

Yelp New York Restaurants Reviews Analysis

Pragati Gupta
Master of science, Information technology
Montclair State University
New Jersey, USA
guptap5@montclair.edu

Abstract— Online reviews posted by the consumers are a great source of information on different websites for the given product or destination. In this digital world, people prefer to check the rating and comments provided by other users instead of obeying the advertisement.

The Online market is a multi-billion market, with Covid-19, this will only increase. Hence the value of reviews is significantly more than ever. Reviews are not only essential for the consumer but play a central role in the functioning and sales of the business. Companies are investing more toward analyzing the reviews because they don't want to mislead the users.

The challenge is how to identify analyze the reviews from the pile of comments and ratings posted on daily bases. This paper focus on the Yelp reviews and rating to understand the reviews opinion and compare the highly rated restaurant vs the low rated restaurant in New York.

Keywords—New York, reviews, machine learning, Online reviews, Yelp, Analysis, Rating.

I. INTRODUCTION

Yelp is the renowned website for registering the business and to find the desirable destination. It could be a restaurant or a place to visit. Yelp allows the user to check the rating and comments of restaurants and consumers can share their experience to provide the feedback. These rating and reviews can be voted by the user as “funny”, “cool” and “useful” which gives more clarity to identify the popular reviews. As the number of customers increasing day by day the number of reviews is becoming the huge set of data.

The fake review can damage the sales of a particular business. The reason can be a competition, revenge, scam; which can affect the owner's reputation. The fake reviews can be paid to attract more customers and the fake reviews also can be a part of self-promotion. Fake reviews mislead the customers and affect the business.

This research paper aims to analyze the reviews from the numerous available comments and ratings.

This paper presents the investigation on reviews based on the number of reviews and rating of the restaurant as well as the comment posted by the customers on the restaurant page. The following research questions are addressed:

RQ1 According to the published research papers, what are the major challenges/issues/risks/methods/ solutions proposed to analyze the reviews?

RQ2 What are the factors involve in investigating the reviews?

RQ3 The comparative analysis of high rated and low rated restaurant on New York area?

II. DATA SOURCE

RQ1 According to the published research papers, what are the major challenges/issues/risks/methods/ solutions proposed to analyze the reviews?

In previous research paper many studies have conducted based on the reviewer details such as location, number of reviews posted, valid IP address, number of friends, followers, active hours. The reviews including stars and the comment posted by the consumer. For example, the reviews are matching with the star rating, the length of the reviews, context information, bag of words tagging the review as positive and negative (sentimental analysis), extremely positive and extremely negative reviews, duplicate reviews, tense used in the comments, most impactful reviews analysis, circle of the reviewers as well as business owners like friends and families.

Some model helps to filter the reviews based on their “Helpful” and “not helpful”. The customer votes the review as helpful or not helpful [7].

Crowd intelligence powered dynamic machine learning framework - RREF (Reliable Review Evaluation Framework), Which follows two approaches. A situation-aware task allocation approach for crowdsourcing quality control, For example worker's track records, workload. Other is RREF enhances the AdaBoost framework for classification accuracy like selection and re-weighting of topic classification models.[2].

The yelp voting system is label as “funny”, “cool” and “useful”. the collected data from the voting system shows that the positive reviews are voted as cool, the negative reviews get funny votes. Early listed business on yelp voted as useful. The longer reviews valued as useful, cool and funny [1].

The quantitatively investigation methods on reviews shows in [6] taking in consideration Rhetorical level including Spontaneous structures and Tenses; Thematic level including

Experiential testimonies, Informative discourse, Sensorial description, Purchase experience, Risk-Reward; Numerosity, Place, Detailed general text, Endorsements and recommendations. Enunciative level including Enunciate interpellation, Spontaneous- effects, Verisimilitude effects.

