# Sentiment Analysis on Consumer Reviews of Amazon Products

Aamir Rashid and Ching-yu Huang

*Abstract*—In today's world, the significance of online shopping is growing day by day. The business ideas have been refashioned and completely transformed by making it so easy for the customers to purchase anything they want at just one click of a mouse button. It is becoming even more popular due to its high level of convenience. The only thing customers must have been the Internet and the appropriate method of payment. Amazon.com is one such widely known E-commerce website and it is being used worldwide. It was initially known for its huge collection of books but later it was expanded to sell electronics and other home appliances and consumer products. At present, Amazon is known to sell millions of products. This growth of E-commerce gave importance to customer needs and opinions which in turn gave rise to an important aspect of online shopping known as 'User Reviews'. User reviews are customer suggestions and opinions about the product which helps other customers make decisions about that product. Such review systems form the backbone of E-commerce. The goal of this project is to understand and analyze the Amazon User Review Dataset with the help of different visualization techniques. These visualization techniques will help showcase various informative statistical trends which will provide us with insights about the Amazon Review system. These insights will help in exploring the possible improvements that can be done to satisfy the customers. Major work will involve empirical analysis for data understanding and exploration by taking into consideration, the various metrics related to the user reviews as opposed to sentimental analysis on the review text which aims at understanding the overall emotion of the reviews which has been done previously.

*Index Terms*—Sentiment analysis, chi-square, tableau, dataset, raw data, correlation, hypothesis, business intelligence.

## I. INTRODUCTION

Recent years have seen a tremendous amount of research in the area of sentimental analysis and efforts expended in the understanding sentiment of the customers [1], [2]. This is because the marketplace for consumer products moves to the Internet, the shopping experience changes in a way that makes much of the information regarding the use of products available online and generated by users [3]. This is different than the traditional marketing practices and the way that product information used to be disseminated: through word of mouth and advertising. In this work, we focus on the amazon product and shed some light on customer behavior. Since its creation as an online bookstore in 1994, Amazon.com has grown rapidly and

been a microcosm for user-supplied reviews [4], [5]. Soon, Amazon opened its reviews to consumers and eventually allowed any user to post a review for anyone of the millions of products on the site. With this increase in anonymous user-generated content, efforts must be made to understand the information in the correct context and develop methods to determine the intent of the author.

Understanding what online users think of its content can help a company market its product as well as manage its online reputation. The purpose of this study is to investigate a small part of this large problem: positive and negative attitudes towards products. Sentiment analysis attempts to determine which features of the text are indicative of its context (positive, negative, objective, subjective, etc.) and build systems to take advantage of these features. The problem of classifying text as positive or negative is not the whole problem in and of itself, but it offers a simple enough premise to build upon further. Much of the work involved in sentiment analysis in content containing personal opinions has been done relatively recently [6], [7]. Pang and Lee [8] used several machine learning systems to classify

A recently enormous amount of research published on sentiment analysis has been increasing for the past years. One of the subtopics of this research is called sentiment analysis or opinion mining, which is, given a bunch of text, we can computationally study people's opinions, appraisals, attitudes, and emotions toward entities, individuals, issues, events, topics, and their attributes. Applications of this technique are diverse. and emotions of existing users before they use a service or purchase a product. Last but not least, researchers [9] uses this information to do an in-depth analysis of market trends and consumer opinions, which could potentially lead to a better prediction of the stock market. However, saying this, to find and monitor opinion sites on the Web and distill the information contained in them remains a formidable task because of the proliferation of diverse sites. Each site typically contains a huge volume of opinionated text that is not always easily deciphered in long forum postings and blogs. The average human reader will have difficulty identifying relevant sites and accurately summarizing the information and opinions contained in them [10]. Besides, to instruct a computer to recognize sarcasm is indeed a complex and challenging task given that at the moment, the computer still cannot think like human beings.

The objective of this paper is to classify the positive and negative reviews of the customers over different electronic products using a correlation between rating for the product and the hopeful numbers and also evaluate the null hypothesis for different categories within the Amazon electronics product. We used Tableau Software [11] in

conjunction with Excel to compute P-value that determines whether our null hypothesis has some creditability. We found out that Amazon data that we used for our analysis was largely biased towards higher rating and null hypothesis in almost all sub-categories in the electronics product miserably failed, and we rather chose the alternative hypothesis.

