



Big Data Analytics Symposium 2022

Yelp Data Analysis - NYC Eating Team

PRESENTED BY Junhua Liang, Virginia Wu, Xin Lu, Wenhao Li

05/05/2022

Yelp Data Analysis

Abstract

Feast is important for everyone in the world, especially in this era as the development of the society's productivity. People who travel to a new city usually want to have a taste of the best restaurant in this area, and those who want to repair their house's roof would also check the rating of the store. For merchant, they also want to open their restaurants in the best location in the city. Thus, it's important to have related analytics to help them. One direction to have has analytics is using people's reviews and their check-in information to get the pattern of consuming in the area and build related models to help recommend restaurants for people and find location for merchant. In this project, we will use Yelp's open datasets to visualize specific information and build model to for recommendations and rating predictions.

Platform(s) where the application runs: NYU Peel Cluster

Yelp Data Analysis

Motivation

Who are the users of this analytics?

People who wrote their reviews and provided feedback to the stores they went to will be the users of this analytics.

Who will benefit from this analytics?

People who need to go to a high rate restaurant.

People who use Yelp app, the system will provide the best recommended stores for them.

People who go to explore a new city and want to have a taste of the best stores.

Why is this analytics so important?

1. The analytics can retrieve people's consuming habits, and recommend the most perfect one for them.
2. The analytics can obtain consuming distribution in each city, and business man could open stores in the good location.

Yelp Data Analysis

Goodness

What steps were taken to access the ‘goodness’ of the analytics?

1. For the recommendation systems, our team verified the recommended items, compared it to the user’s past consuming items. If the recommended one is meaningful, then we could state that our analytics is good.
2. For the data visualization and the pattern we retrieve from it, our team will regard it as a good result if the pattern match normal people’s thinking.

Data Sources

Data Sources

Dataset name	Description	Size of Data
yelp_academic_dataset_user.json	Users' information, including user_id and metadata.	4.34 GB
yelp_academic_dataset_business.json	Restaurant's information, including business_id and metadata.	118.9 MB
yelp_academic_dataset_review.json	Review information on specific user and restaurant.	5.34 GB
yelp_academic_dataset_checkin.json	Restaurant's business ID and their check-in date.	287 MB
yelp_academic_dataset_tip.json	Tip's information on specific user and restaurant.	180.6 MB

