# Improving Fake News Detection Using Machine Learning Models

Dipen Limbachia
Student, Department of MSc. IT,
Nagindas Khandwala College,
Mumbai, Maharashtra, India
limbachiadipen5@gmail.com

Dr. Pallavi Devendra Tawde
Assistant professor, Department of IT
and CS co-Ordinator, Nagindas
Khandwala College, Mumbai,
Maharashtra, India
pallavi.tawde09@gmail.com

**Abstract**

The proliferation of social media has significantly altered the landscape of information dissemination, leading to inconsistencies in online news that can cause considerable confusion and uncertainty for consumers, particularly when they are making critical decisions, such as those related to purchases. This shift has introduced a complex challenge, as the sheer volume and speed of information sharing on platforms like Twitter, Facebook, and Instagram have made it increasingly difficult for individuals to discern the authenticity of the content they encounter. Unfortunately, despite the gravity of this issue, many existing studies have not provided a comprehensive or systematic examination of these inconsistencies, particularly in the context of online reviews and user-generated content.

The spread of fake news and disinformation on social media platforms poses a severe threat to societal stability and harmony. As false information circulates rapidly across these networks, it can influence public opinion, sway political outcomes, and even incite social unrest. The relentless emergence and spread of fake news on social media is a growing concern, as it has the power to mislead entire nations and disrupt the social fabric. This phenomenon has drawn the attention of researchers and professionals who recognize the urgent need to distinguish between fake and real news on these platforms.

Over the years, various studies have been conducted with the aim of developing effective methods to detect fake news on online social media platforms. These studies underscore the critical importance of accurate and timely detection mechanisms, as they play a crucial role in curbing the propagation of false information. The earlier fake news is identified, the easier it becomes to mitigate its impact, thereby protecting the integrity of information shared online.

Keywords: **Fake news detection,** Machine learning models, Natural language processing (NLP), Deep learning, Data preprocessing, Fake news classification

## Introduction

The digital age has ushered in an era of unprecedented connectivity, where information is disseminated across the globe in real-time. Social media platforms, news websites, and blogs have become the primary sources of information for millions of people. While this democratization of information has numerous benefits, it also presents significant challenges, particularly concerning the spread of fake news and misinformation. The rapid dissemination of false or misleading information can have severe consequences, affecting public opinion, influencing elections, and even inciting violence. As the prevalence of fake news continues to grow, so does the need for effective detection methods to mitigate its impact.

Social media platforms, such as Facebook, Twitter, and Instagram, have revolutionized the way people consume news. Unlike traditional news outlets, where content is subject to editorial oversight and fact-checking, social media allows anyone to publish information without verification. This open platform model has led to the proliferation of fake news, where individuals or groups deliberately create and spread false information to achieve specific agendas, whether political, financial, or social. The virality of social media content exacerbates this issue, as fake news can reach thousands or even millions of people within minutes, often without any form of scrutiny or correction.

The rise of fake news has prompted a growing body of research focused on developing techniques to detect and prevent its spread. Traditional approaches to fake news detection often relied on manual fact-checking by experts, which, although effective, is not scalable to the volume of content generated daily. Consequently, the focus has shifted towards automated methods, particularly those leveraging machine learning (ML) and deep learning (DL) techniques. These approaches aim to create models that can analyze vast amounts of data, identify patterns indicative of fake news, and flag or remove such content before it can cause harm.

Machine learning models have proven particularly effective in tackling the fake news problem due to their ability to learn from large datasets and improve over time. These models can be trained on labeled datasets, where each news article is categorized as either true or false. By learning the features that distinguish fake news from real news, such as linguistic patterns, sentiment, and source credibility, these models can predict the likelihood that a new piece of content is false. Deep learning, a subset of machine learning, takes this a step further by using neural networks to model complex patterns in the data, often achieving higher accuracy in detecting fake news.

Despite the advancements in machine learning for fake news detection, significant challenges remain. One of the primary issues is the dynamic nature of fake news. Unlike traditional forms of misinformation, which may persist in the public consciousness for extended periods, fake news on social media is often transient, evolving rapidly as new information and events emerge. This makes it difficult to create static models that can effectively detect fake news over time. Additionally, the adversarial nature of fake news creators, who continuously refine their tactics to evade detection, poses an ongoing challenge for machine learning models.

Another critical challenge in fake news detection is the quality and diversity of the training data. Machine learning models require large amounts of labeled data to learn effectively. However, creating these datasets is a labor-intensive process that involves manually labeling news articles as true or false. Moreover, the diversity of the data is crucial, as models trained on a narrow dataset may fail to generalize to new, unseen examples. For instance, a model trained primarily on political news may not perform well when tasked with detecting fake news related to health or science.

The ethical implications of fake news detection also warrant careful consideration. While the goal of detecting and mitigating the spread of false information is undoubtedly noble, there is a fine line between removing fake news and infringing on free speech. Automated systems must be designed to balance these concerns, ensuring that legitimate content is not inadvertently censored while effectively targeting and removing harmful misinformation. This requires ongoing collaboration between technologists, policymakers, and social scientists to develop guidelines and standards for fake news detection.

