

政治學的資料探勘與機器學習期末報告

以主題塑模看不同論壇對香港反送中議題的立場觀點
-以 **Dcard**、**PTT** 為例

組員：胡慶右、陳禹瑄、曾筑翎、曾淳榆

壹、摘要

香港反送中議題自 2019 年六月以來，一直是我國關注度於討論度皆高的議題，因此本組以主題塑模分析 Dcard 與 PTT 的使用者不同，是否造成對香港反送中議題輿情態度的不同，並探討軟分群（LDA）與硬分群（K-means）在網路論壇輿情分析中的應用比較。研究發現 Dcard 和 PTT 使用者雖然都對香港反送中抱持支持態度，然而 Dcard 使用者的態度較為同理；PTT 使用者則涵蓋小部分的偏激群眾。此外，研究發現由於論壇的發言為一來一往的討論，相較一般文章論壇發言缺乏完整性與脈絡，因此分析論壇輿情時較適用軟分群，另外，由於論壇充滿了顏表情與各種負面語助詞，因此對於論壇的斷詞分析需要謹慎考量停用詞。

貳、研究動機與目的

香港於 2019 年 3 月 15 日開始發起反對《逃犯條例修訂草案》運動（簡稱反送中），並在 6 月 9 日起爆發大規模的社會運動。此議題在當時引起廣泛討論，而台灣人民由於擔心香港事件在台灣上演的「亡國感」對反送中議題更是關注，並且由於鄰近台灣總統大選，香港反送中議題在各大社群媒體上有十分高的關注度及討論度，然而傳統民調甚少直接調查民眾對於香港反送中議題之態度。今周刊〈以台灣人的角度看反送中對台灣的影響〉一文中提到「民調顯示，超過七成的台灣民眾支持香港示威運動，以年輕族群居多」，卻未清楚說明民調調查方式與資料來源，因此本組欲以主題塑模之方式分析輿情，並選擇 Dcard 與 PTT 作為研究對象。本組選擇此兩論壇作為研究對象的原因為，國家發展委員會於 2018 年委託聯合行銷研究股份有限公司進行的〈107 年公民網路參與行為調查報告〉發現，PTT 與 Dcard 的使用者相較其他論壇較集中，且多為青壯年人口（如表一），此外，PTT 與 Dcard 也被認為是最多大學生使用的社團，然而此兩個論壇在使用者年齡趨勢上仍有一些差異（PTT 主要分布在 20-29、30-39 兩群；Dcard

則主要分布在 20-29)，加上兩個論壇在帳號申請上有不同限制，因此本組認為雖然 PTT 與 Dcard 的使用者皆為青壯年，但主要使用者存在差異，因此本組想探討 Dcard 與 PTT 的使用者不同，是否造成對於香港反送中議題的輿情態度的不同，此外，本組試圖於研究中比較軟分群（LDA）與硬分群（K-means）在網路論壇輿情分析中的應用差異。

透過分析不同論壇的輿情，希望了解對於重大社會議題的態度，是否會因為使用者不同，而造成不同的輿論。此外，有別於傳統民調已發展許久，輿情分析是屬於較新穎的研究方式，因此透過比較軟分群與硬分群在輿情分析上的應用，希望提供未來發展輿情分析之參考。

表一、各論壇使用比例

單位：人次/百人

	LINE	FB	IG	Google ⁺	Twitter	Ptt	Dcard
12-14 歲	34.9	59.2	25.9	0.0	0.0	0.0	0.0
15-19 歲	34.7	62.6	19.4	6.7	7.1	3.4	2.9
20-29 歲	36.8	59.5	14.3	3.5	0.8	14.6	7.0
30-39 歲	30.7	58.9	4.1	4.8	1.9	10.5	1.2
40-49 歲	40.6	53.2	6.0	7.0	3.4	6.3	2.9
50-59 歲	42.8	47.7	6.5	8.7	2.6	4.7	0.5
60-64 歲	52.8	46.8	2.1	11.8	2.3	3.7	0.0
65 歲以上	65.1	39.0	1.8	15.4	0.9	1.9	0.9

