# OSEMN Process for Working on Extracted Fields from the Yelp Reviews Data Set

## Rahul Kumar Nalubandhu and Sandra Estrada

## Introduction

Conducting a review analysis gives businesses the opportunity to improve customer experience, identify service gaps and gain real time insights amongst many other benefits. This analysis dives into customer reviews entered in Yelp for businesses in different states throughout the United States.

## Problem

The issue with having little to no customer reviews is that this can negatively impact sales. Customer reviews are an important channel to attract customers and increase sales. A benefit in analyzing reviews at the business, city and state level will provide insight into which states and/or businesses have the least customer engagement and allow for proper intervention.
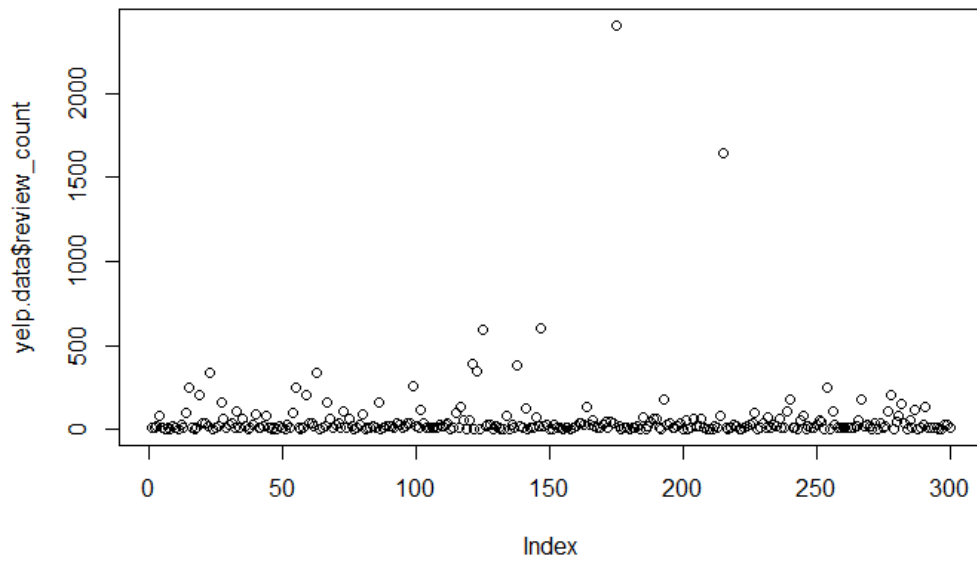
## Obtain Data

Yelp is a one-stop platform which enables customers to connect with businesses. More than 80 million people visit this platform in a month to find businesses and service providers. Customers are given the ability to leave reviews and request quotes from local businesses amongst many other things. In return, local business owners are given the ability to communicate with their customers and respond to reviews to build trust with their customers. The customer review data set is acquired directly through Yelp. The data set is 4.04GB (1 point) and split into multiple Json files (2 points) which contain businesses, reviews, and user data. In addition, the data has punctuation (1 point) and has more than one type of related data (2 points). Based on the point system requirements provided, the yelp data is a 6-point data set.

