

Homework 2

Project Report

Yelp Data Ananlysis

FNU KARAN
NetID: kx361

March 12, 2017

The dataset is freely available at https://www.yelp.com/dataset_challenge/dataset

We just need to fill out a small form and we can access the data. We get a .tar file and we can extract the json files from it. This homework did not require us to use the images dataset. But it made us use the other five datasets, which consisted of reviews, users, business, tips.

Here I did some bery minimal statistics to get some sense out of the data. Mainly it was implemented in FIG. And the observations were plotted using Tableau.

Before performing any of the operations, I made sure that the HPC is working correctly, and flawlessly. Also, the most challenging task was to load the json dataset into Pig variables. This was resolved using the libraries by Twitter, called elephant-bird.

The report contains each question described in sequence. Each section contains the respective questions and the code for each of them and the plots for them.

Following are the operations that were performed:

1 Summarize the number of reviews by US city, by business category

Here is the Pig Script for this problem,

```
1 SET elephantbird.jsonloader.nestedLoad 'true';
2
3 --LOAD FROM THE JSON
4 loadJson = LOAD '/user/kx361/business.json' USING com.twitter.elephantbird.pig
   .load.JsonLoader('-nestedLoad') AS (json:map []);
5
6 --GENERATE THE REQUIRED DATA
7 businessData = FOREACH loadJson GENERATE (int)json#'review_count' as
   review_count, json#'city' as city, json#'categories' as categories;
8
9 --FILTER THE DATA BY REQUIRED CITIES
10 us_business = FILTER businessData BY (
11 (city matches '.*Pittsburgh.*') OR (city matches '.*Charlotte.*') OR (city
   matches '.*Urbana-Champaign.*') OR (city matches '.*Phoenix.*') OR (city
   matches '.*Las Vegas.*') OR (city matches '.*Madison.*') OR (city matches
   '.*Cleveland.*'));
```

```

12
13 --FLATTEN THE DATA OUT
14 flattenedBusinessData = FOREACH us_business GENERATE review_count, city,
    FLATTEN(categories);
15
16 --GROUP THE DATA BY CITY AND CATEGORIES
17 groupedBusinessData = GROUP flattenedBusinessData BY (city, categories);
18
19 --PERFORM THE REQUIRED SUM OPERATION ON REVIEW_COUNT
20 finalData = FOREACH groupedBusinessData GENERATE group.city as city , group.
    categories as category, SUM(flattenedBusinessData.review_count);
21
22 --STORE THE DATA IN THE DIRECTORY
23 store finalData into './p1' using PigStorage(',') ;

```

In the script, elephant-bird was used to load the data in pig variables. The required fields are then extracted from the data with aliases to work on them. Since the data is grouped together, we also flattened out the data to separate out all record to perform the operations that are required.

Once the script is run, and the required output is obtained, then we plot the visualizations in Figure 1,2 and 3.

As we can see from the figures, yelp is mainly being used for restaurant reviews, since it had the maximum number of reviews.

2 Rank all cities by number of stars descending, for each category

Here is the Pig Script for this problem,

```

1 SET elephantbird.jsonloader.nestedLoad 'true';
2
3 --LOAD BUSINESS DATA
4 yelp_business_data = LOAD '/user/kx361/business.json' USING com.twitter.
    elephantbird.pig.load.JsonLoader('-nestedLoad') as (json:map[]);
5
6 --FLATTEN OUT THE CATEGORIES
7 yelp_business_data_category = FOREACH yelp_business_data GENERATE (float)json#
    'stars' AS stars, json#'city' as city, FLATTEN(json#'categories') AS
    categories;
8
9 --GROUP BY CATEGORIES AND CITY
10 yelp_business_group= GROUP yelp_business_data_category BY (categories,city);
11
12 --CALCULATE THE AVERAGE NUMBER OF STARS FOR EACH GROUP
13 yelp_business_group_data= FOREACH yelp_business_group GENERATE group.
    categories as category,group.city AS city, AVG(yelp_business_data_category
    .stars) AS stars;
14
15 --GROUP AND ORDER THIS DATASET BY CATEGORIES
16 yelp_data_order= ORDER yelp_business_group_data BY category ASC;
17 grouped_categories = group yelp_data_order by category;
18
19 --NOW THE EACH CATEGORY CONSISTS CITY AND AVERAGE STARS. SO WE SORT THAT FOR
    EACH SINGLE CATEGORY.