One research paper propose the OneReview method based on the comparison of Yelp and TripAdvisor reviews to detect the fraud and suspicious reviews OneReview focuses on the change point analysis method where evaluate the reviews from different website for the same business, detect those that do not match across the websites and identify them as suspicious and crowd-labelling on the review of every business independently on every website. Data is collected from multiple sources for the business those match on different websites. For each business, the change point analyzer identifies suspicious reviews depending on time frame. Then its Textual (TF-IDF and Sentiments) and contextual features are extracted the details. Contextual feature focus on Review-based features, Business-based features, User-based features. Then The classifier analyses suspicious reviews to detect fraudulent reviews by creating a ground truth dataset, building the K cross validation model, Classify suspicious reviews into fraudulent and genuine reviews. Spammers' data is used to detect fraud campaigns like Campaigns by a set of spammers, Campaigns over certain businesses, Socially networked campaigns.[4]

To understand the user behavior on online book reviews, [] research paper uses the LDA method by extracting common topic in the comments. The Topic terms categorize in eleven categories. This covers the part of DC core elements (Dublin Core) metadata, which consists of 15 core elements.) as well as MARC fields ((Machine-Readable Cataloguing record)) and distinct categories concerning emotion and evaluation. Then considering the reviews factors from the descriptive metadata. Which helps to discover and selected the children's books. [3].

A. Different methods for analysis

i. Information Extraction:

Generally, the first step to do any kind of user opinion analysis is information extraction. It is the process of extracting information from unstructured textual sources to enable finding entities as well as classifying and storing them in a database. Here we are looking for specific information from textual sources. One of the most trivial examples is when you email extract only data from the message for you to add in your calendar. In our use case we can use this method to extract specific information from the user opinions on different platform like positives, negatives, quality etc. about any business or app. The main focus while performing this method is the Accuracy and Usefulness of the extracted data.

ii. Classification:

Classification is a task that is performed with machine learning algorithms that learn how to assign a class label to

examples from the problem domain. It is also considered as predictive modeling problem where a class label is predicted for given input data. Some of the predefined categories which we figure out with our research are:

Containing noise and Fake information

- New feature
- User intention
- Installation problem in case of apps
- Price/Quality/Privacy

This method covers accuracy, efficiency, informative and usefulness of the data collected for reviews.

iii. Clustering:

Clustering is the task of organizing the reviews, rating and comments in to groups called a cluster. The members of the cluster are more similar to each other than to those in other groups. As opposed to Classification this does not required any prior knowledge of the topics. In this we try to find out the similar opinions. For example, opinions talking about the same dish in the restaurants, reviews explaining the hospitality of the restaurants. It also includes similar opinion about the theme or interior. This method helps in determining Accuracy and Usefulness of the reviews.

iv. Search and Information Retrieval:

Search and Information retrieval is the method of finding and tracing user opinions that matches the needed information. This task can be used to find out variety of information or trace the reviews complaining about certain feature, food item or trace the restaurant best reviewed items. This task can also be used to detect some of the key words common across the opinions and can be flag as suspicious or fake. This method helps in accuracy and retrieval of the data.

v. Sentimental Analysis:

This refers to the task of interpreting user emotions in the reviews. It detects the sentiments polarity of the review (e.g., positive, neutral or negative) whether it is a whole review, a sentence or a phrase. This task is the best approach to understand the trend about the place or business. It helps in determining the areas of improvement and also tells whether the reviews associated with them are not repeatable and genuine. Once you understand the preference of users then that can be highlighted for marketing. Also, for our use case of opinion analyses this can use to determine the trend based on the sentiments and understand the business or app if the performance is matching with the sentiments of the reviews.

Sentiment analysis is a type of natural language processing for tracking the mood of the public about a particular product or topic. Sentiment analysis, which is also called opinion mining, involves in building a system to collect and examine opinions about the product made in blog posts, comments, reviews or tweets. [4]

vi. Content Analysis:

The content analysis is study of data to find certain words, themes, concepts and pattern in the data. The content is then analyzed to characterize and quantify the presence and

connection between the certain words, reviews and patterns. This help in find out the relation between the rating and reviews length where they have significance. The content analysis can also help the business to find any recurring issues mentioned in the reviews which can be identified and taken care.

vii. Recommendation:

