



# Fake News Detection Using Machine Learning Techniques

**DR. Orchu Aruna**

Associate Professor,

Vasireddy Venkatadri Institute of Technology,  
Nambur

**Koncha Manoj**

Student,

Vasireddy Venkatadri Institute of Technology,  
Nambur

**Lam Venkata Krishna**

Student,

Vasireddy Venkatadri Institute of Technology,  
Nambur

**Maddula Yasaswini**

Student,

Vasireddy Venkatadri Institute of Technology,  
Nambur

**Kolukula Eswar**

Student,

Vasireddy Venkatadri Institute of Technology,  
Nambur

**Abstract** - Fake news has become a significant issue in today's digital age, spreading misinformation and influencing public opinion. This project aims to develop a fake news prediction using machine learning techniques to identify and classify news articles or statements as either fake or genuine. The system follows a step-by-step process, starting with data collection from reliable sources and preprocessing to clean and prepare the dataset. Feature extraction techniques like term frequency-inverse document frequency (TF-IDF) are applied to convert the textual data into numerical representations. Various machine learning algorithms, such as logistic regression, neural networks, are explored for classification. The chosen model is trained using labelled data, and its performance is evaluated using metrics like accuracy. Once a satisfactory model is obtained, it is deployed in a production environment, enabling users to input news articles for authenticity predictions. Continuous monitoring, maintenance, and updates are essential to ensure the system remains effective in identifying emerging forms of fake news. This project contributes to combating the spread of fake news by leveraging the power of machine learning to distinguish between trustworthy and deceptive information sources.

**Keywords:** Natural Language Processing, TfidfVectorizer, Logistic Regression, Stemming, RegularExpression, stopwords, fake news

## Introduction

The proliferation of false or misleading news stories, commonly referred to as "fake news", has become a serious concern in the digital age. The wide reach and fast dissemination enabled by social media make it easy for inaccurate information to spread quickly without adequate verification (Mustafaraj & Metaxas, 2017). Recent studies have shown that fake news articles often gain more traction than factual reporting, influencing public discourse and opinion (Vosoughi et al., 2018). This underscores the need for automated fake news detection systems to combat the spread of misinformation.

Fake news is frequently produced with the objective to mislead or control viewers for sensationalist, political, financial, or ideological reasons. Fake news has grown more widely as a result of the emergence of digital media and the simplicity with which it can be shared online. Its impact is far-reaching, influencing public opinion, shaping perceptions, and sometimes even causing social and political upheavals. Fake news can take different forms, including false headlines, manipulated images or videos, misleading content, conspiracy theories, and biased reporting. Fake news is created and spread by taking advantage of holes in our information ecosystem. Misinformation spreads because of things like confirmation bias, which makes people accept information that supports their opinions. Furthermore, in the digital era, false narratives spread quickly because of the

ease with which information can now be verified and how quickly it can flow. A diverse strategy is needed to tackle the problem of fake news, including the use of technology, media literacy campaigns, fact-checking activities, and critical thinking abilities. Though the fight against disinformation is still ongoing, numerous organizations, fact-checkers, and digital businesses have developed tools and ways to detect, validate, and combat fake news.

Recent advances in machine learning techniques present an opportunity to develop predictive models that can identify and flag potential fake news with high accuracy. While prior work has applied machine learning for fake news classification, there remain challenges in extracting meaningful features from textual data and improving generalizability across different news domains (Singhania et al., 2017). This project aims to build an effective fake news classification model using a combination of machine learning algorithms and natural language processing techniques for feature engineering. The key objectives are: 1. Compile a robust labeled dataset of verified real and fake news articles. 2. Explore different feature extraction methods like TF-IDF to convert text into informative numerical representations. 3. Evaluate machine learning models including logistic regression, neural networks to classify news articles. 4. Select the best performing model and optimize its hyperparameters. 5. Deploy the model in a web application for real-time fake news prediction.

## Literature Review

Fake news, referring to false or fabricated information disguised as news, has become a serious concern with the rise of social media. The rapid spread of misinformation enabled by online platforms poses risks to societal discourse, public opinion, and policymaking. Developing automated fake news detection systems is therefore an emerging research priority. However, this remains challenging due to the evolving tactics used to generate fake content. Related Work - Existing computational approaches for fake news detection primarily rely on machine learning using news content and metadata as features. Kong et al. (2020) applied natural language processing methods like tokenization, lemmatization, and TF-IDF to convert news titles and text into vector representations. They trained convolutional and recurrent neural networks using these features to classify fake news, achieving up to 83% accuracy.