資料來源：國家發展委員會

參、資料說明

一、開發環境

選擇操作簡單並且以 `python` 語言為主，同時可以下載多種套件的 Jupyter Notebook。

二、 安裝套件

Ptt 安裝套件為：`pandas`(資料視覺化應用套件), `selenium`(操控瀏覽器來擷取 HTML), `time`(管理、控制時間套件)；Dcard 安裝套件為：`Requests`(網路資源 URLs 擷取套件), `json`(將 JSON 物件轉為 Python 資料類型), `pandas`(資料視覺化應用套件), `re`(正則表示式套件)。

三、 爬取方式

有鑑於 Ptt 之網頁格式屬於 HTML 與 CSS 形式，故使用網頁爬蟲，以搜尋引擎爬蟲為主，根據網頁上的超連結進行遍歷爬取；而有鑑於 Dcard 提供網頁之 API，故使用介面爬蟲，通過精準構造特定 API 介面的請求資料，而獲得大量資料資訊。

四、 實際操作

實際爬取 Ptt Gossiping 版、Ptt HatePolitics 版中以「反送」、「香港」為關鍵字的所有討論內容，總共在 Ptt Gossiping 蒐集到 15008 篇文章，在 Ptt HatePolitics 蒐集到 621 篇文章。以下以 Ptt Gossiping 版的爬取過程作為範例說明。

(一) Ptt

設定網址路徑到 Ptt Gossiping 版，其中因為 Ptt 有網站內容分級規定處理(警告：您即將進入之看板內容需滿十八歲方可瀏覽)，因此需先定點取滿 18 歲之選項，該選項寫於 `.over18-button-container.btn-big` 的 class 中，再將「反送」兩字放入搜尋欄中，始完成初始化路徑設定。

```

1 import pandas as pd
2 from selenium import webdriver
3 from selenium.webdriver.common.keys import Keys
4 import time
5
6 driver_path = "c:/Users/Sherry/Desktop/geckodriver.exe"
7 road_home = "https://www.ptt.cc/bbs/Gossiping/index.html"
8 driver = webdriver.Firefox(executable_path=driver_path) # Use Firefox
9 driver.get(road_home)
10 element=driver.find_elements_by_css_selector('.over18-button-container .btn-big')
11 element[0].click()
12 element = driver.find_element_by_css_selector('.query') # select searcher
13 element.send_keys('反送', Keys.ENTER)
14 time.sleep(5)
15 elements = driver.find_elements_by_css_selector('.title a')
16 dates = driver.find_elements_by_css_selector('.date')

```

搜尋結果為 94 頁之資料內頁，因此設定迴圈範圍為 94，選取每個頁面中文章標題和時間，分別將其存入 title 列表和 date 列表，click()進入內容頁，選取所有網民留言，此存在.push-content 的 class 中，回到上一頁重複該動作直到無資料為止，跳到下一頁。

```

18 title=[]
19 comment_all=[]
20 comment=[]
21 date=[]
22 url=[]
23
24 for k in range(94):
25     for i in range(len(elements)):
26         title.append(elements[i].text)
27         date.append(dates[i].text)
28         elements[i].click()
29         comments=driver.find_elements_by_css_selector('.push-content')
30         comment=[]
31         for j in range(len(comments)):
32             comment.append(comments[j].text)
33         comment_all.append(comment)
34         driver.back()
35         elements = driver.find_elements_by_css_selector('.title a')
36         dates = driver.find_elements_by_css_selector('.date')
37         next_page=driver.find_element_by_css_selector('.wide:nth-child(2)') # next page
38         next_page.click()
39         elements = driver.find_elements_by_css_selector('.title a')
40         dates = driver.find_elements_by_css_selector('.date')
41
42 driver.close()

```

利用 pandas 套件將資料整理並且存入 csv 檔中。

```

1 dict={'title':title,'date':date,'content':comment_all}
2 df=pd.DataFrame(dict)

