

備份 pchome 相簿和其照片文字敘述

```

1 import requests, bs4
2 import os
3 import re
4 import html5lib
5
6 with open('pchome_album.txt', mode='rt', encoding='utf-8') as f:
7     selections = f.readlines()
8     user_url = selections[0].strip().upper()
9     # check if url of the 1st line is valid or not
10    if not user_url.startswith('HTTP'):
11        print('請以http開始')
12        quit()
13    if not 'PHOTO.PCHOME.COM.TW' in user_url:
14        print('可能不是pchome相簿')
15        quit()
16
17 URL_HOME = user_url + '{}'
18
19 options = []
20 for selection in selections[1:5]:
21     try:
22         album_range = int(selection.strip())
23     except ValueError:
24         album_range = None
25     options.append(album_range)
26
27 headers = {'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; WOW64) AppleWebKit/
28 /537.36 (KHTML, like Gecko) Chrome/50.0.2661.102 Safari/537.36'}
29
30 DOMAIN = 'http://photo.pchome.com.tw'
31
32 user_page = 1
33
34 # collect links even if spanning multiple pages
35 links = []
36 while True:
37     link_res = requests.get(URL_HOME.format(user_page), headers=headers)
38     link_res.raise_for_status()
39     link_soup = bs4.BeautifulSoup(link_res.content, 'html5lib')
40     links_tmp = link_soup.select('div#alb > a') (1)
41     if not links_tmp:
42         break
43     links.extend(links_tmp)
44     user_page = user_page + 1
45
46 pattern = re.compile("[【】?!-顧。： “ ” ‘ ’ ; \w\s,.-]")
47
48 album_index = options[3]
49 for link in links[options[0]:options[1]:options[2]]:
50     album_url = DOMAIN + link['href'][::-1] + '{}' # discard extra '/'
51     # first page
52     link2_res = requests.get(album_url.format(1), headers=headers)
53     link2_res.raise_for_status()
54     link2_soup = bs4.BeautifulSoup(link2_res.content, 'html5lib')
55     title = link2_soup.select_one('div.tit a').text.strip()
56
57     valid_name = "".join(ch for ch in title if pattern.match(ch)) (2)
58     title = valid_name
59     title = title.strip()
60
61     abstract = link2_soup.select_one('div.ab').text.strip() (3)

```

先選稍後要用到的模組

粗略檢查網址是否合理，若不則離開程式

對使用者指定的選項做些處理

收集該網站的所有相簿的位址連結

一次只處理一本相簿

過濾相簿命名和敘述

```

63     # non-first pages if any (4)
64     page = 2
65     while True:
66         link2b_res = requests.get(album_url.format(page), headers=headers)
67         link2b_res.raise_for_status()
68         link2b_soup = bs4.BeautifulSoup(link2b_res.content, 'html5lib')
69         pics_next = link2b_soup.select('div#pic > a')
70         if not pics_next:
71             break
72         pics.extend(pics_next)
73         page += 1
74
75 folder_album = '{:03}'.format(album_index) + title #new 2 to 3
76 os.makedirs(folder_album, exist_ok=True)
77
78 album_index += 1
79
80 filename_descriptor = folder_album + '/' + title + '.txt'
81 with open(filename_descriptor, mode='wt', encoding='utf-8') as text_fp_w:
82     # save album name & description in the specific file
83     text_fp_w.write(title)
84     delimit_title = '\n' + ' '*len(title) + '\n'
85     text_fp_w.write(delimit_title)
86     text_fp_w.write(abstract)
87     text_fp_w.write('\n\n')
88
89 pic_index = 1
90 for pic in pics:
91     link3_res = requests.get(DOMAIN + pic['href'], headers=headers)
92     link3_res.raise_for_status()
93     #link_soup = bs4.BeautifulSoup(link_res.content, 'html.parser')#bad
94     link3_soup = bs4.BeautifulSoup(link3_res.content, 'html5lib')
95
96     subtitle = link3_soup.find('b', {'id': 'PTSection'}).text
97
98     valid_name2 = "".join(c for c in subtitle if pattern.match(c))
99     subtitle = valid_name2 (6)
100
101     content = link3_soup.find('div', {'id': 'PDSection'}).text
102     pic_big = link3_soup.find('img', {'id': 'PhotoAreaEns'}) (7)
103     # save photo name & description in the specific file
104     text_fp_w.write(subtitle)
105     delimit_subtitle = '\n' + ' '*len(subtitle) + '\n'
106     text_fp_w.write(delimit_subtitle)
107     text_fp_w.write(content)
108     text_fp_w.write('\n\n')
109
110     image_res = requests.get(DOMAIN + pic_big['src'], headers=headers)
111     #image_res.raise_for_status()
112     # some jpg files seem to get lost
113     if image_res.ok == False:
114         print('failed:', image_res.url)
115         continue
116     image_filename = folder_album + '{:03}' + subtitle + '.jpg'
117     with open(image_filename.format(pic_index), mode='wb') as img_fp_w:
118         for chunk in image_res.iter_content(100000):
119             img_fp_w.write(chunk)
120
121     pic_index += 1
122     print('Got', folder_album)

