# Undergraduate AI Capstone NYCU Spr2022 Final Group Project

0816175  高靖  0816006 郭冠德  0816081  張以廉

Our project is about NLP. We try to collect reviews on google map & try to do some prediction of the reviews rating.



Our original slides of our short present:

## Outcomes

- Positive/negative predictions.
- Furthermore, do the predict rating of a comment.
- Dataset:

1. Yelp reviews, which is an open source dataset, containing vast variety of comments to train the model.

2. Crawler from Google map reviews, using API or selenium to collect data from Google to train the model.
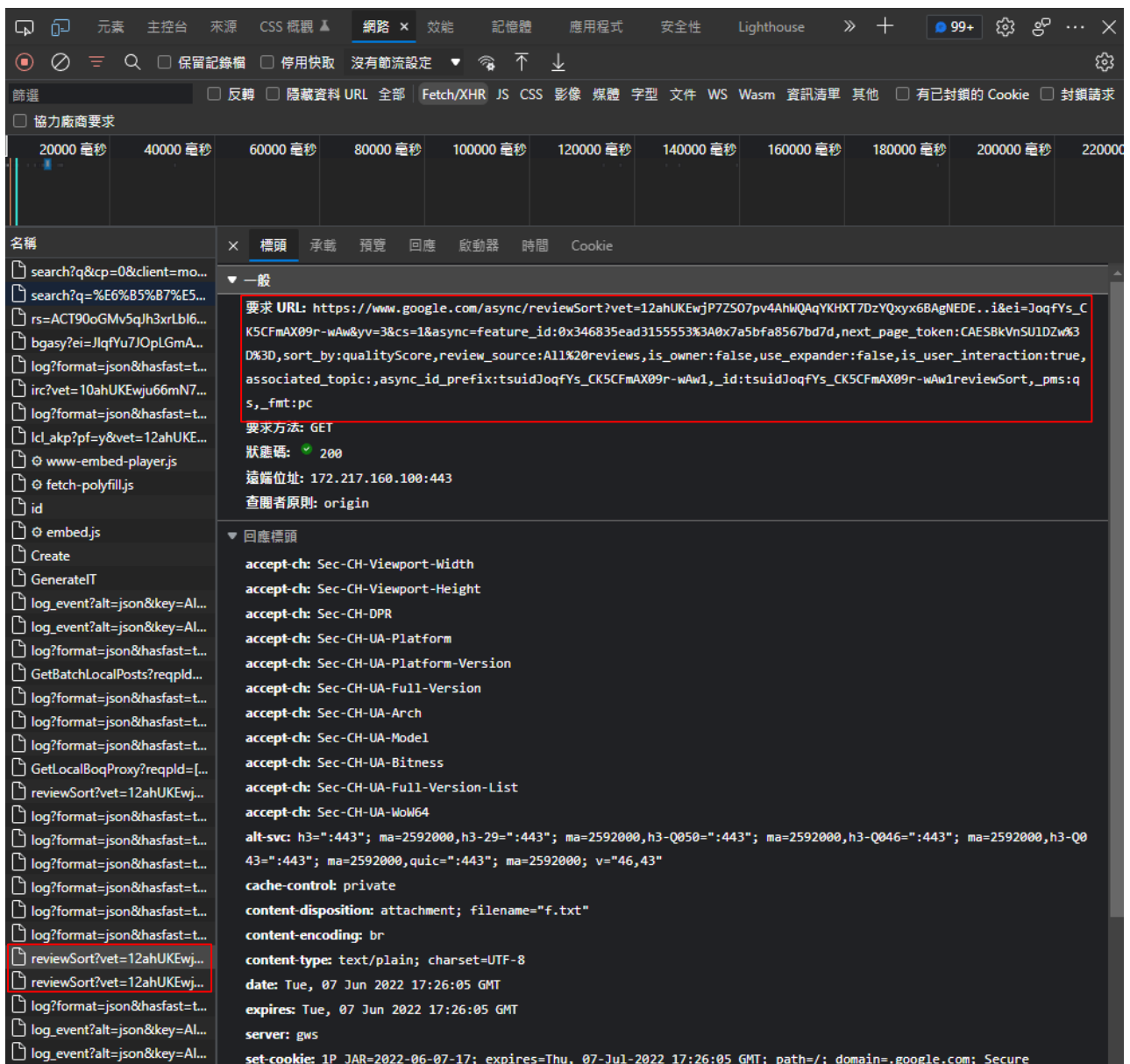
# Data collection:

So, originally, we try to use some online open-source dataset. But in the end, we try to make it our own way. So instead to use some existing data on the Internet, we decide to do crawler on our own.

