

Combatting Deceptive Reviews: Leveraging RoBERTa for Authenticity Assessment and Collaborative Spammer Detection in Buyer-Seller Networks

Kashish Padhiar
21BCE5704

Santosh R
21BCE5058

Samiksha Rajesh Nair
21BCE5832

Abstract— This manuscript presents an innovative methodology designed to combat the pervasive issue of fake reviews plaguing online platforms, posing significant risks to both consumers and businesses by perpetuating misinformation and potentially causing substantial financial losses. By harnessing the latest advancements in natural language processing (NLP) and deep learning, particularly leveraging the cutting-edge RoBERTa model, our approach seeks to significantly enhance the efficacy of fake review detection mechanisms. Through meticulous fine-tuning of the RoBERTa model using meticulously annotated review datasets, our system endeavours to refine its classification capabilities, aiming to achieve heightened precision and robustness in identifying fraudulent reviews. Through rigorous experimentation and validation, we showcase the effectiveness of our proposed methodology in accurately discerning fake reviews, demonstrating notably elevated levels of precision and recall compared to existing approaches. Furthermore, our methodology stands as a testament to the potential of advanced NLP techniques and deep learning models in addressing complex challenges within online ecosystems. By leveraging the inherent capabilities of RoBERTa to capture subtle linguistic nuances and contextual cues, we aim to fortify the defense against deceptive practices, thereby safeguarding the integrity and reliability of online review systems. Through this concerted effort, we aspire to foster a more trustworthy and transparent digital marketplace, where consumers can make informed decisions with confidence, while businesses can compete fairly and ethically on a level playing field.

I. INTRODUCTION

The exponential growth of online reviews has become fundamental to consumers' decision-making processes. Yet, this landscape is marred by the rampant proliferation of fake reviews, which erode the trustworthiness of online platforms and jeopardise consumer confidence. Accurate detection of these fraudulent reviews is imperative to uphold the integrity of online review systems and foster equitable competition among businesses. In response to this pressing challenge, our paper introduces a pioneering methodology aimed at bolstering fake review detection. We propose an innovative approach leveraging advanced natural language processing (NLP) techniques and deep learning models. By harnessing the capabilities of these cutting-edge technologies, we seek to significantly

enhance the efficacy of fake review identification mechanisms. Through the fusion of sophisticated NLP algorithms and deep learning frameworks, our methodology aims to discern subtle patterns and linguistic nuances indicative of fraudulent activity within review datasets. This novel approach holds the promise of revolutionising fake review detection, paving the way for a more transparent and trustworthy online marketplace where consumers can make informed decisions with confidence, and businesses can compete on a level playing field characterised by integrity and fairness.

II. RELATED WORKS

The content provided discusses the increasing prevalence of fake reviews on platforms like Amazon due to the rise in online shopping. To combat this issue, a methodology utilizing sentiment analysis, support vector machine (SVM) algorithm, and logistic regression (LR) algorithm is proposed. Sentiment analysis helps uncover patterns indicative of fake reviews, while SVM and LR algorithms are employed for classification based on sentiment features. The paper outlines steps for data collection, preprocessing, feature extraction, and model training and evaluation. Results show SVM achieving 86% accuracy, while LR achieves 87% accuracy in detecting fake reviews. The conclusion suggests future work to enhance accuracy by combining the proposed methods. [1] This research proposes a new method to fight fake reviews with machine learning. They test two models, LSTM and BERT, on identifying fake reviews. BERT achieved higher accuracy than LSTM. To improve reliability, they introduce Monte Carlo Dropout (MCD). MCD runs the model multiple times with slight variations and averages the results, making the model more certain about its predictions. In conclusion, this approach combines machine learning models with MCD to enhance fake review detection accuracy and reliability. [2] This is a very interesting research paper on using a convolutional neural network to detect fake reviews. It proposes a novel approach to detecting fake reviews that combines text analysis with other factors. The system is able to achieve high accuracy on a dataset of real and fake reviews. The paper also mentions that the Yelp dataset used to train the system is filtered, which means that some of the fake reviews may have already been removed. Methodology: The system uses a convolutional neural network (CNN) to analyze the text of the review. A CNN is a type of artificial intelligence that is good at identifying patterns in data. The system

