

# Predicting Online Reviewer Popularity: A Comparative Analysis of Machine Learning Techniques

Samadrita Bhattacharyya<sup>(✉)</sup>, Shankhadeep Banerjee,  
and Indranil Bose

Management Information Systems, Indian Institute of Management Calcutta,  
Calcutta, India

{samadritabl4, shankhadeepbl5, bose}@iimcal.ac.in

**Abstract.** Online customer reviews have been found to vary in their level of influence on customers' purchase decisions depending on both review and reviewer characteristics. It is logical to expect reviews written by popular reviewers to wield more influence over customers, and therefore an investigation into factors which can help explain and predict reviewer popularity should have high academic and practical implications. We made a novel attempt at using machine learning techniques to classify reviewers into high/low popularity based on their profile characteristics. We compared five different models, and found the neural network model to be the best in terms of overall accuracy (84.2%). Total helpfulness votes received by a reviewer was the top determinant of popularity. Based on this work, businesses can identify potentially influential reviewers to request them for reviews. This research-in-progress can be extended using more factors and models to further enhance the accuracy rate.

**Keywords:** Online reviews · Reviewer popularity · Machine learning techniques · Predictive analytics

## 1 Introduction

Widespread access to the internet is changing the way modern consumers make their purchase decisions. This change is facilitated by e-commerce platforms like Amazon and other online review websites like Yelp which host customer reviews of various products and services. Online customer reviews, also known as electronic word-of-mouth (eWOM) can be defined as “*any positive or negative statement made by potential, actual, or former customers about a product or company, which is made available to a multitude of people and institutions via the Internet*” [1]. According to a survey<sup>1</sup>, for 90% of customers, their buying decisions were influenced by online reviews. Academic studies too have confirmed the importance of online reviews on customers' purchase decisions [2]. However, consumers find some reviews more

---

<sup>1</sup> <http://marketingland.com/survey-customers-more-frustrated-by-how-long-it-takes-to-resolve-a-customer-service-issue-than-the-resolution-38756>

helpful than the others. This could be either because of the characteristics of the review content (review length, review polarity, content and style) [3, 4] and/or the characteristics of the reviewer.

It has been found that apart from information quality, source credibility is an important aspect for information seeking and adoption [5, 6]. In the context of online reviews, ‘source’ would imply the reviewer, hence reviewer credibility should impact review adoption. An accumulation of reviewer credibility over time should lead to more customers following a reviewer, leading to more reviewer popularity. Thus, popularity of reviewers should be associated with the impact of their reviews on the customers. The importance of word-of-mouth of influential reviewers on consumers’ purchase decision is well established in previous research [7, 8]. However, little attention has been given to study factors which make reviewers popular, and hence more influential. These factors could be used to predict reviewer popularity, which could be useful for businesses. Hence, in this research-in-progress, we attempt to identify popular reviewers based on their online profile characteristics. We use data from Yelp website which has an extensive reviewer community and detailed reviewer attributes. One of such attributes is the number of followers of a reviewer, which we use as a proxy for popularity. Other information regarding the reviewer include the number of reviews written, the number of friends a reviewer has, average rating the reviewer provides, years of experience in writing reviews for the website, etc.

We used five different machine learning techniques to classify reviewers as high or low on popularity based on their profile characteristics provided, and compared their performances. Also, we identified the factors which were most impactful in predicting reviewer popularity.

Insights shared in this study might help businesses in targeting popular reviewers for writing reviews about their offerings. Again, by predicting popularity of a reviewer, review websites might prioritize the display of a new review which is yet to get a helpfulness vote. The study also contributes to the growing research on online customer reviews and to the best of our knowledge, this is a novel attempt of using predictive analytics in the context of reviewer popularity.