```

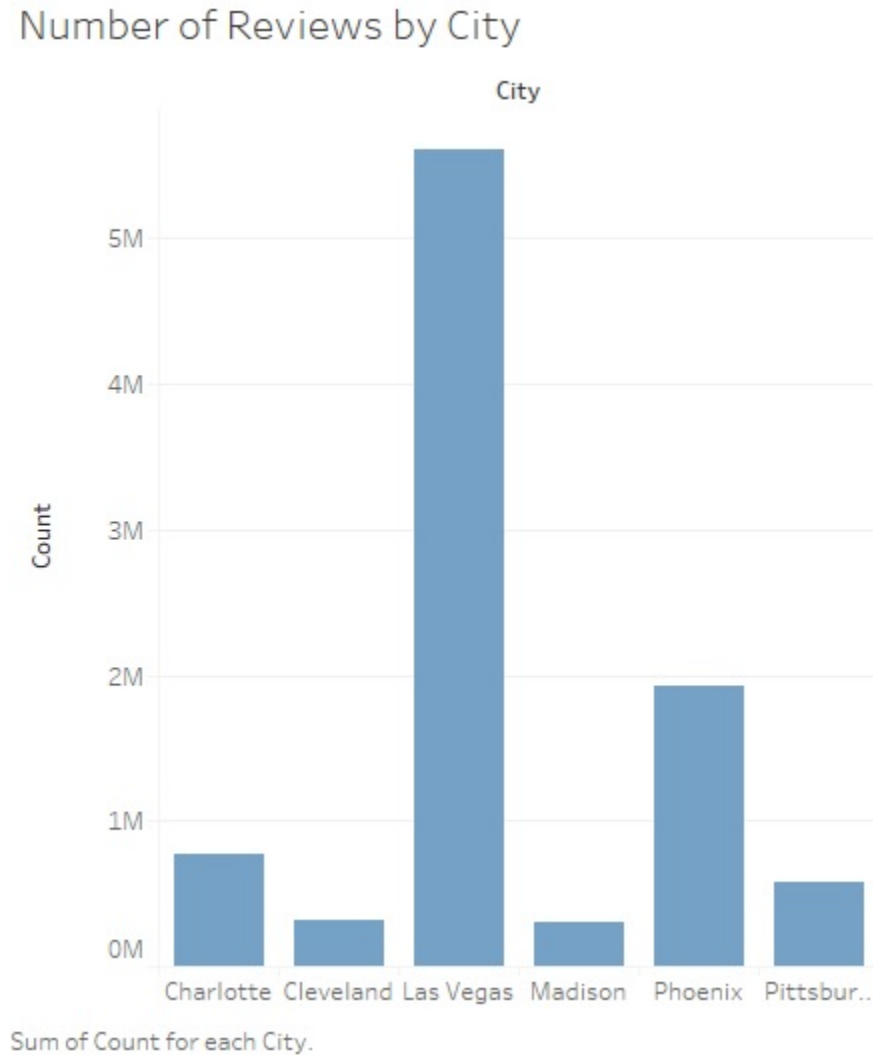


Figure 1: Number of reviews by city.

```

20 sorted_stars_data = FOREACH grouped_categories{
21     sorted= ORDER yelp_data_order by stars desc;
22     GENERATE FLATTEN(sorted);
23 };
24
25 --AND FINALLY STORE THE RESULT.
26 STORE sorted_stars_data INTO './p2' using PigStorage(',') ;

```

The code is self explanatory with all the comments indicating the steps that were taken to work on the problem. Here the catch was to sort the data by two fields. This is indicated in the line number 20-23. What I did, is that I generated a sorted list of stars for each category and then I flattened it so that each of the records is indicated as individual. I did that for each of the category and hence I was able to sort the cities with number of stars for each category.

Once the script is run, and the required output is obtained, then we plot the visualizations in Figure 4,5, and 6

```
> p1 = read.csv('p1.csv')
> summary(p1)
```

City	Category	Number.Of.Reviews
Charlotte : 749	: 6	Min. : 3
Cleveland : 583	Accessories : 6	1st Qu.: 24
Las Vegas :1002	Active Life : 6	Median : 123
Madison : 628	Acupuncture : 6	Mean : 2060
Phoenix : 936	Adult : 6	3rd Qu.: 623
Pittsburgh: 713	Adult Education: 6	Max. :761661
	(Other) :4575	

Figure 2: Summary of the dataset in R.

