文本分析與程式設計 Week02

本課程由卓騰語言科技贊助

學習目標

- ▶ 1. 了解什麼是「特徵詞」
- ▶ 2. 如何使用不同的工具來輔助你找到特徵詞

- 本課程使用的文本是從
- ▶ 1. https://news-taiwan.xyz/uncategorized/39053.html
- ▶ 2. https://www.ctwant.com/article/111388 (有經過編輯)



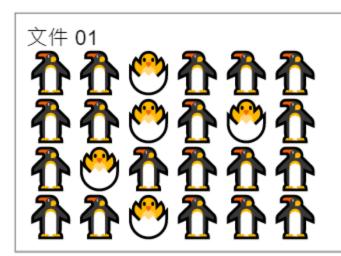
文本分析任務

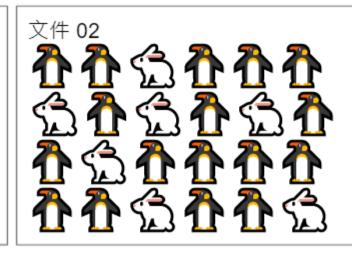
- ▶ 文本分析是常見的自然語言處理(Natural Language Processing, NLP) 任務。
- > 有三種常見的取特徵詞手法
 - 一是取 TF-IDF 特徵詞。
 - 二是依詞性取特徵詞。
 - 三是依「人、事、時、地、物」取特徵詞。
- ▶本週課程會以TF-IDF 和詞性為主
- ▶依「人、事、時、地、物」取特徵詞會是下週內容, 不過這週會先開始介紹一點



- ▶特徵詞可以做為代表文本內容的一種參考維度。
- ▶ 換句話說,一篇文本之中,那些字詞是可以拿來區 辨不同的文件

- ▶以下面這三個由emoji 的文件來舉例,請問什麼圖 案可以拿來區分這三個文件呢?
- ▶ 也就是說,哪一個「詞」是以下文件的「特徵詞」 呢?











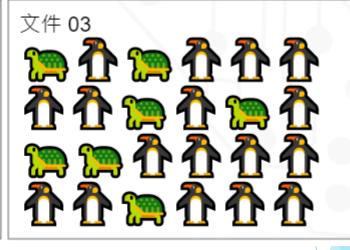


- ▶ 我們會說:
- ▶ 第一個文件的特徵詞是 🖰
- ▶ 第二個文件的特徵詞是 🏠
- 第三個文件的特徵詞是



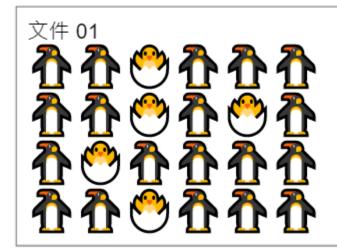


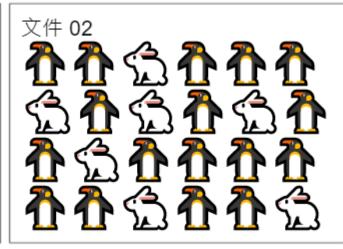




▶ ② □ 是特徵詞,因為他們是分別這三個文件中比較特別的存在

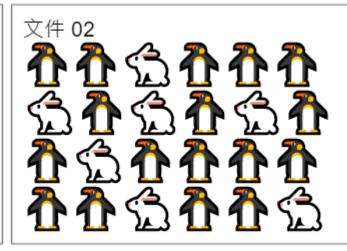








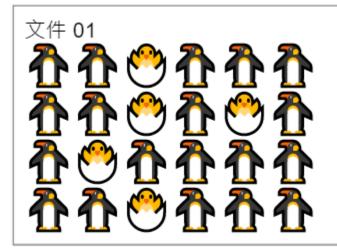


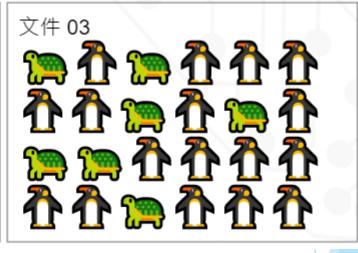




- 所以從以上的例子,我們可以知道一個詞出現的「頻率最高」不能當作判斷特徵詞的唯一指標。
- ▶ 那麼,該怎麼找到屬於人類判斷的「特徵詞」呢?

