```
1    Data Scientist Role Play: Profiling and Analyzing the Yelp Dataset Coursera Worksheet
2
3    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.
4
5    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.
6
7    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.
8    In this case, you may want to save as a PDF to ensure your formatting remains intact
     for you reviewer.
9
10
11
12   Part 1: Yelp Dataset Profiling and Understanding
13
14   1. Profile the data by finding the total number of records for each of the tables below:
15
16   i. Attribute table = 10000
17
18   SELECT COUNT(*)
19   FROM attribute
20
21   ii. Business table = 10000
22
23   SELECT COUNT(*)
24   FROM business
25
26   iii. Category table = 10000
27
28   SELECT COUNT(*)
29   FROM category
30
31   iv. Checkin table = 10000
32
33   SELECT COUNT(*)
34   FROM checkin
35
36   v. elite_years table = 10000
37
38   SELECT COUNT(*)
39   FROM elite_years
40
41   vi. friend table = 10000
42
43   SELECT COUNT(*)
44   FROM friend
45
46   vii. hours table = 10000
47
48   SELECT COUNT(*)
49   FROM hours
50
51   viii. photo table = 10000
52
53   SELECT COUNT(*)
```

```
54    FROM photo
55
56    ix. review table = 10000
57
58    SELECT COUNT(*)
59    FROM review
60
61    x. tip table = 10000
62
63    SELECT COUNT(*)
64    FROM tip
65
66    xi. user table = 10000
67
68    SELECT COUNT(*)
69    FROM user
70
71
72    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.
73
74    i. Business = id:10000
75
76    SELECT COUNT(DISTINCT(id))
77    FROM business
78
79    ii. Hours = business_id:1562
80
81    SELECT COUNT(DISTINCT(business_id))
82    FROM hours
83
84    iii. Category = business_id:2643
85
86    SELECT COUNT(DISTINCT(business_id))
87    FROM category
88
89    iv. Attribute = business_id:1115
90
91    SELECT COUNT(DISTINCT(business_id))
92    FROM attribute
93
94    v. Review = id: 10000, business_id: 8090, user_id: 9581
95
96    SELECT COUNT(DISTINCT(id)), COUNT(DISTINCT(business_id)), COUNT(DISTINCT(user_id))
97    FROM review
98
99    vi. Checkin = business_id:493
100
101   SELECT COUNT(DISTINCT(business_id))
102   FROM checkin
103
104   vii. Photo = id:10000 business_id:6493
105
106   SELECT COUNT(DISTINCT(id)), COUNT(DISTINCT(business_id))
107   FROM photo
108
109   viii. Tip = user_id:537 business_id:3979
110
111   SELECT COUNT(DISTINCT(user_id)), COUNT(DISTINCT(business_id))
112   FROM tip
113
114   ix. User = id:10000
115
116   SELECT COUNT(DISTINCT(id))
117   FROM user
118
119   x. Friend = user_id:11
120
121   SELECT COUNT(DISTINCT(user_id))
```

```
122    FROM friend
123
124    xi. Elite_years = user_id:2780
125
126    SELECT COUNT(DISTINCT(user_id))
127    FROM elite_years
128
129    Note: Primary Keys are denoted in the ER-Diagram with a yellow key icon.
130
131
132
133    3. Are there any columns with null values in the Users table? Indicate "yes," or "no."
134
135        Answer: "NO"
136
137
138        SQL code used to arrive at answer:
139
140    SELECT COUNT(*)
141    FROM user
142    WHERE id IS NULL
143    OR name IS NULL
144    OR review_count IS NULL
145    OR yelping_since IS NULL
146    OR useful IS NULL
147    OR funny IS NULL
148    OR cool IS NULL
149    OR fans IS NULL
150    OR average_stars IS NULL
151    OR compliment_hot IS NULL
152    OR compliment_more IS NULL
153    OR compliment_profile IS NULL
154    OR compliment_cute IS NULL
155    OR compliment_list IS NULL
156    OR compliment_note IS NULL
157    OR compliment_plain IS NULL
158    OR compliment_cool IS NULL
159    OR compliment_funny IS NULL
160    OR compliment_writer IS NULL
161    OR compliment_photos IS NULL
162
163
164    4. For each table and column listed below, display the smallest (minimum), largest
       (maximum), and average (mean) value for the following fields:
165
166        i. Table: Review, Column: Stars
167
168            min:    1    max:    5    avg: 3.7082
169
170            SELECT MIN(stars), MAX(stars), AVG(stars)
171            FROM review
172
173        ii. Table: Business, Column: Stars
174
175            min: 1.0    max: 5.0    avg: 3.6549
176
177            SELECT MIN(stars), MAX(stars), AVG(stars)
178            FROM business
179
180        iii. Table: Tip, Column: Likes
181
182            min:    0    max:    2    avg:    0.0144
183
184            SELECT MIN(likes), MAX(likes), AVG(likes)
185            FROM tip
186
187        iv. Table: Checkin, Column: Count
188
189            min:    1    max:    53    avg: 1.9414
```

```
190
191            SELECT MIN(count), MAX(count), AVG(count)
192            FROM checkin
193
194       v. Table: User, Column: Review_count
195
196            min:    0   max:    2000    avg:  24.2995
197