Yelp Data Analysis

Data Sample: yelp_academic_dataset_user



NYU

Yelp Data Analysis

Data Sample: yelp_academic_dataset_business

address	attributes	business_id	categories	city	hours is_open	latitude	longitude	name postal_code review_count stars state
[1616 Chapala St, ... , , , , , , , True... Pns214eNsf08kk83d... Doctors, Traditio...			Santa Barbara	null	0 34.4266787 -119.7111968 Abby Rappoport, L...	93101	7 5.0 CA	
[87 Grasso Plaza S... , , , , , , , True,... mpf3x-BjTdTEA3yCZ... Shipping Centers,...			Affton [8:0-18:30, 0:0-0...	1 38.551126 -90.335695 The UPS Store 63123	15 3.0 MO			
[5255 E Broadway Blvd , , , , , , , True,, T... tUfFrWirkIKI-TansV... Department Stores...			Tucson [8:0-23:0, 8:0-22...	0 32.223236 -110.880452 Target 85711	22 3.5 AZ			
935 Race St [], u'none', , , , MTSW4McQd7CbVtyjq... Restaurants, Food...			Philadelphia [7:0-21:0, 7:0-20...	1 39.955502 -75.1555641 St Honore Pastries 19107	80 4.0 PA			
101 Walnut St [], , , , , , , True,, T... mWMc6_wTDE0EUBKIG... Brewpubs, Breweri...			Green Lane [12:0-22:0, 12:0...	1 40.3381827 -75.4716585 Perkiomen Valley ... 18054	13 4.5 PA			
615 S Main St [], u'none', None, , , CF33F8-E6oudUQ46H... Burgers, Fast Foo...			Ashland City [9:0-0:0, 0:0-0:0...	1 36.269593 -87.058943 Sonic Drive-In 37015	6 2.0 TN			
[8522 Eager Road, ... , , , , , , , True,, T... n_0UpQx1hsNbPU1... Sporting Goods, F...			Brentwood [10:0-18:0, 0:0-0...	1 38.627695 -90.340465 Famous Footwear 63144	13 2.5 MO			
400 Pasadena Ave S [], null qKRm_2X51Yqxk3bt1... Synagogues, Religi...			St. Petersburg [9:0-17:0, 9:0-17...	1 27.76659 -82.732983 Temple Beth-El 33707	5 3.5 FL			
8025 Mackenzie Rd [], u'full_bar', , , , k0hLbqXX-Bt0vflop... Pubs, Restaurants...			Affton [null	0 38.5651648 -90.3210868 Tsevi's Pub And G... 63123	19 3.0 MO			
2312 Dickerson Pike [], u'none', , , , bBDDegkFA10tx9Lf... Ice Cream & Froze...			Nashville [6:0-16:0, 0:0-0:...	1 36.2081924 -86.7681696 Sonic Drive-In 37207	10 1.5 TN			
[21705 Village Lak... , , , , , , , True,, T... UjsufbvfyfONHeWdv... Department Stores...			Land O' Lakes [9:30-21:30, 9:30...	1 28.1904587953 -82.4573802199 Marshalls 34639	6 3.5 FL			
[], 'none', ('tou... eEOYSGkmpB90uNA71... Vietnamese, Food,...			Tampa Bay [11:0-14:0, 11:0-...	1 27.9552692 -82.4563199 Vietnamese Food T... 33602	10 4.0 FL			
8901 US 31 S [], 'none', ('tou... il_Ro8jwPlHresjw9... American (Traditi...			Indianapolis [6:0-22:0, 6:0-22...	1 39.6371332838 -86.127217412 Denny's 46227	28 2.5 IN			
15 N Missouri Ave , , , , , , , True... jaxMSoInw8Poo3xEm... General Dentistry...			Clearwater [], 7:30-15:30,, , ,	1 27.966235 -82.787412 Adams Dental 33755	10 5.0 FL			
2575 E Bay Dr [], u'none', ('ro... 0bPLkL0QhhP05kt1... Food, Delis, Ital...			Largo [10:0-20:0, 10:0-...	0 27.9161159 -82.7604608 Zio's Italian Market 33771	100 4.5 FL			
205 Race St [], 'full_bar', {'ro... MUTTqe8ugyMdBl186... Sushi Bars, Resta...			Philadelphia [13:30-23:0,, 13:...	1 39.953949 -75.1432262 Tuna Bar 19106	245 4.0 PA			
625 N Stone Ave , , , , , , , True,... rbMypr_Y1UbBx8ggh1... Automotive, Auto ...			Tucson [8:0-17:0, 0:0-0:...	1 32.2298719 -110.9723419 Arizona Truck Out... 85705	10 4.5 AZ			
712 Adams St [], , , , , , , True,, T... MOXSSHqrASOnhgbWD... Vape Shops, Tobac...			New Orleans [10:0-19:0, 10:0-...	1 29.9414679565 -90.129952757 Herb Import Co 70118	5 4.0 LA			
1241 Airline Dr [], u'none', {'to... R0eacJQwBebh05Rqg7... Automotive, Car R...			Kenner [8:0-17:0, 8:0-17...	1 29.981183 -90.2540123 Nifty Car Rental 70062	14 3.5 LA			
1224 South St [], u'none', {'to... R0eacJQwBebh05Rqg7... Korean, Restaurants			Philadelphia [11:30-20:30, 11:...	1 39.943223 -75.162568 BAP 19147	205 4.5 PA			

Yelp Data Analysis

Data Sample: yelp_academic_dataset_tip

business_id	compliment_count	date	text	user_id
3uLgwr0qeCNMjKenH...	0	2012-05-18 02:17:21	Avengers time wit...	AGNUgVwnZUey3gcPC...
QoezRbYQncpRqyrLH...	0	2013-02-05 18:35:10	They have lots of...	NBN4MgHP9D3cw--Sn...
MYoRNLb5chwjQe3c...	0	2013-08-18 00:56:08	It's open even wh...	-cop0vldyKh1qr-vz...
hv-bABTK-glh5wj31...	0	2017-06-27 23:05:38	Very decent fried...	FjMQVZjSqvY8syIO-5...
_uN0OudeJ3Zl_tf6n...	0	2012-10-06 19:43:09	Appetizers.. plat...	ld0AperBXk1h6Ubqm...
7Rm9Ba50bw23KTA8R...	0	2012-03-13 04:00:52	Chili Cup + Singl...	trf3Qcz8qvCDKXiTg...
kh-0iXqkL7b8UXNpg...	0	2013-12-03 23:42:15	Saturday, Dec 7th...	SMGAI1RjyfuYu-c-22...
jtri188kuhe_AuEOJ...	0	2016-11-22 22:14:58	This is probably ...	YVBB9g23nuVJ0u44z...
x0DBZmX4Em1VvbqtK...	0	2012-07-27 01:48:24	Tacos	VL12EhEdT40WqGq0n...
picJRCyqW1cF96Q3X...	0	2012-06-09 22:57:04	Starbucks substit...	4ay-fdVks5WMerYL...
clwjLY7PdYJpe7IP9...	0	2014-06-17 01:20:14	Order the Tortill...	0ttfcRwgRrYsTg9EV...
wlHodvVFLTgK3n12X...	0	2017-03-23 22:01:41	Very good will de...	JsXhBw6MntzTJjh_U...
wUMuvdUeVZODZk7Tj...	0	2013-02-28 02:05:54	If the Hotlight i...	Y0JfJh4B-jrtGc_AH...
MDr7KLYSPkEonvGoj...	0	2011-07-20 21:52:57	Let's go Yankees!	MlnuJ7T14CE0JDK2Z...
ak6R2akvIK9ijw3Fv...	0	2014-06-12 17:34:20	Basically same fo...	ffWWVlmsrN5lZ6sjA...
EXYbKA1tocvOK_1tX...	0	2011-10-13 03:15:15	Don't go for dinn...	j2sEA3hiUcwHfq9M1...
H9fkf4Xkj_j7Zxs1F...	0	2012-03-11 23:16:12	30 mins for take ...	jsaN4TDygu76AGTiB...
ReX09lhufLTax19kr...	0	2013-06-10 20:18:41	Got the grilled c...	kjFgyrCvmVVGSlgWz...
c5nLy7YgXG-IIr0mq...	0	2016-04-23 02:44:03	This is the bomb ...	I6aRZ4sE7ixv0_2r3...
LJaR65ALpz261_d1V...	0	2012-06-02 14:39:28	Helping Mona find...	L1514WTKPH7zWQWA6...