## 2 Literature Review

Past literature one-WoM has primarily focused on identifying the factors related to helpfulness of reviews [9, 10]. A study by [11] found that perceived value of review is influenced by reviewer’s expertise and reputation. [10] used average helpfulness votes received per review and personal information disclosure for finding impact on review helpfulness. Another research has found that reviewer characteristics like the number of reviews posted by a reviewer and the number of helpful votes received by the reviewer on the whole, impacts the helpfulness vote of a review [12]. Study also found that reviews written by a self-described expert are more helpful than those that are not [13]. Some reviewer characteristics such as reviewer quality and reviewer exposure are found to impact sales by reducing perceived uncertainty of buyers [14]. However, to the best of our knowledge there has been no research to identify the dominant factors responsible for making a reviewer popular.

In our study we attempt to differentiate relatively more popular and less popular reviewers using a predictive analytics approach. The determining factors are selected based on support from extant literature and availability of data. Since we did not find much literature directly examining the factors related to reviewer popularity, so we had to use nearest available proxy which is review popularity to justify our variables. Review popularity is based on the perceived helpfulness of the reviews. Review helpfulness has been found to be influenced by reviewer characteristics [13, 15, 16], and thus justifies its use as the proxy. Hence, for predicting reviewer popularity, we identify factors from literature which influence helpfulness of the reviews.

[17] using TripAdvisor data studied the effect of review polarity on helpfulness of the review. They found that reviewers who posted more positive reviews are more likely to receive helpful votes than those who stressed on the negative aspects. Star rating is an indicator of reviewer's polarity. Hence, we include average review rating as a factor of reviewer popularity.

Another important factor found to be influential in deciding review helpfulness is review length [4, 18, 19]. Studies have found that review length provides important cues regarding reviewer characteristics [20]. Prior literature has already established that number of reviews are associated with review helpfulness [11, 12]. Drawing on these we can say that more the number of reviews a user writes on a forum, more popularity she gains.

Also, it could be said that with increased experience, a reviewer writes more useful reviews and hence become more popular. Similarly, the helpfulness votes a reviewer receives should be associated with her popularity, since it validates her credibility. The research by [10] also confirms that average helpfulness votes received by reviewer as one of the possible reviewer characteristics that might affect review helpfulness.

Being a well-reputed reviewer with certification from the website ('Elite' in case if Yelp website) also establishes the reviewer credibility and in turn should influence popularity. Finally, more the number of friends a reviewer has, more popular she is expected to be [11] used number of friends as a measuring variables for reviewers' reputation, hence we use number of friends as a factor of reviewer popularity.

### 3 Data and Methodology

A large dataset with 552,339 records was collected from [Yelp.com](https://www.yelp.com) which was made public as a part of the Yelp Dataset Challenge 2016. After processing data and removing outliers we had 69,612 records which we used for analysis. Yelp hosts customer reviews on local businesses. The reason behind selecting the website is that it provides information regarding the reviewers and their followers. User attributes such as number of followers, number of friends, average review rating, number of reviews written, total helpfulness votes, years of experience, years of reputation, and average review length for each reviewer were provided in the data. Number of followers was used as a proxy for reviewer popularity. Description on the data is shown in Table 1.

**Table 1.** Data descriptive

Variable	Range	Mean	SD
Number of followers	1–23	2.042	1.934
Average review rating per user	1–5	3.816	0.587
Number of reviews written	3–284	45.026	46.541
Years of experience	1–12	5.214	2.084
Years of reputation	0–3	0.133	0.435
Average review votes per user	0–8.5	2.085	1.414
Number of friends	1–49	7.186	9.361
Average review length	1–58	23.114	13.130

We used clustering technique (2-stage clustering) to decide on the number of segments appropriate for classification. The results showed two distinguished clusters. On observation of the clustered data we found mean value of number of followers of a reviewer to be the demarcation value. The reviewers having followers more than mean value are said to be high on popularity and vice-versa. The outcome variable is binary with 1 representing high and 0 representing low.

Data was partitioned in 70:30 ratio for training and testing. Five different models were used for classification: C5, Neural network, Bayesian network, CHAID, and Logistic Regression. We used IBM SPSS modeler as the analytical tool. The models were compared based on overall accuracy, lift, and costs. The agreement of all the models were checked to ensure their comparability.

## 4 Results and Analysis

Different models had different values of accuracy, however, they didn't differ much. The overall accuracy was around 83%–84%. There was 83.8% agreement among the classification techniques. Table 2 summarizes the results.