## II. DATASET

To understand the structure and schema of the dataset, let us look at a sample Amazon Review. As seen in Fig. 1, an Amazon User Review consists of four important aspects:
• Summary: The title of the review
• Review text: The actual content of the review.
• Rating: User rating of the product on a scale of 1 to 5.
• Helpfulness: The number of people who found the review useful. These aspects will help us understand and analyze the reviews to derive insights.
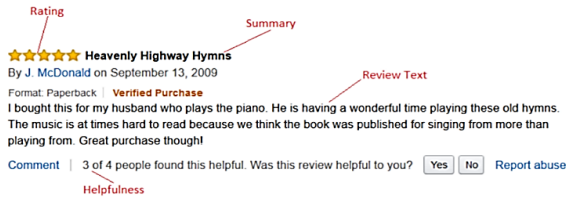


Fig. 1. Actual amazon customer review sample.

### A. Raw Data

Our raw data was downloaded from datafinti's website, which has a total of 3500 records. There are more than 20 columns such as product id, category, sub-category, review title, review text, reviewer name, etc. raw dataset has a lot of null values also data and price format were not correct, there were also outliers and noise in the dataset. Fig. 2 shows a few columns from the Original amazon consumer product review dataset. We can see that in a few columns Values are null and in the wrong format.

### B. Processed Data

Raw data were trimmed using excel because there were more than 20 columns, but we trim it to 10 columns for our analysis [12]. We also format a few data and price columns with proper format using built-in excel functions. We use a built-in excel function to handle outliers. Null and duplicates were removed using tableau filter, include and exclude functions. Fig. 3 shows the normalized data that we used for analysis. It has fewer columns than the then original dataset.
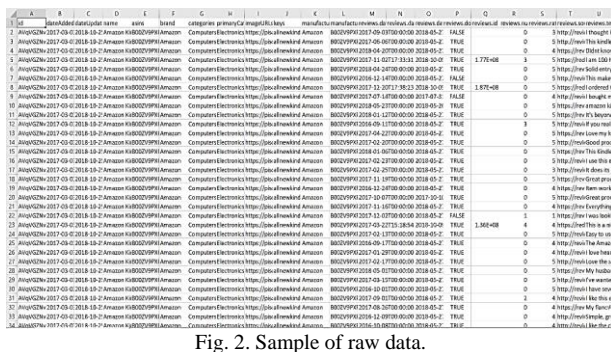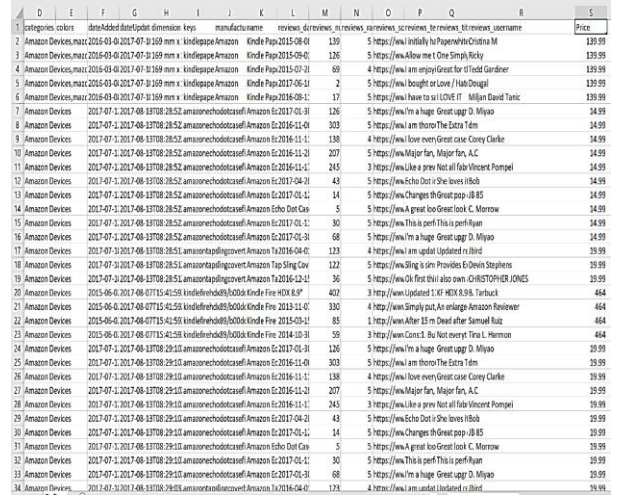


Fig. 2. Sample of raw data.



Fig. 3. Processed data for analysis.

## III. SENTIMENT ANALYSIS

Sentiment analysis refers to the use of natural language processing, text analysis, computational linguistics, and biometrics to systematically identify, extract, quantify, and study affective states and subjective information [13].

### A. Monkey Learn API

We have used monkey learn API for a python programming language to do sentiment analysis. we have trained our models with more than 150 comments which were positive, negative and neutral. We used ratings 4 and 5 for positive comments, 1 and 2 for negative and 3 for neutral comments.

