



## ▼ Review JSON

[illegible]

	user_id		business_id	rating	date
	0	hG7b0MtEbXx5QzbzE6C_VA	ujmEBvifdJM6h6RLv4wQlg	1.0	2013-05-07 04:34:36
	1	yXQM5uF2jS6es16SJzNHfg	NZnhc2sEQy3RmzKTZnqtWQ	5.0	2017-01-14 21:30:33
	2	n6-Gk65cPZL6Uz8qRm3NYw	WTqjgwHIXbSFevF32_DJVw	5.0	2016-11-09 20:09:03
	3	dacAlZ6fTM6mqwW5uxkskg	ikCg8xy5JIg_NGPx-MSIDA	5.0	2018-01-09 20:56:38
	4	ssoyf2_x0EQMed6fqHeMyQ	b1b1eb3uo-w561D0ZfCEiQ	1.0	2018-01-30 23:07:38

 (1637138, 192606)

<https://colab.research.google.com/drive/1UcSfWf7-wLd65GxUSmF5Ye13syHPC-Ge#printMode=true>

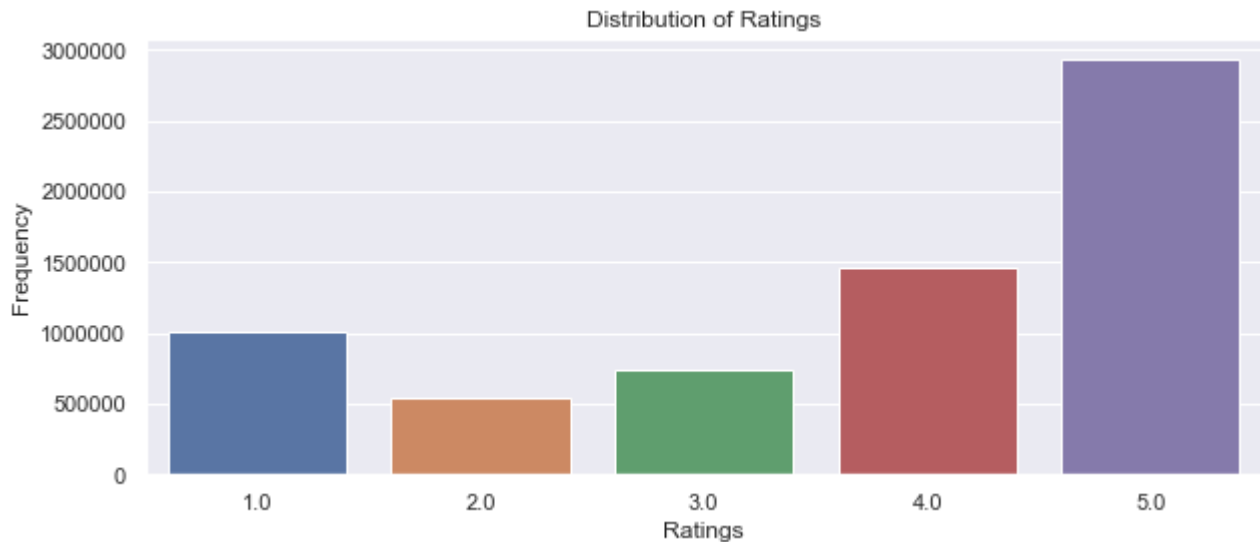
```
len(ratings_active.user_id.unique()), len(ratings_active.business_id.unique()))
```

(1351008, 172970)

## ▼ Distribution of Ratings

```
plt.figure(figsize=(10,4))
sns.set(style='darkgrid')
ax = sns.countplot(ratings['rating'])
# plt.title('Distribution of rating')
plt.title('Distribution of Ratings')
plt.ylabel('Frequency')
plt.xlabel('Ratings')
```

