

PBC

解析合購版

B06302201 陳品儒 B07303050 周鈺淇

B07701236 柯昱丞 B09610020 黃韻文

B09702032 林子昕 P09323028 吳金擇

CONTENTS

01. 專題理念

02. 網頁簡介

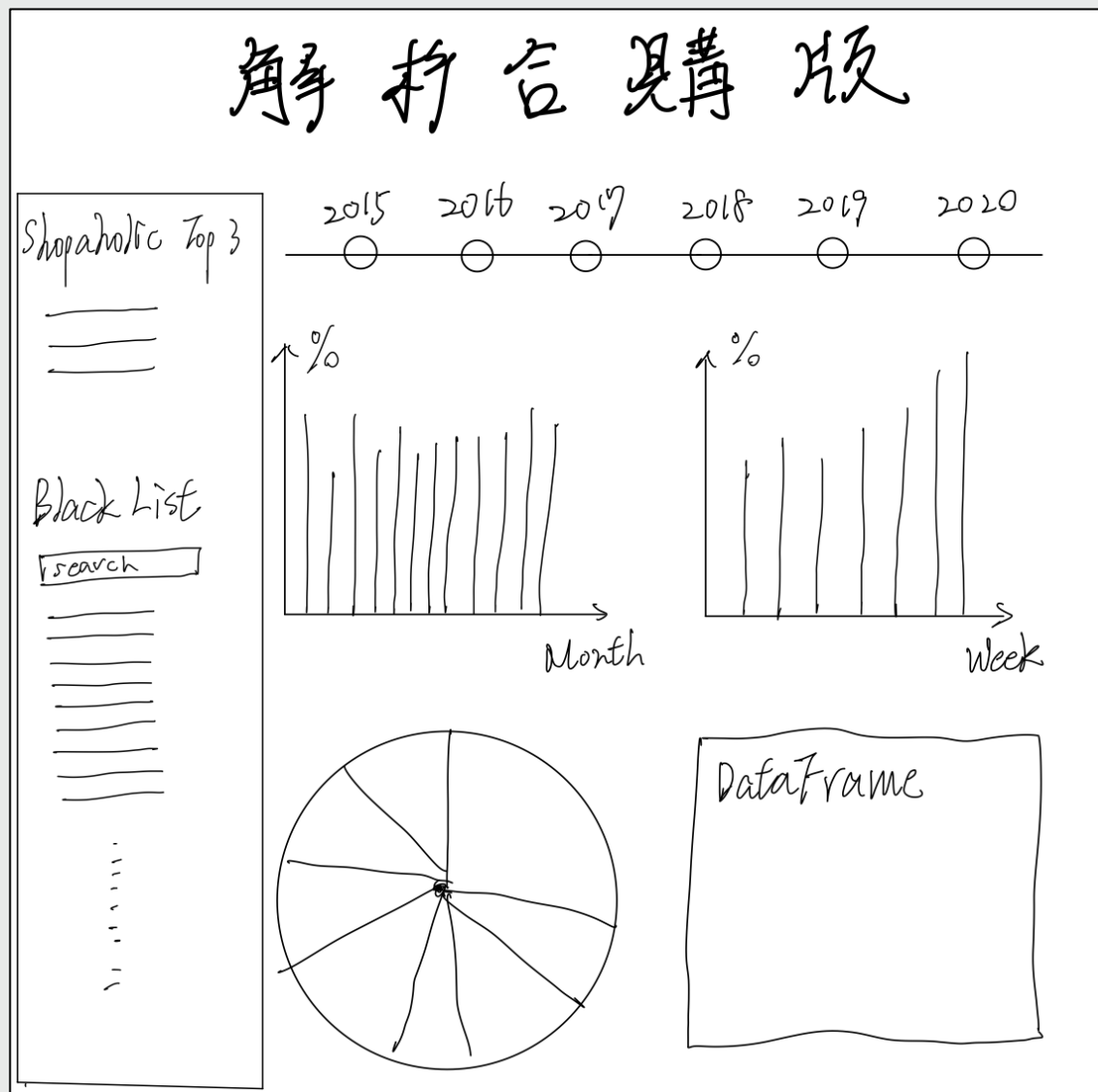
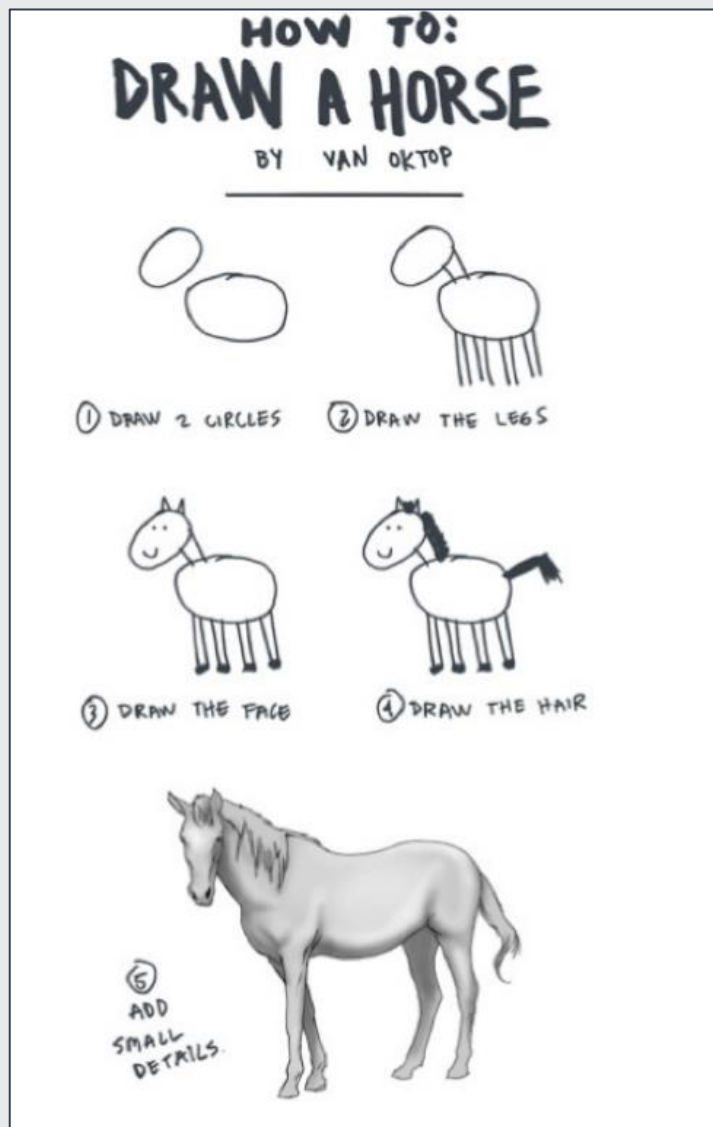
03. 程式簡介

