

# Analysis of YELP dataset

1<sup>st</sup> Jagath Sai Narayana Kakaraparty  
Dept. of ES Data Science  
University at Buffalo  
Buffalo, NY, USA  
jagathsa@buffalo.edu

2<sup>nd</sup> Nisarg Negi  
Dept. of ES Data Science  
University at Buffalo  
Buffalo, NY, USA  
nisargne@buffalo.edu

3<sup>rd</sup> Siddhi Satish Jadhav  
Dept. of ES Data Science  
University at Buffalo  
Buffalo, NY, USA  
sjadhav3@buffalo.edu

**Abstract**—Consumer reviews have become an important deciding factor for businesses to flourish in today’s world. Consumers frequently rely on certified sources for trustworthy reviews of various establishments such as restaurants, hotels, and so on. Yelp is a local business directory and review site that also has social networking capabilities. It enables people to rate and review businesses. The goal is to create a relational database system that will integrate existing Yelp data to give insightful analytics and assist existing and future business owners in making critical decisions about a new business or business expansion.

## I. INTRODUCTION

Every day, the internet generates an exponential amount of data. It is difficult to store massive volumes of data and extract insight from it. Yelp is a website that publishes reviews submitted by the public about local businesses. It allows business owners to enhance their services and users to select the best business accessible. The dataset consists of a subset of Yelp’s businesses, reviews, and user information. It was originally created for the Yelp Dataset Challenge, allowing students to conduct research or analyze Yelp’s data and report their findings. The collection contains information on businesses in 11 metropolitan areas across four nations. We intend to create an effective database system based on this dataset that can manage vast amounts of data and deliver insights such as reviews supplied by other consumers, ratings, tips, and so on which will be valuable to both consumers and businesses.

### A. Targeted Audience

Owners that are prepared to monitor their firm and make significant decisions based on the visual plots, such as modifying the infrastructure or concentrating on an area that will help them expand, will be the system’s users.

The dataset’s real-world goal is to offer relevant insights and assist current and prospective business owners in making crucial choices about the

## II. DATA DESCRIPTION

The source data set is a JSON file, which is transformed into a CSV file using Python programs into tabular form. We chose to store the data in a database system for the reasons listed below: (I) The dataset is 4GB in size, making it challenging to store and access the data using CSV files. Postgres SQL, on the other hand, can store a large amount of data. Therefore, a database system would be a superior

strategy for our instance. (II) The project’s main goal is to establish connections between the tables in order to learn more about the businesses. Through the use of primary keys and foreign keys, Postgre SQL aids in achieving this whereas a relational schema cannot be established using Excel. (III) The yelp data has multiple groups that are there are multiple categories for each business. Using a database system, we can enlist the business based on the groups and perform further analysis. Such functionality is not present in excel. (IV) The main reason behind choosing a database system over CSV for the current data is faster access as the size of the data is quite huge. Indexing can be performed on the database tables for faster lookup whereas CSV depends on the size of the ram to access the data.

## III. DATA IMPORT AND PROCESSING

The dataset is an academic challenge distributed by Yelp. The dataset is in a compressed form. The uncompressed data is in form of JSON and hence requires extraction. Using python, the data is extracted into CSV and imported into Postgres via scripting.

The data that we imported was in JSON format, so we extracted and scraped the data into CSV format using python before loading it into Postgres SQL.

## IV. DATABASE AND TABLE FORMATION

After extracting the data into CSV files, we had a total of 7 CSV files. We create a database called ‘Yelp’ in Postgres SQL and the tables created are as mentioned below.

‘yelp\_business’ table is one of the base tables created using the yelp\_business.csv by first extracting the CSV into a data frame object. We inserted the data into this table using the “to\_sql” library which is used to insert the records present in the data frame object to a SQL database. We used this library for quicker insertion of data into the database. We then altered the table to add a primary key constraint on the column ‘business\_id’.

The ‘yelp\_user’ table is one of the base tables created using the yelp\_user.csv by first extracting the CSV into a data frame object. We inserted the data into this table using the “to\_sql” library which is used to insert the records present in the data frame object to a SQL database. We used this library for quicker insertion of data into the database. We then altered the table to add a primary key constraint on the column ‘user\_id’.

```
CREATE TABLE YELP_BUSINESS
(BUSINESS_ID VARCHAR(50) NOT NULL,
NAME CHAR(30) NOT NULL,
NEIGHBORHOOD VARCHAR(30) NOT NULL,
ADDRESS VARCHAR(100) NOT NULL,
CITY VARCHAR(30) NOT NULL,
STATE CHAR(5) NOT NULL,
POSTAL_CODE INTEGER NOT NULL,
LATITUDE DOUBLE PRECISION NOT NULL,
LONGITUDE DOUBLE PRECISION NOT NULL,
STARS INTEGER NOT NULL,
REVIEW_COUNT INTEGER NOT NULL,
IS_OPEN INTEGER NOT NULL,
CATEGORIES VARCHAR(150) NOT NULL);
```

#### Query 1: Create YELP\_BUSINESS Table

```
ALTER TABLE YELP_BUSINESS
ADD PRIMARY KEY (BUSINESS_ID);
```

#### Query 2: Alter YELP\_BUSINESS Table

```
CREATE TABLE YELP_USER
(USER_ID VARCHAR(50),
NAME CHAR(50),
REVIEW_COUNT INTEGER DEFAULT 0,
YELPING_SINCE CHAR(10) DEFAULT 'NONE');
```