Recommendation method is designed to provide the course of action that the business should follow. The above approaches are designed in such a way that they recommend some opinions to be highlighted depending on the need. We can also find the reviews which have influence on the customers so that we can find the factors affecting the feedback.

viii. Summarization:

User opinion summarization aims to provide a concise and precise summary of the one or more opinions. In this task first we group the reviews together using Clustering technique and then provide a proposal summarizing thousands of reviews in that group. Based on the summary business can work on the improvements or suggested based on the summary.

ix. Visualization:

Visualization can help businesses to identifying pattern and trends which making them easier to interpret information from the data extracted from the reviews. We can use different type of graphical representation like chart, tables. Simple statistics can help us in finding the reviews which may be paid or intentionally spreading negativity or positivity.

B. Mining Techniques

i. Manual Analysis:

This technique is quite famous among the scholars to develop a data set or training records to do perform different kind of opinion analysis. For example, using manual analysis we can find different words, concepts and pattern in the data to group them together and do Content Analysis. Manual analysis typically takes a form of tagging a group of sample reviews with one or more meaningful tags. For instance, tags might indicate type of complaints or sentiment users expresses towards the business or specific product or item getting traction. There are below steps involved in the manual analysis:

- a) Formulate analysis objective
- b) Select reviews for analysis
- c) Specify unit of analysis
- d) Perform coding process
- e) Analyze dataset or use of evaluation

Manual analysis is time consuming and requires lot of human effort.

ii. Natural Language Processing:

The interaction between human and computer. Program a computer to process and analyze the huge amount of data; where data is the information available in any natural language for example English. NLP is one of the branches of Artificial Intelligence. AI makes a computer perform the task which a human can do. NLP shows how to deal with the text data. In NLP the scrape the data from the website, put data into standard formats for future analysis using text pre-processing techniques, organize the cleaned data standard text format. When dealing with numerical data, data cleaning often involves removing null values and duplicate data, dealing with outliers, etc. With text data, there are some common data cleaning techniques, which are also known as text pre-processing techniques. the text must be tokenized, meaning broken down into smaller pieces. The most common tokenization technique is to break down text into words. It can do this using scikit-learn's Count Vectorizer, where every row will represent a different document and every column will represent a different word. EDA is a visualizing method to get the trend to data. Sentimental analysis and Topic modeling are the powerful techniques to text analysis. NLP techniques has Text similarity techniques, Pattern matching techniques, Collocation finding techniques.

iii. Machine Learning:

Field of study that gives computers the ability to learn without being explicitly programmed (Arthur Samuel). Machine learning can be assigned to one of two broad classifications one is Supervised learning and other is Unsupervised learning.

Supervised learning has given a data set and already know what our correct output should look like, having the idea that there is a relationship between the input and the output. Supervised learning is categorized into "regression" and "classification". In a Regression, it is trying to predict results within a continuous output, it trying to map input variables to some continuous function. In a classification, it is trying to predict results in a discrete output. In other words, we are trying to map input variables into discrete categories.

Unsupervised learning: The Other Machine learning category is Unsupervised learning which approach the problem with little or no idea what result will look like. It uses clustering the data based on relationships among the variables in the data.

In a machine learning based classification, two sets of documents are required: training and a test set. A training set is used by an automatic classifier to learn the differentiating characteristics of documents, and a test set is used to validate the performance of the automatic classifier.[5]

A number of machine learning techniques have been adopted to classify the reviews. Machine learning techniques like Naive Bayes (NB), maximum entropy (ME), and support

vector machines (SVM). The other most well-known machine learning methods in the natural language processing area are K-Nearest neighbourhood, ID3, C5, centroid classifier, winnow classifier, and the N-gram model.

Naive Bayes is a simple but effective classification algorithm. The Naive Bayes algorithm is widely used algorithm for document classification. The basic idea is to estimate the probabilities of categories given a test document by using the joint probabilities of words and categories [5].

Support vector machines (SVM), seeks a decision surface to separate the training data points into two classes and makes decisions based on the support vectors that are selected as the only effective elements in the training set [5].

Centroid classification algorithm, s, document is assigned to the class corresponding to the most similar centroid.