Yelp Data Analysis

Data Sample: yelp_academic_dataset_review

business_id cool	date funny	review_id stars	text useful	user_id
XQfwVwDr-v0ZS3_Cb... 0 2018-07-07 22:09:11 0 KU_05udG6zpxOg-Vc... 3.0 If you decide to ... 0 mh_-eMZ6K5RLWhZyI...				
7ATYjTIgM3jUlt4UM... 1 2012-01-03 15:28:18 0 BiTunyQ73aT9WBnpR... 5.0 I've taken a lot ... 1 OyoGAe7OKpv6SyGZT...				
YjUWPpI6HXG5301wP... 0 2014-02-05 20:30:30 0 saUsX_uimxR1CVr67... 3.0 Family diner. Had... 0 8g_iMtfSiwikVnbP2...				
kx2S0es4o-D3ZQBk... 1 2015-01-04 00:01:03 0 AqPFMlE6RsU23_au... 5.0 Wow! Yummy, diff... 1 _7bHUi9Uuf5__HHc...				
e4Vwtrqf-wpJfwesg... 1 2017-01-14 20:54:15 0 Sx8TMOWLNUJBWer-0... 4.0 Cute interior and... 1 bcjbaE6dDog4jkNY9...				
04UD14gamNjLY0IDY... 1 2015-09-23 23:10:31 2 JrIx1S1TzJ-iCu79u... 1.0 I am a long term ... 1 eUta8W_HdHMXPzLBB...				
gmjsEdUsKpj9Xu6p... 0 2015-01-03 23:21:18 2 6AxgBCNX_PNT0xmbR... 5.0 Loved this tour! ... 0 r3zeYsv1XFBR4dJp...				
LHSTtnW3YHCeUkRDG... 0 2015-08-07 02:29:16 0 _ZeMknuYdlQcUqnq... 5.0 Amazingly amazing... 2 yFzsLmaWF2d4Sr0U...				
B5XSoSG3SfvQGtKEG... 0 2016-03-30 22:46:33 1 ZKvDG2sBvHVdF5oBN... 3.0 This easter inste... 1 wSTuiTk-sKNdcFypr...				
gebiRewfieSdtt17P... 0 2016-07-25 07:31:06 0 pUyc0fUwM8vqX7KjR... 3.0 Had a party of 6 ... 0 59MxRhNVhU9MYndMk...				
uMvVYRgGNXF5boolA... 0 2015-06-21 14:48:06 0 rGQRf8UafX70T1MNN... 5.0 My experience wit... 2 1WHRWwQmZOZDAhp2Q...				
EQ-TZ2eeD_E0BHuv0... 0 2015-08-19 14:31:45 0 13Wk_mvAog6XANIuG... 4.0 Locals recommende... 0 ZbqSHbgCjzVAqaa7N...				
lj-E32x9_FA7GmUrB... 0 2014-06-27 22:44:01 0 XW_LfMv0fV2119c6x... 4.0 Love going here f... 0 90AtfnWag-ajVxRbU...				
RZtGWDLCAtuipwaZ... 0 2009-10-14 19:57:14 0 8JFGBuHMoiNDyfcxu... 4.0 Good food--loved ... 0 sm0v0ajNG01S4Pq7d...				
otQS34_MymjPTdNB... 0 2011-10-27 17:12:05 2 UBp0zWyH60Hmw6Fsa... 4.0 The bun makes the... 0 4Uh27DgGzsp6PqrH9...				
BVndHaLiHxEYbr76Z0... 0 2014-10-11 16:22:06 0 0AhBYw8IQ6wlfw1ow... 5.0 Great place for b... 0 1C2lxzUo1Hyye4RFI...				
YtSqYv1Q_p0ltsVPS... 0 2013-06-24 11:21:25 0 oyaMhzBSwfGgemSGu... 5.0 Tremendous servic... 0 Dd1jQj7S-BFGqRbAp...				
rBdG_23USc7DletfZ... 0 2014-08-10 19:41:43 0 LnGZB0ffjgeVDVz5I... 4.0 The hubby and I h... 1 j2wlzrntrbKwy0c0i...				
CLEWowfkj-wKYJlQD... 1 2016-03-07 00:02:18 0 u2vza0Qj2feRshaa... 5.0 I go to blow bar ... 2 NDZvyYHTUWWu-kqgQ...				
eFvzHawVJofxSnD7T... 0 2014-11-12 15:30:27 0 Xs8Z8lmKkosqW5mw... 5.0 My absolute favor... 0 IQsF3Rc6IgCzjVV9D...				