Baarir and Djeflal (2021) also extracted various stylistic, semantic and user-based features, evaluated using a support vector machine classifier. Their approach combining n-grams, part-of-speech tags, source details, and propagation patterns achieved 87% accuracy on a curated real/fake news dataset. However, deep neural networks were not explored. A key challenge highlighted in these works is the lack of labelled fake news data spanning diverse topics and contexts. Ahmed et al. (2017) noted most datasets comprise political news and are limited in size. Ruchansky et al. (2017) proposed a hybrid CSI model combining labelled news content with unlabeled user response data to improve classification accuracy. Proposed Work - This project aims to develop an advanced fake news detection system using deep neural networks. The objectives are:

Incorporate unlabeled news data through pre-training

Compare performance with machine learning models using engineered features.

Evaluate model robustness across diverse news topics and sources.:

## Methodology

### A. Proposed Methodology

The proposed approach is implemented in the following phases.

#### 1) Data Loading :

- The data sets are taken from the Kaggle website
- A total of two datasets are taken.
- In which one dataset consists of only fake news-related data while the other dataset consists of only real news.
- Each dataset consists total of four fields as title, text, subject, date
- We have a total of 21 thousand rows of real news datasets and 23 thousand rows of fake news dataset
- The dataset is loaded as pandas DataFrame.
- The Fake news is labeled as 0 and Real news is labeled as 1.
- The Data is displayed through tables using pandas DataFrame object.
- We merge both the real and fake news datasets and reset their index values for performing operations on each row.
- For manual testing we take out limited rows out of the dataset.

#### 2) Data Pre-processing Phase:

Data preprocessing is a fundamental stage in the data analysis process that involves transforming raw data into a clean, structured format suitable for analysis or modelling. This stage is crucial as it directly impacts the quality and effectiveness of subsequent analytical tasks.

This process includes following steps

- Drop the columns which are not needed.
- Check Whether the dataset consists of null values or not
- Convert the text to lowercase
- Remove the text that is enclosed in square brackets.
- Replace the numbers with a space.
- Remove the URLs present in the text.
- Remove html tags in the text which does not give any meaning to the text.
- Remove the punctuations and all the special character in the text
- And we finally obtain the preprocessed text.
- The preprocessed text is vectorized into feature representations that can be used by machine learning models. Two approaches are evaluated:
  - TF-IDF (Term Frequency - Inverse Document Frequency): Calculates word importance based on frequency in the document and corpus. Scikit-Learn's Tfidf Vectorizer is used with n-gram range (1,3) to vectorize text.

- The text vectors are combined with any additional metadata like source, author, publish date etc. to create the final input feature matrix.
- The matrix is transformed to be compatible with selected machine learning models. Examples - one-hot encoding for deep learning, scaling for Logistic Regression.
- The processed datasets are ready for training supervised learning models to detect fake news.
- This covers the major data preprocessing steps including cleaning, formatting, text vectorization and feature engineering to prepare the raw news data for feeding into machine learning algorithms.

**3) Training and Testing Phase:**

- The dataset is split into train dataset and test dataset.
- Here on splitting we have taken seventy five percent for train dataset and twenty five percent for the test dataset.
- The spitted dataset with their corresponding labels are now fed to the logistic regression algorithm.
- The Logistic Regression algorithm fits the data and build the logistic Regression model for training and testing the data.
- After training and testing the data with the model. Now our model is ready for prediction.
- We get approximately 98 percent of accuracy score.
- Now we can build the prediction system based on the results given by our model.

**B. Algorithms:**

**Logistic Regression:** Logistic regression is a simple linear classification algorithm that will be leveraged for predicting if a given news article is fake or real. It calculates the probability that the input text belongs to the ‘fake’ class based on the logistic/sigmoid function:

$$P(fake|x) = 1 / (1 + e^{-(\beta T x)})$$

Where x is the TF-IDF feature vector for the input news text and β are the model weights learned during training. The probability score is threshold to classify the news as fake or real.