The K-nearest neighbor (KNN), the system finds the k nearest neighbors among training documents. The similarity score of each nearest neighbor document to the test document is used as the weight of the classes of the neighbor document.

Maximum entropy, it creates a model that best accounts for the available data but with a constraint that without any additional information the model should maximize entropy. In other words, the model prefers a uniform distribution by maximizing the conditional entropy.[6]

III. RELATED WORKS

Type	Machine Learning Techniques	No. Studies
Supervised	Naïve Bayes	35
	Support Vector Machine	31
	Decision Tree	26
	Logistic Regression	20
	Random Forest	12
	Neural Network	9
	Linear Regression	5
	K-Nearest Neighbor	4
Unsupervised	Latent Dirichlet Allocation	30
	K-Means	4

Table 1: Machine Learning Techniques

Early efforts have investigated several sequences of machine learning techniques and characteristics to analyzing reviews. They used regression models, heterogeneous graphs, unsupervised anomaly detection, and behavioral models. Some works observed the time of correlations between the time of writing reviews, and the location of their authors, bursts in the number of reviews to locate suspicious patterns, correlations between the time of writing reviews and the location of their authors, length of the reviews, sentiment attached with the reviews etc. [4]

This Research paper analysis the rating, number of reviews of highly rated restaurant and low rated restaurant in New York. One main concern of review analysis or the spam or fake review detection is absence of ground truth. The duplicate reviews and reviews with unnecessary details can we detected based on the other user's response on 'funny', 'cool', 'useful' votes of that reviews. On Yelp, funny, cool, and useful votes are good measures of quality [1]. As each review has his own experience and put that experience in his own words so it is difficult to put all the reviews in one category.

IV. DATA

Yelp's website is Yelp.com which is a crowdsourced local business review and social networking site. In this paper we focused on the restaurant listed in New York location on yelp website. We collected 1152 restaurant data and randomly choose 10 restaurants. five restaurants have highly rated, other five have low rating. As per Wikipedia 78% of businesses listed on the site had a rating of three stars or better, but some negative reviews were very personal or extreme, sosome of the reviews are written in an entertaining or creative manner. We use filters to select the all kind of high and low rated restaurants for better analysis.

V. PROPOSED APPROACH

Natural language processing (NLP) is an exciting branch of artificial intelligence (AI) that allows machines to break down and understand human language. NLP techniques, to interpret text data for analysis. We use text pre-processing techniques, machine learning techniques and Python libraries for NLP. Text pre-processing techniques include tokenization, text normalization and data cleaning. Once in a standard format, various machine learning techniques can be applied to better understand the data. This includes using popular modeling techniques to classify spam, or the sentiment of reviews.

Newer, more complex techniques can also be used such as topic modeling, word embeddings or text generation with deep learning.

This paper analyze all reviews step by step approach to understand the comments, using several NLP libraries in Python including NLTK, TextBlob, spaCy and gensim along with the



Figure 1: Process of reviews analysis

standard machine learning libraries including pandas and scikit-learn.

Our work focuses on analyzing product reviews of high rates and low rated restaurants to analysis the differences.

Our goal is to look at reviews of various restaurants and note their similarities and differences. The steps are web scraping data, cleaning data, organizing data. Getting the data - scraping data from a website. The output will be clean, organized data in two standard text formats:

- Corpus - a collection of text
- Document-Term Matrix - word counts in matrix format

Cleaning the data - Text pre-processing techniques.

Organizing the data - Organize the cleaned data into a way that is easy to input into other algorithms

A) Getting The Data

The data is collected from the Yelp website. Data contain review details of 1152 restaurants of New York. Five random high rated restaurant (rated 4.5 to 5 stars) and five random low rated restaurants (rated 1 to 1.5 stars) selected for the further analysis. For web scraping request and BeautifulSoup used to get data from the web. The Data cleaning is performed to covert the data into the standard formats (Corpus, Document Term Matrix) for the analysis. Corpus is a collection of texts. The Panda's data frame is use to create the corpus format with the restaurant name and the reviews. Panda and Data Frame are the python libraries for data analysis. Data frame is an object in panda. Which looks like a table. Every row of a data frame has an ID and every column of the data frame ha the same data type

For document Term Matrix first data gets clean (remove excess, unnecessary text) then data have to tokenized (break the text into smaller parts) and then all data puts into the matrix format.