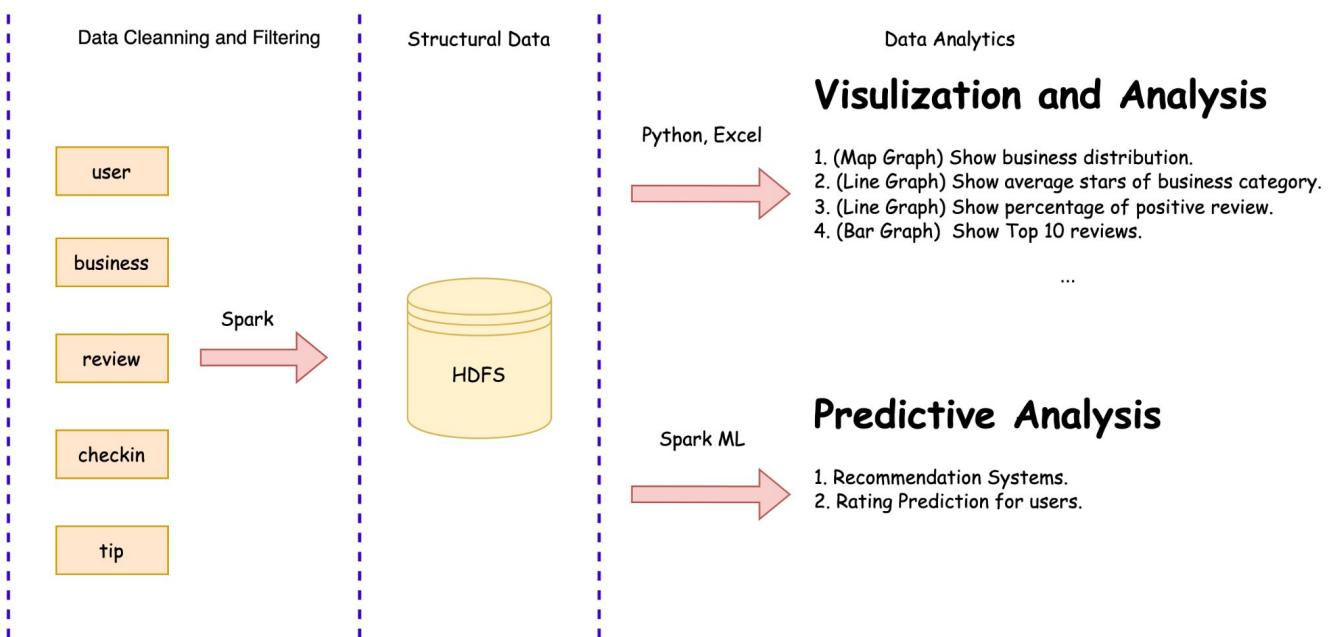
Yelp Data Analysis

Data Sample: yelp_academic_dataset_checkin

business_id	date
--kPU91CF4Lq2-W1...	2020-03-13 21:10:...
--0iUa4sNDFiZFrAd...	2010-09-13 21:43:...
--30_8IhuyMHbSOcN...	2013-06-14 23:29:...
--7PUidqRWpRSpXeb...	2011-02-15 17:12:...
--7jw19RH9JKXgFoh...	2014-04-21 20:42:...
--8Ib0sAAxjKRoYsB...	2015-06-06 01:03:...
--9osgUCSDUWUkoTL...	2015-06-13 02:00:...
--ARBQr1WMsTWiwOK...	2014-12-12 00:44:...
--FWWsIwxRwuw9vIM...	2010-09-11 16:28:...
--FcbSxK1AoEtEAx0...	2017-08-18 19:43:...
--LC8cIrALInl2vyo...	2017-01-12 19:10:...
--MbOh201pATkXa7x...	2013-04-21 01:52:...
--N9yp3ZWqQIm7DqK...	2012-10-06 20:46:...
--03ip9NpXTKD4oBS...	2010-04-17 21:07:...
--OS_I7dnABrXvRCC...	2018-05-11 18:23:36
--S43ruInmIsGrnnk...	2010-08-29 01:17:...
--SJXpAa0E-GCp2sm...	2014-04-06 22:23:...
--Sd930FWITqDHifM...	2013-01-09 17:42:...
--ZVrH2X2QXBFDcIl...	2010-08-12 18:21:...
--ZWv8kG1M2YL58uK...	2010-10-13 18:41:...