```

1	df				
		title	date		content
0		[新聞] 匯集反送中意象作品 港青展出苦水穿石抗	6/20		[我怎覺得中共是水的角色，欸，解放軍黑警是水，港女是石...; 推]
1		[新聞] 反送中週年 港人悼念首位預命抗爭者	6/16		[維尼：挖鼻，：推，：反送懶叫港臭飛去甲囊，：還是會記得一年前那個畫面，：...
2		[新聞] 疫情、反送中成敏感詞？社群媒體被質疑	6/16		[反正臺灣派賤畜還是爽用 哈哈，：可悲綠燦大內宣，：油管師啊師 誰叫RMB真香]
3		Re: [新聞] 反送中周年 台北7千人集會撐港	6/14		[：你也能去，廢話一堆，：7千人誇張耶，有什麼喝的吃的嗎，：柯韓翼大動員也才70...
4		[新聞] 反送中周年 台北7千人集會撐港	6/14		[：下禮拜五，：請洽香港版，：天滅中共，：意淫香港，：一群白癡，：柯韓...
...	
1858		Re: [新聞] 50萬港人今上街反送中惡法	6/09		[：有一天我夢到香港沒有黑社會，：人暴動 坦克開出來 香港發大財，：韓粉要...
1859		[新聞] 50萬港人今上街反送中惡法	6/09		[：港狗死好，：推，：沒有含種這勢多，：比韓粉少了一點，：很多人應該是覺得...
1860		[爆卦] 香港反送中遊行現場	6/09		[：擠爆，：香港沒救了 被綁進中國就只能等死，：五樓送終，：來我們後山花東打拼...
1861		[爆卦] 反送中遊行提前開始了	6/09		[：：坦克待命中，：柯文哲，：台灣都救不了了還管香港，：給他們看看昨天花...
1862		[問卦] 有沒有反送中大遊行的卦？	6/09		[他國事務，：台北也又，：八卦是中共華南軍區的兵力已經部署在深圳了，：干我屁...

(二) Dcard

實際爬取 Dcard entertainer 版、Dcard hkmacdaily 版、Dcard trending 版、Dcard mood 版、Dcard hktrending 版中在標題或內文有「反送」、「香港」兩關鍵字的所有討論內容以及內文，總共在 Dcard trending 蒐集到 4797 篇文章，在 Dcard entertainer 蒐集到 7 篇文章，Dcard mood 版 577 篇文章，Dcard hkmacdaily 版 68 篇文章，Dcard hktrending 蒐集到 13 篇文章。以下以 Dcard trending 版的爬取過程作為範例說明。

設定網址路徑到 Dcard 提供的 API，即可以次以 JSON 的型態瀏覽最多一百筆資料內容，本組選取 id(文章 id)、title(文章標題)、excerpt(文章內文)、createdAt(發文時間)進行 10000 次迴圈，如無資料即跳出迴圈。以下以 Dcard trending 版的爬取過程為範例說明。

```
1 # Dcard-Trending-香港
2 import requests, json
3 import pandas as pd
4 import re
5
6 article_id=[]
7 title_i_want=[]
8 content_for_each=[]
9 date=[]
10 article_all_id=[]
11 excerpt=[]
12
13 for j in range(1000):
14     try:
15         last=str(article_all_id[-1])
16     except:
17         last=str(233914588)
18     url="https://www.dcard.tw/_api/forums/trending/posts?popular=false&limit=100&before=" + last
19     requ = requests.get(url)
20     rejs=requ.json()
21     for i in range(len(rejs)):
22         article_all_id.append(rejs[i]['id'])
23         if ('香港' in rejs[i]['title']) or ('香港' in rejs[i]['excerpt']):
24             title_i_want.append(rejs[i]['title'])
25             date.append(rejs[i]['createdAt'])
26             article_id.append(rejs[i]['id'])
27             excerpt.append(rejs[i]['excerpt'])
```

根據上面步驟得到含有關鍵字之文章 id，放入 API 網址中，前往文章內容，爬取文章的內文與網民留言，分別為 excerpt 和 content，最後利用 pandas 將所有資料進行整理視覺化。

```
28 for id in article_id:
29     url_content="https://www.dcard.tw/_api/posts/{id}/comments".format(id=id)
30     requ = requests.get(url_content)
31     rejs = requ.json()
32     content=[]
33     for i in range(len(rejs)):
34         try:
35             content.append(rejs[i]['content'])
36         except:
37             continue
38     content_for_each.append(content)
39
40 dict={'article':article_id,'createdAt':date,'title':title_i_want,'excerpt':excerpt,'content':content_for_each}
41 df_trending=pd.DataFrame(dict)
42 df_trending
43
```