## Scrub Data

Prior to conducting the analysis for the yelp data set, the data is extracted and consolidated into a csv and Json file for the following fields: business name, city, state, and review count. Visual studio is the primary application utilized to read and extract data in the python programming language. In addition, all null values are deleted from the data set. Null values are deleted to not compromise the integrity of the analysis.
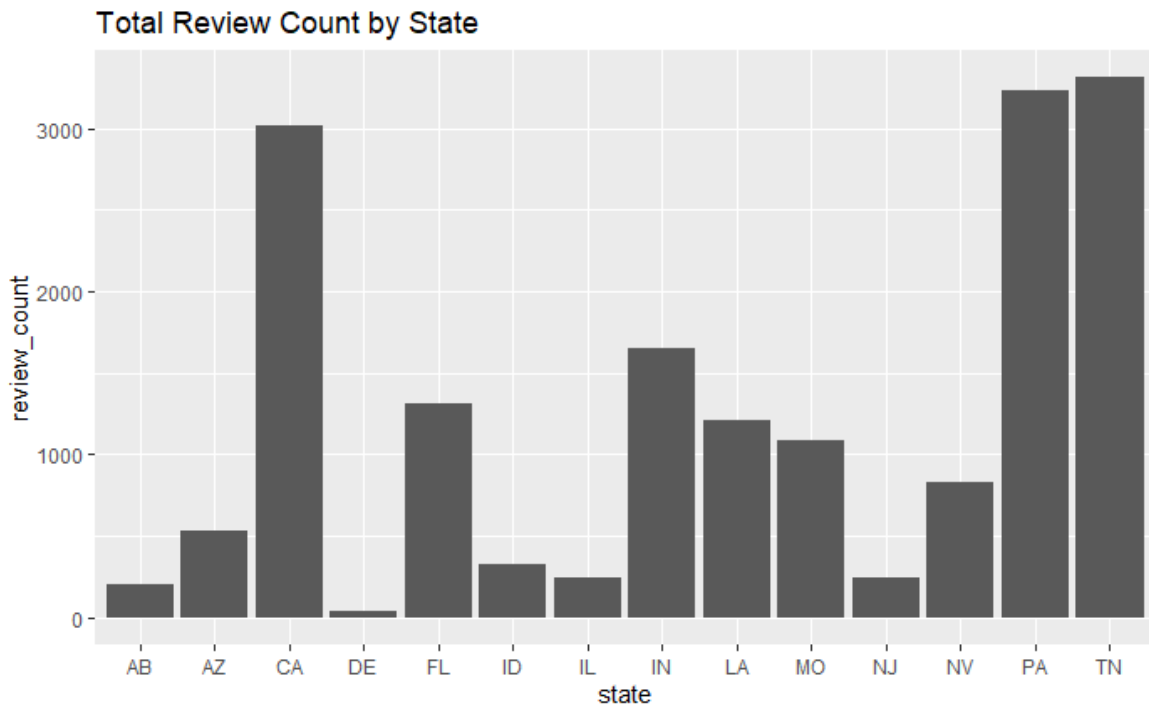
## Explore Data

For the exploration stage, the csv data is loaded into R-studio. To begin the exploration stage the review count field is plotted. There are two evident outliers for businesses having more than 1,000 reviews on yelp. Most of the businesses have a review count of less than 500 reviews.

*Graph 1*

The review count variable is also examined at the state level for the total and average number of reviews. These observations indicate the states with the most reviews are Pennsylvania and Tennessee followed by California. On the opposite end of the spectrum, Delaware is the state with the least total and average number of reviews. On average, the state whose businesses have more reviews compared to other states is California.



*Graph 2*

Graph 3

**Model Data**

After reviewing the dataset, a decision tree is the most suitable modeling technique for this analysis. A decision tree will contain the different factors that determine if businesses are likely to face less customer reviews due to the state they are in or even the demographics of their target customers. While this is deemed the best approach for this analysis, more data needs to be gathered and analyzed to properly complete this.

**Interpret Data**

The business with the most reviews in this data set is Santa Barbara Shellfish Company located in CA with 2,404 reviews. For the companies with the lowest count of reviews, there are a few with only 5 reviews. An interesting observation is that most businesses are in FL. This observation is aligned with the results from the Average Review Count by State plot.

Description: df [259 × 4]

| state <chr> | name <chr> | total_count <int> | review_count <int> |
|---|---|---|---|
| CA | Santa Barbara Shellfish Company | 1 | 2404 |
| TN | Gaylord Opryland Resort & Convention Center | 1 | 1639 |
| NV | Romano's Macaroni Grill | 2 | 678 |
| MO | Budweiser Brewery Experience | 1 | 605 |
| TN | Mike's Ice Cream | 1 | 593 |
| PA | Tuna Bar | 2 | 490 |
| PA | BAP | 2 | 410 |
| CA | Helena Avenue Bakery | 1 | 389 |
| LA | Mahony's Po-Boys & Seafood | 1 | 382 |
| LA | Copper Vine | 1 | 350 |

| | state<br><chr> | name<br><chr> | total_count<br><int> | review_count<br><int> |
|---|---|---|---|---|
| 1 | AB | River City Games | 1 | 5 |
| 2 | AZ | Ballistic Fabrication | 1 | 5 |
| 3 | AZ | Desert Design Center | 1 | 5 |
| 4 | FL | 1-275 Rest Area Manatee County Mile 7 | 1 | 5 |
| 5 | FL | Bay Area Appliance | 1 | 5 |
| 6 | FL | PDQ Temple Terrace | 1 | 5 |
| 7 | FL | Thach Used Tires | 1 | 5 |
| 8 | FL | Zesty Tsunami | 1 | 5 |
| 9 | IL | All In Shipping | 1 | 5 |
| 10 | IL | K-9 Groom Room | 1 | 5 |

1-10 of 10 rows

## Tasks Completed per Team Member

Rahul Kumar Nalubandhu – Read, extracted data from the original data set, removed null values, converted the Json data into a csv format file, and assisted in portion of the r analysis.

