



GoPic

Go to the Destination with your Imagination.

GOPICT FOR YOUR TRAVEL

- 報告日期: 2016.08.04
- 受訓梯次: AB101
- 報告組別: Group 2

成員介紹

組長 鄭云鈞



組員 李鈺祥



組員 王珮涵



組員 倪裕程



組員 黃郁棻

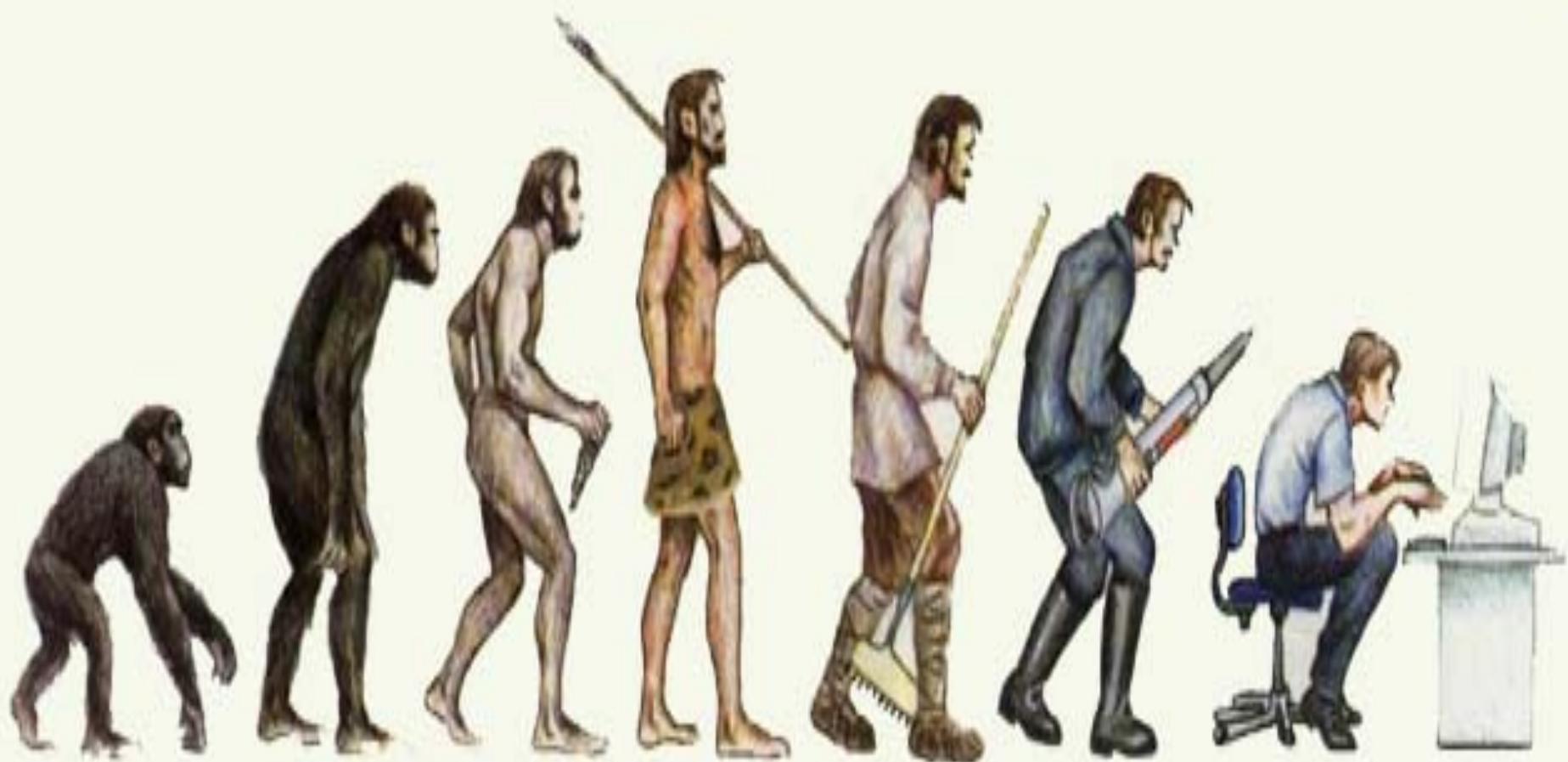


組員 郭南宏

OUTLINE

- 專題源起
- 專題架構
- 資料收集
- 專題核心
- 未來展望
- Q&A





專題源起 | 黃郁棻

原想做旅行社行程/價格比較...

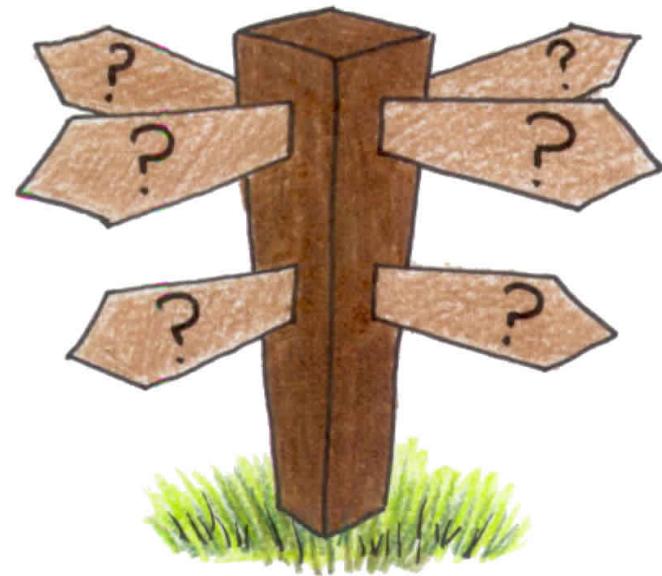


但旅遊咖和FUN TIME...

景點	精選東京限量晴空塔展望台.橫濱紅磚倉庫.台場鋼彈.迪士尼.箱根遊船.溫泉饗宴5日	靜岡綠茶之鄉～東京迪士尼.久能東照宮.忍野八海.猿劇場.溫泉五日	《線上旅展2000》東京四大主題～迪士尼.鬼太郎.猿劇場.藤子不二雄.保證上晴空塔箱根.海盜船溫泉5日
	景點總數：16處 極力推薦：9處 旅遊天數：5天	景點總數：15處 極力推薦：7處 旅遊天數：5天	景點總數：14處 極力推薦：5處 旅遊天數：5天

共同景點 ?	 富士箱根國立公園 ★★★★★ 國家公園^自然與公園	 唉呀這景點我沒有	 富士箱根國立公園 ★★★★★ 國家公園^自然與公園
	 箱根蘆之湖海盜船 ★★★★★ 遊船^搭船遊覽與水上運動^遊覽^戶外活動	 唉呀這景點我沒有	 箱根蘆之湖海盜船 ★★★★★ 遊船^搭船遊覽與水上運動^遊覽^戶外活動
	 忍野八海 ★★★★★ 景點和地標^景點與地標	 唉呀這景點我沒有	 忍野八海 ★★★★★ 景點和地標^景點與地標
	 免稅店 ★★★★★	 免稅店 ★★★★★	 免稅店 ★★★★★

只有圖卻不知道在哪..... ?



一張圖，勝過許多關鍵字！



V.S

Google 搜尋

全部 圖片 地圖 新聞 影片 更多 ▾ 搜尋工具

約有 189,000 項結果 (搜尋時間 : 0.88 秒)

[勇者鬥惡龍創世小玩家全挑戰詳細攻略 - 哈啦區 - 巴哈姆特](#)
forum.gamer.com.tw/C.php?bsn=00251&sna=18998 ▾

2016年5月1日 - 13 篇文章 - 0 位作者

或把牆壁地板升級成石壁板或木壁板來衝到五級吧！... 第一個：首張地圖於擊敗主線病原體大鍋牛的正西方山丘上有**一棟**猶謎 ... A：藍門的**城堡**附近四周有浮在天上的石塊，用聖水解咒就會恢復成鐵塊與史萊姆：初始地圖（主角剛到拉達托姆的第一章圖）再石化公主的石像上方，可以找到史萊姆，交付物品為**白色花**。

【攻略】主線任務攻略(全) @勇者鬥惡龍系列哈啦板 - 巴哈姆特 2016年5月4日

【森林】暮光森林模組(TwilightForest Mod) (圖片 ...) 2013年5月15日

forum.gamer.com.tw 的其他相關資訊

[台北親子餐廳 大樹先生的家 城堡球池,樹屋溜滑梯,沙坑,室內戶外都 ...](#)
cline1413.pixnet.net/blog.../395724047 ▾ [台北親子餐廳 大樹先生的家 城堡球池](#) ▾

2014年10月29日 - 我們大人也無法好好吃頓飯說說話選擇安全又有寬敞遊戲環境的親子餐廳,大人小孩都會 ... 主建築是一棟**白色**羅馬柱的洋房,很難想像裡面是親子餐廳 ... 走進室內,第一眼看到的是讓孩子們瘋狂的木製**城堡**溜滑梯球池,這座球 ... 沖繩自行-沖繩必住♥聖瑪莉娜海濱飯店SUN MARINA♥網友大推。c/p值高有美麗沙灘。

[【台東民宿】布拉諾城堡@ 灰頭土臉:: 隨意窩Xuite日誌](#)

blog.xuite.net/nike4859/mind/23200764-【台東民宿】布拉諾城堡 ▾

200902192313 【台東民宿】布拉諾城堡 ?旅遊食記 ... 強烈的白色建築，以**城堡**的概念搭造而成的民宿。布拉諾的 ... 地中海風格的色調，有著很大一面的**白牆**。顏色藍色 ...

該如何達成？



保證成團

【時尚東京】**迪士尼樂園**・晴空塔・杜莎夫人蠟像館・明治神宮・淺草觀
◎出發日期：8/04E、8/15E、8/21E

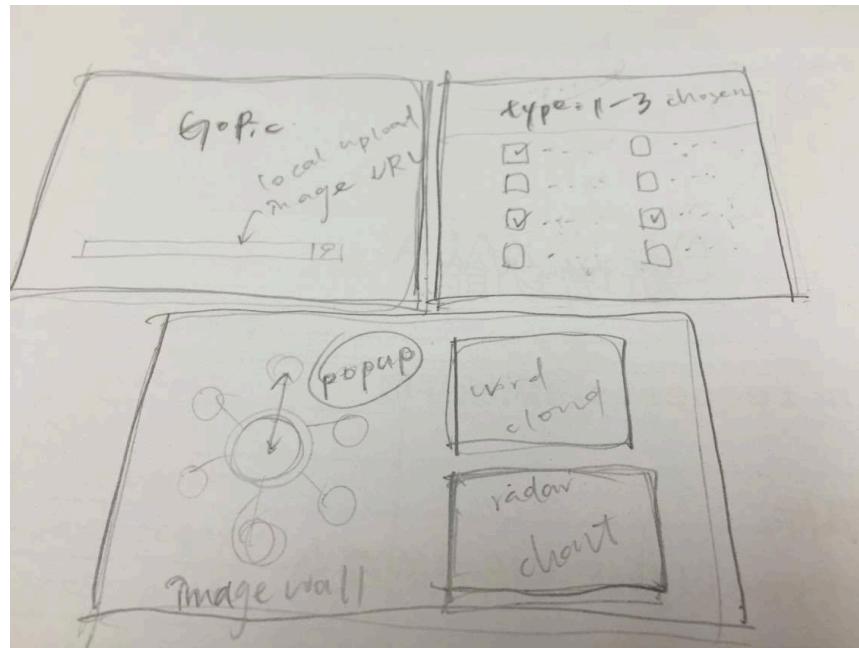
【東京雙樂園】**迪士尼**・KITTY樂園・海盜船・忍野八海・OUTLET・溫!
◎出發日期：8/22G、8/25C

【富士鐵道・靜岡+東京】天上山纜車・淺間大社・**迪士尼**・湯瑪士列車
◎出發日期：8/11

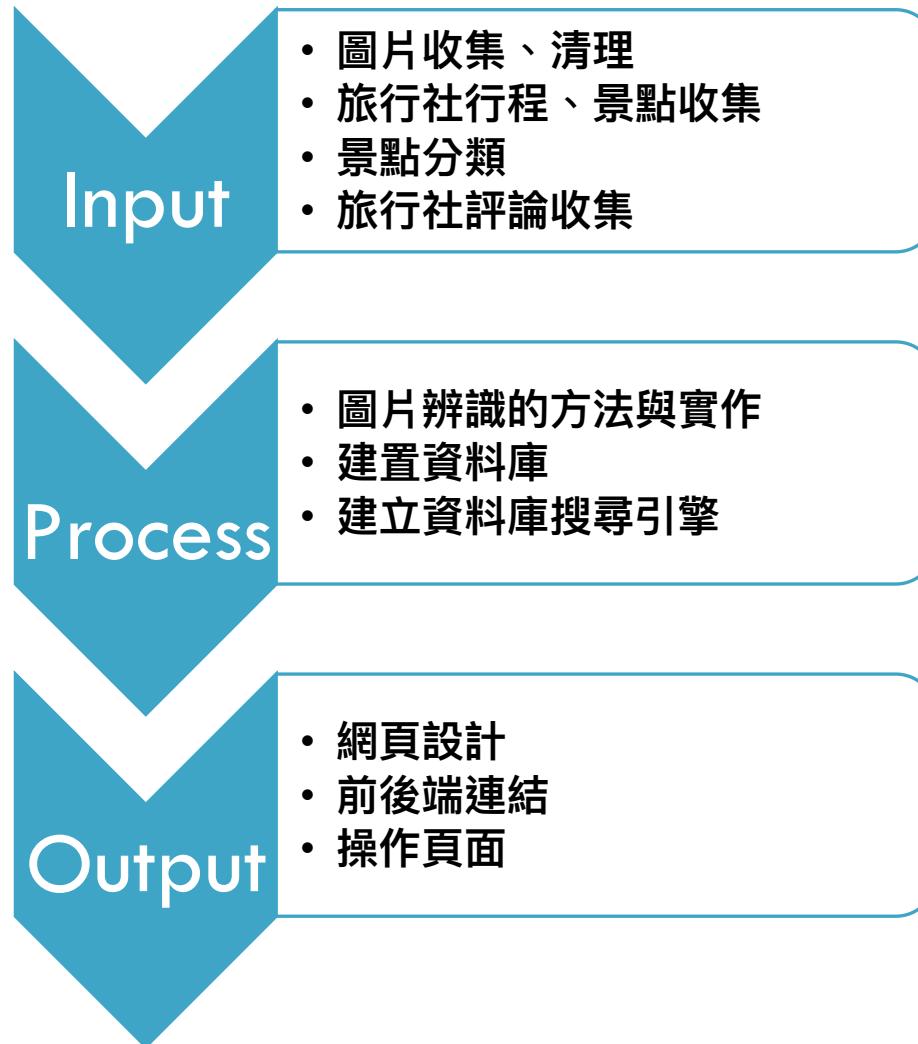
【FUN暑假・東京雙樂園】**迪士尼**・KITTY樂園・箱根遊船・淺草寺・O
5天

初步構想

- ◆ 以圖推薦旅行社行程
- ◆ 考量使用者偏好的景點種類
- ◆ 提供PTT/背包客棧對於旅行社的評價



專題架構



專題過程—工作配置

項目 成員 \	專題綜合 處理	圖像辨識 & Hadoop	文字探勘	資料庫 & 前後端連接	前端網頁 開發	機動協助
鄭云鈞	✓	✓		✓		
倪裕程	✓	✓		✓		
李鈺祥	✓		✓	✓		
王玟涵	✓	✓		✓		✓