also considers non-textual factors such as the rating of the review, the date it was written, and the reviewer's location. The authors test their system on a dataset of real and fake reviews. They find that their system is able to detect fake reviews with high accuracy. The paper also does not discuss the computational cost of training and running the system. Overall, this is a promising piece of research that could help to improve the quality of online reviews. [3] This article provides a comprehensive analysis of features used to detect fake reviews on online review sites. This can mislead both human readers and automated analysis systems. Supervised machine learning approaches are prevalent for detecting fake reviews. Reviewer-Centric Features: These features capture the reviewer's behavior and can expose patterns indicative of fake reviews. Examples include— Textual features: Maximum/Average Content Similarity, Word number average (to detect review duplication). Rating Features: Total number of reviews, ratios of positive/negative/extreme reviews, average deviation from entity's average rating, rating entropy, rating variance (new). Temporal Features: Activity time, maximum rating per day, date entropy, date variance (new). Overall, the article emphasizes the importance of using a comprehensive set of features, including newly proposed ones like burst features and reviewer temporal features with variance, to improve fake review detection. [4] This article proposes a new method to detect fake reviews based on multiple feature fusion and rolling collaborative training. This method combines review text features (sentiment, Doc2vec representation, part-of-speech), user behavior features (rating deviation, number of reviews), and a rolling collaborative training approach. Feature Extraction—Text features: Sentiment analysis using a sentiment dictionary to calculate sentiment intensity. Doc2vec model to represent text as a semantic vector. Rolling Collaborative Training—Use C1 and C2 to label unlabeled reviews with high confidence, add these high-confidence labeled reviews to the training sets for text and user behavior features (L1 and L2), update the classifiers C1 and C2 iteratively. The experiment used reviews from Yelp with fake reviews marked by the platform. Achieved higher accuracy (84.45%) compared to baseline methods (81%). Overall, this method shows promise for improving fake review detection accuracy by leveraging multiple features and unlabeled data through rolling collaborative training. [5] Fake reviews are a major concern for e-commerce businesses as they can mislead customers. Machine learning techniques are widely used to detect fake reviews, categorized into supervised and unsupervised learning. Supervised learning uses labeled data (real vs. fake reviews) to train algorithms for identifying fake reviews. This approach can be effective but requires a lot of labeled data which can be expensive and time-consuming to obtain. Unsupervised learning identifies fake reviews based on reviewer behavior and characteristics of their reviews, without labeled data. Deep learning techniques, such as Recurrent Neural Networks (RNNs) and Convolutional Neural Networks (CNNs), are also showing promise in fake review detection. They can capture the semantic meaning of text and identify patterns in reviewer behavior. [6] This is a very informative article about detecting fake reviewer groups