B)Cleaning Data

First covert the all data into the lower case, remove the punctuations, digits, special symbols. Regular expression from python library is used for the data cleaning. Regular expression (Re) is the tool to search for patterns.

Tokenization: Split the text into smaller pieces for instance the size of a word or sentence. The combination of two word called bigram. For this analysis data is tokenized into words where each word is a item. After tokenization stop words are removed from the data. Stop words are the word with the less meaning like the, a. After that the remaining structured data is called bag of words models. Document Term Matrix's each row is a different restaurant, each Colum is a different term (word, but is can also be bi-gram) and the value inside the matrix is the word counts. Count Vectorizer is use to make the DTM. Now the data is in the structed form to perform the analysis. This Analysis follow the MVP (minimum viable product) approach - start simple and iterate.

C) Exploratory Data Analysis

After the data cleaning step, the next step is to take a look at the data explore. When working with numerical data, some of the exploratory data analysis (EDA) techniques include finding the average of the data set, the distribution of the data, the most common values, etc. The idea is the same when working with text data.

In this analysis EDA is use to find patterns before identifying the hidden patterns with machines learning (ML) techniques. This paper focus on the following for each restaurant:

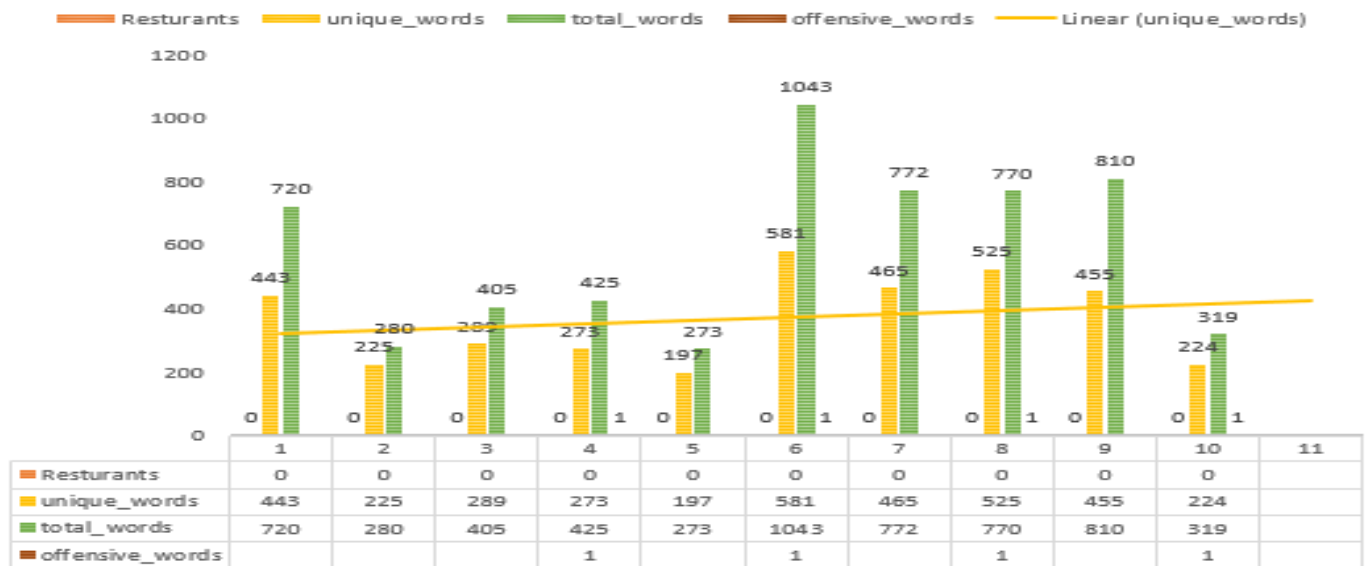


Figure2: Exploratory Data Analysis

- Most common words - find these and create word clouds
- Size of vocabulary - look number of unique words in the comments
- Number of words - use for appreciation and deprecation in the comments and most common terms.

