Data Scientist Role Play: Profiling and Analyzing the Yelp Dataset Coursera Worksheet

This is a 2-part assignment. In the first part, you are asked a series of questions that will help you profile and understand the data just like a data scientist would. For this first part of the assignment, you will be assessed both on the correctness of your findings, as well as the code you used to arrive at your answer. You will be graded on how easy your code is to read, so remember to use proper formatting and comments where necessary.

In the second part of the assignment, you are asked to come up with your own inferences and analysis of the data for a particular research question you want to answer. You will be required to prepare the dataset for the analysis you choose to do. As with the first part, you will be graded, in part, on how easy your code is to read, so use proper formatting and comments to illustrate and communicate your intent as required.

For both parts of this assignment, use this "worksheet." It provides all the questions you are being asked, and your job will be to transfer your answers and SQL coding where indicated into this worksheet so that your peers can review your work. You should be able to use any Text Editor (Windows Notepad, Apple TextEdit, Notepad ++, Sublime Text, etc.) to copy and paste your answers. If you are going to use Word or some other page layout application, just be careful to make sure your answers and code are lined appropriately.
In this case, you may want to save as a PDF to ensure your formatting remains intact for you reviewer.

**Part 1: Yelp Dataset Profiling and Understanding**

1. Profile the data by finding the total number of records for each of the tables below:

```
Select Count(*) AS Total_Number
From Table
```

i. Attribute table = 10000
ii. Business table = 10000
iii. Category table = 10000
iv. Checkin table = 10000
v. elite_years table = 10000
vi. friend table = 10000
vii. hours table = 10000
viii. photo table = 10000
ix. review table = 10000
x. tip table = 10000
xi. user table = 10000

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.

```
Select Count(Distinct (key)) AS Total_Key
From Table
```

i. Business = id: 10000
ii. Hours = business_id: 1562
iii. Category = business_id: 2643
iv. Attribute = business_id: 1115
v. Review = business_id: 8090, id: 10000, user_id: 9581
vi. Checkin = business_id: 493
vii. Photo = business_id: 6493, id: 10000
viii. Tip = business_id: 3979
ix. User = id: 10000
x. Friend = user_id: 11
xi. Elite_years = user_id: 2780

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

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 Count(*)
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:

```
Select min(ColumnName) ,max(ColumnName), avg(ColumnName)
From Table
```

  i. Table: Review, Column: Stars

        min:    1       max:    5       avg: 3.708

  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 city, sum(review_count) as Reviews
From business
Group by city
Order by Reviews DESC
```

Copy and Paste the Result Below:

```
+-----------------+---------+
| city            | Reviews |
+-----------------+---------+
| Las Vegas       |   82854 |
| Phoenix         |   34503 |
| Toronto         |   24113 |
| Scottsdale      |   20614 |
| Charlotte       |   12523 |
| Henderson       |   10871 |
| Tempe           |   10504 |
| Pittsburgh      |    9798 |
| Montréal        |    9448 |
| Chandler        |    8112 |
| Mesa            |    6875 |
| Gilbert         |    6380 |
| Cleveland       |    5593 |
| Madison         |    5265 |
| Glendale        |    4406 |
| Mississauga     |    3814 |
| Edinburgh       |    2792 |
| Peoria          |    2624 |
| North Las Vegas |    2438 |
| Markham         |    2352 |
| Champaign       |    2029 |
| Stuttgart       |    1849 |
| Surprise        |    1520 |
| Lakewood        |    1465 |
| Goodyear        |    1155 |
+-----------------+---------+
(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, sum(review_count) AS count
From business
where city is 'Avon'
group by stars
```

Copy and Paste the Resulting Table Below (2 columns – star rating and count):

```
+-------+-------+
| stars | count |
+-------+-------+
|   1.5 |    10 |
|   2.5 |     6 |
|   3.5 |    88 |
|   4.0 |    21 |
|   4.5 |    31 |
|   5.0 |     3 |
+-------+-------+
```

ii. Beachwood

```sql
Select stars, sum(review_count) AS count
From business
where city is 'Beachwood'
group by stars
```

Copy and Paste the Resulting Table Below (2 columns – star rating and count):