3 Average number of stars for businesses within 10 miles of the University of Wisconsin, Madison, by business category

Here is the Pig Script for this problem,

```
1 SET elephantbird.jsonloader.nestedLoad 'true';
2
3 --LOAD BUSINESS DATA
4 businessData = LOAD '/user/kx361/business.json' USING com.twitter.elephantbird.
   pig.load.JsonLoader('-nestedLoad') as (json:map[]);
5
6 --GET THE REQUIRED DATA
7 requiredData = FOREACH businessData GENERATE FLATTEN(json#'categories') AS
   type_of_business, (float)json#'stars' AS stars, (float)json#'latitude' as
   latitude, (float)json#'longitude' as longitude;
8
9 --FILTER THE DATA BY LONGITUDE AND LATITUDE
10 wisconsinBusiness = FILTER requiredData BY (latitude>42.9083) AND (
   latitude<43.2417) AND (longitude>-89.5839) AND (longitude<-89.2506);
11
12 -- GROUP THE DATA BY BUSINESS TYPE
13 groupedByCategory = GROUP wisconsinBusiness BY type_of_business;
14
15 -- FOR EACH GROUP, GENERATE THE AVG NUMBER OF STARS
16 finalData = FOREACH groupedByCategory GENERATE group as categories, AVG(
   wisconsinBusiness.stars) AS stars;
17
18 --FINALLY SORT THE DATA BY CATEGORIES/TYPE
19 orderedData = ORDER finalData BY categories;
20
21 --STORE IT AS CSV
22 STORE orderedData INTO './p3' using PigStorage(',');
```

The code is self explanatory with all the comments indicating the steps that were taken to work on the problem.

Once the script is run, and the required output is obtained, then we plot the visualizations in Figure 7 and 8.