We were trying to use Yelp dataset(Yelp Dataset) which is a famous Open dataset which includes personal ,educational, and some academic purpose and reviews & user data dataset. It has JSON & SQL form. However, we think do the Chinese part is more interesting. So we collect the google reviews of some famous store in Taiwan. So, how wee collect our data?

By developer tool, we can check the console of Internet/Fetch/XHR/Headers, to see how the google map backend call the backend api. So as we scroll down the reviews, we can check there's an request URL which is called repeatedly.

So we can cut in this point, just to get the reviews by calling the same URL.(next page)

We can compare the 2 URL of the request, to do the loop

crawler. However, I found that

URL1:https://www.google.com/async/reviewSort?vet=12ahUKEwjP7ZSO7pv4AhWQAqYKHXT7Dz

YQxyx6BAgNEDE..i&ei=JoqfYs_CK5CFmAX09r-

wAw&yv=3&cs=1&async=feature_id:0x346835ead3155553%3A0x7a5bfa8567bd7d,next_page_toke

n:CAESBkVnSUlDZw%3D%3D,sort_by:qualityScore,review_source:All%20reviews,is_owner:fa

lse,use_expander:false,is_user_interaction:true,associated_topic:,async_id_prefix:t

suidJoqfYs_CK5CFmAX09r-wAw1,_id:tsuidJoqfYs_CK5CFmAX09r-

wAw1reviewSort,_pms:qs,_fmt:pc

URL2:`https://www.google.com/async/reviewSort?vet=12ahUKEwjP7ZSO7pv4AhWQAqYKHXT7Dz`

`YQxyx6BAgNEDE..i&ei=JoqfYs_CK5CFmAX09r-`

`wAw&yv=3&cs=1&async=feature_id:0x346835ead3155553%3A0x7a5bfa8567bd7d,`<span style="color:red">`next_page_toke`</span>

<span style="color:red">`n`</span>`:CAESBkVnSUl`<span style="color:red">`GQQ`</span>`%3D%3D,sort_by:qualityScore,review_source:All%20reviews,is_owner:fa`

`lse,use_expander:false,is_user_interaction:true,associated_topic:,async_id_prefix:t`

`suidJoqfYs_CK5CFmAX09r-wAw1,_id:tsuidJoqfYs_CK5CFmAX09r-`

`wAw1reviewSort,_pms:qs,_fmt:pc`

The crucial part is about the param: **next_page_token.**

Rest part are the same. So here comes the question, how can we decide find the pattern of the URL? Because the token seems to be random. After some research, I found that it do have regular pattern is that, **each store start with DZw, GQQ,IZw,LQQ**……. So the token seems no rules, but can we get the next page token before do the crawler? The answer is negative, we still need to find the token manually and, then, apply to every store.