```

收集該本相簿內的照片網址連結 (一頁一頁收集)

建立該本相簿的目錄

準備檔案並寫入該本相簿的大略文字敘述

根據前面收集的連結，開始請求下載照片了

搜尋單一照片的文字

敘述並寫入檔案

若成功抓到照片，則將其存成檔案。若失敗，則處理下一張

列個訊息說該本相簿處理結束

備份 pchome 相簿程式的執行結果

名稱	修改日期	類型	大小
001_空難紀錄	2020/9/10 上午 09:51	檔案資料夾	
002_空難救援	2020/9/10 上午 09:51	檔案資料夾	
003_深海打撈	2020/9/10 上午 09:51	檔案資料夾	
004_高山深海打撈	2020/9/10 上午 09:51	檔案資料夾	
005_海難百態	2020/9/10 上午 09:51	檔案資料夾	
pchome_album.txt	2020/9/10 上午 09:44	文字文件	1 KB
pchome_album_ReadFile.exe	2020/9/10 上午 09:50	應用程式	12,365 KB

產生 5 本相簿

按這個檔案執行

原則上右面設定除了
第一行網址改成你的
外，其餘不須變動

```
pchome_album.txt - 記事本
檔案(F) 編輯(E) 格式(O) 檢視(V) 說明
http://photo.pchome.com.tw/salvage
-1
1
# 主程式只讀以上五行
# line 1: pchome相簿使用者的網址
# line 2: 起始相本 (第一本的值設0)
# line 3: 截止相本 (不含該相本，只抓到前一本)
# line 4: 跳本 (負值代表從後往前抓，一般後面的日期較舊)
# line 5: 起始相本的編號 (1代表從001開始，下一本依序為002)
# 注意line 2和line 3不填而line 4設-1，代表由最舊往最新日期下載。
```

魔鬼藏在細節裡

- 利用瀏覽器的開發者模切將網頁切換成人容易看的樹狀圖模式，並使用強大的 `beautifulsoup4` 模組將一團純文字的網頁轉換成有組織的物件
 - (1) `links_tmp = link_soup.select('div#alb > a')`
 - (2) `title = link2_soup.select_one('div.tit a').text.strip()`
 - (3) `abstract = link2_soup.select_one('div.ab').text.strip()`
 - (4) `pics = link2_soup.select('div#pic > a')`
 - (5) `subtitle = link3_soup.find('b', {'id':'PTSection'}).text`
 - (6) `content = link3_soup.find('div', {'id':'PDSection'}).text`
 - (7) `pic_big = link3_soup.find('img', {'id':'PhotoAreaEns'})`

salvage

歡迎大家參觀我的相簿

首頁



span | 52 x 16.8

海難奇聞

歷年來在台...
共18張



高山深海打撈

深海打撈
共12張



深海打撈

深海打撈紀錄
共8張



空難救援

空難救援紀錄
共33張



空難紀錄

空難精華剪輯
共16張

熱門指數

累積人氣: 8,812
 推薦數: 0 個
 訪客人數: 1 人
 被訂閱數: 0 人
 相簿總數: 5 本
 最近修改時間: 2011-01-05

相簿分類

Ctrl + Shift + C
Inspect element
(firefox browser)

```

HTML 搜尋 HTML
(1) <div id="left" class="left">
  <div id="main-area" style="position: relative;">
    <!--MyActivity-->
    <!--我的主打-->
    <div id="MainListSection">
      <div id="alb" class="MarkSet">
        <a href="/salvage/05/" target="_parent">
          <span>...</span>
          <div class="name">
            <span>海難奇聞</span>
          </div>
          </a>
          <div class="dis">歷年來在台...</div>
          共18張
        </div>
        <div id="alb" class="MarkSet">...</div>
        <div id="alb" class="MarkSet">...</div>
        <div id="alb" class="MarkSet">...</div>
      </div>
    </div>
  </div>
  </div>
  
```

過濾樣式

body { font-family: "微軟正黑體", "Microsoft JhengHei", Arial, "文泉驛正黑", "WenQuanYi Zen Hei", "麗黑 Pro", "LiHei Pro", sans-serif; }

body { font-family: Verdana, Arial, Helvetica, sans-serif; -webkit-text-size-adjust: none; }

版面

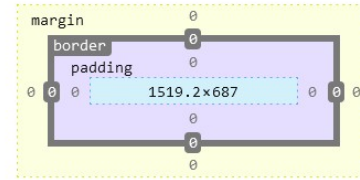
Flexbox

選擇 Flex 容器或項目繼續

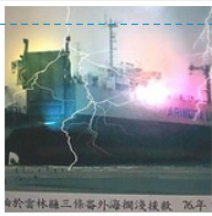
格線

此頁面沒有使用 CSS 格線

Box Model



1519.2x687 static



海難百態 共18張

點閱數：981

歷年來在台灣各附近海域擱淺及沉沒之船舶海難照片

相簿標籤：
全站分類：[拍賣市集](#) 個人分類：[未分類](#)