### B. Positive Sentiments

Fig. 4 shows that when you try to run a review text from our dataset on our trained model its shows positive sentiment with accuracy.
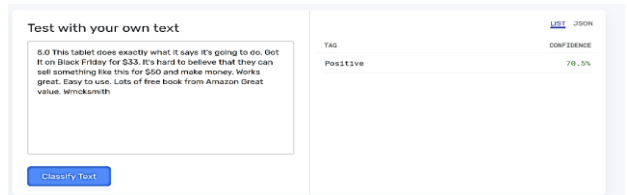


Fig. 4. Positive sentiment.

### C. Negative Sentiment

Fig. 5 shows that negative sentiment on a product review with accuracy and confidence. We can improve our confidence by training our model with more negative words.
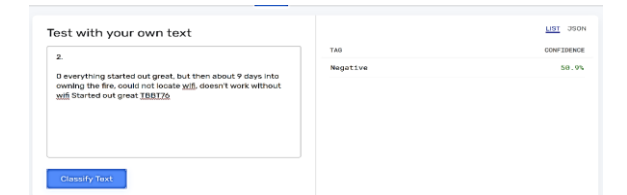


Fig. 5. Negative sentiment.

### D. Neutral Sentiment

Fig. 6 shows neutral sentiment with accuracy and confidence. We can improve accuracy and confidence by

training our model with more and more sentences either, positive, negative or neutral.



Fig. 6. Neutral sentiment.

## IV. CORRELATION

Correlation means association - more precisely it is a measure of the extent to which two variables are related. There are three possible results of a correlational study: a positive correlation, a negative correlation, and no correlation [14].

### A. Correlation between Number of Helpful Reviews and Price

Fig. 7 describes the correlation between the number of Helpful reviews and price using a trendline, as we can see that from trendline it shows that both variables have a negative correlation between them. The color shows details about the correlation. We also calculated the correlation between rating and the hopeful number and found -0.09597 for a sample size of 2800. The correlation calculated is not conclusive with this sample set. It shows overall rating increases, but a helpful number does not increase rather decrease. But the graph between rating and the helpful number is showing a mixed reaction. At rating 5 we see an increase in helpful number but at 2 we see decrease compare to 1. Rating 1 shows more increase in helpful numbers than rating 2 and 3. This is the most likely reason that we have a negative correlation in this scatter plot [15].

### B. Hypothesis Testing Using P-Value

In this section, we want to take a look at the data from the null hypothesis angle and see how much we learn from the data. We look at the rating attribute and the hopefulness number in the Amazon data set for our sentimental analysis. Fig. 8 shows that helpful number vs the rating. There is little information that can be obtained from the scatter plot except for the fact that higher ratings getting higher [16].
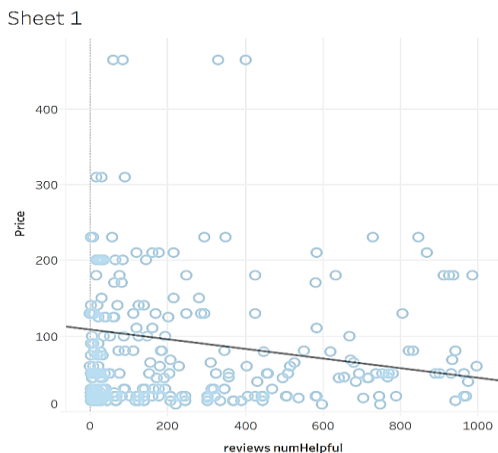


Fig. 7. Correlation between price and helpful reviews.

