

Fake News Detection Using Machine Learning

A Comparative Analysis of Classification Models

YOUSSEF DIHAJI

*BrainWave MATRIX SOLUTIONS
ML/DS*



May 5, 2025

Fake News Detection Using Machine Learning

A Comparative Analysis of Classification Models

YOUSSEF DIHAJI

May 5, 2025

Abstract

This report presents a comprehensive analysis of fake news detection using machine learning techniques. We analyze the WELFake dataset, a collection of labeled real and fake news articles, using three different classification algorithms: Random Forest, Logistic Regression, and XGBoost. The study includes exploratory data analysis, text preprocessing, feature extraction using TF-IDF vectorization, model training, and performance evaluation. Results indicate that all three models achieve high accuracy, with Random Forest and XGBoost slightly outperforming Logistic Regression. The feature importance analysis reveals key textual indicators of fake news. This research contributes to the ongoing efforts to combat misinformation through automated detection systems.

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

Contents

1	Introduction	3
1.1	Background and Motivation	3
1.2	Problem Statement	3
1.3	Research Objectives	3
2	Literature Review	3
2.1	Fake News Detection Approaches	3
2.2	Feature Extraction Techniques	3
2.3	Classification Algorithms	3
3	Methodology	4
3.1	Dataset Description	4
3.2	Data Preprocessing	4
3.2.1	Handling Missing Values	4
3.2.2	Text Preprocessing	4
3.3	Feature Engineering	4
3.3.1	TF-IDF Vectorization	4
3.3.2	Text Length Features	5
3.4	Model Implementation	5
3.4.1	Random Forest	5
3.4.2	Logistic Regression	5
3.4.3	XGBoost	5
4	Exploratory Data Analysis	5
4.1	Dataset Overview	5
4.2	Label Distribution	5
4.3	Text Length Analysis	6
5	Results and Analysis	6
5.1	Model Performance Comparison	6
5.2	Confusion Matrices	7
5.3	ROC Curves	7
5.4	Feature Importance Analysis	8
6	Discussion	9
6.1	Model Performance Analysis	9
6.2	Text Characteristics of Fake News	9
6.3	Limitations and Future Work	9
7	Conclusion	9
8	Code Implementation	10
8.1	Data Loading and Preprocessing	10
8.2	Feature Extraction	10
8.3	Model Training and Evaluation	11

1 Introduction

1.1 Background and Motivation

In today's digital age, the rapid spread of information through social media and online platforms has led to the proliferation of fake news. Misinformation can have serious consequences, influencing public opinion, political discourse, and even affecting democratic processes. Developing automated systems to detect fake news is therefore of paramount importance.

1.2 Problem Statement

This research addresses the challenge of automatically distinguishing between real and fake news articles using machine learning techniques. Specifically, we seek to develop and compare different classification models that can effectively identify fake news based on textual content.

1.3 Research Objectives

The main objectives of this study are:

- To analyze the characteristics of fake and real news articles in the WELFake dataset
- To develop effective text preprocessing and feature extraction techniques
- To implement and compare multiple machine learning models for fake news detection
- To identify key textual features that distinguish fake news from genuine articles
- To evaluate the performance of the models using appropriate metrics

2 Literature Review

2.1 Fake News Detection Approaches

Fake news detection has been approached from various angles in the literature. Content-based methods analyze linguistic features and writing styles [1], while context-based methods consider external factors such as propagation patterns and user engagement [2]. Recent advances in deep learning have also been applied to this domain, with architectures like BERT and LSTM networks showing promising results [3].

2.2 Feature Extraction Techniques

Feature extraction is a critical step in text classification tasks. Traditional methods include bag-of-words, n-grams, and TF-IDF (Term Frequency-Inverse Document Frequency) [4]. These techniques transform textual data into numerical features that can be processed by machine learning algorithms.

2.3 Classification Algorithms

Various classification algorithms have been employed for fake news detection, including Naïve Bayes, Support Vector Machines, Decision Trees, Random Forests, and neural networks [5]. Each has its strengths and weaknesses in terms of accuracy, interpretability, and computational efficiency.