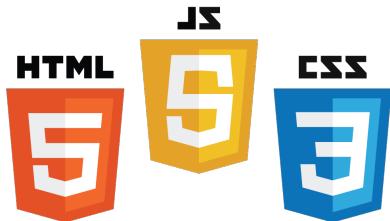


WEB DESIGN

黃郁棻

使用工具

前端網頁



視覺呈現



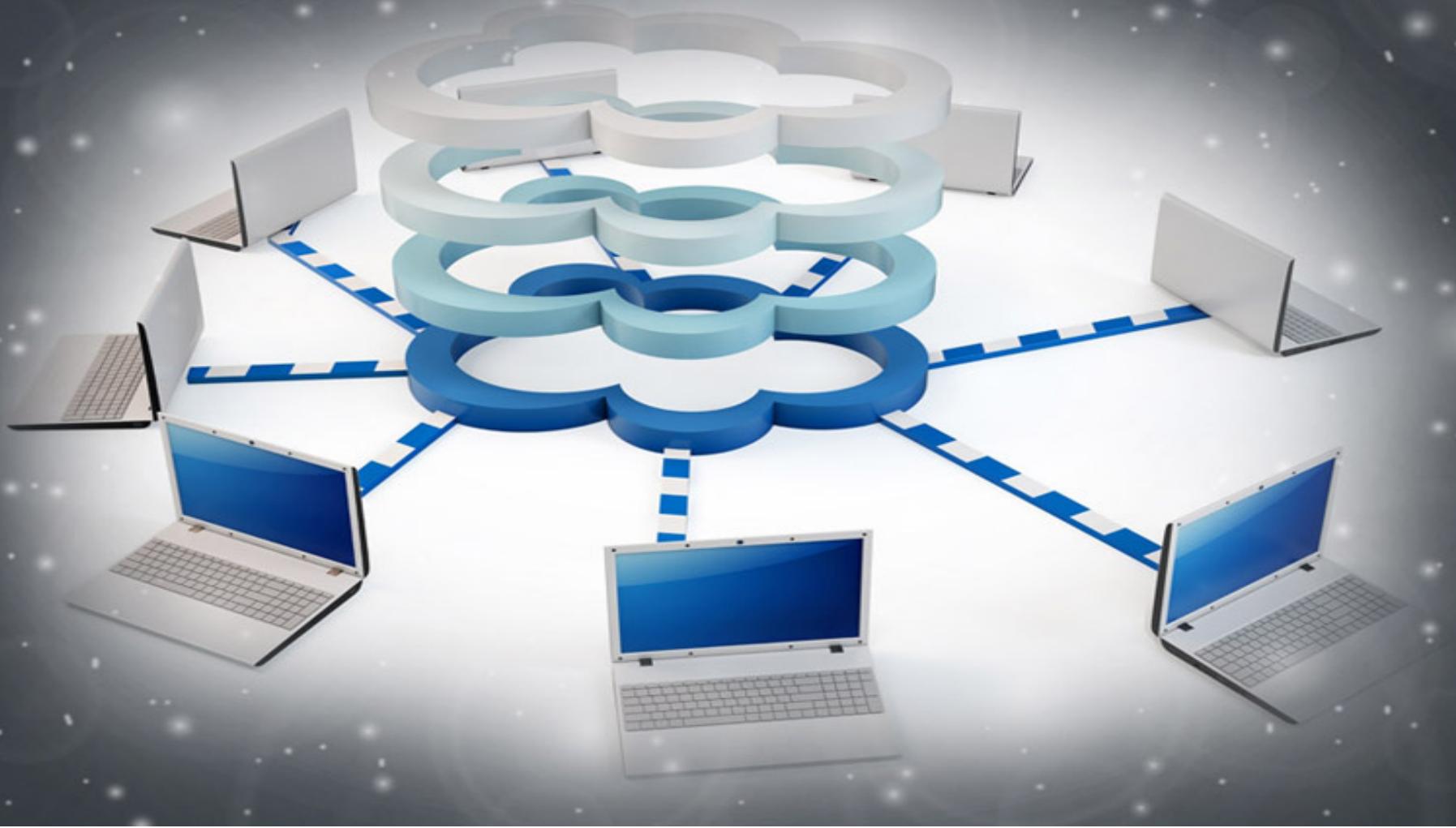
amCharts

後端連線





DEMO | 黃郁棻



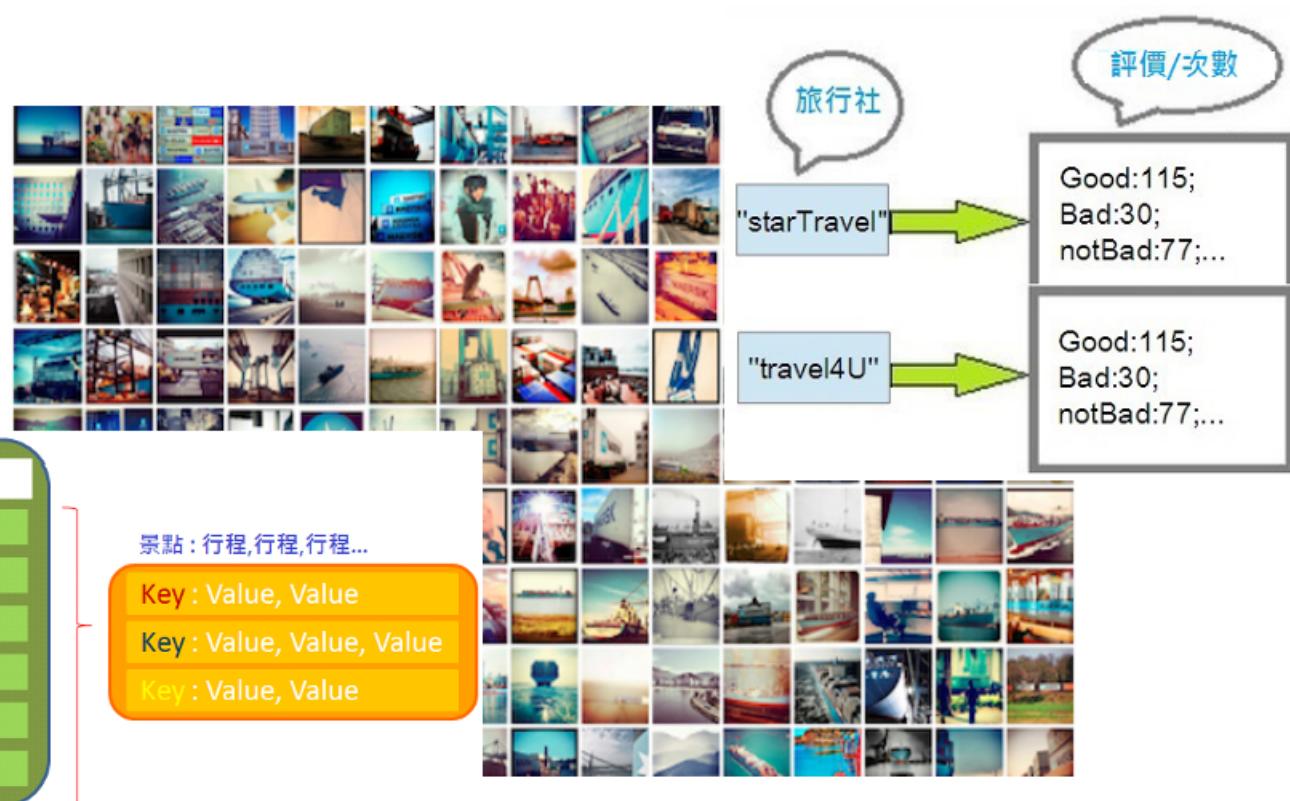
前後端連結與資料庫

王玫涵

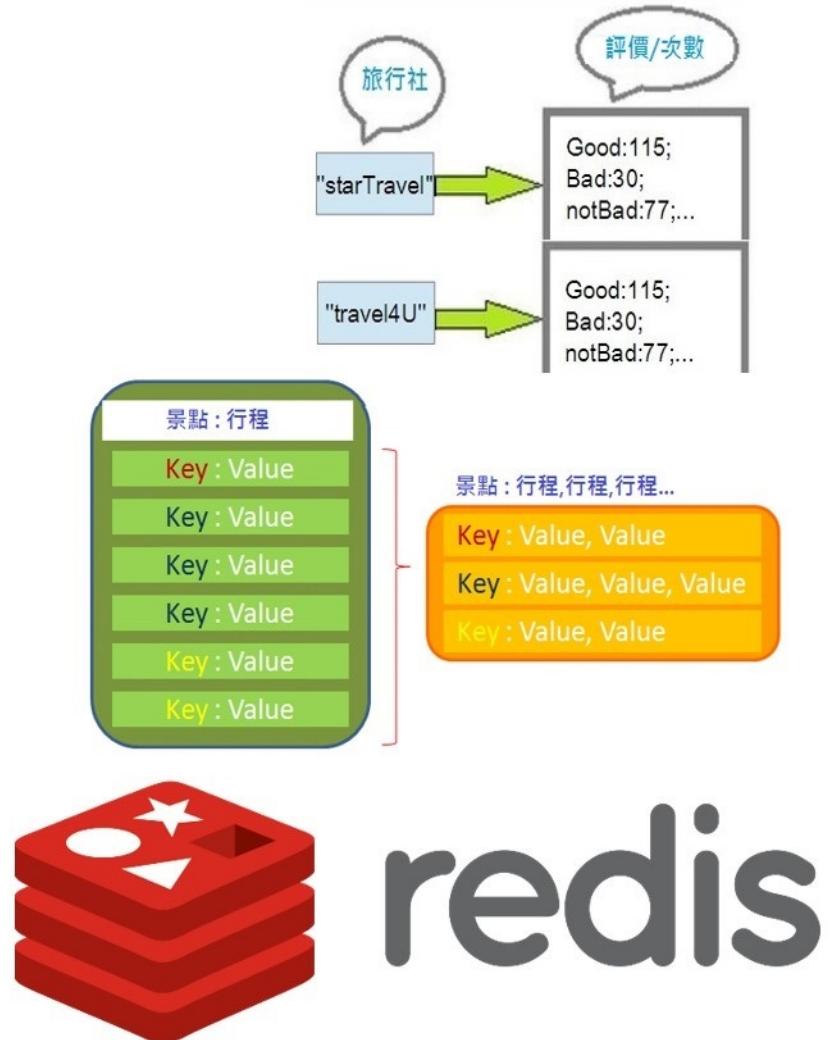
資料庫

資料型態？

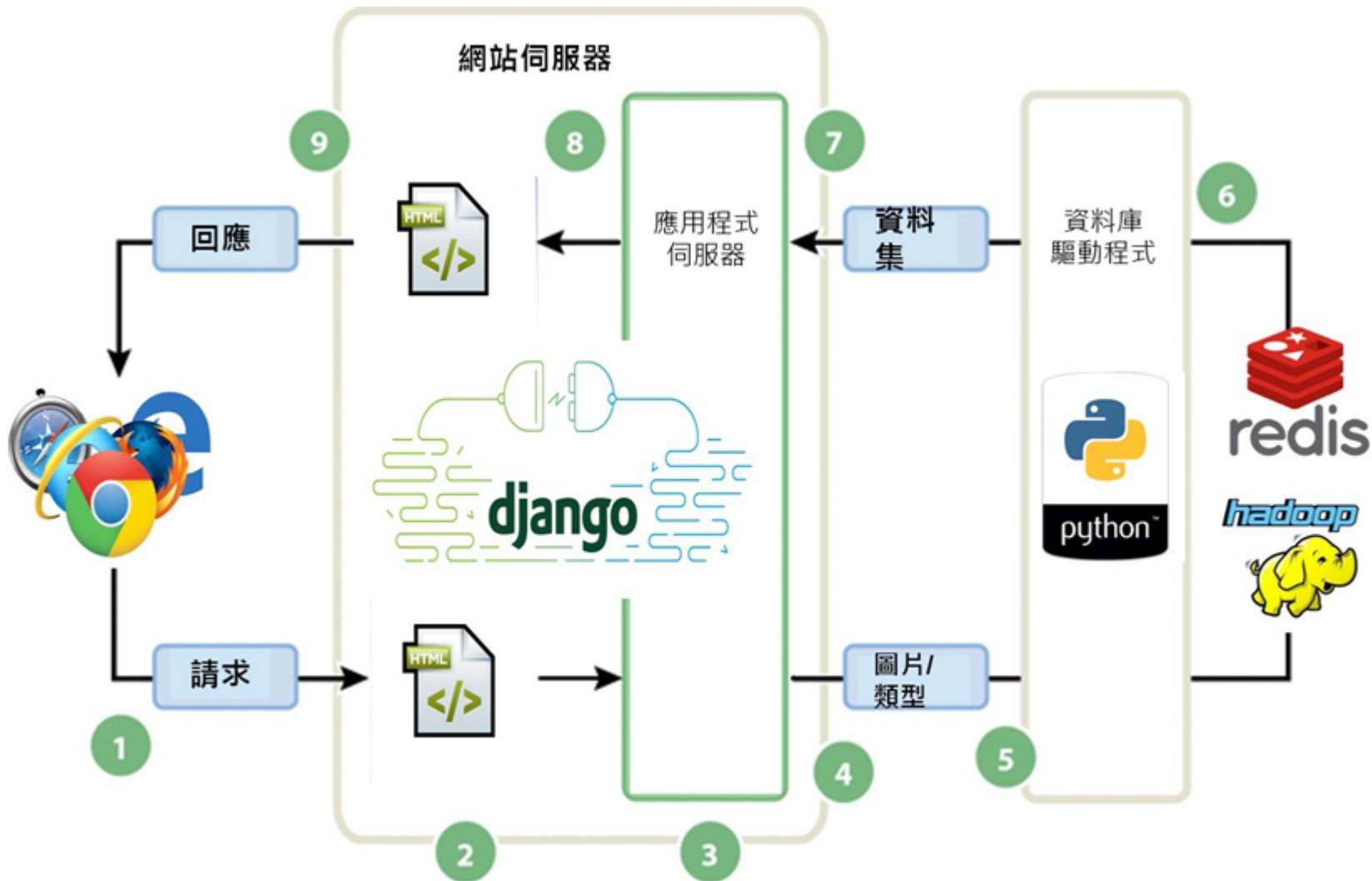
資料操作？
送出圖片,類型



資料庫



前後端連結





資料收集與清理

行程收集與景點分類

郭南宏、鄭云鈞

旅行行程—行程爬取

- ◆ BeautifulSoup
- ◆ Selenium

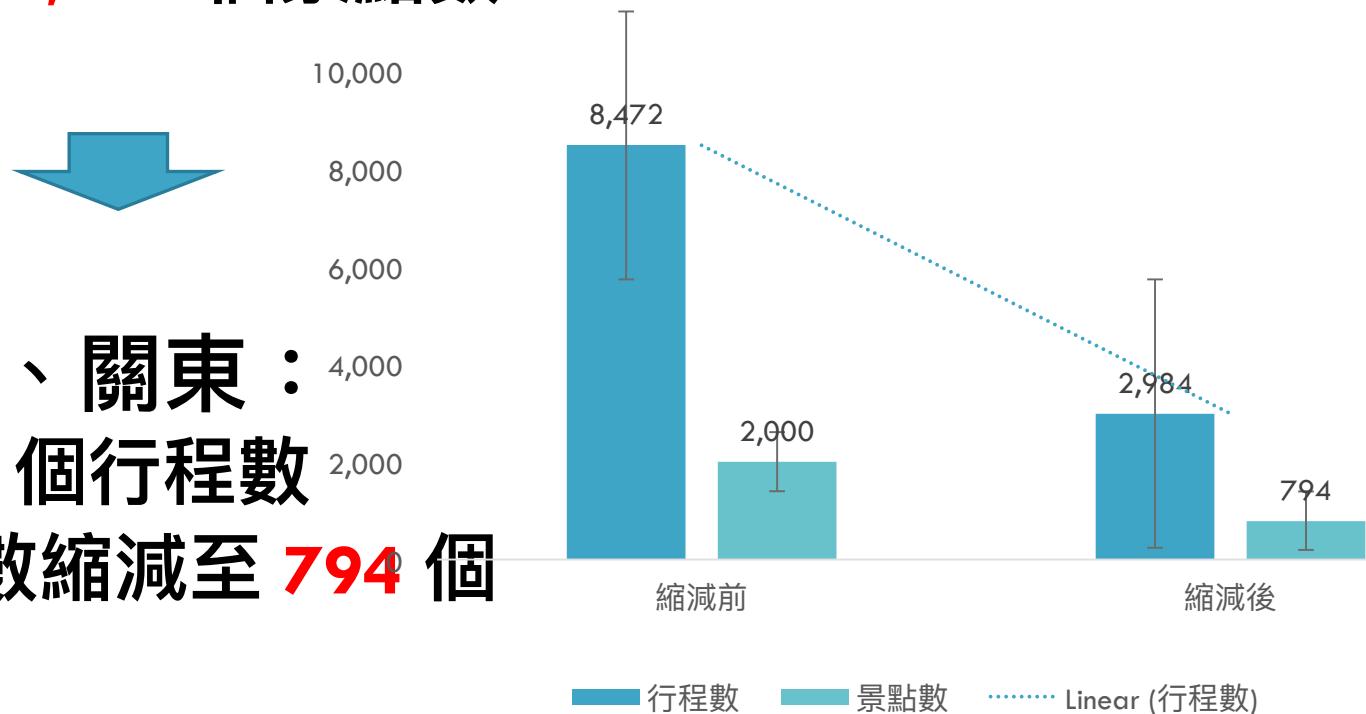


area 北海道、星野滑雪場、Club Med
href http://www.colatour.com.tw/C10A_TourSell/C10A06_TourItinerary.aspx?PatternNo=50331&GASource=國外旅遊E點通
price 35,900 趟
prodNo 50331
subtitle 北海道好個秋(函館進千歲出)～函館小樽情、海洋尼克斯、札幌摩天輪、紅葉溫泉香三日(含稅) - colatour 國外旅遊
title 北海道好個秋(函館進千歲出)～函館小樽情、海洋尼克斯、札幌摩天輪、紅葉溫泉香三日(含稅)
第01天 桃園／函館～小丘幸運漢堡(車上享用)～全森會龐群～函館百葉夜景(世界三大之夜景+纜車次回)～函館(或大沼)
第01天住宿 函館平成館溫泉旅館 或麗都皇家溫泉 或 大沼王子溫泉或大沼greenpia溫泉 或 函館國際 或 函館溫泉IMAGINE RESORT 或 函館啄木亭溫泉旅館+自助晚餐 或 同級旅館
第01天餐食 (早餐) 溫暖的家 (午餐) 機上輕食 (晚餐) 飯店內用自助晚餐 或 日式御膳海鮮料理
第02天 飯店～鶴居市～女子修道院～五稜郭城跡～大、小沼國立公園(紅葉)～洞爺湖展望台～昭和新山～熊牧場(每人發送兩顆蘋果或一袋餅乾可供餵食熊)～洞爺湖或登別溫泉(飯店內自助餐)
第02天住宿 洞爺湖觀光溫泉 或 登別PARK 雅享溫泉飯店(大廳房間皆無WIFI) 或 登別萬世閣(大廳房間皆無WIFI) 或 登別石水亭 或 同級旅館
第02天餐食 (早餐) 飯店早餐 (午餐) 日式海鮮定食 (晚餐) 飯店內用自助晚餐
第03天 飯店～登別地獄谷(紅葉)～尼克斯海洋公園(企鵝遊行)～支笏湖(紅葉)～定山溪溫泉(或KIRORO渡假村)
第03天住宿 定山溪養之湯(大廳有WIFI 房間無) 或 定山溪HOTEL(大廳有WIFI 房間無) 或 定山溪GRAND 或 KIRORO渡假村 或 同級旅館
第03天餐食 (早餐) 飯店早餐 (午餐) 日式壽喜鍋 (晚餐) 飯店內用自助晚餐 或 飯店內會席料理
第04天 飯店～壹平吹(紅葉)～免稅商店～小樽運河放浪(運河工藝館、音樂鐘館、北一硝子館)～銀之鐘咖啡座(贈送小樽風味咖啡杯及點心)～NORBESA摩天輪(眺望札幌風情)～豬肉涮涮鍋+三大盤(吃到飽)+創意自助餐+飲飲無限～札幌
第04天住宿 RESOL 札幌(大廳房間皆無WIFI) 或 札幌城市飯店 或 東急INN(大廳房間有WIFI) 或 札幌 ASPEN(大廳房間有WIFI) 或 札幌ART 或 同級旅館
第04天餐食 (早餐) 飯店早餐 (午餐) 日式石狩鍋+烤花魚料理 (晚餐) 豬肉涮涮鍋+三大盤(吃到飽)+創意自助餐+飲飲無限
第05天 飯店～寒庭溪谷(紅葉)～OUTLET購物賣場～千歲／桃園
第05天住宿 甜蜜的家
第05天餐食 (早餐) 飯店早餐 (午餐) 方便逛街~欲請自理 (晚餐) 機上輕食
PatternNo 50307
area 北海道、星野滑雪場、Club Med
href http://www.colatour.com.tw/C10A_TourSell/C10A06_TourItinerary.aspx?PatternNo=50307&GASource=國外旅遊E點通
price 35,900 趟
prodNo 50331
subtitle 北海道好個秋(千歲進函館出)～函館小樽情、海洋尼克斯、札幌摩天輪、紅葉溫泉香三日(含稅) - colatour 國外旅遊
title 北海道好個秋(千歲進函館出)～函館小樽情、海洋尼克斯、札幌摩天輪、紅葉溫泉香三日(含稅)
第01天 桃園／千歲～寒庭溪谷(紅葉)～支笏湖(紅葉)～定山溪溫泉(或KIRORO渡假村)
第01天住宿 定山溪養之湯(大廳有WIFI 房間無) 或 定山溪HOTEL(大廳有WIFI 房間無) 或 定山溪GRAND 或 KIRORO渡假村 或 同級旅館
第01天餐食 (早餐) 溫暖的家 (午餐) 機上輕食 (晚餐) 飯店內用自助晚餐 或 飯店內會席料理

旅行行程—行程爬取

- ◆ 原行程數量：8,472 個行程
- ◆ 超過 2,000 個景點數

- ◆ 關西、關東：
- ◆ 2,984 個行程數
- ◆ 景點數縮減至 794 個



旅行行程—資料清理

●遇到問題

桃園機場或松山機場／東京成田空港或東京羽田空港 → 飯店 → 成田(購物廣場~自由購物) 或 東京或千葉(自由活動~逛街購物)
(限航班為午前出發)箱根蘆之湖復古海盜船 · 遊覽富士山湖光山色 · 關東古老神社之一 · 箱根神社 → 御殿場OUTLET ·
人氣品牌匯集~可眺望富士山 → 富士五湖溫泉鄉或石和溫泉鄉日本動畫大師 · 宮崎駿 · 三鷹之森吉卜力美術館 → 淺草觀音寺 ·
江戶繁華老街(仲見世通) · 可遠眺東京晴空塔 → 免稅店 → 台場繽紛城 · DiverCity Tokyo Plaza~巨大鋼彈戰士模型全日迪士尼樂園
或 海洋迪士尼(一票到底)~搭乘電車體驗（選擇樂園或海洋，請於出發前七個工作日告知客服人員，未事先告知的貴賓，
視同以樂園為主）自由活動~飯店→東京成田空港或東京羽田空港／桃園機場或松山機場

●每個行程區分景點、說明的表點符號不一定一樣，本專題只要留景點名稱

旅行行程—資料清理

原始行程
資料



行程資料
前處理

台北松山機場東京羽田空港→東京自由活動 ·
建議您可自費搭乘電車前往(1)有樂町、銀座時尚百貨名店街
(2)新宿 大型知名日系百貨、美食餐飲店以及歌舞伎町 (3)橫濱未來21～紅磚赤瓦倉庫群等
淺草觀音寺·雷門～仲見世通→東京晴空塔～
晴空街道商店街→★特別安排【宮崎駿·三鷹之森吉普力美術館】→富士五湖區
日本百選名水～國家天然紀念物〈忍野八海〉→藝猴雜技表演～河口湖猿劇場→
路經彩虹大橋→免稅店→台場DiverCity Tokyo Plaza～超大鋼彈模型立像
全日 東京迪士尼樂園 或 迪士尼海洋(盡情暢遊一票到底) ～搭捷運電車、地下鐵或專車接送
自由活動 · 建議您可自費搭乘電車前往(1)築地場外市場
(2)上野阿美橫町傳統商店街→東京羽田空港台北松山機場

台北松山機場東京羽田空港
東京自由活動 · 建議您可自費搭乘電車前往(1)有樂町、銀座時尚百貨名店街
自由活動(2)新宿 大型知名日系百貨、美食餐飲店以及歌舞伎町 (3)橫濱未來21～紅磚赤瓦倉庫群等
淺草觀音寺·雷門～
仲見世通
東京晴空塔～
晴空街道商店街
★特別安排【宮崎駿·三鷹之森吉普力美術館】
富士五湖區
日本百選名水～國家天然紀念物〈忍野八海〉
藝猴雜技表演～河口湖猿劇場
路經彩虹大橋
免稅店
台場DiverCity Tokyo Plaza～超大鋼彈模型立像
全日 東京迪士尼樂園 或 迪士尼海洋(盡情暢遊一票到底) ～搭捷運電車、地下鐵或專車接送
自由活動 · 建議您可自費搭乘電車前往(1)築地場外市場 (2)上野阿美橫町傳統商店街
東京羽田空港台北松山機場

旅行行程—資料清理

利用Python進行去除不必要資料及景點名字轉換

```
import os

for root, dirs, files in os.walk('attraction\\'):
    for f in files:

        a = os.path.join(root, f)
        title = a.split('\\')[1]
        no = a.split('\\')[2]

        with open(a, 'r') as r:
            dir_name = "tourname\\" + title
            if not os.path.exists(dir_name):      #先確認資料夾是否存在
                os.makedirs(dir_name)
            file = open(dir_name + "\\\" + no, "w")
            try:
                for i in r:
                    ele = i.decode('utf-8')
                    if u'機場' in ele:
                        continue
                    elif u'自由活動' in ele:
                        continue
                    else:
                        m = changename(ele)
                        file.write(m.encode('utf-8'))
            except:
                print a
                continue
            file.close()
    print "done"

def changename(i):
    dictour = {
        u'淺草觀音寺':u'Senso_ji_Temple',
        u'仲見世':u'Shin_Nakamise_Shopping_Street',
        u'晴空街道商店街':u'Tokyo_Solamachi',
        u'晴空塔':u'Tokyo_Skytree',
        u'吉普力美術館':u'Ghibli_Museum',
        u'富士五湖':u'Notenburo_Tensui',
        u'忍野八海':u'Oshino_Hakkai',