Hopeful number. But a lower rating of 1 also gets a significant hopeful number. Now we look at another figure that shows a number of items vs rating. Fig. 8 clearly shows a large bias towards a higher rating. Fig. 9 shows a lot of electronics product is getting a higher rating [17].
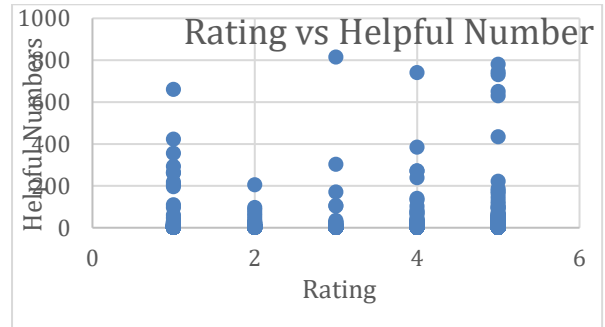


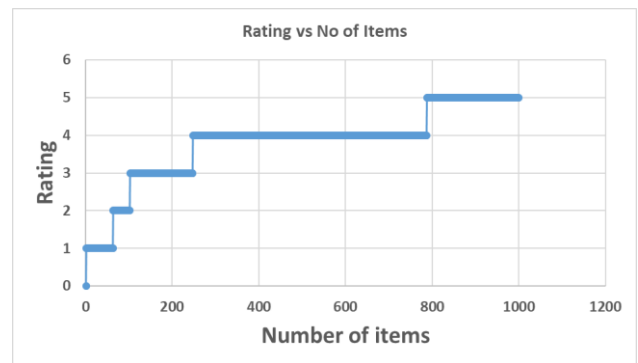Fig. 8. Rating vs hopeful number.



Fig. 9. Rating vs no of electronics items.

After taking a look at the scatter plot we want to perform some hypothesis analysis for different attributes of our data that would be the category of the data vs the rating again. Rating is one of the most important attributes in our Amazon data sentimental analysis.

First, we pick the electronics supply category and we want to see if our null hypothesis passes the p-value test that is computed using a chi-test in excel. The table for the electronics supply with a sample size of 100 is shown in Table I. We have not that many samples of electronics supply and therefore we just randomly picked a sample size of 100 [18].

TABLE I: THE CHI-SQUARE TEST FOR ELECTRONICS SUPPLY

| Electronics Supply Sample Size =100 | | |
|---|---|---|
| Rating | Expected | Observed |
| 1 | 20 | 0 |
| 2 | 20 | 0 |
| 3 | 20 | 4 |
| 4 | 20 | 34 |
| 5 | 20 | 42 |
| P-Value Calculated=6.2951e-18 | | |

We took a liberal approach and use equal distribution for all 5 categories in the sample size of 100 i.e. each category got 20 samples. But when we observed the data, we were kind of surprised which his clearly shown by the p-value [19]. The p-value is very small which means our null hypothesis cannot be accepted and hence we go by the alternative hypothesis. Now we selected another category called office supply and picked randomly 250 samples since

we don't have that many samples in this category either. A summary is shown in Table II.

TABLE II: THE CHI-SQUARE TEST FOR ELECTRONICS OFFICE SUPPLY

| Electronics Office Supply Sample Size =250 | | |
|---|---|---|
| Rating | Expected | Observed |
| 1 | 10 | 3 |
| 2 | 10 | 2 |
| 3 | 20 | 11 |
| 4 | 90 | 66 |
| 5 | 120 | 168 |
| P-Value Calculated=2.75e-08 | | |

Like the previous case, the p-value calculated shown in Table III is very small too and thus we have to discard our null hypothesis and find an alternative hypothesis [19].

TABLE III: THE CHI-SQUARE TEST FOR ELECTRONICS HARDWARE

| Electronics Hardware Sample Size =1000 | | |
|---|---|---|
| Rating | Expected | Observed |
| 1 | 10 | 1 |
| 2 | 50 | 12 |
| 3 | 150 | 40 |
| 4 | 290 | 262 |
| 5 | 500 | 685 |
| P-Value Calculated=9.6e-40 | | |

Now we have to take a very liberal approach and gave equal distribution to all ratings. This case randomly picked different categories of electronics products [20]. The sample size is again in 1000. Each rating is given 200 i.e. 1/5 probability. But when we look at the actual Data, we saw again a huge bias towards higher rating but not as much as the previous examples [20]. Now when we computed the p-value using the chi-squared test then we got an extremely small value of P that means we have to again discard our null hypothesis and stay with the alternative hypothesis. This kind of Data, which is highly biased, it is hard to make a hypothesis that can be true just based on expectation. We have also plotted the rating vs the number of products in the following.