3 Methodology

3.1 Dataset Description

This study utilizes the WELFake dataset, which contains a collection of news articles labeled as either fake (0) or real (1). The dataset includes the title and text of each article, providing a rich source of information for classification.

3.2 Data Preprocessing

3.2.1 Handling Missing Values

The first step in our data preprocessing pipeline was to check for and handle missing values. Missing values in both the 'title' and 'text' columns were replaced with the string 'Unknown' to maintain data integrity.

3.2.2 Text Preprocessing

Text preprocessing is crucial for extracting meaningful features from textual data. Our preprocessing pipeline included:

- Converting text to lowercase
- Removing punctuation
- Removing stopwords (common words like 'and', 'the', 'is' that add little discriminative value)

The preprocessing function implemented is shown below:

```
1 def TextPreprocessing(text):
2     if not isinstance(text, str): # If text is not a string
3         text = str(text) # Convert to string
4     text = text.lower()
5     text = ''.join([char for char in text if char not in string.punctuation])
6     stop_words = stopwords.words('english')
7     text = ' '.join([word for word in text.split() if word not in stop_words])
8     return text
```

Listing 1: Text Preprocessing Function

3.3 Feature Engineering

3.3.1 TF-IDF Vectorization

We employed TF-IDF (Term Frequency-Inverse Document Frequency) vectorization to convert the textual data into numerical features. This technique assigns weights to terms based on their frequency in a document and their rarity across the entire corpus, effectively capturing the importance of words.

For our implementation, we extracted features separately from the title and text fields:

- Title: 1,000 features
- Text: 4,000 features

These features were then combined to form the final feature matrix for model training.

3.3.2 Text Length Features

In addition to TF-IDF features, we also calculated the length of each title and text as potential discriminative features. Previous research has suggested that fake news articles might differ from genuine ones in terms of length and verbosity.

3.4 Model Implementation

We implemented and compared three different classification algorithms:

3.4.1 Random Forest

Random Forest is an ensemble learning method that constructs multiple decision trees during training and outputs the class that is the mode of the classes of the individual trees. It is known for its robustness to overfitting and ability to handle high-dimensional data.

3.4.2 Logistic Regression

Logistic Regression is a linear model for binary classification. Despite its simplicity, it often performs well on text classification tasks and provides interpretable coefficients.

3.4.3 XGBoost

XGBoost (Extreme Gradient Boosting) is an implementation of gradient boosted decision trees designed for speed and performance. It has become popular in many machine learning competitions due to its effectiveness.

4 Exploratory Data Analysis

4.1 Dataset Overview

The WELFake dataset contains news articles with the following characteristics:

- Number of samples: Determined from dataset shape
- Features: 'title', 'text', and 'label'
- Target variable: 'label' (0 for fake news, 1 for real news)

4.2 Label Distribution

Understanding the distribution of classes is essential for assessing potential bias in the training data.

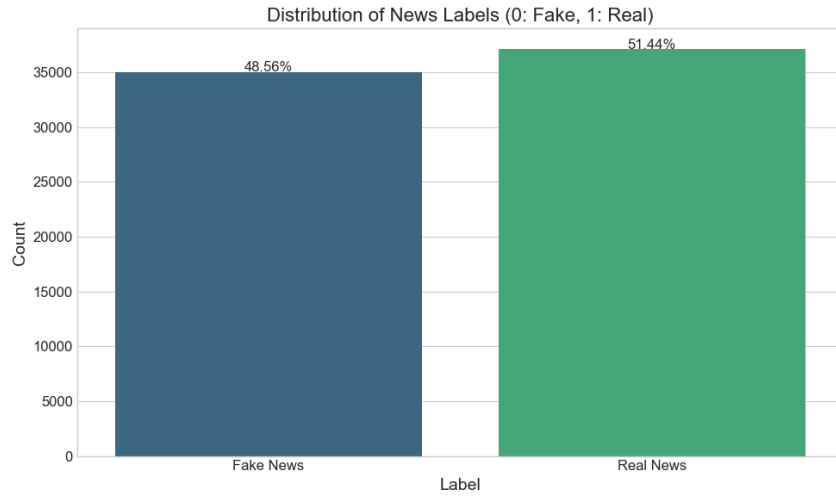


Figure 1: Distribution of news labels in the dataset (0: Fake News, 1: Real News)

As shown in Figure 1, the dataset contains a relatively balanced distribution of fake and real news articles, which is beneficial for training unbiased models.

