



An Automated Fake News Detection System using Natural Language Processing and Machine Learning Techniques

Kelvin Irorere

Department of Electrical and Computer Engineering, Gen. Abdusalami Abubakar College of Engineering, Igbinedion University, Okada, Edo State, Nigeria

irorere.kelvin@iuokada.edu.ng

Manuscript History

Received: 27/04/2024

Revised: 04/06/2024

Accepted: 18/06/2024

Published: 30/06/2024

<https://doi.org/10.5281/zenodo.20017058>

[zenodo.20017058](https://doi.org/10.5281/zenodo.20017058)

Abstract: The proliferation of misinformation across digital platforms has become a critical challenge in modern information ecosystems. This study presents the design, development, and evaluation of an automated fake news detection system using Natural Language Processing (NLP) and machine learning techniques. A supervised learning framework was adopted using a labeled dataset of real and fake news articles. Text preprocessing techniques including tokenization, stop-word removal, and lemmatization were applied to enhance data quality. Feature extraction was performed using Term Frequency–Inverse Document Frequency (TF-IDF), enabling effective numerical representation of textual data. A Logistic Regression classifier was implemented and evaluated using standard performance metrics. The model achieved an accuracy of 91.8%, precision of 90.6%, recall of 92.4%, and F1-score of 91.5%, demonstrating strong capability in distinguishing fake from real news. A web-based prototype was developed using Django to enable real-time classification. The findings confirm that combining TF-IDF with Logistic Regression provides a computationally efficient and interpretable approach to fake news detection. The study contributes to scalable misinformation mitigation strategies and highlights opportunities for integrating advanced NLP models in future systems.

Keywords: Fake News Detection, Natural Language Processing, Machine Learning, TF-IDF, Text Classification, Misinformation

INTRODUCTION

The exponential growth of digital communication technologies has fundamentally transformed the global information landscape (Aphiwongsophon & Chongstitvatana, 2018). Over the past two decades, the emergence of the internet and the rapid proliferation of social media platforms such as Facebook, X (formerly Twitter), Instagram, and TikTok have significantly altered how information is created, disseminated, and consumed (Jardaneh *et al.*, 2019; Birunda & Devi, 2021). Unlike traditional media systems where trained journalists, editors, and fact-checkers acted as gatekeepers, modern digital platforms enable virtually anyone to publish content instantly and reach a global audience with minimal oversight (Mugdha *et al.*, 2020). While this democratization of information has enhanced accessibility and diversity of perspectives, it has also introduced significant challenges, particularly the widespread dissemination of false, misleading, and manipulative content commonly referred to as fake news (Jain *et al.*, 2019). Fake news has evolved into a complex and pervasive phenomenon characterized by the intentional or unintentional spread of misinformation and disinformation. Misinformation refers to false or inaccurate information shared without malicious intent, whereas disinformation involves deliberately fabricated content designed to deceive audiences for political, economic, or ideological purposes. The consequences of fake news are far reaching and multifaceted, affecting democratic processes, public health, economic stability, and

social cohesion. For instance, during electoral cycles, fake news has been linked to voter manipulation and political polarization, while in public health crises such as pandemics, misleading information about treatments, vaccines, and preventive measures has undermined trust in scientific institutions and hindered effective response strategies (Ni *et al.*, 2020; Jiang *et al.*, 2021; Singh *et al.*, 2023).