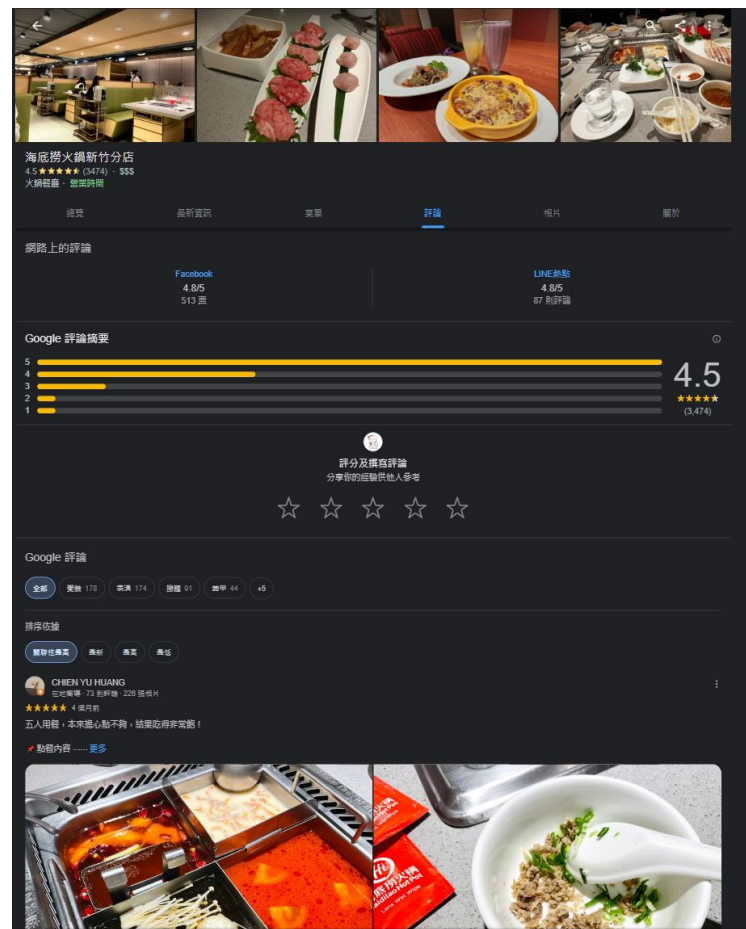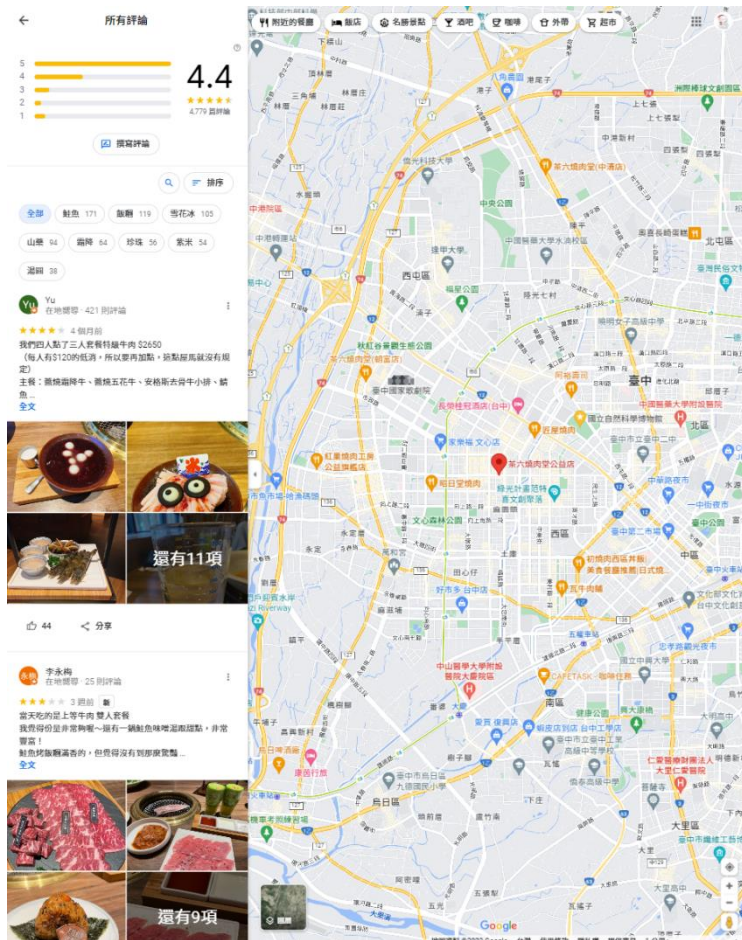So, it may be not a 'smart' way to collect way to collect data but a 'not bad' way, if you do have time to do so. For us, we want to find more 'intelligent' way that not to do any manual work. So here's work around we find:

Instead of using the reviews on web page of the google store We change to using the map application interface, because the

api called here is different from web interface.



The left hand side is  listentitiesreviews?authuser=0&hl=zh-TW&gl=tw&pb=!...4b1!5b1!6b1!7b1!5m2!1s4!at

While the right hand side is  reviewSort?vet=12ahUKEwj...

So the map version has rules to follow to do the crawler, by

iterating one of the params i, which loop as we scroll the page.

So next part is about how to do the crawler.

Initially, we try to use selenium/webdriver, but we finally decide

to call the api of backend. And, google's anti-crawler still block

my continuous request in about 1 min. So, to avoid such

circumstance, I add the header(language for not to get reviews

other than Chinese, or we would receive English review)

```
headers = {
    'user-agent':
    'Mozilla/5.0 (Macintosh; Intel Mac OS X 10_15_6) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/85.0.
    'Accept-Language':'zh-TW'
}
```

Hope not be blocked again, however, it's of no use. So I use

VPN, try not to request from same ip for too many times. But it

still of no use, because our demand is so big that these method

can't prevent we from being blocked. So, we look up for

another third party tool. Scraper API.

