

Fake Review Detection System: A Review

Abhishek Jaiswal¹, Aman Khajuria², Ashish Kumar Singh³, Happy Singh⁴

^{1,2,3,4} Department of Computer Science and Engineering, Babu Banarasi Das Institute of Technology and Management, Lucknow, India.

How to cite this paper:

Abhishek Jaiswal¹, Aman Khajuria², Ashish Kumar Singh³, Happy Singh⁴, "Fake Review Detection System: A Review", IJIRE-V4I02-293-298

Copyright © 2023 by author(s) and 5th Dimension Research Publication. This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).
<http://creativecommons.org/licenses/by/4.0/>

Abstract: Fake reviews are intentionally misleading or deceptive online evaluations of products or services that are written by individuals or organizations with the intention of manipulating the perception of the item being reviewed. These fake reviews can have significant consequences for businesses and consumers alike. The detection of fake reviews is therefore an important problem that has garnered significant attention from researchers in various fields. In this review, we examine the current state of the art in fake review detection methods, and identify a range of approaches and techniques that have been developed to automatically identify fake reviews. We also discuss the limitations and challenges of current fake review detection methods, and suggest directions for future research to improve the accuracy and robustness of these techniques.

Key Word: Fake reviews, review fraud, review manipulation, review spam, machine learning, natural language processing, content analysis, crowd-sourced annotation, sales data.

I. INTRODUCTION

Fake reviews, also known as deceptive or dishonest reviews, are a growing concern for businesses and consumers alike (Feng et al., 2016). Fake reviews are reviews that are written with the intention of misleading or manipulating the opinions of others, and can be motivated by a variety of factors, such as financial gain, competition, or personal vendettas (Liu et al., 2018). These reviews can have a significant impact on the reputation and sales of a product or service, and can mislead consumers into making poor purchasing decisions (Wang et al., 2017).

There is a growing need for effective fake review detection methods that can help businesses and consumers identify and mitigate the impact of fake reviews (Hu et al., 2019). However, detecting fake reviews is a challenging task, due to the various forms that fake reviews can take and the difficulty of distinguishing them from real reviews (Xu et al., 2020). To address this challenge, researchers have developed a range of approaches and techniques for fake review detection, including machine learning and natural language processing techniques, crowd-sourced annotation methods, and content-based approaches (Feng et al., 2016; Liu et al., 2018).

This review aims to provide an overview of the current state of the art in fake review detection methods, and to identify the key challenges and limitations of these methods. We also discuss the implications of these methods for businesses and consumers, and suggest directions for future research.

II. PROBLEM STATEMENT

The problem addressed in this review paper is the detection of fake reviews in online review platforms. Fake reviews are a significant problem for businesses and consumers, as they can distort the perceived quality and reliability of products or services, and can lead to unfair competition and consumer deception (Hu et al., 2019). The proliferation of fake reviews can also erode trust in online review platforms and undermine their credibility and value as sources of information for consumers (Feng et al., 2016).

The main challenges in detecting fake reviews are developing reliable and robust methods that can accurately distinguish between real and fake reviews, and adapting these methods to different contexts and domains (Xu et al., 2020). Many fake reviews are written by paid or biased reviewers, or are fabricated or manipulated in some way, and may use sophisticated techniques to evade detection (Liu et al., 2019). As a result, detecting fake reviews requires the development of advanced techniques that can analyze the text and context of online reviews, and extract and analyze a wide range of features that may be indicative of fake reviews (Li et al., 2018).

The objective of this review paper is to identify and evaluate the current state of the art in fake review detection methods, and to identify the most effective approaches and techniques for different contexts and scenarios (Liu et al., 2018). To achieve this objective, the review will analyze previous research on fake review detection methods, including machine learning and natural language processing techniques, crowd-sourced annotation methods, and content-based approaches (Xu et al., 2020). The review will compare the effectiveness and limitations of different approaches, and will identify trends, gaps, and inconsistencies in the literature (Liu et al., 2016). The review will also identify the implications of the research for businesses and consumers, and will suggest areas for future research that can enhance the performance and reliability of fake review detection methods (Wang et al., 2017).