Text(0.5, 0, 'Ratings')



```
ratings_wordcloud = ratings.sample(n = 100000)
```

```
from wordcloud import WordCloud
cloud = WordCloud(width=1440, height= 1080,max_words= 200, background_color = 'white').genera
plt.figure(figsize=(20, 15))
plt.imshow(cloud)
plt.axis('off');
```





```
# sample input
# {"business_id":"1SWh84yJXfytoVILX0AQ","name":"Arizona Biltmore Golf Club","address":"2818
```

```

import pandas as pd
import json
from tqdm import tqdm
line_count = len(open("data/business.json", encoding="utf8").readlines())
business_ids, names, stars, review_counts = [], [], [], []
with open("data/business.json", encoding="utf8") as f:
    for line in tqdm(f, total=line_count):
        blob = json.loads(line)
        business_ids += [blob["business_id"]]
        names += [blob["name"]]
        stars += [blob["stars"]]
        review_counts += [blob["review_count"]]
#        compliment_counts += [blob["compliment_count"]]
business = pd.DataFrame(
    {"business_id": business_ids, "name": names, "star": stars, "review_count": review_counts}
)

```



100% | 192609/192

## ▼ Restaurants' Average Ratings

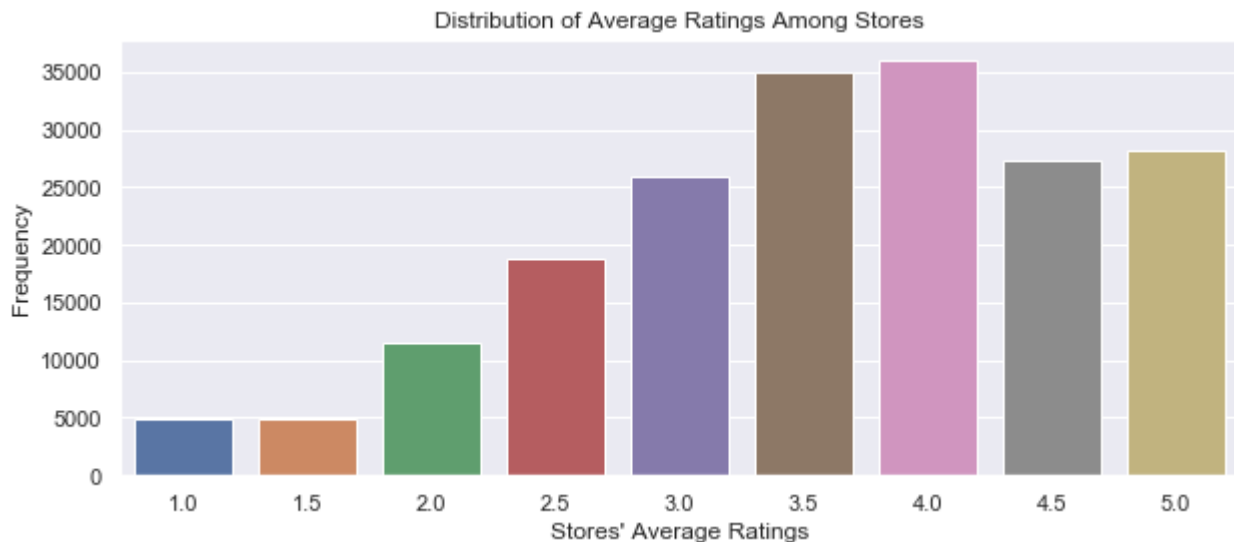
```

plt.figure(figsize=(10,4))
sns.set(style='darkgrid')
ax = sns.countplot(business.star)
# plt.title('Distribution of rating')
plt.title('Distribution of Average Ratings Among Stores')
plt.ylabel('Frequency')
plt.xlabel("Stores' Average Ratings")

```



Text(0.5, 0, "Stores' Average Ratings")