利用 pandas 套件將資料整理並且存入 csv 檔中。

	article	createdAt	title	excerpt	content
	0	233913443	2020-06-20T14:51:25.400Z	《今日香港，明日台灣，誰說的底話》（文長）	我叫b仔，是一名應屆中學文憑試的考生。原本應該在學校無憂無慮地上課，與同學有空打一下LOL，...
	1	233909230	2020-06-19T23:37:27.630Z	大同涉嫌法中資？黃國昌：最好的解藥就是脫權透明化	從這集法律白話文楊貴智的podcast 法客電台中，黃國昌從他競選大同獨立董事的事件，想帶領...
	2	233907618	2020-06-19T15:25:03.474Z	【公子快報】七國集團聯合反對香港國安法	【公子快報】七國集團聯合反對香港國安法，薩德會談了7個小時只剩一個共識，中國一意孤行與西方世...
	3	233902816	2020-06-18T23:57:03.503Z	協助港人 即刻啟動	台灣援港專案，即刻啟動！，今天陸委會正式公布 #香港人道援助關懷行動專案，未來對於香港朋友的...
	4	233900242	2020-06-18T12:58:38.222Z	為什麼要抗議釣魚台被占領 不抗議香港、大陸被占領	我一直覺得很奇怪，最近日本宣布釣魚台島改名 宣示主權，一群人喊著要向日本抗議，可是釣魚台就這...
...
	4440	31276	2014-09-01T13:39:08.533Z	【新聞】黃國昌：香港人夢碎、台灣人夢醒了嗎？	中共人大常委會在八月的最後一天，正式敲碎了港人在2017年依民主原則普選行政長官的期望，依照...
	4441	22948	2014-06-19T11:41:24.747Z	【新聞】臉書掛貼前 中國3波數位「導彈」襲美	http://imgur.com/gabOHuVn/n社群網站臉書(facebook) 下午...
	4442	22338	2014-06-15T05:17:50.535Z	【新聞】林飛帆臉書聲援港人：差點掉下淚來	【新唐人亞太台2014年6月14日訊】針對昨天(13日)香港立法會爆發衝突，示威學生遭到...
	4443	21364	2014-06-06T11:21:00.224Z	香港人也攻進立法院了	反對新界東北前期撥款的東北村民，部份成功衝入立法會大堂圍困靜坐，要求強制派議員出來交代；另有...
	4444	20265	2014-05-13T04:45:37.277Z	香港洗碗工月薪16000港元...	16 000港幣 =n62 257.7746 台幣/nhttp://imgur.com...

4445 rows x 5 columns

肆、資料分析結果

以下分為兩小節分別說明 Dcard 與 PTT 使用軟硬分群的分析結果。

一、Dcard 資料分析說明

（一）硬分群處理

在 Dcard 資料處理，首先以結巴(Jieba-JS)斷詞分析器進行字典法(Dictionary)斷詞。並分別以詞袋模型(Bag of Words)與詞頻模型(TF)進行文字像量化，後，使用 Weka 中的 Add Cluster 套件進行分群，分群演算法則選用層疊式 K 平均法(Cascade Simple K Means)，設定 3 至 7 群，由演算法選出最適群數。