TextBlob Module: Linguistic researchers have labeled the sentiment of words based on their domain expertise. Sentiment of words can vary based on where it is in a sentence. The TextBlob module allows us to take advantage of these labels.

Sentiment Labels: Each word in a corpus is labeled in terms of polarity and subjectivity. A corpus' sentiment is the average of these.

Polarity: How positive or negative a word is. -1 is very negative. +1 is very positive.

Subjectivity: How subjective, or opinionated a word is. 0 is fact. +1 is very much an opinion.

The sentiment performed on the various reviews of the restaurants, both overall and throughout the time line. The result of sentimental analyse s is plot using matplotlib.pyplot.

The data is converted in the list of ten elements, one for each restaurant and each review has been split into ten pieces of text. Calculate the polarity for each piece of text. Plot the result for each restaurant.

D) Topic Modeling

Another popular text analysis technique is called topic modeling. The ultimate goal of topic modeling is to find various

topics that are present in your corpus. Each document in the corpus will be made up of at least one topic, if not multiple topics. Latent Dirichlet Allocation (LDA), which is one of many topic modeling techniques. It was specifically designed for text data.

To use a topic modeling technique, provide

- A document-term matrix and
- The number of topics you would like the algorithm to pick up.

Once the topic modeling technique is applied, keep changing the number of topics. the terms in the document-term matrix, model parameters, or even try a different model, till the results the mix of words in each topic make sense. The genism module for LDA is used. Further noun and adjective filter are used to create a new document-term matrix using only nouns and then adjective then both together. After LDA the topic which makes more sense are pulled down and run some more iterations to get more fine-tuned topics.

E) Text Generation

Markov chains can be used for very basic text generation. Think about every word in a corpus as a state. By making a simple assumption that the next word is only dependent on the previous work - which is the basic assumption of a Markov chain.

Build a Markov Chain Function, A dictionary created form Markov chain function. The keys should be all of the words in the corpus. The values should be a list of the words that follow the keys. Create a Text Generator, Create a function that generates sentences. It will take two things as inputs:

- The dictionary you just created
- The number of words you want generated

VI. RESULTS AND CONCLUSION

- The highly-rated restaurants have more reviews posted by a consumer, but the lower-rated restaurant has fewer reviews. It shows the low-rated restaurants also have good reviews. i.e. all these restaurants serve tasty food and have good services too.
- Among 10 restaurants the amount of profanity words (socially offensive use of language) are not commonly used (rarely used).

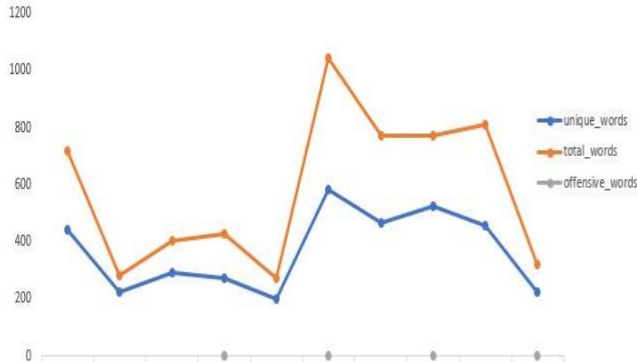


Figure3: Profanity words analysis

- The vocabulary of the reviews is depending on the number of reviews. i.e., vocabulary of unique world used by the reviews depends on the numbers of reviews that particular restaurant gets.
- Sentimental Analysis (Opinion Mining) it is observed that the low rated restaurants are have the negative polarity more than the high rated restaurants i.e. In the reviews analysis we get to know that the customer emotions are positive in writing comment for high rated restaurants.

