikGAr1T6lXkEgynRnsZCofv87ioaB6F2Jx"
```

```
In [4]: res = content(GET(url))
```

```
In [5]: res
```

```
Out[5]: $posts
$data
1. $message
  '#最新'
  $created_time
  '2015-07-26T03:56:53+0000'
  $id
  '352962731493606_729357620520780'

2. $message
  '用影片讓人更了解過去'
  $created_time
  '2015-07-26T03:50:00+0000'
  $id
  '352962731493606_729335113856364'

3. $message
```

先把 蘋果新聞的 po文抓下來

```
In [6]: tbl = lapply(res$posts$data, data.frame, stringsAsFactors=FALSE)
```

```
In [7]: tbl
```

```
Out[7]: 1.
```

	message	created_time	id
1	#最新	2015-07-26T03:56:53+0000	352962731493606_729357620520780

```
2.
```

	message	created_time	id
1	用影片讓人更了解過去	2015-07-26T03:50:00+0000	352962731493606_729335113856364

```
3.
```

	message	created_time	id
1	曹錦輝下放3A	2015-07-26T03:40:11+0000	352962731493606_729354070521135

```
4.
```

	message	created_time	id
1	好奇拍攝出來會是什麼樣子	2015-07-26T03:20:01+0000	352962731493606_729334647189744

```
5.
```

```
In [8]: tbl = rbind_all(tbl)
```

```
In [9]: head(tbl) 整理成表格
```

```
Out[9]:
```

	message	created_time	id
1	#最新	2015-07-26T03:56:53+0000	352962731493606_729357620520780
2	用影片讓人更了解過去	2015-07-26T03:50:00+0000	352962731493606_729335113856364
3	曹錦輝下放3A	2015-07-26T03:40:11+0000	352962731493606_729354070521135
4	好奇拍攝出來會是什麼樣子	2015-07-26T03:20:01+0000	352962731493606_729334647189744
5	該架小型飛機事發時載著5個人，機上兩人以及地面一名居民可能走避不及....	2015-07-26T03:05:36+0000	352962731493606_729340873855788

時間最好要切開來儲存 ...
以後要 Query 比較方便

ex. 20150725 , 130513

可以回想前面的正規表達，
我們使用 stringr 套件的 str_replace function

把 created_time 欄位內的 - : +0000 給拿掉

```
In [10]: tbl$created_time = str_replace_all(tbl$created_time, "-|:|[+]0000", "")  
  
In [11]: tbl$created_time[1:10]  
  
Out[11]: '20150726T035653' '20150726T035000' '20150726T034011' '20150726T032001' '20150726T030536' '20150726T025000'  
         '20150726T024401' '20150726T022000' '20150726T015001' '20150726T005001'  
  
In [12]: str_split(tbl$created_time, "T")[1:10]  
  
Out[12]: 1. '20150726' '035653'  
         2. '20150726' '035000'  
         3. '20150726' '034011'  
         4. '20150726' '032001'  
         5. '20150726' '030536'  
         6. '20150726' '025000'  
         7. '20150726' '024401'  
         8. '20150726' '022000'  
         9. '20150726' '015001'  
        10. '20150726' '005001'
```

利用 T 這個字，把日期時間給切開

```
In [13]: time_tbl = do.call(rbind, str_split(tbl$created_time, "T"))
```

```
In [14]: time_tbl = as.data.frame(time_tbl, stringsAsFactors = FALSE)
```

```
In [15]: names(time_tbl) = c('date', 'time')
```

```
In [16]: head(time_tbl)
```

Out[16]:

	date	time
1	20150726	035653
2	20150726	035000
3	20150726	034011
4	20150726	032001
5	20150726	030536
6	20150726	025000

把切開的時間日期, 合併成一張表格

```
In [17]: tbl = cbind(tbl, time_tbl )[-2]
```

把 date 和 time 合併到原本的 tbl 表格

```
In [18]: head(tbl)
```

[-2] 是指原本 tbl 的第二欄 create_time 我不要了

Out[18]:

	message	id	date	time
1	#最新	352962731493606_729357620520780	20150726	035653
2	用影片讓人更了解過去	352962731493606_729335113856364	20150726	035000
3	曹錦輝下放3A	352962731493606_729354070521135	20150726	034011
4	好奇拍攝出來會是什麼樣子	352962731493606_729334647189744	20150726	032001
5	該架小型飛機事發時載著5個人，機上兩人以及地面一名居民可能走避不及....	352962731493606_729340873855788	20150726	030536
6	嗯哼~BJ4 #曹錦輝	352962731493606_729317440524798	20150726	025000