04. 心得展望

01

PART ONE
專題理念



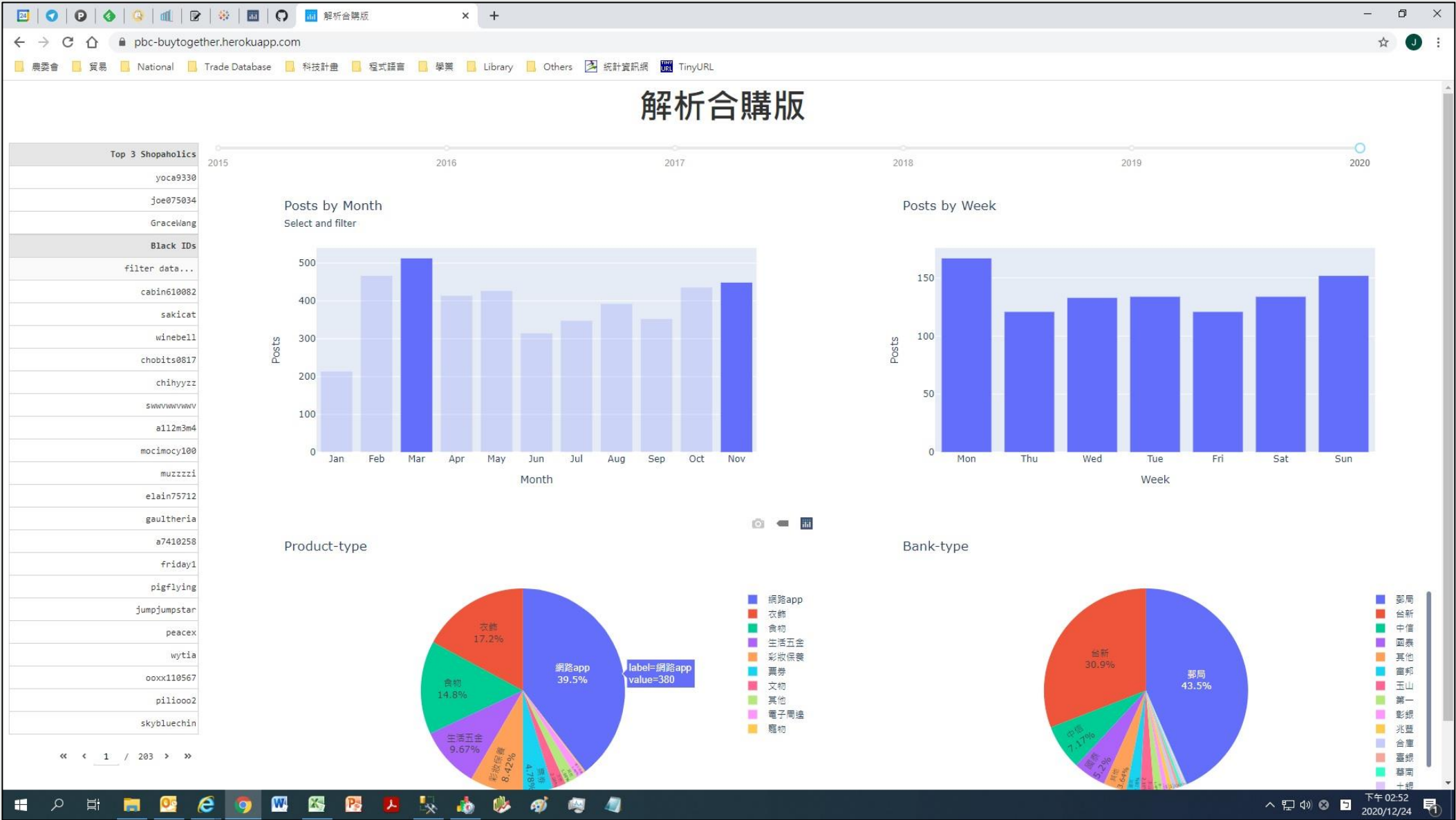


02

PART TWO

網頁簡介

<https://pbc-buytogether.herokuapp.com/>



03

PART THREE

程式簡介

<https://github.com/PinJu-Chen/PTT-Analysis>

01 爬取資料

02 資料整理

03 製作圖表

04 網頁部署

程式簡介 - 01 爬取資料

PART THREE

```
PTT crawler.py x data_manipulation.py x app.py x
1  #coding=utf-8
2  import re
3  import sys
4  import json
5  import requests
6  import time
7  from datetime import datetime
8  from bs4 import BeautifulSoup

174 def store(data):
175     with open(fileName, 'a', encoding="utf-8") as f:
176         f.write(data.encode(sys.stdin.encoding, "replace").decode(sys.stdin.encoding))
177
178 if __name__ == "__main__":
179     PttName, ParsingPage = 'BuyTogether', int(3999)
180     start_time = time.time()
181     print('Start parsing ' + PttName + '....')
182     fileName = 'D:\PCB\data-' + PttName + '-' + datetime.now().strftime('%Y%m%d%H%M%S') + '.json'
183     # 檢查看板是否為18禁,有些看板為18禁
184     soup = over18(PttName)
185     ALLpageURL = soup.select('.btn.wide')[1]['href']
186     # 得到本看板全部的index數量
187     ALLpage = int(getPageNumber(ALLpageURL)) + 1
188     index_list = []
189     for index in range(ALLpage, ALLpage - int(ParsingPage), -1):