## V. VISUALIZATIONS AND ANALYSIS

Here, several figures generated by the Tableau software are used to show and discuss the distribution of ratings across the total number of reviews.

### A. Common Words Used in Reviews

Fig. 10 shows the common words used in the reviews which have good ratings [21]. The words with higher frequency are greater in font size and are darker in color. The font size of the word decreases, and the color becomes lighter as the frequency of the words in the reviews.

As can be seen, the words great, good, like, love, very are common in reviews with positive ratings (4 - 5 stars). Fig. 11 shows the common words used in the reviews which have bad ratings. The words with higher frequency are greater in font size and are darker in color. The font size of the word decreases, and the color becomes lighter as the frequency of the words in the reviews decreases. As can be seen, the word not is the most common in reviews with

negative ratings (1 - 2 stars) followed by other words like no, don't. The words great, good, easy which had a high frequency for reviews with positive ratings have a comparatively lower frequency in reviews with bad ratings.
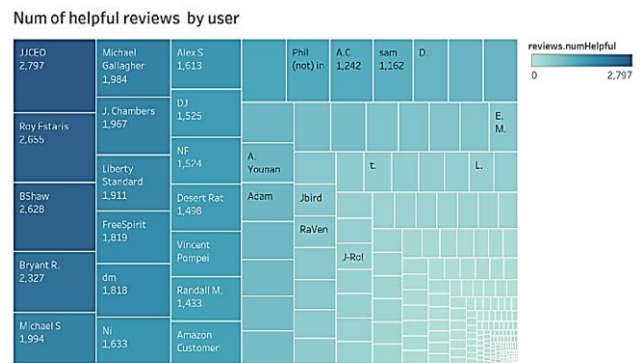

Fig. 10. Word cloud for good rating.


Fig. 11. Word cloud for bad rating.

### B. Number of Helpful Reviews by User

In Fig. 12 below we can see that users NF, Bashaw, Adam, Alex S and so on have a higher number of helpful reviews as compared to other users, so we can assume that these users reviews about any amazon product will be more helpful for any consumer who is reading reviews before buying any product.


Fig. 12. Most numbers of helpful reviews by users.

### C. Price of Product and Reviews

As we see in Fig. 13, there are many numbers of reviews for the products with a low-price range. As the price range of the products increases, the number of reviews decreases [21], [22]. This indicates that a smaller number of users buy expensive products online. Almost 60 percent of the total reviews are for the products that have a price range of $0-150. This shows that users prefer to buy products that are not extremely expensive online.

### D. Predictive Analysis Using Forecasting

Forecasting review year and review rating is shown in Fig. 14. The trend of the sum of Reviews. Rating (actual & forecast) for Reviews. Date Year. The color shows details about the Forecast indicator. The marks are labeled by the sum of Reviews [23], [24]. Rating (actual & forecast) is shown in Table IV.
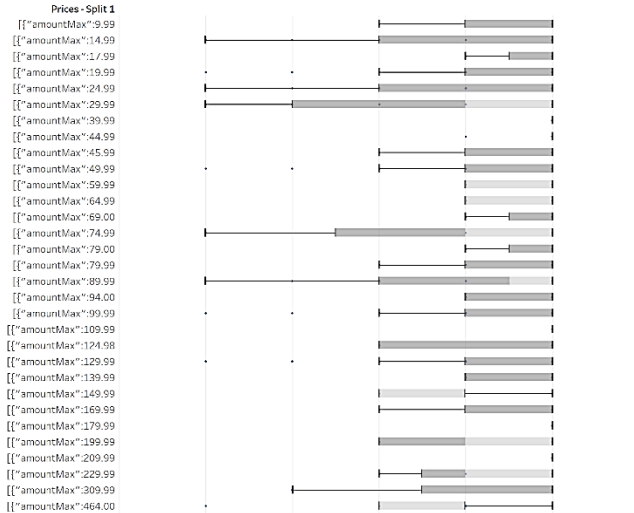


Fig. 13. Price of product vs reviews.