```
In [19]: tbl$id = str_replace_all(tbl$id, "352962731493606_", "") 順便把 id 多餘部分拿掉
```

```
In [20]: head(tbl)
```

Out[20]:

	message	id	date	time
1	#最新	729357620520780	20150726	035653
2	用影片讓人更了解過去	729335113856364	20150726	035000
3	曹錦輝下放3A	729354070521135	20150726	034011
4	好奇拍攝出來會是什麼樣子	729334647189744	20150726	032001
5	該架小型飛機事發時載著5個人，機上兩人以及地面一名居民可能走避不及....	729340873855788	20150726	030536
6	嗯哼~BJ4 #曹錦輝	729317440524798	20150726	025000

蘋果新聞半夜也 po 文
好認真喔 + _ +

資料庫操作

```
In [21]: library(RSQLite)
library(dplyr)
```

```
>Loading required package: DBI
```

```
In [22]: setwd('/tmp//RCrawler')
```

```
In [23]: getwd()
```

在 /tmp/RCrawler 目錄下，創造一個資料庫檔案 my_fb.db

```
Out[23]: '/tmp/RCrawler'
```

```
In [24]: my_fb <- src_sqlite("my_fb.db", create = TRUE)
```

```
In [27]: dir("/tmp/RCrawler")
```

```
Out[27]: 'my_fb.db'
```

```
In [25]: dim(tbl)
```

```
Out[25]: 25 4
```

tbl 裡面目前有 25 筆料

```
In [26]: tbl[1:5, ]
```

```
Out[26]:
```

	message	id	date	time
1	駕照是怎麼考到的？	729354443854431	20150726	042000
2	黃國昌也要投入立委選戰了	729364717186737	20150726	041859
3	(頂眼鏡)	729360497187159	20150726	040755
4	#最新	729357620520780	20150726	035653
5	用影片讓人更了解過去	729335113856364	20150726	035000

```
In [27]: dbWriteTable(my_fb$con, "my_fb_table", tbl[1:5, ])
```

```
Out[27]: TRUE
```

把前五筆資料寫入到資料庫裡面，
表格名稱為 my_fb_table

```
In [28]: #dbRemoveTable(my_fb$con, "my_fb_table")
```

```
In [29]: tbl(my_fb, sql("SELECT * FROM my_fb_table"))
```

查詢看看有沒有資料

```
Out[29]: Source: sqlite 3.8.6 [my_fb.db]  
From: <derived table> [?? x 4]
```

	message	id	date	time
1	駕照是怎麼考到的？	729354443854431	20150726	042000
2	黃國昌也要投入立委選戰了	729364717186737	20150726	041859
3	(頂眼鏡)	729360497187159	20150726	040755
4	#最新	729357620520780	20150726	035653
5	用影片讓人更了解過去	729335113856364	20150726	035000
..

```
In [30]: dbWriteTable(my_fb$con, "my_fb_table", tbl[6:10, ], append=TRUE)
```

Out[30]: TRUE

再把第六到第十筆資料寫入到資料庫裡面

```
In [32]: tbl[6:10, ]
```

Out[32]:	message	id	date	time
6	曹錦輝下放3A	729354070521135	20150726	034011
7	好奇拍攝出來會是什麼樣子	729334647189744	20150726	032001
8	該架小型飛機事發時載著5個人，機上兩人以及地面一名居民可能走避不及....	729340873855788	20150726	030536
9	嗯哼~BJ4 #曹錦輝	729317440524798	20150726	025000
10	#最新	729331047190104	20150726	024401

```
In [33]: tbl(my_fb, sql("SELECT * FROM my_fb table WHERE rowid >= 6"))
```

Out[33]: Source: sqlite 3.8.6 [my_fb.db]
From: <derived table> [?? x 4]

透過rowid 來查詢第六筆以後的資料

1 message
2 曹錦輝下放3A
3 好奇拍攝出來會是什麼樣子
4 該架小型飛機事發時載著5個人，機上兩人以及地面一名居民可能走避不及....
5 懿哼-BJ4\n\n#曹錦輝
#最新
..
Variables not shown: id (chr), date (chr), time (chr) ...

如果我想另外開一個 R 來連結剛才的資料庫 並對資料庫進行關鍵字或時間查詢

