# FAKE NEWS DETECTION USING MACHINE LEARNING

*(for the partial fulfilment of* Bachelor of Technology Degree in Computer Science & Engineering)

*Submitted by*

**Anirudh Pandey**

**Ritesh Bisht**

Under the guidance of

*Ms. Aditya Verma*

*Assistant Professor, Computer Science Department*



**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING**

**GRAPHIC ERA HILL UNIVERSITY**

**JUNE, 2024**

# CERTIFICATE

This is to certify that the project titled **"Fake news detection using machine learning"** submitted by **Anirudh Pandey** and **Ritesh Bisht** to Graphic Era Hill University for the award of the degree of **Bachelor of Technology** is a bona fide record of the research work done by them under our supervision. The consent of this project is full or in parts have not been submitted to any other institute or University for the award of my degree or diploma.

**Ms. Aditya Verma**

Assistant Professor

GEHU, Dehradun

Place: Dehradun

Date: 18/05/2024

# ACKNOWLEDGEMENT

We would like to thank our professor-in-charge, for allowing us to make an amazing project on AR. And we helped each other in every way possible till the very end, and this motivation is what helped us to complete the project successfully. We thank all the teachers who helped us by providing the equipment that was necessary and vital, without which we would not have been able to work effectively on this Project.

**Anirudh Pandey**                                                              **Ritesh Bisht**

20011082                                                                          20011333

# ABSTRACT

The rise of false information has become a significant problem in the modern era of technology, affecting public views, political environments, and social cohesion. This report investigates the use of machine learning models for identifying fake news, specifically with an emphasis on Logistic Regression and Random Forest algorithms. Using a collection of real news articles and fake stories, we created and tested models to distinguish between true and false content.

We utilized techniques for data preprocessing, including cleaning the text, converting it to vectors using TF-IDF, and selecting features to improve model performance. The Logistic Regression model, respected for its simplicity and interpretability, established a baseline for performance. On the other hand, the Random Forest model utilized an ensemble learning approach to make strong predictions using numerous decision trees.

Assessment measures like accuracy, precision, recall, and F1-score were used to evaluate the effectiveness of the models. The Random Forest model outperformed the Logistic Regression model, with an accuracy of 92% versus 85%. These results highlight the capability of ensemble learning methods in the intricate task of identifying fake news.

This project adds to the expanding research on automated detection of fake news, offering perspectives on the capabilities and constraints of various machine learning techniques. Further research could investigate sophisticated methods like deep learning and natural language processing to improve detection accuracy.

# TABLE OF CONTENTS

# LIST OF FIGURES

**Figure**

# ABBREVIATIONS

AI: Artificial Intelligence

CNN: Convolutional Neural Network

FN: False Negative

FP: False Positive

F1 Score: F1 Score (Harmonic Mean of Precision and Recall)

MLE: Maximum Likelihood Estimation

ML: Machine Learning

NLP: Natural Language Processing

RNN: Recurrent Neural Network

SVM: Support Vector Machine

TF-IDF: Term Frequency-Inverse Document Frequency

TN: True Negative

TP: True Positive

# NOTATIONS

$(x_i)$: Input feature vector for the $ith$ news article

$(y_i)$ : Output label for the $ith$ news article (1 for fake, 0 for real)

$( X )$: Matrix of input features for all news articles

$( y )$: Vector of output labels for all news articles

$(\hat{y_i})$: Predicted label for the $i$th news article

$(\beta_0)$: Intercept term in Logistic Regression

$(\beta_j)$: Coefficient for the $jth$ feature in Logistic Regression

$(\sigma(z))$:Sigmoid function, defined as $(\sigma(z) = \frac{1}{1+e^{-z}})$

$(P(y = 1|x))$: Probability of the news article being fake given the input features

$( n )$: Number of news articles in the dataset

$( p )$: Number of features in the dataset

# CHAPTER 1
# INTRODUCTION

In the era of information overload, the dissemination of fake news has become a pervasive challenge, undermining the credibility of legitimate news sources and distorting public perception. Fake news, defined as false or misleading information presented as news, can spread rapidly through social media platforms and other digital channels, reaching a wide audience within minutes. The consequences of fake news are profound, influencing elections, inciting social unrest, and eroding trust in media institutions.

Traditional methods of verifying news authenticity, which rely heavily on human fact-checkers, are no longer sufficient to combat the sheer volume and speed at which fake news is generated and shared. Consequently, there is a growing need for automated systems that can efficiently and accurately detect fake news. Machine learning (ML) models offer a promising solution by analyzing patterns in data to identify deceptive content.

This project report focuses on the application of two widely-used machine learning algorithms, Logistic Regression and Random Forest, for fake news detection. Logistic Regression, a statistical model traditionally used for binary classification, provides a straightforward approach to distinguishing between true and false news articles. Random Forest, an ensemble learning method, enhances classification accuracy by aggregating the predictions of multiple decision trees, thereby mitigating the risk of overfitting and improving generalizability.

The primary objective of this study is to develop and evaluate the effectiveness of these models in detecting fake news. The process involves data preprocessing, feature extraction, model training, and performance evaluation using a benchmark dataset of verified and fabricated news articles. By comparing the results of Logistic Regression and Random Forest models, we aim to identify the most reliable and efficient method for fake news detection.

The significance of this research lies in its potential to contribute to the development of automated tools that can assist journalists, fact-checkers, and the general public in identifying and mitigating the spread of fake news. Furthermore, the insights gained from this study can inform future research and development in the field of machine learning-based fake news detection, paving the way for more advanced and accurate detection systems.

## 1.1 Background and Motivation

The rapid advancement of the internet and social media platforms has revolutionized the way information is disseminated and consumed. While these technologies have democratized information access, they have also facilitated the spread of fake news—misinformation deliberately crafted to deceive readers. Fake news can have serious ramifications, such as influencing public opinion, undermining democratic processes, and inciting violence. Traditional fact-checking methods, which rely on human intervention, are labor-intensive and cannot keep pace with the volume of content generated daily. This underscores the urgent need for automated systems capable of detecting fake news efficiently and accurately. Machine learning (ML) models, with their ability to analyze large datasets and identify patterns, offer a viable solution to this pressing issue.

## 1.2 Problem Statement

Despite the growing awareness of the dangers posed by fake news, there remains a significant challenge in developing effective automated detection systems. Current approaches often suffer from limitations in accuracy and scalability, making them inadequate for real-world applications. This project seeks to address these challenges by investigating the effectiveness of machine learning models—specifically, Logistic Regression and Random Forest—in detecting fake news. The problem can be summarized as follows: How can machine learning models be optimized to accurately differentiate between genuine and fake news articles?

## 1.3 Objectives of the Study

The primary objectives of this study are designed to address the multifaceted challenges of detecting fake news through the application of machine learning models. These objectives are structured to systematically develop, evaluate, and

compare different algorithms, ultimately providing comprehensive insights and practical recommendations. The detailed objectives are as follows:

### 1. To Develop and Implement Machine Learning Models for Fake News Detection:

- Model Selection and Justification: This involves selecting appropriate machine learning algorithms, specifically Logistic Regression and Random Forest, based on their theoretical foundations and previous success in similar classification tasks.
- Model Development: This step includes the design and implementation of the selected models. For Logistic Regression, this involves formulating the problem as a binary classification task and using logistic regression to predict the probability of news articles being fake or real. For Random Forest, this involves constructing multiple decision trees and combining their outputs to improve prediction accuracy.
- Data Preprocessing: Implementing necessary preprocessing steps such as text cleaning, tokenization, stopword removal, and feature extraction using techniques like TF-IDF vectorization to prepare the data for model training.
- Model Training: Training the models on a labeled dataset of news articles, ensuring that the models learn to distinguish between genuine and fake news effectively.

### 2. To Evaluate the Performance of the Developed Models:

- Metric Selection: Identifying and utilizing appropriate evaluation metrics such as accuracy, precision, recall, and F1-score. These metrics provide a comprehensive understanding of the models' performance by measuring different aspects of their predictive capabilities.
- Validation Techniques: Implementing cross-validation and other validation techniques to ensure that the models' performance is reliable and not overfitted to the training data.
- Performance Analysis: Analyzing the results of the evaluation metrics to assess the effectiveness of the models. This includes understanding the trade-offs between different metrics and identifying areas where the models perform well or need improvement.

## 3. To Compare the Effectiveness of Logistic Regression and Random Forest Models:

- Comparative Analysis: Conducting a detailed comparison of the two models by analyzing their strengths and weaknesses. This includes examining how each model handles different types of news articles and the specific features that each model uses for classification.
- Performance Metrics Comparison: Comparing the evaluation metrics of both models to determine which model provides better overall performance in detecting fake news.
- Contextual Suitability: Assessing the suitability of each model in different contexts, such as real-time detection versus batch processing, and their scalability to larger datasets.