#### Query 3: Create YELP\_USER Table

```
ALTER TABLE YELP_USER
ADD PRIMARY KEY (USER_ID);
```

#### Query 4: Alter YELP\_USER Table

yelp\_business table consists of a 'categories' column in the form of lists; hence, we need to remove the partial dependency. So, we created a new table 'yelp\_categories' which consisted of the categories from the 'yelp\_business' table, making it follow 1 NF. Insertion into the table is achieved using python code. We then altered the table to add a foreign key constraint on the column 'business\_id' which is referenced from the 'yelp\_business' table and we also add a delete cascade constraint on the table.

```
CREATE TABLE YELP_CATEGORIES
(BUSINESS_ID VARCHAR(50) NOT NULL,
CATEGORIES CHAR(50) NOT NULL);
```

#### Query 5: Create YELP\_CATEGORIES Table

The 'yelp\_checkin' table created using the yelp\_checkin.csv by first extracting the CSV into a data frame object. We inserted the data into this table using the "to\_sql" library which is used to insert the records present in the data frame object to a SQL database. We used this

```
ALTER TABLE YELP_CATEGORIES
ADD FOREIGN KEY (BUSINESS_ID)
REFERENCES YELP_BUSINESS(BUSINESS_ID)
ON DELETE CASCADE;
```

#### Query 6: Alter YELP\_CATEGORIES Table

library for quicker insertion of data into the database. We then altered the table to add a foreign key constraint on the column 'business\_id' which is referenced from the 'yelp\_business' table and we also add a delete cascade constraint on the table.

```
CREATE TABLE YELP_CHECKIN
(BUSINESS_ID VARCHAR(50) NOT NULL,
WEEKDAY VARCHAR(20) NOT NULL,
HOUR TIME NOT NULL,
CHECKINS INTEGER NOT NULL);
```

#### Query 7: Create YELP\_CHECKIN Table

```
ALTER TABLE YELP_CHECKIN
ADD FOREIGN KEY (BUSINESS_ID)
REFERENCES YELP_BUSINESS(BUSINESS_ID)
ON DELETE CASCADE;
```

#### Query 8: Alter YELP\_CHECKIN Table

The 'yelp\_business\_hours' table is created using the yelp\_business\_hours.csv by first extracting the CSV into a data frame object. We inserted the data into this table using the "to\_sql" library which is used to insert the records present in the data frame object to a SQL database. We used this library for quicker insertion of data into the database. We then altered the table to add a foreign key constraint on the column 'business\_id' which is referenced from the 'yelp\_business' table and we also add a delete cascade constraint on the table.

```
CREATE TABLE YELP_BUSINESS_HOURS
(BUSINESS_ID VARCHAR(50),
MONDAY VARCHAR(20) DEFAULT 'NONE',
TUESDAY VARCHAR(20) DEFAULT 'NONE',
WEDNESDAY VARCHAR(20) DEFAULT 'NONE',
THURSDAY VARCHAR(20) DEFAULT 'NONE',
FRIDAY VARCHAR(20) DEFAULT 'NONE',
SATURDAY VARCHAR(20) DEFAULT 'NONE',
SUNDAY VARCHAR(20) DEFAULT 'NONE');
```

#### Query 9: Create YELP\_BUSINESS\_HOURS Table

The 'yelp\_tip' table is created using the yelp\_tips.csv by first extracting the CSV into a data frame object. We inserted the data into this table using the "to\_sql" library which is used to insert the records present in the data frame object to a SQL database. We used this library for quicker insertion of data into the database. We then altered the table to add a

```
ALTER TABLE YELP_BUSINESS_HOURS
ADD FOREIGN KEY (BUSINESS_ID)
REFERENCES YELP_BUSINESS(BUSINESS_ID)
ON DELETE CASCADE;
```

#### Query 10: Alter YELP\_BUSINESS\_HOURS Table

foreign key constraint on the column 'business\_id' which is referenced from the 'yelp\_business' table and we also add a delete cascade constraint on the table.

```
CREATE TABLE YELP_TIP
(TEXT VARCHAR(200) DEFAULT 'NONE',
LIKES INTEGER DEFAULT 0,
BUSINESS_ID VARCHAR(50) NOT NULL,
USER_ID VARCHAR(50) NOT NULL,
FOREIGN KEY (BUSINESS_ID)
REFERENCES YELP_BUSINESS(BUSINESS_ID),
FOREIGN KEY (USER_ID) REFERENCES YELP_USER(USER_ID));
```

#### Query 11: Create YELP\_TIP Table

```
ALTER TABLE YELP_TIP
ADD FOREIGN KEY (BUSINESS_ID)
REFERENCES YELP_BUSINESS(BUSINESS_ID)
ON DELETE CASCADE;
```

