Part 1: Yelp Dataset Profiling and Understanding

- 1. Profile the data by finding the total number of records for each of the tables below:
- i. Attribute table =
- ii. Business table =
- iii. Category table =
- iv. Checkin table =
- v. elite years table =
- vi. friend table =
- vii. hours table =
- viii. photo table =
- ix. review table =
- x. tip table =
- xi. user table =

Solution=

- i. Attribute table =1000
- ii. Business table =1000
- iii. Category table =1000
- iv. Checkin table =1000
- v. elite_years table =1000
- vi. friend table = 1000
- vii. hours table =1000
- viii. photo table =1000
- ix. review table = 1000

- x. tip table = 1000
- xi. user table =1000
- 2. Find the total distinct records by either the foreign key or primary key for each table. If two foreign keys are listed in the table, please specify which foreign key.
- i. Business =
- ii. Hours =
- iii. Category =
- iv. Attribute =
- v. Review =
- vi. Checkin =
- vii. Photo =
- viii. Tip =
- ix. User =
- x. Friend =
- xi. Elite years =

Note: Primary Keys are denoted in the ER-Diagram with a yellow key icon.

Solution=

- i. Business =1000 distinct records for primary key 'id' of business table
- ii. Hours =1562 distinct records for foreign key 'business id' of hours table
- iii. Category =2643 distinct records for foreign key 'business_id' of category table
- iv. Attribute =1115 distinct records for foreign key 'business_id' of attribute table
- v. Review =1000 distinct records for primary key 'id',8090 distinct records for primary key 'business_id' and 9581 distinct records for primary key 'user_id' from the review table
- vi. Checkin = 493 distinct records for primary key 'user id' from the checkin table
- vii. Photo =6493 distinct records for foreign key 'business_id' and 1000 distinct records for primary key 'id' from the photo table

viii. Tip =3979 distinct records for foreign key 'business_id' and 537 distinct records for foreign key 'user_id' from the tip table

- ix. User = 1000 distinct records for primary key 'id' of user table
- x. Friend = 11 distinct records for foreign key 'id' from the friend table
- xi. Elite_years =2780 distinct records for foreign key 'user_id' from the elite_years table
 - 3. Are there any columns with null values in the Users table? Indicate "yes," or "no."

Answer: "no"

SQL code used to arrive at answer:

SELECT*

FROM user

WHERE

id IS NULL OR

name IS NULL OR

review_count IS NULL OR

yelping_since IS NULL OR

useful IS NULL OR

funny IS NULL OR

cool IS NULL OR

fans IS NULL OR

average_stars IS NULL OR

compliment_hot IS NULL OR

compliment_more IS NULL OR

compliment profile IS NULL OR

compliment cute IS NULL OR

```
compliment_list IS NULL OR
compliment_note IS NULL OR
compliment_plain IS NULL OR
compliment_cool IS NULL OR
compliment_funny IS NULL OR
compliment_writer IS NULL OR
compliment_photos IS NULL;
```

- 4. For each table and column listed below, display the smallest (minimum), largest (maximum), and average (mean) value for the following fields:
 - i. Table: Review, Column: Stars

min: 1 max: 5 avg: 3.7082

ii. Table: Business, Column: Stars

min:1 max: 5 avg:3.6549

iii. Table: Tip, Column: Likes

min: 0 max: 2 avg: 0.0144

iv. Table: Checkin, Column: Count

min: 1 max: 53 avg: 1.9414

v. Table: User, Column: Review count

min:0 max:2000 avg: 24.2995

5. List the cities with the most reviews in descending order:

SQL code used to arrive at answer: SELECT SUM(review_count) AS reviews,city

FROM business

GROUP BY city

ORDER BY SUM(review_count) DESC;

Copy and Paste the Result Below:

reviews city					
+-	+	-+			
	82854 Las Vegas	1			
I	34503 Phoenix	I			
I	24113 Toronto	I			
	20614 Scottsdale	1			
	12523 Charlotte	1			
l	10871 Henderson	1			
	10504 Tempe	1			
ı	9798 Pittsburgh	Ī			

```
9448 | Montréal
  8112 | Chandler
  6875 | Mesa
  6380 | Gilbert
  5593 | Cleveland
  5265 | Madison
  4406 | Glendale
  3814 | Mississauga
  2792 | Edinburgh
  2624 | Peoria
                     I
  2438 | North Las Vegas |
  2352 | Markham
  2029 | Champaign
  1849 | Stuttgart
  1520 | Surprise
   1465 | Lakewood
   1155 | Goodyear
                       1
+----+
(Output limit exceeded, 25 of 362 total rows shown)
  6. Find the distribution of star ratings to the business in the following cities:
i. Avon
SQL code used to arrive at answer:SELECT stars,COUNT(stars)
                   FROM business
                   WHERE city='Avon'
```

GROUP BY stars;

