

美食公道伯



第二組

巨資四B 07170203 陳永珅

巨資四B 07170213 陳政彣

巨資四B 07170232 何詩彥



Table of contents



04 Algorithm

02 定義 Graph

05 落地實現

03 資料處理

06 心得分享

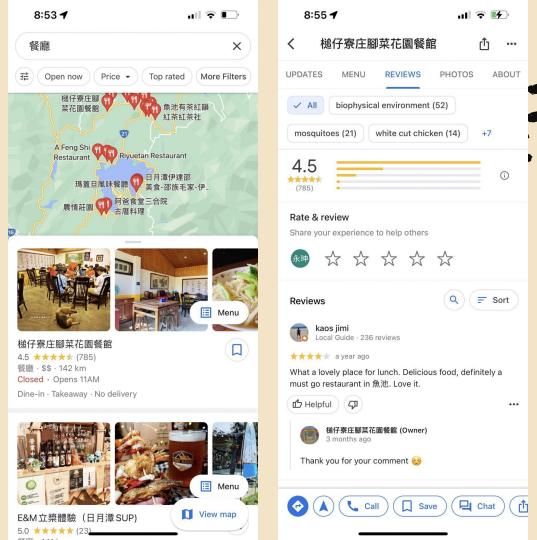
01

發想&資料探索



情境 & 提問

大家是否曾經與家人、朋 友出遊,使用Google Map 搜尋餐廳?除了參考其他 網頁, 在 Google Map 中 哪些資訊會列入你的考



參考資訊



以下是我們團隊會參考的資訊:

星度、餐廳被留言數量、TopN留言、近期留言、留言照片、該篇留言讚數、留言者有沒有頭貼、留言內容、留言者聲量(留言次數、留言被按讚數)。

總結來說我們除了看他的評分高低以外也會關注這個評分是否真實。



真實度

我們參考的資訊大多與用戶留言相關(用戶與餐廳的互動),那是否可以透過社群網絡來找出餐廳評分的真實度?

餐廳名稱	Stars	Truth
ABC bar	3	1(?)
CDE bar	5	1.5(?)
FGH bar	5	6(?)



我們嘗試在 yelp 資料集中找到以下資料:

餐廳被留言數量、TopN留言、近期留言、留言的照片、該篇留言讚數、留言者有沒有頭貼、留言內容、留言者聲量(留言次數、留言被按讚數)。

使用了review、business、user資料集



名稱 个 business.json --checkin.json 🚢 Dataset_Challenge_Dataset_Agreement.pdf photo.json --review.json 🚢 tip.ison -user.json 🚢 Yelp_Dataset_Challenge_Round_13.pdf 45

誰反映了真實度?





資料篩選

基於Top N留言、近期留言、留言者質量, 我們篩選了什麼資料?

- 1. Las Vegas 當地屬於 restaurants & bars、有營業的店家(運算考量)
- 2. 2018 年有 elite 資格的用戶
- 3. 20180101 以後的評論







date	text	cool	funny	useful	stars	business_id	user_id	review_id	
2018/4/27 18:53	Went in on a Friday since I was strolling arou	0	0	0	4	gOOfBSBZlffCkQ7dr7cpdw	1RQaXb0xSLsUzwJ0u4uZ-w	C0onMemq7n2bMx3fMQXF9w	0
	Came to have dinner and they were pretty busy	0	0	0	3	H5Y6o9H4rtxbHJFgLql6Fw	qXcTj0AyR6BF5wQpzRSYaQ	V_Wd_Ybjvr47o3utbEa9dg	1
2018/5/28 21:25	I love Q Karaoke so I decided to eat at QBistr	1	1	1	5	XsSgv3vBOyOBXn3Co8EVIg	qEjm0_ivRn8WlfJipeFAgg	g2s2D_gDldt1SHA9_03DHQ	2
2018/8/14 4:17	Great service at the bar for happy hour and in	0	0	0	4	pH0BLkL4cbxKzu471VZnuA	kdDTqKBbfNZeKEbWAcqZWQ	T0Hs0KAAuHEzHXq00baJUQ	3
2018/5/29 2:12	I went on national burger day ha ha so I had a	0	0	0	4	isw3cS3hOKdKeBgi3lF3-A	IPILcdNG426qOopHH4EwkQ	H9iVzGZklifXAHBSanqWOw	4

誰反映了真實度?