        u'河口湖':u'Kawaguchiko_Lake',
        u'彩虹大橋':u'Rainbow_Bridge',
        u'免稅':u'Aeon_Mall_Ayagawa',
        u'DiverCity':u'Diver_City_Tokyo_Plaza',
        u'迪士尼':u'Tokyo_Disneyland'
    }

    for ele in dictour:
        if ele in i:
            return dictour[ele] + '\n'
            break
    return i
```

旅行行程—資料清理

已整理
資料

Senso_ji_Temple
Shin_Nakamise_Shopping_Street
Tokyo_Skytree
Tokyo_Solamachi
Ghibli_Museum
Notenburo_Tensui
Oshino_Hakkai
Kawaguchiko_Lake
Rainbow_Bridge
Aeon_Mall_Ayagawa
Diver_City_Tokyo_Plaza
Tokyo_Disneyland

旅行行程一分數制定

● 拉手分數

● 計算每個景點在旅行社行程出現的熱門度

```
import os
score = {}
dicscore = {}
num = 0
for root, dirs, files in os.walk('tournname\\'):
    for f in files:
        a = os.path.join(root, f)
        title = a.split('\\')[1]
        no = a.split('\\')[2]
        with open(a, 'r') as r:
            num += 1
            for ele in r:
                b = repr(ele).strip("'").split('\\')[0]
                b = b.decode('utf=8')
                if b not in score:
                    score[b] = 1
                else:
                    score[b] += 1
```

```
[('Aeon_Mall_Ayagawa', 1932), ('Tokyo_Disneyland', 497), ('Sannenzaka_Ninenzaka', 494), ('Kiyomizu_dera_Temple', 486), ('Togetsukyo', 441), ('Shinsaibashi', 429), ('Senso_ji_Temple', 422), ('Tokyo_Skytree', 403), ('Osaka_Castle', 386), ('Dotonbori', 360), ('Diver_City_Tokyo_Plaza', 353), ('Universal_Studios_Japan', 346), ('Hakone_Pirate_Ship', 342), ('Todai_ji_Temple', 321), ('Fushimi_Inari_taisha_Shrine', 318), ('Shin_Nakamise_Shopping_Street', 316), ('Sagano', 305), ('Nara_Park', 292), ('Hakone_Shrine_Kuzuryu_Shrine_Singu', 254), ('Otonashi_Falls', 227), ('Fuji_Hakone_Izu_National_Park', 223), ('Nonomiya_Shrine', 215), ('Rokuon_ji_Temple', 210), ('Akashi_Kaikyo_Bridge', 204), ('Oshino_Hakkai', 200), ('Machikado_Amusement_Park_Mosaic_Garden', 195), ('Rainbow_Bridge', 192), ('Hanamikoji_Street', 186), ('Jisyu_Shrine', 184), ('Maiko_Park', 168), ('Meiji_Jingu', 156), ('Gion', 151), ('Sagano_Romantic_Train', 149), ('Isawa_Onsen', 145), ('Kobe_Kitano_Museum_Ijinkan_gai', 138), ('Yasaka_Shrine', 119), ('Arima_Onsen', 117), ('Notenburo_Tensui', 114), ('Tokyo_Solamachi', 111), ('Naritasan_Shinshoji_Temple', 98), ('Yokohama_Red_Brick_Warehouse', 94), ('St_Paul_s_Catholic_Church_Karuizawa', 90), ('Byodoin_Temple', 89),
```

● 計算每個景點在所有行程出現的次數

```
import os
dicnum = {}
num = 0
for root, dirs, files in os.walk('tournname\\'):
    for f in files:
        dic = {}
        a = os.path.join(root, f)
        title = a.split('\\')[1]
        no = a.split('\\')[2].split('.txt')[0]
        no = no.decode('utf-8')
        dic['agency'] = title
        with open(a, 'r') as r:
            list = []
            num += 1
            for ele in r:
                for i in score:
                    b = repr(ele).strip("'").split('\\')[0]
                    b = b.decode('utf=8')
                    if b == i:
                        dic[b] = score[i]
            dicnum[no] = dic
```

```
{"16JT731B": {"Isawa_Onsen": 145, "Aeon_Mall_Ayagawa": 1932, "Shin_Nakamise_Shopping_Street": 316, "agency": "lionTravel", "Oshino_Hakkai": 200, "Tokyo_Solamachi": 111, "Meiji_Jingu": 156, "Hakone_Pirate_Ship": 342, "Tokyo_Skytree": 403, "Senso_ji_Temple": 422, "Rainbow_Bridge": 192, "Hakone_Shrine_Kuzuryu_Shrine_Singu": 254, "Tokyo_Disneyland": 497}, "PACTT-260 CX-TT160724A": {"Hakone_Kowakien_Yunessun": 17, "Shuzenji_Onsen": 25, "Minato_Mirai_21": 63, "agency": "pacificTour", "Byodoin_Omotesando": 32, "Tsurugaoka_Hachimangu_Shrine": 29, "Enoshima_Island": 13, "The_Hakone_Open_Air_Museum": 10, "Izu_Kyoko_Line": 25, "Aeon_Mall_Ayagawa": 1932}, "VDR0000001943_OSA05CI0611B": {"Daiganbo": 17, "Kanazawa_Castle": 2, "World_Heritage_Shirakawa_go_Gassho_Frame_Housing_Community": 21, "Midosuji_Street": 78, "Mikurigaike_Pond": 17, "Sanmachi_Suji": 16, "Kappa_Bashi": 12, "Yuki_no_Otani": 1, "Tsuzumi_Gate": 1, "Kenrokuen_Garden": 14, "Kurobe_Dam": 17, "Aeon_Mall_Ayagawa": 1932, "Kurokabe_Square": 1, "agency": "eztravel", "Bijodaira": 12, "Rinku_Premium_Outlets": 18, "Pond_Taisho": 10, "Higashichaya_Old_Town": 7, "Yamanaka_Onsen": 3, "Dotonbori": 360, "Kamikochi": 17, "Shinsaibashi": 429, "Tojinbo_Cliff": 4, "Miho_Museum": 10}, "TY005CI160724G": {"Isawa_Onsen": 145, "Aeon_Mall_Ayagawa": 1932, "Shin_Nakamise_Shopping_Street": 316, "agency": "bestTour", "Fuji_Hakone_Izu_National_Park": 223, "Oshino_Hakkai": 200, "Yokohama_Hakkeijima_Sea_Paradise": 42, "Meiji_Jingu": 156, "Hakone_Pirate_Ship": 342, "Tokyo_Skytree": 403, "Senso_ji_Temple": 422, "Hakone_Shrine_Kuzuryu_Shrine_Singu": 254, "Tokyo_Disneyland": 497}, "OSA04CI16830
```

