

WEB CRAWLER

張家瑋 博士

新漢股份有限公司創新工業4.0中心顧問

PTT Introduction

作者 oppo5566 (5566)
標題 [問卦] 發錢 預測大谷翔平本日打擊
時間 Thu Apr 12 07:26:39 2018

大谷翔平 本日要先發打擊對上遊騎兵左投Matt moore，

這是上次鄉民預測中的發錢名單

以及收到P幣之後的感謝回信

<https://i.imgur.com/vEm8tme>. 推 badbadook: 好屌
推 robinyu85: 2/0/1

那這次要預測的推文格式為：安

範例：3/1/1

前十位預測中的鄉民稅後各100F

PS.不用擔心錢不夠發，有朋友贊助

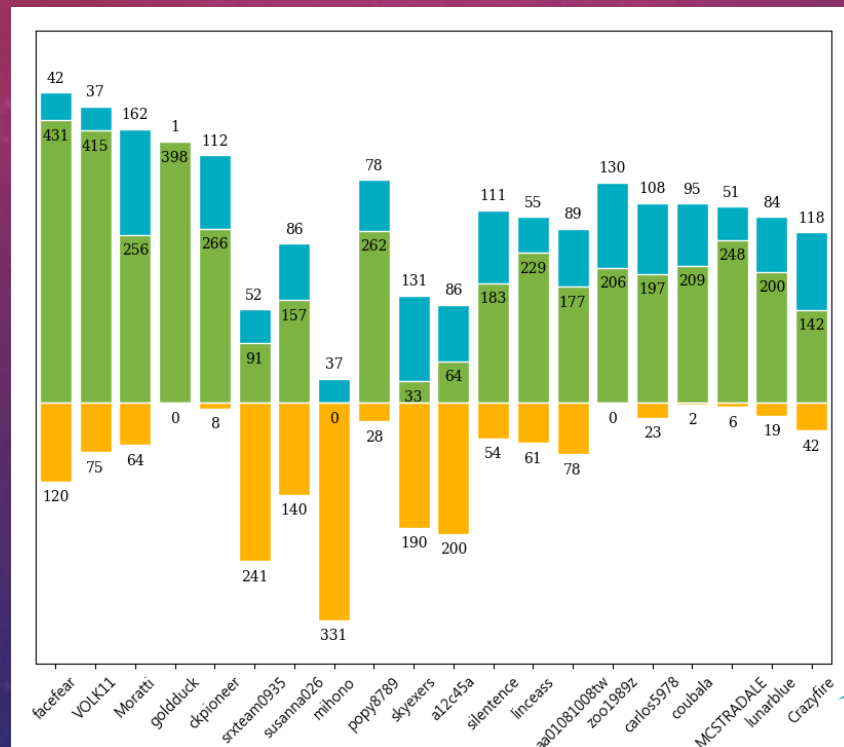
【板主:Kay731/RS556/Ra...】【八卦板】置底協尋文，請大家幫忙！				看板 (Gossiping)
編號	日期	作者	文章標題	回覆人數:14099
786165	+ 4	4/12	ahuang80919	R: 新聞 台灣人島國心態 柯P:電視很少國際新聞
786164	+ 2	4/12	hachilou	R: 問卦 說中國古代文化在日本的是什麼人?
786165	+X3	4/12	Safin	R: 新聞 陳菊任府院:現階段選舉本是重要工作
786166	+X2	4/12	borondawon	R: F B 陳沂:很多台男放任自己變肥宅
786167	+10	4/12	mike901003	R: 新聞 法航開航 桃園機場飛巴黎已漫百分百
786168	+ 4	4/12	RonaldReagan	R: 問卦 有沒有台南奇美醫院科謝俊民醫生的八卦
786169	+ 4	4/12	lianpig5566	R: 問卦 勵! 且幾歲才到臉
786170	+ 7	4/12	aghs386690	R: 新聞 南部茶棧車版連3天超時駕駛 遊覽車GPS還
786171	+ 7	4/12	1ElITS	R: 問卦 公司最近提供沐浴設備
786172	+ 5	4/12	xiaopen74269	R: 問卦 原來一箭可以切屎巴
786173	+ 8	4/12	CORSA	R: 問卦 為啥越南/北韓不是中國不可分割的一部分?
786174	+ 4	4/12	bora	R: 問卦 要怎麼把川普這顆棋子效用最大化?
786175	+ 5	4/12	Cocochia	R: 問卦 臉書是不是快掛了?
786176	+99	4/12	Eliphalet	R: 新聞 台北LOL夢到被人偷摸! KTV醒來怒告男同事
786177	+ 4	4/12	vmlinux	R: 問卦 台大校長與斯隆撤令 是不是很像?
786178	+ 4	4/12	penisman	R: 問卦 有沒有台南奇美醫院科謝俊民醫生的八
786179	+95	4/12	arrenwu	R: F B 陳沂:很多台男放任自己變肥宅
786180	+ 9	4/12	Pattaya	R: 新聞 葉俊榮塔中誦禪經 還不止一次
786181	+ 4	4/12	waymayday	R: 新聞 黃網掛禁網12億 妙天也曾被告訴詐欺
786182	+ 4	4/12	qno5566	R: 問卦 發給 預測大谷理平本月打臉

