

Exploratory Data Analysis of Yelp Data

Sophie Wang, Elyse Cheung-Sutton, Aneri Dand, Vaishnavi Kashyap, Qianyun(April) Li

Background

5 collections:

1. Business
2. Users
3. Check-ins
4. Reviews
5. Tips

Initial EDA:

1. Only 37 states including Canadian states
2. 2011 had the most reviews

```
1 db.business.findOne()
2
3 {
4   "_id" : ObjectId("5fff2a5243bfad3bc4878a2e"),
5   "business_id" : "YzvJg0SayhoZgCljUJRF9Q",
6   "name" : "Carlos Santo, NMD",
7   "address" : "8880 E Via Linda, Ste 107",
8   "city" : "Scottsdale",
9   "state" : "AZ",
10  "postal_code" : "85258",
11  "latitude" : 33.5694041,
12  "longitude" : -111.8902637,
13  "stars" : 5.0,
14  "review_count" : 4.0,
15  "attributes" : {
16    "GoodForKids" : "True",
17    "ByAppointmentOnly" : "True"
18  },
19  "categories" : "Health & Medical, Fitness",
20  "hours" : null
21 }
```

```
12 db.user.findOne()
1
2 {
3   "_id" : ObjectId("5ffe9726a7cec2080b58cea"),
4   "user_id" : "nltvfPzc8eg1qv92iDIaw",
5   "name" : "Rafael",
6   "review_count" : 553.0,
7   "yelping_since" : "2007-07-06 03:27:11",
8   "useful" : 628.0,
9   "funny" : 225.0,
10  "cool" : 227.0,
11  "elite" : "",
12  "friends" : "oeMvJh94PiGQnx_6GIndPQ, wm1z1Pa3KvHgSDRKfwhfDg, IkR1b6Xs91PPW7pon7VW1g, A8Aq8f0-XvLBcyMk2G3dJ0, eEZM",
13  "fans" : 14.0,
14  "average_stars" : 3.57,
15  "compliment_hot" : 3.0,
16  "compliment_more" : 2.0,
17  "compliment_profile" : 1.0,
18  "compliment_cute" : 0.0,
19  "compliment_list" : 1.0,
20  "compliment_note" : 11.0,
21  "compliment_plain" : 15.0,
22  "compliment_cool" : 22.0,
23  "compliment_funny" : 22.0,
24  "compliment_writer" : 112,
25  "compliment_photos" : 1
26 }
```

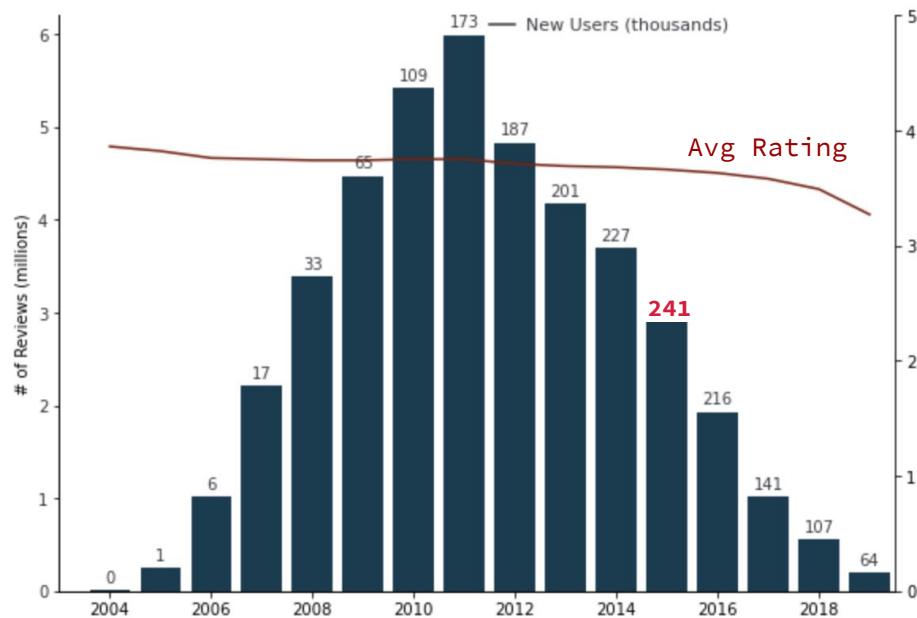
```
112 db.review.findOne()
1
2 {
3   "_id" : ObjectId("5ffe5e24b0b5b2799bcbcf90d"),
4   "review_id" : "-DQb3fBYScdqy1_9irsJtA",
5   "user_id" : "CWenIpiWvvcBJPXF5A60Q",
6   "business_id" : "oS96aJ1HFWcFALGHKKXjAw",
7   "stars" : 5.0,
8   "useful" : 5.0,
9   "funny" : 2.0,
10  "cool" : 3.0,
11  "text" : "So happy to have found a little piece of France in Pittsburgh! As a francophile and food-lover, i was ex",
12  "date" : "2014-04-19 17:37:58"
13 }
```

How many Reviews and new Users between 2004-2019?