#### Query 12: Alter YELP\_TIP Table

The 'yelp\_review' table is created using the yelp\_review.csv by first extracting the CSV into a data frame object. We inserted the data into this table using the "to\_sql" library which is used to insert the records present in the data frame object to a SQL database. We used this library for quicker insertion of data into the database. We then altered the table to add a primary key constraint on the column 'review\_id' and a foreign key constraint on the column 'business\_id' which is referenced from the 'yelp\_business' table and on the column 'user\_id' which is referenced from the 'yelp\_user' table. We also add a delete cascade constraint on the table.

```
CREATE TABLE YELP_REVIEW
(REVIEW_ID VARCHAR(50) NOT NULL,
USER_ID VARCHAR(50) NOT NULL,
BUSINESS_ID VARCHAR(50) NOT NULL,
STARS INTEGER DEFAULT 0,
TEXT VARCHAR(300) DEFAULT 'NONE');
```

#### Query 13: Create YELP\_REVIEW Table

### V. PROBLEMS FACED AND THEIR SOLUTION

When we tried to create tables by defining the constraints and types for each of the attributes in the tables, when we inserted data into the tables, the data was inserted into the

```
ALTER TABLE YELP_REVIEW
ADD PRIMARY KEY (REVIEW_ID);
ALTER TABLE YELP_REVIEW
ADD FOREIGN KEY (BUSINESS_ID)
REFERENCES YELP_BUSINESS(BUSINESS_ID)
ON DELETE CASCADE;
ALTER TABLE YELP_REVIEW
ADD FOREIGN KEY (USER_ID)
REFERENCES YELP_USER(USER_ID)
ON DELETE CASCADE;
```

#### Query 14: Alter YELP\_REVIEW Table

table the with the original data types and with no constraints. Hence, we altered each of the tables to define the data types and constraints separately.

```
ALTER TABLE YELP_BUSINESS
ALTER COLUMN BUSINESS_ID TYPE VARCHAR(500),
ALTER COLUMN NAME TYPE CHAR(200),
ALTER COLUMN ADDRESS TYPE VARCHAR(200),
ALTER COLUMN CITY TYPE VARCHAR(50),
ALTER COLUMN STATE TYPE CHAR(5),
ALTER COLUMN POSTAL_CODE TYPE CHAR(10),
ALTER COLUMN LATITUDE TYPE DOUBLE PRECISION,
ALTER COLUMN LONGITUDE TYPE DOUBLE PRECISION,
ALTER COLUMN STARS TYPE DECIMAL(1,0),
ALTER COLUMN REVIEW_COUNT TYPE NUMERIC(100,0),
ALTER COLUMN IS_OPEN TYPE NUMERIC(1,0);
```

#### Query 15: Alter YELP\_BUSINESS Table to change data types

```
ALTER TABLE YELP_BUSINESS
ALTER COLUMN BUSINESS_ID SET NOT NULL,
ALTER COLUMN ADDRESS SET NOT NULL,
ALTER COLUMN CITY SET DEFAULT 'NONE',
ALTER COLUMN STATE SET DEFAULT 'NONE',
ALTER COLUMN POSTAL_CODE SET DEFAULT '0',
ALTER COLUMN LATITUDE SET DEFAULT '0.0000',
ALTER COLUMN LONGITUDE SET DEFAULT '0.0000',
ALTER COLUMN STARS SET NOT NULL,
ALTER COLUMN REVIEW_COUNT SET NOT NULL,
ALTER COLUMN IS_OPEN SET NOT NULL;
```

#### Query 16: Alter YELP\_BUSINESS Table to add constraints

```
ALTER TABLE YELP_USER
ALTER COLUMN USER_ID TYPE VARCHAR(500),
ALTER COLUMN NAME TYPE CHAR(200),
ALTER COLUMN REVIEW_COUNT TYPE NUMERIC(100,0),
ALTER COLUMN YELPING_SINCE TYPE DATE
USING TO_DATE(YELPING_SINCE, 'YYYY-MM-DD');
```

#### Query 17: Alter YELP\_USER Table to change data types

### VI. FUNCTIONAL DEPENDENCIES AND BCNF CHECK

In order to maintain data integrity by removing redundancies in a database, it is necessary to normalize the database relations. The database is already normalized since we handled the partial dependency that one of the table was holding, tables