Number of reviews by Categories

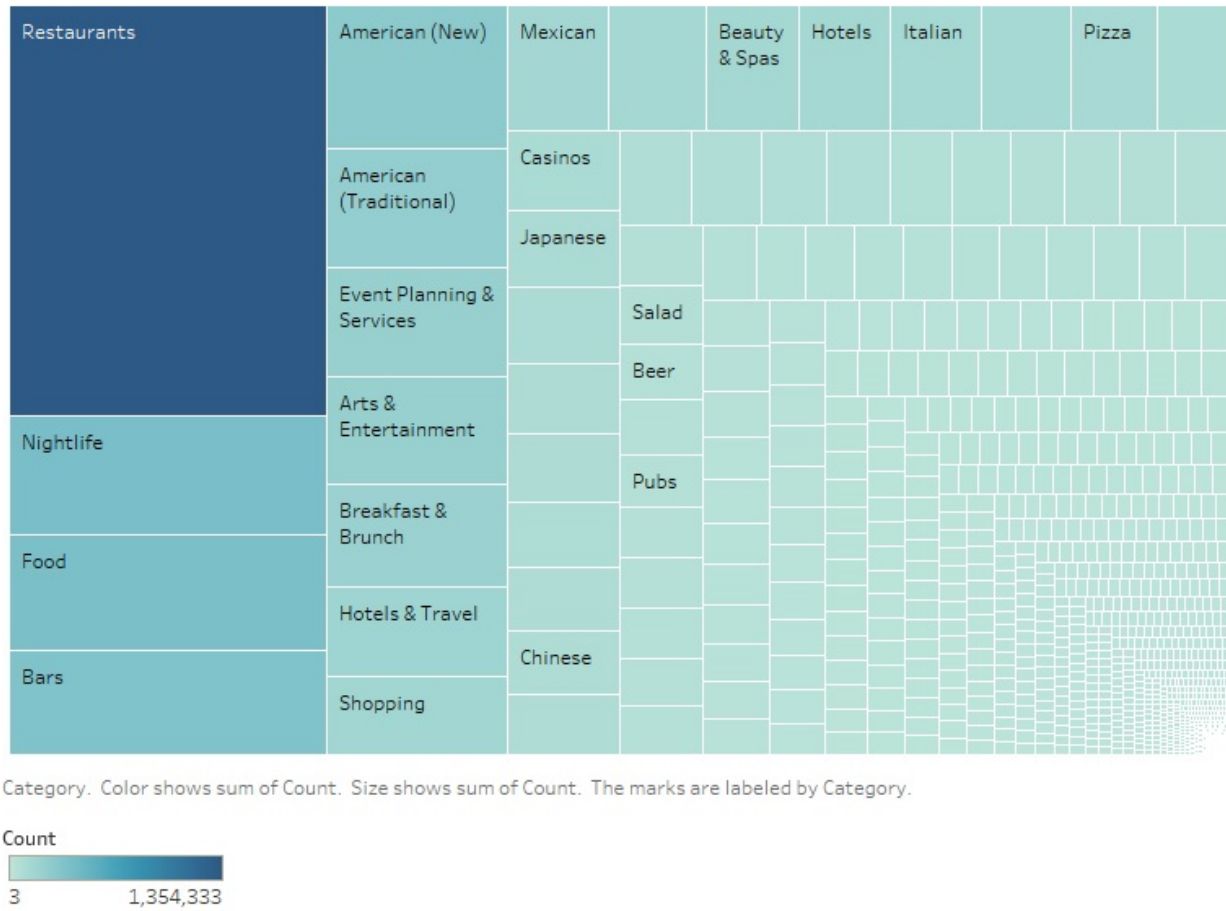


Figure 3: Number of reviews by categories.

4 Average number of stars for the top 10 reviewers

Here is the Pig Script for this problem,

```

1 SET elephantbird.jsonloader.nestedLoad 'true';
2
3 --Fetch business data
4 businessData = LOAD './business.json' USING com.twitter.elephantbird.pig.load.
    JsonLoader('-nestedLoad=true') AS (yelp_business: map[]);
5 business = FOREACH businessData GENERATE yelp_business#'business_id' as
    business_id, yelp_business#'categories' as categories;
6
7 --Fetch review data
8 reviewData = LOAD './yelp_academic_dataset_review.json' USING com.twitter.
    elephantbird.pig.load.JsonLoader('-nestedLoad=true') AS (yelp_review: map
    []);
9 review = FOREACH reviewData GENERATE yelp_review#'user_id' as user_id1,
    yelp_review#'business_id' as business_id, (int)yelp_review#'stars' as

