

2019

嚴選台北地區五星飯店之評論

高子淇、賴永琪、劉萱

目录 Contents

01



動機與目的



動機與目的

01 動機

隨著時間累積大量的顧客評論也成為重要的研究資料，利用文字探勘的方式，找出有用的關鍵字來推薦使用者最佳訂房選擇，以縮短使用者訂房時間。

02 介紹

由於技術上的困難，在嘗試過後，選擇以評分與評論之間的關係做成分類。

03 目的

了解消費者對於旅館的評論與評分之間的關聯，進而能使網站平台，飯店了解消費者的需求，最終能使消費者更佳便利的瀏覽訂房網站。



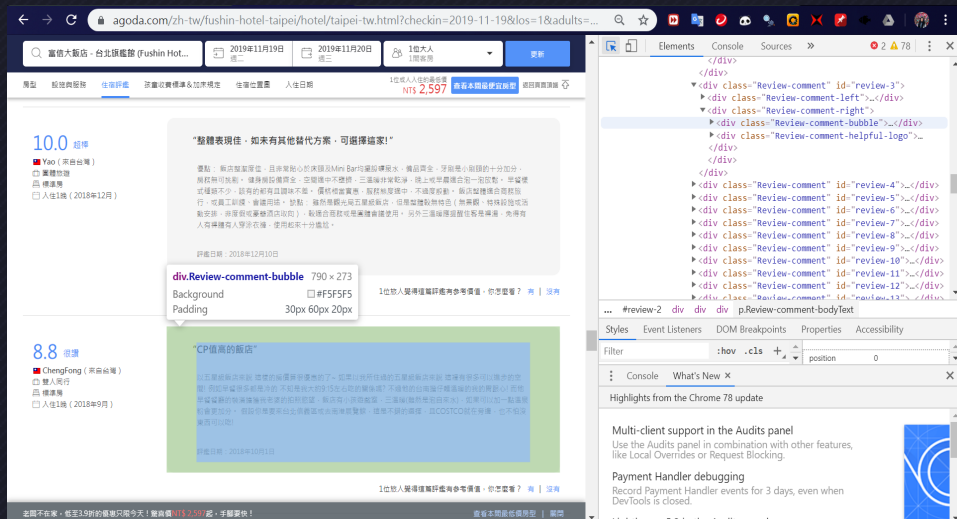
02



資料收集



使用(<https://www.agoda.com/zh-tw/>)提取台北地區五星級飯店的顧客評鑑





資料收集

02 使用方法：selenium

直接運行在瀏覽器中，通過一系列命令來模擬操作，可以將這些命令轉化成實際的HTTP請求在瀏覽器中運行

```
In [1]: import requests
        from bs4 import BeautifulSoup
        import selenium
        from selenium import webdriver
        from selenium.common.exceptions import NoSuchElementException
        from selenium.webdriver.common.keys import Keys
        import time
        from bs4 import BeautifulSoup
        from selenium import webdriver
        from selenium.webdriver.support.ui import WebDriverWait
        from selenium.webdriver.common.by import By
        from selenium.webdriver.support import expected_conditions as ec
        import pandas as pd
```



資料來源

03 爬下飯店評論及評分將 兩項製程資料表做分類分析

[illegible]

03



模型與分析步驟



分析步驟

Agoda | 台北市住宿 | 價格保證

hkg.agoda.com/zh-tw/search?asq=4ol%2FMUcWtaTA3fKh8%2F12VysLs8HBOYxs%2BeCgsfL...