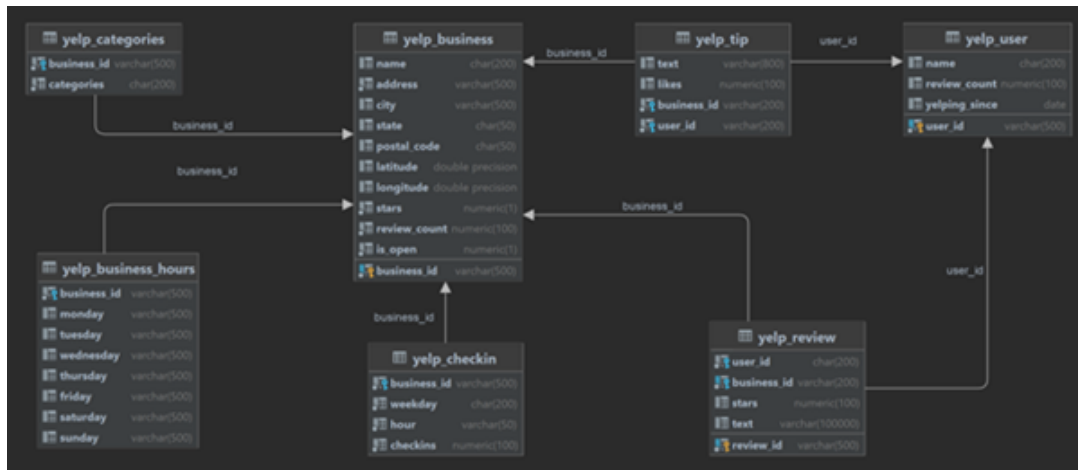


Fig. 1: Entity Relationship Diagram

```

ALTER TABLE YELP_USER
ALTER COLUMN USER_ID SET NOT NULL,
ALTER COLUMN NAME SET DEFAULT 'UNKNOWN',
ALTER COLUMN REVIEW_COUNT SET DEFAULT 0,
ALTER COLUMN YELPING_SINCE SET DEFAULT CURRENT_DATE;
  
```

Query 18: Alter YELP\_USER Table to add constraints

```

ALTER TABLE YELP_CATEGORIES
ALTER COLUMN BUSINESS_ID TYPE VARCHAR(500),
ALTER COLUMN CATEGORIES TYPE CHAR(200);
  
```

Query 19: Alter YELP\_CATEGORIES Table to change data types

```

ALTER TABLE YELP_CATEGORIES
ALTER COLUMN BUSINESS_ID SET NOT NULL,
ALTER COLUMN CATEGORIES SET NOT NULL;
  
```

Query 20: Alter YELP\_CATEGORIES Table to add constraints

```

ALTER TABLE YELP_CHECKIN
ALTER COLUMN BUSINESS_ID SET NOT NULL,
ALTER COLUMN WEEKDAY SET NOT NULL,
ALTER COLUMN HOUR SET NOT NULL,
ALTER COLUMN CHECKINS SET NOT NULL;
  
```

Query 21: Alter YELP\_CHECKIN Table to change data types

```

ALTER TABLE YELP_CHECKIN
ALTER COLUMN BUSINESS_ID TYPE VARCHAR(500),
ALTER COLUMN WEEKDAY TYPE CHAR(200),
ALTER COLUMN HOUR TYPE VARCHAR(50),
ALTER COLUMN CHECKINS TYPE NUMERIC(100,0);
  
```

Query 22: Alter YELP\_CHECKIN Table to add constraints

```

ALTER TABLE YELP_BUSINESS_HOURS
ALTER COLUMN BUSINESS_ID TYPE VARCHAR(500),
ALTER COLUMN MONDAY TYPE VARCHAR(500),
ALTER COLUMN TUESDAY TYPE VARCHAR(500),
ALTER COLUMN WEDNESDAY TYPE VARCHAR(500),
ALTER COLUMN THURSDAY TYPE VARCHAR(500),
ALTER COLUMN FRIDAY TYPE VARCHAR(500),
ALTER COLUMN SATURDAY TYPE VARCHAR(500),
ALTER COLUMN SUNDAY TYPE VARCHAR(500);
  
```

Query 23: Alter YELP\_BUSINESS\_HOURS Table to change data types

```

ALTER TABLE YELP_BUSINESS_HOURS
ALTER COLUMN BUSINESS_ID SET NOT NULL,
ALTER COLUMN MONDAY SET DEFAULT 'NONE',
ALTER COLUMN TUESDAY SET DEFAULT 'NONE',
ALTER COLUMN WEDNESDAY SET DEFAULT 'NONE',
ALTER COLUMN THURSDAY SET DEFAULT 'NONE',
ALTER COLUMN FRIDAY SET DEFAULT 'NONE',
ALTER COLUMN SATURDAY SET DEFAULT 'NONE',
ALTER COLUMN SUNDAY SET DEFAULT 'NONE';
  
```

Query 24: Alter YELP\_BUSINESS\_HOURS Table to add constraints

```

ALTER TABLE YELP_TIP
ALTER COLUMN TEXT TYPE VARCHAR(800),
ALTER COLUMN LIKES TYPE NUMERIC(100,0),
ALTER COLUMN BUSINESS_ID TYPE VARCHAR(200),
ALTER COLUMN USER_ID TYPE VARCHAR(200);
  
```

Query 25: Alter YELP\_TIP Table to change data types

```

ALTER TABLE YELP_TIP
ALTER COLUMN TEXT SET DEFAULT 'NONE',
ALTER COLUMN LIKES SET DEFAULT 0,
ALTER COLUMN BUSINESS_ID SET NOT NULL,
ALTER COLUMN USER_ID SET NOT NULL;
  
```

Query 26: Alter YELP\_TIP Table to add constraints

are already in 1NF. Since the primary keys alone are able

### Query 27: Alter YELP\_REVIEW Table to change data types

Query 28: Alter YELP\_REVIEW Table to add constraints

Since checkins can be defined by business\_id, weekday, and hour there exists only one FD. For the functional dependency shown above, the attributes on the left side i.e. business id,

Since all the attributes can be defined by review\_id alone, there exists only one FD with others being subsets of this FD.

- For the functional dependency shown above, the attribute on the left side i.e. review\_id is a key and there exists no transitive dependencies which means the table is in BCNF.