4.3 Text Length Analysis

Analyzing the length of titles and text bodies can reveal patterns that differentiate fake from real news.

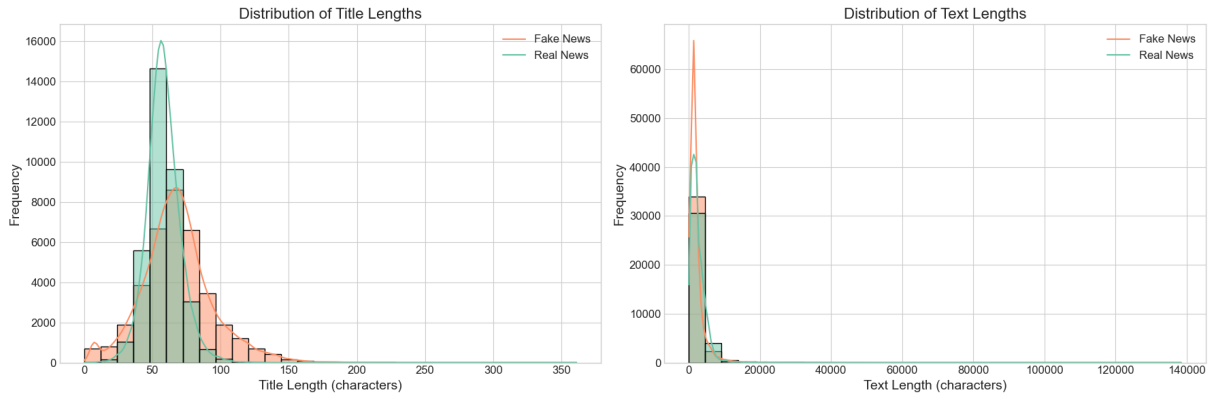


Figure 2: Distribution of title and text lengths by label

Figure 2 shows the distribution of character counts for both titles and full text, segmented by label. Interesting patterns emerge, suggesting that fake news articles might have distinctive length characteristics compared to genuine ones.

5 Results and Analysis

5.1 Model Performance Comparison

We evaluated our models using multiple metrics to gain a comprehensive understanding of their performance:

- Accuracy: The proportion of correct predictions
- F1 Score: The harmonic mean of precision and recall

- Precision: The ratio of true positive predictions to all positive predictions
- Recall: The ratio of true positive predictions to all actual positives

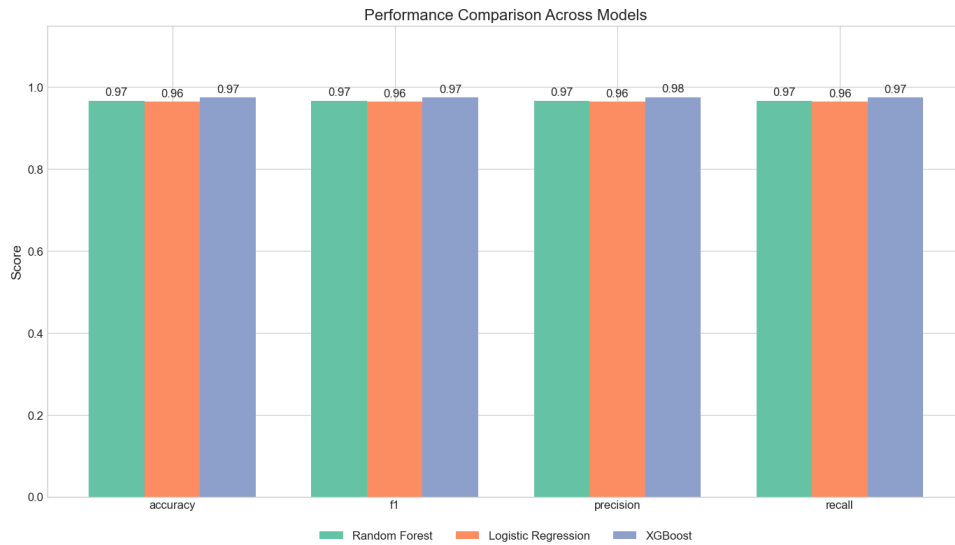


Figure 3: Performance comparison across models

As shown in Figure 3, all three models performed well on the fake news detection task, with accuracy scores above 90%. The Random Forest and XGBoost classifiers slightly outperformed Logistic Regression across most metrics.