(y) 回應 (X) 推文 (X) 轉錄 (E) 相關主題 (V) 投標題 (V) 進板畫面

0.82), 04/12/2018 07:28:49

是打我就

Top 20 踴躍留言者分析的圖表實作



推

→

噓

Top 20
的鄉民ID

情緒辨識

輸入留言！偵測你的情緒狀態：罪有英德
罪有英德 ---你很不開心喔？

輸入留言！偵測你的情緒狀態：賴清德
賴清德 ---你很不開心喔？

輸入留言！偵測你的情緒狀態：蔡英文
蔡英文 ---感覺還不錯哦！

輸入留言！偵測你的情緒狀態：勞動部
勞動部 ---你很不開心喔？

輸入留言！偵測你的情緒狀態：勞基法
勞基法 ---你很不開心喔？

輸入留言！偵測你的情緒狀態：健保
健保 ---你很不開心喔？



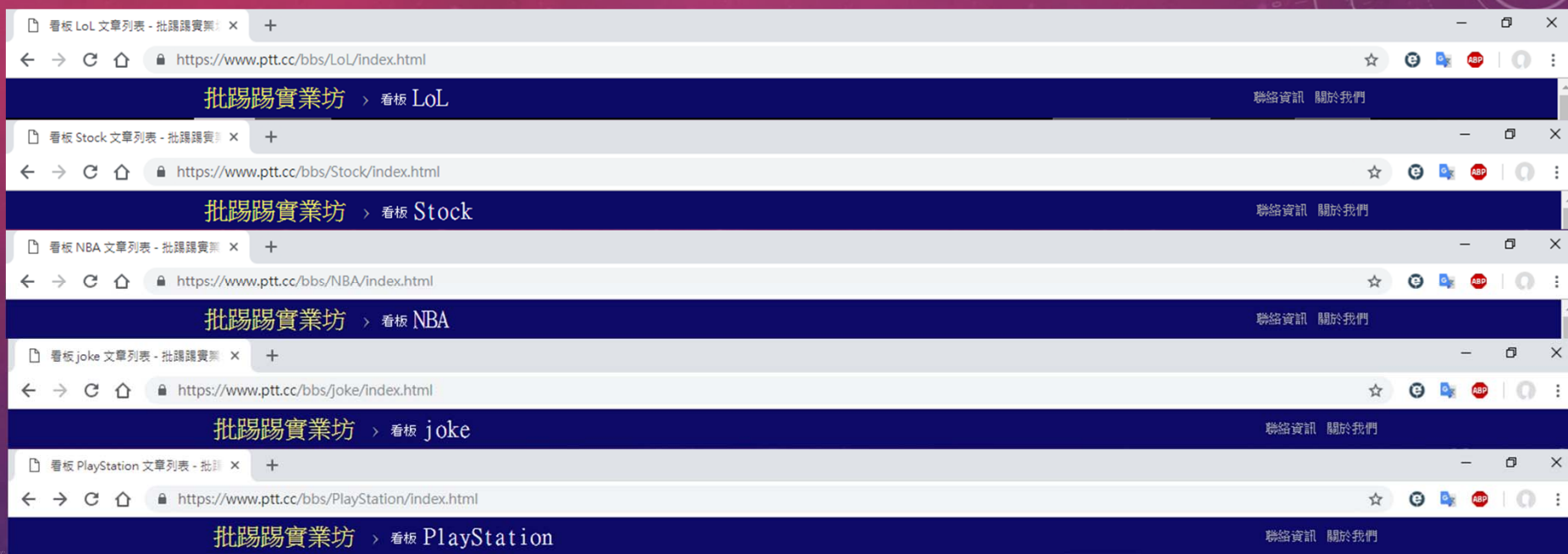
「今天來帶各位用PYTHON爬PTT」

做個友善的英雄聯盟玩家，一起支持G-REX

<https://www.ptt.cc/bbs/LoL/index.html>

批踢踢實業坊 > 看板 LoL		聯絡資訊 關於我們	
看板	精華區	最舊	最新
搜尋文章...			
9	[閒聊] 法洛士的武器怎麼不會跟他講話?	lycs0908	10/07 ...
5	[問題] 今年跟S4，哪一年的台灣比較強?	wyner	10/07 ...
5	[閒聊] S6和今年的MMD哪個比較難受?	bygamantou	10/07 ...
19	[閒聊] GBM 推特	HaiTurtle	10/07 ...
9	[閒聊] LMS終於可以和其它四賽區平起平坐了嗎?	lovealgebra	10/07 ...
44	[閒聊] LMS 今年全明星賽要派哪 2 個選手去?	S890127	10/07 ...
1	[閒聊] 考特 Scott FB	lovealgebra	10/07 ...
71	[外絮] Toyz Instagram	iamwhoim	M 10/07 ...
	Re: [閒聊] LMS 今年全明星賽要派哪 2 個選手去?	FeiWenKing	10/07 ...

逛PTT的過程中你會發現PTT的URL網址都有一個固定格式
`https://www.ptt.cc/bbs/<看板名稱>/index.html`



事前作業

- 開始爬蟲前，我們必須先安裝一些Python套件
- `$ pip install requests`
- `$ pip install bs4`

Get hrefs



The screenshot shows a web browser window with the address bar displaying `https://www.ptt.cc/bbs/Lol/index.html`. The page is the "批踢踢實業坊" (PTT) forum, specifically the "看板 LoL" (LoL board). The page has a dark blue header with navigation links like "看板", "精華區", "最舊", "上頁", "下頁", and "最新". A search bar is present with the text "搜尋文章...". The main content area lists forum posts. A context menu is open over the 14th