Fig. 2: Explain query before indexing

```

-- Query Plan
1  Gather (cost=8886.22..20867.26 rows=8268 width=108) (actual time=0.099..4293.921 rows=8099 loops=1)
2   Buckers Planned: 2
3   Buckers Launched: 2
4   → Parallel Hash Join (cost=8886.22..16236.76 rows=8002 width=108) (actual time=18.391..50.517 rows=1766 loops=1)
5     Hash Cond: (t.business_id = b.business_id)
6     → Parallel Seq Scan on yelp_categories t (cost=0.00..11287.31 rows=21767 width=74) (actual time=0.683..84.363 rows=6418 loops=1)
7     Filter: (categories = 'Restaurants')
8     Rows Removed by Filter: 41049
9     → Parallel Hash Join (cost=8886.22..16666.28 rows=11889 width=80) (actual time=18.116..18.116 rows=8929 loops=1)
10      Buckers: 37968 Buckets: 1 Memory Usage: 384000
11      → Parallel Seq Scan on yelp_business b (cost=0.00..16666.28 rows=11889 width=80) (actual time=0.437..43.961 rows=26775 loops=1)
12      Filter: (city = 'Las Vegas')
13      Rows Removed by Filter: 147992
14 Planning Time: 13.139 ms
15 Execution Time: 4293.921 ms

```

Fig. 3: Explain query after indexing

```

SELECT T.TEXT, B.NAME
FROM YELP_BUSINESS AS B
INNER JOIN YELP_CATEGORIES AS C
ON B.BUSINESS_ID = C.BUSINESS_ID
INNER JOIN YELP_TIP AS T
ON B.BUSINESS_ID = T.BUSINESS_ID
WHERE C.CATEGORIES = 'RESTAURANTS'
AND B.CITY = 'LAS VEGAS' LIMIT 10;

```

Query 30: Select query 1

```

-- Query Plan
1  Limit (cost=10409.00..10409.00 rows=10 width=80) (actual time=77.965..426.868 rows=10 loops=1)
2   → Nested Loop (cost=10409.00..10409.00 rows=10 width=80) (actual time=77.962..426.868 rows=10 loops=1)
3     Join Filter: (t.business_id = b.business_id)
4     → Gather (cost=10409.00..10706.39 rows=10000 width=110) (actual time=100.872..1009.193 rows=10 loops=1)
5       Buckers Planned: 1
6       Buckers Launched: 1
7       → Parallel Hash Join (cost=10409.00..10809.79 rows=10000 width=110) (actual time=100.871..10.468 rows=10 loops=1)
8         Hash Cond: (t.business_id = b.business_id)
9         Hash Join (cost=10409.00..10706.39 rows=10000 width=110) (actual time=100.862..1.371 rows=10 loops=1)
10          → Parallel Seq Scan on yelp_tip t (cost=0.00..11287.31 rows=21767 width=74) (actual time=0.467..16.868 rows=10 loops=1)
11          → Parallel Hash Join (cost=10409.00..11287.31 rows=21767 width=110) (actual time=100.401..11.894 rows=10 loops=1)
12            Filter: (categories = 'Restaurants')
13            Rows Removed by Filter: 41049
14            → Hash Join (cost=10409.00..11287.31 rows=21767 width=110) (actual time=100.401..11.894 rows=10 loops=1)
15              Hash Cond: (t.business_id = b.business_id)
16              Hash Join (cost=10409.00..11287.31 rows=21767 width=110) (actual time=100.401..11.894 rows=10 loops=1)
17              → Parallel Seq Scan on yelp_business b (cost=0.00..16666.28 rows=11889 width=80) (actual time=0.437..43.961 rows=26775 loops=1)
18              Filter: (city = 'Las Vegas')
19              Rows Removed by Filter: 1
20 Planning Time: 1.888 ms
21 Execution Time: 1009.193 ms

```

Fig. 4: Select query 1 before indexing

```

-- Query Plan
1  Limit (cost=0.00..10409.00 rows=10 width=80) (actual time=71.058..117.960 rows=10 loops=1)
2   → Nested Loop (cost=0.00..10409.00 rows=10 width=80) (actual time=71.057..117.960 rows=10 loops=1)
3     Join Filter: (t.business_id = b.business_id)
4     → Nested Loop (cost=0.00..12735.16 rows=10 width=80) (actual time=71.060..71.060 rows=10 loops=1)
5       → Seq Scan on yelp_categories c (cost=0.00..11287.31 rows=21767 width=74) (actual time=0.000..0.000 rows=10 loops=1)
6       Filter: (categories = 'Restaurants')
7       Rows Removed by Filter: 41049
8       → Index Scan using yelp_business_idx on yelp_business b (cost=0.00..10409.00 rows=10 width=80) (actual time=71.057..71.057 rows=10 loops=1)
9       Index Scan using yelp_tip_idx on yelp_tip t (cost=0.00..12735.16 rows=10 width=80) (actual time=0.000..0.000 rows=10 loops=1)
10       Filter: (city = 'Las Vegas')
11       Rows Removed by Filter: 1
12       → Index Scan using yelp_business_idx on yelp_business b (cost=0.00..10409.00 rows=10 width=80) (actual time=71.057..71.057 rows=10 loops=1)
13       Index Scan using yelp_tip_idx on yelp_tip t (cost=0.00..12735.16 rows=10 width=80) (actual time=0.000..0.000 rows=10 loops=1)
14 Planning Time: 275.867 ms
15 Execution Time: 117.960 ms