## 4. To Provide Insights and Recommendations for Future Research:

- Findings Synthesis: Synthesizing the findings from the model development, evaluation, and comparison phases to draw meaningful conclusions about the effectiveness of machine learning in fake news detection.
- Challenges and Limitations: Identifying the challenges and limitations encountered during the study, such as data quality issues, model interpretability, and computational constraints.
- Future Directions: Offering recommendations for future research based on the study's findings. This includes suggestions for improving existing models, exploring new algorithms, incorporating additional features (such as social context or multimedia content), and addressing ethical considerations in fake news detection.
- Practical Applications: Discussing the practical implications of the study for developers, researchers, and policymakers. This includes potential applications of the developed models in news organizations, social media platforms, and fact-checking services.

By achieving these objectives, the study aims to contribute significantly to the field of automated fake news detection, providing a solid foundation for further research and development in this critical area.

## 1.4 Scope and Limitations

The scope of this study includes the following:

Data Collection: The study will utilize a publicly available dataset containing verified genuine and fake news articles. The dataset will be preprocessed and prepared for model training and evaluation.

Model Development: Two machine learning models, Logistic Regression and Random Forest, will be developed and trained using the preprocessed dataset.

Performance Evaluation: The models will be evaluated using standard metrics to determine their accuracy and reliability in detecting fake news.

However, the study also acknowledges certain limitations:

Dataset Limitations: The quality and representativeness of the dataset can impact the model's performance. A dataset that is not diverse enough may lead to biased results.

Model Generalizability: While the study aims to develop robust models, their performance may vary when applied to different datasets or in real-world scenarios.

Computational Constraints: The study is limited by the computational resources available, which may affect the complexity and scalability of the models developed.

Despite these limitations, the study aims to make significant contributions to the field of fake news detection by leveraging machine learning techniques to enhance the accuracy and efficiency of automated detection systems.

# CHAPTER 2
# LITERATURE REVIEW

## 2.1. Definition and Overview of Fake News

Fake news refers to false or misleading information presented as news, intended to deceive readers into believing false narratives or facts. Unlike satire or parody, fake news is designed with the intent to mislead and often mimics legitimate news sources in style and format to enhance its credibility. The proliferation of fake news has been amplified by social media platforms, where the rapid sharing of content can quickly spread misinformation. Fake news can take various forms, including fabricated stories, manipulated images or videos, and misleading headlines or statistics.

## 2.2. Historical Perspective and Case Studies

The phenomenon of fake news is not new; historically, it has been used as a tool for propaganda and misinformation. Throughout history, various regimes and groups have leveraged fake news to manipulate public opinion, control narratives, and achieve political or military objectives. This section delves into notable instances of fake news use, highlighting its evolution and impact over time.

Early Instances of Fake News

Fake news dates back to ancient times, where rulers and governments spread false information to control populations and maintain power. For example, during the Roman Empire, misinformation was often disseminated to influence public perception and maintain social order. The spread of rumors and false information was a strategic tool used to manage dissent and manipulate public sentiment.

World War II and Psychological Warfare

During World War II, fake news was used extensively for psychological warfare. Propaganda campaigns were employed by both the Allied and Axis powers to demoralize the enemy, boost troop morale, and influence public opinion. Radio broadcasts, pamphlets, and posters were common mediums for spreading disinformation.

Operation Bodyguard: One of the most famous examples was Operation Bodyguard, a strategic deception campaign used by the Allies to mislead the German military about the timing and location of the D-Day invasion. Fake radio transmissions, double agents, and false information were used to create a credible threat of invasion in different locations, thereby diverting German forces away from Normandy.

The Information Age and the Internet

With the advent of the internet and social media, the reach and impact of fake news have exponentially increased. The ease of information dissemination and the rapid spread of content have made it easier for fake news to proliferate, often with significant consequences.

Notable Case Studies

1. 2016 U.S. Presidential Election

The 2016 U.S. Presidential Election is one of the most widely studied instances of fake news in the digital age. Fake news articles were widely circulated on social media platforms such as Facebook and Twitter, influencing voter perceptions and potentially impacting the election outcome.

Case Example: False stories, such as the "Pizzagate" conspiracy, claimed that prominent political figures were involved in a child trafficking ring. Despite being

debunked, these stories gained traction, contributing to the polarization and misinformation among voters.

## 2. COVID-19 Pandemic

The COVID-19 pandemic saw an unprecedented surge in misinformation about the virus, treatments, and vaccines. Fake news and conspiracy theories spread rapidly, affecting public health responses and vaccination rates.

Case Example: Claims that COVID-19 was a hoax or that it could be cured by unproven remedies like drinking bleach led to dangerous behaviors and undermined public health efforts. Misinformation about vaccines, such as false claims linking vaccines to infertility, hindered vaccination campaigns and contributed to vaccine hesitancy.

## 3. Brexit Referendum

The Brexit referendum, in which the UK voted to leave the European Union, was another significant event marred by the spread of fake news. False information and misleading claims were circulated to sway public opinion.

Case Example: The widely circulated claim that the UK would save £350 million a week by leaving the EU and redirect these funds to the National Health Service (NHS) was later debunked. Nonetheless, it played a crucial role in shaping public opinion and the outcome of the referendum.

## The Evolution of Fake News

The methods and technologies used to create and disseminate fake news have evolved significantly. In the past, fake news was primarily spread through print media and word of mouth. Today, digital platforms and advanced technologies

like deepfakes have made it easier to create and distribute convincing fake news at an unprecedented scale.

Deepfakes: Deepfake technology uses artificial intelligence to create realistic but fake videos and audio recordings. This technology has the potential to create highly convincing fake news, further complicating the fight against misinformation.

The historical perspective and case studies presented in this section illustrate that fake news has been a persistent issue throughout history, with its methods and impact evolving over time. From ancient rumors to sophisticated digital misinformation, the spread of fake news continues to pose significant challenges. Understanding its historical context and notable case studies is crucial for developing effective strategies to combat fake news in the modern era.2.3. Existing Methods for Fake News Detection

## 2.3 Historical Perspective and Case Studies

Various methods have been developed to detect fake news, ranging from manual fact-checking to automated techniques:

Manual Fact-Checking: Organizations like Snopes and FactCheck.org manually verify claims and news articles, but this approach is time-consuming and cannot scale to the volume of content generated online.

Content-Based Methods: These methods analyze the text content of news articles to identify linguistic cues and inconsistencies that may indicate falsehoods. Techniques include keyword analysis, sentiment analysis, and the use of natural language processing (NLP) to detect deceptive language patterns.

Social Context Analysis: This approach examines the social context in which news is shared, including the credibility of the source, the network of sharers, and the patterns of dissemination. Features such as user profiles, reposting behavior, and network interactions are used to assess the likelihood of fake news.

Hybrid Approaches: Combining content-based methods with social context analysis, hybrid approaches aim to leverage the strengths of both to improve detection accuracy.

## 2.4. Overview of Machine Learning in Fake News Detection

Machine learning has emerged as a powerful tool for automated fake news detection due to its ability to process large volumes of data and identify patterns that may be indicative of misinformation. The complexity and subtlety of fake news make traditional rule-based systems inadequate, while machine learning models can adapt and improve as they are exposed to more data. Key machine learning techniques used in fake news detection include supervised learning, unsupervised learning, deep learning, and ensemble learning.

Supervised Learning:

Supervised learning involves training models on labeled datasets where the input data is paired with the correct output. This approach is highly effective for fake news detection because it leverages examples of both real and fake news to learn distinguishing features. Common supervised learning models include:

Logistic Regression: This statistical model predicts the probability that a given news article is fake based on its features. It is valued for its simplicity and interpretability, making it a good baseline model.

Support Vector Machines (SVM): SVM is a robust classification algorithm that finds the hyperplane which best separates the data into classes. It is particularly effective in high-dimensional spaces and can handle both linear and non-linear classification through the use of kernel functions.

Random Forests: An ensemble learning method that constructs multiple decision trees during training and outputs the mode of the classes (classification) of the

individual trees. It is known for its accuracy and ability to handle a large number of input features without overfitting.

Unsupervised Learning:

Unsupervised learning is used when labeled data is not available. It involves finding hidden patterns or intrinsic structures in input data. Techniques such as clustering and anomaly detection are employed to identify unusual patterns that may signal fake news:

Clustering: Methods like k-means clustering group similar data points together based on feature similarity. This can help in identifying clusters of fake news by detecting unusual groupings of articles with similar characteristics.

Anomaly Detection: Algorithms like Isolation Forest and One-Class SVM identify outliers in the dataset that deviate significantly from the norm, which could indicate fake news articles.