198            SELECT MIN(review_count), MAX(review_count), AVG(review_count)
199            FROM user
200
201
202    5. List the cities with the most reviews in descending order:
203
204        SQL code used to arrive at answer:
205
206    SELECT SUM(review_count), city
207    FROM business
208    GROUP BY city
209    ORDER BY SUM(review_count) DESC
210
211        Copy and Paste the Result Below:
212
213        +-------------------+-----------------+
214    | SUM(review_count) | city            |
215    +-------------------+-----------------+
216    |             82854 | Las Vegas       |
217    |             34503 | Phoenix         |
218    |             24113 | Toronto         |
219    |             20614 | Scottsdale      |
220    |             12523 | Charlotte       |
221    |             10871 | Henderson       |
222    |             10504 | Tempe           |
223    |              9798 | Pittsburgh      |
224    |              9448 | Montréal        |
225    |              8112 | Chandler        |
226    |              6875 | Mesa            |
227    |              6380 | Gilbert         |
228    |              5593 | Cleveland       |
229    |              5265 | Madison         |
230    |              4406 | Glendale        |
231    |              3814 | Mississauga     |
232    |              2792 | Edinburgh       |
233    |              2624 | Peoria          |
234    |              2438 | North Las Vegas |
235    |              2352 | Markham         |
236    |              2029 | Champaign       |
237    |              1849 | Stuttgart       |
238    |              1520 | Surprise        |
239    |              1465 | Lakewood        |
240    |              1155 | Goodyear        |
241    +-------------------+-----------------+
242    (Output limit exceeded, 25 of 362 total rows shown)
243
244
245
246    6. Find the distribution of star ratings to the business in the following cities:
247
248    i. Avon
249
250    SQL code used to arrive at answer:
251
252    SELECT stars AS 'Star_Rating', count(stars) AS Count
253    FROM business
254    WHERE city = 'Avon'
255    GROUP BY stars;
256
257    Copy and Paste the Resulting Table Below (2 columns - star rating and count):
258
```

```
259    +-------------------+-------+
260    | SUM(review_count) | stars |
261    +-------------------+-------+
262    |                10 |   1.5 |
263    |                 6 |   2.5 |
264    |                88 |   3.5 |
265    |                21 |   4.0 |
266    |                31 |   4.5 |
267    |                 3 |   5.0 |
268    +-------------------+-------+
269
270    ii. Beachwood
271
272    SQL code used to arrive at answer:
273
274    SELECT stars AS 'Star_Rating', count(stars) AS Count
275    FROM business
276    WHERE city = 'Beachwood'
277    GROUP BY stars;
278
279    Copy and Paste the Resulting Table Below (2 columns – star rating and count):
280
281    +-------------------+-------+
282    | SUM(review_count) | stars |
283    +-------------------+-------+
284    |                 8 |   2.0 |
285    |                 3 |   2.5 |
286    |                11 |   3.0 |
287    |                 6 |   3.5 |
288    |                69 |   4.0 |
289    |                17 |   4.5 |
290    |                23 |   5.0 |
291    +-------------------+-------+
292
293    7. Find the top 3 users based on their total number of reviews:
294
295        SQL code used to arrive at answer:
296
297    SELECT id, name, review_count
298    FROM user
299    ORDER BY review_count DESC
300    LIMIT 3
301
302
303        Copy and Paste the Result Below:
304
305    +------------------------+--------+--------------+
306    | id                     | name   | review_count |
307    +------------------------+--------+--------------+
308    | -G7Zkl1wIWBBmD0KRy_sCw | Gerald |         2000 |
309    | -3s52C4zL_DHRK0ULG6qtg | Sara   |         1629 |
310    | -8lbUNlXVSoXqaRRiHiSNg | Yuri   |         1339 |
311    +------------------------+--------+--------------+
312
313
314    8. Does posing more reviews correlate with more fans?
315
316        Please explain your findings and interpretation of the results:
317
318    Posing more reviews does not correlate with more fans. Amy has the most fans with the
       number of 503 and she has 609 reviews. Yuri has only 76 fans while he has 1339 reviews.
       Jeb has 0 fans while he has 57 reviews.
319
320    I used 2 SQL codes for this analysis.
321
322    SELECT name, review_count, fans
323    FROM user
324    ORDER BY fans DESC
325
```

```
326    +-----------+--------------+------+
327    | name      | review_count | fans |
328    +-----------+--------------+------+
329    | Amy       |          609 |  503 |