相簿分類

相片回應

訂閱的相簿 [更多>>](#)

我的收藏 [更多>>](#)



倒插沉沒



擱淺



擱淺



外傘頂洲擱淺



西子灣擱淺

```

<!--相簿資訊 開始-->
<div id="ainfo">
  <div class="cover">...</div>
  <div class="dis">
    <div class="tit"> (2)
      <a href="/salvage/05/">海難百態</a>
      <span class="am">共18張</span>
    </div>
    <div class="t">...</div>
    <div class="ab" style="height:57px;overflow:hidden">...</div>
    <!--隱藏標籤-->
    <div id="hidden_tags" style="position:absolute; width:320px; z-index:999; background:#fef..px; line-heig
    0 50px; visibility: hidden;" onmouseout="HiddenTags()" onmouseover="ShowAllTags()">...</div> <event>
    <!--隱藏標籤-->
    <div class="cb" style="display:block;">...</div>
    <div class="cb" style="display:none;">...</div>
    <div class="cb">...</div>
  </div>
</div>

```

```

元素 {
}
html, body {
  margin: 0;
  padding: 0;
}
body {
  font-family: "微軟正黑體", "Microsoft JhengHei",
  Arial, "文泉驛正黑", "WenQuanYi Zen Hei", "麗黑
  Pro", "LiHei Pro", sans-serif;
}
BODY {
  FONT-SIZE: 12px;
  COLOR: #000;
  FONT-FAMILY: Arial;
  background: #alalal;
  margin: 0px;
}
body {
  margin: 0;
  FONT-FAMILY: Arial;
}

```

版面 計算樣式 變更 字型 動畫

Flexbox
選擇 Flex 容器或項目繼續

格線
此頁面沒有使用 CSS 格線

Box Model

margin: 0
border: 0
padding: 0
content: 1519.2x1150

1519.2x1150 static

ktsalvage

歡迎大家參觀我的相簿

首頁



海難百態 共18張
相簿建立時間: 點閱數: **div.ab** 400 x 57

歷年來在台灣各附近海域擱淺及沉沒之船舶海難照片

- 相簿分類
- 相片回應
- 訂閱的相簿 更多>>

搜尋 HTML

```
<div class="ab" style="height:57px;overflow:hidden">
  歷年來在台灣各附近海域擱淺及沉沒之船舶海難照片
</div>
<!--隱藏標籤-->
<div id="hidden_tags" style="position:absolute; width:320px; z-index:999; background:#fef...px; line-heig
0 50px; visibility: hidden;" onmouseout="HiddenTags()" onmouseover="ShowAllTags()">...</div> [event]
<!--隱藏標籤-->
<div class="cb" style="display:block;">...</div>
<div class="cb" style="display:none;">...</div>
<div class="cb">
  全站分類:
  <a href="/class3.html?books_class=017">拍賣市集</a>
  | 個人分類:
  <a href="/photo_user_class.html?n=salvage&c=">未分類</a>
</div>
```

版面 計算樣式 變更 字型 動畫

元素 { height: 57px; overflow: hidden; }

#ainfo .ab { height: 90px; font-size: 12px; color: #777; }

繼承自 div #ainfo .dis { font-size: 13px; line-height: 20px; color: #2b2b2b; }

繼承自 div .w950 { text-align: left; }

繼承自 body

Flexbox 選擇 Flex 容器或項目繼續

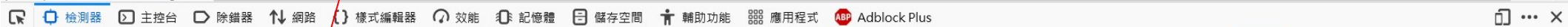
格線 此頁面沒有使用 CSS 格線

Box Model margin: 0, border: 0, padding: 0, 400x57, 0, 0, 0, 0

400x57 static



photo.pchome.com.tw/salvage/106



搜尋 HTML

(4)

```

</div>
</div>
<!--相簿資料 結束-->
<div id="MainListSection">
  <div id="pic" class="MarkSet">
    <a href="/salvage/106">
      <span id="img"/>
        
        
      </span>
      <div class="name">
        <span>倒插沉沒</span>
      </div>
    </a>
  </div>
  <div id="pic" class="MarkSet">

```

濾選樣式

Flexbox

Grid

Box Model

400x57

static


```

<!--照片資料 start-->
<div class="picinfo">
  <div class="jumppic">...</div>
  <h1 style="margin-top:10px">
    <a href="/salvage/05/" title="海難百態">海難百態</a>
    <span class="ar">></span>
    <b id="PTSection" style="font-size:15px;" title="倒插沉沒">倒插沉沒</b>
    <span>(1/18)</span>
  </h1>

```



(7)



```

<div id="PhotoArea">
  

```

日本工作船倒插沉沒

(6)

```

<!--照片資料 start-->
<div class="picinfo">
  <div class="dis">
    <div id="PDSection">日本工作船倒插沉沒</div>
  </div>
</div>
<!--照片資料 end-->

```