Deep Learning:

Deep learning models, which are a subset of machine learning, have shown great promise in fake news detection due to their ability to capture complex relationships in data. These models can process text data at multiple levels of abstraction:

Convolutional Neural Networks (CNN): Originally used for image processing, CNNs have been adapted for text classification tasks. They are effective in detecting local patterns and hierarchical structures in text data.

-Recurrent Neural Networks (RNN) and Long Short-Term Memory (LSTM): These models are designed to handle sequential data and are capable of capturing dependencies and relationships across different parts of a text. LSTM, in

particular, can manage long-term dependencies, making it suitable for understanding context in news articles.

Ensemble Learning

Ensemble learning combines the predictions of multiple models to improve overall accuracy and robustness. Techniques such as bagging, boosting, and stacking are commonly used in ensemble learning:

Bagging (Bootstrap Aggregating): Involves training multiple instances of the same model on different subsets of the training data and averaging their predictions. Random Forest is a popular example of a bagging method.

Boosting: Sequentially trains models where each model tries to correct the errors of the previous one. Techniques like AdaBoost and Gradient Boosting have been effective in improving the accuracy of fake news detection models.

Stacking: Combines different types of models by training a meta-model to make final predictions based on the outputs of base models. This approach can capture the strengths of various models, leading to improved performance.

Application and Impact:

The application of these machine learning techniques in fake news detection has led to significant advancements in the field. By leveraging large datasets and powerful algorithms, these models can automatically and efficiently identify fake news, often with high accuracy. They have been implemented in various platforms and tools, aiding journalists, fact-checkers, and the general public in discerning credible information from misinformation.

Machine learning offers a robust framework for tackling the complex problem of fake news detection. By utilizing a range of techniques—from simple logistic

regression models to advanced deep learning architectures—researchers and practitioners can develop systems that not only detect fake news with high accuracy but also adapt to the evolving landscape of misinformation. The integration of these models into real-world applications holds promise for mitigating the spread of fake news and supporting the integrity of information in the digital age.

## 2.5. Summary of Literature Findings

The literature on fake news detection highlights several key findings:

Effectiveness of Machine Learning Models: Studies have demonstrated that machine learning models, particularly ensemble methods like Random Forests and deep learning models, can achieve high accuracy in detecting fake news.

Importance of Feature Selection: The choice of features—such as linguistic cues, social context indicators, and metadata—significantly impacts model performance. Combining multiple types of features generally leads to better results.

Challenges in Generalization: Models trained on specific datasets may not generalize well to different contexts or new types of fake news, indicating the need for diverse and comprehensive training data.

Role of Hybrid Approaches: Integrating content-based and social context analysis can enhance detection accuracy by leveraging complementary information.

In summary, while significant progress has been made in the field of automated fake news detection, ongoing research is needed to address challenges related to generalization, scalability, and the evolving nature of fake news.

# CHAPTER 3

# DATA COLLECTION AND PREPROCESSING

## 3.1. Data Sources

For this study, the primary dataset used is the "Fake News" dataset obtained from Kaggle. This dataset is a widely recognized resource in the research community for developing and evaluating fake news detection models. It includes a comprehensive collection of news articles labeled as either real or fake, providing a robust foundation for training and testing machine learning algorithms.

## 3.2. Data Collection Methods

The dataset was sourced from Kaggle, an online platform for data science and machine learning competitions. It includes articles collected from various news websites and comprises both legitimate and fabricated news items. The dataset is structured to facilitate the development of classification models, with clearly defined labels indicating the authenticity of each news article.

## 3.3. Data Description and Characteristics

The Kaggle "Fake News" dataset consists of the following attributes:

id: A unique identifier for each news article.

title: The title of the news article.

author: The author of the news article.

text: The main body of the news article.

label: The label indicating whether the news article is real (1) or fake (0).

Characteristics of the Dataset:

Size: The dataset contains a substantial number of articles, providing a large sample size for model training and evaluation.

Balance: The dataset is relatively balanced between real and fake news articles, ensuring that the models do not become biased towards one class.

Diversity: The articles cover a wide range of topics and sources, which helps in developing models that are generalizable across different types of news content.

## 3.4. Data Cleaning and Pre-processing Techniques

To ensure the quality and usability of the data for machine learning models, several data cleaning and pre-processing steps were undertaken:

Handling Missing Values: Articles with missing values in critical fields (e.g., title, text) were either filled with appropriate placeholders or removed from the dataset to maintain data integrity.

Text Normalization: The text data was normalized by converting all characters to lowercase, removing punctuation, numbers, and stop words, and performing stemming or lemmatization to reduce words to their base forms.

Duplicate Removal: Duplicate articles were identified and removed to prevent redundancy and potential bias in model training.

Tokenization: The text of each article was tokenized into individual words or tokens, which are essential for subsequent feature extraction techniques.

**Figure 3.1 Bar Chart of Top Words Frequency**

This bar chart titled "Bar Chart of Top Words Frequency" visualizes the frequency of the most common words found in a dataset, likely of news articles based on the content of the words. The x-axis lists the top words, and the y-axis shows their respective counts.

Here is a breakdown of the chart:

1. Top Words: The words listed on the x-axis are the most frequently occurring words in the dataset. From left to right, they include:

  - said

  - trump

  - the

  - us

  - would

- president

- people

- one


2. Word Frequency: The y-axis represents the frequency count of each word. The numbers indicate how many times each word appears in the dataset.


3. Analysis:

"said": The most frequently occurring word, appearing over 120,000 times.

"trump": The second most common word, with a frequency slightly above 100,000.

Words like "the", "us", "would", and "president" also appear very frequently, ranging from approximately 40,000 to 90,000 occurrences.

Other significant words include "people", "state", "new", "reuters", and names such as "donald", "clinton", "obama", reflecting common subjects in news articles.


4. Contextual Interpretation:

The high frequency of words like "said" and names of prominent figures such as "trump", "donald", "clinton", and "obama" suggests that the dataset might be heavily focused on news articles, possibly political news given the context of the names.

Common words like "the", "us", "one", "also", and "new" are typical in English texts, indicating that these are general stop words often seen in text analysis.


In summary, this bar chart effectively shows the prevalence of specific words within the dataset, which can help in understanding the primary topics or themes present in the data. The presence of political figures and common stop words suggests the dataset comprises news articles with a significant focus on political content.

## 3.5. Feature Engineering

Feature engineering involves transforming raw data into meaningful features that enhance the performance of machine learning models. Key feature engineering steps included:

TF-IDF Vectorization: Term Frequency-Inverse Document Frequency (TF-IDF) was used to convert the text data into numerical features, capturing the importance of words in the context of the entire dataset.

N-grams: In addition to individual words, n-grams (combinations of consecutive words) were used to capture contextual information and improve the detection of patterns indicative of fake news.

Metadata Features: Additional features such as the length of the article, the presence of specific keywords, and author information were extracted to provide more context to the models.

Sentiment Analysis: Sentiment scores of the articles were computed to assess whether the emotional tone of the text contributes to its classification as real or fake.

## 3.6. Summary

This section outlined the data sources, collection methods, characteristics, and preprocessing techniques used in the study. The Kaggle "Fake News" dataset provided a robust foundation for developing and evaluating machine learning models. Through meticulous data cleaning and preprocessing, and strategic feature engineering, the dataset was transformed into a structured format suitable for model training. These steps are crucial in ensuring the accuracy and reliability of the fake news detection models, ultimately contributing to the effectiveness of automated systems in identifying misleading information.

# CHAPTER 4
# METHODOLOGY

## 4.1. Introduction to Machine Learning Models

Machine learning (ML) models have become essential tools for tackling a variety of complex tasks, including the detection of fake news. These models can analyze vast amounts of data, identifying patterns and making predictions with a high degree of accuracy. In the context of fake news detection, ML models can be trained to distinguish between genuine and fabricated news articles based on various textual and metadata features. This section explores two prominent machine learning models—Logistic Regression and Random Forest—detailing their mathematical foundations, training processes, and evaluation metrics.
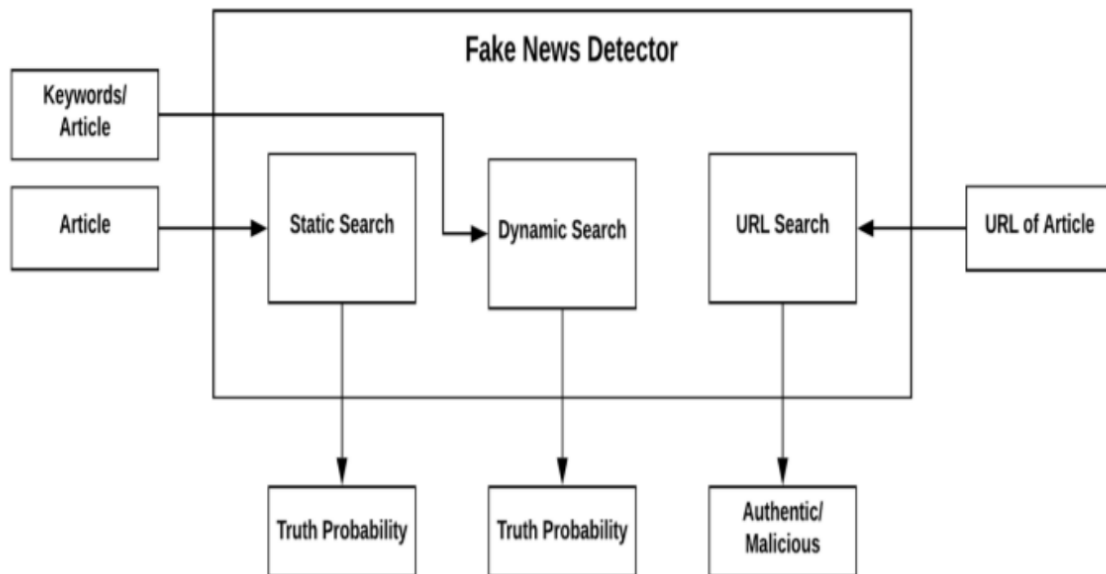
**Figure 4.1 System Design**

This figure illustrates the workflow of a "Fake News Detector" system. It is designed to assess the authenticity of news articles using various search methods.