```

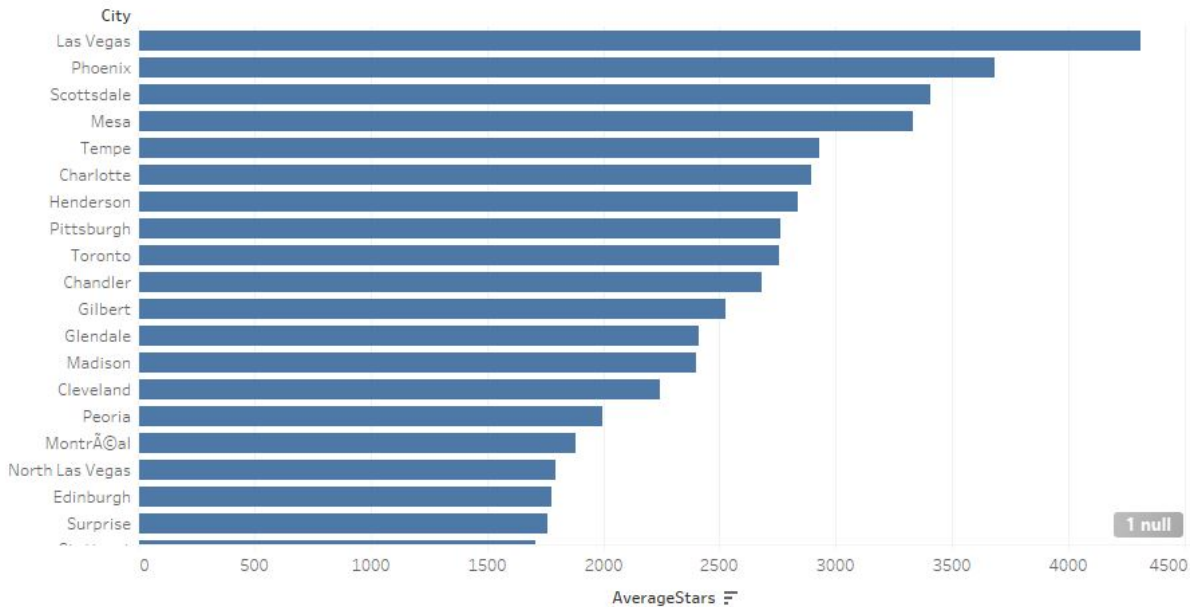


Figure 4: Total of Average stars for each category by city.

```

stars;
10
11 --Fetch user data
12 userData = LOAD './yelp_academic_dataset_user.json' USING com.twitter.
    elephantbird.pig.load.JsonLoader('-nestedLoad=true') AS (yelp_user: map[])
    ;
13 user = FOREACH userData GENERATE yelp_user#'user_id' as user_id, (int)
    yelp_user#'review_count' as review_count;
14
15 --Sort the order in descending order
16 sortedUsers = ORDER user BY review_count DESC;
17
18 --Get the top 10
19 top10Users = LIMIT sortedUsers 10;
20
21 --Join the users with reviews on by user id
22 userReviewJoined = JOIN top10Users BY user_id, review BY user_id1;
23
24 --Join the joined part with businesses
25 userReviewBusinessJoined = JOIN business BY business_id, userReviewJoined BY
    business_id;
26
27 --Flatten the data
28 flattenedData = FOREACH userReviewBusinessJoined GENERATE FLATTEN(categories),
    stars, user_id;
29
30 --group the data by users and also categories
31 grouped = GROUP flattenedData BY (user_id, categories);
32
33 --calculate the average and sort it

```

```

> p2 = read.csv('p2.csv')
> summary(p2)
      Category      City      Average.Stars
Restaurants: 674   Las Vegas : 1002      4      : 6383
Food         : 480   Phoenix   : 942     3.5     : 5567
Shopping     : 394   Scottsdale: 845     4.5     : 4785
Nightlife    : 391   Mesa      : 759     5       : 4735
Pizza        : 380   Charlotte : 749     3       : 3787
Bars         : 379   Toronto  : 742     2.5     : 2169
(other)      :45007  (other)  :42666   (other):20279
> |

```

Figure 5: Summary of the dataset in R.

```

34 grouped_result = FOREACH grouped GENERATE group, AVG(flattenedData.stars) as
    avg_stars;
35 ordered_result = ORDER grouped_result BY user_id;
36
37 --store the result into the directory
38 STORE grouped_result INTO './p4' using PigStorage(',');

```

The code is self explanatory with all the comments indicating the steps that were taken to work on the problem.

Once the script is run, and the required output is obtained, then we plot the visualizations in Figure 9, 10 and 11.

5 Summarize star rating for reviews in January through May For the top 10 and bottom 10 food business near UWM (in terms of stars)

Here is the Pig Script for this problem,

```

1 SET elephantbird.jsonloader.nestedLoad 'true';
2
3 --
4 business_data = LOAD './business.json' USING com.twitter.elephantbird.pig.load
    .JsonLoader('-nestedLoad') AS (json:map []);
5
6 required_business_data = FOREACH business_data GENERATE json#'business_id' as
    business_id, json#'city' as city, json#'name' as name, (double) json#'
    latitude' as latitude, (float) json#'longitude' as longitude, (double)
    json#'stars' as stars, FLATTEN(json#'categories') as categories;
7
8 filtered_data = FILTER required_business_data BY (latitude <= 43.2467) AND (
    latitude >= 42.90833) AND (longitude >= -89.58389) AND (longitude <=
    -89.25056);
9
10 filtered_food_business = FILTER filtered_data by categories matches '.*Food.*'
    ;
11