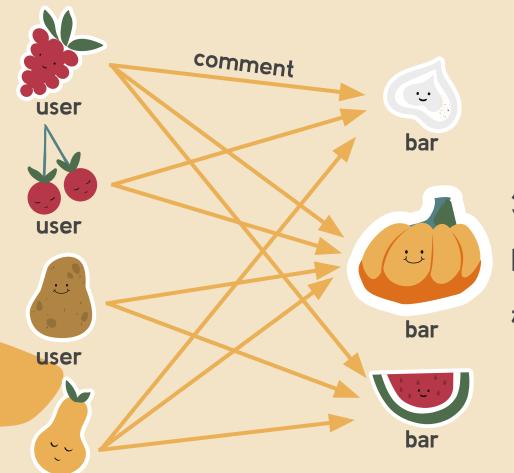






定義 Graph





user



節點:菁英用戶、餐廳&酒吧

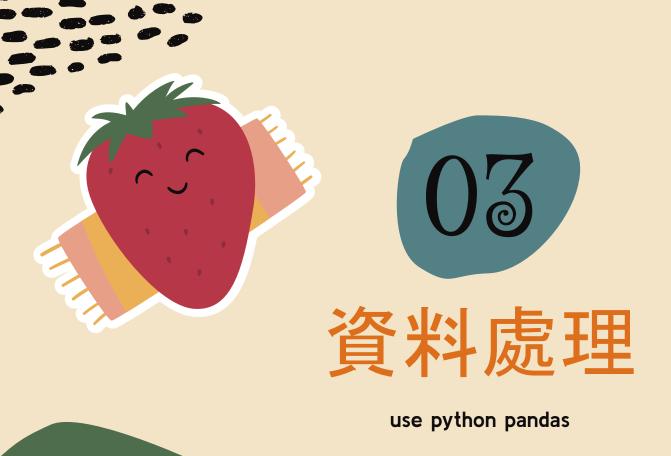
關係:用戶「評論」餐聽

權重:useful、funny、cool 相加

誰反映了真實度?







1. 以這張表當主鍵, 去合併user及business資料集

	review_id	user_id	business_id	stars	useful	funny	cool	text	date
0	C0onMemq7n2bMx3fMQXF9w	1RQaXb0xSLsUzwJ0u4uZ-w	gOOfBSBZlffCkQ7dr7cpdw	4	0	0	0	Went in on a Friday since I was strolling arou	2018/4/27 18:53
1	V_Wd_Ybjvr47o3utbEa9dg	qXcTj0AyR6BF5wQpzRSYaQ	H5Y6o9H4rtxbHJFgLql6Fw	3	0	0	0	Came to have dinner and they were pretty busy	2018/2/11 3:30
2	g2s2D_gDldt1SHA9_03DHQ	qEjm0_ivRn8WlfJipeFAgg	XsSgv3vBOyOBXn3Co8EVIg	5	1	1	1	I love Q Karaoke so I decided to eat at QBistr	2018/5/28 21:25
3	T0Hs0KAAuHEzHXq00baJUQ	kdDTqKBbfNZeKEbWAcqZWQ	pH0BLkL4cbxKzu471VZnuA	4	0	0	0	Great service at the bar for happy hour and in	2018/8/14 4:17
4	H9iVzGZklifXAHBSanqWOw	IPILcdNG426qOopHH4EwkQ	isw3cS3hOKdKeBgi3lF3-A	4	0	0	0	I went on national burger day ha ha so I had a	2018/5/29 2:12



- 2. 將重複評論的列合併。
- 3. 篩選要使用的欄位: business_name、 user_name、 stars、useful、funny、cool
- 4. 將 useful、funny、cool相加當作權重(weight)





結果



餐廳數量:663

評論數:9997

用戶數:1298

04 Algorithm





Algorithm









PageRank - define method

「PageRank」是 Google 用於「評等網頁重要性」的一種方式。

	Google	Our case		
數量假設	一個網頁收到的其他網頁指向的 入 鏈結(in-link)越多, 說明該網頁越 重要。	一間餐廳 被評論的數量越多 , 說明 該家餐廳評分 越趨近於真實 。		
質量假設	當一個 質量高的網 頁指向 (out-link)的網頁也質量也高。	評論餐廳的人質量越高 ,說明這家 餐廳的 評分質量也越高 。		

誰反映了真實度?