Fig. 14. Forecasting review rating vs review years.
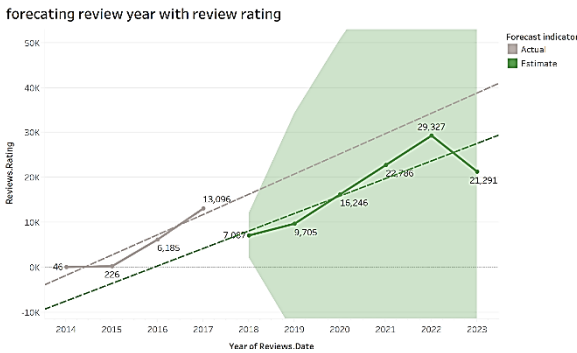
TABLE IV: REVIEW RATING

| Rows: | Reviews.Rating |
|---|---|
| Columns: | Year of Reviews. Date |
| Text: | Sum of Reviews.Rating (actual & forecast) |
| Color: | Forecast indicator |

### E. Trend Lines Model

A linear trend model is computed for the sum of Reviews. Rating (actual & forecast) given Reviews. Date Year. The model may be significant at $p <= 0.05$ as shown in Table V.

TABLE V: P VALUE COMPUTATION FOR TREND LINES

| Model formula: | Forecast indicator*( Year of Reviews. Date + intercept ) |
|---|---|
| p-value (significance): | 0.0031656 |

If we take a look at the dimensions of Forecast, we can easily see that in Fig. 14 we have two forecast indicators which are actual and estimate forecast. Where actual forecast Represent the actual data presented by our data set and estimated forecast represents that predictive data when is calculated using tableau forecast algorithm [25], [26] which represents the prediction of future reviews on different amazon products. Where actual years are between 2014 to 2018 and the estimated forecast is between the year 2019 to 2023 as shown in Table VI.

TABLE VI: SUM OF REVIEW RATING

| Initial | Changes from Initial | Seasonal Effect | | Contribution | | | |
|---|---|---|---|---|---|---|---|
| Sep 2018 | Sep 2018 – Aug 2023 | High | Low | Trend | Season | Quality | |
| 335 ± 1,169 | 2,400 | Dec 2022 1,673 | Apr 2023 -707 | 57.5% | 42.5% | Poor | |

### Options Used to Create Forecasts

| Time-series: | Year of Reviews. Date |
|---|---|
| Measures: | Sum of Reviews. Rating |

| Forecast forward: | 60 months (Sep 2018 – Aug 2023) |
|---|---|
| Forecast based on: | Oct 2014 – Aug 2018 |
| Ignore last: | 1 month (Sep 2018) |
| Seasonal pattern: | 12-month cycle |

### F. Product Reviews from the Last Five Years

Fig. 15 below shows the total number of reviews on different amazon products between the Year 2014 to 2018. The different color shows detail about the different product name. as we can see from Fig. 15 that for the year 2014 total number of reviews is very less than just under a thousand reviews on all products and reviews are a lot more during the year 2016 to 2017 which almost cross 12 Thousand on different amazon products [27].
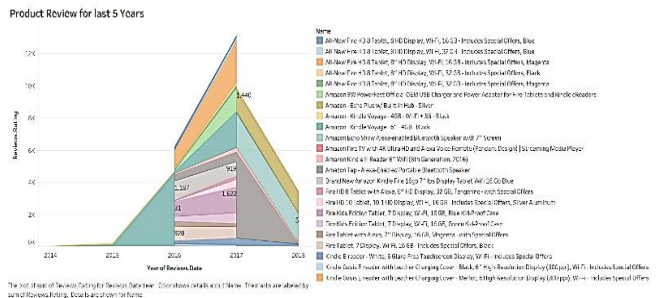


Fig. 15. Product yearly reviews.

## VI. CONCLUSION

From our study, we see that most of the amazon products have either four or five ratings in this dataset, so that's why in our plot we can't show much of information on 1, 2, and 3 rating products. We also see that product with high prices has less rating then product at low prices, it means most of the consumer has more review on low price product. And at last, we also see that for the first year that was 2014 total number of reviews on amazon products were very less, but

between the year 2016-2017, it was much higher.