Yelp Data Analysis

Design Diagram



Yelp Data Analysis

Core Challenge 1

Data visualization:

Peel cluster doesn't support tableau for data visualization. We export analytics results in csv format then use python and excel to visualize our results

Yelp Data Analysis

Core Challenge 2

Recommended Systems:

For the recommended system, the mainly core challenge is training the model using large volume of dataset (4G+). This scale of data requires more training iterations for our model. Plus, our model will have lower accuracy since and bias because of the lack of complete data. Thus, it's important to do the future engineering of the training set.

Yelp Data Analysis

Core Challenge 3

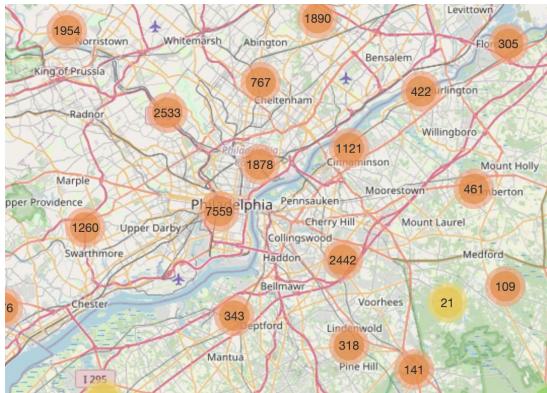
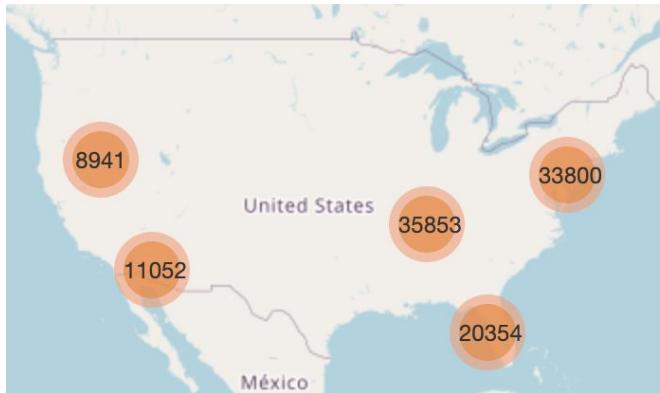
In friends relationship analysis:

Each user has a list of his friends' id. However, some of those ids are not included in our primary key. In other words, we don't have data for that user. But while rendering the Force Atlas Graph, these data points are necessary. To solve this problem, we created a dummy point for that user and assigned it with 0 weight.

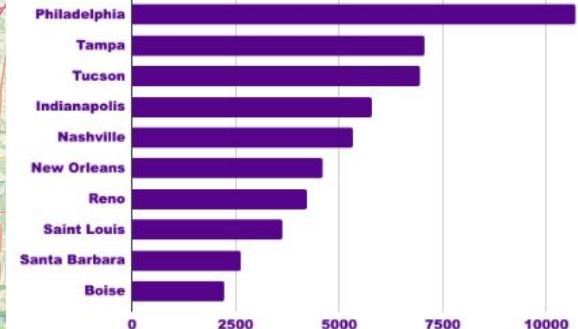
Yelp Data Analysis

Results - Overview:

We show our data of each business and number of business in different area on the map