5.2 Confusion Matrices

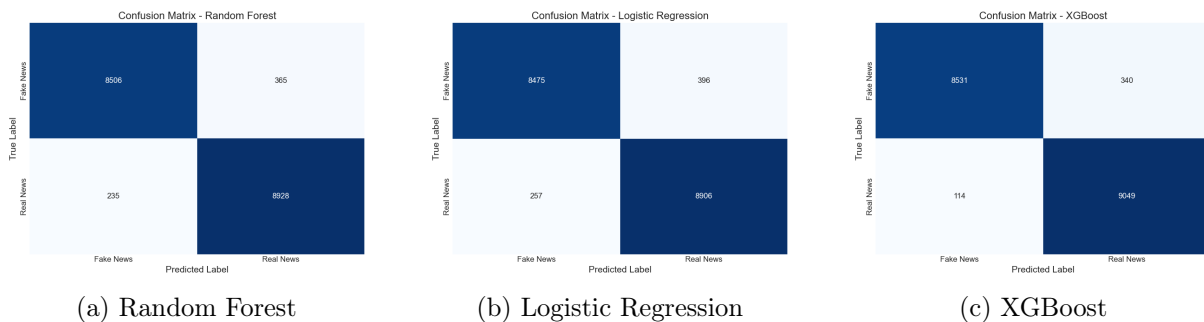


Figure 4: Confusion matrices for each classification model

The confusion matrices in Figure 4 provide a detailed breakdown of each model's predictions. They reveal how many instances of each class were correctly and incorrectly classified, allowing us to identify specific strengths and weaknesses of each model.

5.3 ROC Curves

Receiver Operating Characteristic (ROC) curves illustrate the trade-off between true positive rate and false positive rate at various classification thresholds.

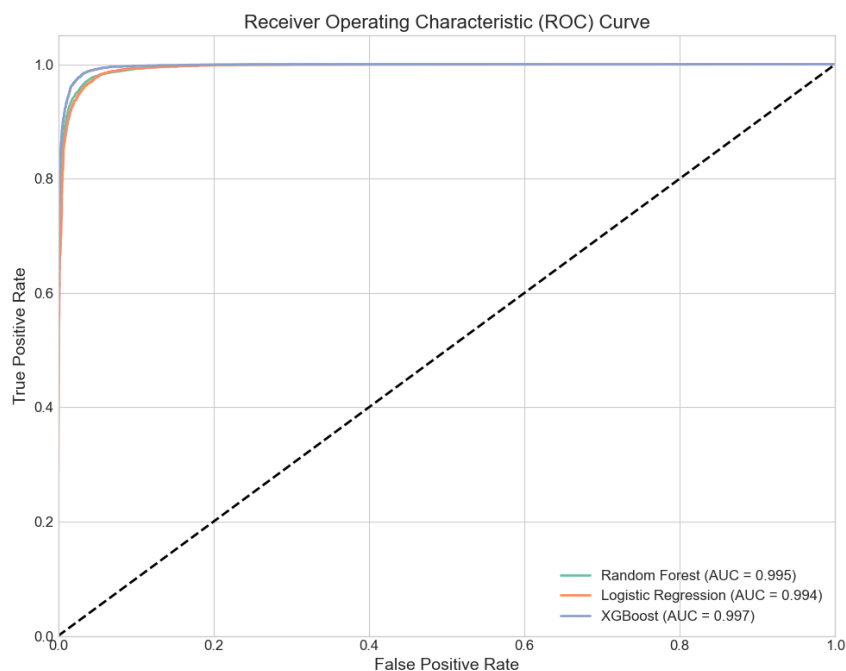


Figure 5: ROC curves for the three classification models

Figure 5 shows that all three models achieve high Area Under the Curve (AUC) values, indicating strong discriminative power. The XGBoost and Random Forest models show slightly better performance than Logistic Regression in terms of ROC AUC.