Sandra Estrada – Converted the csv format into a Json format file, completed the r analysis, and wrote the OSEMN report.

## Appendix

Initial Data:

Python code for conversions:

```python
6    import csv
7    import json
8    import pandas as pd
9    from pathlib import Path
10
11   #read the dataset
12   with open("yelp_academic_dataset_business.json", "r", encoding='utf-8') as infile:
13       data = infile.readlines() #read all lines
14       json_data = []
15       for line in data :
16           line_data = json.loads(line)
17           # handles all Null or empty values. It doest read the whole row
18           if any(val == "None" or val == "" for val in line_data.values()):
19               continue
20           json_data.append(line_data)
21
22   json_dataframe = pd.DataFrame.from_records(json_data[:300])
23   json_dataframe_new = json_dataframe[['name', 'state', 'city', 'review_count']]
24
25   with open("format_json.json", "w") as outfile:
26       outfile.write('[')
27       for i, row in json_dataframe_new.iterrows():
28           json.dump(row.to_dict(), outfile)
29           if i < len(json_dataframe_new) - 1:
30               outfile.write(',')
31               outfile.write('\n')
32       outfile.write(']')
33
34   #convert from json to csv
35   json_dataframe_new.to_csv('format.csv',index = False)
36
37   #convert from csv to json
38   with open("format.csv", "r", encoding='utf-8') as csvFile:
39       csv_read = csv.DictReader(csvFile)
40       conv_json_data = {}
41       for line, rows in enumerate(csv_read, start=1):
42           conv_json_data.update({"Business {:02}".format(line):rows})
43       with open("format.json", "w", encoding='utf-8') as jsonFile:
44           json.dump(conv_json_data, jsonFile, indent=4)
45
```

CSV Sample:

| name | state | city | review_count |
|---|---|---|---|
| Abby Rappoport, LAC, CMQ | CA | Santa Barbara | 7 |
| The UPS Store | MO | Affton | 15 |
| Target | AZ | Tucson | 22 |
| St Honore Pastries | PA | Philadelphia | 80 |
| Perkiomen Valley Brewery | PA | Green Lane | 13 |
| Sonic Drive-In | TN | Ashland City | 6 |
| Famous Footwear | MO | Brentwood | 13 |
| Temple Beth-El | FL | St. Petersburg | 5 |
| Tsevi's Pub And Grill | MO | Affton | 19 |
| Sonic Drive-In | TN | Nashville | 10 |
| Marshalls | FL | Land O' Lakes | 6 |
| Denny's | IN | Indianapolis | 28 |
| Adams Dental | FL | Clearwater | 10 |
| Zio's Italian Market | FL | Largo | 100 |
| Tuna Bar | PA | Philadelphia | 245 |
| Arizona Truck Outfitters | AZ | Tucson | 10 |
| Herb Import Co | LA | New Orleans | 5 |
| Nifty Car Rental | LA | Kenner | 14 |
| BAP | PA | Philadelphia | 205 |
| Roast Coffeehouse and Wine Bar | AB | Edmonton | 40 |
| Barnes & Noble Booksellers | IN | Indianapolis | 38 |
| Hibachi Express | IN | Indianapolis | 20 |
| Romano's Macaroni Grill | NV | Reno | 339 |
| Super Dog | TN | Nashville | 6 |
| Indian Walk Veterinary Center | PA | Newtown | 15 |
| H&M | CA | Santa Barbara | 24 |
| The Green Pheasant | TN | Nashville | 161 |

format (+)