in online review systems. Fake reviewer groups can significantly affect the reputation of products online. Existing methods focus on individual reviewers, while fake reviewers often work in groups. Overall, REAL is a promising approach for detecting fake reviewer groups and improving the trustworthiness of online review systems. The paper discusses the limitations of using only textual features to detect fake reviewers. The authors acknowledge that REAL may not perform as well on very large datasets. They plan to improve the accuracy of REAL for detecting larger fake reviewer groups. [7] Fake reviews are a significant problem for e-commerce platforms like Yelp. They can mislead consumers and damage businesses' reputations. Machine learning models can be used to identify fake reviews based on textual features of the review text. This study compared tree-based models (Random Forest and XGBoost) with transformer-based models (BERT and GPT-3) for fake review detection on Yelp data. Transformer-based models, specifically GPT-3, outperformed tree-based models in terms of accuracy, precision, recall, and F1 score. The study acknowledges that tree-based models may still be preferable in some cases due to their faster training times and interpretability. [8] This document discusses fake reviews and how Naive Bayes can be used to identify them. Fake reviews mislead consumers and harm businesses. Naive Bayes, a machine learning algorithm, can analyze review text to identify features that suggest fakeness. Data is collected (e.g., reviews with labels of real or fake). Reviews are preprocessed (e.g., removing stop words, stemming words). The model is trained on labeled data to learn the characteristics of real and fake reviews. New reviews are classified as real or fake based on the trained model. Benefits of Naive Bayes— Simple and easy to implement, fast training time, works well with text data, provides insights into why a review is classified as fake. Naive Bayes generally outperforms SVM in this application in terms of accuracy, training speed, and interpretability. Naive Bayes assumes features are independent, which may not always be true. The document proposes a system that combines review content analysis with reviewer behavior analysis for potentially better detection. Overall, the document provides a good explanation of how Naive Bayes can be used as a tool to combat fake reviews. [9] This paper proposes a method for classifying online product reviews as real or fake using Dempster-Shafer (D-S) theory. Fake reviews mislead consumers and harm businesses. The D-S theory, a mathematical framework for reasoning with uncertainty, is used to analyze features of reviews and classify them. Preprocess review text (remove stop words and punctuation, convert to vectors). Use BERT for sentiment analysis (positive or negative). Train-test split the data (80% training, 20% testing). Define classifications (real or fake review). Assign mass functions to each feature (degree of belief). Combine mass functions using Dempster's rule. Classify the review based on the highest mass value. It can handle uncertainty inherent in fake review detection and achieves higher accuracy compared to other methods (based on the paper's findings). This research shows promise for using D-S theory to combat fake reviews and improve online shopping experiences. [10]

III. EXISTING SYSTEM

The prevailing systems devised for fake review detection predominantly hinge upon traditional methodologies characterized by feature engineering and shallow learning algorithms. Commonly employed techniques include support vector machines (SVM) or random forests, which rely on manually crafted features extracted from textual data. These features typically encompass various linguistic attributes such as word frequencies, syntactic patterns, or semantic cues. Through the training of classifiers on these engineered features, these systems aim to discern between genuine and fake reviews within online platforms. While these conventional approaches may exhibit reasonable accuracy under controlled conditions, they face notable challenges in adapting to the dynamic and evolving landscape of fake review tactics. One prominent limitation lies in their reliance on static feature sets, which may struggle to capture the intricate nuances and evolving strategies employed by malicious actors. As fake review tactics become increasingly sophisticated and nuanced, traditional systems may falter in detecting these subtleties, thus compromising their effectiveness. Moreover, the manual feature engineering process is inherently labor-intensive and may not scale well to large datasets or diverse domains. The need for domain-specific feature crafting further exacerbates this issue, potentially limiting the generalizability of the detection system across different contexts. Additionally, the performance of these systems may be hindered by the inherent subjectivity and ambiguity present in textual data, leading to suboptimal classification outcomes. In light of these challenges, there arises a pressing need for more advanced and adaptable methodologies in fake review detection. By leveraging cutting-edge techniques in natural language processing (NLP) and deep learning, such as transformer-based models like RoBERTa, it becomes possible to overcome the limitations of traditional systems. These modern approaches have demonstrated remarkable capabilities in capturing complex linguistic patterns and contextual nuances inherent in textual data, thereby offering a promising avenue for enhancing the efficacy of fake review detection mechanisms in online platforms.

IV. PROBLEM STATEMENT

The central problem addressed in this paper pertains to the precise identification of fraudulent reviews within online platforms. Faced with a corpus of reviews, the primary objective is to discern between genuine and fake instances based on the content provided. Fake reviews often manifest through a myriad of deceptive tactics, such as embellished claims, unnatural language usage, or skewed sentiment expressions, thereby posing considerable challenges to accurate detection. The overarching aim is to devise a system that is both resilient and scalable, capable of effectively distinguishing between authentic and fraudulent reviews with a high degree of precision and recall. By achieving this goal, the integrity of online review systems can be safeguarded, fostering trust among consumers and ensuring fair competition among businesses in the digital marketplace.