In our statistical analysis, we calculate the correlation between product price vs the number of helpful reviews. We find that when the price goes up helpful reviews are less in number and vice versa, it means they have a negative correlation between them. We also perform a chi-square test to calculate p-value on products from different categories like office supply and hardware. For office supply, we took a sample size of 250 values since there are only a few products in that category [28]. We find that our expected values are different from observed values and the p-value is also very small. Hence it does not satisfy our null hypothesis.

For hardware we tool sample size of 1000 products, our results show that expected and observed values are again different, which again do not satisfy our null hypothesis.

Because our data is biased towards high ratings like 4 and 5 that's why it is hard to find a good hypothesis to calculate the chi-square test [29]. For future work, we want to expand our work with more products and more category and also compare amazon products with other electronics products from different companies and use advance visualization techniques to show more about the dataset [30].

## CONFLICT OF INTEREST

Both Authors do not have any conflict of interest.

## AUTHOR CONTRIBUTIONS

A. R. analyzed the raw data, transform the raw data into processed data Using SQL queries, applied BI techniques to create different Visualizations, applied sentiment analysis and Natural language processing Algorithms, used machine learning and data Mining techniques to preform correlation and chi-square test, wrote the paper. C.H. guided the overall research and revised the paper. Both authors approved the final draft of the paper.