台北市
47間住宿還有空房

2020年1月7日
週二

2020年1月8日
週三

1位大人
1間客房

搜出好價

篩選條件
人氣選項
價格
5 星級評等
地點
顯示全部

想搜什麼就搜什麼

飯店+民宿
飯店
家庭旅遊專區
民宿
機票

家庭旅遊模式
關閉

地圖找房
查看台北市地圖
顯示附免費早餐的住宿

台北市推薦住宿

恭喜！您今日獲得了加碼5%的折扣優惠。
只要啟用折扣碼即可解鎖優惠。
啟用折扣碼

排序方式
最契合優先
最低價優先
距離
住客評鑑優等
A級機密價

這些是我們特別為你精選的住宿選項。

免費取消
高CP值
富信大飯店 - 台北旗艦館 (Fushin Hotel-)
2,439 份評論
很讚
8.5

台北市

五星級

飯店



技術說明

```
browser = webdriver.Firefox()
browser.get("https://www.agoda.com/zh-tw/search?city=4951&languageId=20&userId=fbea6c98-bc10-4d2c-a45f-38aa37091ae0&sessionId=qi
#while len(soup.select(".btn_pagination2_next"))>0:
browser.execute_script("window.scrollTo(0,336)")
browser.execute_script("window.scrollTo(0,354.6666564941406)")
browser.execute_script("window.scrollTo(0,469.3333435058594)")
browser.execute_script("window.scrollTo(0,696.6666870117188)")
browser.execute_script("window.scrollTo(0,813.3333129882812)")
browser.execute_script("window.scrollTo(0,928.6666870117188)")
browser.execute_script("window.scrollTo(0,1046.6666259765625)")
browser.execute_script("window.scrollTo(0,1160.6666259765625)")
browser.execute_script("window.scrollTo(0,1277.3333740234375)")
browser.execute_script("window.scrollTo(0,1396.6666259765625)")
browser.execute_script("window.scrollTo(0,1508.6666259765625)")
browser.execute_script("window.scrollTo(0,1814)")
browser.execute_script("window.scrollTo(0,1858)")
browser.execute_script("window.scrollTo(0,2134.666748046875)")
browser.execute_script("window.scrollTo(0,2204.666748046875)")
browser.execute_script("window.scrollTo(0,2440.666748046875)")
browser.execute_script("window.scrollTo(0,2553.333251953125)")
browser.execute_script("window.scrollTo(0,2985.333251953125)")
browser.execute_script("window.scrollTo(0,3047.333251953125)")
browser.execute_script("window.scrollTo(0,3364)")
browser.execute_script("window.scrollTo(0,3667.333251953125)")
browser.execute_script("window.scrollTo(0,3831.333251953125)")
browser.execute_script("window.scrollTo(0,4114.66650390625)")
browser.execute_script("window.scrollTo(0,4524)")
browser.execute_script("window.scrollTo(0,4576.66650390625)")
browser.execute_script("window.scrollTo(0,5190.66650390625)")
browser.execute_script("window.scrollTo(0,5340.66650390625)")
browser.execute_script("window.scrollTo(0,5727.33349609375)")
```



技術説明

```
browser.execute_script("window.scrollTo(0,14500)")
soup = BeautifulSoup(browser.page_source)
http=[]
for ele in soup.select(".PropertyCardItem.ssr-search-result a"):
    http.append("https://www.agoda.com"+ele.get("href"))
browser.close()
comment=[]
score=[]
for i in http:
    url=i
    browser = webdriver.Firefox()
    browser.get(url)
    browser.execute_script("window.scrollTo(0, document.body.scrollHeight)")
    time.sleep(5)
    if browser.find_element_by_css_selector("span.Searchbox__searchButton__text"):
        browser.find_element_by_css_selector("span.Searchbox__searchButton__text").click()
        soup = BeautifulSoup(browser.page_source)

        for ele in soup.select(".Review-comment"):
            for ele1 in ele.select(".Review-comment-leftScore"):

                for ele2 in ele.select(".Review-comment-bodyText"):
                    if len(ele1.text)*len(ele2.text)!=0:

                        score.append(ele1.text)
                        comment.append(ele2.text)

browser.close()
```



模型與分析步驟



用評分將評論分成三類

大於 9 為一類，在 7、8 區間呈另一類，不足者即為最後一類



分類完後，將分成訓練集與測試集



利用課堂上的分類方法



模型與分析步驟

01 評論與評分之高低關係

02 問題出現

03 模型訓練集有.95準確率

04 模型測試集有.55的準確率

evaluation on train data

confusion matrix:

```
[[242 10 9]
 [ 0 188 0]
 [ 0 0 51]]
```

accuracy:

0.962

evaluation on test data

confusion matrix:

```
[[230 113 42]
 [ 35 18 4]
 [ 1 2 1]]
```

accuracy:

0.5582959641255605

04



實驗與分析結果



實驗

01 過程中有使用jieba進行斷詞

02 然而，有過多不必要的單字使真正提及優點反而被掩蓋。

	support	itemsets
0	0.65	(\n)
1	0.50	(也)
2	0.55	(很)
3	0.50	(早餐)
4	0.60	(是)
5	0.65	(有)
6	0.90	(的)
7	0.55	(都)
8	0.65	(飯店)
9	0.50	(,)
10	0.55	(有, \n)
11	0.60	(的, \n)
12	0.50	(\n, 飯店)
13	0.50	(的, 也)
14	0.55	(的, 很)
15	0.60	(是, 的)
16	0.60	(的, 有)
17	0.55	(都, 的)
18	0.65	(的, 飯店)
19	0.50	(的, 有, \n)
20	0.50	(的, \n, 飯店)



實驗困難



1.selenium動態網站自動爬取下一頁，只能到第二頁，第三頁不能再繼續

2.訂房網站會出現沒有飯店名稱、詳細內容的推薦，解決方案if沒有詳細內容（程式碼）就跳過。



3.爬下來的comment 包含許多unicode string



模型與分析步驟

01 更多資料時測試集的準確率下降

02 問題又出現

03 模型訓練集有.95準確率

04 模型測試集有.48的準確率

evaluation on train data

confusion matrix:

```
[[370 14 12]
```

```
[ 1 242  2]
```

```
[ 0  0 59]]
```

accuracy:

0.9585714285714285

evaluation on test data

confusion matrix:

```
[[496 391 134]
```

```
[ 48  62  16]
```

```
[ 0  0  0]]
```

accuracy:

0.4864864864864865



分析結果

只有近五成的準確率，可以再更探討更多因素，可以再了解情緒字詞。

目录 Contents

05



結論



結論

單純只看評論與評分的分類可能準確度較低，還需搭配情緒字詞或是個人對於分數的滿意程度上感覺的不同來深入了解

文字探勘專題計畫
書

2019

感謝一路有你