**Model Training:**

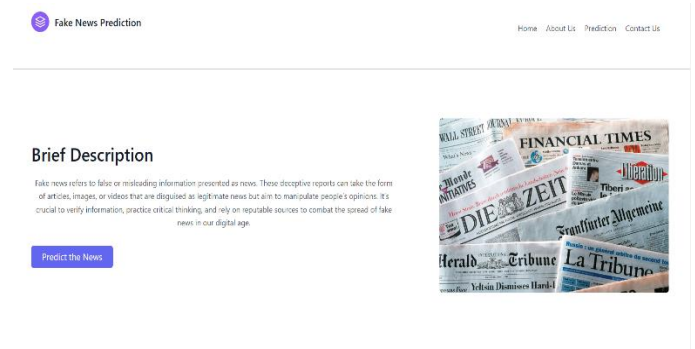
- β weights are randomly initialized
- Training data with TF-IDF vectors and fake/real labels is fed to the model
- Loss is calculated using cross-entropy between predicted and true labels
- β is updated to minimize loss using gradient descent optimization
- L2 regularization is used to prevent overfitting

**Model Evaluation:**

- Model is evaluated on unseen test set
- Classification accuracy, precision, recall metrics are tracked
- AUC-ROC curve shows model’s ability to distinguish fake vs real news
- Errors are analyzed to improve model performance

The logistic regression algorithm provides an interpretable baseline model before exploring more complex neural network architectures. Its performance on imbalanced datasets is also relatively robust. The probability outputs allow ranking fake news likelihood rather than binary classification.

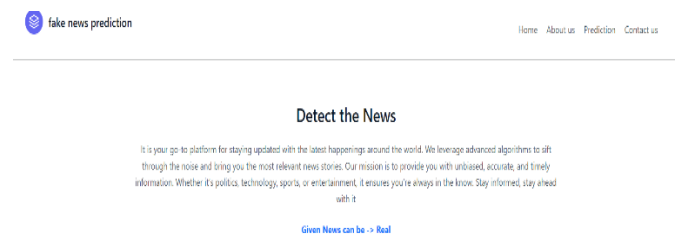
**C. Results And Discussions:**



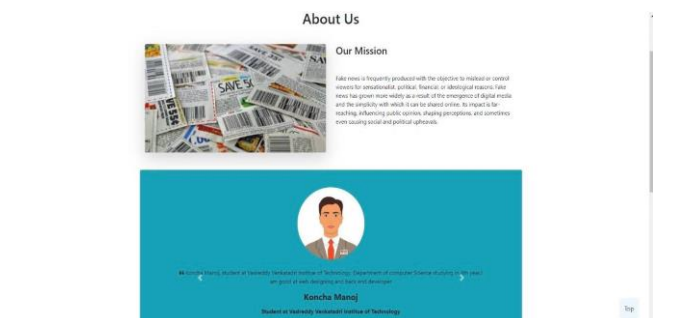
The result includes a website which gives a brief description about the project. To predict the news as a real or a fake article should go to the Prediction page



The Prediction page includes the input box to enter the text of the news article which have to be predicted. On submitting the article the prediction will be done



The Result will be displayed on the screen which depicts the misinformation.



This page is about us. Which tells the details of our group members.

## Conclusion

In conclusion, this project successfully developed a robust fake news prediction system using machine learning techniques. Through meticulous data preprocessing and feature engineering, the system achieved high accuracy in classifying news articles as fake or genuine. Its utilization of logistic regression as a baseline model demonstrated effectiveness in discerning between authentic and deceptive content. Moving forward, continuous monitoring and updates will be essential to ensure the system's reliability in identifying emerging forms of fake news. Overall, this project contributes to combating the spread of misinformation by leveraging advanced technology to distinguish trustworthy information sources from deceptive ones.