It can automatically change from millions of

proxy server, and do the retry also. So I use the

tool and do implementation in python. (all codes will be on

github also).

| author | grade | comment |
|--------|-------|---------|
| 高慧嘉 | 5 | 這家當地人超推<br>珍珠奶茶是奶蓋，很奶很棒<br>火鍋料非常多　料多超實在<br>用餐時間人很多，需要等一下，用餐有兩樓，廁所在一樓。<br><br>很少小火鍋能吃這麼飽的，很推 |
| Scottie Zhan | 5 | 每鍋價格180-200，料好實在，比起同性質的小火鍋，　CP值超高。肉都有6-7片，火鍋料至少有七種，品質都還不錯，鍋底鋪滿高麗菜，上面都會附上 |

So I use the tool to get data, and here comes the last question,

It's too slow. Only 2500 reviews takes few hours to craw. So I change the way we craw to multi-thread. Use the python module, `import concurrent.futures`,we create all URL as a list, and view them as a pool, then the workers will execute the crawler from pool once they're free.

```python
with concurrent.futures.ThreadPoolExecutor(
        max_workers=NUM_THREADS) as executor:
    executor.map(scrape_url, list_of_urls)
```
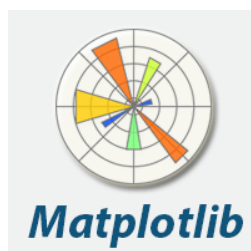
After crawling, export the data as excel.

So, that's about the whole process of the data collection.

## NLP part:

For data analysis, preprocessing, I use 3 famous tool in Machine Learning & in NLP field. The Pandas & matplotlib & numpy.
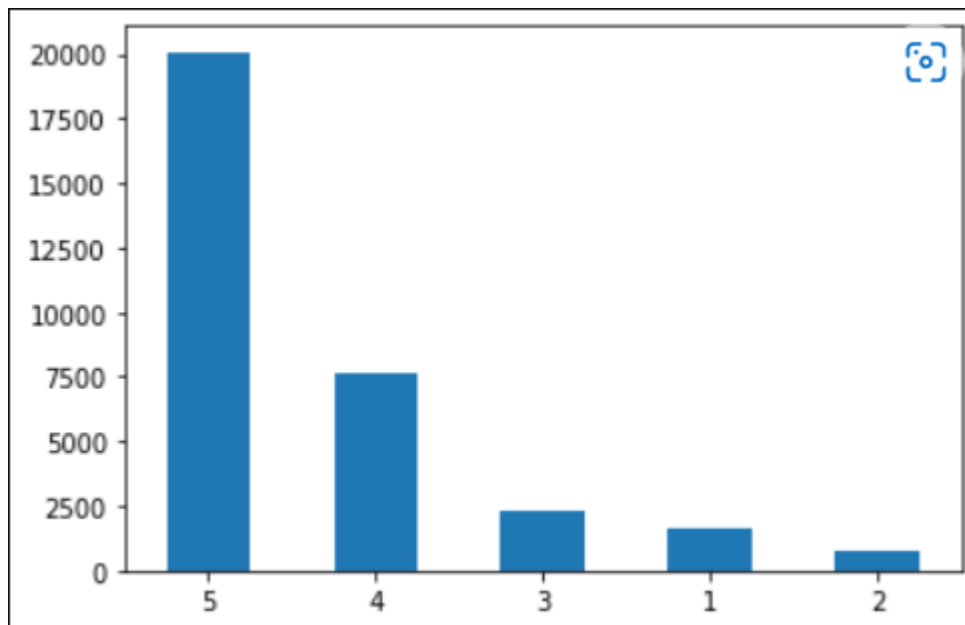
So, first, let's import our data using pandas.read_excel()