V. PROPOSED ARCHITECTURE

The proposed system harnesses the power of a fine-tuned RoBERTa model for the detection of fake reviews within online platforms. RoBERTa, a transformer-based model, is initially pretrained on vast text corpora, endowing it with the ability to comprehend intricate linguistic structures and semantic nuances inherent in textual data. Through the fine-tuning process, the parameters of RoBERTa are adapted to the specific task of fake review detection by training it on a meticulously annotated dataset of reviews. This fine-tuning mechanism enables RoBERTa to effectively discern between genuine and fake reviews, leveraging its comprehensive understanding of language patterns and contextual cues. By exposing RoBERTa to labeled review data during the fine-tuning phase, the model learns to recognize subtle indicators indicative of fraudulent activity, thus enhancing its discriminatory capabilities. Furthermore, the fine-tuning process augments RoBERTa's performance and generalization capabilities, enabling it to extrapolate learnings from the training dataset to accurately classify unseen reviews. Through the utilization of a fine-tuned RoBERTa model, the proposed system offers a sophisticated and adaptable solution to the challenge of fake review detection. By leveraging the inherent strengths of transformer-based architectures and fine-tuning techniques, the system aims to achieve heightened accuracy and robustness in identifying fraudulent reviews, thereby fortifying the integrity of online review systems and preserving consumer trust in the digital marketplace.

VI. MODULES AND MODULE DESCRIPTION

1. Data Collection Module:
 - Utilizes web scraping techniques or API calls to gather a diverse corpus of reviews from various online platforms.
 - Implements mechanisms to ensure data integrity and compliance with platform terms of service.
 - Aggregates reviews across different products or services to create a comprehensive dataset for training and evaluation purposes.
2. Preprocessing Module:
 - Performs text cleaning operations to remove noise and irrelevant information from the raw review data.
 - Implements tokenization to break down the text into individual tokens or words for further analysis.
 - Conducts stop words removal to eliminate common words that may not contribute to the classification task.
 - Normalizes the text by converting it to lowercase and handling issues like spelling variations or abbreviations.
3. Fine-Tuning Module:
 - Adapts the pretrained RoBERTa model to the specific task of fake review detection by fine-tuning its parameters.
 - Utilizes techniques such as transfer learning to leverage the knowledge gained from the pretrained model.
 - Implements strategies to prevent overfitting and ensure optimal performance on the target task.
 - Adjusts hyperparameters and optimization algorithms to fine-tune the model effectively.
4. Evaluation Module:

- Assesses the performance of the fine-tuned model on a separate test dataset to evaluate its effectiveness.
 - Calculates metrics such as accuracy, precision, recall, and F1-score to measure the model's performance.
 - Conducts statistical analysis to identify strengths and weaknesses of the model and potential areas for improvement.
 - Generates visualizations or reports to communicate evaluation results and insights to stakeholders.
5. Deployment Module:
- Integrates the trained model into a production environment to enable real-time fake review detection on online platforms.
 - Implements scalable and efficient deployment strategies to handle large volumes of review data.
 - Monitors model performance in production and implements mechanisms for continuous improvement and updates.
 - Collaborates with platform administrators to ensure compliance with relevant regulations and guidelines regarding automated review detection systems.

VII. APPROACH USED

The central algorithm employed in this project is fine-tuning, a technique that entails adjusting the parameters of the pretrained RoBERTa model to suit the task of fake review detection. Fine-tuning involves leveraging a labeled dataset specifically curated for this purpose, enabling the model to learn task-specific features and nuances inherent in identifying fraudulent reviews. By updating the parameters of RoBERTa based on the characteristics of fake reviews present in the labeled dataset, the model can adapt its representations to better capture the subtle linguistic cues and patterns indicative of fraudulent activity. Through the fine-tuning process, RoBERTa undergoes iterative adjustments that optimize its performance for the target task, effectively enhancing its discriminative capabilities and enabling more accurate classification of genuine and fake reviews. This iterative refinement enables the model to continually refine its understanding of fake review characteristics, leading to improved performance and robustness in real-world applications.

VII. TOOLS USED

- PyTorch: Deep learning framework for building and training neural network models.
- Transformers Library: Provides implementations of transformer-based models, including RoBERTa.
- Pandas: Python library for data manipulation and analysis.
- Scikit-learn: Machine learning library for implementing classification algorithms and evaluation metrics.
- Google Colab: Cloud-based platform for running Python code and training machine learning models.