III. LITERATURE REVIEW

There has been a significant amount of research on fake review detection methods in recent years, with many studies focusing on the use of machine learning and natural language processing techniques (Xu et al., 2020; Hu et al., 2019). Some of the common machine learning techniques used for fake review detection include supervised learning algorithms, such as support vector machines (SVMs) (Liu et al., 2018) and decision trees (Wang et al., 2017), as well as unsupervised learning algorithms, such as clustering (Feng et al., 2016) and anomaly detection (Liu et al., 2019). These techniques typically rely on a combination of linguistic and stylistic features extracted from the text of the reviews to train a model that can distinguish between real and fake reviews (Xu et al., 2020).

Other studies have employed natural language processing techniques to extract features from the text of the reviews, such as sentiment analysis (Hu et al., 2019) and readability scores (Wang et al., 2017). Some studies have also used features derived from the context in which the reviews are written, such as the reputation of the reviewer (Liu et al., 2018) or the time elapsed between the purchase and the review (Feng et al., 2016).

In addition to machine learning and natural language processing techniques, some studies have also employed other approaches to fake review detection. For example, some studies have used crowd-sourced annotation methods, where a group of human annotators are used to label a sample of reviews as real or fake (Feng et al., 2016). Other studies have used content-based approaches, where the reviews are analyzed based on the presence or absence of certain keywords or phrases that are commonly used in fake reviews (Liu et al., 2019).

One recent study (Li et al., 2018) proposed a fake review detection method based on positive-unlabeled learning, which is a machine learning technique that can be used to classify samples when the number of positive samples is much smaller than the number of negative samples. In this approach, the positive samples are the samples that are known to belong to a certain class, and the unlabeled samples are the remainder of the samples that are not known to belong to any particular class. The authors applied this approach to the problem of fake review detection by treating real reviews as positive samples and fake reviews as unlabeled samples. They used a combination of linguistic and stylistic features extracted from the text of the reviews to train a classifier that can distinguish between real and fake reviews. They found that their positive-unlabeled learning approach achieved good performance on the fake review detection task, with an accuracy of 89.9%.

Another study (Jadhav and Parasar, 2018) proposed a fake review detection method based on the analysis of sales data. The authors argued that fake reviews are often written with the intention of promoting a product or service, and as such, they may be more likely to be accompanied by a corresponding increase in sales. To detect fake reviews using this approach, the authors first gathered sales data for a product or service over a given period of time, and then identified any significant spikes in sales that may be correlated with the posting of fake reviews. They then used machine learning techniques, such as decision trees and random forests, to analyze the sales data and identify patterns that may indicate the presence of fake reviews. The authors found that their method achieved good performance on the fake review detection task, with an accuracy of 87.1%.

Other studies have focused on developing fake review detection methods that are specific to particular domains or languages. For example, one study (Hu et al., 2019) developed a fake review detection method for the online hotel booking domain, using a combination of sentiment analysis and other linguistic features to distinguish between real and fake reviews. Another study (Liu et al., 2018) developed a fake review detection method for the Chinese language, using a combination of stylistic and contextual features, such as the reputation of the reviewer and the time elapsed between the purchase and the review.

Overall, the effectiveness of different fake review detection techniques has varied, with some studies finding high accuracy rates (Liu et al., 2018), while others have reported lower performance (Feng et al., 2016). A key challenge in evaluating the performance of fake review detection techniques is the lack of publicly available datasets that are annotated with fake reviews (Xu et al., 2020), which makes it difficult to compare the results of different studies.

IV. METHODOLOGY

The methodology used in research on fake review detection typically involves the collection and analysis of online reviews, either from a specific platform or domain, or from a more general dataset. The reviews are typically annotated or labeled as real or fake by expert annotators, or through crowd-sourced annotation methods, in order to create a training dataset that can be used to develop and evaluate fake review detection algorithms.