特徵詞取得方法





- ▶ 我們有三種方法
 - 取 TF-IDF
 - **看**詞性
 - 看「人、事、時、地、物」

用TF-IDF來取得特徵詞

- ► TF-IDF 是一種統計方法,用在計算某一詞在一篇文章中的重要性。
- ► 而TF 和 IDF 各有自己的意思



- ► TF 是 Term Frequency 的縮寫,意思是某個字在本文件裡出現的加權後的頻率
- ▶ 所以計算方法可以從計算每一個字的數量開始
- > 我們會以這個emoji 文件為例



- ▶ TF 的計算方法如下,以右邊的emoji 文件為例
 - 1. 計算每個詞在本文件中的頻率 例如:



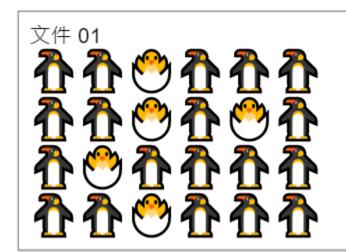


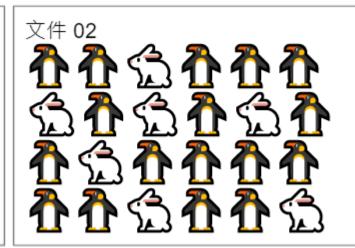


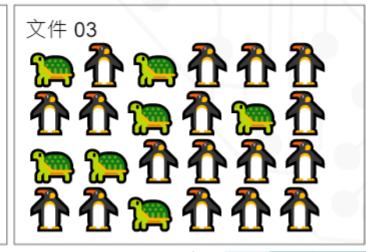
- ► 從上頁的例子,我們可以知道企鵝的TF比較大
- ► 所以從TF 的角度來說,企鵝比較重要

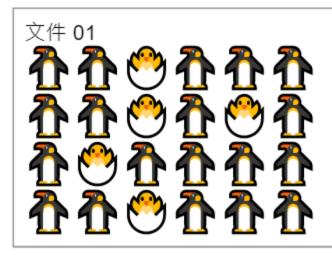


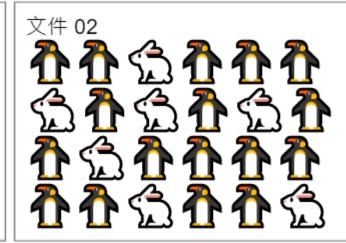
- ► IDF 是 Inverse Document Frequency的縮寫,也就是某個字在所有的文件中出現的頻率
- > 我們再來看看這個emoji 文件

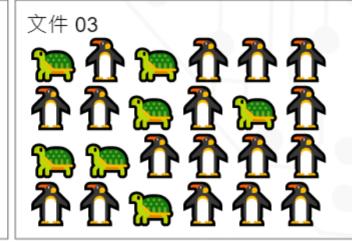




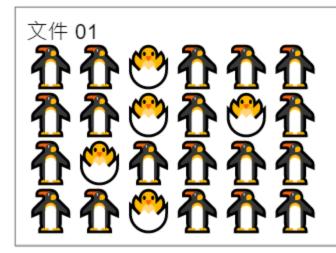


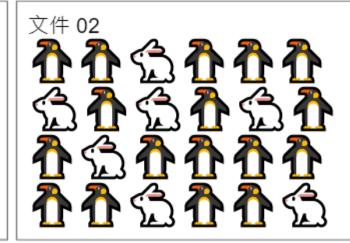


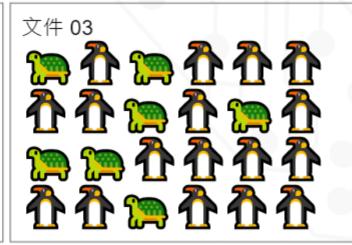




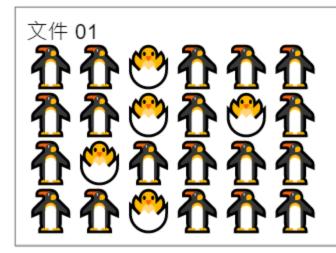
- ▶ 計算IDF時,我們要先計算這些內容
 - 1. 我們需要知道我們有多少文件 → 目前我們有三個文件
 - 2. 某個字出現在文件中的次數 例如
 - → ↑ 出現在這三個文件中
 - → 🔥 😭 各出現在一個文件中 (我們用 🖰 當例子)