```
In [2]: library(RSQLite)  
library(dplyr)
```

```
In [3]: my_fb <- src_sqlite("/tmp/RCrawler/my_fb.db")
```

```
In [4]: res = tbl(my_fb, sql("SELECT * FROM my_fb_table WHERE message LIKE '%曹錦輝%'"))
```

```
In [5]: collect(res)
```

Out[5]:

	message	id	date	time
1	曹錦輝下放3A	729354070521135	20150726	034011
2	嗯哼~BJ4 #曹錦輝	729317440524798	20150726	025000

```
In [6]: res = tbl(my_fb, sql("SELECT * FROM my_fb_table WHERE time LIKE '02%'"))
```

```
In [7]: collect(res)
```

Out[7]:

	message	id	date	time
1	嗯哼~BJ4 #曹錦輝	729317440524798	20150726	025000
2	#最新	729331047190104	20150726	024401

練習 5分鐘

自己寫入資料庫
查詢自己有興趣的關鍵字

Data Wrangling with dplyr and tidyverse

Cheat Sheet

R Studio

Syntax - Helpful conventions for wrangling

`dplyr::tbl_df(iris)`

Converts data to `tbl` class. `tbl`'s are easier to examine than data frames. R displays only the data that fits onscreen:

```
Source: local data frame [150 x 5]
  Sepal.Length Sepal.Width Petal.Length Petal.Width Species
 1          5.1         3.5          1.4      0.2   setosa
 2          4.9         3.0          1.4      0.2   setosa
 3          4.7         3.2          1.3      0.2   setosa
 4          4.6         3.1          1.5      0.2   setosa
 5          5.0         3.6          1.4      0.4   setosa
 6          5.4         3.9          1.7      0.4   setosa
 7          4.6         3.4          1.4      0.3   setosa
 8          5.0         3.4          1.5      0.2   setosa
...
Variables not shown: Petal.Width (dbl), Species (fctr)
```

`dplyr::glimpse(iris)`

Information dense summary of `tbl` data.

`utils::View(iris)`

View data set in spreadsheet-like display (note capital V).

iris					
	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
1	5.1	3.5	1.4	0.2	setosa
2	4.9	3.0	1.4	0.2	setosa
3	4.7	3.2	1.3	0.2	setosa
4	4.6	3.1	1.5	0.2	setosa
5	5.0	3.6	1.4	0.4	setosa
6	5.4	3.9	1.7	0.4	setosa
7	4.6	3.4	1.4	0.3	setosa
8	5.0	3.4	1.5	0.2	setosa

`dplyr::%>%`

Passes object on left hand side as first argument (or, argument) of function on righthand side.

```
x %>% f(y) is the same as f(x, y)
y %>% f(x, ., z) is the same as f(x, y, z)
```

"Piping" with `%>%` makes code more readable, e.g.