旅行行程一分數制定

- 拉手分數
- 行程中每個景點的分數加總當作拉手分數

```
import os
dictotal = {}
num = 0
for root, dirs, files in os.walk('tourname\\'):
    for f in files:
        dic = {}
        a = os.path.join(root, f)
        title = a.split('\\')[1]
        no = a.split('\\')[2].split('.txt')[0]
        no = no.decode('utf-8')
        dic[u'agency'] = title
        with open (a, 'r') as r:
            num = 0
            for ele in r:
                for i in score:
                    b = repr(ele).strip("").split('\\')[0]
                    b = b.decode('utf=8')
                    if b == i:
                        num += score[i]
        dictotal[no] = num
```

```
{"16JT731BR": 6902, "PACTT-260 CX-TT160724A": 2171, "VDR000001943_OSA05CI0611B": 2999, "TY005CI160724G": 4932, "OSA04CI16830Y": 3599, "FRN0000010835": 4686, "0BA3FED93711B971781D8FFBDEF22268": 3770, "OSA05BR160802HN": 11714, "16JX725BRR": 4987, "VDR0000001843_OSAGE16060305A": 5804, "16JX725BRB": 4800, "OSA05160828TB": 6472, "48893": 6918, "16JX725BR4": 5502, "VDR0000007614_JTP100160603IT": 6359, "GFG0000004368": 7165, "GFG0000003409": 6070, "PACTT-260 NH-TT160807A": 2112, "FRN0000014116": 8885, "PACTT-222
```

旅行行程一分數制定

● 類型分數

● 使用TripAdvisor的景點所屬類型

```
url = 'https://www.tripadvisor.com.tw/Attractions-g294232-Activities-Japan.html'
import requests
from bs4 import BeautifulSoup as bs
import time, re
num = 0
rs = requests.session()
res = rs.get(url)
soup = bs(res.text, 'lxml')
for k in soup.select('div#CHILD_GEO_FILTER div.filter_list.al_border div.filter')[0:-1]:
    a = 'https://www.tripadvisor.com.tw' + k.select('a')[0]['href']
    res3 = rs.get(a)
    soup3 = bs(res3.text, 'lxml')
    if len(soup3.select('div.pageNumbers a')) > 0:
        if re.search(u'Activities', a) is not None:
            ccc = re.sub(u'Activities', u'Activities-oa{}', a)
            pagenum = int(soup3.select('div.pageNumbers a')[-1]['data-page-number'])

for b in range(0, pagenum):
    res4 = rs.get(ccc.format(b*30))
    soup4 = bs(res4.text, 'lxml')
    try:
        for ele in soup4.select('div#FILTERED_LIST'):
            for i in ele.select('div.element_wrap'):
                if len(i.select('div.pi3n_reasoning_v2')) > 0:
                    for m in i.select('div.pi3n_reasoning_v2 a'):
                        dickobe.setdefault(m.select('div.property_title a')[0]['href'].split('-')[4], []).append(m.text)
                    num += 1
            else:
                f = 'https://www.tripadvisor.com.tw' + i.select('div.property_title a')[0]['href']
                res5 = rs.get(f)
                soup5 = bs(res5.text, 'lxml')
                if len(soup5.select('div.pageNumbers a')) == 0:
                    for g in soup5.select('div#FILTERED_LIST'):
                        for h in g.select('div.element_wrap'):
                            for j in h.select('div.pi3n_reasoning_v2'):
                                for p in j.select('span.matchedTag.noTagImg'):
                                    dickobe.setdefault(h.select('div.property_title a')[0]['href'].split('-')[4], []).append(p.text)
                else:
                    for q in range(0, pagenum):
                        res6 = rs.get(ddd.format(q*30))
                        soup6 = bs(res6.text, 'lxml')
                        for r in soup6.select('div#FILTERED_LIST'):
                            for s in r.select('div.element_wrap'):
                                for t in s.select('div.pi3n_reasoning_v2'):
                                    for u in t.select('span.matchedTag.noTagImg'):
                                        dickobe.setdefault(s.select('div.property_title a')[0]['href'].split('-')[4], []).append(u.text)
                            time.sleep(2)
                        except:
                            print ccc.format(b*30)
                else:
                    for c in soup3.select('div#FILTERED_LIST'):
                        for d in c.select('div.element_wrap'):
                            if len(d.select('div.pi3n_reasoning_v2')) > 0:
                                for e in d.select('div.pi3n_reasoning_v2 a'):
                                    dickobe.setdefault(d.select('div.property_title a')[0]['href'].split('-')[4], []).append(e.text)
                            time.sleep(2)
            time_end= time.time()
```

旅行行程一分數制定

- 類型分數
- 將每個行程景點轉換為TripAdvisor的景點類型

```
import os
dictype = {}
num = 0
for root, dirs, files in os.walk('tourname\\'):
    for f in files:
        dic = {}
        a = os.path.join(root, f)
        title = a.split('\\')[1]
        no = a.split('\\')[2].split('.txt')[0]
        no = no.decode('utf-8')
        dic[u'agency'] = title
        with open(a, 'r') as r:
            num = 0
            add = {}
            type = []
            for ele in r:
                for i in dicname:
                    b = repr(ele).strip("").split('\\')[0]
                    b = b.decode('utf=8')
                    if b == i:
                        add[b] = dicname[i]
            dictype[no] = add
```

```
{"16JT731BR": {"Isawa_Onsen": ["Spa 與養生"], "Aeon_Mall_Ayagawa": ["購物"], "Shin_Nakamise_Shopping_Street": ["購物"], "Oshino_Hakkai": ["景點與地標"], "Tokyo_Solamachi": ["購物"], "Meiji_Jingu": ["景點與地標"], "Hakone_Pirate_Ship": ["遊覽", "戶外活動", "搭船遊覽與水上運動"], "Tokyo_Skytree": ["景點與地標"], "Senso_ji_Temple": ["景點與地標"], "Rainbow_Bridge": ["景點與地標"], "Hakone_Shrine_Kuzuryu_Shrine_Singū": ["景點與地標"], "Tokyo_Disneyland": ["水上樂園和遊樂場"]}, "PACTT-260 CX-TT160724A": {"Hakone_Kowakien
```

旅行行程一分數制定

● 類型分數

● 計算每個行程各類型的個數

```
import os
dicscoresum = {}
for root, dirs, files in os.walk('tourname\\'):
    for f in files:
        dic = {}
        num = 0
        a = os.path.join(root, f)
        title = a.split('\\')[1]
        no = a.split('\\')[2].split('.txt')[0]
        no = no.decode('utf-8')
        with open(a, 'r') as r:
            for ele in r:
                b = repr(ele).strip("").split('\\')[0]
                b = b.decode('utf-8')
                for i in dicname:
                    if b == i:
                        for m in range(0, len(dicname[i])):
                            if dicname[i][m] not in dic:
                                dic[dicname[i][m]] = 1
                            else:
                                dic[dicname[i][m]] += 1
                for n in score:
                    if b == n:
                        num += score[n]
        dic[u'score'] = num

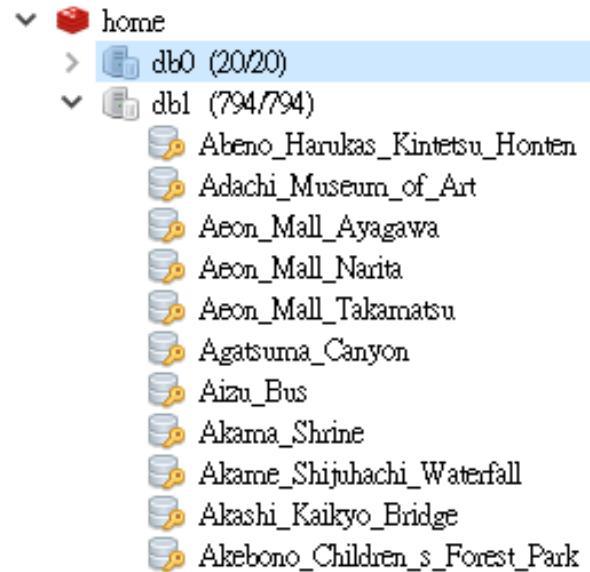
dicscoresum[no] = dic
```

```
{"16JT731BR": {"購物": 4, "遊覽": 1, "搭船遊覽與水上運動": 1, "景點與地標": 6, "Spa 與養生": 1, "水上樂園和遊樂場": 1, "戶外活動": 1, "score": 6902}, "PACTT-260 CX-TT160724A": {"購物": 1, "自然與公園": 1, "博物館": 1, "景點與地標": 3, "交通": 2, "Spa 與養生": 2, "score": 2171, "旅客資源": 1, "遊覽": 1}, "VDR0000001943_OSA05CI0611B": {"購物": 3, "自然與公園": 7, "博物館": 2, "景點與地標": 1, "Spa 與養生": 1, "score": 2999, "戶外活動": 2}, "TY005CI160724G": {"購物": 2, "自然與公園": 1, "動物園和水族館": 1, "搭船遊覽與水上運動": 1, "景點與地標": 5, "Spa 與養生": 1, "水上樂園和遊樂場": 1, "戶外活動": 1, "score": 4932, "遊覽": 1}, "OSA04CI6830Y": {"購
```

旅行行程—匯入資料庫

● 特定景點在那些行程出現

```
import os
dicattr = {}
listattr = []
num = 0
for root, dirs, files in os.walk('tourname\\'):
    for f in files:
        dic = {}
        a = os.path.join(root, f)
        title = a.split('\\')[1]
        no = a.split('\\')[2].split('.txt')[0]
        no = no.decode('utf-8')
        dic[u'agency'] = title
        with open(a, 'r') as r:
            num = 0
            add = {}
            type = []
            for ele in r:
                for i in dicname:
                    b = repr(ele).strip("").split('\\')[0]
                    b = b.decode('utf-8')
                    if b == i:
                        dicattr.setdefault(i, []).append(no)
```



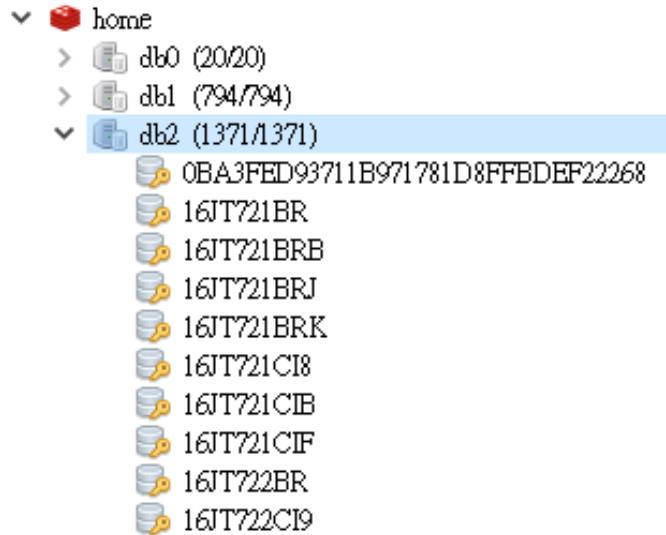
```
{"Korakukan": ["VDR0000001888_SDJ05160702A"], "Marukin_Soysouce_Memorial": ["OSA05CI160726BP", "OSA05CI160802DP", "OSA05CI160814AP", "OSA05CI160816BP", "OSA05CI160820BP", "OSA05CI160904AP", "VDR0000001888_NPD061705CI6E", "VDR0000001888_HIJ05160614C", "VDR000001888_HIJ06160612B", "VDR0000001888_OSA05160604Y", "VDR0000001888_OSA05160614H", "VDR0000001888_OSA05160701N", "VDR0000001975_OSA05160628TS", "JHIJ5-CIK3", "GFG0000009417", "GFG0000009418"], "Hozugawa_River_Boat_Ride": ["VDR0000001846_16JX604CIH", "VDR0000001846_16JX629BRC", "VDR0000001943_OSA05BR6721J", "JOSBA5", "OSA05B6828J", "OSA05BR6725J", "OSA05BR6822J", "OSA05BR6904E", "OSA05BR6909E", "16JX724CIC", "16JX727BRC", "16JX731CIC", "16JX801BRH"], "Kotobukiyo": ["VDR0000001888_HIJ06160609L"], "Kasamat su_Park": ["VDR0000001843_OSAGE16060305A", "VDR0000001888_OSA05160616V", "VDR0000001943_OSA05BR6721J", "VDR0000001975_OSA05160703TC", "534", "OSA05B6828J", "OSA05BR6725J", "OSA05BR6822J", "OSA05BR6904E", "OSA05BR6909E", "OSA05CI6907J", "16JX726CI9", "GFG0000002670", "GFG0000007050", "OSA05CI132", "OSA05CI1406", "OSA05CI16D"], "Goryokaku_Park": ["VDR0000001888_SDJ05160812B"], "SHIS
```

旅行行程—匯入資料庫

● 行程所有類型、拉手分數、旅行社、旅行社團號

```
import os
dictotalsum = {}
for root, dirs, files in os.walk('tourname\\'):
    for f in files:
        dic = {}
        list = []
        num = 0
        a = os.path.join(root, f)
        title = a.split('\\')[1]
        no = a.split('\\')[2].split('.txt')[0]
        no = no.decode('utf-8')
        dic['tourno'] = no
        dic['agency'] = title
        with open(a, 'r') as r:
            b = repr(ele).strip("").split('\\')[0]
            b = b.decode('utf=8')
            list.append(b)
            for i in dicname:
                if b == i:
                    for m in range(0, len(dicname[i])):
                        if dicname[i][m] not in dic:
                            dic[dicname[i][m]] = 1
                        else:
                            dic[dicname[i][m]] += 1
            for p in name:
                if p not in dic:
                    dic[p] = 0
            for n in score:
                if b == n:
                    num += score[n]
        dic['score'] = num
        dictotalsum[no] = dic
```

```
{"16JT731BR": {"課程活動及工作坊": 0, "搭船遊覽與水上運動": 1, "購物": 4, "自然與公園": 0, "score": 6902, "博物館": 0, "動物園和水族館": 0, "agency": "lionTravel", "景點與地標": 6, "tourno": "16JT731BR", "演唱會和表演": 0, "交通": 0, "Spa 與養生": 1, "tour": ["Aeon_Mall_Ayagawa", "Oshino_Hakkai", "Hakone_Pirate_Ship", "Hakone_Shrine_Kuzuryu_Shrine_Singu", "Isawa_Onsen", "Meiji_Jingu", "Rainbow_Bridge", "Aeon_Mall_Ayagawa", "Senso_ji_Temple", "Shin_Nakamise_Shopping_Street", "Tokyo_Skytree", "Tokyo_Solamachi", "Tokyo_Disneyland"], "水上樂園和遊樂場": 1, "戶外活動": 1, "旅客資訊": 0, "飲食": 0, "休閒與娛樂": 0, "遊覽": 1}, "PACTT-260 CX-TT1607
```