- 2015: most new Users
- 2011: most reviews
- Average review rating constant
- Expect 2019 to have most reviews and new users given that Yelp is always growing

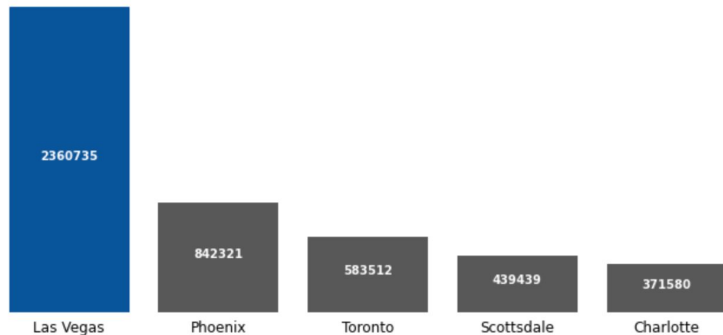


Yelp Reviews through the Years



Top 5 Cities w/ Most Reviews & Categories Most Reviewed in Top City

Top 5 Cities with Most Yelp Reviews



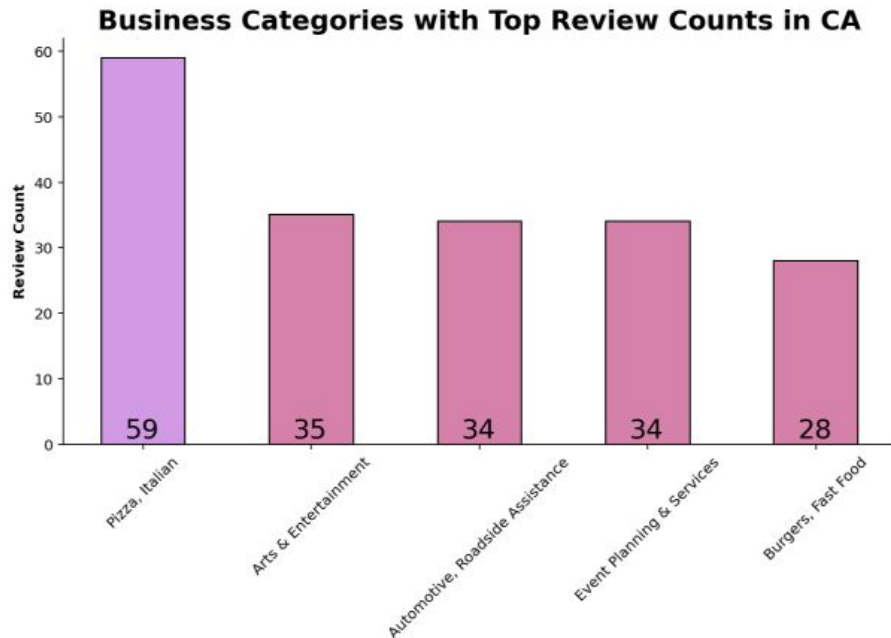
- Las Vegas ranked #1 with 2,360,735 reviews - not surprising
- Different versions of “Mexican Restaurant” and “Nail Salon, Beauty” made up the Top 5 Categories reviewed in Las Vegas - a bit surprising, would have thought “Bars/Clubs” would be in the top 5 reviewed

```
{
  "category" : "Mexican, Restaurants",
  "num_reviews" : 143.0
}
{
  "category" : "Restaurants, Mexican",
  "num_reviews" : 133.0
}
{
  "category" : "Food, Coffee & Tea",
  "num_reviews" : 114.0
}
{
  "category" : "Nail Salons, Beauty & Spas",
  "num_reviews" : 108.0
}
{
  "category" : "Beauty & Spas, Nail Salons",
  "num_reviews" : 107.0
}
```

Which business categories in California have the most review counts?

