

驗證碼 (captcha) 的光學文字辨識

- 實現網路爬蟲過程中，會遇到有些網站需要你輸入驗證碼
 - Completely Automated Public Turing test to tell Computers and Humans Apart，全自動區分電腦和人類的公開圖靈測試
- 使用 tesseract OCR(光學文字辨識) 獲得驗證碼
 - 透過容器：`docker run --rm -v $(pwd):/data -w /data vimagick/tesseract test.png stdout -l eng --psm 13`
 - 在容器內替瀏覽器安裝中文字型 `apt-get install fonts-arphic-ukai fonts-arphic-uming`
 - 直接將工具安裝在作業系統上：`apt-get install tesseract-ocr ; apt-get install imagemagick`
- 測試結果
 - 簡單的4位數字小圖片(國家森林遊樂區網路訂房系統)：**八成以上成功率**

驗證碼 (國家森林遊樂區網路訂房)

回東勢林管處 網站地圖 回首頁 設為首頁 加到我的最愛 遊樂區門票電子票證 訂房Q&A

行政院農業委員會
林務局東勢林區管理處

歡迎使用網路訂房系統，[登入/加入會員](#)

目前系統時間 10時42分14秒
(若長時間未使用本系統，請按重新整理更新時間)

國家森林遊樂區網路訂房系統

輸入欲檢索之字串

訂房流程

線上訂房

LOCATION 首頁 > 訂房流程

登入會員

請輸入帳號： (勿使用身分證字號作為帳號)

舊會員請輸入密碼： [忘記密碼](#)

訂房遊客請於三日內(含當日、例假日)匯訂金，否則將依規取消訂單(外國遊客除外)

驗證碼： 6815

待辨識和輸入的驗證碼

測試樣本

3.692

9.477

3.513

3.262

9.198

9.898

2.571

0.288

0.534

9.587

4.520

1.813

7.705

7.337

0.373

光學文字辨識後的結果

```
001 3692
[]
002 []
003 []
004 []
[]
005 []
[]
006 []
007 []
[]
008 []
009 []
[]
010 []
9587
[]
011 4520
[]
[]
012 []
013 77205
[]
014 []
[]
015 []
```

最初範例的設定



```
001 3692
[]
002 9477
[]
003 3513
[]
004 3262
[]
005 9198
[]
006 9898
[]
007 2571
[]
008 0288
[]
009 0534
[]
010 9587
[]
011 45720
[]
012 '1813
[]
013 77205
[]
014 7332
[]
015 0.3.72'3
[]
```

將 psm 改為 13



```
001 3692
[]
002 9477
[]
003 3513
[]
004 3262
[]
005 9198
[]
006 9898
[]
007 2571
[]
008 0288
[]
009 0534
[]
010 9587
[]
011 45720
[]
012 1813
[]
013 77205
[]
014 7332
[]
015 03723
[]
```

限定只辨識數字

```
3692
9477
3513
3262
9198
9898
2571
0288
0534
9587
4520
1813
7705
7337
0373
```

正解

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
1	3692									9587	4520		77205		
2															
3	3692									9587	4520		77205		
4	3692									9587			77205		03723
5	2 9 46 4 3	7 7 4 9	9 1 5 3	42 6 2 3	8 9 1 9	0 9 8 9	1 7 2	2 8 6 0	4 3 5 0	7 8 0 9	0 2 0 4	3 1 8 1	5 0 7 7	7 0 3 7	3 7 3 3
6	3692	9477	3513	3262	9198	9898	2571	0288	0534	9587	4520	1813	77205	7332	03723
7	3692	9477	3513	3262	9198	9898	2571	0288	0534	9587	4520	1813	77205	7332	03723
8	3692	9477	3513	3262	9198	9898	2571	0288	0534	9587	45720	1813	77205	7332	03723
9	3692	9477	3513	3262	9198	9898	2571	0288	0534	9587	45720	1813	77205	7332	03723
10	3692	9477	3513	3262	9198	9898	2571	0288	0534	9587	4520	1813	77205	7332	03723
11	3692	9477	3513	62	91	9898	2571	0288	0534	9587	452	1813	77205	733	03723
12	3692	9477	3513	62	91	9898	71	0288	05	9587	452	1813	77205	733	0373
13	3692	9477	3513	3262	9198	9898	2571	0288	0534	9587	45720	1813	77205	7332	03723
14	3692	9477	3513	3262	9198	9898	2571	0288	0534	9587	4520	1813	7705	7337	0373