- ▶ 計算IDF時,我們要先計算這些內容
 - 3. 接下來就是將剛剛的文件數取log
 - 4. 最後按照這個算式: log(文件數/該單詞出現在幾個文件)
 - → 所以 **1** 是log(3/3) = 0
 - → 🖰 是 log (3/1) 大約等於 0.477

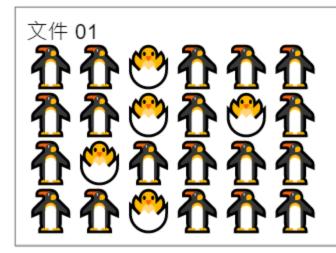




- ▶計算之後我們發現 ②的算出來比較大, 然後 1 比較小
- ▶ 這個代表什麼呢?
- ▶ 這是代表 ②在IDF的角度中,比較重要

課間練習1

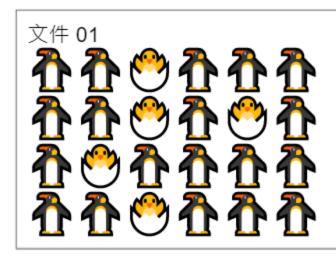
- ▶ 想想看, 你是怎麼判斷一個文章的大意呢?
- ▶ 想想看特徵詞的定義,你覺得特徵詞等同於文章大意嗎?
- ► 看完TF 和 IDF 的內容,你覺決TF 和 IDF 哪個比較接近我們人類思考之中的特徵詞呢?為什麼?
- ► 想一想TF和IDF是如何計算的,你覺得和人類判斷語意一樣嗎?為什麼?

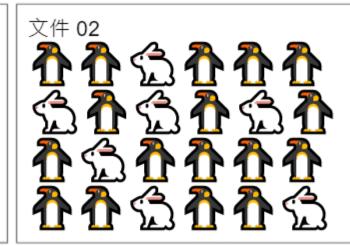


文# 02 **孙孙公孙孙 孙公孙公孙 孙公孙 孙 孙 孙**

- ▶最後,我們可以將我們計算出來的TF 和 IDF 相乘
- ▶ 例如 🏠 是0.792 * 0 = 0
- ▶ 而 🔥 是 0.2 * 0.477 = 0.0954 → 比較重要
- ▶ 0 和 0.0954 就是所謂的「 TF-IDF 權重值」

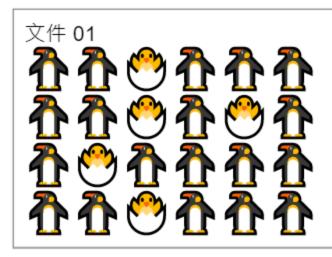
TF-IDF 可以告訴我們什麼?

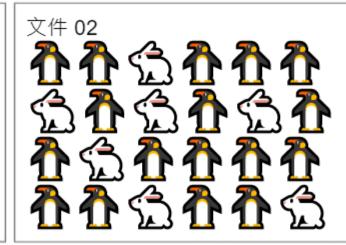


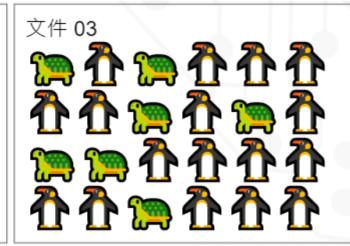


▶ 現在我們知道TF 和 IDF 背後的原理,那TF-IDF 可以告訴我們什麼呢?

TF-IDF 可以告訴我們什麼?







- ► TF-IDF 可以告訴我們文件數量和詞的頻率之間的關係
- ► 所以我們可以發現TF-IDF 可以幫助我們過濾常見的 詞。常見字詞為什麼要過濾呢?我們先來做課間練 習二來想想看

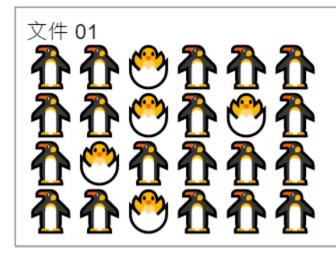
課間練習2

▶請看下一頁的新聞,並利用前一堂課學習到的 Articut 斷詞方法斷詞,並計算頻率,什麼字詞頻率 最高?高頻率的那些字是可以代表那篇新聞的字詞 嗎?