## REFERENCES

[1] M. Hu and B. Liu, "Mining and summarizing customer reviews," in *Proc. the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2004, pp. 168–177.
[2] S. Mukherjee and P. Bhattacharyya, *Feature Specific Sentiment*.
[3] A. Bhatt, A. Patel, H. Chheda, and K. Gawande, "Amazon review classification and sentiment analysis," *International Journal of Computer Science and Information Technologies*, vol. 6, no. 6, pp. 5107–5110, 2015.
[4] X. Fang and J. Zhan, "Sentiment analysis using product review data," *Journal of Big Data*, vol. 2, pp. 1–14, 2015.
[5] J. McAuley. (2016). Amazon product data. [Online]. Available: http://jmcauley.ucsd.edu/data/amazon/
[6] J. McAuley, R. Pandey, and J. Leskovec, "Inferring networks of substitutable and complementary products," *Knowledge Discovery and Data Mining*, 2015.
[7] J. Mcauley, C. Targett, and A. Hengel, "Image-based recommendations on styles and substitutes," *SIGIR*, 2015.
[8] B. Pang, L. Lee, and S. Vaithyanathan, "Thumbs up? Sentiment classification using machine learning techniques," in *Proc. the ACL-02 Conference on Empirical Methods in Natural Language Processing*, 2002, pp. 79–86.
[9] S. Mukherjee and P. Bhattacharyya, "Feature specific sentiment analysis for product reviews," *Computational Linguistics and Intelligent Text Processing*, pp. 475–487, 2012.
[10] Tableau. (2016). [Online]. Available: http://interworks.co.uk/business-intelligence/why-tableau/
[11] T. Pinch and F. Kesler, *How Aunt Amy Gets Her Free Lunch: A Study of the Top-Thousand Customer Reviewers at Amazon.com*, 2011.
[12] Thelwall, K. Buckley, G. Paltoglou, D. Cai, and A. Kappas, "Sentiment strength detection in short informal text," *Journal of the American Society for Information Science and Technology*, vol. 61, no. 12, pp. 2544–2558, 2010.
[13] Yessenov and S. Misailovic, "Sentiment analysis of movie review comments," *Methodology*, pp. 1–17, 2009.
[14] Apache AsterixDB. (2019). [Online]. Available: https://asterixdb.apache.org/
[15] Integrate your data with cross-database joins in Tableau 10. (2019). [Online]. Available: https://www.tableau.com/about/blog/2016/7/integrate-your-data-cross-database-joins-56724
[16] Join Your Data - Tableau. (2019). [Online]. Available: https://onlinehelp.tableau.com/current/pro/desktop/en-us/joining_tables.htm#about-queries-and-crossdatabase-joins
[17] PostgreSQL: Documentation: 10: F.34.Ââpostgres_fdw. (2019). [Online]. Available: https://www.postgresql.org/docs/10/postgres-fdw.html
[18] PostgreSQL: The world's most advanced open source database. (2019). [Online]. Available: https://www.postgresql.org/
[19] A. Abbasi, S. Sarker, and R. Chiang, "Big data research in information systems: Toward an inclusive research agenda," *Journal of the Association for Information Systems*, vol. 17, no. 2, pp. i–xxxii, Feb. 2016.
[20] C. Alonso-Fernandez, A. Calvo, M. Freire, I. Martinez-Ortiz, and B. Fernandez-Manjon, "Systematizing game learning analytics for serious games," in *Proc. the Global Engineering Education Conference (EDUCON)*, 2017, pp. 1111–1118.
[21] D. Angrave, A. Charlwood, I. Kirkpatrick, M. Lawrence, and M. Stuart, "HR and analytics: why HR is set to fail the big data challenge," *Human Resource Management Journal*, vol. 26, no. 1, pp. 1–11, Jan. 2016.
[22] I. Benbasat and A. Dexter, "Value and events approaches to accounting: An experimental evaluation," *Accounting Review*, vol. 54, no. 4, pp. 735–749, 1979.
[23] M. Brehmer and T. Munzner, "A multi-level typology of abstract visualization tasks," *IEEE Transactions on Visualization and Computer Graphics*, vol. 19, no. 12, pp. 2376–2385, 2013.
[24] E. K. Choe, B. Lee, H. Zhu, N. H. Riche, and D. Baur, "Understanding self-reflection: How people reflect on personal data through visual data exploration," in *Proc. the 11th EAI International Conference on Pervasive Computing Technologies for Healthcare*, 2017, pp. 173–182.
[25] B. F. Chorpita, A. Bernstein, and E. L. Daleiden, "Driving with roadmaps and dashboards: Using information resources to structure the decision models in service organizations," *Administration and Policy in Mental Health and Mental Health Services Research*, vol. 35, pp. 114–123, Nov. 2007.
[26] J. Coffman, T. Beer, P. Patrizi, and E. H. Thompson, "Benchmarking evaluation in foundations: Do we know what we are doing?" *The Foundation Review*, vol. 5, no. 2, pp. 36–51, July 2013.
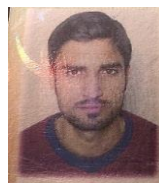[27] R. Crooks, "Representationalism at work: Dashboards and data analytics in urban education," *Educational Media International*, pp. 1–14, Dec. 2017.
[28] S. Few. (2017). There's nothing mere about semantics. [Online]. Available: https://www. perceptualedge.com/blog/?p=2793
[29] D. Filonik, R. Medland, M. Foth, and M. Rittenbruch, "A customisable dashboard display for environmental performance visualisations," *Persuasive Technology*, vol. 7822, pp. 51–62, 2013.
[30] K. Rall, M. L. Satterthwaite, A. V. Pandey, J. Emerson, J. Boy, O. Nov, and E. Bertini, "Data visualization for human rights advocacy," *Journal of Human Rights Practice*, vol. 8, no. 2, pp. 171–197, July 2016.

**Aamir Rashid** was a graduate student of School of Computer Science, Kean University, Union, New Jersey. Prior to joining, Rashid completed his bachelor of science in computer science from COMSATS University, Islamabad, Pakistan. His research interests are in the area of data science, data mining, machine learning and python programming. He graduated with a master of science in computer information systems from Kean University in May 2019.

**Ching-yu Huang** is an assistant professor of the School of Computer Science at Kean University since September 2014. Dr. Huang received a Ph.D. in computer & information science from New Jersey Institute of Technology, Newark, New Jersey, USA.

Prior to joining Kean University, Dr. Huang had more than 16 years of experience in the industry and academics in software development and R&D in bioinformatics. His research focuses SNP genotype calling and cluster detection; image processing and pattern recognition, especially in microarray and fingerprint; geotagged images and location information reconstruction; database application development; data processing automation; e-learning, educational multimedia, methodology, and online tools for secondary schools and colleges. Dr. Huang has more than 30 publications in journals and conferences and more than 40 presentations in workshops and invited lectures.