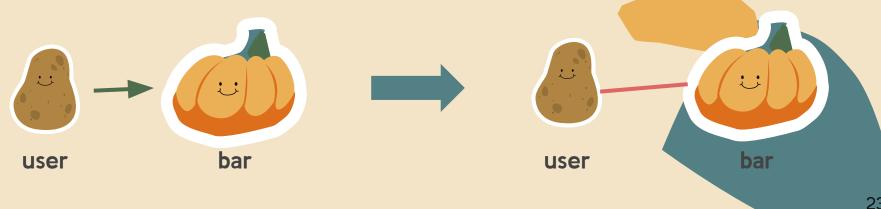


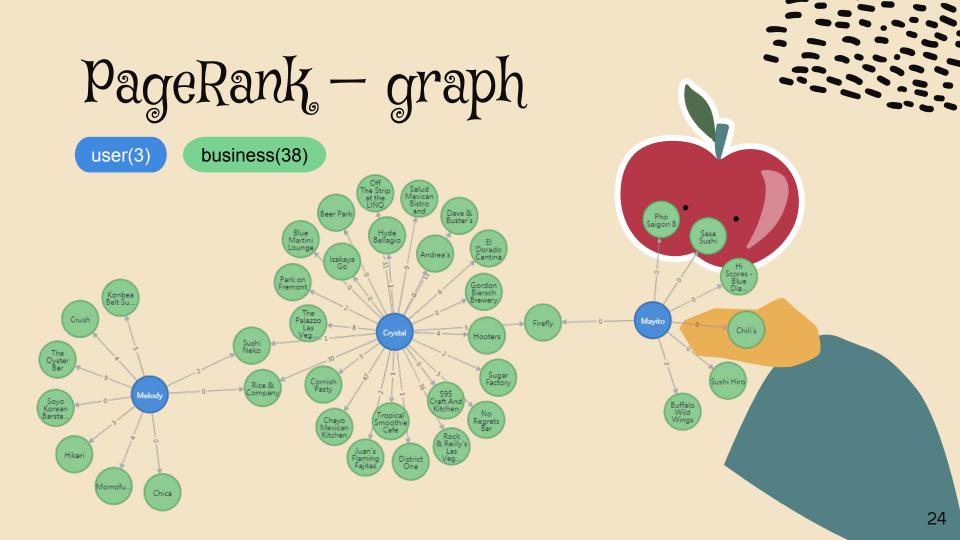


PageRank — define method

調整成無向圖

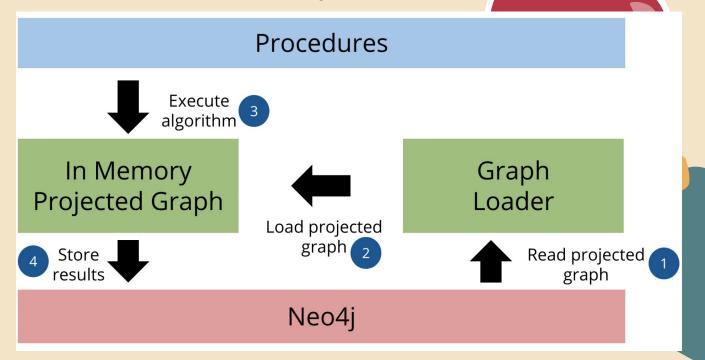
註:計算完 PageRank 分數後再將用戶節點的分數移除。





PageRank - gds info

Graph Data Science library



PageRank - code_add_graph

```
CALL gds.graph.create(
      'myGraph',
     ['user', 'business'],
 4
 5
       Related: {
          orientation: 'UNDIRECTED',
          properties: 'weight'
 9
10
```

PageRank - code_PageRank -

- 1 CALL gds.pageRank.stream('myGraph') YIELD nodeId, score
- 2 WITH gds.util.asNode(nodeId) AS n, score
- 3 MATCH (n)-[r:Related]-()
- 4 RETURN n.Name AS name, score, count(r) AS interactions
- 5 ORDER BY score DESC





"name"	"score"	"interactions"
"Momofuku Las Vegas"	25.641347145149098	204
"HEXX kitchen + bar"	16.85194591643035	136
"Lotus of Siam"	15.666573188885478	132
"Sapporo Revolving Sushi"	12.962212150692016	124
"Gordon Ramsay Pub & Grill"	12.777973644927615	102
"Holsteins"	12.272574440683506	103
"SkinnyFATS"	11.958860096266198	108
"Black Tap"	11.931775801865307	105
"Yard House"	11.036061905182281	98
"Beauty & Essex"	10.984403766069727	89
"The Peppermill Restaurant & Fireside Lounge"	10.752671837482556	90

PageRank - remove_user

	name	score
0	Momofuku Las Vegas	25.641347
1	HEXX kitchen + bar	16.851946
2	Lotus of Siam	15.666573
3	Sapporo Revolving Sushi	12.962212
4	Gordon Ramsay Pub & Grill	12.777974
659	Holo Holo	0.237453
660	Agua El Manantial	0.237397