```
iris %>%
  group_by(Species) %>%
  summarise(av = mean(Sepal.Width)) %>%
  arrange(av)
```

來源網址

<https://www.rstudio.com/wp-content/uploads/2015/02/data-wrangling-cheatsheet.pdf>

Tidy Data - A foundation for wrangling in R

In a tidy data set:



Each **variable** is saved in its own column

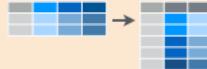


Each **observation** is saved in its own row

Tidy data complements R's **vectorized operations**. R will automatically preserve observations as you manipulate variables. No other format works as intuitively with R.

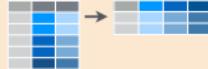


Reshaping Data - Change the layout of a data set



`tidyr::gather(cases, "year", "n", 2:4)`

Gather columns into rows.



`tidyr::spread(pollution, size, amount)`

Spread rows into columns.

`tidyr::separate(storms, date, c("y", "m", "d"))`

Separate one column into several.

`tidyr::unite(data, col, ..., sep)`

Unite several columns into one.

`dplyr::data_frame(a = 1:3, b = 4:6)`
Combine vectors into data frame (optimized).

`dplyr::arrange(mtcars, mpg)`
Order rows by values of a column (low to high).

`dplyr::arrange(mtcars, desc(mpg))`
Order rows by values of a column (high to low).

`dplyr::rename(tb, y = year)`
Rename the columns of a data frame.

Subset Observations (Rows)



`dplyr::filter(iris, Sepal.Length > 7)`

Extract rows that meet logical criteria.

`dplyr::distinct(iris)`

Remove duplicate rows.

`dplyr::sample_frac(iris, 0.5, replace = TRUE)`

Randomly select fraction of rows.

`dplyr::sample_n(iris, 10, replace = TRUE)`

Randomly select n rows.

`dplyr::slice(iris, 10:15)`

Select rows by position.

`dplyr::top_n(storms, 2, date)`

Select and order top n entries (by group if grouped data).

Subset Variables (Columns)



`dplyr::select(iris, Sepal.Width, Petal.Length, Species)`

Select columns by name or helper function.

Helper functions for select - ?select

`select(iris, contains("x"))`

Select columns whose name contains a character string.

`select(iris, ends_with("Length"))`

Select columns whose name ends with a character string.

`select(iris, everything())`

Select every column.

`select(iris, matches("t.*"))`

Select columns whose name matches a regular expression.

`select(iris, num_range("x", 1:5))`

Select columns named x1, x2, x3, x4, x5.

`select(iris, one_of("Species", "Genus"))`

Select columns whose names are in a group of names.

`select(iris, starts_with("Sepal"))`

Select columns whose name starts with a character string.

`select(iris, Sepal.Length:Petal.Width)`

Select all columns between Sepal.Length and Petal.Width (inclusive).

`select(iris, -Species)`

Select all columns except Species.

Logic in R - ?Comparison, ?base::Logic

<	Less than	!=	Not equal to
>	Greater than	%in%	Group membership
==	Equal to	is.na	Is NA
<=	Less than or equal to	!is.na	Is not NA
>=	Greater than or equal to	&, , !, xor, any, all	Boolean operators

來源網址

<https://www.rstudio.com/wp-content/uploads/2015/02/data-wrangling-cheatsheet.pdf>

Summarise Data



`dplyr::summarise(iris, avg = mean(Sepal.Length))`

Summarise data into single row of values.

`dplyr::summarise_each(iris, funs(mean))`

Apply summary function to each column.

`dplyr::count(iris, Species, wt = Sepal.Length)`

Count number of rows with each unique value of variable (with or without weights).



Summarise uses **summary functions**, functions that take a vector of values and return a single value, such as:

`dplyr::first`

First value of a vector.

`dplyr::last`

Last value of a vector.

`dplyr::nth`

Nth value of a vector.

`dplyr::n`

of values in a vector.

`dplyr::n_distinct`

of distinct values in a vector.

`IQR`

IQR of a vector.

`min`

Minimum value in a vector.

`max`

Maximum value in a vector.

`mean`

Mean value of a vector.

`median`

Median value of a vector.

`var`

Variance of a vector.

`sd`

Standard deviation of a vector.

Group Data

`dplyr::group_by(iris, Species)`

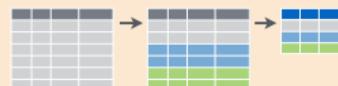
Group data into rows with the same value of Species.

`dplyr::ungroup(iris)`

Remove grouping information from data frame.

`iris %>% group_by(Species) %>% summarise(...)`

Compute separate summary row for each group.



Make New Variables



`dplyr::mutate(iris, sepal = Sepal.Length + Sepal.Width)`

Compute and append one or more new columns.

`dplyr::mutate_each(iris, funs(min_rank))`

Apply window function to each column.

`dplyr::transmute(iris, sepal = Sepal.Length + Sepal.Width)`

Compute one or more new columns. Drop original columns.



Mutate uses **window functions**, functions that take a vector of values and return another vector of values, such as:

`dplyr::lead`

Copy with values shifted by 1.

`dplyr::lag`

Copy with values lagged by 1.

`dplyr::dense_rank`

Ranks with no gaps.

`dplyr::min_rank`

Ranks. Ties get min rank.

`dplyr::percent_rank`

Ranks rescaled to [0, 1].

`dplyr::row_number`

Ranks. Ties get to first value.

`dplyr::ntile`

Bin vector into n buckets.

`dplyr::between`

Are values between a and b?

`dplyr::cume_dist`

Cumulative distribution.

`dplyr::cumall`

Cumulative all

`dplyr::cumany`

Cumulative any

`dplyr::cummean`

Cumulative mean

`cumsum`

Cumulative sum

`cummax`

Cumulative max

`cummin`

Cumulative min

`cumprod`

Cumulative prod

`pmax`

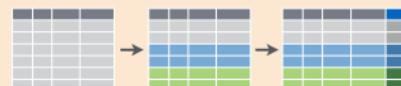
Element-wise max

`pmin`

Element-wise min

`iris %>% group_by(Species) %>% mutate(...)`

Compute new variables by group.



Combine Data Sets



Mutating Joins

`dplyr::left_join(a, b, by = "x1")`

Join matching rows from b to a.

`dplyr::right_join(a, b, by = "x1")`

Join matching rows from a to b.

`dplyr::inner_join(a, b, by = "x1")`

Join data. Retain only rows in both sets.

`dplyr::full_join(a, b, by = "x1")`

Join data. Retain all values, all rows.

Filtering Joins

`dplyr::semi_join(a, b, by = "x1")`

All rows in a that have a match in b.

`dplyr::anti_join(a, b, by = "x1")`

All rows in a that do not have a match in b.



Set Operations

`dplyr::intersect(y, z)`

Rows that appear in both y and z.

`dplyr::union(y, z)`

Rows that appear in either or both y and z.

`dplyr::setdiff(y, z)`

Rows that appear in y but not z.

Binding

`dplyr::bind_rows(y, z)`

Append z to y as new rows.

`dplyr::bind_cols(y, z)`

Append z to y as new columns.

Caution: matches rows by position.

Source: https://rstudio.com/wp-content/uploads/2015/02/data-wrangling-cheatsheet.pdf
Version: 2015-02-02
Last updated: 2015-02-02
R version: 3.2.3
Platform: x86_64-pc-linux-gnu
RStudio: 0.98.1101
tidyverse: 0.4.0
dplyr: 0.2.0
Updated: 1/15

我的爬蟲想自動化 ...

可以每天每小時就去抓一次嗎 (￣▽￣)~*