VIII. IMPLEMENTATION RESULTS

The experimental findings validate the efficacy of the proposed approach in effectively identifying fake reviews with a remarkable level of accuracy and robustness. Through rigorous testing, the fine-tuned

RoBERTa model exhibits state-of-the-art performance on benchmark datasets, surpassing the capabilities of traditional machine learning algorithms and heuristic-based methods commonly employed in fake review detection. Evaluation metrics, including accuracy, precision, recall, and the F1-score, provide quantitative evidence of the superior performance achieved by the proposed system compared to existing approaches. Specifically, the fine-tuned RoBERTa model demonstrates a significant improvement in accuracy and precision, accurately discerning between genuine and fake reviews with heightened sensitivity and specificity. Moreover, its elevated recall ensures that a substantial proportion of fraudulent reviews are correctly identified, minimizing false negatives and enhancing overall detection effectiveness. The superior F1-score further corroborates the system's ability to strike a balance between precision and recall, reflecting its robust performance across various evaluation criteria. These findings underscore the transformative potential of leveraging advanced natural language processing techniques and deep learning models, exemplified by the fine-tuned RoBERTa model, in addressing the pervasive challenge of fake review detection. By achieving unprecedented levels of accuracy and robustness, the proposed system promises to bolster the integrity of online review systems, instill confidence among consumers, and foster fair competition among businesses in the digital marketplace.

```
device="cuda"
query = ""My old bot was wearing this to the Macy's in January.
This is the first one I've ever had. I'm a 320, and the first pair I bought were just a little tight.
I'm a bit disappointed.
This is my second pair.
I'm looking forward to wearing them to the Macy's in the fall.
I like the way they look. Love these! These are my favorite.
I have a hard time finding jeans that fit me comfortably, but I have a hard time finding jeans that don't fit.
These jeans are super comfortable and have a great price point.
I have some great jeans to wear for work, but these are the only jeans that I wear for work or for my family.
I will be buying more! I have a lot of compliments on them. I love these shoes.
I love the color and the fit. They fit my body well and are comfortable. I have a wide foot and these fit me well.

I'm 5'4", 130lbs and these fit well.
I would recommend them. I wear a size 11.5 in jeans and this fits perfect.
I have a narrow foot and this fits perfect. It is very comfortable and fits great.
I bought a small and it fit perfectly. I will order another size up. I bought these for my husband, he loves them and he loves them!
This is the best pair of sunglasses for the price! They are so comfortable and easy to use.
I wear them all the time and they don't hurt my feet. I wear them everyday and my feet are so happy with them!""
tokens = tokenizer.encode(query, return_tensors="pt")
all_tokens = len(tokens[0])
mask = torch.ones_like(tokens)

with torch.no_grad():
    logits = model(tokens.to(device), attention_mask=mask.to(device))[0]
    probs = logits.softmax(dim=-1)

fake, real = probs.detach().cpu().flatten().numpy().tolist()

print(f"Real Probability: {real}\nFake Probability: {fake}")

Real Probability: 0.000283153121173382
Fake Probability: 0.999716991233826
```

```
def predict(query, model, tokenizer, device="cuda"):
    tokens = tokenizer.encode(query)
    all_tokens = len(tokens)
    tokens = tokens[:tokenizer.model_max_length - 2]
    used_tokens = len(tokens)
    tokens = torch.tensor([tokenizer.bos_token_id] + tokens + [tokenizer.eos_token_id]).unsqueeze(0)
    mask = torch.ones_like(tokens)

    with torch.no_grad():
        logits = model(tokens.to(device), attention_mask=mask.to(device))[0]
        probs = logits.softmax(dim=-1)

    fake, real = probs.detach().cpu().flatten().numpy().tolist()
    return real

query = ""Worth the money Best mobile phone Camera quality is very nice Battery backup is very good Sound quality is amazing.""
predict(query, model, tokenizer)

0.911635484693989
```

	precision	recall	f1-score	support
CG	0.97	0.96	0.97	4010
OR	0.97	0.97	0.97	4077
accuracy			0.97	8087
macro avg	0.97	0.97	0.97	8087
weighted avg	0.97	0.97	0.97	8087