Json Sample :

```json
C: > Users > estra > OneDrive > Desktop > Python > CS512Module3 > {} format.json > ...
1    {
2        "Business 01": {
3            "name": "Abby Rappoport, LAC, CMQ",
4            "state": "CA",
5            "city": "Santa Barbara",
6            "review_count": "7"
7        },
8        "Business 02": {
9            "name": "The UPS Store",
10           "state": "MO",
11           "city": "Affton",
12           "review_count": "15"
13       },
14       "Business 03": {
15           "name": "Target",
16           "state": "AZ",
17           "city": "Tucson",
18           "review_count": "22"
19       },
20       "Business 04": {
21           "name": "St Honore Pastries",
22           "state": "PA",
23           "city": "Philadelphia",
24           "review_count": "80"
25       },
26       "Business 05": {
27           "name": "Perkiomen Valley Brewery",
28           "state": "PA",
29           "city": "Green Lane",
30           "review_count": "13"
31       },
32       "Business 06": {
33           "name": "Sonic Drive-In",
34           "state": "TN",
35           "city": "Ashland City",
36           "review_count": "6"
37       },
```

R-Studio code:

```r
1   ---
2   title: "datawrangling"
3   output: pdf_document
4   date: "2023-01-28"
5   ---
6
7   ```{r setup, include=FALSE}
8   knitr::opts_chunk$set(echo = TRUE)
9
10  library(ggplot2)
11  library(magrittr)
12  library(dplyr)
13  library(tidyverse)
14  library(data.table)
15  ```
16
17
18  ```{r}
19  #load data
20  yelp.data <- read.csv("format.csv", head = T)
21
22  #print data - only prints 10 first and 10 last rows
23  head(yelp.data, n=20)
24  tail(yelp.data, n=20)
25
26  #plot for review count
27  plot(yelp.data$review_count) #a couple outliets can be see in the plot
28
29  ```
30
31  ```{r}
32  #total review count by state barchart
33  sum.plot <- ggplot(yelp.data, aes(x = state, y = review_count)) + stat_summary(fun = sum, geom = "bar")
34
35  print(sum.plot + ggtitle("Total Review Count by State"))
36
37  ```
38
39  ```{r}
40  #average review count by state barchart
41  avg.plot <- ggplot(yelp.data, aes(x = state, y = review_count)) + stat_summary(fun.y = mean, geom = "bar")
42
43  print(avg.plot + ggtitle("Average Review Count by State"))
44  ```
45
46  ```{r}
47  sum.plot.businessname <- ggplot(yelp.data, aes(x = name, y = review_count)) + stat_summary(fun = sum, geom =
    "bar")
48
49  print(sum.plot.businessname + ggtitle("Total Review Count by Business"))
50  ```
51
52
53  ```{r}
54  #Count for business name in data set
55  table(yelp.data$name)
56  ```
57
58  ```{r}
59  #Count of business name in each state
60  yelp.data.table <- yelp.data %>% group_by(state, name) %>%
61    summarise(total_count=n(),.groups = 'drop') %>%
62    as.data.frame()
63
64  yelp.data.table
65  ```
66
67
68  ```{r}
69  yelp.data.table3 <- yelp.data %>% group_by(state, name) %>%  select(review_count) %>%
    summarise(total_count=n(),review_count=sum(review_count), .groups = 'drop') %>% arrange(review_count) %>%
    as.data.frame()
70  yelp.data.table3
71
72
73  ```
74
75  ```{r}
76  head(yelp.data.table3, n=10)|
77
78  tail(yelp.data.table3, n=10)
79  ```
80
```

# References

e-satisfaction. (n.d.). *7 reasons why customer reviews are important*. Retrieved January 27, 2023, from https://www.e-satisfaction.com/7-reasons-why-customer-reviews-are-important/#:~:text=Analyzing%20reviews%20left%20by%20your,what%20your%20customers%20truly%20want

*Convert CSV to JSON using python with new headers*. Stack Overflow. (1969, July 1). Retrieved January 27, 2023, from https://stackoverflow.com/questions/73089933/convert-csv-to-json-using-python-with-new-headers

Rana, D. S. (2022, September 12). *Learn Review Analysis: Why, how, data sources with Free Trial*. Learn Review Analysis | Why, How, Data Sources with Free Trial. Retrieved January 29, 2023, from https://www.repustate.com/blog/review-analysis/#:~:text=Review%20analysis%20allows%20you%20to,new%20sales%20opportunities%2C%20and%20more.

*Indexing and selecting data - pandas 1.5.3 documentation*. (n.d.). Retrieved January 29, 2023, from https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy.

Yelp for Business. (2023, January 23). *Plan, start, grow, and advertise your small business*. Yelp for Business. Retrieved January 29, 2023, from https://business.yelp.com/