資料收集與清理

圖片收集與清理

黃郁棻

圖片來源：取得圖片URL

- ◆ 東京、京都、大阪、神戶共**794個**景點
- ◆ Google 圖片搜尋：圖檔大小共約**77G**

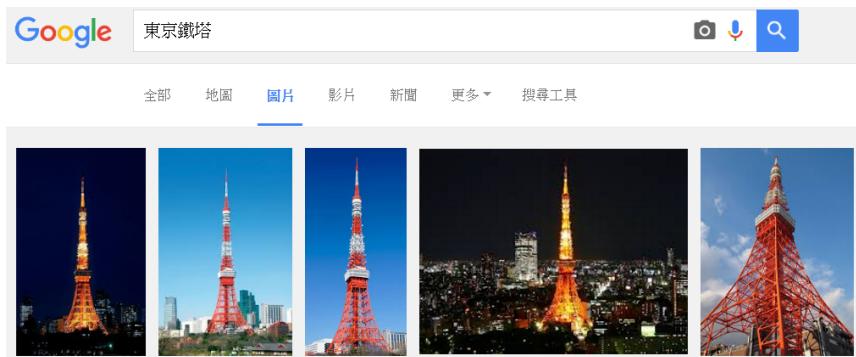


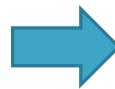
Image URL List

http://www.gotokyo.org/tc/kanko/minato/spot/images/02_12.jpg
http://www.gotokyo.org/tc/kanko/minato/spot/Images/01_12.jpg
<https://upload.wikimedia.org/wikipedia/commons/e/ed/TaroTokyo20110213-TokyoTower.jpg>
http://travelheadline.com/wp-content/uploads/2014/01/241460_12984258180chX.jpg
http://travelheadline.com/wp-content/uploads/2014/01/1_edited1.jpg
<http://pic.pimg.tw/handazuku/1382187426-2482834366.jpg>
https://fastjapan.com/en/wp-content/uploads/2016/01/shutterstock_129611225.jpg
http://b.share.photo.xuite.net/twgoq07/1be6c1b/9757958/432214805_m.jpg
http://www.wall001.com/human/tokyo-tower_01/mxxx01/t-tower-228.jpg
http://farm5.static.flickr.com/4048/4646822023_4804fd47ea_o.jpg
http://sib.tw/article_vv/2015042901/007.jpg
http://farm5.static.flickr.com/401/4651115794_689bd15455_o.jpg
http://www.backpackers.com.tw/forum/gallery/images/97962/1_DSC05922.JPG
http://3.share.photo.xuite.net/yolmg/13d7ece/4671419/194688019_m.jpg
<http://pic.pimg.tw/handazuku/1382187427-699339213.jpg>
http://farm9.static.flickr.com/8037/8032928423_ac7e4deee2_o.jpg
http://farm9.static.flickr.com/8458/8035512745_f2e9fb2f3_o.jpg
<http://cdn2.ettoday.net/images/592/d592675.jpg>
https://farm2.staticflickr.com/1632/25761356326_cfa49d2d19_b.jpg
<http://pic.pimg.tw/pella/1338436610-3571660953.jpg>
http://farm6.static.flickr.com/5018/5428439899_0f9fda1ecc_o.jpg
http://img103.myspd.com.cn/20130108/1/Myspd_22847_201301081645010004B.jpg
http://pic.pimg.tw/r50130/1430743787-2385792713_n.jpg
<http://pic.pimg.tw/uniquevera/4bd048f6d4c0d.jpg>
<https://i.ytimg.com/vi/apXtvJvP9Lo/maxresdefault.jpg>
<http://pic.pimg.tw/barbara0726/1366633624-2309163822.jpg>
http://www.backpackers.com.tw/forum/gallery/images/217553/1_東京鐵塔S130.jpg

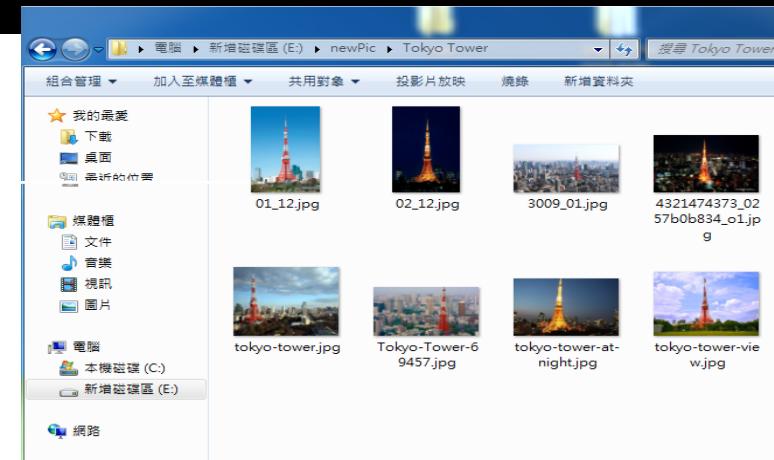
圖片下載與清理

Image URL List

```
http://www.gotoko.org/tc/kanko/minato/spot/images/02_12.jpg  
http://www.gotoko.org/tc/kanko/minato/spot/images/01_12.jpg  
https://upload.wikimedia.org/wikipedia/commons/e/ed/TaroTokyo20110213-TokyoTower  
http://travelheadline.com/wp-content/uploads/2014/01/241460_12984258180chX.jpg  
http://travelheadline.com/wp-content/uploads/2014/01/1_edited1.jpg  
http://pic.pimg.tw/handazyuku/1382187426-2482854366.jpg  
https://fastjapan.com/en/wp-content/uploads/2016/01/shutterstock_129611225.jpg  
http://b.share.photo.xuite.net/twggoo7/1be6c1b/9757958/432214805_m.jpg  
http://www.wall001.com/human/tokyo-tower_01/mxxx01/t-tower-228.jpg  
http://farm5.static.flickr.com/4048/4646822023_4804fd47ea_o.jpg  
http://sib.tw/article_vw/2015042901/007.jpg  
http://farm5.static.flickr.com/4017/4651115794_689bd15455_o.jpg  
http://www.backpackers.com.tw/forum/gallery/images/97962/1_DSC05922.JPG  
http://3.share.photo.xuite.net/yoimg/13d7ece/4671419/194688019_m.jpg  
http://pic.pimg.tw/handazyuku/1382187427-699339213.jpg  
http://farm9.static.flickr.com/8037/8032928423_ac7e4dee2_o.jpg  
http://images2.gamme.com.tw/news2/2012/52/23/p5_Wn6Scj5_Z.jpg  
http://farm9.static.flickr.com/8458/8035512745_f2e9fb2f3_o.jpg  
http://cdn2.ettoday.net/images/592/d592675.jpg  
https://farm2.staticflickr.com/1632/25761356326_cfa49d2d19_b.jpg  
http://pic.pimg.tw/joella/1338436610-3571660953.jpg  
http://farm6.static.flickr.com/5018/5428439899_0f9fd1ecc_o.jpg  
http://img103.mypsd.com.cn/20130108/1/Mypsd_22847_201301081645010004B.jpg  
http://pic.pimg.tw/ra50130/1430743787-2385792713_n.jpg  
http://pic.pimg.tw/uniquevera/4bd048f6d4c0d.jpg  
https://i.ytimg.com/vi/apXtvJvF9Lo/maxresdefault.jpg  
http://pic.pimg.tw/barbara0726/1366633624-2309163822.jpg  
http://www.backpackers.com.tw/forum/gallery/images/217553/1_東京鐵塔S130.jpg
```



```
命令提示字元 - python auto-picDown.py  
Microsoft Windows [版本 6.1.7601]  
Copyright (c) 2009 Microsoft Corporation. All rights reserved.  
  
C:\Users\BIG DATA>cd /d e:\  
  
e:\>python auto-picDown.py  
http://www.gotoko.org/tc/kanko/minato/spot/images/01_12.jpg  
http://www.japan-guide.com/g9/3009_01.jpg  
http://yokosojapan.co.jp/wp/wp-content/uploads/2015/12/tokyo-tower.jpg  
http://cdn.tokyotimes.com/wp-content/uploads/2010/01/4321474373_0257b0b834_o1.jpg  
http://www.destination360.com/asia/japan/tokyo/images/s/tokyo-tower.jpg  
https://images.trvl-media.com/media/content/shared/images/travelguides/destination/179900/Tokyo-Tower-69457.jpg  
http://www.gotoko.org/tc/kanko/minato/spot/images/02_12.jpg  
http://jpninfo.com/wp-content/uploads/2015/10/tokyo-tower-view.jpg  
http://mobiquitous.org/2013/media/carousel-images/tokyo-tower.jpeg  
https://media-cdn.tripadvisor.com/media/photo-s/01/3e/f3/2a/tokyo-tower-at-night.jpg  
http://www.theodora.com/wfb/photos/japan/tokyo_tower_environstokyo_japan_photo_y_shimizu_jnto.jpg
```



◆ 可突破Google一次下載100張圖片的限制



旅行社評論與文字探勘

李鈺祥

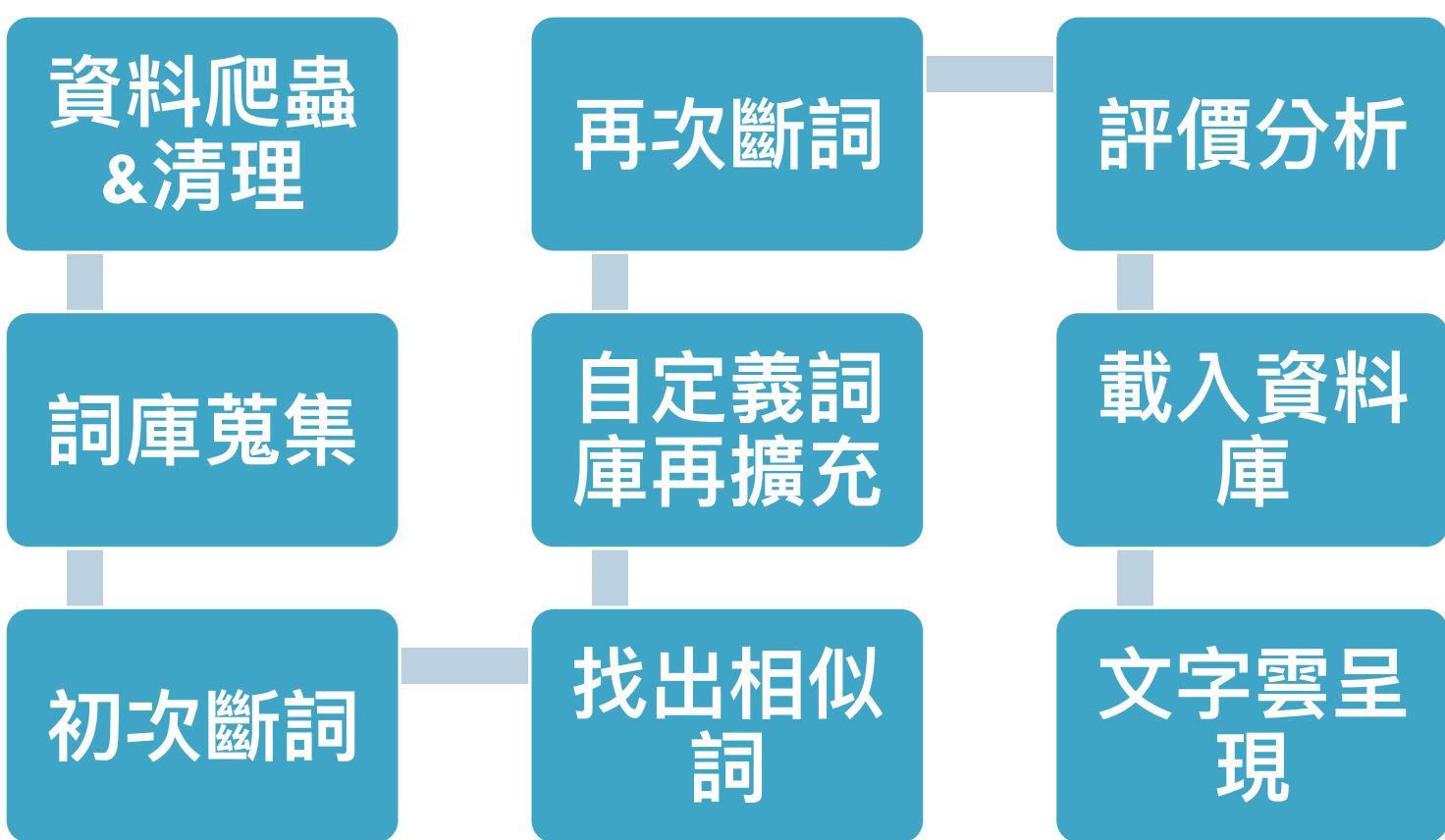
文字探勘—功能簡介

- 負責功能：旅行社評價分析
- 用途：
 - 選擇旅行社的評估參考
 - 行程規劃之參考指標



EVALUATION	
<input type="radio"/>	Excellent
<input type="radio"/>	Very good
<input type="radio"/>	Good

文字探勘—處理流程



文字探勘—資料爬蟲

- 資料來源：
 - PTT(旅行社板)
 - 背包客棧
- 資料數量：
 - 共 13,410 篇
- 資料年份：
 - 自 2011 ~ 至今



文字探勘—資料清理

- 使用語言：Python
- 去除標的物：
 - 去除不要的html標籤
 - 去除不要的文字內容
- 透過：
 - Python中的函式庫
 - 正規表示法



他：那我幫你查一下，然後就動動電腦

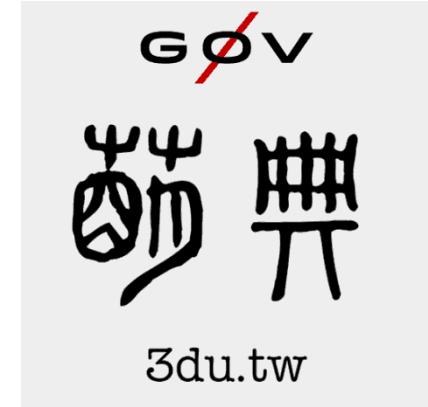
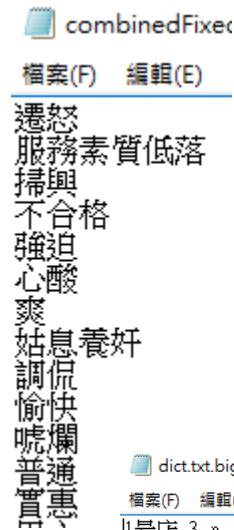
文字探勘—詞庫蒐集