So the data is import as dataFrame like this:

```
        grade                                                       comment
0           5   這家當地人超推\n珍珠奶茶是奶蓋，很奶很棒\n火鍋料非常多   料多超實在\n用餐時間人很多，...
1           5   每鍋價格180-200，料好實在，比起同性質的小火鍋， CP值超高。肉都有6-7片，火鍋料至...
2           5   來台中必吃的愛店\n聽說牛奶鍋超讚\n但我是點泰式酸辣鍋\n香料味很重很喜歡\n那個酸跟辣很...
3           5                 泡菜鍋真的很夠味，魚酥給的大方，吃了開心，下次想買回家自己煮加料
4           5   以前就常常跟朋友聚餐來吃～今天吃吃辣味起司豬肉鍋～對於不敢吃單吃起司鍋的我 這個辣味起司鍋真...
...        ...                                                      ...
44735       4   吃過滿多口味的，整體CP值都滿高的，用料實在，價格算平實，值得推薦，唯一缺點就是買的人很多要排隊。
44736       4       熱熱吃好吃，外面奶粉撒超多，如果要排隊我是不會買的  ，不排還可以買來吃一下
44737       4                        人多常常會排隊，口味多樣但有些提早賣完！
44738       4   這個好吃!超脫了國外那種甜甜圈的感覺!我周六,日下午三點多去的,排隊好長!國內外旅客都有,我...
44739       4       味道不錯，冷的或是熱的都很好吃，缺點就是要花太多時間去排隊，最少都要排一小時起跳。

[44740 rows x 2 columns]
```

Let's using matplotlib to see how our data distribution:



There's huge difference between amounts of 5 starts & others,

Actually it's normal, because real trend on Google maps is that

Taiwanese seem love to give 5 stars on revies.

Let's hand over it in next version, let's move on first.

Next step, we're going to deal with the row without comment.

Because some people have no comments but only stars like
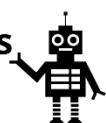
Sasa Chen:

(this step is address reviews below)

So our crawler will only get None value in excel. So we replace 'None' to Nan, then use pandas.dropna() to get rid of those empty value. And after dealing with the some basic operation. We can start to train our model.

Let's see positive/negative analysis first, So we change the reviews into 2 categories, rating >3 are positive, others are negative. And we choose to use bert as our model, which is a bidirectional encoder from transformers. So using pyTorch & CUDA, I build the model on my own environment. Finishing the fine tune part of bert is to set the label as int, and use f1score, confusion matrix to present the outcome. I use sklearn's package.

**Simple Transformers**

**Transformers**
Documentation

| | comment | status |
|---|---|---|
| 0 | 這家當地人超推\n珍珠奶茶是奶蓋，很奶很棒\n火鍋料非常多 料多超實在\n用餐時間人很多，... | 1 |
| 1 | 每鍋價格180-200，料好實在，比起同性質的小火鍋，CP值超高。肉都有6-7片，火鍋料至... | 1 |
| 2 | 來台中必吃的愛店\n聽說牛奶鍋超讚\n但我是點泰式酸辣鍋\n香料味很重很喜歡\n那個酸跟辣很... | 1 |
| 3 | 泡菜鍋真的很夠味，魚酥給的大方，吃了開心，下次想買回家自己煮加料 | 1 |
| 4 | 以前就常常跟朋友聚餐來吃～今天吃吃辣味起司豬肉鍋～對於不敢吃單吃起司鍋的我 這個辣味起司鍋真... | 1 |
| ... | ... | ... |
| 32315 | None | 1 |
| 32316 | None | 1 |
| 32317 | None | 1 |
| 32318 | None | 1 |
| 32319 | None | 1 |

So the original data would look like this after doing 2-label.

Using train_test_split in sklearn.model_selection, I use test data

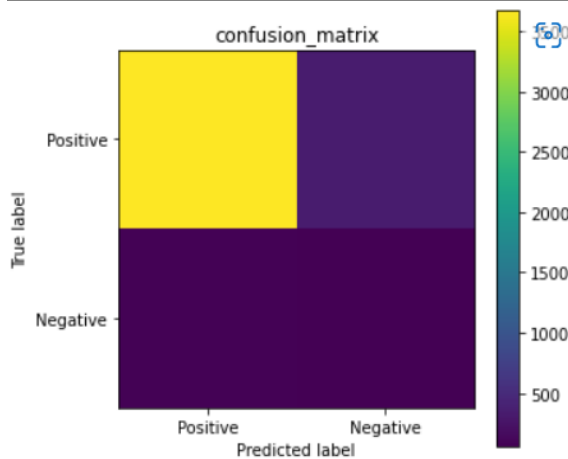Size of 4200. And most of data are positive, which is quite

normal and fulfill original trend. So let's do 5-label version.