```

Fig. 5: Select query 1 after indexing

```

SELECT CAT, COUNT(CAT) FROM
(SELECT C.CATEGORIES AS CAT
FROM YELP_BUSINESS
AS B INNER JOIN YELP_CATEGORIES AS C
ON B.BUSINESS_ID = C.BUSINESS_ID
WHERE B.CITY = 'LAS VEGAS' )
AS TEMP GROUP BY CAT HAVING COUNT(CAT) =1;

```

Query 31: Select query 2

```

-- Query Plan
1  Sort Key: (categories)
2  Filter: (count(categories) = 1)
3  Rows Removed by Filter: 306
4  → Gather Merge (cost=10010.01..19770.16 rows=10 width=80) (actual time=476.609..4566.803 rows=10 loops=1)
5    Buckers Planned: 6
6    Buckers Launched: 6
7    → Sort (cost=10010.01..19770.16 rows=10 width=80) (actual time=476.609..79.404 rows=10 loops=1)
8      Sort Key: (categories)
9      Sort Method: quicksort Memory: 2048
10     Buckers: 3 Sort Method: quicksort Memory: 2048
11     Buckers: 3 Sort Method: quicksort Memory: 2048
12     → Parallel HashAggregate (cost=10010.01..16565.16 rows=10 width=80) (actual time=476.110..79.106 rows=10 loops=1)
13       Hash Key: (categories)
14       Buckers: 1 Memory Usage: 2168
15       Buckers: 1 Memory Usage: 2168
16       Buckers: 1 Memory Usage: 2168
17       Buckers: 1 Memory Usage: 2168
18       → Parallel Hash Join (cost=10010.01..16565.16 rows=10 width=80) (actual time=476.107..41.106 rows=10 loops=1)
19         Hash Cond: (t.business_id = b.business_id)
20         → Parallel Seq Scan on yelp_categories t (cost=0.00..11287.31 rows=21767 width=74) (actual time=0.000..0.000 rows=10 loops=1)
21         → Parallel Seq Scan on yelp_business b (cost=0.00..16666.28 rows=11889 width=80) (actual time=0.437..43.961 rows=26775 loops=1)
22         Filter: (city = 'Las Vegas')
23         Rows Removed by Filter: 1
24 Planning Time: 8.284 ms
25 Execution Time: 4566.803 ms

```

Fig. 6: Select query 2 before indexing

```

-- Query Plan
1  Sort Key: (categories)
2  Filter: (count(categories) = 1)
3  Rows Removed by Filter: 1000
4  → Gather Merge (cost=10010.01..19770.16 rows=10 width=80) (actual time=456.907..4566.713 rows=10 loops=1)
5    Buckers Planned: 6
6    Buckers Launched: 6
7    → Sort (cost=10010.01..19770.16 rows=10 width=80) (actual time=456.907..79.014 rows=10 loops=1)
8      Sort Key: (categories)
9      Sort Method: quicksort Memory: 2048
10     Buckers: 3 Sort Method: quicksort Memory: 2048
11     Buckers: 3 Sort Method: quicksort Memory: 2048
12     → Parallel HashAggregate (cost=10010.01..16565.16 rows=10 width=80) (actual time=456.904..41.106 rows=10 loops=1)
13       Hash Key: (categories)
14       Buckers: 1 Memory Usage: 2168
15       Buckers: 1 Memory Usage: 2168
16       Buckers: 1 Memory Usage: 2168
17       → Parallel Hash Join (cost=10010.01..16565.16 rows=10 width=80) (actual time=456.904..41.106 rows=10 loops=1)
18         Hash Cond: (t.business_id = b.business_id)
19         → Parallel Seq Scan on yelp_categories t (cost=0.00..11287.31 rows=21767 width=74) (actual time=0.000..0.000 rows=10 loops=1)
20         → Parallel Seq Scan on yelp_business b (cost=0.00..16666.28 rows=11889 width=80) (actual time=0.437..43.961 rows=26775 loops=1)
21         Filter: (city = 'Las Vegas')
22         Rows Removed by Filter: 1
23 Planning Time: 8.213 ms
24 Execution Time: 4566.713 ms

```

Fig. 7: Select query 2 after indexing

```

SELECT R.STARS, R."TEXT"
FROM YELP_REVIEW R INNER JOIN
YELP_CATEGORIES C
ON C.BUSINESS_ID = R.BUSINESS_ID
WHERE C.CATEGORIES = 'RESTAURANTS' LIMIT 1000;