330    | Mimi      |          968 |  497 |
331    | Harald    |         1153 |  311 |
332    | Gerald    |         2000 |  253 |
333    | Christine |          930 |  173 |
334    | Lisa      |          813 |  159 |
335    | Cat       |          377 |  133 |
336    | William   |         1215 |  126 |
337    | Fran      |          862 |  124 |
338    | Lissa     |          834 |  120 |
339    | Mark      |          861 |  115 |
340    | Tiffany   |          408 |  111 |
341    | bernice   |          255 |  105 |
342    | Roanna    |         1039 |  104 |
343    | Angela    |          694 |  101 |
344    | .Hon      |         1246 |  101 |
345    | Ben       |          307 |   96 |
346    | Linda     |          584 |   89 |
347    | Christina |          842 |   85 |
348    | Jessica   |          220 |   84 |
349    | Greg      |          408 |   81 |
350    | Nieves    |          178 |   80 |
351    | Sui       |          754 |   78 |
352    | Yuri      |         1339 |   76 |
353    | Nicole    |          161 |   73 |
354    +-----------+--------------+------+
355    (Output limit exceeded, 25 of 10000 total rows shown)
356
357
358    SELECT name, review_count, fans
359    FROM user
360    ORDER BY fans ASC
361
362    +---------+--------------+------+
363    | name    | review_count | fans |
364    +---------+--------------+------+
365    | Joe     |            2 |    0 |
366    | Jeb     |           57 |    0 |
367    | Jed     |            8 |    0 |
368    | Rae     |            2 |    0 |
369    | Ryan    |            2 |    0 |
370    | Joe     |            1 |    0 |
371    | Scott   |            7 |    0 |
372    | John    |            3 |    0 |
373    | Ron     |            9 |    0 |
374    | Bryan   |            5 |    0 |
375    | Patti   |            2 |    0 |
376    | Gary    |           23 |    0 |
377    | Kristin |           28 |    0 |
378    | Cynthia |            4 |    0 |
379    | Mrme    |            2 |    0 |
380    | Austin  |            2 |    0 |
381    | Mesut   |           25 |    0 |
382    | Lissa   |            3 |    0 |
383    | Tara    |            3 |    0 |
384    | Lyndsey |            1 |    0 |
385    | Annie   |           11 |    0 |
386    | Daniece |            2 |    0 |
387    | Alex    |            7 |    0 |
388    | Mary    |            2 |    0 |
389    | Garen   |            3 |    0 |
390    +---------+--------------+------+
391    (Output limit exceeded, 25 of 10000 total rows shown)
392
393    9. Are there more reviews with the word "love" or with the word "hate" in them?
394
```

```
395     Answer:
396
397 Yes. There are more rewievs with the word "love" than with the word "hate". There are
    1780 reviews with the word 'love' and 232 reviews with the word 'hate'.
398
399     SQL code used to arrive at answer:
400
401 SELECT COUNT(*)
402 FROM review
403 WHERE text LIKE '%love%'
404
405 SELECT COUNT(*)
406 FROM review
407 WHERE text LIKE '%hate%'
408
409
410 10. Find the top 10 users with the most fans:
411
412     SQL code used to arrive at answer:
413
414 SELECT name, fans
415 FROM user
416 ORDER BY fans DESC
417 LIMIT 10
418
419     Copy and Paste the Result Below:
420
421 +-----------+------+
422 | name      | fans |
423 +-----------+------+
424 | Amy       |  503 |
425 | Mimi      |  497 |
426 | Harald    |  311 |
427 | Gerald    |  253 |
428 | Christine |  173 |
429 | Lisa      |  159 |
430 | Cat       |  133 |
431 | William   |  126 |
432 | Fran      |  124 |
433 | Lissa     |  120 |
434 +-----------+------+
435
436 11. Is there a strong relationship (or correlation) between having a high number of
    fans and being listed as "useful" or "funny?" Out of the top 10 users with the highest
    number of fans, what percent are also listed as "useful" or "funny"?
437
438 Key:
439 0% - 25% - Low relationship
440 26% - 75% - Medium relationship
441 76% - 100% - Strong relationship
442
443     SQL code used to arrive at answer:
444
445 SELECT name, fans, useful, funny
446 FROM user
447 ORDER BY fans DESC
448 LIMIT 10
449
450     Copy and Paste the Result Below:
451
452 +-----------+------+--------+--------+
453 | name      | fans | useful |  funny |
454 +-----------+------+--------+--------+
455 | Amy       |  503 |   3226 |   2554 |
456 | Mimi      |  497 |    257 |    138 |
457 | Harald    |  311 | 122921 | 122419 |
458 | Gerald    |  253 |  17524 |   2324 |
459 | Christine |  173 |   4834 |   6646 |
460 | Lisa      |  159 |     48 |     13 |
```

```
461  | Cat      |  133 |   1062 |    672 |
462  | William  |  126 |   9363 |   9361 |
463  | Fran     |  124 |   9851 |   7606 |
464  | Lissa    |  120 |    455 |    150 |
465  +----------+------+--------+--------+
```