+----+

| stars | COUNT(stars) |

+----+

| 1.5 | 1 |

| 2.5 | 2 |

| 3.5 | 3 |

| 4.0 | 2 |

| 4.5 | 1 |

| 5.0 | 1 |

+----+

ii. Beachwood

SQL code used to arrive at answer: SELECT stars, COUNT(stars)

FROM business

WHERE city='Beachwood'

GROUP BY stars;

+----+

| stars | COUNT(stars) |

+----+

| 2.0 | 1 |

| 2.5 | 1 |

3.0 2

| 3.5 | 2 |

| 4.0 | 1 |

7. Find the top 3 users based on their total number of reviews:

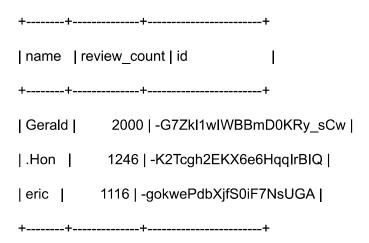
SQL code used to arrive at answer: SELECT name,review_count FROM user

GROUP BY name

ORDER BY review_count DESC

LIMIT 3;

Copy and Paste the Result Below:



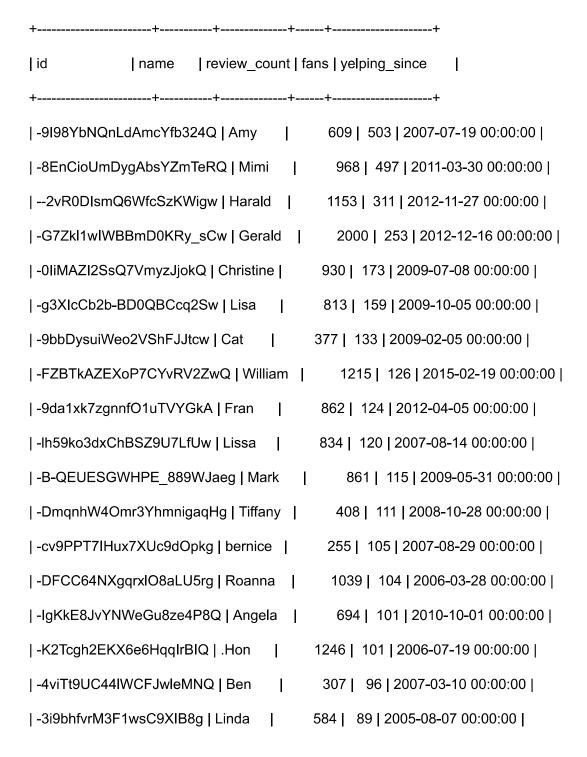
8. Does posing more reviews correlate with more fans?

Please explain your findings and interpretation of the results:

No,Posing more reviews is not directly related with more fans.It also depends on the yelping since factor.

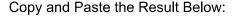
The longer they have been yelping, the more number of genuine reviews they give, hence resulting in more number of fans.

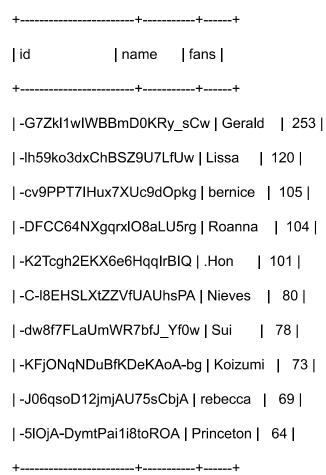
In most of the cases,we see that with increase in the number of reviews,a corresponding increase in the number of fans is not seen.



```
|-ePh4Prox7ZXnEBNGKyUEA|Jessica| 220| 84|2009-01-12 00:00:00|
 9. Are there more reviews with the word "love" or with the word "hate" in them?
      Answer: Number of reviews with the word 'love' in them =1780
       Number of reviews with the word 'hate' in them =232
       Hence, there are more reviews with the word love in it.
  SQL code used to arrive at answer:
SELECT COUNT(text)
FROM review
WHERE text LIKE '%love%';
SELECT COUNT(text)
FROM review
WHERE text LIKE '%hate%';
  10. Find the top 10 users with the most fans:
SQL code used to arrive at answer:
SELECT id,name,fans
FROM user
GROUP BY name
ORDER BY fans DESC
```

LIMIT 10;





Part 2: Inferences and Analysis

1. Pick one city and category of your choice and group the businesses in that city or category by their overall star rating.

Compare the businesses with 2-3 stars to the businesses with 4-5 stars and answer the following questions. Include your code.

i. Do the two groups you chose to analyze have a different distribution of hours?

I chose 'Toronto' as city and 'Food' as category. Yes, they have different distribution of working hours.

Working hours for 2-3 star rating businesses are 7 and working hours for 4-5 star rating businesses are 13.