```

Query 32: Select query 3

```

-- Query Plan
1  Limit (cost=10409.00..10409.00 rows=10 width=80) (actual time=100.872..426.868 rows=10 loops=1)
2   → Nested Loop (cost=10409.00..10409.00 rows=10 width=80) (actual time=100.872..426.868 rows=10 loops=1)
3     Join Filter: (t.business_id = b.business_id)
4     → Gather (cost=10409.00..10706.39 rows=10000 width=110) (actual time=100.872..1009.193 rows=10 loops=1)
5       Buckers Planned: 1
6       Buckers Launched: 1
7       → Parallel Hash Join (cost=10409.00..10809.79 rows=10000 width=110) (actual time=100.871..10.468 rows=10 loops=1)
8         Hash Cond: (t.business_id = b.business_id)
9         Hash Join (cost=10409.00..10706.39 rows=10000 width=110) (actual time=100.862..1.371 rows=10 loops=1)
10          → Parallel Seq Scan on yelp_tip t (cost=0.00..11287.31 rows=21767 width=74) (actual time=0.467..16.868 rows=10 loops=1)
11          → Parallel Hash Join (cost=10409.00..11287.31 rows=21767 width=110) (actual time=100.401..11.894 rows=10 loops=1)
12            Filter: (categories = 'Restaurants')
13            Rows Removed by Filter: 41049
14            → Hash Join (cost=10409.00..11287.31 rows=21767 width=110) (actual time=100.401..11.894 rows=10 loops=1)
15              Hash Cond: (t.business_id = b.business_id)
16              Hash Join (cost=10409.00..11287.31 rows=21767 width=110) (actual time=100.401..11.894 rows=10 loops=1)
17              → Parallel Seq Scan on yelp_business b (cost=0.00..16666.28 rows=11889 width=80) (actual time=0.437..43.961 rows=26775 loops=1)
18              Filter: (city = 'Las Vegas')
19              Rows Removed by Filter: 1
20 Planning Time: 8.438 ms
21 Execution Time: 1009.193 ms

```

Fig. 8: Select query 3 before indexing

```

-- Query Plan
1  Limit (cost=0.00..10409.00 rows=10 width=80) (actual time=71.058..117.960 rows=10 loops=1)
2   → Nested Loop (cost=0.00..10409.00 rows=10 width=80) (actual time=71.057..117.960 rows=10 loops=1)
3     Join Filter: (t.business_id = b.business_id)
4     → Nested Loop (cost=0.00..12735.16 rows=10 width=80) (actual time=71.060..71.060 rows=10 loops=1)
5       → Seq Scan on yelp_categories c (cost=0.00..11287.31 rows=21767 width=74) (actual time=0.000..0.000 rows=10 loops=1)
6       Filter: (categories = 'Restaurants')
7       Rows Removed by Filter: 41049
8       → Index Scan using yelp_business_idx on yelp_business b (cost=0.00..10409.00 rows=10 width=80) (actual time=71.057..71.057 rows=10 loops=1)
9       Index Scan using yelp_tip_idx on yelp_tip t (cost=0.00..12735.16 rows=10 width=80) (actual time=0.000..0.000 rows=10 loops=1)
10       Filter: (city = 'Las Vegas')
11       Rows Removed by Filter: 1
12       → Index Scan using yelp_business_idx on yelp_business b (cost=0.00..10409.00 rows=10 width=80) (actual time=71.057..71.057 rows=10 loops=1)
13       Index Scan using yelp_tip_idx on yelp_tip t (cost=0.00..12735.16 rows=10 width=80) (actual time=0.000..0.000 rows=10 loops=1)
14 Planning Time: 275.867 ms
15 Execution Time: 117.960 ms

```

Fig. 9: Select query 3 after indexing

## Indexes created

```

CREATE INDEX REVIEW_IDX
ON YELP_REVIEW(REVIEW_ID, BUSINESS_ID);

```

Query 33: Indexing query 1

```

CREATE INDEX BUSINESS_IDX
ON YELP_BUSINESS(BUSINESS_ID);

```

Query 34: Indexing query 2

```

CREATE INDEX CATEGORY_IDX
ON YELP_CATEGORIES(BUSINESS_ID);

```

Query 35: Indexing query 3



```
CREATE INDEX TIPS_IDX
ON YELP_TIP (BUSINESS_ID, USER_ID);
```

Query 36: Indexing query 4

## VIII. CRUD OPERATIONS

### a. Insert query

```
INSERT INTO yelp_business (business_id, name, address, city, state, postal_code, latitude, longitude, stars, review_count, is_open)
VALUES ('ABCDEF', 'ABCDEF', 'ABCDEF', 'ABCDEF', 'ABCDEF', 'ABCDEF', 'ABCDEF', 'ABCDEF', 'ABCDEF', 'ABCDEF', 'ABCDEF');
```

Fig. 10: Insert query

### b. Select query

```
SELECT * FROM yelp_business;
```

Fig. 11: Select query

### c. Index creation query

```
CREATE INDEX review_idx ON yelp_review (review_id, business_id);
```

Fig. 12: Indexing query

### d. Alter query

```
ALTER TABLE yelp_business
ALTER COLUMN business_id TYPE varchar(500),
ALTER COLUMN name TYPE char(200),
ALTER COLUMN address TYPE varchar(500),
ALTER COLUMN city TYPE varchar(500),
ALTER COLUMN state TYPE char(50),
ALTER COLUMN postal_code TYPE char(50),
ALTER COLUMN latitude TYPE double precision,
ALTER COLUMN longitude TYPE double precision,
ALTER COLUMN stars TYPE decimal(1,0),
ALTER COLUMN review_count TYPE numeric(100,0),
ALTER COLUMN is_open TYPE Numeric(1,0);
```

Fig. 13: Alter query

### e. Delete query

```
delete from yelp_business where business_id = 'ABCDEF';
```

Fig. 14: Delete query

### f. Order by query

```
SELECT state, COUNT(*) AS business_count
FROM yelp_business
GROUP BY state, city
ORDER BY business_count DESC
LIMIT 1000;
```

Fig. 15: Order by query

### g. Join query

```
SELECT B.name, B.latitude, B.longitude, C.categories, B.city, B.state, B.review_count
FROM yelp_business AS B INNER JOIN yelp_categories AS C
ON B.business_id = C.business_id where B.city = 'Las Vegas' and C.categories = 'Restaurants';
```

Fig. 16: Join query

## IX. WEB APPLICATION

Users can visualize states with most number of businesses by selecting the number of states they want to view the data for which is shown by a Bubble plot.