    Please explain your findings and interpretation of the results:

All of the top 10 users with the highest number of fans are also listed as "useful" and "funny". I believe there is a
strong correlation (76% - 100% - Strong relationship) between having a high number of fans and being listed as "useful" or "funny".

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 analyzed Las Vegas and Food on this question. Yes, two groups have different distribution of hours. The food place with 2.5 stars open between 8:00-22:00 on Saturday while the food place with higher rating 4.0 opens late on Saturday.

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

I analyzed Las Vegas and Food on this question. Yes, two groups have different number of reviews. The food place with 2.5 stars have 6 reviews while the food place with higher rating 4.0 have 30 reviews.

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

The food places in two different groups are located in different postal codes.

SQL code used for analysis:

```
SELECT business.name, business.city, category.category, business.stars ,hours.hours,
business.review_count, business.address, business.postal_code
FROM (business INNER JOIN category ON business.id =
category.business_id) INNER JOIN hours ON hours.business_id =
business.id
WHERE business.city = 'Las Vegas' AND category.category = "Food"
GROUP BY business.stars;
```

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 more reviews on average than the businesses that are closed.

ii. Difference 2:have more reviews on average than
    The businesses that are open listed as 'funny' compared to the businesses that are closed.

| AVG(business.stars) | SUM(business.review_count) | AVG(business.review_count) | COUNT(review.funny) | is_open |
|---|---|---|---|---|
| 3.52039473684 | 35261 | 23.1980263158 | 1 | 0 |
| 3.67900943396 | 269300 | 31.7570754717 | 13 | 1 |

```
511   +--------------------+--------------------------+--------------------------+--------
      ------------+---------+

512
513   SQL code used for analysis:
514
515   SELECT
      AVG(business.stars),SUM(business.review_count),AVG(business.review_count),COUNT(review.fu
      nny), business.is_open
516   FROM business
517   LEFT JOIN review
518   ON business.id = review.id
519   GROUP BY business.is_open
520
521   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.
522
523   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:
524
525   i. Indicate the type of analysis you chose to do:
526      Comparison of the average number of reviews and stars of the restaurants among
         different cities.
527
528   ii. Write 1-2 brief paragraphs on the type of data you will need for your analysis and
      why you chose that data:
529
530      Getting a decision to open a restaurant in one of the cities used in the analysis.
         The city which has lowest average number of reviews and stars would be a target
         city to open a restaurant.
531
532   iii. Output of your finished dataset:
533
534      +-----------+------------+----------------------------+--------------------+
535   | city      | category   | AVG(business.review_count) | AVG(business.stars) |
536   +-----------+------------+----------------------------+--------------------+
537   | Charlotte | Restaurants |                        5.5 |               4.25 |
538   | Las Vegas | Restaurants |                      265.5 |              3.875 |
539   | Phoenix   | Restaurants |                126.166666667 |                3.5 |
540   | Tempe     | Restaurants |                        5.0 |                2.5 |
541   | Toronto   | Restaurants |                       29.9 |                3.4 |
542      +-----------+------------+----------------------------+--------------------+
543
544   iv. Provide the SQL code you used to create your final dataset:
545
546   SELECT business.city, category.category, AVG(business.review_count), AVG(business.stars)
547   FROM business
548   LEFT JOIN category
549   ON category.business_id=business.id
550   WHERE business.city IN ('Phoenix', 'Toronto', 'Charlotte', 'Las Vegas') AND
      category='Restaurants'
551   GROUP BY business.city
552
```