With the same processing before, we change the None to Nan

then remove them from data. And minus 1 from original

score, Cause the bert's label need to start from 0, instead of 1.
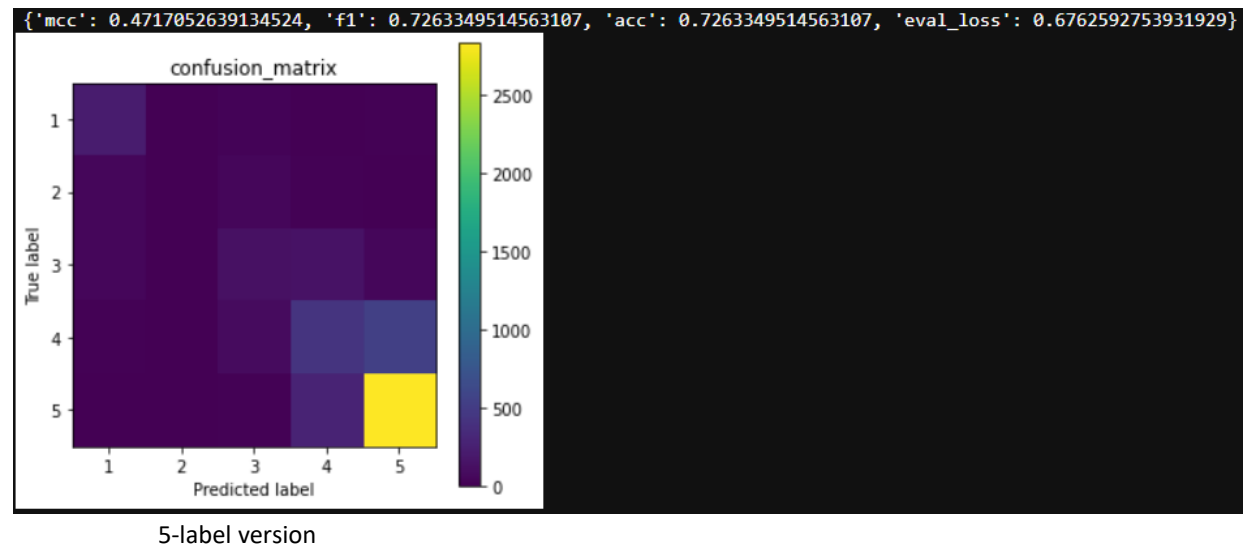
{'mcc': 0.8087241072066188, 'tp': 3669, 'tn': 329, 'fp



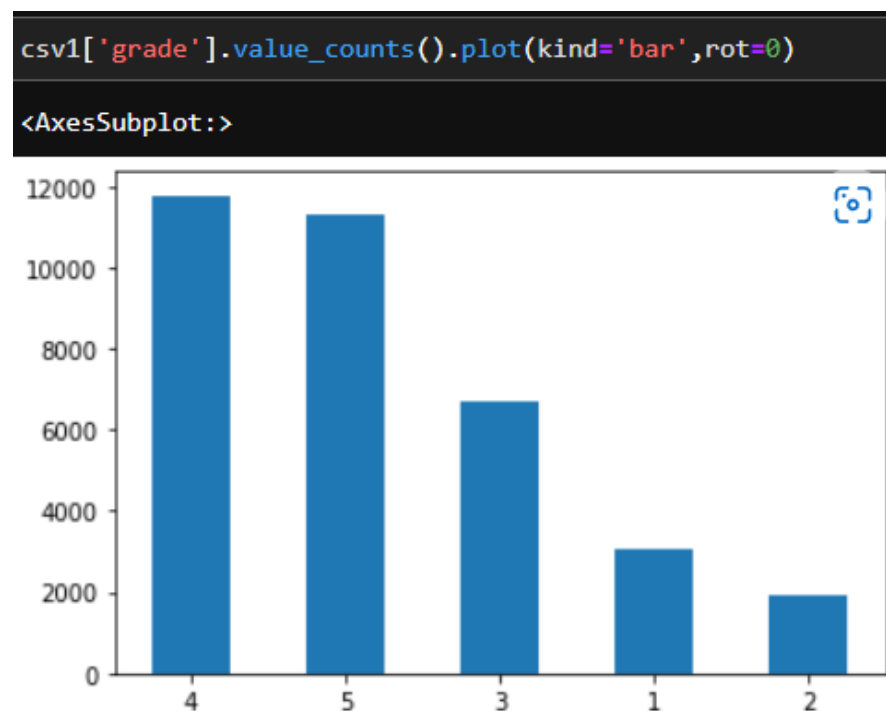2-label version

```
csv1[['grade']] = csv1[['grade']]-1
```



scikit learn

Here's the outcome of 5-label:

{'mcc': 0.4717052639134524, 'f1': 0.7263349514563107, 'acc': 0.7263349514563107, 'eval_loss': 0.6762592753931929}



5-label version

We can see that most of the data still appears in label 5-stars.