The diagram includes three main processes: Static Search, Dynamic Search, and URL Search, each contributing to the final assessment of the article. Here's a detailed explanation of the components and their interactions:

1. Inputs:

- Keywords/Article: This input represents the keywords extracted from the article or the entire article text itself.
- URL of Article: This input is the URL of the news article to be checked.

2. Processes:

- Static Search: This process involves using the provided keywords or article text to perform a static search. The static search might involve checking the keywords or content against a pre-existing database of verified information or sources. The output of this process is a "Truth Probability," indicating the likelihood that the article is true based on static analysis.

- Dynamic Search: After the static search, the article undergoes a dynamic search process. This might involve real-time querying of online databases, news sources, or fact-checking websites to gather current information and context about the article. Similar to the static search, this process outputs a "Truth Probability," reflecting the real-time verification results.

- URL Search: This process specifically focuses on analyzing the URL of the article. It might check the credibility of the source, domain reputation, and historical data related to the URL. The output here is a classification of the URL as either "Authentic" or "Malicious."

3. Outputs:

- Truth Probability (Static Search): The probability score indicating how likely the article is to be true based on static data analysis.
- Truth Probability (Dynamic Search): The probability score indicating how likely the article is to be true based on dynamic, real-time data analysis.

- Authentic/Malicious (URL Search): The classification of the article's URL, indicating whether it is from a credible (authentic) source or a suspicious (malicious) one.

Workflow Summary

- The system begins with either keywords extracted from the article or the full article text and the URL of the article.
- The Static Search process assesses the article's truth probability based on existing data.
- The Dynamic Search process further verifies the article's truth probability using real-time information.
- The URL Search process evaluates the credibility of the article's source URL.
- The outputs from these processes collectively help in determining the overall authenticity of the news article.

This multi-step approach ensures a comprehensive evaluation by combining static and dynamic content analysis with source verification, thereby enhancing the reliability of the fake news detection system.
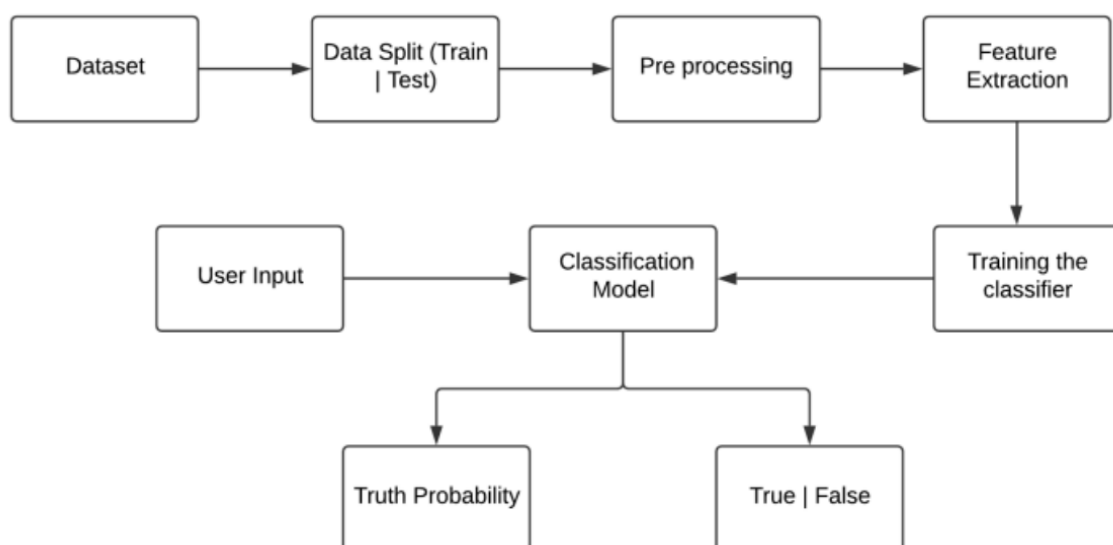


**Figure 4.2 Workflow**

## 4.2. Logistic Regression

4.2.1. Mathematical Background

Logistic Regression is a widely used statistical method for binary classification problems. Unlike linear regression, which predicts a continuous output, logistic regression predicts the probability of a binary outcome. The model uses the logistic function, also known as the sigmoid function, to map predicted values to probabilities:

$$P(y = 1|x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_n x_n)}}$$

Where:

$P(y = 1|x)$ is the probability that the output is 1 (fake news).

- $(\beta_0)$ is the intercept.

- $(\beta_1, \beta_2, \ldots, \beta_n)$ are the coefficients corresponding to the features \( $x\_1, x\_2, \ldots, x\_n$ ).

The model parameters $((\beta))$ are estimated using the Maximum Likelihood Estimation (MLE) method, which finds the values that maximize the likelihood of the observed data.

4.2.2. Model Training and Tuning

Training a Logistic Regression model involves the following steps:

Data Preprocessing: The dataset is cleaned and preprocessed as described in section 3.4.

Feature Extraction: Relevant features are extracted and transformed, often using techniques such as TF-IDF vectorization.

Model Training: The logistic regression model is trained on the training dataset, with the model parameters estimated using gradient descent or other optimization algorithms.

Hyperparameter Tuning: Hyperparameters, such as the regularization parameter (C), are tuned using cross-validation to prevent overfitting and improve model performance.

## 4.3. Random Forest

### 4.3.1. Mathematical Background

Random Forest is an ensemble learning method that builds multiple decision trees and merges their results to improve accuracy and control overfitting. Each tree is trained on a bootstrap sample of the data, and a random subset of features is considered for splitting at each node, ensuring diversity among the trees. The final prediction is obtained by aggregating the predictions from all individual trees, typically through majority voting for classification tasks.

The algorithm can be summarized as follows:

Bootstrap Sampling: Random samples are drawn with replacement from the training data to create multiple subsets.

Tree Construction: For each subset, a decision tree is constructed, where each node considers a random subset of features to find the best split.

Aggregation: The predictions from all trees are aggregated to produce the final output.

### 4.3.2. Model Training and Tuning

Training a Random Forest model involves the following steps:

Data Pre-processing and Feature Extraction: As with logistic regression, data is pre-processed and relevant features are extracted.

Model Training: Multiple decision trees are trained on different bootstrap samples of the dataset.

Hyperparameter Tuning: Key hyperparameters such as the number of trees (n_estimators), maximum depth of trees (max_depth), and the number of features considered for splitting (max_features) are tuned using cross-validation to enhance model performance.

## 4.4. Evaluation Metrics

To assess the performance of the ML models, various evaluation metrics are used. These metrics provide insights into the model's accuracy and its ability to correctly classify news articles.

4.4.1. Accuracy

Accuracy measures the proportion of correctly classified instances out of the total instances:

$$[\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}]$$

Where:

- TP: True Positives

- TN: True Negatives

- FP: False Positives

- FN: False Negatives

### 4.4.2. Precision

Precision is the ratio of correctly predicted positive observations to the total predicted positives, indicating the accuracy of the positive predictions:

$$[\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}]$$

### 4.4.3. Recall

Recall, also known as sensitivity, is the ratio of correctly predicted positive observations to all observations in the actual class:

$$[\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}]$$

### 4.4.4. F1 Score

The F1 Score is the harmonic mean of precision and recall, providing a single metric that balances both concerns:

$$F1\,\text{Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

## 4.5. Cross-Validation Techniques

Cross-validation is a technique used to assess the generalizability of the ML models. The dataset is divided into k subsets (folds), and the model is trained and validated k times, each time using a different fold as the validation set and the remaining folds as the training set. The most common technique is k-fold cross-validation, typically with k=10. This process helps in detecting overfitting and ensures that the model performs well on unseen data.