（二）硬分群資料分析結果

使用詞袋模型進行向量化後的分群結果呈現 3 群，分別為第一群的 3 筆、第二群的 319 筆與第三群的 403 筆。依此分群結果所製作的文字雲如下：

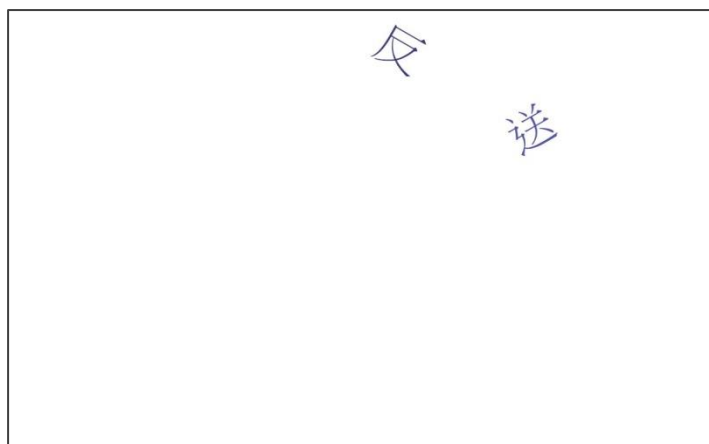


圖 1 第一群



圖 2 第二群

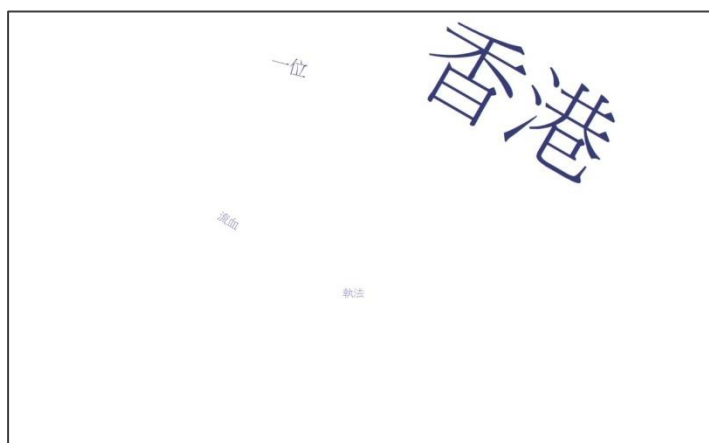


圖 3 第三群

使用詞頻模型進行向量化後的分群結果呈現 6 群，分別為第一群的 50 筆、

第二群的 1 筆、第三群的 3 筆、第四群的 669 筆、第五群的 1 筆與第六群的 1 筆。依此分群結果所製作的文字雲如下：



圖 4 第一群



圖 5 第二群



圖 6 第三群

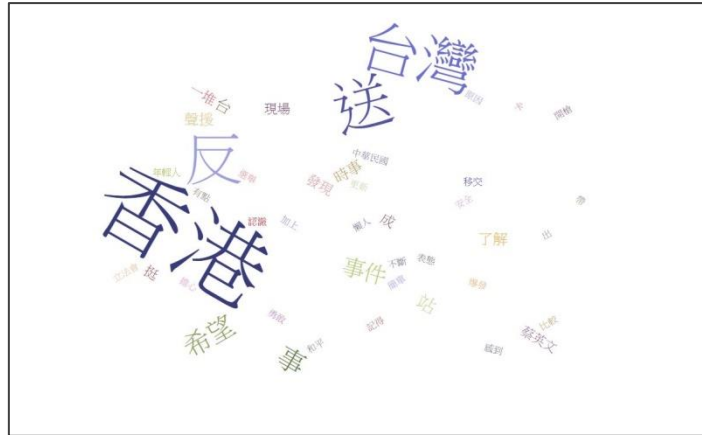


圖 7 第四群

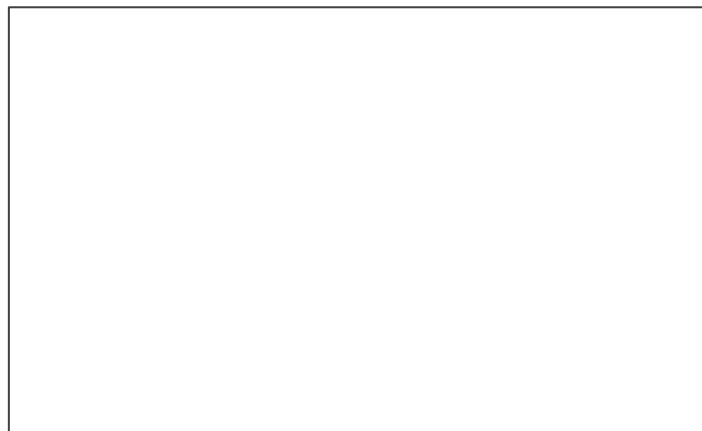


圖 8 第五群



圖 9 第六群

從上述兩節的分群結果來看，其解釋效果並不理想。無論是以詞袋或詞頻模型的結果來進行分群，都出現了多數資料集中於其中一、兩個分群的結果，無法促進對 Dcard 上輿論的理解，故本文於後續接著嘗試軟分群的分群方式。

(三) 軟分群處理

以結巴（Jieba-JS）斷詞分析器進行字典法（Dictionary）斷詞。本部分無須進行文字的向量化。可直接執行主題塑模（Topic Modeling）分析。於群數的設定，不同於硬分群由演算法找出最適群數，主題塑模須手動設定，故本文依序由3至6群進行嘗試。後續選出解釋效果較佳的群數進行分析。

(四) 軟分群資料分析結果

依 3 至 6 群分別執行主題塑模後，本文認為分為 3 群之結果較佳，分別為第一群的 172 筆、第二群的 130 筆與第三群的 420 筆。依此分群結果所製作的文字雲如下：