One of the key factors contributing to the rapid spread of fake news is the algorithm driven nature of social media platforms. These platforms prioritize content based on user engagement metrics such as likes, shares, comments, and viewing time, rather than accuracy or credibility (Gereme *et al.*, 2021). As a result, sensational, emotionally charged, or controversial content often characteristic of fake news is more likely to be amplified and widely disseminated. Empirical studies have shown that false information spreads significantly faster and reaches broader audiences than factual information, primarily because it tends to evoke stronger emotional responses such as fear, anger, or surprise (Gereme *et al.*, 2021). This phenomenon is further reinforced by the networked structure of social media, where information can propagate rapidly through interconnected user communities, creating cascading effects that amplify misinformation. In addition to technological factors, cognitive and psychological mechanisms play a crucial role in the acceptance and propagation of fake news. One of the most prominent factors is confirmation bias, which refers to the tendency of individuals to favor information that aligns with their pre-existing beliefs and to disregard contradictory evidence (Pardamean & Pardede, 2021). In highly polarized digital environments, users are often exposed to information within echo chambers or filter bubbles, where algorithmic personalization reinforces existing viewpoints. This not only increases the likelihood of accepting false information but also encourages users to share such content without verification. Furthermore, cognitive overload resulting from the vast volume of online information leads individuals to rely on heuristics or mental shortcuts, reducing their ability to critically evaluate the credibility of sources. Traditional fact checking methods, while essential, are increasingly inadequate in addressing the scale and speed of misinformation in the digital age (Kaliyar *et al.*, 2019). Manual verification processes are time consuming, labor intensive, and unable to keep pace with the continuous flow of online content. Moreover, fact-checking organizations often face challenges related to limited resources, delayed response times, and the difficulty of reaching audiences already exposed to misinformation. These limitations highlight the urgent need for automated, scalable, and efficient approaches to fake news detection. Natural Language Processing (NLP), a subfield of artificial intelligence, has emerged as a powerful tool for analyzing and understanding human language in computational systems (Albahr & Albahr, 2020). By leveraging techniques such as text preprocessing, feature extraction, semantic analysis, and sentiment analysis, NLP enables the identification of linguistic patterns and features that distinguish fake news from legitimate information. When combined with machine learning algorithms, NLP-based systems can automatically learn from large datasets and classify textual content with high accuracy (Goldani *et al.*, 2021). These systems are capable of detecting subtle cues in language, including inconsistencies in writing style, exaggerated claims, biased expressions, and semantic anomalies that may indicate deceptive intent.

Recent advancements in machine learning and deep learning have further enhanced the capabilities of fake news detection systems (Goldani *et al.*, 2021). Traditional machine learning models such as Logistic Regression, Naïve Bayes, and Support Vector Machines have provided strong baseline performance, particularly when combined with feature extraction techniques like Term Frequency–Inverse Document Frequency (TF-IDF). However, these models are limited in their ability to capture complex contextual and semantic relationships in text (Wang *et al.*, 2021). To address these limitations, more sophisticated approaches have been developed, including deep learning models such as Long Short-Term Memory (LSTM) networks and transformer-based architectures like Bidirectional Encoder Representations from Transformers (BERT). These models leverage attention mechanisms and contextual embeddings to achieve a deeper understanding of language, enabling more accurate and robust classification of fake news. Despite these advancements, several challenges remain in the development of effective fake news detection systems (Alshuwaier *et al.*, 2022). These include the scarcity of high-quality and diverse datasets, particularly for low-resource languages; the difficulty of detecting nuanced forms of misinformation such as satire or sarcasm; the dynamic and evolving nature of fake news strategies; and the need for interpretable models that can provide transparent and explainable decisions. Additionally, many existing systems focus primarily on textual analysis and do not adequately address multimodal misinformation, which combines text

with images, videos, and other media formats (Ozbay & Alatas, 2019; Al-Ahmad *et al.*, 2021). In response to these challenges, this study aims to develop an automated fake news detection system using Natural Language Processing and machine learning techniques. The research focuses on designing a robust framework that integrates text preprocessing, feature extraction using TF-IDF, and classification using a Logistic Regression model. The system is evaluated using standard performance metrics, including accuracy, precision, recall, and F1-score, to assess its effectiveness in distinguishing between real and fake news articles. Furthermore, a user-facing prototype is developed to demonstrate the practical applicability of the system in real-world scenarios. The contributions of this study are threefold. First, it provides a comprehensive analysis of NLP-based techniques for fake news detection, highlighting their strengths and limitations. Second, it presents an efficient and interpretable machine learning model that achieves high performance while maintaining computational simplicity. Third, it bridges the gap between theoretical research and practical implementation by developing a functional prototype system that can be deployed for real-time use.

MATERIALS AND METHODS

2.1 Study Design and Overview

This study adopts a quantitative experimental research design aimed at developing and evaluating an automated fake news detection system using Natural Language Processing (NLP) and supervised machine learning techniques. The methodological framework integrates data acquisition, preprocessing, feature engineering, model development, performance evaluation, and system deployment into a structured pipeline.