## ▼ Top Name of Restaurant in Yelp

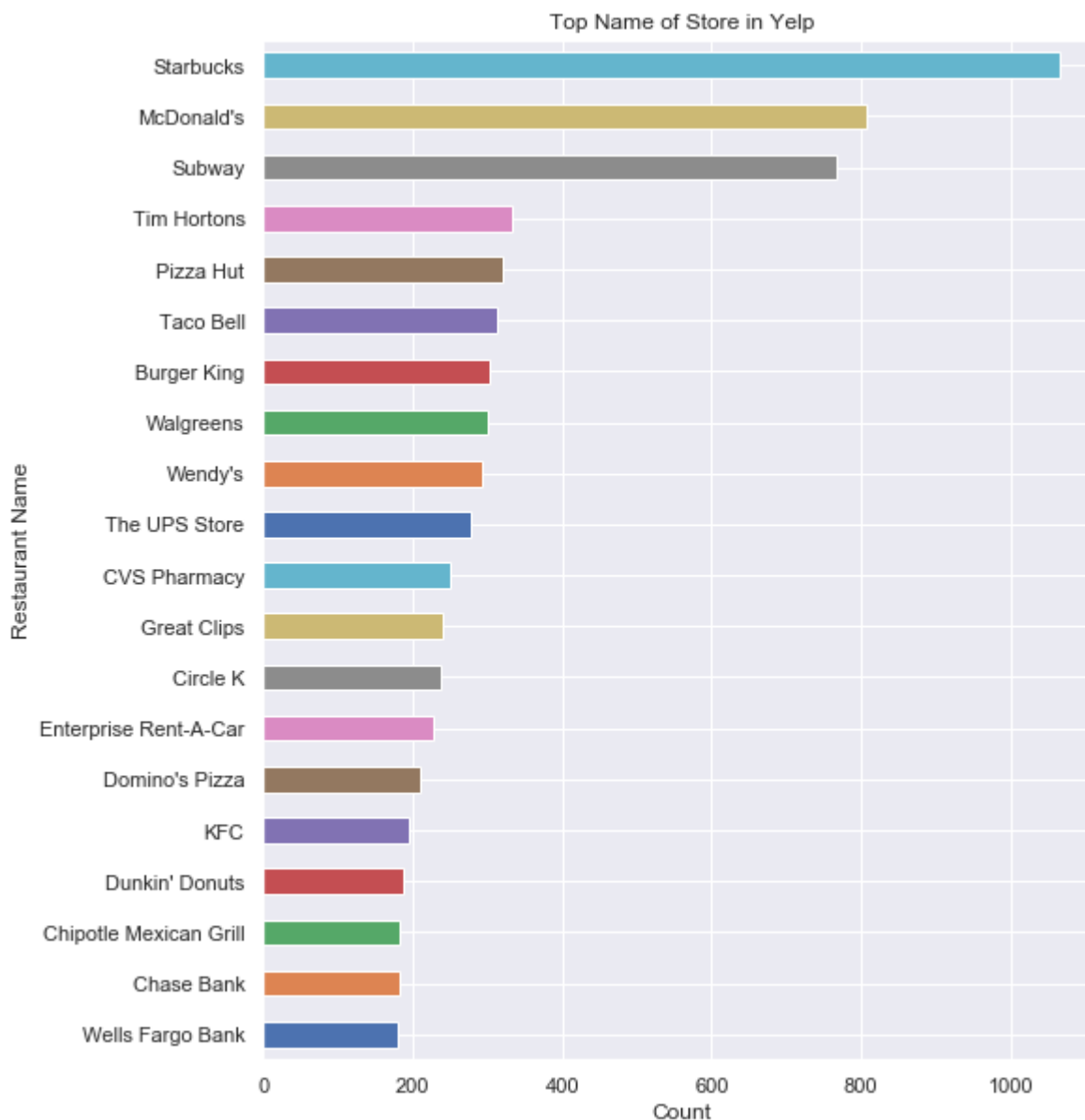
```

a = business.groupby('name').size().sort_values(ascending = False)[:20]
sns.set(style='darkgrid')
plt.figure(figsize=(8,10))
a.sort_values(ascending=True).plot(color=sns.color_palette(), kind='barh')
plt.title("Top Name of Store in Yelp")
plt.ylabel("Restaurant Name")
plt.xlabel("Count")

```



Text(0.5, 0, 'Count')



## ▼ Tip JSON

```

line_count = len(open("data/tip.json", encoding="utf8").readlines())
user_ids, business_ids, texts, compliment_counts = [], [], [], []
with open("data/tip.json", encoding="utf8") as f:
    for line in tqdm(f, total=line_count):
        blob = json.loads(line)

```

100% | 1223094/12230