課間練習2

- ► The news is from https://www.cna.com.tw/news/amov/202106270091.aspx

TF-IDF 可以告訴我們什麼?





- ▶ 通過剛剛的課堂練習2,還有以上我們很熟悉emoji 文件,以及我們前面討論的討論詞,我們會發現反 而最常出現在文件中的詞,通常都並不是最可以突 出文件特點的詞,反而只有出現在某些文件的才是 特徵詞
- ► 如果還想更多了解TF-IDF 可以參考以下內容: https://bit.ly/tfidf_explain



- ▶ 從以上的操作,我們會發現TF-IDF ,也就是單頻 文章字詞的頻率來判斷,其實解讀內容有限。
- ▶ 為什麼呢?
- ▶ 1. TF 中所得到的高頻詞,有可能無法解讀
 - 從以上的練習題,我們會發現,一篇文章如果我們直接 只看高頻詞那我們可能會完全不了解這篇文章再講什麼, 因為這些高頻詞,很多都是「功能詞」,也就是語法架 構的一部分,並非「實詞」(content word),也就是具 備比較多「語意」的詞彙

什麼是功能詞 content words

- ► Content words,或是又稱「實詞」,指的是句子中有重要意義的意思。
- ▶ 例如以下句子 (取次於 Snow White 的維基條目)
- "Snow White" is a 19th-century German fairy tale that is today known widely across the Western world.
- ▶ 其中實詞就包括 Snow White, German, fairy tale
- ▶ 而相對於實詞,如果只有文法上功能的就是虛詞, 也就是功能詞



什麼是功能詞 content words

- ▶ 而相對於實詞,如果只有文法上功能的就是虛詞
- ▶ 例如在剛剛的句子中
 "Snow White" is a 19th-century German fairy tale that is today known widely across the Western world.
- ▶ a, that, the 就是虛詞

▶ 取自於https://en.wikipedia.org/wiki/Snow_White

課間練習3

▶ 可能你會說,我們可以把功能詞濾掉啊?然後再做 TF-IDF 分析,這樣TF-IDF 會不會比較有意義呢? 請和同學討論看看。

- ▶ 其實課間練習3 可以從兩個角度來看
- ▶ 1. 可以,但是要知道賦予TF-IDF 意義的理解語言、 懂這些字詞意義的你
 - 如果只剩下功能詞,的確從TF-IDF或許可以解構出更多的內容。不過光是從「頻率」中,其實不一定可以看出所有端倪,甚至頻率也有可能誤導,為什麼呢?

- ▶ 因為當電腦在計算TF-IDF 時其實並沒有注意語句的順序
- 如果說有兩篇文章
 - 1. 小明向小美說我愛你,你是我的最愛。
 - 2. 小明向小美說你愛我,我是你的最愛。
 - 根據上述TF-IDF 的算法,如果我們想算出第一篇和第二篇文章「你」和「我」的差異,因為字數完全一樣,所以從TF-IDF 來看,這是兩篇特徵一樣的文章。不過我們因為有上下文,也就是語境(context) 所以我們知道不一樣

- ▶ 2. 可能不具什麼特別意義
- ▶ IDF 或許很多人就會覺得已經可以抓到整篇文章的 精華,因為IDF 所計算出來的結果是某篇文章獨有 的內容
- ► 不過「某篇文章獨有內容」可以代表整篇文章想要 講的概念嗎?

- ▶ 請看以下比喻:
- ► 便當街裡,每家店都有排骨飯,但只有「好好吃便 當店」的排骨飯有加辣椒。
- ▶ 從這個敘述中,「辣椒」是一個特徵,這個好比 IDF 計算出來的結果
- ▶ 不過我們會把「好好吃便當店」的這份有加辣椒的 飯的排骨飯,叫做「辣椒飯」嗎?其實並不會

- ▶ 同理,如果我們有兩篇文章,當中IDF所計算出來的,其實也頂多算是我們辨別某些文章的特徵,但是是不是這個特徵一定等於整篇文章的大意,並不一定
- ► 所以當我們在理解TF-IDF的內容時,我們需要清楚TF-IDF的原理,以及分清楚我們想要透過文本分析想要分析什麼內容
- ▶ 我們需要清楚理解TF-IDF僅是工具,TF-IDF結果其 實是不能直接等同於文本解釋