圖 10 第一群



圖 11 第二群

其比例差異進行分群，而呈現較佳的解釋效果。

二、PTT 資料分析

PTT 資料分析由於資料量過大，現有設備尚無法獲致結果。硬分權方式分為兩種：一、2020 年 6 月到 2019 年 6 月的所有觀察樣本；二、2020 年 6 月到 2019 年 6 月樣本中，單一樣本回應數超過五十則以上。以下分別說明。

（一）所有觀察樣本分析結果

以關鍵字「反送」所有樣本作為分析對象，以結巴（Jieba-JS）斷詞分析器進行字典法（Dictionary）斷詞。以「詞袋」模式作為向量分析，使用 Weka 中的 Add Cluster 套件進行分群，分群演算法則選用層疊式 K 平均法（Cascade Simple K Means），設定 3 至 7 群，由演算法選出最適群數。

其結果如下：第一群 1377 筆、第二群 96 筆、第三群 419 筆，僅第三群形成文字雲（419 筆），其餘兩群文字雲空白，也就是說僅 22% 的資料能形成「廣泛一致」的內容。



圖 13 PPT 觀察樣本第三群文字雲

由於資料過度「廣泛」，無明顯主題性存在，其原因如下：一、PTT 資料過度零碎，前後言無邏輯慣性；二、僅三兩篇回應的文章也納入分析對象；三、比之於 Dcard 回文包含較多情緒用詞或顏文字（例：■ ㄈ \$@ 齣 屎 哈）。基於此上理由進行資料分析修改：首先，以超過五十則回應的單筆資料作為樣本，此原因有二：一、輿論不僅是單一情緒性發言，而是長久討論結果，充分討論的內容，更能代表「意見」；二、每十筆資料約略兩筆有 50 則討論，將得到 300-500 筆資料，內容不致過少；第二個改善方式為增加停用詞——■、\$、@、齣、屎、哈，；第三，刪除英文，因英文多代表連結網址或圖片網址。

（二）五十則回應的樣本

經篩選後資料餘 507 筆，且僅七筆資料發生於 2020 年，也將使得分析結果與 Dcard 分析結果更具可比性。第一群 306 筆（圖 14），第二群 51 筆（圖 15），第三群 150 筆（圖 16），文字雲分別如下：