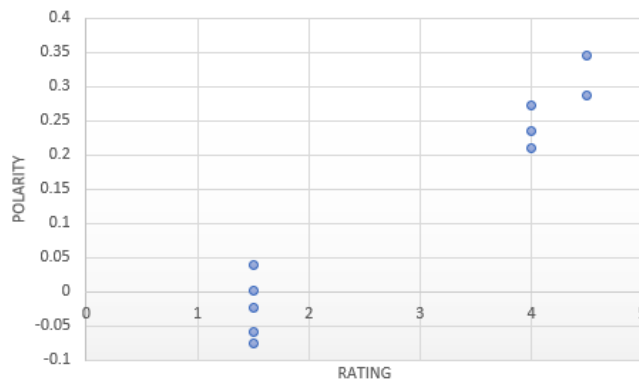


Figure4: Polarity Analysis

- Subjective sentences generally refer to personal opinion, emotion or judgment whereas objective refers to factual information. Range is from 0 to +1, where 0 is fact. +1 is very much an opinion. An objective sentence presents some factual

information about the world, while a subjective sentence expresses some personal feelings or beliefs. We observed that all the reviews of listed restaurant are towards opinion (subjectivity). All the comments are meaningful.

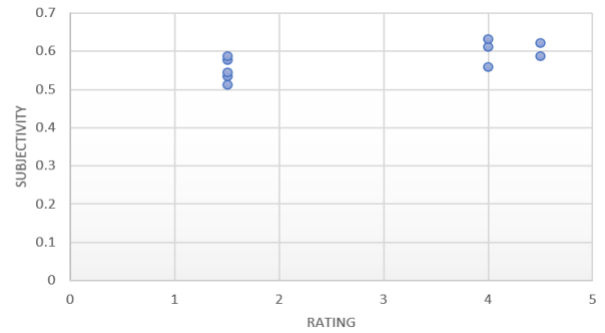


Figure5: Subjectivity Analysis

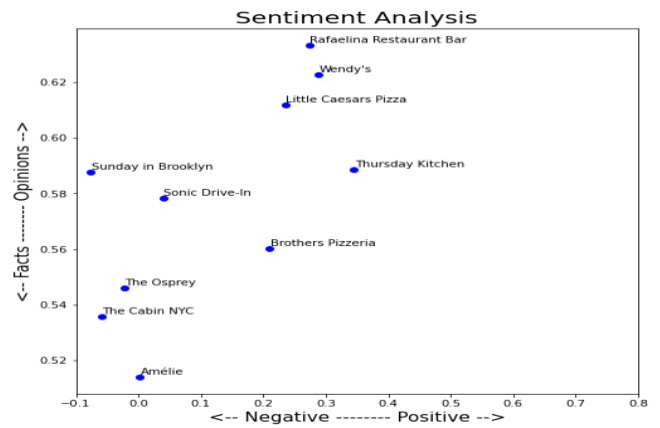


Figure6: Sentiment Analysis

- Then we created a list which holds all the text and have 10 elements for each transcript. After that we calculated the polarity of each text to visualized the change in polarity over the time. It is observed that some of the high rated as well as some low rated restaurants have positive to negative and then negative to positive variation over the time. That show that each restaurant reviews get polarity changes over the time. /seasons.

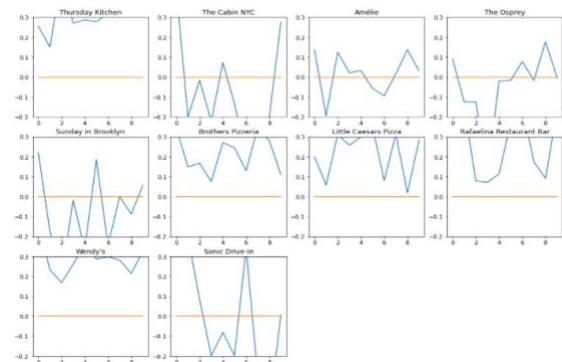


Figure7: Polarity Analysis over the time

- Text Modeling: Using LDA method, the transcript divides in to four topics which are snacks, Interior, food, brunch. This topic is generated from topic modeling by using customer reviews. This shows highly rated restaurant comments are about food, interior and snacks, but low rated restaurant comments topics are brunch, interior, and snacks. The finding is that most of the topics of comments in review section are same.