• 詞庫來源：

- Jieba 內建繁體詞庫
- 萌典詞庫
- 自定義詞庫

• 用途：

- 初次斷詞之前處理
- 提升斷詞精準度
- 提升Word2Vec之訓練成效

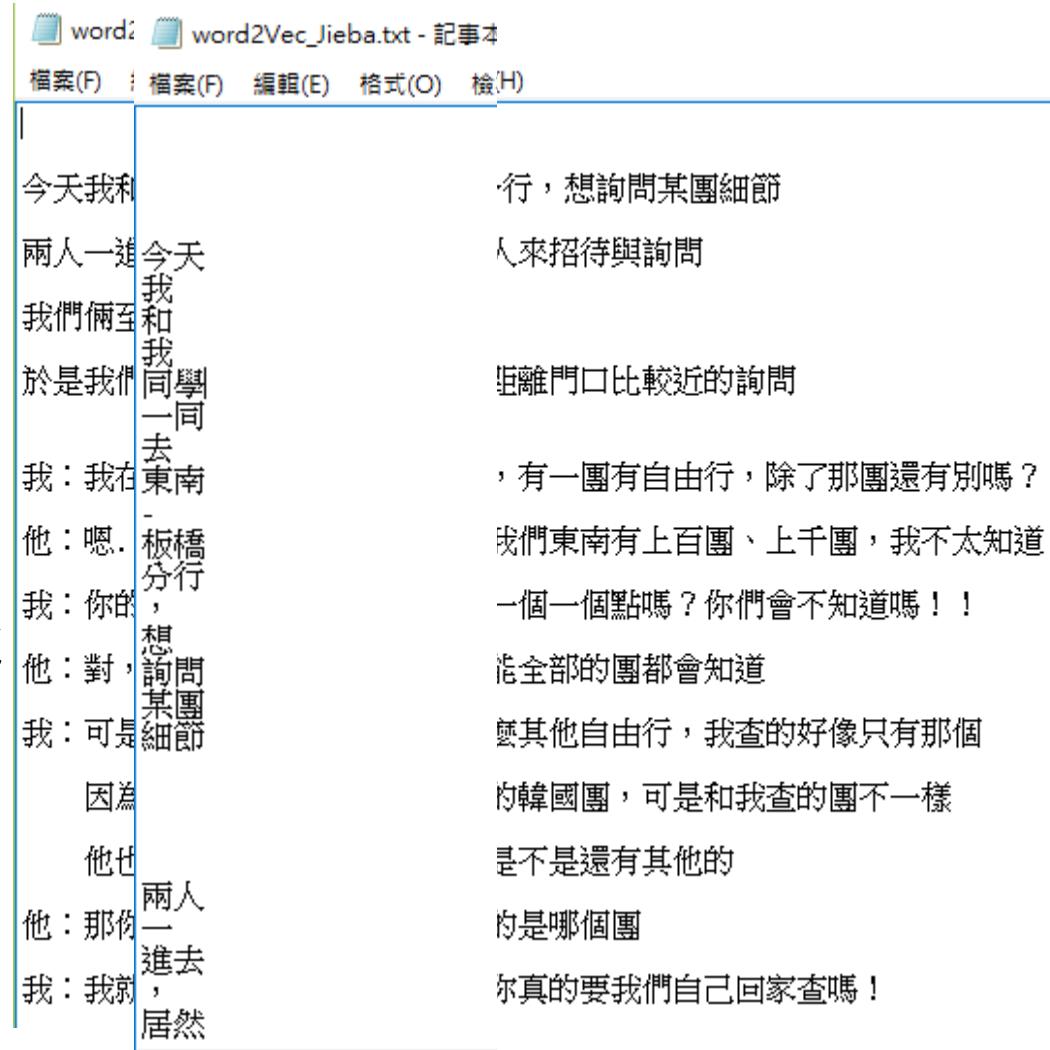


A screenshot of a Windows context menu. At the top, it says 'dict.txt.big.txt - 記事本'. Below that are standard Windows options: '檔案(F)', '編輯(E)', '格式(O)', '檢視(V)', and '說明(H)'. The menu lists many words with their definitions:

- 1號店 3 n 1號店 3 n 4S店 3 n 4S店 3 n AA制 3 n AB型 3
- 十五册 4 m 一万七千 5 m 一万七千余 2 m 一万七千多 2 m
- 三五 6 m 一三六八 2 m 一三四团 3 m 一三四團 3 m 一上
- 两回 8 m 一两场 10 m 一两块 4 m 一两声 24 m 一两处 14
- 2 m 一个五岁 12 m 一个五百 2 m 一个五袋 2 m 一个亿 24
- 2个桶 5 m 一个点 77 m 一个班 84 m 一个飘 4 m 一个甲子
- 一九九五年 20 m 一九九八 3 m 一九九八年 47 m 一九九六
- 万 8 m 一二十两 3 m 一二十个 5 m 一二十件 2 m 一二十
- 下 26 i 一人之交 3 i 一人份 3 m 一人传虚 3 i 一人传虚
- 一来 854 d 一来一去 3 i 一来一往 3 d 一来二去 48 l 一
- 一個十五歲 8 m 一個十個 3 m 一個十八九歲 7 m 一個十八
- 论 3 i 一停 3 d 一側 729 f 一傳 72 b 一傳十 3 i 一億
- m 一兩聲 24 m 一兩萬 40 m 一兩萬元 6 m 一兩萬塊 4 m -
- 则以喜 3 l 一則以優 3 l 一則以慎 3 i 一刪而空 3 z 一
- 24 m 一十五万 2 m 一十五个 3 m 一十五位 3 m 一十五個
- m 一千二百 18 m 一千二百七十二处 2 m 一千二百七十二處
- 百多年 2 m 一千八百萬噸 5 m 一千八百餘 2 m 一千八百餘
- 13 m 一千間 3 m 一千间 3 m 一千隻 10 m 一千零一夜 30
- mrr 一古脑 5 l 一古脑儿 108 l 一古脑 5 l 一古脑兒 108

文字探勘—初次斷詞

- 使用工具：**Jieba**
- 用途：
 - Word2Vec之前處理
- 斷詞範疇：
 - 涵蓋PTT與背包客棧
 - 不分各家旅行社
 - 共 13,410 篇文章



The screenshot shows a Windows Notepad window with two tabs: 'word2Vec.txt' and 'word2Vec_Jieba.txt'. The 'word2Vec_Jieba.txt' tab is active and displays a list of words from a text file, with each word followed by its corresponding Jieba tokenization. The tokens are separated by vertical lines. The right column contains the original text and the left column contains the tokenized text.

原句	斷詞	說明
今天我和	今天 我 和	行，想詢問某團細節
兩人一起	兩 人 一 起	人來招待與詢問
我們倆至	我 倆 至	
於是我們	於 是 我 倆	距離門口比較近的詢問
我：我在	我 在	
他：嗯..	他 愠 ..	
我：你的	我 你 的	，有一團有自由行，除了那團還有別嗎？
他：板橋	他 板 橋	我們東南有上百團、上千團，我不太知道
我：分行	我 分 行	一個一個點嗎？你們會不知道嗎！！
他：對，	他 對 ，	這全部的團都會知道
我：可是	我 可 是	
因為	因 為	要其他自由行，我查的好像只有那個
他也	他 也	的韓國團，可是和我查的團不一樣
他：那個	他 那 倆	是不是還有其他的
我：我就	我 就	是哪個團
		你真的要我們自己回家查嗎！

文字探勘—找出相似詞

- 使用工具：**Word2Vec**
- 用途：
 - 找出相似並且可能遺漏的字詞
- 投入訓練量：
 - 751,490 句 / 訓練單位

```
...
這邊是讀取訓練結果，可以不用再重新訓練。
使用most_similar這個方法做字詞的查詢，並且顯示前N個相似字詞
...
這邊是讀取訓練結果，可以不用再重新訓練。
使用most_similar這個方法做字詞的查詢，並且顯示前N個相似字詞
...
load_model = Word2Vec.load(r'D:\BigData\Python\works\chinese_stopwords.model')
for name, weight in load_model.most_similar(u'不爽', topn=10):
    print name, weight
```

不高興 0.823772668839
貪小便宜 0.813368439674
唬爛 0.808220982552
不划算 0.806670665741
被騙 0.805783569813
活該 0.805532753468
笨蛋 0.803521513939
很慘 0.800624847412
想想 0.800486445427
得了 0.798668861389
難看 0.796682596207
很煩 0.839766919613
好笑 0.837970972061

文字探勘—自定義詞庫再擴充

- 遺漏詞來源：

- Word2Vec訓練出的相似詞

- 作法：

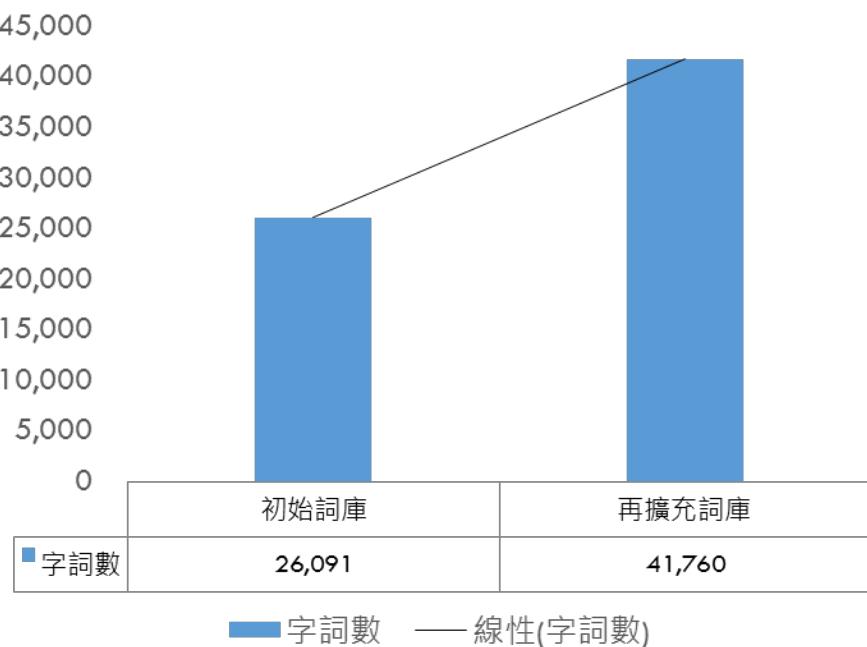
- 篩選出符合實際需求之遺漏詞

- 並建立於原有詞庫中

- 數量差異：

- 擴充前詞量 **26,091字**
 - 擴充詞後量 **41,760字**

詞庫數量擴充前後比較



文字探勘—再次斷詞

- 使用工具：**Jieba**
- 用途：
 - 評價分析前處理
- 斷詞範疇：
 - 涵蓋PTT與背包客棧
 - 分為**20家旅行社**
 - 共 13,410 篇文章

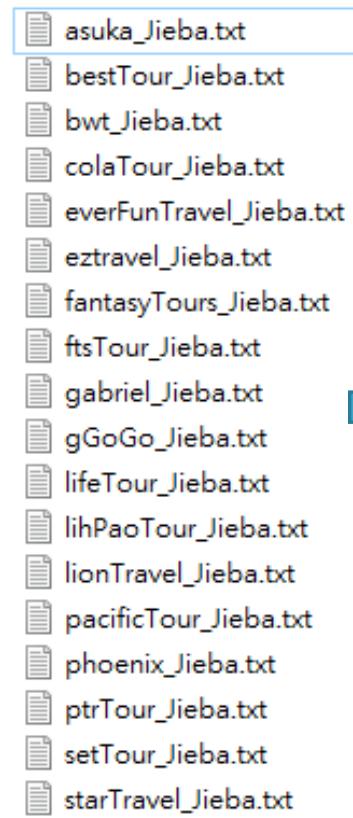
asuka.txt
bestTour.txt
bwt.txt
colaTour.txt
everFunTravel.txt
eztravel.txt
fantasyTours.txt
ftsTour.txt
gabriel.txt
gGoGo.txt
lifeTour.txt
lihPaoTour.txt
lionTravel.txt
pacificTour.txt
phoenix.txt
ptrTour.txt
setTour.txt
starTravel.txt



asuka_Jieba.txt
bestTour_Jieba.txt
bwt_Jieba.txt
colaTour_Jieba.txt
everFunTravel_Jieba.txt
eztravel_Jieba.txt
fantasyTours_Jieba.txt
ftsTour_Jieba.txt
gabriel_Jieba.txt
gGoGo_Jieba.txt
lifeTour_Jieba.txt
lihPaoTour_Jieba.txt
lionTravel_Jieba.txt
pacificTour_Jieba.txt
phoenix_Jieba.txt
ptrTour_Jieba.txt
setTour_Jieba.txt
starTravel_Jieba.txt

文字探勘—評價分析

- 使用方法：
 - 共現詞頻分析
- 用途：
 - 找出符合該旅行社之真實評價
 - 提供使用者在挑選旅遊方案的參考

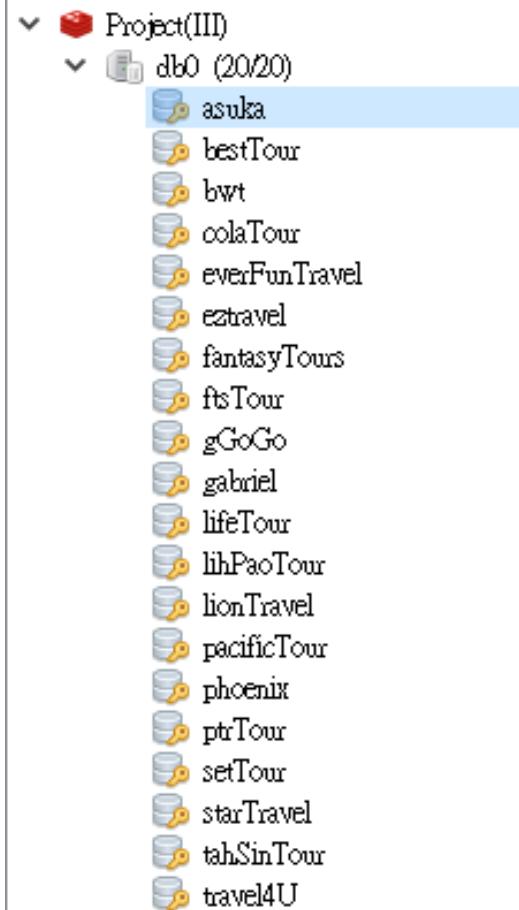


asuka_result.txt - 記事本
檔案(F) 編輯(E) 格式(O) 檢

評價	相似度
差	0.328701
便宜	0.274981
比較簡單	0.234626
只有蜻蜓點水	0.197447
抱怨	0.169635
太趕	0.148218
不夠好	0.142587
佳	0.131311
不錯	0.127974
很好	0.118758
超爽	0.117677
比較推薦	0.090384
也不錯	0.083083
比較好	0.077478
也不少	0.073037
真的不錯	0.071499
有點誇張	0.066846
比較划算	0.066046
真的很推	0.056228
鳥	0.053387
沒有特別推薦	0.051242

文字探勘—載入資料庫

- 使用資料庫：**Redis**
- 儲存型態：
 - Hash (Key/Value of 2 layers)
- 儲存資料：
 - 20家旅行社之評價分析
- 用途：
 - 為前端網頁文字雲之呈現資料





圖片辨識

倪裕程、鄭云鈞

圖片辨識—圖片資料狀態

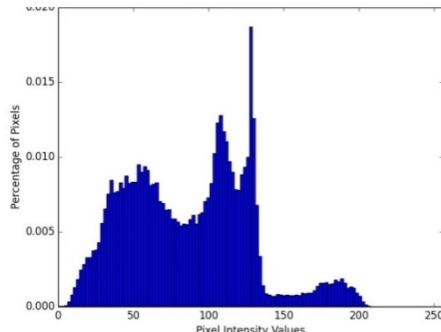
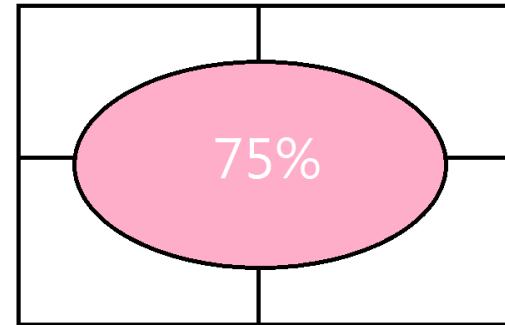
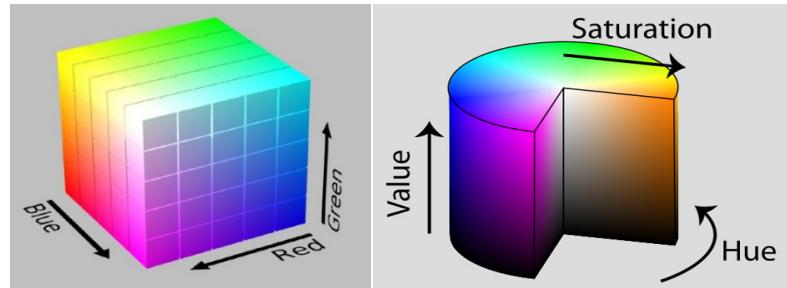
- 景點圖片比對遇到的問題