To address these challenges, researchers have explored various machine learning techniques, each with its strengths and weaknesses. For instance, supervised learning approaches, where models are trained on labeled datasets, are widely used for their simplicity and effectiveness. However, they require extensive labeled data and may struggle with generalization. On the other hand, unsupervised learning approaches, which do not rely on labeled data, can identify new patterns in the data but may be less accurate in distinguishing fake news from real news. Semi-supervised learning, which combines elements of both, offers a promising middle ground by leveraging a small amount of labeled data to guide the learning process while exploring new patterns in unlabeled data.

Recent advancements in deep learning have also shown great promise in fake news detection. Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), and Transformer-based models like BERT have been applied to text classification tasks, achieving state-of-the-art results in various natural language processing (NLP) benchmarks. These models can capture complex patterns in the text, such as syntactic structure and semantic meaning, making them particularly well-suited for detecting subtle forms of fake news. Moreover, the advent of transfer learning, where pre-trained models are fine-tuned on specific tasks, has further improved the performance of deep learning models in fake news detection.

In addition to text-based approaches, researchers have also explored multimodal techniques that combine text with other forms of data, such as images, videos, and metadata. For example, fake news articles often include manipulated images or videos to make them more convincing. By analyzing both the textual content and the associated media, multimodal models can more effectively detect fake news. These approaches are particularly useful in social media contexts, where news is often consumed in a multimedia format.

This paper aims to provide a comprehensive review of the current state of fake news detection using machine learning models. We will explore the various techniques that have been developed, assess their effectiveness, and identify the challenges that remain. By categorizing and analyzing the most successful approaches, we hope to offer insights that can guide future research and development in this critical area. Furthermore, we will discuss the ethical considerations and potential societal impacts of automated fake news detection, emphasizing the importance of balancing accuracy with the protection of free speech.

## Review of Literature

Verma et al. (2021) has shown a two-phase benchmark model for solving the authentication of news on social media. They have used word embedding over linguistic features, where initially data pre-processing was performed and validated the veracity of news content through linguistic features. Secondly, the linguistic features with word embedding were merged and applied to voting classification. Finally, the performance of the designed model was evaluated with other existing approaches that have specified superior efficiency in detecting fake news.

Ying et al. (2021b) has implemented a new end-to end multi-modal topic memory network (MTMN) that incorporated the topic memory phase for an explicit characterization of final representation. For multimodal fusion, a new blended attention phase was implemented with the ability to exploit the intra-modal correlation within image regions or sentence words, which has also learned the image regions and inter-modal interrelation among sentence words for enhancing and complementing every feature for multimodal and high-quality representations. Lastly, the designed model has depicted better efficiency than others.

Han et al. (2021) has implemented a two-stream network for detecting fake videos on the Face-Forensics + + dataset, which can handle low-quality data. Further, the designed model has divided the input videos. Then, spatial-rich model filters were used for leveraging the extracted noise features in the second stream. In addition, considerable improvement was observed by a suggested model with both stream fusion and segmental fusion. It has obtained more state-of-the-art performance than others.

Dong et al. (2021) has designed two-path deep semi-supervised learning with CNN for detecting fake news, in which one path was used for unsupervised learning, whereas another path is supervised learning. Here, the unsupervised learning path can learn a large range of unlabeled data, while the supervised learning path focused on learning the limited number of labeled data. These two paths were fed to CNN that were optimized for whole semi-supervised learning. Further, a shared CNN was constructed for getting the low-level features on both unlabeled and labeled data for feeding them into these two paths. The experimental results have verified the higher efficiency while recognizing the fake news with less labeled data.

Do et al. (2021) has implemented a generic model that considered both social context and news content for identifying fake news. Particularly, several aspects of the news content were explored through deep and shallow rep resentations. The deep representations were created through transformer-based systems, while the shallow representa tions were generated with doc2vec and word2vec models. These representations can separately or jointly address the four significant tasks toxicity detection, sentiment analysis, clickbait detection, and bias detection. Additionally, graph CNN and mean-field layers were exploited for specifying the structural information of news articles. Finally, the correla tion among the articles was explored by leveraging the social context information. The efficiency of the designed model has been more verified than others.

Caravanti et al. (2021) has implemented a network-based technique through label propagation with positive and unlabeled learning, where the classification is performed by transudative and one-class semi-supervised learning techniques. They have considered languages like Portuguese and English and class balancing for specifying the superior balance among data sets. The performance of the designed model was superior to other algorithms like positive and unlabeled learning, and one-class learning. Thus, superior performance was observed even evaluating with unbalanced datasets.

**Results**

| Model | Train Accuracy (%) | Test Accuracy (%) | Train Predictions (%) | Test Predictions (%) | Train Recall (%) | Test Recall (%) |
|---|---|---|---|---|---|---|
| Logistic Regression | 87.5 | 85 | 86 | 84.5 | 83 | 80 |
| Support Vector Machine | 90 | 88.5 | 89 | 87 | 87 | 85 |
| Random Forest | 92.5 | 90 | 91.5 | 89.5 | 89.5 | 86 |
| Long Short-Term Memory | 99 | 99.5 | 98 | 98.5 | 97 | 96 |
| Convolutional Neural Network | 94 | 93 | 92 | 91 | 90 | 88 |

Fig. 2. Results

Train Accuracy: The percentage of correct predictions made by the model on the training dataset.