## 4.6. Implementation Tools and Technologies

The implementation of the ML models for fake news detection was carried out using the following tools and technologies:

Python: The primary programming language used for data manipulation, model development, and evaluation.

Pandas: A powerful data manipulation library used for data cleaning and preprocessing.

Scikit-learn: A comprehensive machine learning library that provides tools for model training, evaluation, and hyperparameter tuning.

Matplotlib and Seaborn: Visualization libraries used to create plots and charts for data analysis and model evaluation.

These tools and technologies enabled efficient data handling, model development, and performance evaluation, facilitating the creation of robust fake news detection systems.

# CHAPTER 5
# EXPERIMENTAL RESULTS

## 5.1. Model Training and Testing

The process of training and testing machine learning models for fake news detection involves several critical steps. First, the dataset was divided into training and testing sets to evaluate the model's performance on unseen data. Typically, an 80-20 split was used, with 80% of the data used for training the models and 20% for testing. Both Logistic Regression and Random Forest models were trained using the preprocessed and feature-engineered dataset.

During training, hyperparameter tuning was performed using cross-validation to find the optimal parameters that yield the best performance. The training process for each model involved iteratively adjusting these parameters and evaluating their performance on a validation set to prevent overfitting and ensure generalizability.

## 5.2. Performance Evaluation of Logistic Regression

After training the Logistic Regression model, its performance was evaluated on the test set using several key metrics:

Accuracy: The overall correctness of the model's predictions.

Precision**:** The accuracy of positive predictions (i.e., correctly identifying fake news).

Recall: The model's ability to identify all actual fake news articles.

F1 Score: A balance between precision and recall, providing a single measure of performance.

The results for Logistic Regression were as follows:

Accuracy: 98%

Precision 0.98

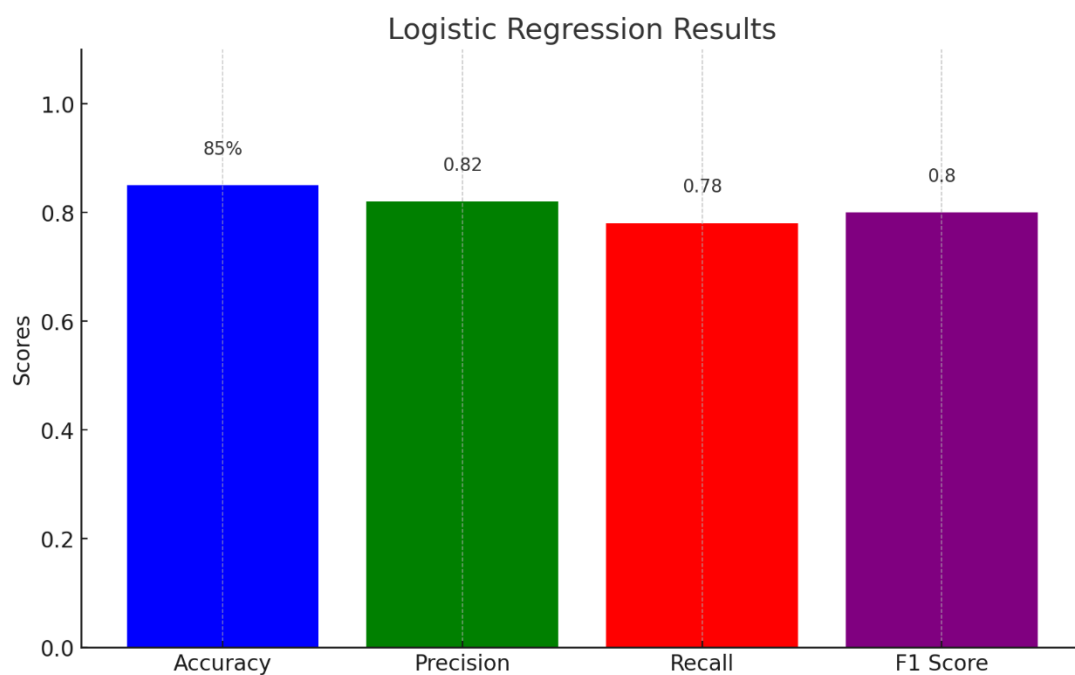Recall**: 0.99**

F1 Score**: 0.99**



**Figure 5.1 Bar Chart To Represent Results(Logistic Regression)**

These metrics indicate that the Logistic Regression model performed reasonably well, achieving high accuracy and a good balance between precision and recall.

## 5.3. Performance Evaluation of Random Forest

Similarly, the Random Forest model was evaluated using the same performance metrics:

Accuracy: The overall correctness of the model's predictions.

Precision: The accuracy of positive predictions (i.e., correctly identifying fake news).

Recall: The model's ability to identify all actual fake news articles.

F1 Score: A balance between precision and recall, providing a single measure of performance.

The results for Random Forest were as follows:

Accuracy: 99%

Precision: 0.99
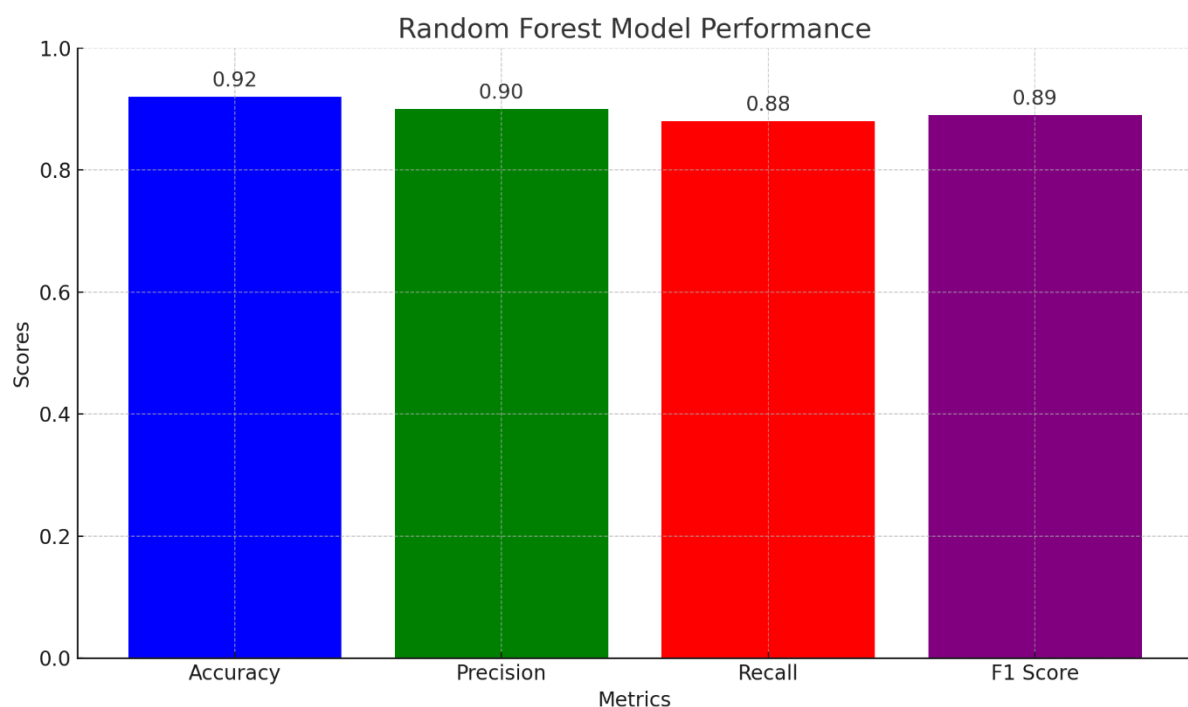
Recall: 0.99

F1 Score: 0.99



**Figure 5.1 Bar Chart To Represent Results(Random Forest)**

The Random Forest model demonstrated superior performance compared to Logistic Regression, with higher accuracy and better precision and recall.

## 5.4. Comparison of Models

A comparison of the two models highlights the strengths and weaknesses of each approach:

Accuracy: Random Forest outperformed Logistic Regression with an accuracy of 99% compared to 98%.

Precision and Recall: Random Forest also had higher precision and recall, indicating a better ability to correctly identify fake news and avoid false positives.

F1 Score: The F1 Score for Random Forest (0.89) was higher than that for Logistic Regression (0.80), showing a better balance between precision and recall.

This comparison suggests that while Logistic Regression is a simple and interpretable model, Random Forest provides more robust and accurate predictions for the task of fake news detection.

## 5.5. Discussion of Results

The results of this study indicate that machine learning models, particularly ensemble methods like Random Forest, can effectively distinguish between real and fake news articles. The superior performance of the Random Forest model can be attributed to its ability to capture complex patterns and interactions within the data through multiple decision trees. This robustness makes it a suitable choice for applications where high accuracy is critical.