延伸資料

► 如果還想更多了解TF-IDF 可以參考以下內容: https://www.cc.ntu.edu.tw/chinese/epaper/0031/2 0141220_3103.html

▶以下這個連結是有關Jieba以及Articut的TF-IDF運 算以及優缺點比較:

https://blog.droidtown.co/post/186883773617/tf-idf

如何利用TF-IDF 來取出特徵詞

- ► Articut 內建的 analyse 工具包裡的 extract_tags() 函式來取得經TF-IDF計算的特徵詞
- 延續上週課程中的棒球與籃球的例子

```
baseball_TFIDF = articut.analyse.extract_tags(baseballResultDICT)
print(baseball_TFIDF)
```

```
basketball_TFIDF = articut.analyse.extract_tags(baseballResultDICT)
print(baseball_TFIDF)
```



如何利用TF-IDF 來取出特徵詞

棒球語料的 TF-IDF 存在 baseball_TFIDF 裡,印出後內容如下:

執行結果

['\n', '大都會', '馬林魚', '投手', '敲出', '安打', '2', '1', '巴斯', '局面', '康福托', '比賽', '被', '週三', '紐約', '此戰', '使用', '車輪戰', '4名', '輪番', '上陣', '壓制', '打線']

► 籃球語料的 TF-IDF 存在 baseball_TFIDF 裡,印出後內容如下:

執行結果

['\n', '兩', '米契爾', '康利', '罰', '勝利', '太陽', '攻勢', '爵士', '籃板', '3', '助攻', '11分', '投', '一', '昨晚', '紐約', '西區', '霸王']

TF-IDF 的計算結果:觀察

▶ 從兩篇文章的 TF-IDF 特徵詞裡,如果沒有相關的 背景知識知道「大都會」、「馬林魚」是棒球大聯 盟的球隊,而「爵士」是 NBA 的球隊名稱的話, 其實很難確認究竟哪些詞彙是可以做為「棒球類」 或是「籃球類」的文本分類特徵詞。

TF-IDF 的計算結果:觀察

► 但我們可以確認的是「投手」、「安打」、「打線」 這三個 名詞,應該是只有棒球類的文本才會用到。 而「籃板」這個 名詞,則是籃球類的文本才會用到。 同理,我們也能觀察到「敲出」、「保送」、「打擊」這幾個動詞 同時出現的文章,有很大的機率是在描寫棒球類的文本。而「助攻」、「上籃」則常見於描寫籃球比賽過程的動詞。

TF-IDF 的計算結果:觀察

▶ 這個觀察最重要的是讓我們發現「詞性」(前文提到的「<mark>名詞</mark>」、「<mark>動詞</mark>」)是可以做為抽取特徵詞的一種方法。接下來,我們利用 Articut 的詞性篩選工具來取出「<mark>名詞</mark>」和「動詞」。

用詞性來取得特徵詞

名詞比較-輸入

▶ 我們使用 Articut 內建的 .getNounStemLIST() 函式 來取得文本中各句的名詞。

```
baseballNounLIST = articut.getNounStemLIST(baseballResultDICT)
print(wordExtractor(baseballNounLIST))
```

basketballNounLIST = articut.getNounStemLIST(basketballResultDICT)
print(wordExtractor(basketballNounLIST))

名詞比較-輸出

棒球報導文本的名詞

執行結果

['主審', '優勢', '分', <mark>'勝利'</mark>, '安打', '局面', '意識地', '手', '打線', '打者', '投手', '此戰', '比賽', '滿壘', '球', '球隊', '登板', '眼', '角滑球', '觸身球', '身肘', '身體', '車輪戰', '追平比數', '這球', '陽春砲', '領先']

▶ 籃球報導文本的名詞:

執行結果

['中', '人', <mark>'勝利', </mark>'單場', '場', '布克', '延長賽', '戰', '攻勢', '機會', '次', '波格丹諾維奇', '籃板', '階段', '霸王']

▶ 從名詞裡看出來,有些詞彙(例如「勝利」)是重覆的,要得到更好的效果的話,還可以把重覆的詞彙去除以便讓不同的類別都具有各自獨立的特徵詞表。

名詞比較-輸出

棒球報導文本的名詞

執行結果

▶ 籃球報導文本的名詞:

執行結果

就觀察名詞的情形,黃色部分的標記有助於我們理解文本所討論的主題是棒球類,相較於籃球綠色標記的部分有助於我們理解文本主題為籃球。

動詞比較-輸入

▶ 我們使用使用 Articut 內建的 .getVerbStemLIST() 函式來取得文本中各句的名詞。

```
baseballVerbLIST = articut.getVerbStemLIST(baseballResultDICT)
print(wordExtractor(baseballVerbLIST))
```

basketballVerbLIST = articut.getVerbStemLIST(basketballResultDICT)
print(wordExtractor(basketballVerbLIST))

動詞比較-輸出

棒球報導文本的動詞

執行結果

['上場', '上陣', '伸進', '使用', '保送', '再見', '判定', '到', '壓制', '失', '帶', '形成', '打出', '打 擊', '投', '投出', '拿下', '接下來', '推出', '敲', '有下', '看', '碰觸', '讓', '贏得', '越來', '輪到', ' 進入', '關', '隨', '靠近', '面對']

▶ 籃球報導文本的動詞:

執行結果

['上籃', '分', '力圖', '助攻', '合計', '回應', '得分', '得手', '打出', '抄截', '投', '拿到', '施展', '比', '犯規', '狂轟', '用', '留給', '罰', '讀秒', '追分', '追平', '造成', '錯失', '鎖定', '開始', '隨']

► 從動詞裡看出來,有些詞彙 (例如「打出」、「投」、「隨」) 是重覆的,要得到更好的效果的話,還可以把重覆的詞彙去除以便讓不同的類別都具有各自獨立的特徵詞表。



動詞比較-輸出

棒球報導文本的動詞

執行結果

['上場', '上陣', '伸進', '使用', '<mark>保送</mark>', '再見', '判定', '到', '壓制', '失', '帶', '形成', '打出', '<mark>打擊</mark>', '投', '投出', '拿下', '接下來', '推出', '<mark>敲</mark>', '有下', '看', '碰觸', '讓', '贏得', '越來', '輪到', '進入', '關', '隨', '靠近', '面對']

▶ 籃球報導文本的動詞:

執行結果

['<mark>上籃</mark>', '分', '力圖', '助攻', '合計', '回應', '得分', '得手', '打出', '抄截', '投', '拿到', '施展', ' 比', '犯規', '狂轟', '用', '留給', '罰', '<mark>讀秒</mark>', '追分', '追平', '造成', '錯失', '鎖定', '開始', '隨'**]**

▶ 針對動詞的部分,則是在以上標記的部分顯示出文本類別的線索,要特別注意的是,再區分文本時需要多多考慮不同的詞性(例如:名詞與動詞),如此一來才能正確找出文本的區分關鍵。



用「人、事、時、地、物」取特徵詞

取得人事時地物的工具

- 如果你想知道的特徵詞可以以「人事時地物」來分類,那可以嘗試使用以下的工具來幫助你
- ▶「人」、「事」、「地」、「時」和「物」因為通常都是名詞,所以可以使用下面的工具
 - 可以用剛剛取名詞的.getNounStemLIST()
 - 可以用 articut.getContentWordLIST(),這個是用來取得 content word
- ~ 不過更詳細的內容,下週會再介紹



取得人事時地物的工具

- 如果你想知道的特徵詞可以以「人事時地物」來分類,那可以嘗試使用以下的工具來幫助你
- ▶「地」
 - Articut 也可以直接擷取地方的名稱
 - 可以用 articut.getLocationStemLIST()

課間練習4

- ▶請使用籃球和棒球的文章來練習
- ▶ articut.getLocationStemLIST() 和用 articut.getContentWordLIST() 和前面取名 詞和動詞的方法是一樣的,請嘗試自己摸索使用看 看。
- ▶ 觀察看看 content word (內容詞) 和 location word (地方名稱) 是否具有比「動詞」或「名詞」更好/更差的特徵表現能力。說說看你的觀察。
- ▶ 請直接使用之前的basketball 和 baseball 的 resultDICT



- ▶ 作業敘述:
- ► 2021年由於疫情關係遠距教學十分盛行,很多學者對遠距教學裡頭的學生端以及教師端進行不少研究, 而在這個任務中我們將從另一個新的角度來探討遠 距教學對社會的影響,我們蒐集了10篇從聯合新聞 網的新聞,利用關鍵字*遠距教學*進行檢索,並取出 前十篇進行分析。