- 日夜
- 季節
- 建築物



圖片辨識—顏色辨識

- 色塊直方
 - 將圖片從RGB色彩轉為
 - H (色相)
 - S (飽和度)
 - V (明度)
- 將圖片分成五大區域
 - 中間區域佔75%
- 分別針對各區域計算
 - 顏色直方圖



圖片辨識—特徵值選取

- 計算特徵值

- 每個區塊都會分成H(8)S(12)V(3)

- 分別有288個每個顏色出現的累積值

- 有5個區塊，共1440個累積值做為特徵值

dataset\\100000.jpg	5.15E-06	4.17E-05	3.57E-06	1.19E-06	0	0	0	0.995165	0	0	0	0	0	0	0	0
dataset\\100001.jpg	3.72E-05	3.51E-05	1.24E-05	0	0	0	0	0.95551	0	2.06E-06	0	1.24E-05	0	0	0	0

圖片檔名

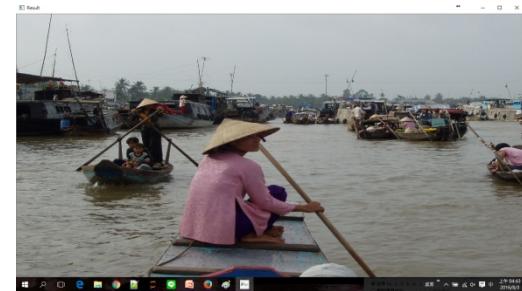
共儲存有1440個累積值

圖片辨識—顏色辨識範例

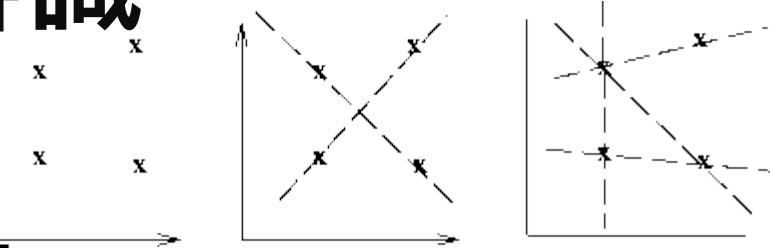
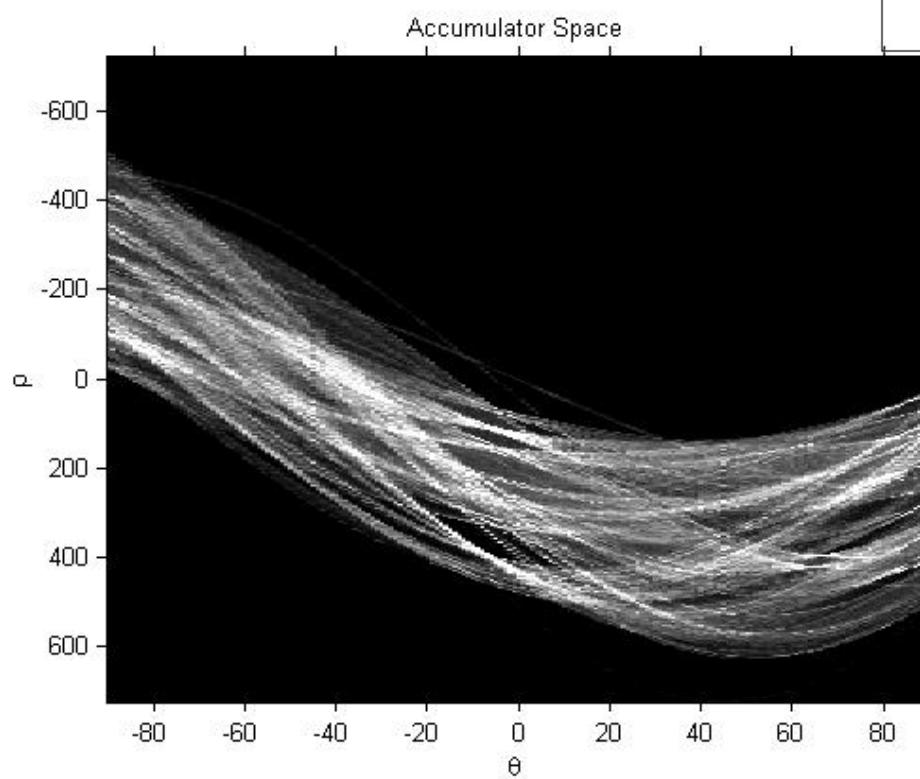
query



result



圖片辨識—形狀辨識

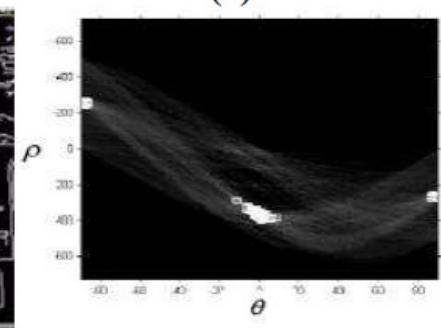


(a)

(b)



(c)



(d)

圖片辨識—特徵值選取

- 計算特徵值

- The peak percentage 從0至179度
- The distance ratio feature 從0至179度
- 共360個維度

dataset100301.jpg	0.001548	0.019436	0.00479	0.000945	0.022681	0.007074	0.000268	0.009835	0.004457	8.31E-05	0.003626	0.00101	9.23E-06	0.000357	3.69E-05	6.16E-06
dataset100400.jpg	0.000411	0.000776	0.262822	8.96E-05	2.24E-05	4.48E-05	4.85E-05	1.87E-05	7.47E-06	8.21E-05	2.24E-05	3.36E-05	3.36E-05	1.49E-05	0	1.87E-05

圖片檔名

共儲存有360特徵值

圖片辨識—特徵值選取

- **The peak percentage**

The peak percentage is defined by

$$p_\theta = \frac{n_\theta}{n}$$

where n_θ represents the number of Hough peaks in the θ^{th} bin and n denotes the total number of Hough peaks in the Hough parameter space.

圖片辨識—特徵值選取

• The distance ratio feature

Also the distance ratio feature is defined by

$$d_\theta = \frac{\sqrt{\sum_{j=1}^{n_\theta} (\rho_j)^2}}{C_\theta}$$

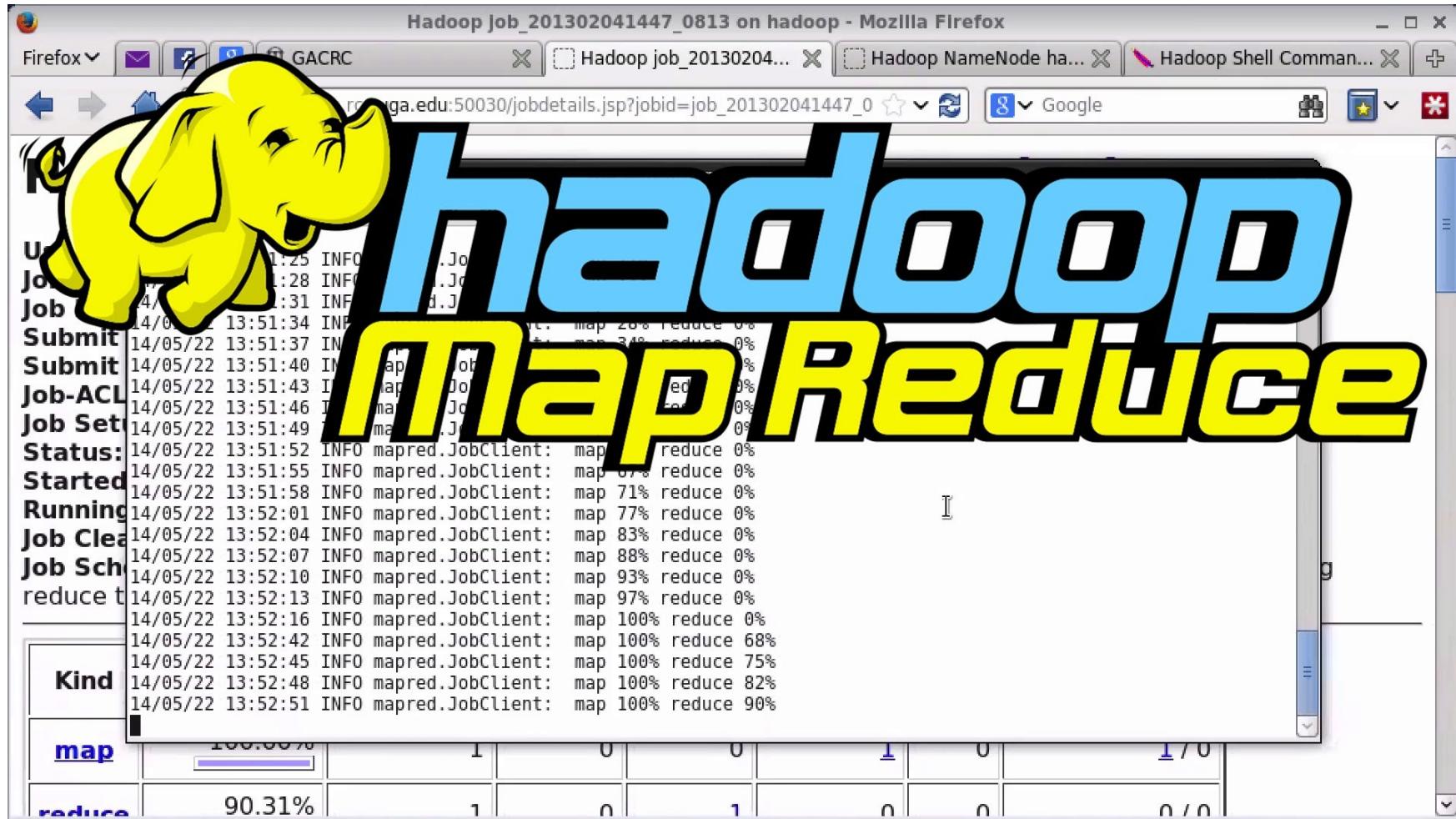
where ρ_j represents the distance from the θ^{th} Hough cell to the horizontal line and C_θ signifies the weighted-centroid of the θ^{th} cell defined by

$$C_\theta = \frac{\sum_{j=1}^{n_\theta} w_j \cdot \rho_j}{\sum_{j=1}^{n_\theta} w_j}$$

with w_j representing the number of points in the edge map that form the j^{th} Hough peak. The distance ratio vectors of the

圖片辨識—演算法選擇

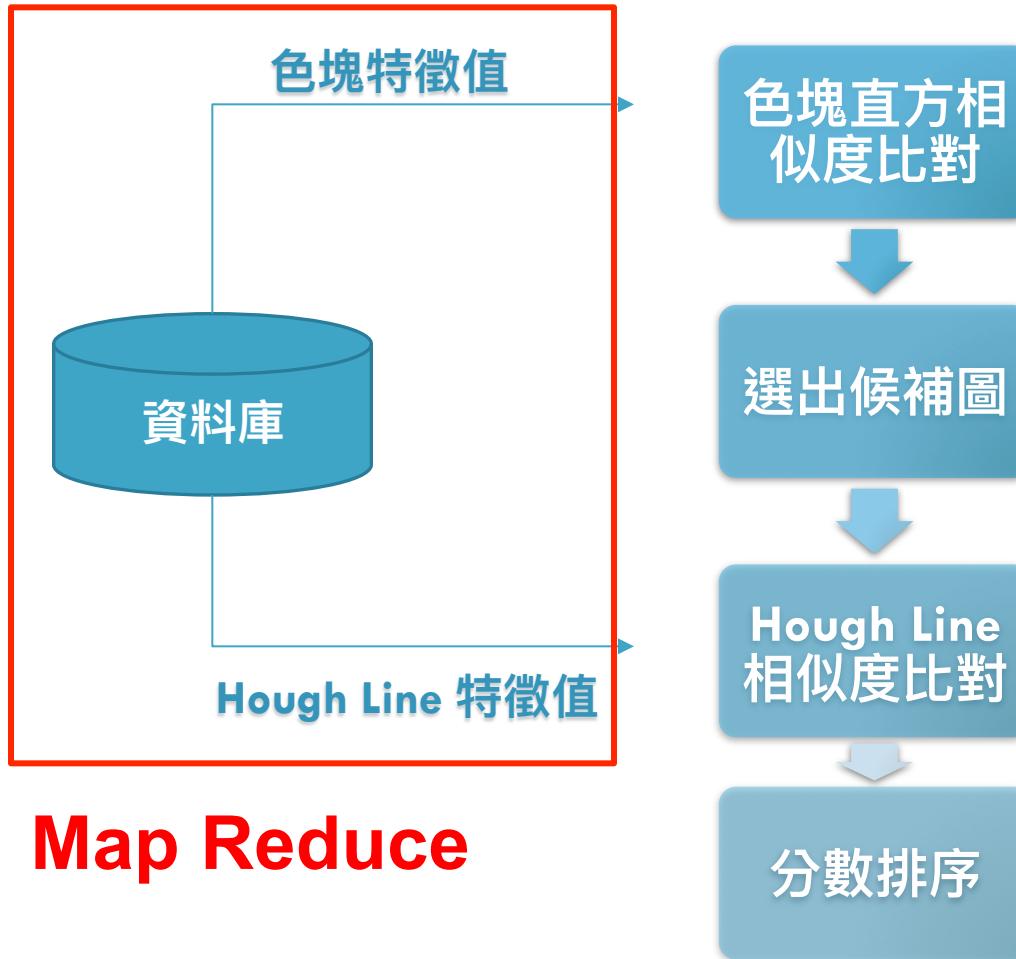
- 色塊直方圖相似比對
 - 解決日夜或是季節顏色問題
 - 減少Hough Line特徵值比對運算時間
- Hough Line 特徵值相似度比對
 - 解決建築物形狀比對
 - 解決色塊搜尋到不相關圖片



分散式運算應用

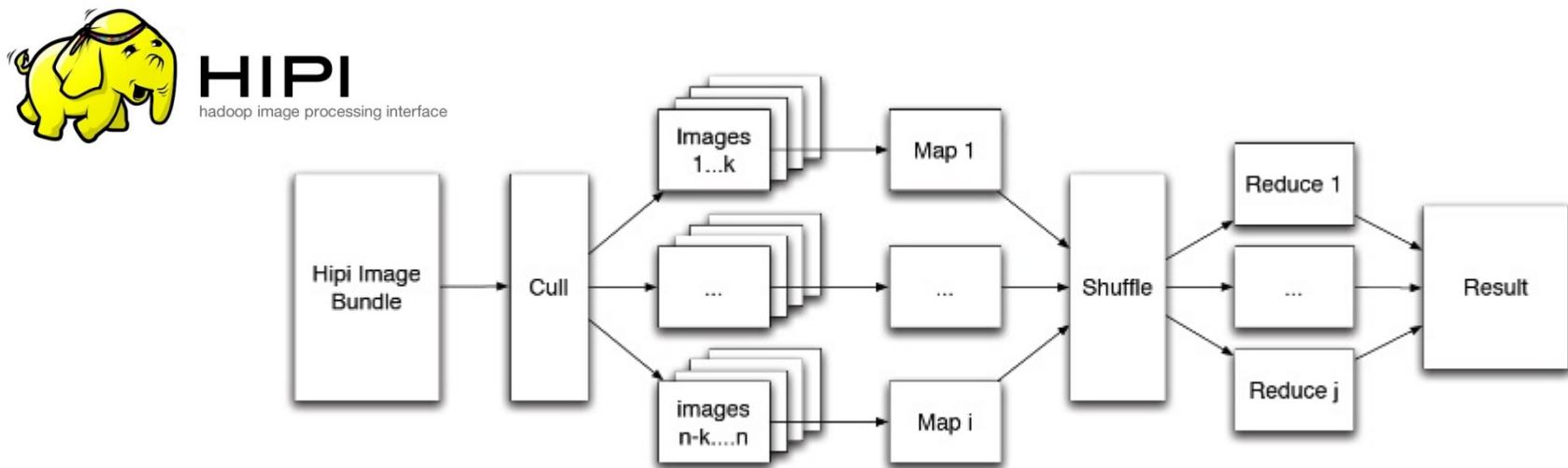
倪裕程

圖片辨識—分散式應用



圖片辨識—資料庫特徵值

- 為何使用 Hadoop map reduce
 - 行程景點眾多，網路爬取圖片量大
 - 藉由 Hadoop 進行分散式儲存及運算
 - 可節省利用單機運算的時間