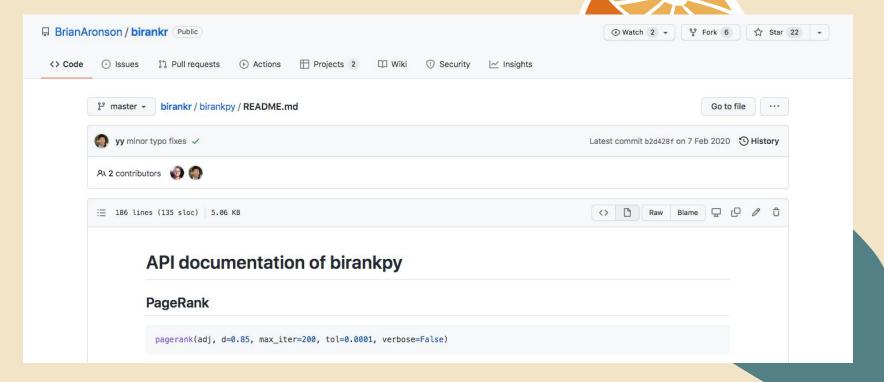
Birank - method intro

Bipartite networks are commonly reduced to unipartite networks for further analysis, such as calculating node centrality (e.g. PageRank,see Figure 1(c)).

To overcome the issues of one-mode projection, we present BiRank, an R and Python package that performs PageRank on bipartite networks directly.



Birank – package



Birank — data

	business_name	user_name	stars	weight
0	Chica	Melody	4	0
1	Chili's	Mayito	3	0
2	Q Bistro	Crystal	5	3
3	SUSHISAMBA - Las Vegas	Lulu	4	0
4	Eureka!	Theresa	4	0
7	SUSHISAMBA - Las Vegas	Juliet	4	0
8	Sammy's Beach Bar & Grill	Rodrigo	3	2
9	Hawthorn Grill	George	5	3
10	Battista's Hole In the Wall	Jhordan	3	7

Birank - code

To performing BiRank on this bipartite network, just:

import birankpy

bn = birankpy.BipartiteNetwork()

bn.set_edgelist(df, top_col='business_name', bottom_col='user_name',

weight_col=weight)

top_birank_df, bottom_birank_df = bn.generate_birank()



Birank — output

	business_name	business_name_birank
0	Holsteins	3.352381e-01
1	Malena's Yogurt Plus	3.101665e-01
2	Blue Martini Lounge	3.219785e-02
3	Rí Rá Irish Pub	2.012090e-02
4	Lazy Dog Restaurant & Bar	1.937041e-02
658	Al-Hayat Hookah Lounge	2.869219e-07
659	Agua El Manantial	2.869219e-07
660	Margaritaville Casino	2.869219e-07
661	Oliva Minimart & Cafeteria	2.869219e-07
662	OTORO Robata Grill & Sushi	2.869219e-07

落地實現



→ 將透過 PageRank 計算得 到的排序轉換為 PR 值

```
output['rank'] = output.index + 1

N = len(output)
def pr(s):
    pr_val = int((N - s)*100 / N)
    return pr_val

output['pr'] = output['rank'].apply(pr)
```

0
1
2
3
4
658
659





PR 分數的疑慮



惡性競爭

無法爭取絕對分數的高分 -> 攻擊比自己高分的店家



使用者判讀

找到兩家 PR 不高的餐廳 -> 感覺 star 都是假的





商家流失

平台上一定有 PR 低的店家 -> 不願待於己不利的平台









資料觀察

從互動角度切入觀察





心得重點



分析結果

客觀、有效率

分析方法

透過 Graph 看到隱藏的資訊

SNA 是以一個與過往完全不同的角度來解讀資料, 互動行為常常隱匿在資料集堂 中。起初訂定題目的時候常常 盯著資料集的欄位左思右想, 每每有不錯的 陳永珅 卻發現用不到 graph, 完全把 graph 當作視覺化工具在用!但我們真正用 graph 來解 決問題的時候,就能發現考量到整個網絡互動的結果更為客觀。 平常使用Google Map在查找附近的餐廳時,可能會找到數間星數都差不多的店家。 這時要從這些店家中找出一間前往拜訪,往往都需要從下面的評論區仔細觀察後, 何詩彥 綜合比較出一間自己認為最好的餐廳。這個過程其實需要耗費不少時間, 如果我們 的這個想法可以實現的話, 這個問題就可以更客觀、更有效率的解決, 相信能成為 有選擇障礙朋友的好夥伴。 這學期首次接觸 SNA, 以前都只是將資料集的每個欄位去做分析, 看看有哪些關聯 ,但沒想過平凡無奇的資料集,原來可以有那麼多連結,各自的節點看似沒關聯,但 陳政彣 卻影響著整個整個生態, 這次主題範圍訂在餐廳及 Bar, 但用相同的模式, 只要找得

到資料集,可以將這個主題運用在更多場景上,讓踩雷的人越來越少。

小組分工

	分工項目
陳永珅	Neo4j 演算法、主題發想、簡報製作
陳政彣	資料處理、主題發想、簡報製作
何詩彥	Neo4j 製圖、主題發想、簡報製作