— — —

- The business categories with top review counts are Italian food, entertainment, automotive, event Planning and fast food in California (based on this dataset).
- Caveat: This dataset does not represent the comprehensive yelp data within US. (Number of business in CA is a lot smaller than other states). So it's hard to say in reality which category in CA has the most reviews. Also high review count does not equal high popularity.



What do people complain most about restaurants?

Top 5 Common Nouns in Negative Reviews

Food

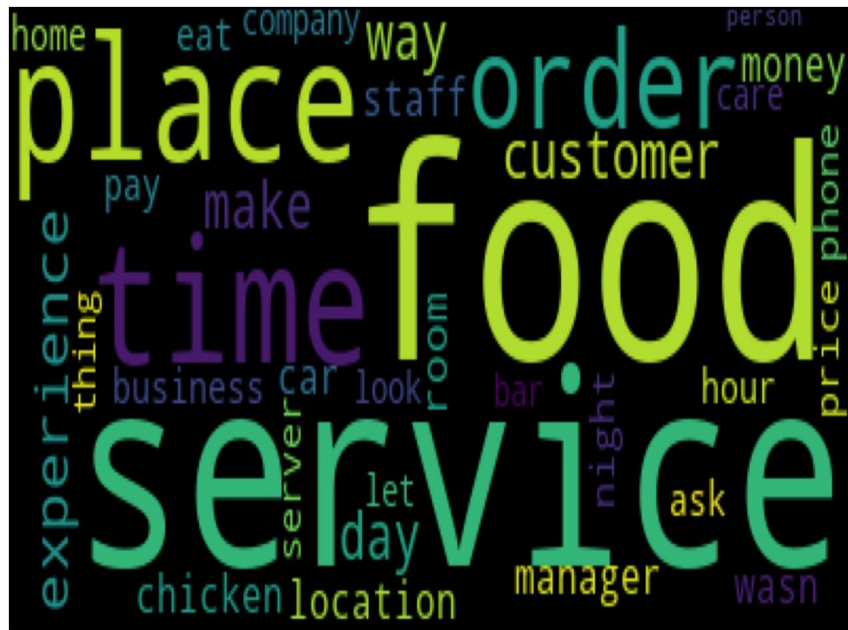
Service

Place

Time

Order

- For restaurants, no wonder terrible food is the most frequent reason for negative reviews.
- Besides, restaurants should also pay attention to service, place, waiting time and order service in case of negative reviews.



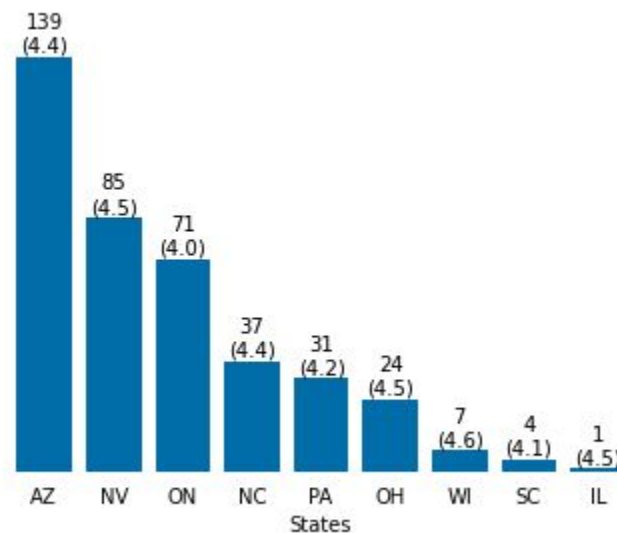
Which states have the highest number of pilates studios?

— — —

- Based on EDA, it appears that Arizona has the highest number of Pilates studios.
- However, the veracity of the trend shown in the adjacent graph comes into question as the Yelp dataset contains information only about 9 US states.
- California and Texas (US states with the highest number of gyms*) are not included in the data.

*Source: <https://www.exercise.com/learn/ten-fittest-states-in-the-us/>
<https://www.ibisworld.com/industry-statistics/number-of-businesses/gym-health-fitness-clubs-united-states/>

Pilates Studio Count and Average Rating by State in the US



Note: Figures in parenthesis represent average ratings

Data Pipeline Efficiency

- From pymongo (not distributed computing setting) to spark, we achieve 95% runtime reduction.
- In Spark, using techniques to boost efficiency we were able to run and execute code in ~40% less time.

— — —

Key Takeaways

- Running multiple jupyter notebooks on the same cluster was very slow owing to memory constraints
- Having the flexibility to control the number of partitions helps improve query efficiency
- Large data does not always translate to better / accurate insights