圖片辨識—資料庫特徵值

```
@Override
public void map(HipiImageHeader header, FloatImage image,
                 Context context) throws IOException, InterruptedException {
    String output = null;
    if (header == null) {
        output = "Failed to read image header.";
    } else if (image == null) {
        output = "Failed to decode image data.";
    } else {
        System.loadLibrary(Core.NATIVE_LIBRARY_NAME);
        Mat images = this.convertFloatImageToOpenCVMat(image);
        houghLineTransform houghlineChar = new houghLineTransform(
                images);
        Map<Double, Double> DistanceRatio = houghlineChar
                .getDistanceRatio();
        Map<Double, Double> Percentage = houghlineChar.getPercentages();
        String filename = header.getMetaData("filename");
        output = filename;
        // Write in to map
        for (double key : Percentage.keySet()) {
            output = output + "," + Percentage.get(key);
        }
        for (double key : DistanceRatio.keySet()) {
            output = output + "," + DistanceRatio.get(key);
        }
    }
    context.write(new IntWritable(1), new Text(output));
}
```

```
@Override
public void map(HipiImageHeader header, FloatImage image,
                 Context context) throws IOException, InterruptedException {
    String output = null;
    if (header == null) {
        output = "Failed to read image header.";
    } else if (image == null) {
        output = "Failed to decode image data.";
    } else {
        // Load OpenCV native library
        try {
            System.loadLibrary(Core.NATIVE_LIBRARY_NAME);
        } catch (UnsatisfiedLinkError e) {
            System.err.println("Native code library failed to load.\n" + e
                    + Core.NATIVE_LIBRARY_NAME);
            System.exit(1);
        }
        // convert to OpenCV Mat
        Mat images = this.convertFloatImageToOpenCVMat(image);
        // Hsv hist feature
        MatOfInt bin = new MatOfInt(8, 12, 3);
        ColorDescriptor cd = new ColorDescriptor(bin);
        List<Mat> feature = cd.describe(images);
        String name = header.getMetaData("filename");
        output = name;
        for (int i = 0; i < 5; i++) {
            for (int j = 0; j < 288; j++) {
                String elementOfMat = "" + feature.get(i).get(j, 0)[0];
                output = output + "," + elementOfMat;
            }
        }
    }
    context.write(new IntWritable(1), new Text(output));
}
```

圖片辨識—特徵值提取時間差

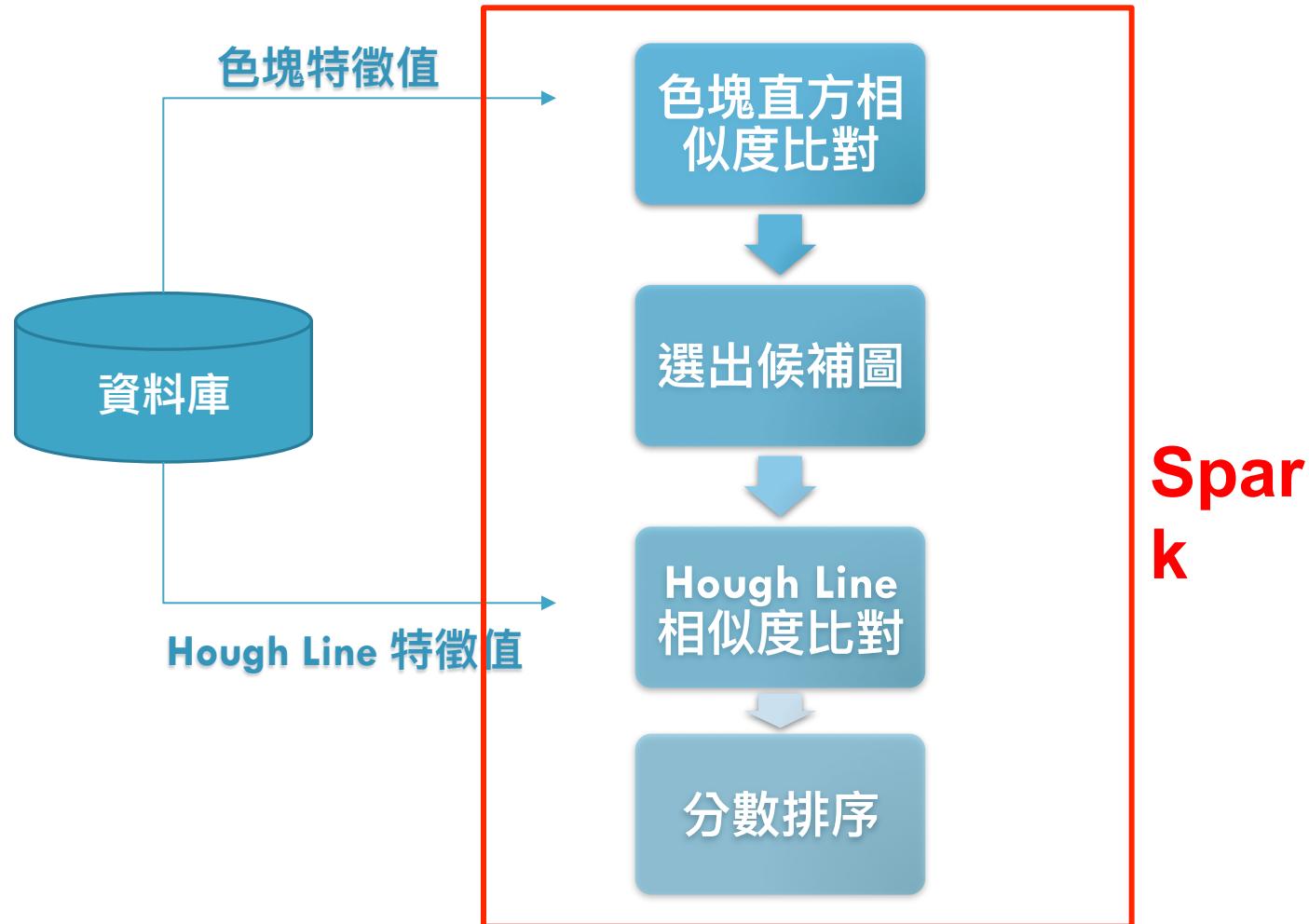
- 以10張圖比較Hough line特徵值計算時間
 - 單機- 75s
 - Map reduce- 33s
- 以100張圖比較Hough line特徵值計算時間
 - 單機- 88m42s
 - Map reduce- 56m5s



圖片資料庫搜尋引擎

鄭云鈞

圖片辨識—流程與分散式應用



圖片辨識— SPARK

- 呼叫後可一直在記憶體中持續運行
- 在記憶體中運行計算，計算速度快
- 提供彈性分佈式數據集(RDD)
- 利用Spark計算被查詢的圖片及特徵值比對

```
bins = (8, 12, 3)
query = url_to_image(image_seqfile_path)
queryFeatures = describe(query,bins)

text = taxiFile.map(lambda k: k.split(",")).map(lambda k: colorfeature(k,queryFeatures)).sortByKey().map(lambda k: k[1])
colorquery = text.take(105)

houghquery = hough(query)
result = taxiFile.map(lambda k: k.split(",")).filter(lambda k: k[0].encode('utf-8') in colorquery).map(lambda k:rankingFunction(houghquery,k)).reduceByKey(max).map(lambda k: (k[1],k[0])).map(lambda k: k[1])

final = result.take(3)

connection = redis.Redis(host = '36.236.48.175', port = 6379, db = 4)
connection.flushdb()
connection.lpush('query', *final)
```

圖片辨識—色塊直方相似度比對

- 先進行圖片色塊比對
- 計算被查詢圖片與資料庫圖片的色塊卡方相似度

```
text = taxiFile.map(lambda k: k.split(",")).map(lambda k: colorfeature(k,queryFeatures)).sortByKey().map(lambda k: k[1])
```

$$\chi^2(x, y) = \sqrt{\sum_{i=1}^p \left(\frac{x_i - E(x_i)}{E(x_i)} \right)^2 + \sum_{i=1}^p \left(\frac{y_i - E(y_i)}{E(y_i)} \right)^2}$$



```
(2.1774022620351747e-07, '123300.jpg'),  
(6.4339739536879405e-07, '128300.jpg'),  
(2.4936570199980801e-06, '119602.jpg'),  
(3.7364309595940535e-06, '145300.jpg'),  
(8.0154014764149456e-06, '141601.jpg'),  
(9.6966978644122727e-06, '122700.jpg'),
```

- 數值越小，相似度越高
- 只要找出前105張圖來做Hough Line比對

圖片辨識—搜尋範例

Query



候補圖



圖片辨識— HOUGH LINE 相似度比對

●再進行Hough Line形狀比對

```
result = taxiFile.map(lambda k: k.split(",")).filter(lambda k: k[0].encode('utf-8') in colorquery).map(lambda k: rankingFunction(houghquery,k)).reduceByKey(max).map(lambda k: (k[1],k[0])).map(lambda k: k[1])
```

The m^{th} circular correlation between \mathbf{p}^q and \mathbf{p}^k is defined as

$$s_p^k(m) = \frac{\sum_{j=1}^{J-m} (P_j^q P_{j+m}^k) + \sum_{j=J-m+1}^J (P_j^q P_{j+m-J}^k)}{\sqrt{\sum_{j=1}^J (P_j^q)^2} \sqrt{\sum_{j=1}^J (P_j^k)^2}}$$

which is normalized between 0 and 1. We calculate the maximum circular correlation $\max_m\{s_p^k(m)\}$ between \mathbf{p}^q and \mathbf{p}^k , and use to retrieve the candidate images.

```
u'137901.jpg',
u'0.0830464',
u'0.0272187',
u'0.38288',
u'0.0509615',
u'0.0277534',
u'0.0314432',
u'0.0886613',
u'0.0373789',
u'0.00251332',
u'0.220584',
u'0.015882',
```



Similarly, the m^{th} circular correlation between \mathbf{r}^q and \mathbf{r}^k is defined as

$$s_r^k(m) = \frac{\sum_{j=1}^{J-m} (r_j^q r_{j+m}^k) + \sum_{j=J-m+1}^J (r_j^q r_{j+m-J}^k)}{\sqrt{\sum_{j=1}^J (r_j^q)^2} \sqrt{\sum_{j=1}^J (r_j^k)^2}}$$

which is normalized between 0 and 1. Also we calculate the maximum circular correlation $\max_m\{s_r^k(m)\}$ between \mathbf{r}^q and \mathbf{r}^k , and use to retrieve the candidate images.

```
('137901.jpg', 0.775079348140713),
('147500.jpg', 0.4231559767303664),
('141801.jpg', 0.5377062084315848),
('134503.jpg', 0.5447344501367279),
('148601.jpg', 0.3893849151577147),
('118601.jpg', 0.47478382659779406),
('117801.jpg', 0.6186679239688069),
('123402.jpg', 0.6685508610019886),
```

圖片辨識— HOUGH LINE 相似度比對

```
result = taxiFile.map(lambda k: k.split(",")).filter(lambda k: k[0].encode('utf-8') in colorquery).map(lambda k:rankingFunction(oughquery,k)).reduceByKey(max).map(lambda k: (k[1],k[0])).map(lambda k: k[1])  
  
('137901.jpg', 0.775079348140713),  
('147500.jpg', 0.4231559767303664),  
('141801.jpg', 0.5377062084315848),  
('134503.jpg', 0.5447344501367279),  
('148601.jpg', 0.3893849151577147),  
('118601.jpg', 0.47478382659779406),  
('117801.jpg', 0.6186679239688069),  
('123402.jpg', 0.6685508610019886),  
  
'118201.jpg',  
'132700.jpg',  
'147500.jpg',  
'112001.jpg',  
'113402.jpg',  
'119600.jpg',  
'145300.jpg',  

```

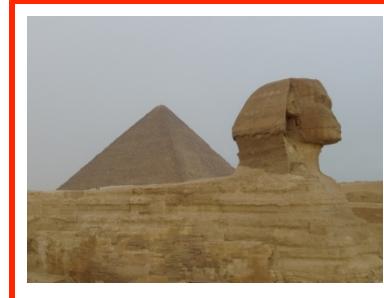
- 數值越大，表示建築物形狀越相近
- 選出前三個數值最大的景點作行程推薦

圖片辨識—搜尋範例

Query



候補圖



圖片辨識—行程推薦

- 將Spark查詢結果從Redis資料庫叫出

The_Tsukiji_Market
Daigoji_Temple
Ginkaku_ji_Temple

- 從Redis資料庫拉出相對應相關行程

```
[ 'GFG0000002920' ]  
[ '469_1', '469' ]  
[ 'VDR0000001838_N00061305BR6B', 'VDR0000001838_N00061005CI6EB', '50499', '50398', '43408', '41639', '39404' ]
```

- 從使用者選擇的類型計算出最適合行程
- 計算每個行程中特定類型所占分數，找出最大值
- 若有同分，則再比較拉手分數
- 若無選擇喜歡類型，只比較拉手分數

圖片辨識—行程推薦

- 三種結果訊息呈現
- 該行程所屬旅行社 `setTour`

- 該行程有哪些景點

[Fuji_Hakone_Izu_National_Park](#)
[Shin_Nakamise_Shopping_Street](#)
[Venetian_Glass_Museum](#)
[Senso_ji_Temple](#)
[Hakone_Pirate_Ship](#)
[Diver_City_Tokyo_Plaza](#)

- 該行程每種類別所佔比例

景點與地標 14
自然與公園 3
購物 2
水上樂園和遊樂場 2
旅客資訊 1
博物館 1
休閒與娛樂 1
飲食 0

未來展望

- 增加風景辨識功能
- 增加自由行功能
- 文字探勘前處理更完善
- 景點行程類型前處理更完善

" I have a dream..."

Q&A時間

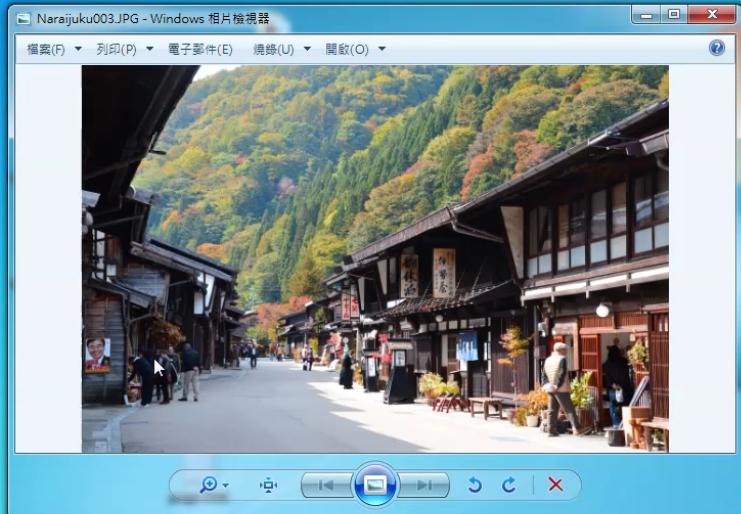


Thanks for watching



BACKUP







Please select 1-3 types of tourist spots that you prefer

- Sights & Landmarks
- Nature & Parks
- Museums
- Spa & Wellness
- Shopping
- Zoos & Aquariums
- Water & Amusement Parks
- Outdoor Activities
like hiking tour, canoeing etc.

Send data

非常專業
也不錯
很好 抱怨很滿意
驚豔 佳
瑕庇 比較划算
相當不錯
很漂亮
真的好 再麻煩比較好