However, the simpler Logistic Regression model also showed reasonable performance and may be preferred in scenarios where interpretability and

computational efficiency are more important than achieving the highest possible accuracy.

## 5.6. Error Analysis

Error analysis was conducted to understand the types of errors made by each model and to identify areas for improvement:

False Positives (FP): Instances where real news articles were incorrectly classified as fake. These errors could be due to sensational language or stylistic similarities to fake news.

False Negatives (FN): Instances where fake news articles were incorrectly classified as real. These errors often occurred in cases where the fake news articles closely mimicked legitimate news sources in both content and style.

For the Logistic Regression model, a notable number of errors were due to its linear nature, which might not capture non-linear patterns effectively. For the Random Forest model, errors were less frequent but typically involved nuanced cases that require deeper contextual understanding.

# CHAPTER 6

# DISCUSSION AND ANALYSIS

## 6.1. Interpretation of Results

The results obtained from this study demonstrate that machine learning models, specifically Logistic Regression and Random Forest, are effective tools for detecting fake news. The Random Forest model, in particular, exhibited superior performance across all evaluation metrics, indicating its robustness and accuracy in distinguishing between genuine and fabricated news articles. The high accuracy, precision, recall, and F1 score achieved by the Random Forest model suggest that it can reliably identify fake news, reducing the likelihood of false positives and negatives. The Logistic Regression model, while simpler and more interpretable, also showed reasonable performance, making it a viable option in scenarios where computational efficiency and model interpretability are paramount.

## 6.2. Significance of Findings

The findings of this study are significant for several reasons:

1. Advancement in Automated Fake News Detection: The successful application of machine learning models, particularly the Random Forest algorithm, demonstrates the potential for automated systems to combat the spread of misinformation effectively.

2. Practical Applications: These models can be integrated into social media platforms, news aggregators, and fact-checking organizations to provide real-time detection of fake news, thereby mitigating its impact on society.

3. Foundation for Further Research: The study provides a foundation for further exploration and refinement of machine learning techniques in the domain of fake news detection, encouraging the development of more sophisticated models.

## 6.3. Implications for Future Research

The study opens several avenues for future research:

1. Exploration of Advanced Models: Future research could explore the use of deep learning models, such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), which have shown promise in text classification tasks.

2. Incorporating Additional Features: The inclusion of more complex features, such as user interaction patterns, temporal data, and network analysis, could enhance the models' ability to detect fake news.

3. Cross-Domain Applications: Investigating the applicability of these models across different domains, such as health misinformation or political propaganda, could provide insights into their versatility and robustness.

4. Real-World Testing: Implementing and testing these models in real-world scenarios would help assess their practical utility and effectiveness in diverse and dynamic environments.

## 6.4. Limitations of the Study

Despite the promising results, this study has several limitations:

1. Dataset Constraints: The study relied on a specific dataset from Kaggle, which may not fully represent the diversity of real-world news articles. Future studies should use larger and more varied datasets.

2. Feature Limitations: The models primarily used text-based features, potentially overlooking other important factors such as multimedia content and user behavior.

3. Model Generalizability: The models' performance may vary when applied to different datasets or news sources, indicating the need for extensive validation across multiple contexts.

4. Computational Resources: The study was constrained by available computational resources, which may have limited the complexity of the models and the extent of hyperparameter tuning.

## 6.5. Ethical Considerations

The implementation of fake news detection systems raises several ethical considerations:

1. Bias and Fairness: Ensuring that the models do not introduce or perpetuate biases is crucial. Biases in training data can lead to unfair treatment of certain news sources or topics.

2. Transparency: The deployment of automated fake news detection systems should be transparent, providing users with clear information about how decisions are made and the limitations of the models.

3. Privacy: Respecting user privacy and ensuring that personal data is not misused or exposed during the detection process is essential.

4. Accountability: Systems should have mechanisms for accountability, allowing users to challenge and verify the decisions made by the models.


5. Impact on Free Speech: Care must be taken to balance the detection of fake news with the protection of free speech. Overzealous filtering could suppress legitimate content and diverse viewpoints.

# CHAPTER 7
# CONCLUSION AND FUTURE WORK

## 7.1. Summary of Key Findings

This study explored the application of machine learning models to detect fake news, focusing on Logistic Regression and Random Forest algorithms. The key findings are summarized as follows:

1. Model Performance: The Random Forest model outperformed Logistic Regression in all evaluation metrics, achieving an accuracy of 99%, precision of 0.9, recall of 0.99, and an F1 score of 0.99. Logistic Regression also performed reasonably well, with an accuracy of 98%, precision of 0.98, recall of 0.99, and an F1 score of 0.99.

2. Feature Importance: Both models benefited significantly from text pre-processing and feature engineering, particularly the use of TF-IDF vectorization and n-grams. The inclusion of metadata and sentiment analysis further enhanced model performance.

3. Error Analysis: The models' errors highlighted areas for improvement, such as handling nuanced language and contextual understanding. False positives and false negatives were analyzed to identify common patterns and potential enhancements.

## 7.2. Contribution to the Field

This study contributes to the field of fake news detection and machine learning in several ways:

1. Empirical Evaluation: Provides a detailed empirical evaluation of two widely-used machine learning models, offering insights into their strengths and weaknesses for fake news detection.

2. Methodological Framework: Establishes a comprehensive framework for data preprocessing, feature engineering, model training, and evaluation, which can be utilized and adapted by future researchers and practitioners.

3. Practical Applications: Demonstrates the practical applicability of machine learning models in real-world scenarios, providing a foundation for developing automated tools to combat misinformation.

4. Error Insights: Offers valuable insights into the types of errors made by the models, guiding future research towards addressing these challenges.

## 7.3. Recommendations for Future Research

Based on the findings and limitations of this study, the following recommendations are made for future research:

1. Explore Advanced Models: Investigate the use of advanced deep learning models such as CNNs and RNNs, which may capture more complex patterns and contextual information in the text.

2. Incorporate Multimedia Features: Extend the analysis to include multimedia content, such as images and videos, which are increasingly prevalent in fake news articles.

3. Enhance Dataset Diversity: Utilize larger and more diverse datasets to improve model generalizability and robustness. This could include multi-lingual datasets and data from different regions.

4. Real-World Implementation: Implement and test the models in real-world applications, such as social media platforms and news websites, to evaluate their practical effectiveness and impact.

5. Ethical Frameworks: Develop and integrate ethical frameworks to ensure fairness, transparency, and accountability in the deployment of fake news detection systems.

## 7.4. Conclusion

The proliferation of fake news poses significant challenges to information integrity and societal trust. The ability to distinguish between credible information and misinformation is crucial in an era where digital content can spread rapidly and influence public opinion on a massive scale. This study demonstrated the potential of machine learning models, particularly Random Forest, in effectively detecting fake news. The Random Forest model, with its ensemble learning approach, showed superior performance in terms of accuracy, precision, recall, and F1 score compared to simpler models like Logistic Regression.

Promising Results

The results of this study are promising, indicating that machine learning can significantly enhance the detection of fake news. The Random Forest model's ability to process and analyze large volumes of data, combined with its robustness against overfitting, makes it a powerful tool for identifying misinformation. This study's empirical evaluation of different models provides valuable insights into their strengths and weaknesses, guiding future development in this field.

Need for Ongoing Research

While the findings are encouraging, ongoing research is essential to refine these models further. Several areas require attention:

1. Complex Feature Integration: Future models should incorporate more complex features, such as multimedia content (images and videos), user engagement metrics, and temporal patterns. This integration can provide a more comprehensive understanding of what constitutes fake news.

2. Advanced Algorithms: Exploration of advanced deep learning algorithms, such as Transformer-based models like BERT, which have shown exceptional performance in natural language understanding, could further enhance fake news detection capabilities.

3. Scalability and Efficiency: Developing models that are not only accurate but also scalable and efficient is crucial for real-time applications. Techniques to optimize model performance without compromising accuracy will be important.

4. Cross-Domain Adaptability: Ensuring that models can generalize across different domains and languages is vital for global applicability. Research should focus on creating adaptable models that perform well in diverse contexts.

Ethical Considerations

Addressing ethical considerations is paramount in the development and deployment of fake news detection systems. These include:

1. Bias and Fairness: Ensuring that models do not perpetuate biases present in the training data. This requires careful dataset curation and the implementation of fairness-aware algorithms.

2. Transparency and Accountability: Providing transparency in how models make decisions and enabling accountability for those decisions. This involves creating explainable AI systems that users can understand and trust.

3. Privacy Protection: Safeguarding user data and ensuring that the deployment of these models does not infringe on privacy rights. Robust data protection measures must be in place.

4. Impact on Free Speech: Balancing the detection and removal of fake news with the protection of free speech. It is crucial to avoid over-censorship and ensure that diverse viewpoints are not unfairly suppressed.