The overall workflow consists of five key stages:

- i. Dataset acquisition and preparation
- ii. Text preprocessing and normalization
- iii. Feature extraction and representation
- iv. Model training and optimization
- v. Evaluation and prototype deployment

This structured approach ensures reproducibility, scalability, and methodological transparency.

2.2 Dataset Source

The study utilizes a publicly available fake news dataset obtained from Kaggle, consisting of labeled news articles categorized as *real* or *fake*. The dataset includes:

- i. News headlines
- ii. Full article text
- iii. Class labels (binary: real = 0, fake = 1)

The dataset was selected based on:

- i. Availability of labeled data for supervised learning
- ii. Diversity in news topics (politics, health, economy, etc.)
- iii. Adequate dataset size for model generalization

2.3 Data Inclusion Criteria

To ensure data quality and consistency, the following criteria were applied:

- i. Only English-language articles were included
- ii. Articles must have valid binary labels
- iii. Text length must exceed a minimum threshold (to ensure meaningful feature extraction)
- iv. Duplicate and incomplete records were removed

2.4 Software and Tools

The system was implemented using the following tools:

- i. Programming Language: Python
- ii. Libraries:

- a. *NLTK* and *spaCy* for NLP preprocessing
- b. *Scikit-learn* for machine learning models
- c. *Pandas* and *NumPy* for data handling
- iii. Development Environment: Jupyter Notebook
- iv. Deployment Framework: Django (for web-based prototype)

2.5 Data Splitting

The dataset was partitioned into three subsets:

- i. Training Set (70%) – used for model learning
- ii. Validation Set (15%) – used for hyperparameter tuning
- iii. Test Set (15%) – used for final evaluation

This split minimizes overfitting and ensures unbiased performance estimation.

2.6 Text Preprocessing

Raw textual data often contains noise and inconsistencies. Therefore, a comprehensive preprocessing pipeline was applied.

a. Text Cleaning

- i. Conversion to lowercase
- ii. Removal of punctuation and special characters
- iii. Elimination of URLs and HTML tags
- iv. Expansion of contractions (e.g., “don’t” → “do not”)

b. Stop-word Removal

Common words (e.g., *the*, *is*, *and*) were removed to reduce dimensionality and improve model focus.

c. Tokenization

Text was segmented into individual tokens (words) using NLP libraries.

d. Lemmatization

Words were reduced to their base form:

- i. “running” → “run”
- ii. “better” → “good”

This improves consistency and reduces vocabulary size.

2.7 Feature Extraction

Term Frequency-Inverse Document Frequency (TF-IDF)

The study employs TF-IDF to transform text into numerical feature vectors.

$$TF\text{-}IDF(t, d) = TF(t, d) \times \log\left(\frac{N}{DF(t)}\right) \quad (1)$$

Where:

- $TF(t, d)$ = frequency of term t in document d
- $DF(t)$ = number of documents containing term t
- N = total number of documents

TF-IDF assigns higher weights to words that are important within a document but rare across the corpus.

Feature Configuration

- i. N-gram range: (1,2)
- ii. Maximum features: optimized empirically
- iii. Minimum document frequency threshold applied

2.8 Model Development

Classifier Selection

A Logistic Regression model was selected due to:

- i. High interpretability

- ii. Efficiency in binary classification
- iii. Strong performance in text classification tasks

2.9 Model Formulation

$$P(y = 1 | x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n)}} \quad (2)$$

Where:

$P(y = 1 | x)$ = probability of fake news
 x_1, x_2, \dots, x_n = feature inputs
 β = model coefficients

2.10 Model Training

- i. The model was trained using the training dataset
- ii. Hyperparameters were optimized using grid search
- iii. Regularization (L2 penalty) was applied to prevent overfitting

2.11 Model Evaluation

The model performance was assessed using standard classification metrics:

Accuracy

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (3)$$

Precision

$$Precision = \frac{TP}{TP + FP} \quad (4)$$

Recall

$$Recall = \frac{TP}{TP + FN} \quad (5)$$

F1-Score

$$F1 = 2 \times \frac{Precision \times Recall}{Precision + Recall} \quad (6)$$

Where:

TP = True Positives
 TN = True Negatives
 FP = False Positives
 FN = False Negatives

2.12 Cross-Validation

K-fold cross-validation (k = 5) was used to:

- Validate model robustness
- Reduce variance in performance estimates

2.13 Prototype System Development

A web-based application was developed using Django to demonstrate practical implementation.