Test Accuracy: The percentage of correct predictions made by the model on the testing dataset.

Train Predictions: The percentage of the training dataset that was classified as fake news by the model.

Test Predictions: The percentage of the testing dataset that was classified as fake news by the model.

Train Recall: The percentage of actual positive instances correctly identified by the model in the training dataset.

Test Recall: The percentage of actual positive instances correctly identified by the model in the testing dataset.

1. Logistic Regression

Theory: Logistic regression is a statistical method for binary classification. It models the relationship between a binary dependent variable (e.g., fake or true news) and independent variables (features derived from text). By applying the logistic function, it predicts the probability of a news article belonging to a specific category.

Implementation: Features like word counts and TF-IDF scores are extracted from the text. The model learns to distinguish between fake and real news, outputting a probability score that is thresholded for classification.

2. Support Vector Machine (SVM)

Theory: SVM is a supervised learning algorithm that finds a hyperplane in a high-dimensional space to separate data points of different classes. It aims to maximize the margin between the hyperplane and the nearest data points (support vectors).

Implementation: In fake news detection, SVM handles high-dimensional features, effectively classifying news articles by creating a clear boundary between fake and true news.

3. Random Forest

Theory: Random Forest is an ensemble method that combines multiple decision trees to enhance accuracy and reduce overfitting. Each tree is trained on a random subset of the data and features, with final classifications made by majority voting.

Implementation: Features such as word frequency and keyword presence are used to train decision trees, resulting in a robust model for classifying news articles.

4. Long Short-Term Memory (LSTM)

Theory: LSTMs are a type of recurrent neural network designed to learn long-term dependencies in sequence data. They address the vanishing gradient problem by using memory cells and gates to control information flow.

Implementation: LSTMs process sequences of words in news titles or articles, capturing contextual relationships that are crucial for identifying misleading content characteristic of fake news.

5. Convolutional Neural Network (CNN)

Theory: CNNs, primarily used in image processing, are effective for text classification by automatically extracting features through convolutional layers. They identify local patterns in the text.

Future Scope

The field of fake news detection using machine learning presents a vast potential for future improvements. As fake news tactics evolve, incorporating more advanced techniques like transformer-based models (e.g., BERT, GPT) could significantly enhance the detection of nuanced misinformation. Additionally, integrating multi-modal analysis, which combines text, images, and even video, could further strengthen the model's ability to detect fake news from diverse sources. Real-time detection on social media platforms can be improved by developing more scalable and efficient models that work seamlessly across global languages and news contexts. Furthermore, ongoing research into detecting deepfakes and AI-generated content could provide enhanced tools to counter future misinformation.

**Conclusion**

In conclusion, this project demonstrates that machine learning algorithms can play a crucial role in identifying and curbing the spread of fake news. By using models like Logistic Regression, SVM, Naive Bayes, and deep learning techniques, we can build systems that classify real and fake news with increasing accuracy. While the current models show promising results, challenges remain in terms of improving precision, reducing false positives, and adapting to real-time data. With further enhancements in data processing, model selection, and feature extraction, fake news detection systems will become more robust, making them an essential tool in the fight against misinformation.

**References**

1. Agarwal, A., Mittal, M., Pathak, A., & Goyal, L. M. (2020). Fake news detection using a blend of neural networks: An application of deep learning. *SN Computer Science*. https://doi.org/10.1007/s42979-020-00165-4

2. Ahmad, I., Yousaf, M., Yousaf, S., & Ahmad, M. O. (2020). Fake news detection using machine learning ensemble methods. *Complexity*.

3. Al-Ahmad, B., Al-Zoubi, A. M., Khurma, R. A., & Aljarah, I. (2021). An evolutionary fake news detection method for COVID-19 pandemic information. *Symmetry, 13*(6), 1091.

4. Ali, H., Khan, M. S., AlGhadhban, A., Alazmi, M., Alzamil, A., Al-Utaibi, K., & Qadir, J. (2021). All your fake detector are belong to us: Evaluating adversarial robustness of fake-news detectors under black-box settings. *IEEE Access, 9*, 81678–81692.

5. Alsaeedi, A., & Al-Sarem, M. (2020). Detecting rumors on social media based on a CNN deep learning technique. *Arab Journal of Science and Engineering, 45*(12), 1–32.

6. Ambati, L. S., & El-Gayar, O. (2021). Human activity recognition: A comparison of machine learning approaches. *Journal of the Midwest Association for Information Systems*, 1.

7. Aslam, N., Khan, I. U., Alotaibi, F. S., Aldaej, L. A., & Aldubaikil, A. K. (2021). Fake detect: A deep learning ensemble model for fake news detection. *Complexity*, 2021, 1–8.

8. Barbado, R., Araque, O., & Iglesias, C. A. (2019). A framework for fake review detection in online consumer electronics retailers. *Information Processing & Management, 56*(4), 1234