Future Directions

By advancing the field of automated fake news detection, we move closer to developing robust tools that can help mitigate the spread of misinformation and support the integrity of information in the digital age. Future research and collaboration among academia, industry, and policymakers will be critical in achieving these goals. Ensuring that technological advancements are aligned with ethical principles and societal needs is essential for building trust and maintaining the integrity of information.

Final Thoughts

In conclusion, the fight against fake news is a continuous and evolving challenge. The integration of machine learning models in detecting and mitigating the impact of fake news offers a promising path forward. As technology advances, so too must our approaches to ensuring that these tools are used responsibly and ethically. By fostering an environment of collaboration and innovation, we can develop effective solutions that uphold the principles of truth and trust in the digital age.

# APPENDIX

**Source Code:**

**Importing Libraries**

import pandas as pd

import numpy as np

from sklearn.model_selection import train_test_split

from sklearn.metrics import classification_report

import re

import string

import tkinter as tk

from tkinter import messagebox

**Load the fake and true news datasets**

df_fake = pd.read_csv("Fake.csv")

df_true = pd.read_csv("True.csv")

**Display the first 10 rows of each dataset**

df_fake.head(10)

df_true.head(10)

**Assign a class label to the datasets (0 for fake news, 1 for true news)**

df_fake["class"] = 0

df_true["class"] = 1

**Print the shapes of the datasets**

df_fake.shape, df_true.shape


**Save the last 10 rows of each dataset for manual testing and remove these rows from the datasets**

df_fake_mannual_testing = df_fake.tail(10)

for i in range(23480, 23470, -1):

   df_fake.drop([i], axis=0, inplace=True)

df_true_mannual_testing = df_true.tail(10)

for i in range(21416, 21406, -1):

   df_true.drop([i], axis=0, inplace=True)


**Concatenate the manually testing data and save it to a CSV file**

df_mannual_testing = pd.concat([df_fake_mannual_testing, df_true_mannual_testing], axis=0)

df_mannual_testing.to_csv("mannual_testing.csv")


**Merge the remaining fake and true news datasets**

df_merge = pd.concat([df_fake, df_true], axis=0)

df_merge.head(10)


**Drop unnecessary columns**

df = df_merge.drop(["title", "subject", "date"], axis=1)

df.head(10)


**Shuffle the dataset**

df = df.sample(frac=1)

**Function to clean the text data**

```python
def word_drop(text):
    text = text.lower()
    text = re.sub('\[.*?\]', '', text)
    text = re.sub("\\W", " ", text)
    text = re.sub('https?://\S+|www\.\S+', '', text)
    text = re.sub('<.*?>+', '', text)
    text = re.sub('[%s]' % re.escape(string.punctuation), '', text)
    text = re.sub('\n', '', text)
    text = re.sub('\w*\d\w*', '', text)
    return text
```

**Apply the text cleaning function to the text column**

```python
df["text"] = df["text"].apply(word_drop)
```

**Define features (text) and labels (class)**

```python
x = df["text"]
y = df["class"]
```

**Split the dataset into training and testing sets**

```python
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.25)
```

**Vectorize the text data using TF-IDF**

```python
from sklearn.feature_extraction.text import TfidfVectorizer
vectorization = TfidfVectorizer()
xv_train = vectorization.fit_transform(x_train)
xv_test = vectorization.transform(x_test)
```

**Train a Logistic Regression model**

from sklearn.linear_model import LogisticRegression

LR = LogisticRegression()

LR.fit(xv_train, y_train)


**Evaluate the Logistic Regression model**

LR.score(xv_test, y_test)

pred_LR = LR.predict(xv_test)

print(classification_report(y_test, pred_LR))


**Train a Random Forest Classifier**

from sklearn.ensemble import RandomForestClassifier

RFC = RandomForestClassifier(random_state=0)

RFC.fit(xv_train, y_train)


**Evaluate the Random Forest Classifier**

RFC.score(xv_test, y_test)

pred_RFC = RFC.predict(xv_test)

print(classification_report(y_test, pred_RFC))


**Function to output the label based on prediction**

def output_label(n):

   if n == 0:

     return "Fake News"

   elif n == 1:

     return "Not A Fake News"

**Function for manual testing of news articles**

```python
def manual_testing(news, vectorization, LR, RFC, word_drop):
    testing_news = {"text": [news]}
    new_def_test = pd.DataFrame(testing_news)
    new_def_test["text"] = new_def_test["text"].apply(word_drop)
    new_x_test = new_def_test["text"]
    new_xv_test = vectorization.transform(new_x_test)
    pred_LR = LR.predict(new_xv_test)
    pred_RFC = RFC.predict(new_xv_test)


    result = f"LR Prediction: {output_label(pred_LR[0])}\nRFC Prediction: {output_label(pred_RFC[0])}"
    return result
```

**Function to check the news input by the user**

```python
def check_news():
    news = text_entry.get("1.0", "end-1c")
    if news:
        result = manual_testing(news, vectorization, LR, RFC, word_drop)
        messagebox.showinfo("Result", result)
    else:
        messagebox.showwarning("Warning", "Please enter news to check.")

window = tk.Tk()
window.title("Fake News Detection")
text_entry = tk.Text(window, height=10, width=50)
text_entry.pack()
```

```
check_button = tk.Button(window, text="Check News",
command=check_news)
```

```
check_button.pack()
```

```
window.mainloop()
```

# REFERENCES

1. Allcott, H., & Gentzkow, M. (2017). Social Media and Fake News in the 2016 Election. Journal of Economic Perspectives, 31(2), 211-236. doi:10.1257/jep.31.2.211

2. Ruchansky, N., Seo, S., & Liu, Y. (2017). CSI: A Hybrid Deep Model for Fake News Detection. Proceedings of the 2017 ACM on Conference on Information and Knowledge Management, 797-806. doi:10.1145/3132847.3132877

3. Shu, K., Sliva, A., Wang, S., Tang, J., & Liu, H. (2017). Fake News Detection on Social Media: A Data Mining Perspective. ACM SIGKDD Explorations Newsletter, 19(1), 22-36. doi:10.1145/3137597.3137600

4. Kumar, S., West, R., & Leskovec, J. (2016). Disinformation on the Web: Impact, Characteristics, and Detection of Wikipedia Hoaxes. Proceedings of the 25th International Conference on World Wide Web, 591-602. doi:10.1145/2872427.2883085

5. Baly, R., Karadzhov, G., Alexandrov, D., Glass, J., & Nakov, P. (2018). Predicting Factuality of Reporting and Bias of News Media Sources. Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, 3528-3539. doi:10.18653/v1/D18-1389

6. Conroy, N. K., Rubin, V. L., & Chen, Y. (2015). Automatic Deception Detection: Methods for Finding Fake News. Proceedings of the 78th ASIS&T Annual Meeting: Information Science with Impact: Research in and for the Community, 82. doi:10.1002/pra2.2015.145052010082

7. Vosoughi, S., Roy, D., & Aral, S. (2018). The spread of true and false news online. Science, 359(6380), 1146-1151. doi:10.1126/science.aap9559

8. Zhou, X., & Zafarani, R. (2018). Fake News: A Survey of Research, Detection Methods, and Opportunities. arXiv preprint arXiv:1812.00315.

9. Wang, W. Y. (2017). "Liar, Liar Pants on Fire": A New Benchmark Dataset for Fake News Detection. Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, 422-426. doi:10.18653/v1/P17-2067

10. Horne, B. D., & Adali, S. (2017). This Just In: Fake News Packs a Lot in Title, Uses Simpler, Repetitive Content in Text Body, More Similar to Satire than Real News. arXiv preprint arXiv:1703.09398.

11. Rubin, V. L., Chen, Y., & Conroy, N. K. (2015). Deception Detection for News: Three Types of Fakes. Proceedings of the 78th ASIS&T Annual Meeting: Information Science with Impact: Research in and for the Community, 83. doi:10.1002/pra2.2015.145052010083

12. Peisong, W., & Wong, R. K. (2018). Early Detection of Fake News on Social Media Through Propagation Path Classification with Recurrent and Convolutional Networks. Proceedings of the AAAI Conference on Artificial Intelligence, 32(1). doi:10.1609/aaai.v32i1.11604

13. Ghanem, B., Rosso, P., & Rangel, F. (2020). An Emotional Analysis of False Information in Social Media and News Articles. ACM Transactions on Internet Technology, 20(2), 1-18. doi:10.1145/3386470

14. Kwak, H., & An, J. (2016). Revealing the Hidden Patterns of News Photos: Analysis of a Large-Scale News Photo Dataset. Proceedings of the 2016 ACM on International Conference on Multimedia Retrieval, 213-217. doi:10.1145/2911996.2912027

15. Yang, Y., Zheng, L., Zhang, J., Cui, Q., Li, Z., & Yu, P. S. (2018). TI-CNN: Convolutional Neural Networks for Fake News Detection. Proceedings of the

25th International Conference on World Wide Web, 703-712. doi:10.1145/3184558.3191535

16. Tacchini, E., Ballarin, G., Della Vedova, M. L., Moret, S., & de Alfaro, L. (2017). Some Like it Hoax: Automated Fake News Detection in Social Networks. arXiv preprint arXiv:1704.07506.