post, which is titled "[閒聊] LMS 今年全明星賽要派哪 2 個選手去?". The context menu options are: "在新分頁中開啟連結(T)", "在新視窗中開啟連結(V)", "在無框架視窗中開啟連結(S)", "另存連結為(K)...", "複製連結網址(E)", "IE Tab Options", and "檢查(N)". The URL bar at the bottom shows the full URL of the selected post: `https://www.ptt.cc/bbs/Lol/M.1538920763.A.52A.html`.

Post Number	Post Title	Author	Date
11	[閒聊] 法洛士的武器怎麼不會跟他講話?	lycs0908	10/07
5	[閒聊] 今年跟S4, 哪一年的台灣比較強?	wyner	10/07
5	[閒聊] S6和今年的MMD哪個比較難受?	bygamantou	10/07
24	[閒聊] GBM 推特	HaiTurtle	10/07
14	[閒聊] LMS 今年全明星賽要派哪 2 個選手去?	lovealgebra	10/07
56	[閒聊] LMS 今年全明星賽要派哪 2 個選手去?	S890127	10/07
1	[閒聊] 考特 S	lovealgebra	10/07
爆	[外網] Toyz Instagram	iamwhoim	10/07
14	Re: [閒聊] LMS 今年全明星賽要派哪 2 個選手去?	FeiWenKing	10/07

Get hrefs

The screenshot shows a web browser window displaying a forum page titled "批踢踢實業坊" (Ptt). The page lists several forum posts. The browser's developer tools are open, showing the "Elements" panel. A red box highlights the following HTML code snippet:

```
<a href="/bbs/Lol/M.1538920763.A.52A.html">[外架] Toyz Instagram</a> == $0
```

The "Styles" panel on the right shows the default link styles from the forum's CSS.

Get hrefs

點進來可以發現每篇文章的網址跟剛剛選取的Href是一樣的



sbug.py - 未命名 (工作區) - Visual Studio Code

檔案 (F) 編輯 (E) 選取項目 (S) 檢視 (V) 前往 (G) 偵錯 (D) 終端機 (T) 說明 (H)

sbug.py x

```
10
11 for item in results:
12     item_href = item.find("a").attrs["href"] #用迴圈去把每個div中a標籤的 'href' 找出來因為每個div只有一個a所以只需要用find("a")
13     article_href.append(item_href) # 將每個解析出來的href 放到list裡面
14     print(item_href)
```