190         page_url = 'https://www.ptt.cc/bbs/' + PttName + '/index' + str(index) + '.html'
191         index_list.append(page_url)
192
193     store('\n')
194     crawler(index_list)
195
196     # 移除最後一個 ", " 號
197     with open(fileName, 'r', encoding="utf-8") as f:
198         content = f.read()
```

程式簡介 - 02 資料整理

PART THREE

```
PTT crawler.py x data_manipulation.py x app.py x
1  # -*- coding: utf-8 -*-
2  import json
3  import pandas as pd
4  # import re
5
6  # 讀取json檔
7  with open('rawdata.json', 'r', encoding="utf-8") as f:
8      data = json.loads(f.read())
9
10 # 以list建立datafram，以方便後續繪圖
11 buy2df = pd.DataFrame(all_list, columns=['id','author','title','date','ip','likes','content'])
12
13 # 全改為小寫，以方便對照
14 df['content'] = df['content'].str.lower()
15
16 for i in b.keys():
17     for j in range(len(b[i])):
18         b[i][j] = b[i][j].lower()
19
20 # 建立bank清單，並以長度做排序
21 bank = []
22
23 for i in b.keys():
24     for j in range(len(b[i])):
25         bank.append((i, b[i][j]))
26
27 # 定義銀行函數
28 def banktype(content):
29     for i in range(len(bank)):
30         if content.find(bank[i][1]) > 0:
31             return bank[i][0]
32     return '其他'
33
34 # 將分類函數套用至df
35 df['bank'] = df['content'].apply(lambda x: banktype(x))
36
37 df.to_csv(r'rawdata.csv', index=False)
```