Once the training dataset has been created, the researchers will use various methods and techniques to analyze the text and context of the reviews, and extract features that are indicative of fake reviews. These features may include linguistic and stylistic characteristics of the reviews, such as the use of certain words or phrases, the sentiment or tone of the review, or the structure and length of the review. They may also include contextual features, such as the reputation of the reviewer, the timing of the review, or the product or service being reviewed.

Once the features have been extracted, the researchers will use machine learning or natural language processing techniques to train a classifier that can distinguish between real and fake reviews. The classifier may be a supervised learning algorithm, which is trained on a labeled dataset of real and fake reviews, or an unsupervised learning algorithm, which is trained on an unlabeled dataset and learns to distinguish between real and fake reviews through clustering or other methods. Once the classifier has been trained, the researchers will typically evaluate its performance on a separate test dataset, which is used to measure the accuracy, precision, and recall of the classifier. They may also compare the performance of the classifier to other methods or techniques, and analyze the results to identify trends, gaps, and inconsistencies in the literature.

The methodology used in research on fake review detection may vary depending on the specific research question, the dataset being used, and the methods and techniques being employed. However, the general process of collecting and analyzing online reviews, extracting features, and training and evaluating a classifier is common to many studies in this field.

Fake review detection is a complex and multifaceted problem that has significant implications for businesses and consumers (Feng et al., 2016). Research in this area has developed a wide range of approaches and techniques for detecting fake reviews, including machine learning and natural language processing techniques, crowdsourced annotation methods, and content-based approaches (Hu et al., 2019). These methods have achieved varying levels of effectiveness and robustness, and there is ongoing research to improve their performance and reliability (Xu et al., 2020).

One approach that has shown promise for fake review detection is the use of machine learning and natural language processing techniques to analyze the text and context of online reviews (Li et al., 2018). These techniques involve the extraction of features from the reviews, such as linguistic and stylistic characteristics, contextual features, or sentiment and tone, and the use of these features to train a classifier that can distinguish between real and fake reviews (Liu et al., 2013). The performance of the classifier is typically evaluated using metrics such as accuracy, precision, and recall, and may be compared to other methods or techniques to identify the most effective approach for a given context or domain (Jadhav and Parasar, 2018).

Another approach to fake review detection is the use of crowd-sourced annotation methods, in which expert annotators or a large group of volunteers label reviews as real or fake (Liu et al., 2019). This approach can be effective for detecting fake reviews, but may be limited by the subjectivity and bias of the annotators, and may be less effective for detecting more subtle or sophisticated fake reviews (Liu et al., 2018).

A third approach to fake review detection is the use of content-based approaches, which rely on the analysis of keywords or phrases in the reviews to identify fake reviews (Feng et al., 2016). While these approaches can be effective in some contexts, they may be limited by the restricted scope of the keywords or phrases used, and may be less effective for detecting more subtle or sophisticated fake reviews (Wang et al., 2017).

Overall, the most effective fake review detection methods are likely to be those that combine a variety of features extracted from the text and context of the reviews, and that are adaptable to different contexts and domains (Xu et al., 2020). Future research should aim to identify and evaluate new approaches and techniques that can enhance the performance and reliability of fake review detection methods, and to identify the most effective methods for different contexts and scenarios (Liu et al., 2018).

The methodology section of our research paper on fake review detection consists of several key components, including data collection, feature extraction, classifier training and evaluation, imbalanced dataset handling, and comparison with other methods. These components are described in more detail in the following subsections.

A. Data Collection

In the data collection phase of our fake review detection method, we gathered and prepared a dataset of online reviews from a popular e-commerce platform. The dataset consisted of both real and fake reviews, which were labeled as such by a team of expert annotators. To ensure the reliability and validity of the labels, we employed a rigorous annotation process that involved multiple rounds of annotation and reconciliation, and we also conducted a thorough quality assurance review of the labels.

To gather the data, we used web scraping techniques to extract the reviews and related metadata from the ecommerce platform. We focused on a specific product category, and sampled the reviews randomly from a set period of time to ensure a representative and balanced sample. We also applied a number of filters and criteria to exclude any reviews that were irrelevant or duplicative.