## References

- [1] E. Z. Mathews and N. Preethi, "Fake News Detection: An Effective Content-Based Approach Using Machine Learning Techniques," 2022 International Conference on Computer Communication and Informatics (ICCCI), Coimbatore, India, 2022, pp. 1-7, doi: 10.1109/ICCCI54379.2022.9741049.
- [2] V. Gupta, R. S. Mathur, T. Bansal and A. Goyal, "Fake News Detection using Machine Learning," 2022 International Conference on Machine Learning, Big Data, Cloud and Parallel Computing (COM-IT-CON), Faridabad, India, 2022, pp. 84-89, doi: 10.1109/COM-IT-CON54601.2022.9850560.
- [3] G. Rawat, T. Pandey, T. Singh, S. Yadav and P. K. Aggarwal, "Fake News Detection Using Machine Learning," 2023 International Conference on Artificial Intelligence and Smart Communication (AISC), Greater Noida, India, 2023, pp. 759-762, doi: 10.1109/AISC56616.2023.10085488.
- [4] Isaac, Hosea & Olalere, J & Adebayo, M. (2023). A Machine Learning Approach to Fake News Detection Using Support Vector Machine (SVM) and Unsupervised Learning Model. *Advances in Multidisciplinary and scientific Research Journal Publication*. 11. 10.22624/AIMS/CSEAN-SMART2023P1. ,TY - BOOK,AU - Jain, Anjali,AU - Shakya, Avinash,AU - Khatter, Harsh,AU - Gupta, Amit,PY - 2019/09/01,SP - 1,EP - 4,T1 - A smart System for Fake News Detection Using Machine Learning,DO - 10.1109/ICICT46931.2019.8977659.
- [5] E. Z. Mathews and N. Preethi, "Fake News Detection: An Effective Content-Based Approach Using Machine Learning Techniques," 2022 International Conference on Computer Communication and Informatics (ICCCI), Coimbatore, India, 2022, pp. 1-7, doi: 10.1109/ICCCI54379.2022.9741049.
- [6] Kaliyar, R.K., Goswami, A. & Narang, P. DeepFakeE: improving fake news detection using tensor decomposition-based deep neural network. *J Supercomput* 77, 1015–1037 (2021)
- [7] Jadhav, Shrutika S. and S. Thepade. "Fake News Identification and Classification Using DSSM and Improved Recurrent Neural Network Classifier." *Applied Artificial Intelligence* 33 (2019): 1058 - 1068.
- [8] Mandeep Singh, Mohammed Wasim Bhatt, Harpreet Singh Bedi, Umang Mishra, Performance of bernoulli's naive bayes classifier in the detection of fake news, *Materials Today: Proceedings*, 2020, ISSN 2214-7853.
- [9] Detecting fake news for reducing misinformation risks using analytics approaches
- [10] A. Dey, R.Z. Rafi, S.H. Parash, S.K. Arko, A. Chakrabarty, in: *Fake News Pattern Recognition Using Linguistic Analysis*, IEEE, 2018, pp. 305–309
- [11] M. Granik, V. Mesyura, in: *Fake News Detection Using Naive Bayes Classifier*, IEEE, 2017, pp. 900–903.
- [12] Lakshmanarao A, Swathi Y, Kiran TSR (2019) An efficient fake news detection system using machine learning. *Int J Innov Technol Exploring Eng (IJITEE)* 8(10).
- [13] Pedro Henrique Arruda Faustini, Thiago Ferreira Covães, Fake news detection in multiple platforms and languages, *Expert Systems with Applications*, vol 158, 2020, 113503, ISSN 0957-4174.
- [14] N. Garg and K. Sharma, "A review of classification and clustering techniques for sentiment analysis based on social network data." *International Journal of u- and e-service, Science and Technology*, vol. 13, no. 2, 2022.
- [15] M. Granik and V. Mesyura, "Fake news detection using naive bayes classifier," in 2017 IEEE first Ukraine conference on electrical and computer engineering (UKRCON). IEEE, 2017, pp. 900–903.
- [16] S. B. Parikh and P. K. Atrey, "Media-rich fake news detection: A survey," in 2018 IEEE conference on multimedia information processing and retrieval (MIPR). IEEE, 2018, pp. 436–441.
- [17] S. Shabani and M. Sokhn, "Hybrid machine-crowd approach for fake news detection," in 2018 IEEE 4th International Conference on Collaboration and Internet Computing (CIC). IEEE, 2018, pp. 299–306
- [18] N. F. Baair and A. Djeflal, "Fake News detection Using Machine Learning," 2020 2nd International Workshop on Human-Centric Smart Environments for Health and Well-being (IHSH), Boumerdes, Algeria, 2021, pp. 125-130, doi: 10.1109/IHSH51661.2021.9378748. keywords: {Support vector machines;Text recognition;Social networking (online);Supervised learning;Machine learning;Feature extraction;Classification algorithms;Fake news;Social media;Web Mining;Machine Learning;Support Vector Machine;TF-IDF}
- [19] Chih-Chung Chang and Chih-Jen Lin. LIBSVM – A Library for Support Vector Machines, July 15, 2018.
- [20] Kaggle. Getting Real about Fake News, 2016.
- [21] Gerard Salton and J Michael. McGill. 1983. Introduction to modern information retrieval, 1983.
- [22] J. C. Reis, A. Correia, F. Murai, A. Veloso, and F. Benevenuto, "Supervised learning for fake news detection," *IEEE Intelligent Systems*, vol. 34, no. 2, pp. 76–81, 2019
- [23] N. Garg and K. Sharma, "Sentiment Analysis of Events on Social Web." *International Journal of Innovative Technology and Exploring Engineering (IJITEE)*, vol. 9, no. 6, 2020.)