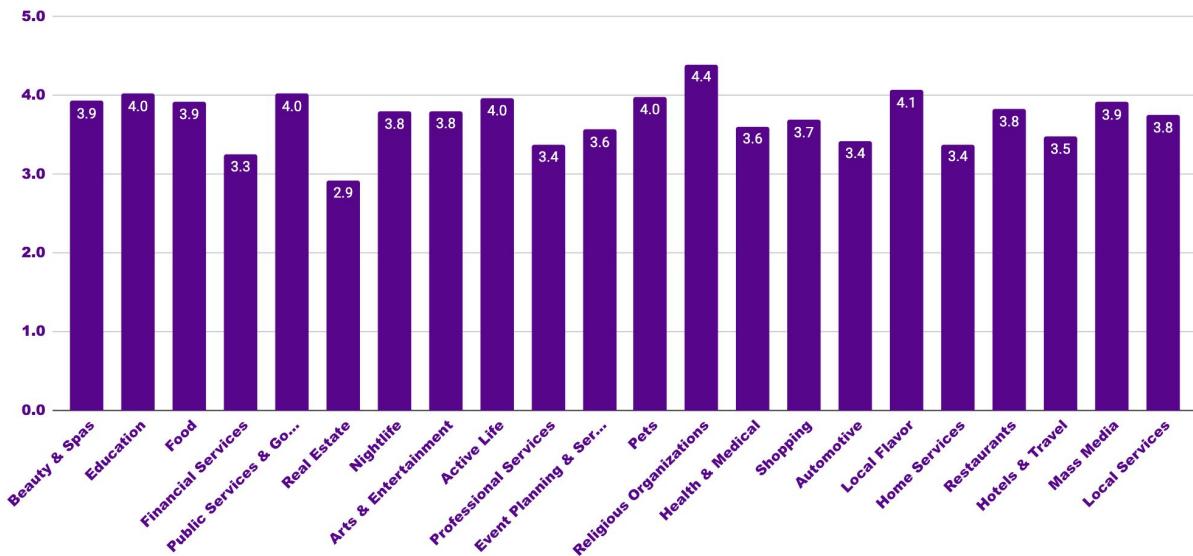
Top 10 Cities With the Most Businesses



Yelp Data Analysis

Result: Average Stars of Business Category

Avg Stars in terms of Business Category



These figures show the Average Stars of Business Category.

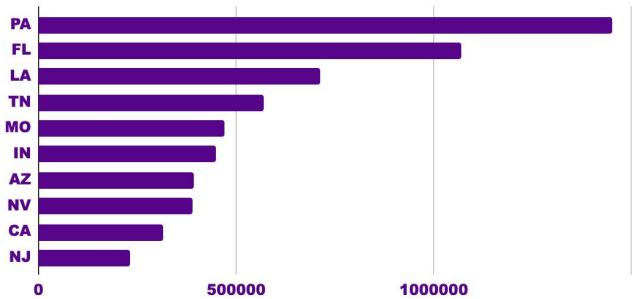
People scores higher on religious organizations.

The worst rating category is Real Estate.

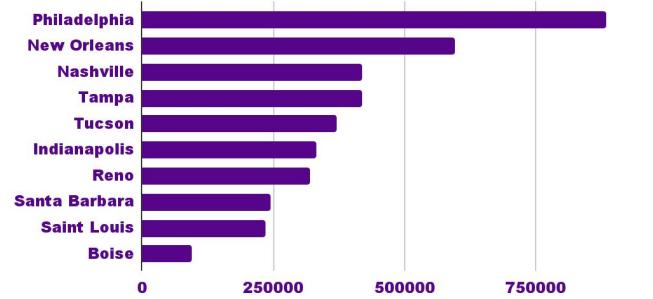
Yelp Data Analysis

Result: Top 10 Review(State, City, Category)

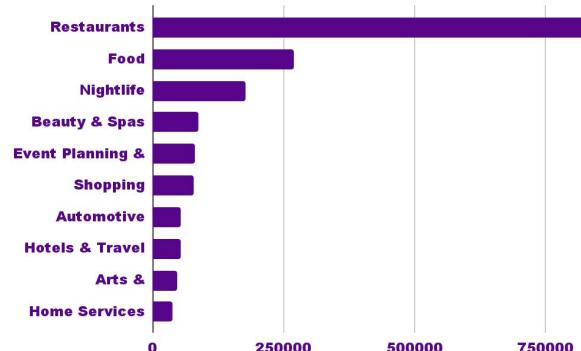
Top 10 Reviews States



Top 10 Reviews cities



Top 10 Categories in terms of Reviews Number



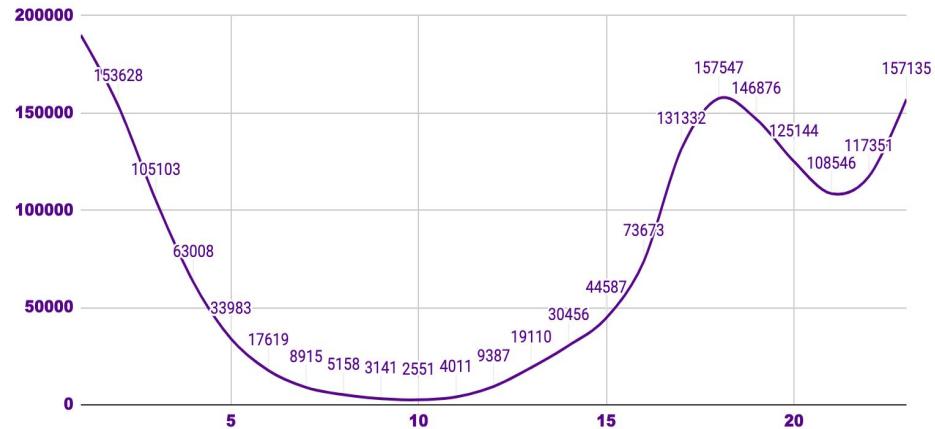
As we can see, restaurants is the most popular business category amount the dataset.