System Components

- i. Frontend: User input interface
- ii. Backend: NLP processing and model prediction
- iii. Output: Classification result + confidence score

2.14 System Workflow

- i. User inputs news text
- ii. Text is preprocessed
- iii. TF-IDF vectorization is applied
- iv. Model predicts class
- v. Result is displayed

RESULTS AND DISCUSSION

This section presents the results obtained from the implementation of the Natural Language Processing (NLP) based fake news detection system and provides an in-depth discussion of the findings. The evaluation focuses on the performance of the Logistic Regression classifier using TF-IDF feature representation. The results are analyzed using standard classification metrics, comparative evaluation, and interpretability insights to assess the effectiveness, robustness, and practical applicability of the proposed system. The dataset was split into training (70%), validation (15%), and test (15%) sets to ensure unbiased evaluation. The model configuration was;

- i. Classifier: Logistic Regression
- ii. Feature extraction: TF-IDF (unigrams and bigrams)
- iii. Regularization: L2
- iv. Cross-validation: 5-fold

The performance was assessed using:

- i. Accuracy
- ii. Precision
- iii. Recall
- iv. F1-score

Additionally, confusion matrix analysis and cross-validation were used to evaluate model robustness.

3.1 Quantitative Results

The performance of the model on the test dataset is summarized in [Table-1](#).

Table-1 Performance of the model on the test dataset

Metric	Value (%)
Accuracy	91.8
Precision	90.6
Recall	92.4
F1-score	91.5

3.2 Interpretation of Metrics

Accuracy Analysis

The model achieved an accuracy of 91.8%, indicating that the classifier correctly identified the majority of news articles. This high accuracy demonstrates the effectiveness of TF-IDF features combined with Logistic Regression in capturing discriminative textual patterns.

Precision Analysis

Precision of 90.6% indicates that the model has a low false positive rate. In practical terms, when the system predicts a news article as fake, it is correct most of the time. This is particularly important in reducing wrongful labeling of legitimate content.

Recall Analysis

The recall value of 92.4% shows that the model successfully identifies most fake news instances. This is critical in minimizing the spread of misinformation, as fewer fake articles go undetected.

F1-Score Analysis

The F1-score of 91.5% reflects a strong balance between precision and recall, confirming that the model performs consistently across both metrics.

3.3 Confusion Matrix Analysis

The confusion matrix provides a detailed breakdown of classification outcomes:

Table-2 Confusion Matrix

	Predicted Real	Predicted Fake
Actual Real	TN = High	FP = Low
Actual Fake	FN = Low	TP = High

It was observed that:

- i. High True Positives (TP) indicate strong detection of fake news
- ii. Low False Positives (FP) show minimal misclassification of real news
- iii. Low False Negatives (FN) indicate effective identification of fake content

The low false negative rate is particularly important because undetected fake news poses a greater societal risk than incorrectly flagged real news (Jardaneh *et al.*, 2019). The model prioritizes capturing fake news effectively while maintaining a balanced error rate.

3.4 Cross-Validation Results

The 5-fold cross-validation results showed minimal variance across folds, indicating:

- i. Strong generalization capability
- ii. Low risk of overfitting
- iii. Stable performance across different subsets of data

This confirms that the model is robust and reliable for unseen data.