```
fb_bot.R
```

```
1 library(httr)
2 library(dplyr)
3 library(logging)
4
5 basicConfig()
6 addHandler(writeToFile, file="/tmp/RCrawler/fb_bot.log")
7
8 url = "https://graph.facebook.com/v2.4/352962731493606/?access_token=CAACEdEose0cBAHXt13FUZCC86TRWhaTu0a4BE6dNI0ePCPH6PQFcN9jzXiRiUZARjeVvJeqfY5YW9TnTr8M"
9 res = content(GET(url))
10 tbl = lapply(res$posts$data, data.frame, stringsAsFactors=FALSE)
11 tbl = rbind_all(tbl)
12 print(head(tbl))
13
14 loginfo("finished %s", "352962731493606")
```

```
apple — root@localhost: /tmp/RCrawler — ssh — 80x29
root@localhost: /tmp/RCrawler# Rscript fb_bot.R
```



可以在 command line
利用 Rscript 執行 R 程式

```
fb_bot.R
```

```
1 library(httr)
2 library(dplyr)
3 library(logging)
4
5 basicConfig()
6 addHandler(writeToFile, file="/tmp/RCrawler/fb_bot.log")
7
8 url = "https://graph.facebook.com/v2.4/352962731493606/?access_token=CAACEdEose0cBAHXt13FUZCC86TRWhaTu0a4BE6dNI0ePCPH6PQFcN9jzXiRiUZARjeVvJeqfY5YW9TnTr8Mw"
9 res = content(GET(url))
10 tbl = lapply(res$posts$data, data.frame, stringsAsFactors=FALSE)
11tbl = rbind_all(tbl)
12 print(head(tbl))
13
14 loginfo("finished %s", "352962731493606")
```

```
apple — root@localhost: /tmp/RCrawler — ssh — 80x29
root@localhost: /tmp/RCrawler
```

```
The following objects are masked from 'package:stats':
  filter, lag

The following objects are masked from 'package:base':
  intersect, setdiff, setequal, union

Loading required package: methods
Source: local data frame [6 x 3]
```

message
1 小編...小編還沒老...
2 一名30歲女子與兒子一起乘扶手梯上樓，快抵達樓層時腳下踏板突然鬆脫，眼看就要被捲進扶手梯時，母親在最後一刻舉起兒子，自己卻不幸被捲入手扶梯下方的機械間...
3 小知識，希望在意外發生時可以幫上大家
4 賀！
5 執行結果
6 末段感言相當值得回味...
請看完內文再下評論