```

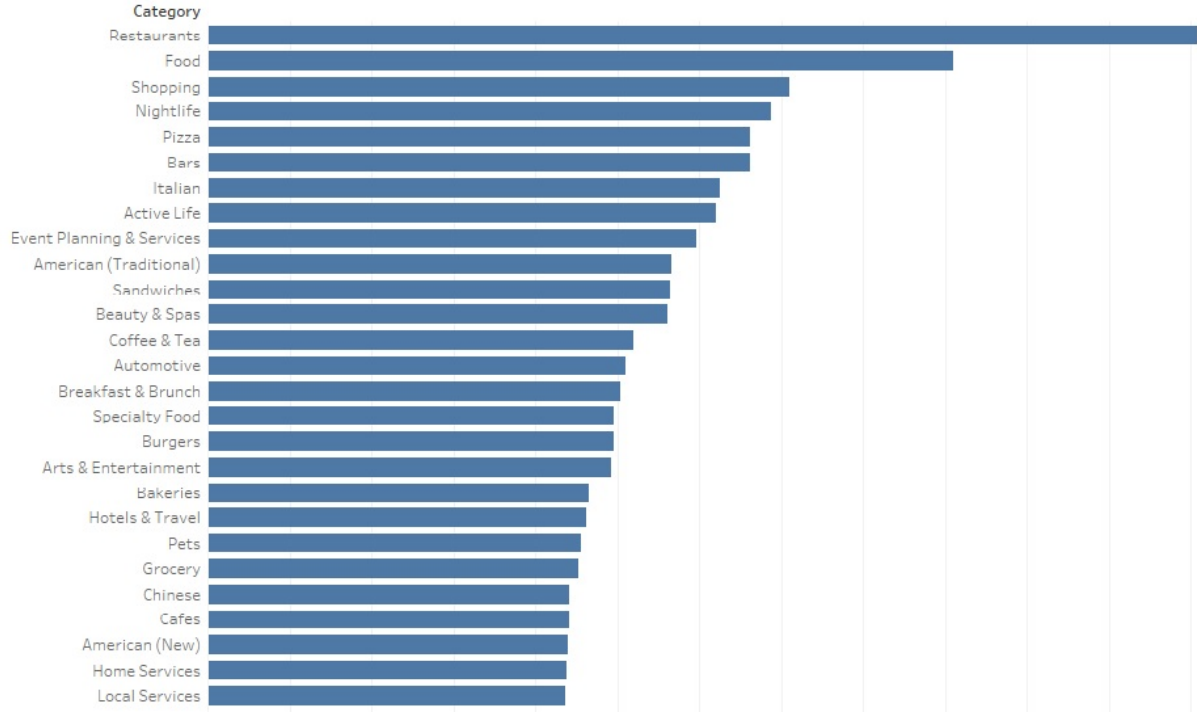


Figure 6: Total of average stars for each city by category.

```

12 business_id_stars = FOREACH filtered_food_business GENERATE business_id, stars
    ;
13
14 ordered_data = ORDER business_id_stars by stars DESC;
15
16 top10 = LIMIT ordered_data 10;
17
18 Ordered_By_Stars_Asc = ORDER business_id_stars by stars;
19
20 bottom10 = LIMIT Ordered_By_Stars_Asc 10;
21
22 top10bottom10 = UNION top10, bottom10;
23
24 review_data = LOAD './yelp_academic_dataset_review.json' USING com.twitter.
    elephantbird.pig.load.JsonLoader('-nestedLoad') AS (json:map []);
25
26 required_review_data = FOREACH review_data GENERATE json#'review_id' as
    review_id, json#'business_id' as business_id, json#'date' as date, json#'
    stars' as stars;
27
28 joined_reviews = JOIN required_review_data by business_id, top10bottom10 by
    business_id;
29
30 final_required_Data = FOREACH joined_reviews GENERATE top10bottom10::
    business_id as bid, (double) required_review_data::stars as star,
    required_review_data::review_id, SUBSTRING(required_review_data::date,5,7)
    as month;