So to get a more accurate outcome, we try to make the

numbers of 5-stars not exceed others too much. So we remove

10000 reviews of 5-stars. And the data distribution now:

```
csv1['grade'].value_counts().plot(kind='bar',rot=0)

<AxesSubplot:>
```

Also, we remove the reviews that contains translation.

Because some foreigner might comment in language other than

Chinese. Like the example of a comment of 海底撈:

（由 Google 提供翻譯）來自中國的一家火鍋連鎖店。
毫無疑問，他們的服務。

說明細節和建議菜。
也是酋長的特別
您最多可以選擇四種湯。
訓練有素的服務員一直都在要求和服務。感到有些尷尬和不安。

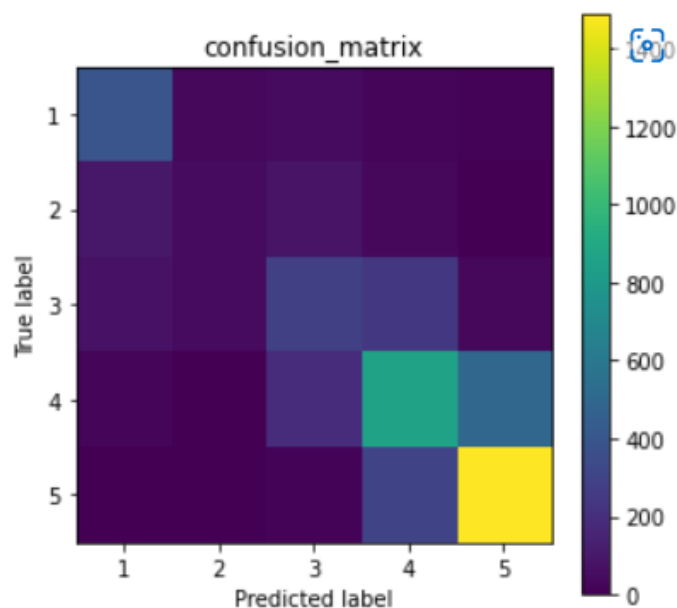附帶麵條搖擺器和換臉秀。
休息和等候區提供小吃和飲料。
還提供兒童房 / 美甲沙龍。
質量很好，$$$相當高的價格。

（原始評論）
A hotpot chain store from China.
No doubt of their service.

Explain details and suggestion dishes.
Also chief's special.
You can choose up to four different soups.
Waiters been trained to come ask and serve all the time. Feel a little embarrassing and disturbing.

Come with noodle swing and change face show.
Rest and wait area with snacks and drinks to serve.

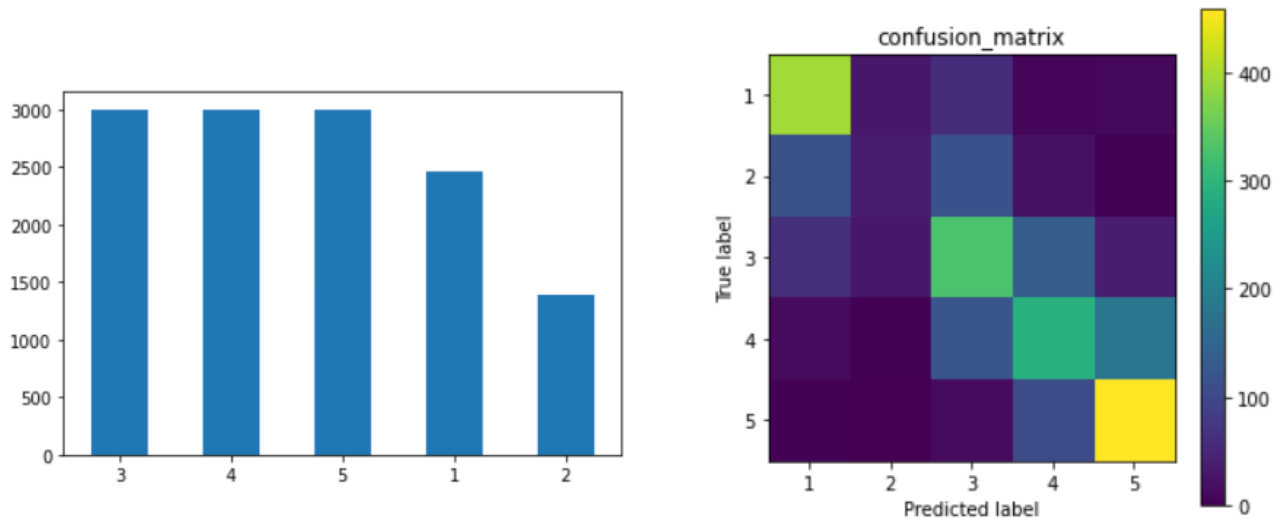So if comment contains 'Google 提供翻譯', we remove it.