```
Variables not shown: created_time (chr), id (chr)
2015-07-26 13:09:58 INFO::finished 352962731493606
root@localhost:/tmp/RCrawler#
```

65

接著來設定一下爬蟲的自動化排程

```
apple — root@localhost: /tmp/RCrawler — ssh — 80x29
root@localhost:/tmp/RCrawler# crontab -e
```

按下 Enter 之後會出現

```
# Edit this file to introduce tasks to be run by cron.
#
# Each task to run has to be defined through a single line
# indicating with different fields when the task will be run
# and what command to run for the task
#
# To define the time you can provide concrete values for
# minute (m), hour (h), day of month (dom), month (mon),
# and day of week (dow) or use '*' in these fields (for 'any').#
# Notice that tasks will be started based on the cron's system
# daemon's notion of time and timezones.
#
# Output of the crontab jobs (including errors) is sent through
# email to the user the crontab file belongs to (unless redirected).
#
# For example, you can run a backup of all your user accounts
# at 5 a.m every week with:
# 0 5 * * 1 tar -zcf /var/backups/home.tgz /home/
#
# For more information see the manual pages of crontab(5) and cron(8)
#
# m h dom mon dow   command
~
~
~
~
~
~/tmp/crontab.bezMwp/crontab" 22L, 888C          1,1      All
```

```
apple — root@localhost: /tmp/RCrawler — ssh — 80x29
root@localhost: /tmp/RCrawler
# Edit this file to introduce tasks to be run by cron.
#
# Each task to run has to be defined through a single line
# indicating with different fields when the task will be run
# and what command to run for the task
#
# To define the time you can provide concrete values for
# minute (m), hour (h), day of month (dom), month (mon),
# and day of week (dow) or use '*' in these fields (for 'any').#
# Notice that tasks will be started based on the cron's system
# daemon's notion of time and timezones.
#
# Output of the crontab jobs (including errors) is sent through
# email to the user the crontab file belongs to (unless redirected).
#
# For example, you can run a backup of all your user accounts
# at 5 a.m every week with:
# 0 5 * * 1 tar -zcf /var/backups/home.tgz /home/
#
# For more information see the manual pages of crontab(5) and cron(8)
#
# m h  dom mon dow   command
*/1 * * * * Rscript /tmp/RCrawler/fb_bot.R
~
~ 每分鐘就用 Rscript 去執行 fb_bot.R
~ 
~/tmp/crontab.AqHNQj/crontab" 23L, 931C          1,1      All
```

<http://www.puritys.me/docs-blog/article-20-cron-jobs-crontab-%E6%8E%92%E7%A8%8B%E6%95%99%E5%AD%B8.html>

```
fb_bot.R
```

```
1 library(httr)
2 library(dplyr)
3 library(logging)
4
5 basicConfig()
6 addHandler(writeToFile, file="/tmp/RCrawler/fb_bot.log")
7
8 url = "https://graph.facebook.com/v2.4/352962731493606/?access_token=CAACEdEose0cBAHxtl3FUZCC86TRWhaTu0a4BE6dNI0ePCPH6PQFcN9jzXiRiUZARjeVvJeqfY5YW9TnTr8M"
9 res = content(GET(url))
10 tbl = lapply(res$posts$data, data.frame, stringsAsFactors=FALSE)
11 tbl = rbind_all(tbl)
12 print(head(tbl))
13
14 loginfo("finished %s", "352962731493606")
```

```
apple — root@localhost: /tmp/RCrawler — ssh — 80x29
root@localhost: /tmp/RCrawler#
root@localhost:/tmp/RCrawler# cat fb_bot.log
2015-07-26 13:09:02 INFO::finished 352962731493606
2015-07-26 13:09:47 INFO::finished 352962731493606
2015-07-26 13:09:58 INFO::finished 352962731493606
2015-07-26 13:10:03 INFO::finished 352962731493606
2015-07-26 13:11:02 INFO::finished 352962731493606
2015-07-26 13:12:02 INFO::finished 352962731493606
2015-07-26 13:13:02 INFO::finished 352962731493606
2015-07-26 13:14:02 INFO::finished 352962731493606
2015-07-26 13:15:02 INFO::finished 352962731493606
root@localhost:/tmp/RCrawler#
```

每分鐘就會產生一筆 log 記錄

練習 5分鐘

觀察看看有沒有 log 自動產生

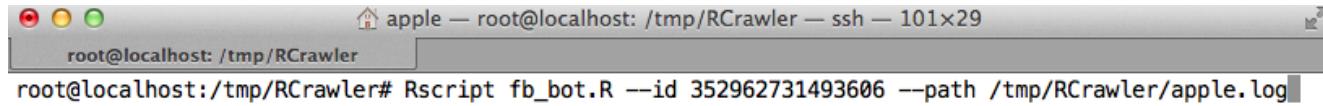
可以不只抓蘋果新聞嗎？我還想抓其他粉絲團 ...

R 語言可以吃參數，利用參數可以達到靈活性

ex. 直接在 command line 執行

`Rscript fb_bot.R --id xxxxx --path xxxx`