IX. CONCLUSION

In summary, this manuscript introduces a pioneering methodology for augmenting fake review detection through the utilization of a fine-tuned RoBERTa model. By harnessing cutting-edge natural language processing techniques and leveraging deep learning models, our system attains unparalleled performance levels in the identification of fake reviews, characterized by heightened precision and recall. The proposed approach represents a significant advancement in combating the pervasive issue of fraudulent reviews plaguing online platforms. By capitalizing on the inherent capabilities of advanced NLP and deep learning, the system achieves a comprehensive understanding of linguistic nuances and contextual cues, enabling it to discern subtle indicators of fraudulent activity within review datasets. This enhanced discriminatory prowess not only bolsters the integrity of online review systems but also fosters trust among consumers, thereby fostering a more transparent and trustworthy digital marketplace. Furthermore, the proposed methodology holds promise as a scalable and adaptable solution to the challenges posed by fake reviews, offering a robust framework for addressing the evolving tactics employed by malicious actors. By contributing to the integrity and trustworthiness of online review systems, our approach plays a pivotal role in safeguarding consumer interests and promoting fair competition among businesses in the digital era.

X. REFERENCES

- [1] S. Akshara, S. Shiva, S. Kubireddy, T. Arun and V. L. Kanthety, "A Small Comparative Study of Machine Learning Algorithms in the Detection of Fake Reviews of Amazon Products," 2023 6th International Conference on Contemporary Computing and Informatics (IC3I), Gautam Buddha Nagar, India, 2023
- [2] S. Zabeen, A. Hasan, M. F. Islam, M. S. Hossain and A. A. Rasel, "Robust Fake Review Detection Using Uncertainty-Aware LSTM and BERT," 2023 IEEE 15th International Conference on Computational Intelligence and Communication Networks (CICN), Bangkok, Thailand, 2023, pp. 786-791, doi: 10.1109/CICN59264.2023.10402342.
- [3] S. Ashraf, F. Rehman, H. Sharif, H. Kirn, H. Arshad and H. Manzoor, "Fake Reviews Classification using Deep Learning," 2023 International Multi-disciplinary Conference in Emerging Research Trends (IMCERT), Karachi, Pakistan, 2023
- [4] J. Fontanarava, G. Pasi and M. Viviani, "Feature Analysis for Fake Review Detection through Supervised Classification," 2017 IEEE International Conference on Data Science and Advanced Analytics (DSAA), Tokyo, Japan, 2017
- [5] J. Wang, H. Kan, F. Meng, Q. Mu, G. Shi and X. Xiao, "Fake Review Detection Based on Multiple Feature Fusion and Rolling Collaborative Training," in IEEE Access, vol. 8, pp. 182625-182639, 2020
- [6] R. Agarwal and D. K. Sharma, "Detecting Fake Reviews using Machine learning techniques: a survey," 2022 2nd International Conference on Advance Computing and Innovative Technologies in Engineering (ICACITE), Greater Noida, India, 2022
- [7] C. Cao, S. Li, S. Yu and Z. Chen, "Fake Reviewer Group Detection in Online Review Systems," 2021 International Conference on Data Mining Workshops (ICDMW), Auckland, New Zealand, 2021
- [8] P. WANG, Y. LIN and J. CHAI, "Unmasking Deception: A Comparative Study of Tree-Based and Transformer-Based Models for Fake Review Detection on Yelp," 2023 IEEE International Conference on Systems, Man, and Cybernetics (SMC), Honolulu, Oahu, HI, USA, 2023
- [9] P. Kalaivani, V. D. Raj, R. Madhavan and A. P. Naveen Kumar, "Fake Review Detection using Naive Bayesian Classifier," 2023 International Conference on Sustainable Computing and Smart Systems (ICSCSS), Coimbatore, India, 2023
- [10] T. R. Sree and R. Tripathi, "Fake Review Detection using Evidential Classifier," 2023 Second International Conference on Advances in Computational Intelligence and Communication (ICACIC), Puducherry, India, 2023