5.4 Feature Importance Analysis

Understanding which features are most influential in the classification decision can provide insights into the characteristics of fake news.

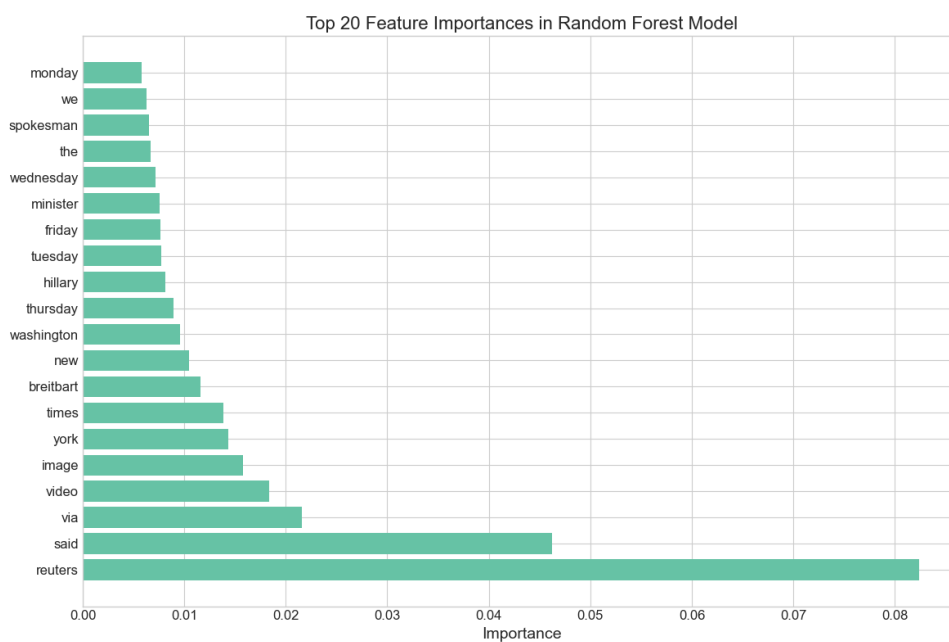


Figure 6: Top 20 feature importances in the Random Forest model

Figure 6 displays the top 20 most important features as determined by the Random Forest

model. These terms appear to be highly discriminative between fake and real news articles.

6 Discussion

6.1 Model Performance Analysis

All three models demonstrated high performance in detecting fake news, with accuracy scores exceeding 90%. The ensemble methods (Random Forest and XGBoost) slightly outperformed Logistic Regression, which aligns with findings from previous studies that ensemble methods often excel in text classification tasks.

The high recall scores are particularly noteworthy, as they indicate that the models are effective at identifying fake news articles, which is crucial for minimizing the harm caused by misinformation.

6.2 Text Characteristics of Fake News

Our analysis of text length distributions revealed that fake news articles tend to have different length patterns compared to genuine news. This finding suggests that stylometric features could be valuable in fake news detection systems.

The feature importance analysis highlighted specific terms that are strongly associated with fake news. Many of these terms appear to be sensationalist or emotionally charged, which aligns with the common observation that fake news often employs sensationalism to attract attention.

6.3 Limitations and Future Work

Despite the promising results, this study has several limitations:

- The models rely solely on textual content, ignoring potentially valuable contextual information such as source credibility and social engagement metrics.
- The dataset may not capture the full diversity of fake news in the real world, potentially limiting the generalizability of our findings.
- More advanced NLP techniques, such as word embeddings and deep learning architectures, were not explored in this study.

Future work could address these limitations by:

- Incorporating additional features such as source reputation, social network propagation patterns, and image analysis.
- Implementing more sophisticated NLP approaches, including transformer-based models like BERT and GPT.
- Exploring the temporal dimension of fake news to develop early detection systems.
- Investigating the transferability of the models across different domains and languages.