So we do some research on this category and get some interesting results

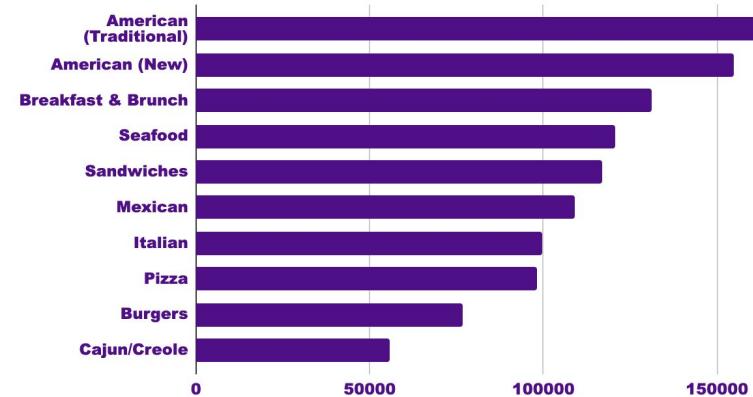
Yelp Data Analysis

Results: Restaurant Overview Analyze

Restaurant Checkin Time Distribution Throughout the Day



Top 10 Cuisines



Yelp Data Analysis

Results - Review : length of positive and negative review



This graph compares the length of positive and negative reviews.

It shows that people were more likely to use more words to express dissatisfaction when they scored low, while they were more stingy with their evaluations when they scored high.

Yelp Data Analysis

Results - Review: the most used words

helpful **available** **clean** **quiet** **fast**
generous **incredibly** **friendly** **prefer**
complimentary **pleasantly**
hot **easy** **attentive** **pleased**
worth **interesting** **rich** **seasoned**
variety **fair** **authentic** **tender**
reasonable **phenomenal**
fresh **warm** **excited** **sweet** **solid**
nicely **soft** **polite**
delicious **prompt** **free**
reasonably **comfortable** **cool** **recommendations**
lovely **recommend**
enjoy **satisfied** **appreciate** **affordable** **beautiful** **consistently**

Indian Food Positive Keywords

frozenstale
issue
strange
weird
smelled
dark
lack
expensive
hung
confused
greasy
sick
break
slow
mistake
dirty
bother
problems
lukewarm
miss
lost
odd
complained
spoiled
doubt
ridiculous
mess
limited
messed
refused
sour
smell
gross
pathetic
fried
downhill
issues
ridiculous
cheap
disgusting
refused
hard
poor
overpriced
poorly
lacked

Indian Food Negative Keywords

friendly *perfection excited beautiful prefer fast consistently*
reasonable *cool pleasant cool prompt variety*
helpful *easy*
crisp *solid*
polite nicely *satisfied*
sweet *loves healthy*
delicious *incredibly solid*
clean *worth affordable reasonably*
generous *modern*
free *fancy*
fresh *attentive tender fairly soft*
hot *available seasoned*
recommend *efficient consistent pleased warm*
authentic

Chinese Food Positive Keywords

A word cloud diagram showing words associated with food texture and taste. The words are arranged in a grid-like shape, with colors ranging from blue to green. The most prominent word is "fried" in yellow, located at the top center. Other large words include "sour" (blue), "dirty" (blue), "nasty" (blue), "issue" (green), "rude" (blue), "mess" (blue), "miss" (blue), "cheap" (yellow), "poor" (yellow), "downhill" (yellow), "stale" (yellow), "slow" (yellow), "fat" (yellow), "lost" (yellow), "lacked" (blue), "smelled" (blue), "issues" (blue), "weird" (blue), "mistake" (blue), "dark" (blue), "hung" (blue), "messed" (blue), "annoyed" (blue), "dim" (purple), "ridiculous" (purple), "overpriced" (purple), "trash" (purple), "strange" (purple), "doubt" (purple), "lack" (purple), "pan" (purple), "old" (purple), "cold" (green), "greasy" (green), "disgusting" (green), and "sticky" (green).

Chinese Food Negative Keywords

This is the word most used by users in positive and negative reviews of Indian Restaurant and Chinese Restaurant.

We can use this method to analyze people's evaluation of the restaurant and make improvements.



Yelp Data Analysis

Results: Recommended Systems

In this project, our team applied **Alternating Least Square (ALS)** model in Spark to train a recommended system.

After finishing the training, our team recommended 5 items for 10 randomly selected people.

Result is listed as follows (just show the first 4 people and their top 2 recommended services):

name	r1_name	r1_categories	r2_name	r2_categories
Benny	It' Sugar	Food, Candy Stores, Specialty Food	El Pastorito	Restaurants, Mexican
Varusha	Not Just Crab	Restaurants, Seafood, Nightlife	Zithers Candy	Candy Stores
Crosby	June's Restaurant	Ice Cream & Frozen Yogurt	Pizza Hut	Italian, Chicken Wings
Daniel	City Grocers	Food, Grocery	Amore Pasta	Italian, Restaurants

Yelp Data Analysis

Results: Recommended Systems - Simple Test

Our team verified our recommended systems by comparing people's past consuming records and the recommended items.