圖 14 第一分群

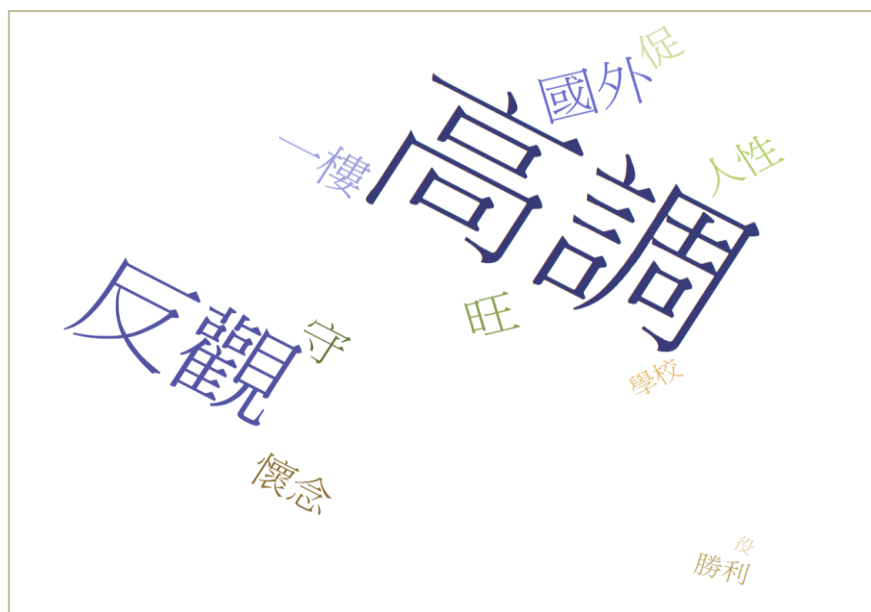


圖 15 第二分群



圖 16 第三分群

分群	第1群 (306) 下載	第2群 (51) 下載	第3群 (150) 下載
屬性為分群中最大	筆數	0: ^ 1: 2: 一下 3: 一中 4: 一個	542: 梁家輝 559: 歌 587: 沒 984: 骨氣 987: 高調
屬性大於全部資料均值	123: 公民 210: 商人 994: 黎	11: 一樓 72: 人性 92: 促 154: 勝利 178: 反觀	57: 之前 58: 之後 59: 乖乖 60: 九二 61: 乾
屬性小於全部資料均值	278: 姑 332: 帥 512: 暴民 564: 歲 607: 消失	253: 夕 542: 梁家輝	筆數 22: 下來 24: 下台 28: 不同 49: 中央
屬性為分群中最小	0: ^ 1: 2: 一下 3: 一中 4: 一個	筆數 278: 姑 332: 帥 559: 歌 984: 骨氣	123: 公民 210: 商人 512: 暴民 564: 歲 607: 消失

圖 17 分析數據

(三) 小結

由以上資料可認為三群分為三種態度：第一群為正義情緒，反送中對此群而言更易聯想到公民權力；第二群為平常態度，也就是容易視反送中為單一事件而已；第三群為偏激態度，其關鍵字包含更多政治情緒——反共、聯想韓國瑜或柯文哲，且「公民」詞彙為分群中較小的屬性，也就是說對第三群而言，他們更帶有偏激政治性的情緒。此外，額外發現為梁家輝在反送中的精神意義，其為第三群屬性最大的資料，此代表網民對於藝人的政治性發表很是在意。

伍、結論

從上述 Dcard 和 PTT 的回文分析發現 Dcard 與 PTT 的使用者確實存在差

異，對於香港反送中議題的輿情態度不同，Dcard 和 PTT 的使用者雖然都對香港反送中抱持著支持態度，然而，Dcard 的使用者抱持著對於該事件較為同理的態度，如較多「聲援」、「關心」等字詞，並且相較於 PTT 也較少激進用詞；PTT 的使用者則是涵蓋小部分的偏激群眾。本組比較軟分群（LDA）與硬分群（K-means）在網路論壇輿情的差異，發現硬分群較不適用於論壇性質的文字，論壇性質的網民討論網站較適用於軟分群（LDA）的方法，因為其中的文章回文就如同人與人之間的口語，表示方式也因此較為支離破碎，因此本組發現，在進行論壇網站的斷詞分析時，需要謹慎考量停用詞，才能更準確切割該名使用者真正想要表達的含意。最後，本組提出一個反思，在網路媒體當道的時代，以往的問卷、民調似乎無法滿足現今大數據能提供的資訊，因此輿情分析成為新的趨勢，然而以論壇做為民調對象是否適當呢？本組透過此研究結果認為，論壇可能僅反映特定的強烈情緒，且可能僅代表某部分群眾之意見，因此分析結果是否能推論母體令人存疑。然而透過此研究本組期望了解不同群體對於社會重大議題是否具有不同輿情態度，並嘗試比較軟分群（LDA）與硬分群（K-means）在應用上的差異，希望透過此研究能使輿情分析更好的反應民意，並且提供未來發展輿情分析之參考。

陸、參考資料

今周刊，2019，〈以台灣人的角度看反送中對台灣的影響〉，

<https://www.businessday.com.tw/article/category/80396/post/201911180011/%E4%BB%A5%E5%8F%B0%E7%81%A3%E4%BA%BA%E7%9A%84%E8%A7%92%E5%BA%A6%20%E7%9C%8B%E5%8F%8D%E9%80%81%E4%B8%AD%E5%B0%8D%E5%8F%B0%E7%81%A3%E7%9A%84%E5%BD%B1%E9%9F%BF>，查閱時間 2020/06/15。

國家發展委員會，2018，〈107 年公民網路參與行為調查報告〉，

<https://ws.ndc.gov.tw/Download.ashx?u=LzAwMS9hZG1pbmlzdHJhdG9yLzEwL2NrZmlsZS9kOTU3MTFjZi02ZGU5LTQ0ZGYtYjliMS1kMTM3NWZlOTk5OTYucGRm&n=MTA35bm05YW5rCR57ay6Lev5Y%2BD6liH6KGM54K66Kq%2F5p%2B15aCx5ZGKLnBkZg%3D%3D&icon=.pdf>，查閱時間 2020/06/15。