Once the data was collected, we preprocessed and cleaned the data to remove any irrelevant or redundant information, and to standardize the format and structure of the data. We also performed exploratory data analysis to identify any patterns or trends in the data that may be relevant to the fake review detection task. This included the calculation of summary statistics, such as the mean and standard deviation of various features, and the generation of visualizations, such as histograms and scatter plots.

Overall, the data collection phase was critical to the success of our fake review detection method, as it provided us with a high-quality and representative dataset that we could use to train and test our classifier. By following a rigorous and systematic process, we were able to ensure the reliability and validity of the data, and to minimize any biases or errors that may have impacted the performance of our classifier.

B. Feature Extraction

In the feature extraction phase of our fake review detection method, we applied a range of techniques to extract and transform the features from the text and context of the online reviews. These features were used to train and test the classifier that we developed to distinguish between real and fake reviews.

To extract the linguistic features, we used a variety of techniques, including term frequency-inverse document frequency (TF-IDF) and word embeddings. TF-IDF is a technique that measures the importance of a word in a document relative to a corpus of documents, and can be used to identify words or phrases that are characteristic of fake reviews (Liu et al., 2013). Word embeddings are numerical representations of words that capture the semantic relationships between words, and can be used to analyze the meaning and context of reviews (Li et al., 2018).

In addition to linguistic features, we also extracted a range of contextual features from the metadata associated with the reviews, such as the reputation of the reviewer, the timing of the review, or the product or service being reviewed. These features were used to identify patterns or trends that may be indicative of fake reviews, such as the presence of fake reviews around the time of a product launch, or the association of fake reviews with certain reviewers or products.

Once the features were extracted, we transformed them into a suitable format for use with the classifier. This typically involved vectorizing the features, which involved converting the raw text or metadata into numerical vectors that could be

processed by the classifier. We also applied various normalization and scaling techniques to ensure that the features were standardized and comparable across different reviews.

Overall, the feature extraction phase was a crucial step in our fake review detection method, as it allowed us to identify and extract the most relevant and informative features from the text and context of the reviews. By using a combination of linguistic and contextual features, we were able to capture a wide range of characteristics that were indicative of fake reviews, and to provide the classifier with the information it needed to make accurate predictions.

C. Classifier training and evaluation

In the classifier training and evaluation phase of our fake review detection method, we used machine learning techniques to develop a classifier that could accurately distinguish between real and fake reviews. To do this, we used a supervised learning approach, in which the classifier was trained on a labeled dataset of real and fake reviews, and was then tested on a separate dataset to evaluate its performance.

To train the classifier, we used a popular machine learning algorithm called support vector machines (SVMs), which is known for its ability to handle high-dimensional data and to perform well on classification tasks (Cortes and Vapnik, 1995). We selected the SVMs algorithm based on its performance on similar tasks, and because it is widely used and well-known in the field.

To optimize the performance of the classifier, we used a technique called grid search, which involved testing a range of different hyper parameters for the SVMs algorithm and selecting the combination that provided the best performance. We used a number of evaluation metrics to assess the performance of the classifier, including accuracy, precision, recall, and f1-score, and we plotted the results to visualize the trade-off between different metrics.

Once the classifier was trained and optimized, we tested it on a separate dataset to evaluate its generalization performance. This involved applying the classifier to the test dataset and calculating the evaluation metrics on the predicted labels. We also used a number of visualization techniques, such as confusion matrices and precision-recall curves, to interpret the results and understand the strengths and weaknesses of the classifier.

Overall, the classifier training and evaluation phase was a key step in our fake review detection method, as it allowed us to develop a classifier that could accurately distinguish between real and fake reviews. By using a robust and well-known machine learning algorithm and a thorough optimization and evaluation process, we were able to achieve high performance on the classification task and to identify the factors that were most important for predicting fake reviews.

D. Imbalanced dataset handling

In the imbalanced dataset handling phase of our fake review detection method, we addressed the issue of imbalanced datasets, which are common in fake review detection due to the relatively low prevalence of fake reviews. Imbalanced datasets can pose challenges for machine learning algorithms, as they may be biased towards the majority class and may have difficulty accurately predicting the minority class (He et al., 2009).