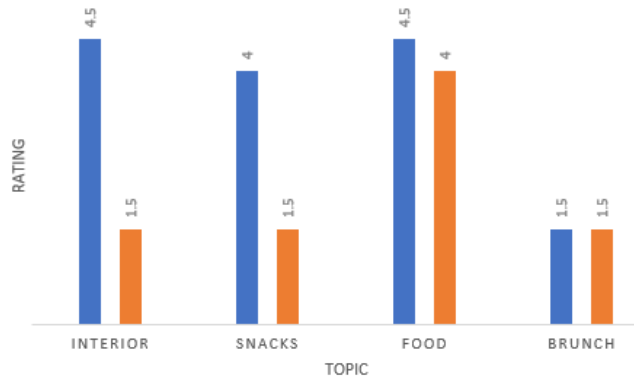


Figure8: Topic Generation

- Text Generation: We built the Markov Chain Function, created the dictionary of comments for one highly rated restaurant and one low rated restaurant. By using the random library, we predicted the next comment for high rated restaurants: Thursday Kitchen the predicted review is 'Rich and had more filling since this was spicy Korean fusion tapas, I've had anything.' This comment shows positive sentiments. The predicted comment for Low rated restaurant Little Caesars Pizza came up as 'African American are raising their own prices store is run by an.' This comment shows the negative sentiments.

VII. FUTURE WORK

Fake reviews start with the reviewer so in future we can analyse the reviewer comments on the restaurants; then filter the comment based on the locations. if a reviewer is posting comments from distinct locations. we can move it as suspicious.

Nowadays the bar-code check-in and order are a must in some of the restaurants. With the metadata. we can write a code to

recheck the comment posted by the person who uses the bar code method are legitimate or not by doing simple supervised machine learning, text, and a bag of words analysis. The "verification of visit" will be same as "verification of purchase" feature used by Amazon.

We can also check the number of reviews posted by a reviewer of a restaurant. For this, we can set a value as a set point. Reviewers who posted reviews more than the set points will come under the suspicious category then we can investigate them.

Acknowledgment

We thank the anonymous reviewers for their valuable feedback, and professor John Jenq, who provide their very helpful comments on this research.

References

- [1] Saeideh Bakhshi, Partha Kanuparth, Partha Kanuparth, "Understanding Online Reviews: Funny, Cool or Useful?" CSCW 2015, March 14-18, 2015, Vancouver, BC, Canada.
- [2] Xinpeng Min, Yuliang Shi, Lizhen Cui, Han Yu, Yuan Miao, "Efficient crowd-powered active learning for reliable review evaluation," ICCSE'17, July 6-9, 2017, Beijing, China.
- [3] Yunseon Choi, Soohyung Joo, "Topic Detection of Online Book Reviews, Topic Detection of Online Book Reviews," 2019 ACM/IEEE Joint Conference on Digital Libraries (JCDL).
- [4] Shirin Nilizadeh, Hojjat Aghakhani, Eric Gustafson, Christopher Kruegel, Giovanni Vigna, "Think outside the dataset: Finding fraudulent reviews using cross-dataset analysis," WWW '19, May 13-17, 2019, San Francisco, CA, USA.
- [5] Michael Crawford, Taghi M. Khoshgoftar, Joseph D Prusa, Aaron N. Richter and Hamzah Al Najada, "Survey of review spam detection using machine learning techniques," Crawford et al. Journal of Big Data (2015) 2:23.
- [6] Paula Almiron-Chamadoira, "Online reviews as a Genre," DTUC '18, October 3-5, 2018, Paris, France.
- [7] Richong Zhang and Thomas Tran, "An entropy-based model for discovering the usefulness of online product reviews," 2008 IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology.
- [8] Petr Hajek, Aliaksandr Barushka, "A comparative study of machine learning methods for detection of fake online consumer reviews," ICEBI '2019, November 9-11, 2019, Prague, Czech Republic.
- [9] I. S. Jacobs and C. P. Bean, "Fine particles, thin films and exchange anisotropy," in Magnetism, vol. III, G. T. Rado and H. Suhl, Eds. New York: Academic, 1963, pp. 271-350.