name	p1_name	p1_categories	p2_name	p2_categories
Benny	The Seed	Food, Juice Bars, Specialty Food	The Cheesecake Factory - Reno	Restaurants, Mexican
Varusha	Roller Kingdom	Active Life, Skating Rinks		
Crosby	Polish Cottage	Polish, Restaurants	Pasco Kitchen & Lounge	American (Traditional)
Daniel	Pokey's	Pubs, Nightlife, Bars		

Yelp Data Analysis

Results: Friends relationship

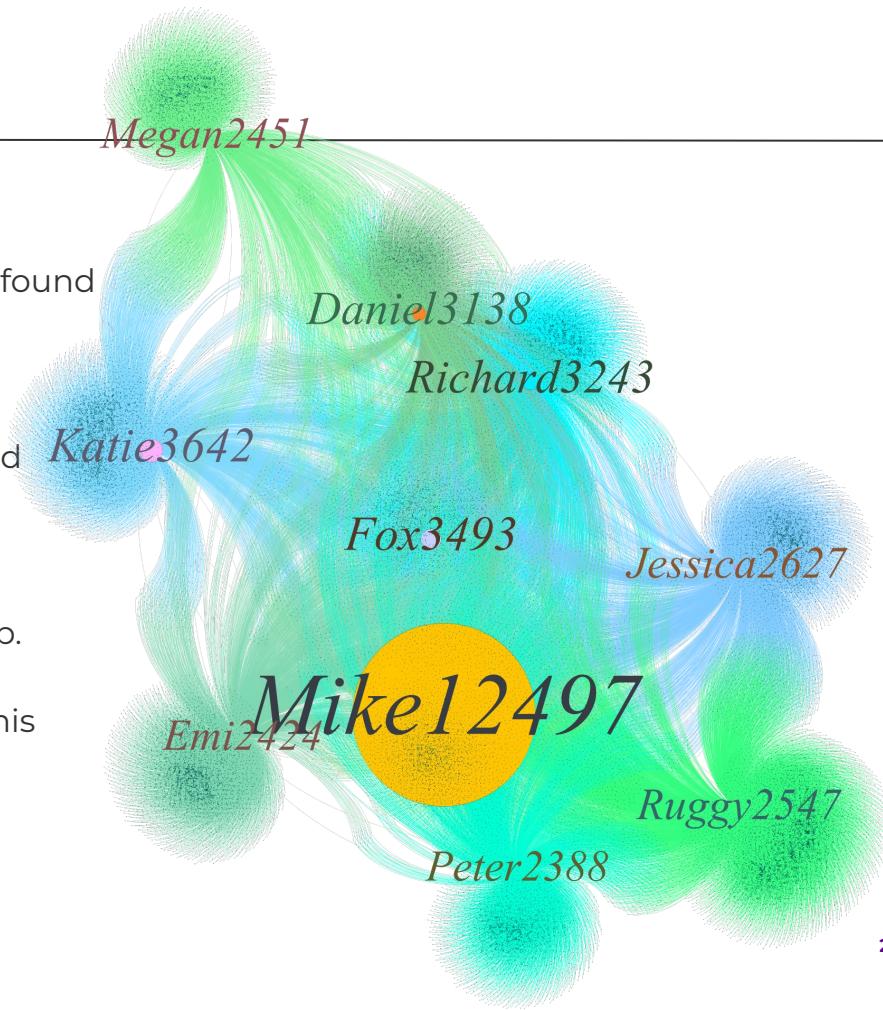
We analyzed yelp_academic_dataset_user.json and found the most popular users, aka influencers, who have the most followers.

The right hand side is a Force Atlas 2 Graph rendered with 57624 users data.

Each vertex represents for a user and an edge represents for a friend relationship.

For example, Mike has 12497 followers and most of his followers have also followed other influencers.

This was a 4K image before compressing.



Yelp Data Analysis

Obstacles

1. The datasets are really large compared to other teams, especially the review and user datasets. Thus, when running our program on the clusters, it may take long time to finishing each step, and it's not convenient to use the Zeppelin, because it's too slow to show the results.
2. The training set for the recommended system is too large to train for more iterations in Spark. It's easy to cause stack overflow when the number of iterations is greater than 25, this situation greatly affects the accuracy of the model. I think our model should have better performance if it had enough training iterations.

Yelp Data Analysis

Summary

Yelp's open datasets could actually explore many interesting patterns.

We analyze what makes a good restaurant and what concerns customers, and then make recommendations of the future improvement and profit growth. Building visualization on this dataset could provide people with potential stores to solve their requirements, sometimes helps them reduce much time to do the decision and exploration. Analyzing friends relationship helps for improving users experience of commercial advertising.

Hence, this data analytics could actually help people a lot.

Yelp Data Analysis

Acknowledgement

1. Thanks NYU Peel Cluster for providing us with machines to run our applications.
2. Thanks Professor Tang Yang for organizing and instructing this course.
3. Thanks Yelp for opening the Yelp datasets for us to perform our analytics.
4. Thanks everyone in this team for providing support for each other.
5. Thanks Tableau for providing us free student license.

Thank you!