```

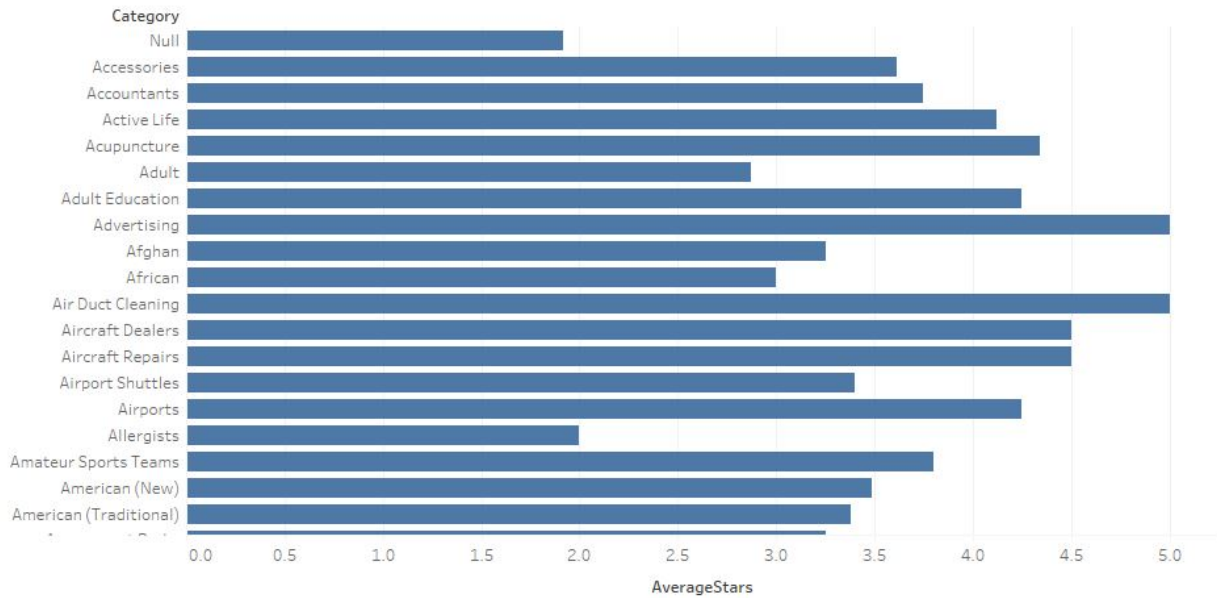



Figure 7: Average number of stars for the categories of businesses near University of Wisconsin.

```
> p3 = read.csv('p3.csv')
> summary(p3)
      Category      AverageStars
      : 1      Min.      :1.000
Accessories: 1      1st Qu.:3.500
Accountants: 1      Median :3.833
Active Life: 1      Mean    :3.817
Acupuncture: 1      3rd Qu.:4.250
Adult       : 1      Max.    :5.000
(other)     :656
> |
```

Figure 8: Summary of the dataset in R.

```
31
32 filtered_data_by_month = FILTER final_required_Data BY (month matches '
33 01|02|03|04|05');
34
35 grouped_data_by_business = GROUP filtered_data_by_month by bid;
36
37 avg_rating = FOREACH grouped_data_by_business GENERATE group, AVG(
38   filtered_data_by_month.star) as avg_stars;
39
40 STORE avg_rating into './p5' using PigStorage ('','');
```

The code is self explanatory with all the comments indicating the steps that were taken to work on the problem.

Once the script is run, and the required output is obtained, then we plot the visualizations in Figure 12 and 13.

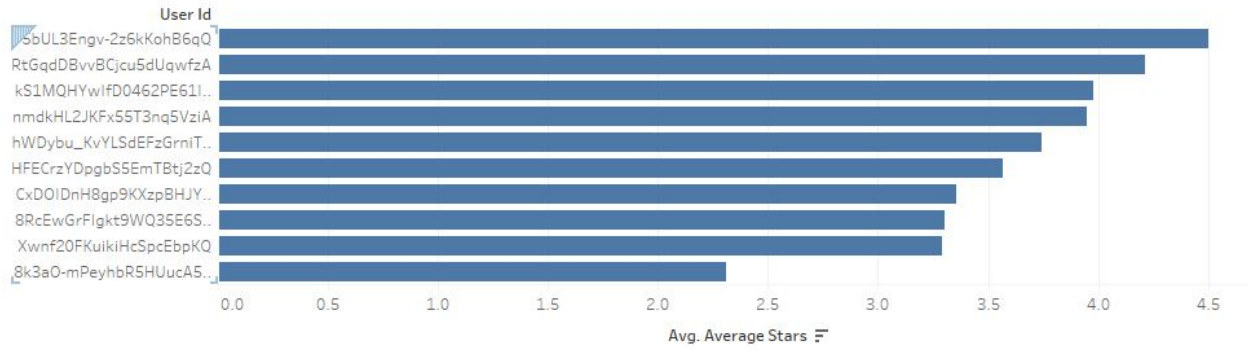


Figure 9: Average number of stars for every user in Top 10.

This might not look much, but it required a good understanding of PIG, and R. Understanding the data and plotting various plots could not have been done without the understanding of the data. Since, the data was huge, so these manipulations could not have been done any local database like SQL, Oracle or for that matter MongoDB.