問題 輸出 偵錯主控台 終端機 1: Python

au@Au: ~/Pythonpractice-/pythonbug\$ /usr/bin/python3 /home/au/Pythonpractice-/pythonbug/sbug.py

```
/bbs/Lol/M.1538971790.A.244.html
/bbs/Lol/M.1538972118.A.203.html
/bbs/Lol/M.1538972249.A.7C2.html
/bbs/Lol/M.1538972300.A.072.html
/bbs/Lol/M.1538973525.A.BAF.html
/bbs/Lol/M.1538974777.A.101.html
/bbs/Lol/M.1538975185.A.61F.html
/bbs/Lol/M.1538975992.A.DF9.html
/bbs/Lol/M.1538976279.A.EFE.html
/bbs/Lol/M.1538976985.A.497.html
/bbs/Lol/M.1538977637.A.999.html
/bbs/Lol/M.1538978203.A.C27.html
/bbs/Lol/M.1538978304.A.C94.html
/bbs/Lol/M.1538978346.A.DFD.html
/bbs/Lol/M.1533315732.A.0CE.html
/bbs/Lol/M.1506551347.A.B64.html
/bbs/Lol/M.1537379839.A.C35.html
/bbs/Lol/M.1538305719.A.7AD.html
/bbs/Lol/M.1538915112.A.889.html
```

master Python 3.6.6 64-bit 第 14 行, 第 21 欄 空格: 4 UTF-8 LF Python

Parser Article

可以看到「作者」、「看板」、「標題」、「時間」資料
都放在 ``

```
▼<div class="article-metaline">
  <span class="article-meta-tag">作者</span>
  <span class="article-meta-value">coolplus (cool)</span>
</div>
▼<div class="article-metaline-right">
  <span class="article-meta-tag">看板</span>
  <span class="article-meta-value">LoL</span>
</div>
▼<div class="article-metaline">
  <span class="article-meta-tag">標題</span>
  <span class="article-meta-value">[閒聊] GREX教練阿WEI：我們最有信心對100T</span>
</div>
▼<div class="article-metaline">
  <span class="article-meta-tag">時間</span>
  <span class="article-meta-value">Mon Oct 8 11:27:35 2018</span>
```

Parser Article

```
author = soup.select('span.article-meta-value')[0].text #作者
board = soup.select('span.article-meta-value')[1].text #看板
title = soup.select('span.article-meta-value')[2].text #標題
time = soup.select('span.article-meta-value')[3].text #時間
push_tag = soup.select('span.push-tag') #推文
push_userid = soup.select('span.push-userid') #推文id
push_content = soup.select('span.push-content') #推文內容
push_ipdatetime = soup.select('span.push-ipdatetime') #推文時間
print('作者:', author)
print(board, ' 看版')
print('標題:', title)
print('時間:', time)
push_list_len = len(push_tag) #計算推文筆數
count = 0
while (count < push_list_len):
    print (push_tag[count].text + push_userid[count].text + push_content[count].text + push_ipdatetime[count].text)
    count = count+1
print("=====分隔線=====")
```


Parser Article

處理內文的資料

```
content = soup.find(id="main-content").text  
target_content = u'※ 發信站: 批踢踢實業坊(ptt.cc), '  
content = content.split(target_content)  
content = content[0].split(time)  
main_content = content[1].replace('--', '  ')
```

#content 文章內文
#去除掉 target_content
#去除掉文末 --

Output

將剛才所整理好的所有資料用print打印出來

```
print('作者:', author)
print(board, ' 版')
print('標題:', title)
print('時間:', time)
print('內文:', main_content)

push_list_len = len(push_tag)          #計算推文筆數
count = 0

while (count < push_list_len):         #利用迴圈印出所有推文
    print (push_tag[count].text + push_userid[count].text + push_content[count].text + push_ipdatetime[count].text)
    count = count+1
```

輸出結果

=====分隔線=====

作者: airiguodala

LoL 版

標題: [公告] LoL 板 開始舉辦樂透!

時間: Sun Oct 7 20:25:11 2018

內文:

特!

請到 LoL 板 按 'f' 參與樂透!

16:00 KT vs TL 17:00 EDG vs MAD 猜TL跟MAD誰撈的比較久(比較慢輸)

1.TL

2.MAD

3.TL/MAD贏 開獎順序:3>2=1

因為間隔都很短 又不像LMS這麼會修 開一次混合的試試看

如果下注狀況不好之後就不開惹

一張 100 Ptt幣 (平民級)

樂透結束時間: 10/10/2018 16:05:17 Wed

※編輯: airiguodala (114.36.7.34), 10/07/2018 20:26:08

推 ilove640 : 特! 10/07 20:25

→ ray221740718: 特! 10/07 20:25

→ asd860079 : 特! 10/07 20:25

推 dinosaur8484: 特!!!!!!!!!!!!!! 10/07 20:25

推 soalzelance : 特! 10/07 20:26

推 Ball0427 : 特! 10/07 20:26

推 hkr91511208 : 3 10/07 20:26

推 Znps : 特! 10/07 20:26

推 marginalFeng: <https://i.imgur.com/RANbro2.jpg> 10/07 20:26

推 tw15 : 三是什麼啊 10/07 20:26