對 15 個樣本，藉由改變 psm，操作在 13 種不同模式下，最後共得到 195 辨識結果，其中以**模式 6 和 7 效果最佳**，模式 8,9,10 和 13 次之，而**模式 2 對數字辨識明顯不行**。

A ~ O : 15 samples
 1 ~ 13: 13 modes
 14 : 正解

← 正確數字 (row # 14)

tesseract 簡介

https://zh.wikipedia.org/wiki/Tesseract



Tesseract [編輯]

維基百科，自由的百科全書

A → 文

此條目可參照英語維基百科相應條目來擴充。*(2018年7月24日)*

若您熟悉來源語言和主題，請協助參考外語維基百科擴充條目。請勿直接提交機械翻譯，也不要翻譯不可靠、低品質內容。依版權協議，譯文需在編輯摘要註明來源，或於討論頁頂部標記 {{Translated page}} 標籤。

Tesseract是一個光學字元辨識引擎，支援多種作業系統。^[1]Tesseract是基於Apache許可證的自由軟體^[2]，自2006年起由Google贊助開發^[3]。

2006年，Tesseract被認為是最精準的開源光學字元辨識引擎之一。^{[2][4]}

透過容器來執行 tesseract

https://registry.hub.docker.com/r/vimagick/tesseract/

100%

tesseract

Tesseract is an Open Source OCR engine, available under the Apache 2.0 license. It can be used directly, or (for programmers) using an API. It supports a wide variety of languages.

Tesseract doesn't have a built-in GUI, but there are several available from the 3rdParty page.

Quick Start

```
$ alias tesseract='docker run --rm -v `pwd`: /data -w /data vimagick/tesseract'
```

```
$ tesseract input.png output -l eng --psm 3
```

```
$ cat output.txt
```

```
The (quick) [brown] {fox} jumps!
```

```
$ tesseract chinese.jpg stdout -l chi_tra --psm 8 --oem 0
```

```
學習
```

Docker Pull Command

```
docker pull vimagick/tesseract
```

Owner

感謝該作者的啟發



vimagick

Source Repository



GitHub

vimagick/dockerfiles

`docker run --rm -v $(pwd):/data -w /data vimagick/tesseract input.png output -l eng --psm 3`

測試程式 (1)

```
1 import os
2 #import subprocess
3
4 #docker_cmd1="docker run --rm -v $(pwd):/data -w /data vimagick/tesseract \
5 #test_{:03}.png stdout -l eng --psm {}"
6
7 #docker_cmd2="docker run --rm -v $(pwd):/data -w /data vimagick/tesseract \
8 #test_{:03}.png stdout --oem 0 -c tessedit_char_whitelist=0123456789 --psm {}"
9
10 docker_cmd3="docker run --rm -v $(pwd):/data -w /data vimagick/tesseract \
11 test_{:03}.png stdout -c tessedit_char_whitelist=0123456789 --psm {}"
12
13 mode = 13# 3, '1st':6,10 '2nd':7,8,9,13 '3rd':11,12
14
15 for number in range(15):
16
17     ocr_result = os.popen(docker_cmd3.format(number+1, mode)).read()
18
19     print('{:03}'.format(number+1), ocr_result)
```

測試程式 (2)

```
1 import os
2 import csv
3 import time
4
5 docker_cmd3="docker run --rm -v $(pwd):/data -w /data vimagick/tesseract \
6 test_{:03}.png stdout -c tesseract_char_whitelist=0123456789 --psm {}"
7
8 tmp = str(time.time_ns()) # used as unique number
9 mode_max = 14
10 sample_max = 16
11
12 for mode in range(1, mode_max):
13     psm = dict()
14
15     for sample in range(1, sample_max):
16
17         ocr_result = os.popen(docker_cmd3.format(sample, mode)).read()
18
19         psm_key = '{:03}'.format(sample)
20         print(psm_key, ocr_result)
21         psm[psm_key] = ocr_result
22
23     fieldnames = ['{:03}'.format(sample) for sample in range(1, sample_max)]
24
25     with open('ocr_' + tmp + '.csv', 'a', newline='', encoding='utf-8') as f_w:
26         writer = csv.DictWriter(f_w, fieldnames=fieldnames)
27         writer.writerow(psm)
28
```