And here's the consequence:



We can see the diagonal is more clear, which means our model

is getting better after the label become more balanced.

So we change the data distribution again, make data more

balanced. And here's our final outcome:



We can see the diagonal is clear but a little bit divergent.

And start-2 seems the most bad, so let's check the original data:

東西好吃是好吃，但是與預約了，還是要稍微等一下，而且建議店家可以改座位結帳，不然連結帳都還要排隊有點無奈

動線很好，干貝好吃，蒸蛋不錯，螃蟹腥味很重完全不行，甜點太普通蛋糕也太少，冰淇淋太軟，其他還好，這樣價位CP 值不太推

流浪漢不會跳舞　　還可以啦！幫家人、朋友慶生會，OK

All of these comments are not 'that' bad I think. But they're

only worth of 2 stars. Basically, the comment is too subjective &

someone will comment in strange viewpoint to affect our

prediction.

So comments seem to be neutral can be low stars also, which affect our model.

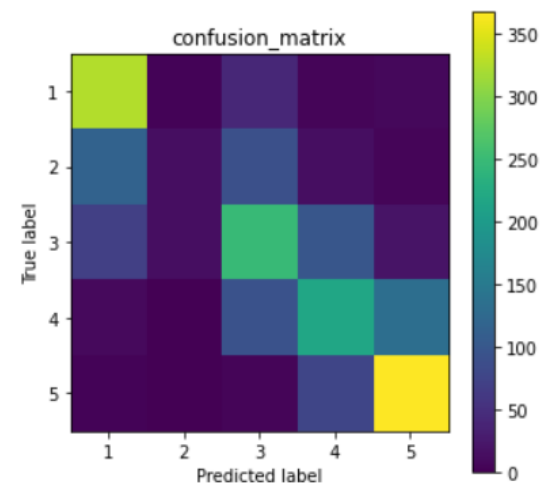To do it more precise, I filter the comment that's too short like: '難吃','還不錯','流浪漢不會跳舞', only leave the comment whose length is more than 15.

And the outcome looks much the same as the previous one.

So I think the most of the reason is that reviews are subjective, and we can analysis in such way to get most of the result correct(70% of the majority people).


confusion_matrix

Last part, I do a manually input reviews, to check the outcome by your own.

```
input your review of 1 star 有夠難吃....菜都臭掉 還有蟑螂 店內衛生不佳 真的很需要加強
input your review of 2 star 用餐體驗不佳 菜也都不太好吃 雖然沒有很貴 但整體還是差評 需要改進的地方有非常多 肉質不好 同樣價位寧願
去吃石二鍋
input your review of 3 star 我認為是還可接受的餐廳 普普通通的用餐環境&體驗 大部分都還可以接受啦 有一點點小貴就是 希望價錢可以調整
一下 會更棒
input your review of 4 star 還不錯的餐廳！位置寬敞舒適 我認為菜品大部分也都沒甚麼問題 只剩價錢有一點點小小貴了 但很適合家庭一起來
用餐！好吃！
input your review of 5 star 超級讚...無可挑剔的餐廳 我認為每一道菜品都非常完美 很好吃 很讚的體驗！
```

```
Running Evaluation: 100%|████████████████████████████████| 1/1 [00:00<00:00, 1.78it/s]

your review of star 1 is rated as 1 stars
your review of star 2 is rated as 2 stars
your review of star 3 is rated as 3 stars
your review of star 4 is rated as 4 stars
your review of star 5 is rated as 5 stars
```

The reviews are all from my own opinion when I share

comment in real condition, and the outcome is good!

So that's all about our project. In the process, most of time

spending on the model, tons of error and fix the crawler bug.

It's a long journey & We learn a lot in the project.

Hope you enjoy our final project!

Group 23

Cute Bert 😊 (easy to learn, hard to utilize & debug, my

environment was a chaos after the project)