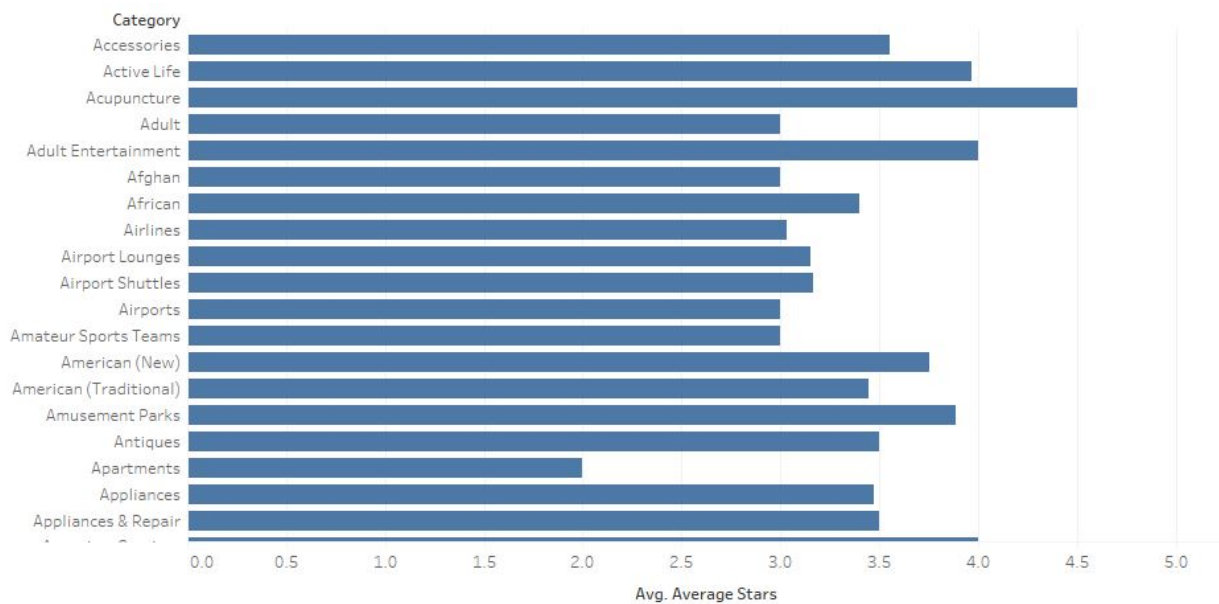


Figure 10: Average number of stars for each category by these top 10 reviewers.

```
> summary(p4)
```

UserId	Category	AverageStars
CxDOIDnH8gp9KXzpBHJYXw:385	Arts & Entertainment	10
hwDybu_KvYLSdEFzGrniTw:347	American (New)	9
Xwnf20FKuikiHCSpCEbpKQ:154	Bakeries	9
nmdkHL2JKFx55T3nq5VziA:114	Bars	9
HFECrZYDpgbS5EmTBtj2zQ: 95	Event Planning & services	9
8RCEWGrFIgkt9wQ35E6SnQ: 74	Food	9
(other)	(other)	1255

```
> |
```

Figure 11: Summary of the dataset in R.

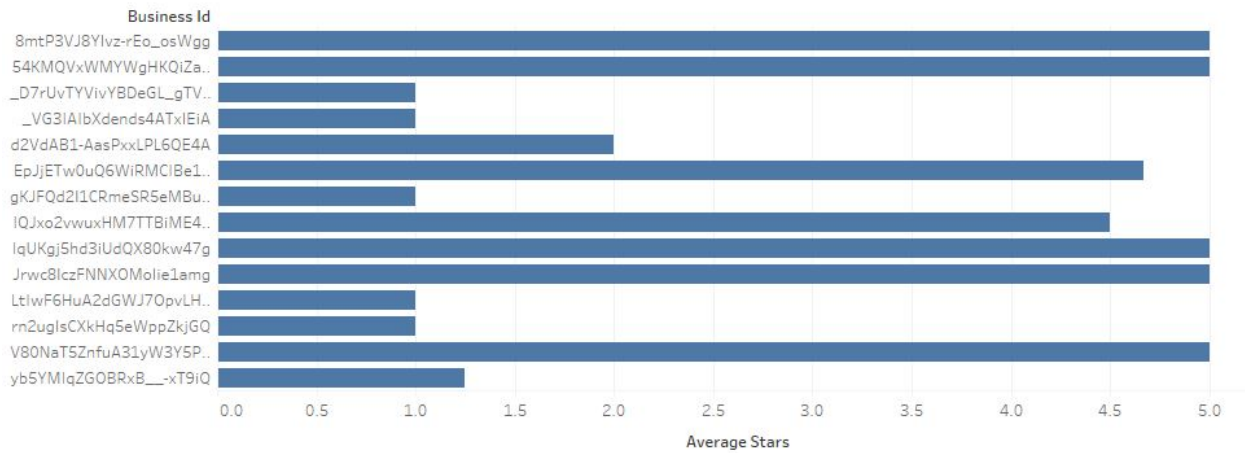


Figure 12: Average number of stars for top 10 and bottom 10 businesses.

```
> p5 = read.csv('p5.csv')
> summary(p5)
```

BusinessId	AverageStars
_D7rUvTYVivYBDeGL_gTVQ:1	Min. :1.00
_VG3IAIbXdends4ATxIEIA:1	1st Qu.:1.00
54KMQVxWMYWgHKQiza_58Q:1	Median :3.25
8mtP3VJ8Ylvz-rEo_oswgg:1	Mean :3.03
d2VdAB1-AasPxxLPL6QE4A:1	3rd Qu.:5.00
EpJjETw0uQ6WiRMCIBe10A:1	Max. :5.00
(Other)	:8

```
> |
```

Figure 13: Summary of the dataset in R.