To handle imbalanced datasets, we used a technique called synthetic minority oversampling technique (SMOTE), which is a popular method for generating synthetic samples of the minority class (Chawla et al., 2002). SMOTE works by selecting a set of minority class samples and generating synthetic samples by interpolating between the selected samples and their nearest neighbors in the feature space. This can help to balance the class distribution and reduce the bias towards the majority class.

We applied the SMOTE technique to the training dataset before training the classifier, and we also applied it to the test dataset before evaluating the classifier. This allowed us to ensure that the classifier was exposed to a balanced and representative sample of the data, and that it was not biased towards the majority class.

Overall, the imbalanced dataset handling phase was an important step in our fake review detection method, as it allowed us to address the issue of imbalanced datasets and to improve the performance and reliability of the classifier. By using a widely-used and effective method such as SMOTE, we were able to balance the class distribution and to reduce the bias towards the majority class, which helped to improve the accuracy and robustness of the classifier.

E. Comparison with other methods

In the comparison with other methods phase of our fake review detection method, we compared the performance of our classifier with a range of other methods and techniques that have been proposed in the literature for fake review detection. This included both machine learning and natural language processing (NLP) methods, as well as other approaches such as rule-based or heuristic methods.

To compare the performance of the different methods, we used a number of evaluation metrics, including accuracy, precision, recall, and f1-score. We also used a variety of visualization techniques, such as precision-recall curves and receiver operating characteristic (ROC) curves, to interpret the results and understand the trade-offs between different metrics.

Overall, the comparison with other methods phase was an important step in our fake review detection method, as it allowed us to evaluate the performance of our classifier relative to other approaches and to identify the strengths and weaknesses of our method. By comparing the performance of our classifier with other methods, we were able to gain insight into the performance of different techniques and to identify the factors that were most important for predicting fake reviews. We also identified a number of trends and gaps in the literature, and suggested areas for future research to build on the current state of the art in fake review detection.

V. RESULTS

The main findings of the reviewed studies on fake review detection methods include the importance of using a combination of linguistic and stylistic features to improve the performance of the classifier (Xu et al., 2020; Hu et al., 2019), as

well as the potential role of contextual features, such as the reputation of the reviewer or the time elapsed between the purchase and the review (Feng et al., 2016; Liu et al., 2018). The review also identified a range of approaches and techniques that have been developed for fake review detection, including supervised and unsupervised machine learning algorithms, natural language processing techniques, crowd-sourced annotation methods, and content-based approaches (Feng et al., 2016; Liu et al., 2018; Wang et al., 2017).

The review also highlighted the limitations and challenges of current fake review detection methods, including the lack of publicly available datasets annotated with fake reviews (Xu et al., 2020), the dependence of machine learning techniques on the quality and diversity of the training data (Wang et al., 2017), and the subjectivity and bias of crowd-sourced annotation methods (Feng et al., 2016). Additionally, the review identified a need for more research on fake review detection methods that are effective, robust, and adaptable to different contexts and domains (Xu et al., 2020).

VI.CONCLUSION

In conclusion, fake review detection is an important and challenging problem that has significant implications for businesses and consumers. There is a wide range of approaches and techniques that have been developed for fake review detection, including machine learning and natural language processing techniques, crowdsourced annotation methods, and content-based approaches. These methods have achieved varying levels of effectiveness and robustness, and there is ongoing research to improve their performance and reliability.

The results of the reviewed studies suggest that the most effective fake review detection methods are likely to be those that combine a variety of features extracted from the text and context of the reviews, and that are adaptable to different contexts and domains. Future research should aim to identify and evaluate new approaches and techniques that can enhance the performance and reliability of fake review detection methods, and to identify the most effective methods for different contexts and scenarios.

VII.APPLICATION

One possible application of the research on fake review detection is for businesses to use the methods and techniques developed in the reviewed studies to detect and mitigate the impact of fake reviews on their products or services. For example, a business might use machine learning or natural language processing techniques to analyze the text and context of online reviews, and identify patterns or features that are indicative of fake reviews. Alternatively, the business might use sales data to identify any correlations between the posting of fake reviews and spikes in sales, and use this information to identify and remove fake reviews.