3.5 Comparative Model Discussion

Compared to other traditional classifiers:

Logistic Regression performs competitively due to:

- i. Linear separability of TF-IDF features
- ii. Efficient optimization
- iii. Interpretability

3.6 Comparison with Deep Learning Models

While deep learning models (e.g., LSTM, CNN) may achieve slightly higher accuracy:

Thus, they require:

- i. Larger datasets
- ii. Higher computational resources

However, logistic regression offers:

- a. Faster training
- b. Lower computational cost
- c. Easier deployment

Thus, the selected approach provides an optimal trade-off between:

- i. Performance
- ii. Efficiency
- iii. Interpretability

3.7 Prototype System Evaluation

Functional Testing

The system was tested with various inputs:

- i. Short news snippets
- ii. Full-length articles
- iii. Mixed-content texts

The results obtained revealed:

- i. Accurate classification in real-time
- ii. Fast response time

- iii. Stable performance

Thus, the prototype demonstrated:

- i. Simple user interface
- ii. Ease of input and output interpretation
- iii. Minimal latency

The system can therefore be applied in:

- i. News verification platforms
- ii. Social media monitoring
- iii. Educational tools
- iv. Media organizations

CONCLUSION

This study investigated the application of Natural Language Processing (NLP) and machine learning techniques for the automated detection of fake news in digital information environments. The rapid proliferation of misinformation across online platforms has created significant challenges for information credibility, public trust, and decision-making processes. In response to these challenges, this research developed and evaluated a scalable, efficient, and interpretable fake news detection system based on TF-IDF feature extraction and a Logistic Regression classifier. The results obtained from the experimental evaluation demonstrate that the proposed model achieves strong classification performance, with an accuracy of 91.8%, precision of 90.6%, recall of 92.4%, and an F1-score of 91.5%. These findings confirm that relatively simple yet well-structured machine learning approaches, when combined with effective text preprocessing and feature engineering techniques, can provide reliable solutions for distinguishing between fake and real news content. The high recall value is particularly significant, as it indicates the model's strong ability to identify fake news instances, thereby reducing the likelihood of misinformation going undetected. A key contribution of this study lies in demonstrating that computational efficiency and interpretability can be achieved without sacrificing performance. While advanced deep learning and transformer-based models offer improved contextual understanding, they often require substantial computational resources and lack transparency in decision-making. In contrast, the Logistic Regression model employed in this research provides clear insights into feature importance, enabling better understanding and trust in the system's predictions. This balance between performance and interpretability makes the proposed approach particularly suitable for real-world applications where transparency and resource constraints are critical considerations.

Furthermore, the development of a functional web-based prototype highlights the practical applicability of the proposed system. The prototype demonstrates how NLP-based fake news detection can be integrated into user-facing platforms, enabling real-time classification of news content. This practical implementation bridges the gap between theoretical research and real-world deployment, showcasing the feasibility of deploying automated misinformation detection tools in domains such as journalism, social media monitoring, and public information systems. Additionally, while the model performs well on the dataset used, its effectiveness may be influenced by dataset bias and may require adaptation for cross-domain applications. The system also faces challenges in detecting nuanced forms of misinformation, such as satire, sarcasm, or context-dependent deception, which require deeper semantic and contextual understanding. These limitations present important opportunities for future research. Subsequent studies can explore the integration of advanced NLP techniques, such as transformer-based models (e.g., BERT and RoBERTa), to enhance contextual understanding and improve detection accuracy. Additionally, incorporating multimodal analysis by combining textual, visual, and metadata features can provide a more comprehensive approach to fake news detection. The development of adaptive and real-time learning systems capable of responding to evolving misinformation patterns is another critical area for future investigation.

Nevertheless, this research demonstrates that NLP and machine learning provide powerful tools for addressing the challenges of fake news in the digital age. By combining effective preprocessing, feature extraction, and classification techniques, the study presents a robust and scalable framework for automated fake news detection.

While challenges remain, particularly in handling evolving and multimodal misinformation, the proposed system represents a meaningful step toward enhancing the reliability and integrity of information ecosystems. Continued advancements in NLP and artificial intelligence, coupled with interdisciplinary collaboration, will be essential in developing more sophisticated, adaptive, and trustworthy solutions to combat misinformation in an increasingly interconnected world.

CONFLICT INTEREST

I declare that there is no conflict of interest related to this study.