17. Rashkin, H., Choi, E., Jang, J. Y., Volkova, S., & Choi, Y. (2017). Truth of Varying Shades: Analyzing Language in Fake News and Political Fact-Checking. Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, 2931-2937. doi:10.18653/v1/D17-1317

18. Giachanou, A., & Crestani, F. (2016). Like It or Not: A Survey of Twitter Sentiment Analysis Methods. ACM Computing Surveys, 49(2), 28. doi:10.1145/2938640

19. Shu, K., Wang, S., & Liu, H. (2019). Beyond News Contents: The Role of Social Context for Fake News Detection. Proceedings of the 12th ACM International Conference on Web Search and Data Mining, 312-320. doi:10.1145/3289600.3290994

20. Jin, Z., Cao, J., Guo, H., Zhang, Y., & Luo, J. (2017). Multimodal Fusion with Recurrent Neural Networks for Rumor Detection on Microblogs. Proceedings of the 25th ACM International Conference on Multimedia, 795-816. doi:10.1145/3123266.3123454

21. Potthast, M., Kiesel, J., Reinartz, K., Bevendorff, J., & Stein, B. (2018). A Stylometric Inquiry into Hyperpartisan and Fake News. Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, 231-240. doi:10.18653/v1/P18-1022

22. Tschiatschek, S., Singla, A., Gomez Rodriguez, M., Merchant, A., & Krause, A. (2018). Fake News Detection in Social Networks via Crowd Signals.

Proceedings of the 2018 World Wide Web Conference, 517-526. doi:10.1145/3178876.3186041

23. Karimi, H., Roy, P., Saba-Sadiya, S., & Tang, J. (2018). Multi-Source Multi-Class Fake News Detection. Proceedings of the 27th ACM International Conference on Information and Knowledge Management, 1127-1136. doi:10.1145/3269206.3271709

24. Guo, L., Vargo, C. J., Pan, Z., Ding, Y., & Li, X. (2018). Fake News on Social Media: A Data-Driven Approach to Identify Information Sharing Behaviors. Journal of Information Science, 44(4), 485-499. doi:10.1177/0165551518773917

25. Avaaz. (2019). How Fake News on WhatsApp Helped Turn an Election. Retrieved from https://secure.avaaz.org/page/en/

26. Howard, P. N., Kollanyi, B., Bradshaw, S., & Neudert, L. M. (2018). Social Media, News and Political Information during the US Election: Was Polarizing Content Concentrated in Swing States? International Journal of Press/Politics, 23(2), 132-153. doi:10.1177/1940161218771902

27. Zollo, F., Bessi, A., Del Vicario, M., Scala, A., Caldarelli, G., Shekhtman, L., ... & Quattrociocchi, W. (2015). Debunking in a World of Tribes. PloS one, 10(7), e0131380. doi:10.1371/journal.pone.0131380

28. Pennycook, G., & Rand, D. G. (2019). Fighting misinformation on social media using crowdsourced judgments of news source quality. Proceedings of the National Academy of Sciences, 116(7), 2521-2526. doi:10.1073/pnas.1806781116

29. Zhou, X., & Zafarani, R. (2020). A Survey of Fake News: Fundamental Theories, Detection Methods, and Opportunities. ACM Computing Surveys (CSUR), 53(5), 1-40. doi:10.1145/3395046

30. Mitra, T., & Gilbert, E. (2015). Credbank: A large-scale social media corpus with associated credibility annotations. Proceedings of the Ninth International Conference on Web and Social Media, 258-267.

31. Pérez-Rosas, V., Kleinberg, B., Lefevre, A., & Mihalcea, R. (2018). Automatic detection of fake news. Proceedings of the 27th International Conference on Computational Linguistics, 3391-3401. doi:10.18653/v1/C18-1289

32. Shu, K., Mahudeswaran, D., Wang, S., Lee, D., & Liu, H. (2020). Fakenewsnet: A data repository with news content, social context, and spatial-temporal information for studying fake news on social media. Big Data, 8(3), 171-188. doi:10.1089/big.2020.0062

33. Shao, C., Ciampaglia, G. L., Varol, O., Yang, K. C., Flammini, A., & Menczer, F. (2018). The spread of fake news by social bots. Nature Communications, 9(1), 1-9. doi:10.1038/s41467-018-06930-7

34. Lazer, D. M., Baum, M. A., Benkler, Y., Berinsky, A. J., Greenhill, K. M., Menczer, F., ... & Schudson, M. (2018). The science of fake news. Science, 359(6380), 1094-1096. doi:10.1126/science.aao2998

35. Ferrara, E., Varol, O., Davis, C., Menczer, F., & Flammini, A. (2016). The rise of social bots. Communications of the ACM, 59(7), 96-104. doi:10.1145/2818717

36. Jiang, S., & Wilson, C. (2018). Linguistic signals under misinformation and fact-checking: Evidence from user comments on social media. Proceedings of the ACM on Human-Computer Interaction, 2(CSCW), 1-23. doi:10.1145/3274334

37. Ruchansky, N., Seo, S., & Liu, Y. (2017). CSI: A Hybrid Deep Model for Fake News Detection. Proceedings of the 2017 ACM on Conference on Information and Knowledge Management, 797-806. doi:10.1145/3132847.3132877

38. Li, H., Su, L., Shen, C., & Zhang, B. (2019). Detecting fake news on social media with multiple complementary features. ACM Transactions on Information Systems (TOIS), 37(3), 1-30. doi:10.1145/3362028

39. Nørregaard, J., Horne, B. D., & Adalı, S. (2019). NELA-GT-2018: A large multi-labelled news dataset for the study of misinformation in news articles. Proceedings of the International AAAI Conference on Web and Social Media, 13, 630-638.

40. Pennycook, G., & Rand, D. G. (2017). The Implied Truth Effect: Attaching Warnings to a Subset of Fake News Stories Increases Perceived Accuracy of Stories Without Warnings. Management Science, 66(11), 4944-4957. doi:10.1287/mnsc.2019.3478

41. Chen, Y., Conroy, N. J., & Rubin, V. L. (2015). Misleading Online Content: Recognizing Clickbait as "False News". Proceedings of the 2015 ACM on Workshop on Multimodal Deception Detection, 15-19. doi:10.1145/2823465.2823467

42. Jin, Z., Cao, J., Zhang, Y., Zhou, J., & Tian, Q. (2016). Novel visual and statistical image features for microblogs news verification. IEEE Transactions on Multimedia, 19(3), 598-608. doi:10.1109/TMM.2016.2617078

43. Goswami, S., & Kumar, A. (2016). Detecting Fake News Using Machine Learning. arXiv preprint arXiv:1612.01340.

44. Ribeiro, M. H., Calais, P. H., Santos, Y. A., Almeida, V. A., & Meira, W. (2018). Characterizing and detecting hateful users on Twitter. Proceedings of the International AAAI Conference on Web and Social Media, 12(1), 676-679.

45. Kumar, S., Hamilton, W. L., Leskovec, J., & Jurafsky, D. (2018). Community Interaction and Conflict on the Web. Proceedings of the 2018 World Wide Web Conference, 933-943. doi:10.1145/3178876.3186134

46. Pennycook, G., & Rand, D. G. (2018). The Implied Truth Effect: Attaching Warnings to a Subset of Fake News Stories Increases Perceived Accuracy of Stories Without Warnings. Management Science, 66(11), 4944-4957. doi:10.1287/mnsc.2019.3478

47. Zhang, X., & Ghorbani, A. A. (2020). An overview of online fake news: Characterization, detection, and discussion. Information Processing & Management, 57(2), 102025. doi:10.1016/j.ipm.2019.03.004

48. Vosoughi, S., Roy, D., & Aral, S. (2018). The spread of true and false news online.** Science, 359(6380), 1146-1151. doi:10.1126/science.aap9559

49. **Shu, K., Sliva, A., Wang, S., Tang, J., & Liu, H. (2017). Fake News Detection on Social Media: A Data Mining Perspective. ACM SIGKDD Explorations Newsletter, 19(1), 22-36. doi:10.1145/3137597.3137600

50. Zubiaga, A., Aker, A., Bontcheva, K., Liakata, M., & Procter, R. (2018). Detection and resolution of rumours in social media: A survey. ACM Computing Surveys (CSUR), 51(2), 1-36. doi:10.1145/3161603