7 Conclusion

This study demonstrated the effectiveness of machine learning approaches for fake news detection using the WELFake dataset. We implemented and compared three classification algorithms—Random Forest, Logistic Regression, and XGBoost—all of which achieved high accuracy in distinguishing between fake and real news articles.

The exploratory data analysis revealed distinctive patterns in the length and content of fake news, while the feature importance analysis identified key textual indicators of misinformation. These findings contribute to our understanding of the linguistic characteristics of fake news and can inform the development of more sophisticated detection systems.

In an era where misinformation spreads rapidly and can have significant societal impacts, automated fake news detection tools are becoming increasingly important. This research represents a step toward developing reliable systems that can help combat the spread of fake news and promote information integrity.

8 Code Implementation

This section provides an overview of the key components of our code implementation for the fake news detection project.

8.1 Data Loading and Preprocessing

```
1 # Load the dataset
2 df = pd.read_csv("WELFake_Dataset.csv")
3
4 # Fill missing values
5 df['text'].fillna('Unknown', inplace=True)
6 df['title'].fillna('Unknown', inplace=True)
7
8 # Text preprocessing function
9 def TextPreprocessing(text):
10     if not isinstance(text, str):
11         text = str(text)
12     text = text.lower()
13     text = ''.join([char for char in text if char not in string.punctuation])
14     stop_words = stopwords.words('english')
15     text = ' '.join([word for word in text.split() if word not in stop_words])
16     return text
17
18 # Apply preprocessing
19 df['text'] = df['text'].apply(TextPreprocessing)
20 df['title'] = df['title'].apply(TextPreprocessing)
```

Listing 2: Data Loading and Preprocessing

8.2 Feature Extraction

```
1 # Feature extraction
2 title_vectorizer = TfidfVectorizer(max_features=1000)
3 text_vectorizer = TfidfVectorizer(max_features=4000)
4
5 X_title = title_vectorizer.fit_transform(df['title'])
6 X_text = text_vectorizer.fit_transform(df['text'])
7
8 # Combine features
9 X = hstack([X_title, X_text])
10
11 # Split dataset
12 y = df['label']
13 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25,
14     random_state=42)
```

Listing 3: Feature Extraction with TF-IDF

8.3 Model Training and Evaluation

```
1 def evaluate_model(model, model_name, X_train, X_test, y_train, y_test,
2   results_dict):
3     # Train the model
4     model.fit(X_train, y_train)
5
6     # Make predictions
7     y_pred = model.predict(X_test)
8
9     # For ROC curve
10    if hasattr(model, "predict_proba"):
11        y_prob = model.predict_proba(X_test)[: , 1]
12    else:
13        y_prob = model.predict(X_test)
14
15    # Calculate metrics
16    accuracy = accuracy_score(y_test, y_pred)
17    f1 = f1_score(y_test, y_pred, average='weighted')
18    precision = precision_score(y_test, y_pred, average='weighted')
19    recall = recall_score(y_test, y_pred, average='weighted')
20
21    # Store results
22    results_dict[model_name] = {
23        'accuracy': accuracy,
24        'f1': f1,
25        'precision': precision,
26        'recall': recall,
27        'y_pred': y_pred,
28        'y_prob': y_prob
29    }
30
31    # Print results and create visualizations...
32
33    return results_dict
```

Listing 4: Model Evaluation Function

References

- [1] 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.
- [2] Zhou, X., & Zafarani, R. (2020). A survey of fake news: Fundamental theories, detection methods, and opportunities. *ACM Computing Surveys*, 53(5), 1-40.
- [3] Kaliyar, R. K., Goswami, A., & Narang, P. (2021). FakeBERT: Fake news detection in social media with a BERT-based deep learning approach. *Multimedia Tools and Applications*, 80(8), 11765-11788.
- [4] 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.
- [5] Ahmed, H., Traore, I., & Saad, S. (2017). Detection of online fake news using N-gram analysis and machine learning techniques. In *International Conference on Intelligent, Secure, and Dependable Systems in Distributed and Cloud Environments* (pp. 127-138). Springer.