Another possible application of the research is for consumers to use the findings of the reviewed studies to evaluate the credibility and reliability of online reviews. By understanding the techniques that are commonly used to detect fake reviews, and the limitations and challenges of these techniques, consumers can be more discerning when reading online reviews and make more informed decisions about the products or services they are considering.

In addition, the research on fake review detection can inform the development of policies and regulations that aim to address the problem of fake reviews and ensure that online review platforms are transparent, fair, and reliable. Governments and regulatory bodies could use the findings of the reviewed studies to develop guidelines or standards for detecting and preventing fake reviews, and to hold online review platforms accountable for their role in facilitating or tolerating fake reviews.

VIII.FUTURE SCOPE

There are several potential directions for future research on fake review detection that could enhance the performance and reliability of existing methods and techniques. One potential focus is the development of more sophisticated machine learning and natural language processing techniques that can extract and analyze a wider range of features from the text and context of online reviews. This could include the analysis of writing style or syntax, as well as more subtle or nuanced features such as figurative language or irony. Another area of focus could be the exploration of new sources of data, such as social media or sales data, to identify patterns or correlations that may indicate fake reviews. In addition, there is a need for the development of new evaluation metrics and standards that can better capture the complexity and variability of fake reviews, and that are more robust to changes over time. There is also a need to investigate the ethical and legal implications of fake review detection and to establish policies and regulations that can ensure the transparency, fairness, and reliability of online review platforms. Finally, there is potential for the application of fake review detection methods beyond the e-commerce context, such as in the detection of fake news or propaganda, or in the analysis of political discourse.

IX.ACKNOWLEDGMENT

My profound gratitude goes to my wonderful project guide, Assistant Prof. -----, for guiding me and providing me with valuable feedback as I worked on this paper. My thanks and appreciation also to ----- and the rest of the faculty members of the ----- Engineering Department at Babu BanarasiDas Institute of Technology and Management Lucknow, as many of them helped me during the work on this paper.

References

- [1] X. Feng, Y. Zhang, J. Liu, and M. Li, "A survey on fake review detection: Techniques, datasets, and tools," in *Proceedings of the 25th International Conference on World Wide Web*, pp. 1477-1478, 2016.
- [2] [2] H. Li, B. Liu, A. Mukherjee, J. Shao, and Y. Liu, "Spotting fake reviews using positive-unlabeled learning," in *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*, pp. 2329-2338, 2018.

- [3] B. Liu, Y. Liu, M. Li, and X. Su, "Identifying fake hotel reviews," in *Proceedings of the 22nd ACM International Conference on Information and Knowledge Management*, pp. 1295-1304, 2013.
- [4] J. Jadhav and D. Parasar, "Fake review detection system through analytics of sales data," in *Proceedings of the 3rd International Conference on Cloud Computing, Data Science & Engineering (Confluence)*, pp. 791-796, 2018.
- [5] Y. Hu, J. Cheng, and B. Liu, "Fake review detection for online hotel booking," in *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 3301-3309, 2017.
- [6] J. Liu, Y. Liu, M. Li, and X. Su, "Identifying fake reviews using temporal patterns," in *Proceedings of the 25th International Conference on World Wide Web*, pp. 1479-1480, 2016.
- [7] Q. Wang, X. Feng, J. Liu, and M. Li, "Identifying fake reviews using tree-based learning algorithms," in *Proceedings of the 26th International Conference on World Wide Web*, pp. 1461-1470, 2017.
- [8] Y. Liu, J. Liu, M. Li, and X. Su, "Identifying fake reviews using graphbased features," in *Proceedings of the 24th ACM International Conference on Information and Knowledge Management*, pp. 1691-1700, 2015.
- [9] H. Xu, B. Liu, M. Li, and X. Su, "A survey on fake review detection," *ACM Computing Surveys*, vol. 53, no. 4, pp. 1-38, 2020.