- 資料概述
- ▶ 新聞來源: 聯合新聞網

	標題	類 別	內容
0	仁寶Q3營收 可望成長	股 市	代工廠大仁寶(2324)今年來持續受惠居家工作、遠距教學等宅 經濟熱潮,推升筆電銷售暢旺,惟第
1	元山 下半年出貨 看增	股 市	隨著歐美疫苗施打率日益提升,車市有望在下半年回溫,國內車 用電子與生活科技系統廠元山(6275
2	宅經濟退燒 Chromebook、平板看淡	產 經	新冠肺炎疫情引發居家工作、遠距教學熱潮,帶來的 Chromebook、平板等宅經濟硬體需求熱
3	署假將至兒盟籲政府助弱勢兒少課	文	全台學校自5月下旬停止到校上課,採遠距學習,兒福聯盟今舉
	輔轉型線上陪伴	教	行「弱勢兒少數位學習困境與對策」線上
4	竹市議員發起愛心筆電助學計畫 支援	文	三級警戒以來全國已停課快兩個月,所有學生需依賴3C設備線上
	弱勢學子署假學習	教	上課,弱勢家庭數位落差大,缺乏手機
5	手機終有訊號了 雙溪區魚行里基地台	地	雙溪區魚行里的坤溪、公館、八股、丁子蘭等聚落,一直以來都
	啟用	方	是區內通訊的隱蔽帶,造成居民手機訊號
6	避免「線上中輟」! 中山工商遠距教	文	疫情讓各級學校停課」,採遠距教學已超過一個月,有學校出現
	學按表操課點名課輔	教	「線上中輟」現象,令家長憂心忡忡。高



- ▶ 任務1: 找出文教類的特徵詞
- ▶ 請運用Articut裡頭extract_tags()這個功能找出文教類的前 20名特徵詞吧!
- ▶ 任務2: 找出股市類的特徵詞
- ▶ 請運用Articut裡頭`extract_tags()`這個功能找出股市類的 前20名特徵詞吧!
- ▶ 任務3: 比較兩類特徵詞
- ▶ 請你比較兩類特徵詞,請問你覺得有那些字詞能代表文教類,而那些字詞能代表股市類?而那些具代表性的字詞反映出遠距教學對這兩方面的影響是什麼呢?



- ► 任務4: 根據特徵詞將文本進行歸類¶
- ▶ 假如今天沒有產經這一個新聞類別,你覺得那些新聞應該要歸到哪文教類還是股市類呢?請試著用特徵詞的方式進行比較。

- ▶ 任務5: 運用名詞比較來理解文本差異
- ▶ 從上面的特徵詞來看感覺名詞較多,因此我們希望 更進一步從名詞下手,希望你能將文教類以及股市 類裡的名詞用`getNounStemLIST()`這個方法找出 名詞,然後用我們之前在week1裡面學到的計算字 詞頻率的方式,做出他的頻率表,之後將頻率表視 覺化成文字雲。如下所示:



- ▶ 任務5: 運用名詞比較來理解文本差異
- ▶ **5.1** 名詞頻率分析表格(文教類為例) 作方法詳見上週作業

	Noun	FREQ
0	學校	17
1	學生	17
2	孩子	12
3	老師	10
4	設備	9
5	家長	7
6	網路	7
7	教師	7
8	資源	7
9	弱勢兒	7
10	兒盟	7
11	疫情	7

- ▶ 任務5: 運用名詞比較來理解文本差異
- ▶ 5.2 文字雲分析(文教類為例)
- ▶ 所需工具解釋請見下一頁



```
wordcloud = WordCloud(width = 800, height = 800, #設定畫布大小
             background color ='white', #背景顏色
             min font size = 10, #最小字體
              font path="TaipeiSansTCBeta-Regular.ttf") #畫中文的文字雲一定要設定字體
路徑(你可以直接把字體放在存這份檔案的同一層)
wordcloud.generate from frequencies(frequencies=data) #繪製文字雲時有很多方法,其中一個
方法就是將檔案轉為dictionary的格式
plt.figure(figsize = (8, 8)) #顯示的圖大小
plt.imshow(wordcloud, interpolation="bilinear") #這部分是有關顏色顯示的方式,參考
plt.axis("off") #不要有x,y軸的座標
plt.show() #顯示圖片
```