```
1 library(httr)
2 library(dplyr)
3 library(logging)
4 library(optparse)
5
6 option_list = list(
7   make_option(c("-i", "--id"), action="store", default="me", type='character',
8     help="please enter the facebook id"),
9   make_option(c("-p", "--path"), action="store", default="/tmp/RCrawler/fb_bot.log", type='character',
10     help="please enter the log path")
11 )
12
13 opt = parse_args(OptionParser(option_list=option_list))
14
15 basicConfig()
16 addHandler(writeToFile, file=opt$path)
17
18 TOKEN = "CAACEEdEose0cBAHxtl3FUZCC86TRWhaTu0a4BE6dNI0ePCPH6PQFcN9jzXiRiUZARjeVvJeqfY5YW9TnTr8Mw2ZBy4Q1T0NK57MoCD6tLiXMKzquv89D2wlNeRg8xssBwKw050CTEQwJ1bP22A"
19
20 url = paste0("https://graph.facebook.com/v2.4/", opt$id, "/?access_token=", TOKEN, "&debug=all&fields=posts&format=json&method=get&pretty=0&suppress_http_c
21 res = content(GET(url))
22 tbl = lapply(res$posts$data, data.frame, stringsAsFactors=FALSE)
23 tbl = rbind_all(tbl)
24 print(head(tbl))
25
26 loginfo_("finished %s", opt$id)
```



A screenshot of a terminal window titled "apple — root@localhost: /tmp/RCrawler — ssh — 101x29". The window shows the command "root@localhost:/tmp/RCrawler# Rscript fb_bot.R --id 352962731493606 --path /tmp/RCrawler/apple.log" being typed. The command is highlighted with red underlines.

這樣就能帶參數執行 R 程式
(其實是我平常跑實驗很愛用 XD)

記得順便到 crontab -e 裡面去做修改 :)

那如果程式跑到一半，網路不穩或連線異常斷掉
我可以讓程式自動重跑嗎 ???

ex. 發現它斷線了，自動重新連線，最多重連3次

```
fb_bot.R
1 library(httr)
2 library(dplyr)
3 library(logging)
4 library(optparse)
5
6
7 option_list = list(
8
9   make_option(c("-i", "--id"), action="store", default="me", type='character',
10             help="please enter the facebook id"),
11
12   make_option(c("-p", "--path"), action="store", default="/tmp/Rcrawler/fb_bot.log", type='character',
13             help="please enter the log path")
14 )
15
16 opt = parse_args(OptionParser(option_list=option_list))
17
18 basicConfig()
19 addHandler(writeToFile, file=opt$path)
20
21 TOKEN = "CAACEdEose0cBAHXtl3FUZCC86TRWhaTu0a4BE6dNI0ePCPH6PQFcN9jzXiRiUZARjeVvJeqfY5YW9TnTr8Mw2ZBy4Q1T0NK57MoCD6tLiXMKzquv89D2w1NeRg8xssBwKwO50CTEqwJ1bP2
22
23 url = paste0("https://graph.facebook.com/v2.4/", opt$id, "/?access_token=", TOKEN, "&debug=all&fields=posts&format=json&method=get&pretty=0&suppress_http
24
25 for(i in 1:3){
26
27   tryCatch({
28     res = content(GET(url))
29     break
30
31   }, error = function(e){
32     logerror("url :%s", url)
33     logerror(e)
34     Sys.sleep(5)
35   })
36 }
37
38 tbl = lapply(res$posts$data, data.frame, stringsAsFactors=FALSE)
39 tbl = rbind_all(tbl)
40 print(head(tbl))
41
42 loginfo("finished %s", opt$id)
43
```