**Table 2.** Summary of results for various predictive models

Models	Overall accuracy	Lift	Factors used	AUC
Neural network	84.2%	2.43	7	0.863
Logistic regression	83.82%	2.39	7	0.856
C5 1	83.62%	2.33	7	0.810
Bayesian network	83.52%	2.33	7	0.822
CHAID 1	83.38%	2.39	5	0.856

All the models show nearly same level of accuracy with neural network giving the best value. We found that number of reviews and average helpfulness votes received by a reviewer were the two most important predictors among all followed by number of friends and average review rating. The least important factors turned out to be average

review length and years of experience. Figure 1 depicts the predictor importance. Table 3 shows the confusion matrix for neural network. 85.9% of reviewers who are low on popularity are predicted correctly, whereas 70.1% of reviewers high on popularity were correctly predicted. Prediction accuracy is higher for less popular class. Businesses would try to minimize the number of less popular reviewers being predicted as more popular ones, since that would incur costs in investing their time and resources on uninfluential reviewers. In our model, this case was found to be just 14.1%, which is on the lower side of error.

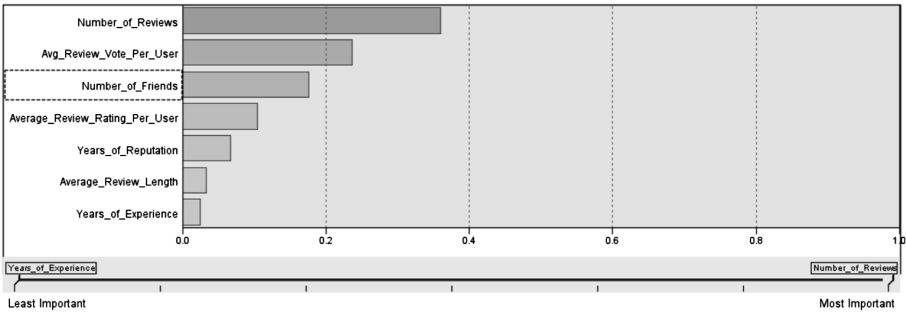


Fig. 1. Predictor importance graph

Table 3. Confusion matrix for neural network

Observed	Predicted	
	0	1
0	85.9%	14.1%
1	29.9%	70.1%

All other models except CHAID used all of the inputs to predict the output variable. CHAID model discarded years of ‘Elite’ and average review length as predictors.

## 5 Conclusion

In this research-in-progress paper, we attempted to use predictive analytics to classify online reviewers into two distinct classes based on their popularity. We compared five different machine learning techniques - C5, Neural network, Bayesian network, CHAID, and Logistic Regression. Among all, the neural network model turned out to be the best with 84.2% accuracy. Number of reviews written was found to be the most important factor.

In future we plan to incorporate few more factors such as review content characteristics, and try to improve the accuracy of prediction. For example, review

subjectivity, polarity, topic relevance, spelling & grammar, etc. could be some more variables to consider. Additionally, we also want to create ensemble models using more than one predictive models (like using a neural network for accuracy and a decision tree for rules), and analyze the results.

Owing to the significant impact of online reviews on customers' purchase decisions, it is important for the businesses to manage the reviews received for their products and services. In order to get more impactful reviews, it is important for businesses to identify influential and popular reviewers. Based on some characteristics or cues about the reviewer, if it is possible to predict reviewer popularity, businesses could leverage the information to target those reviewers and encourage them to write reviews about their products or services. If necessary, they might also incentivize the most popular reviewers. Also, it is advisable for the businesses to keep a track of the issues raised by popular reviewers to proactively address those. Businesses could use extract ideas from their reviews to enhance the offerings if needed.

For e-commerce and online review sites, the insights from this study could be helpful in many ways. They can develop recommender systems based on different characteristics of a reviewer, predict their popularity, and display their reviews as top reviews on their sites. This would be particularly useful for those websites where social interaction (like following) among reviewers and other consumers is not possible.