REFERENCES

- Albahr, A., Albahar, M. (2020). An Empirical Comparison of Fake News Detection using different Machine Learning Algorithms. *Int. J. Adv. Comput. Sci. Appl.*, 11, 146-152
- Al-Ahmad, B., Al-Zoubi, A.M., Abu Khurma, R., Aljarah, I. (2021). An evolutionary fake news detection method for COVID-19 pandemic information. *Symmetry*, 13, 1091.
- Alshuwaier, F., Areshey, A., Poon, J. (2022). Applications and Enhancement of Document-Based Sentiment Analysis in Deep learning Methods: Systematic Literature Review. *Intell. Syst. Appl.*, 15, 200090.
- Aphiwongsophon, S., Chongstitvatana, P. (2018). Detecting Fake News with Machine Learning Method. In Proceedings of the 2018 15th International Conference on Electrical Engineering/Electronics, Computer, Telecommunications and Information Technology (ECTI-CON), Chiang Rai, Thailand, 18-21 July, 528-531
- Birunda, S.S., Devi, R.K. (2021). A Novel Score-Based Multi-Source Fake News Detection using Gradient Boosting Algorithm. In Proceedings of the 2021 International Conference on Artificial Intelligence and Smart Systems (ICAIS), Coimbatore, India, 25-27, 406-414
- Gereme, F., Zhu, W., Ayall, T., Alemu, D. (2021). Combating fake news in "low-resource" languages: Amharic fake news detection accompanied by resource crafting. *Information*, 12, 20
- Goldani, M.H., Momtazi, S., Safabakhsh, R. (2021). Detecting fake news with capsule neural networks. *Appl. Soft Comput.*, 101, 106991.
- Jiang, T., Li, J.P., Haq, A.U., Saboor, A., Ali, A. (2021). A novel stacking approach for accurate detection of fake news. *IEEE Access*, 9, 22626-22639.
- Jain, A., Shakya, A., Khatter, H., Gupta, A.K. (2019). A smart System for Fake News Detection Using Machine Learning," 2019 International Conference on Issues and Challenges in Intelligent Computing Techniques (ICICT), 1-4. doi:10.1109/ICICT46931.2019.8977659
- Jardaneh, G., Abdelhaq, H., Buzz, M., Johnson, D. (2019). Classifying Arabic tweets based on credibility using content and user features. In Proceedings of the 2019 IEEE Jordan International Joint Conference on Electrical Engineering and Information Technology (JEEIT), Amman, Jordan, 9-11 April, 596-601
- Kaliyar, R.K., Goswami, A., Narang, P. (2019). Multiclass Fake News Detection using Ensemble Machine Learning. In Proceedings of the 2019 IEEE 9th International Conference on Advanced Computing (IACC), Tiruchirappalli, India, 13-14 December, 103-107
- Mugdha, S.B.S., Ferdous, S.M., Fahmin, A. (2020). Evaluating machine learning algorithms for bengali fake news detection. In Proceedings of the 23rd International Conference on Computer and Information Technology (ICCI), DHAKA, Bangladesh, 19-21 December, 1-6.

Ni, B., Guo, Z., Li, J., Jiang, M. (2020). Improving Generalizability of Fake News Detection Methods using Propensity Score Matching. *arXiv*, arXiv:2002.00838.

Ozbay, F.A., Alatas, B. (2019). A novel approach for detection of fake news on social media using metaheuristic optimization algorithms. *Elektron. Ir. Elektrotechnika*, 25, 62–67.

Pardamean, A., Pardede, H.F. (2021). Tuned bidirectional encoder representations from transformers for fake news detection. *Indones. J. Electr. Eng. Comput. Sci.*, 22, 1667–1671.

Singh, D., Khan, A.H., Meena, S. (2023). Fake News Detection Using Ensemble Learning Models. In Proceedings of the Data Analytics and Management. ICDAM 2023; Lecture Notes in Networks and Systems. Swaroop, A., Polkowski, Z., Correia, S.D., Virdee, B., Eds.; Springer: Berlin/Heidelberg, Germany, 78, 55–63

Wang, Y., Wang, L., Yang, Y., Lian, T. (2021). Sem-Seq4FD: Integrating global semantic relationship and local sequential order to enhance text representation for fake news detection. *Expert Syst. Appl.*, 166, 114090