```
+-------+-------+
| stars | count |
+-------+-------+
|   2.0 |     8 |
|   2.5 |     3 |
|   3.0 |    11 |
|   3.5 |     6 |
|   4.0 |    69 |
|   4.5 |    17 |
|   5.0 |    23 |
+-------+-------+
```

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

SQL code used to arrive at answer:

```sql
Select name, review_count
From user
order by review_count desc
limit 3
```

Copy and Paste the Result Below:

```
+--------+--------------+
| name   | review_count |
+--------+--------------+
| Gerald |         2000 |
| Sara   |         1629 |
| Yuri   |         1339 |
+--------+--------------+
```

8. Does posing more reviews correlate with more fans?

   Please explain your findings and interpretation of the results:

As table below illustrates, posing more reviews does not necessarily
correlate with more fans. For example, although, Gerald has posed the most
reviews, he has fewer fans in comparison with Mimi. Therefore, sorting the
users in descending order based on their total number of reviews does not
sort the fans in the same order, meaning that there is not a correlation
between the total number of reviews and number of fans.

```
Select name, review_count, fans
From user
Order by review_count DESC
```

```
+-----------+--------------+------+
| name      | review_count | fans |
+-----------+--------------+------+
| Gerald    |         2000 |  253 |
| Sara      |         1629 |   50 |
| Yuri      |         1339 |   76 |
| .Hon      |         1246 |  101 |
| William   |         1215 |  126 |
| Harald    |         1153 |  311 |
| eric      |         1116 |   16 |
| Roanna    |         1039 |  104 |
| Mimi      |          968 |  497 |
| Christine |          930 |  173 |
| Ed        |          904 |   38 |
| Nicole    |          864 |   43 |
| Fran      |          862 |  124 |
| Mark      |          861 |  115 |
| Christina |          842 |   85 |
| Dominic   |          836 |   37 |
| Lissa     |          834 |  120 |
| Lisa      |          813 |  159 |
| Alison    |          775 |   61 |
| Sui       |          754 |   78 |
| Tim       |          702 |   35 |
| L         |          696 |   10 |
| Angela    |          694 |  101 |
| Crissy    |          676 |   25 |
| Lyn       |          675 |   45 |
+-----------+--------------+------+
(Output limit exceeded, 25 of 10000 total rows shown)
```

9. Are there more reviews with the word "love" or with the word "hate" in them?

Answer: love

SQL code used to arrive at answer:

```sql
Select Count(*)
From review
WHERE text LIKE "%love%" ---> 1780


Select Count(*)
From review
WHERE text LIKE "%hate%" ---> 232
```

10. Find the top 10 users with the most fans:

SQL code used to arrive at answer:

```sql
Select name, fans
From user
Order by fans DESC
limit 10
```

Copy and Paste the Result Below:

```
+-----------+------+
| name      | fans |
+-----------+------+
| Amy       |  503 |
| Mimi      |  497 |
| Harald    |  311 |
| Gerald    |  253 |
| Christine |  173 |
| Lisa      |  159 |
| Cat       |  133 |
| William   |  126 |
| Fran      |  124 |
| Lissa     |  120 |
+-----------+------+
```

**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.

**City:** Las Vegas          **Category:** Shopping

```
Select b.city, c.category
From business b inner join category c on b.id = c.business_id
where city is 'Las Vegas'
```

#then I picked Resorts as my category.

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

Yes. They all have a different time distribution each day.

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

Yes, 2-3 star has 6 reviews, on the other hand 4-5 stars have 32 and 4 reviews.

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

Yes, the restaurant address which has 2-3 star is 3808 E Tropicana Ave and the restaurants address which have 4-5 stars are 1000 Sceniss Loop Dr and 3555 W Reno Ave, Ste F.

SQL code used for analysis:
```
Select b.name, b.City, b.stars, b.review_count, c.category, h.hours, b.a
ddres,
CASE
    WHEN hours LIKE "%monday%" THEN 1
    WHEN hours LIKE "%tuesday%" THEN 2
    WHEN hours LIKE "%wednesday%" THEN 3
    WHEN hours LIKE "%thursday%" THEN 4
    WHEN hours LIKE "%friday%" THEN 5
    WHEN hours LIKE "%saturday%" THEN 6
    WHEN hours LIKE "%sunday%" THEN 7
    END AS ord,
 CASE
    WHEN B.stars BETWEEN 2 AND 3 THEN '2-3 Stars'
    WHEN B.stars BETWEEN 4 AND 5 THEN '4-5 Stars'
    END AS star_rating
from business b inner join category c on b.id = c.business_id inner join
hours h on b.id = h.business_id
where city is 'Las Vegas' and category is 'Shopping'
GROUP BY stars,ord
```

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 open ones have much more number of review than the ones that are closed.

ii. Difference 2:

The open restaurants have slightly higher average stars than the closed ones.

SQL code used for analysis:

```
Select Count(DISTINCT(id)), avg(review_count), sum(review_count), avg(stars)
, is_open
FROM business
GROUP BY is_open
```

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:

If an entrepreneur is going to make an attempt to open a new bar, I wanted to suggest to him/her in which city it would be most logical to open this business.

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

First of all, I had to learn about the cities with bars in order to make comparisons. After that, I classified the number of stars of the bars within themselves to make a better comparison. I added the neighborhoods to my query so that I might catch a clue. After my code was finalized, I made my query and interpreted the output as analysis.

iii. Output of your finished dataset:

```
+-------+------------------------------------+-----------------------+----------+------------+------------+
| stars | name                               | neighborhood          | category | city       | star_rating |
+-------+------------------------------------+-----------------------+----------+------------+------------+
|   3.0 | Irish Republic                     |                       | Bars     | Chandler   | 2-3 stars  |
|   4.0 | Nabers Music, Bar & Eats           |                       | Bars     | Chandler   | 3-4 stars  |
|   2.0 | Iron City Grille                   |                       | Bars     | Coraopolis | 1-2 stars  |
|   3.0 | Brubaker's Pub                     |                       | Bars     | Hudson     | 2-3 stars  |
|   3.5 | Hi Scores - Blue Diamond           | Southwest             | Bars     | Las Vegas  | 3-4 stars  |
|   4.0 | TWIISTED Burgers & Sushi           |                       | Bars     | Medina     | 3-4 stars  |
|   4.0 | Eklectic Pie - Mesa                |                       | Bars     | Mesa       | 3-4 stars  |
|   3.0 | The Erin Mills Pump & Patio        |                       | Bars     | Mississauga| 2-3 stars  |
|   3.0 | Restaurant Rosalie                 | Ville-Marie           | Bars     | Montréal   | 2-3 stars  |
|   4.5 | The Wine Mill                      |                       | Bars     | Peninsula  | 4-5 stars  |
|   3.0 | Gallagher's                        |                       | Bars     | Phoenix    | 2-3 stars  |
|   4.0 | Bootleggers Modern American Smokehouse |                   | Bars     | Phoenix    | 3-4 stars  |
|   2.5 | The Fox & Fiddle                   | Greektown             | Bars     | Toronto    | 2-3 stars  |
|   3.5 | The Charlotte Room                 | Entertainment District| Bars     | Toronto    | 3-4 stars  |
|   4.0 | Halo Brewery                       | Wallace Emerson       | Bars     | Toronto    | 3-4 stars  |
|   4.5 | Cabin Fever                        | High Park             | Bars     | Toronto    | 4-5 stars  |
|   4.0 | Cabin Club                         |                       | Bars     | Westlake   | 3-4 stars  |
+-------+------------------------------------+-----------------------+----------+------------+------------
```

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

```sql
Select b.stars, b.name, b.neighborhood, c.category,b.city,
CASE
        WHEN B.stars BETWEEN 0 AND 1 THEN '2-3 stars'
        WHEN B.stars BETWEEN 1 AND 2 THEN '1-2 stars'
        WHEN B.stars BETWEEN 2 AND 3 THEN '2-3 stars'
        WHEN B.stars BETWEEN 3 AND 4 THEN '3-4 stars'
        WHEN B.stars BETWEEN 4 AND 5 THEN '4-5 stars'
        END AS star_rating
from business b inner join category c on b.id = c.business_id
where c.category is 'Bars'
order by city, star_rating
```