## Summary Dashboard

### Max business indexed

Plot a bubble plot of states with the most businesses indexed on Yelp data

Choose top "n" Cities

20

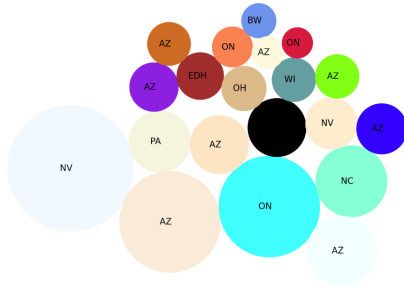


Fig. 17: Bubble plot of states with most businesses

## Reviews

Get a list of reviews for business of a category in a city

Choose state

ON

Choose city

Toronto

Choose category

3D Printing

Positive Reviews Word Cloud

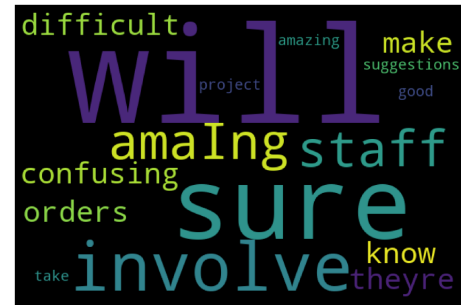


Fig. 19: Positive reviews for a business category in a city

Users can visualize the locations of the businesses with specific categories by selecting the state, business, and the category they wish to view the information about.

## Map

Get a locations of a category of business in a city on a map

Choose state

ON

Choose city

Toronto

Choose category

Accountants

### Map of 'Accountants' in 'Toronto'

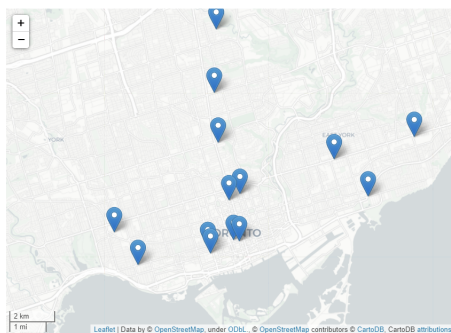


Fig. 18: Locations of business categories in a city

Users can visualize the most used positive and negative reviews by the consumers.

Negative Reviews Word Cloud



Fig. 20: Negative reviews for a business category in a city

Users can visualize the possible business opportunities for a particular category in a city since those cities would have the least number of businesses in that particular category.

## Business opotunities

We will list down categories which have only one business instance in the selected city

Choose state

NLK

Choose city

Livingston

	cat
0	American (New)
1	Arts & Entertainment
2	Buffets
3	Cafes
4	Childrens Clothing
5	Coffee & Tea
6	Dance Clubs
7	Department Stores
8	Diners
9	Donuts

Fig. 21: Business opportunities based on categories in a city



Users can query the database to get any relevant information they would like to fetch.

Top 15 tips for a selected query

```
select text from yelp_tip where business_id in (select yelp_business.business_id from yelp_business wh
```

0	
0	Offers safe disposal of freon appliances (air conditioners, refrigerators, etc.) for \$10.
1	Really fun show!
2	The CLO shows at the cabaret and during the summer are always a fun time!
3	Be sure to check out the mausoleum—it has amazing stained glass windows!
4	Panther Hollow Trail goes through Schenley Park, and connects to trail through ravin
5	Beautiful in springtime, esp. with the flowering trees!
6	Largest Catholic cemetery in PGH.
7	Yelp deal—\$35 for \$50
8	The worlds only in-airport robot repair shop!
9	Two doors down from their old corner location. The guy said someone else is moving

Fig. 22: Top 15 tips

## X. FURTHER ENHANCEMENTS

The database that we designed is good for businesses to get enough insights, but there is a scope for a few improvements. Since this is an already populated database and static in nature, we can improve the time taken to fetch the data by creating views for the queries that are frequently accessed on the front end.

## REFERENCES

[1] <https://www.kaggle.com/code/jagangupta/what-s-in-a-review-yelp-ratings-eda>

[2] <https://www.yelp.com/dataset>

[3] [https://github.com/sai2499/Kryptonite\\_Yelp\\_Dataset](https://github.com/sai2499/Kryptonite_Yelp_Dataset)

[4] <https://www.yelp.com/dataset/documentation/main>

[5] <https://github.com/Yelp/dataset-examples>

[6] <https://towardsdatascience.com/converting-yelp-dataset-to-csv-using-pandas-2a4c8f03bd88>