程式簡介 - 03 製作圖表

PART THREE

```
PTT crawler.py x data_manipulation.py x app.py x
1  # -*- coding: utf-8 -*-
2  import pandas as pd
3  import plotly.express as px
4  import dash
5  import dash_core_components as dcc
6  import dash_html_components as html
7  from dash.dependencies import Input, Output
8  import dash table
32  # 黑名單
33  blackdf = df[(df.title.str.contains('黑人')|
34              df.title.str.contains('灰人')|
35              df.title.str.contains('黑名單')|
36              df.title.str.contains('判決'))&
37              ((df.year == 2020)|
38              (df.year == 2019)|
39              (df.year == 2018)|
40              (df.year == 2017))&
41              (~df.title.str.contains('RE:'))&
42              (~df.title.str.contains('Re:'))]
43  blackdf['black_ID'] = blackdf.title.str.replace('^.+?([a-zA-Z0-9]+).+$', r'\1')
44
45
46  '''-----main-----'''
47  # css
48  external_stylesheets = ['https://codepen.io/chriddyp/pen/bWLwgP.css']
49
50  # 建立app
51  app = dash.Dash(__name__, external_stylesheets=external_stylesheets)
52  # server = app.server
53
54  # 主標
55  header = html.H1(children="解析合購版", style={'font-weight': 'bold'})
```

程式簡介 - 03 製作圖表

PART THREE

```

57 # Top 3 table
58 top3 = dash_table.DataTable(
68
69 # Blacklist table
70 blacklist = dash_table.DataTable(
83
84 # 時間軸
85 year_slider = dcc.RangeSlider(id='year_slider',
97
98 # 以dcc.Graph建立instance，存放月份圖
99 month_bar = dcc.Graph(id='month_bar',
100                        selectedData=None,
101                        className="five columns")
102
103 # 以dcc.Graph建立instance，存放星期圖
104 week_bar = dcc.Graph(id='week_bar', className="five columns")
105
106 # 以dcc.Graph建立instance，存放商品圓餅圖
107 product_pie = dcc.Graph(id='product_pie', className="five columns")
108
109 # 以dcc.Graph建立instance，存放支付銀行圓餅圖
110 bank_pie = dcc.Graph(id='bank_pie', className="five columns")

```

程式簡介 - 03 製作圖表

PART THREE

```
130 '''-----Month callback-----'''
131 @app.callback(
132     Output('month_bar', 'figure'),
133     Output('month_bar', 'selectedData'),
134     Input('year_slider', 'value'))
135 def update_month_bar(selected_year):
136     # 年度
137     dfm = dff[(dff['year'] <= max(selected_year)) & (dff['year'] >= min(selected_year))]
138
139     # 更新 month_bar 月份圖
140     chartm = px.bar(x=dfm.groupby('month').size().index,
141                    y=dfm.groupby('month').size(),
142                    title="Posts by Month" + "<br>" + '<span style="font-size: 12px;">Select and filter</span>',
143                    labels={"x": "Month",
144                           "y": "Posts"},
145                    category_orders={"x":
146                                    ['Jan', 'Feb', 'Mar', 'Apr', 'May', 'Jun',
147                                     'Jul', 'Aug', 'Sep', 'Oct', 'Nov', 'Dec']})
148
149     # 設定clickmode，作為後續的input
150     chartm.update_layout(clickmode='event+select')
151     return chartm, None
```



Jump to Favorites, Apps, Pipelines, Spaces...



Deploy using Heroku Git

Use git in the command line or a GUI tool to deploy this app.

Install the Heroku CLI

Download and install the [Heroku CLI](#).

If you haven't already, log in to your Heroku account and follow the prompts to create a new SSH public key.

```
$ heroku login
```

Clone the repository

Use Git to clone pbc-buytogether's source code to your local machine.

```
$ heroku git:clone -a pbc-buytogether  
$ cd pbc-buytogether
```

Deploy your changes

Make some changes to the code you just cloned and deploy them to Heroku using Git.

```
$ git add .  
$ git commit -am "make it better"  
$ git push heroku master
```

04

PART FOUR

心得展望

複習
基礎觀念

學習
相關套件

BI
商業智慧

短期

- 修正版面小bug
- 即時抓取資料、即時更新圖表

中期

- 使用NLTK及wordcloud製作文字雲

長期

- 應用於其他資料，製作更多商業智慧分析

THANKS

謝 謝 聆 聽