SELECT stars

,CASE WHEN stars>=4 THEN "4-5 stars"

WHEN stars>=2 THEN "2-3 stars"

ELSE "Below 2"

END Star_rank

,business.city,COUNT(DISTINCT hours.business_id) AS Company_count

,COUNT(hours.hours) AS working_hours

,category.category

FROM((business JOIN hours ON business.id=hours.business id)

JOIN category ON business.id=category.business_id)

WHERE city='Toronto' AND category= 'Food'

GROUP BY Star_rank;

ii. Do the two groups you chose to analyze have a different number of reviews?

I chose 'Toronto' as city and 'Food' as category. The two categories have different number of reviews.

Number of reviews for 2-3 star rating businesses are 70 and number of reviews for 4-5 star rating businesses are 272.

SELECT stars

,CASE WHEN stars>=4 THEN "4-5 stars"

WHEN stars>=2 THEN "2-3_stars"

ELSE "Below 2"

END Star rank

,business.city,COUNT(DISTINCT hours.business_id) AS Company_count

,COUNT(hours.hours) AS working_hours

,SUM(review count) AS Number of reviews

,category.category

FROM((business JOIN hours ON business.id=hours.business_id)

JOIN category ON business.id=category.business_id)

WHERE city='Toronto' AND category= 'Food'

GROUP BY Star rank;

iii. Are you able to infer anything from the location data provided between these two groups? Explain.

2-3 star rating businesses are located near one another.4-5 star rating businesses are apart from one anothe.

SQL code used for analysis:

SELECT stars

,CASE WHEN stars>=4 THEN '4-5 Stars'

WHEN stars>=2 THEN '2-3 Stars'

ELSE 'below 2'

END Star_Rank

,category,city,postal_code,address,neighborhood

FROM business INNER JOIN category ON business.id=category.business_id

WHERE city='Las Vegas' AND category='Shopping'

ORDER BY Star_Rank;

2. Group business based on the ones that are open and the ones that are closed.

What differences can you find between the ones that are still open and the ones that are closed?

List at least two differences and the SQL code you used to arrive at your answer.

- i. Difference 1:The businesses that are open have a greater number of reviews when compared to the businesses that are closed.
- ii. Difference 2:The average stars given are very closed to each other 3.68 vs. 3.52.

We can assume that businesses which are closed was not solely due to poor services or poor quality.

SQL code used for analysis:

SELECT

CASE

WHEN is open='0' THEN 'closed'

WHEN is_open='1' THEN 'open'

ELSE 'NO Status'

END Status

,COUNT(DISTINCT id) AS Company num

,SUM(review_count) AS num_of_reviews

,ROUND(AVG(review_count),2) AS AVG_REVIEW

,ROUND(AVG(stars),2) AS AVG_STARS

FROM business

GROUP BY Status;

3. For this last part of your analysis, you are going to choose the type of analysis you want to conduct on the Yelp dataset and are going to prepare the data for analysis.

Ideas for analysis include: Parsing out keywords and business attributes for sentiment analysis,

clustering businesses to find commonalities or anomalies between them,

predicting the overall star rating for a business,

predicting the number of fans a user will have, and so on.

These are just a few examples to get you started, so feel free to be creative and come up with your own problem you want to solve.

Provide answers, in-line, to all of the following:

i. Indicate the type of analysis you chose to do:

The analysis done is to find out which categories of business have higher stars rating and more number of companies.

ii. Write 1-2 brief paragraphs on the type of data you will need for your analysis and why you chose that data:

For this analysis, we will need data such as id, stars, and review count from the business table and category table.

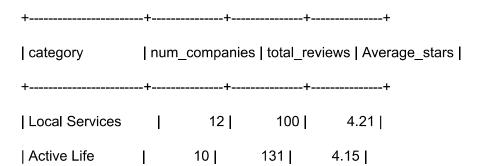
We are counting the numbers of companies within each category,

the average stars given by the consumers to see how they perform,

and the total reviews given to see if the data is relevant and ensure it's not biased.

We are only analysing the categories with at least 10 companies and an average of 3+ stars to reduce any irrelevant data.

iii. Output of your finished dataset:



Health & Medica	al	17	203	4.09
Home Services		16	94	4.0
Shopping	ļ	30	977	3.98
Beauty & Spas	1	13	119	3.88
American (Trad	11	1128	3.82	
Food		23	1781	3.78
Bars	I	17	1322	3.5
Nightlife	1	20	1351	3.48
Restaurants		71	4504	3.46
+	+	+	+	+

iv. Provide the SQL code you used to create your final dataset:

SELECT category,

COUNT(DISTINCT id) AS num_companies,

SUM(review_count) AS total_reviews,

ROUND(AVG(stars),2) AS Average_stars

FROM business INNER JOIN category

ON business.id=category.business_id

GROUP BY category

HAVING Average_stars>=3 AND num_companies >=10

ORDER BY Average_stars DESC;