## References

1. Hennig-Thurau, T., Gwinner, K.P., Walsh, G., Gremler, D.D.: Electronic word-of-mouth via consumer-opinion platforms: what motivates consumers to articulate themselves on the internet? *J. Interact. Mark.* **18**(1), 38–52 (2004)
2. Zhu, F., Zhang, X.: Impact of online consumer reviews on sales: the moderating role of product and consumer characteristics. *J. Mark.* **7**(2), 133–148 (2010)
3. Korfiatis, N., García-Bariocanal, E., Sánchez-Alonso, S.: Evaluating content quality and helpfulness of online product reviews: the interplay of review helpfulness vs. review content. *Electron. Commer. Res. Appl.* **11**(3), 205–217 (2012)
4. Mudambi, S.M., Schuff, D.: What makes a helpful online review? A study of customer reviews on Amazon.com. *MIS Q.* **34**(1), 185–200 (2010)
5. Cheung, C.M.K., Lee, M.K.O.: Information adoption in an online discussion forum. In: 2nd International Conference on e-Business, ICE-B 2007, pp. 322–328 (2007)
6. Zhang, W., Watts, S.: Knowledge adoption in online communities of practice. In: ICIS 2003 Proceedings, vol. 9(1), pp. 96–109 (2003)
7. Kim, Y., Srivastava, J.: Impact of social influence in e-commerce decision making. In: Proceedings of ninth International Conference on Electronic Commerce - ICEC 2007, p. 293 (2007)
8. Li, Y.M., Lin, C.H., Lai, C.Y.: Identifying influential reviewers for word-of-mouth marketing. *Electron. Commer. Res. Appl.* **9**(4), 294–304 (2010)
9. Devi, J.: Estimating the helpfulness and economic impact of product reviews. *Int. J. Innov. Res. Dev.* **1**(5), 232–236 (2012)
10. Ghose, A., Ipeirotis, P.G.: Estimating the helpfulness and economic impact of product reviews: mining text and reviewer characteristics. *IEEE Trans. Knowl. Data Eng.* **23**(10), 1498–1512 (2011)

11. Liu, Z., Park, S.: What makes a useful online review? Implication for travel product websites. *Tour. Manag.* **47**, 140–151 (2015)
12. Otterbacher, J.: Helpfulness' in online communities: a measure of message quality. In: *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pp. 1–10 (2009)
13. Connors, L., Mudambi, S.M., Schuff, D.: Is it the review or the reviewer? A multi-method approach to determine the antecedents of online review helpfulness. In: *Proceedings of the 44th Hawaii International Conference on System Sciences*, pp. 1–10 (2011)
14. Hu, N., Liu, L., Zhang, J.: Do online reviews affect product sales? The role of reviewer characteristics and temporal effects. *Inf. Technol. Manag.* **9**(3), 201–214 (2008)
15. Forman, C., Ghose, A., Wiesenfeld, B.: Examining the relationship between reviews and sales: the role of reviewer identity disclosure in electronic markets. *Inf. Syst. Res.* **19**(3), 291–313 (2008)
16. Ngo-Ye, T.L., Sinha, A.P.: The influence of reviewer engagement characteristics on online review helpfulness: a text regression model. *Decis. Support Syst.* **61**(1), 47–58 (2014)
17. Fang, B., Ye, Q., Kucukusta, D., Law, R.: Analysis of the perceived value of online tourism reviews: influence of readability and reviewer characteristics. *Tour. Manag.* **52**, 498–506 (2016)
18. Kim, S.-M., Pantel, P., Chklovski, T., Pennacchiotti, M.: Automatically assessing review helpfulness. In: *Proceedings of 2006 Conference Empirical Methods Nature Language Processing (EMNLP 2006)*, pp. 423–430, July 2006
19. Schindler, R.M., Bickart, B.: Perceived helpfulness of online consumer reviews: the role of message content and style. *J. Consum. Behav.* **11**(3), 234–243 (2012)
20. Baek, H., Ahn, J., Choi, Y.: Helpfulness of online consumer reviews: readers' objectives and review cues. *Int. J. Electron. Commer.* **17**(2), 99–126 (2012)