

# Mid-Term Project Report: Fake News Detection System

---

Course Title: Machine Learning

Instructor: Gulshan Yasmeen

Institution: NAVTAC AI Training Program

Student Name: Sheryar Sher

Project Title: Fake News Detection System

Submission Date: October 25, 2024

## **Table of Contents**

A. Project Proposal

B. Data Mining and Exploration

C. Data Preprocessing

D. Data Visualization

E. Model Development

F. Model Evaluation

G. Conclusion and Recommendations

H. Documentation

## A. Project Proposal

### Title

Fake News Detection System: A Machine Learning Approach to Identify Misinformation

### Problem Statement

In today's digital age, the rapid spread of misinformation and fake news has become a significant challenge affecting public opinion, political discourse, and social stability. With the exponential growth of social media platforms and online news consumption, distinguishing between authentic and fabricated news articles has become increasingly difficult for the average reader. This project addresses the critical need for automated systems that can accurately identify fake news articles, helping users make more informed decisions about the information they consume.

### Objectives

The primary objectives of this project are:

1. Develop an accurate machine learning model that can distinguish between real and fake news articles with high precision
2. Implement comprehensive text preprocessing techniques including stemming, stopword removal, and feature extraction
3. Compare multiple machine learning algorithms including Logistic Regression, Naive Bayes, Random Forest, and SVM to identify the best-performing approach
4. Create an interactive web application using Streamlit that allows users to input news text and receive real-time predictions
5. Achieve high accuracy and reliability in fake news detection to make the system practically useful
6. Develop comprehensive evaluation tools and visualization functions for model assessment

### Dataset Description

The project utilizes a comprehensive fake news dataset with the following characteristics:

- Source: Publicly available dataset for fake news detection research
- Size: 24,353 training samples, 8,117 test samples, and 8,117 evaluation samples
- Features: id (unique identifier), title (headline), text (article content), label (0=Real, 1=Fake)
- Class Distribution: Approximately balanced with 13,246 fake news articles and 11,107 real news articles

- Data Quality: Clean dataset with no missing values, ensuring reliable model training

df  
✓ 0.1s

		title	text	label
0		Palestinians switch off Christmas lights in Be...	RAMALLAH, West Bank (Reuters) - Palestinians s...	1
1		China says Trump call with Taiwan president wo...	BEIJING (Reuters) - U.S. President-elect Donal...	1
2		FAIL! The Trump Organization's Credit Score W...	While the controversy over Trump s personal ta...	0
3		Zimbabwe military chief's China trip was norma...	BEIJING (Reuters) - A trip to Beijing last wee...	1
4		THE MOST UNCOURAGEOUS PRESIDENT EVER Receives ...	There has never been a more UNCOURAGEOUS perso...	0
...		...	...	...
24348		Mexico Senate committee OK's air transport dea...	MEXICO CITY (Reuters) - A key committee in Mex...	1
24349		BREAKING: HILLARY CLINTON'S STATE DEPARTMENT G...	IF SHE S NOT TOAST NOW THEN WE RE IN BIGGER TR...	0
24350		trump breaks from stump speech to admire beaut...	kremlin nato was created for agresion \r\nru...	0
24351		NFL PLAYER Delivers Courageous Message: Stop B...	Dallas Cowboys star wide receiver Dez Bryant t...	0
24352		NORDSTROM STOCK TAKES NOSEDIVE After Trump Twe...	UPDATE: Nordstrom stock closed up slightly tod...	0

Dataset insights:

```

... Dataset Shape: (24353, 4)
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 24353 entries, 0 to 24352
Data columns (total 4 columns):
#   Column  Non-Null Count  Dtype
---  -
0    id      24353 non-null   int64
1   title   24353 non-null   object
2   text    24353 non-null   object
3   label   24353 non-null   int64
dtypes: int64(2), object(2)
memory usage: 761.2+ KB
...

```

## B. Data Mining and Exploration

### Initial Data Analysis

The dataset exploration revealed several key insights through comprehensive analysis:

#### Dataset Structure

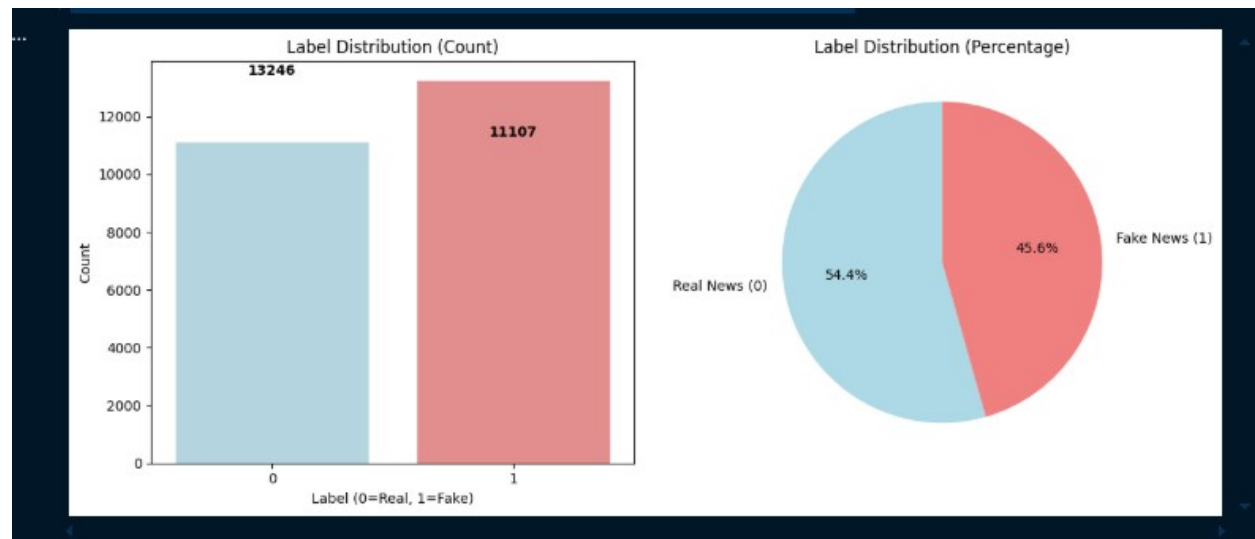
- Training Set: 24,353 articles (54.4% fake, 45.6% real)
- Test Set: 8,117 articles (53.8% fake, 46.2% real)
- Evaluation Set: 8,117 articles (53.1% fake, 46.9% real)

#### Data Quality Assessment

- No missing values detected across all features
- Consistent data format across all three datasets
- Text content varies significantly in length and complexity
- Memory usage: 79.4 MB for text content, 3.7 MB for titles

	id	label
count	24353.000000	24353.000000
mean	12176.000000	0.543917
std	7030.249889	0.498078
min	0.000000	0.000000
25%	6088.000000	0.000000
50%	12176.000000	1.000000
75%	18264.000000	1.000000
max	24352.000000	1.000000

### Class Representation



## Key Findings from Exploration

1. **Balanced Dataset:** The relatively balanced class distribution prevents bias toward either real or fake news classification
2. **Rich Text Content:** The combination of titles and article text provides comprehensive information for analysis
3. **Diverse Topics:** The dataset covers various domains, making the model more robust and generalizable
4. **Clean Data:** No preprocessing required for missing values, allowing focus on feature engineering

## C. Data Preprocessing

### Text Preprocessing Pipeline

A comprehensive text preprocessing pipeline was implemented to prepare the data for machine learning:

#### 1. Data Cleaning

- Removed the id column as it's not relevant for classification
- Handled any potential null values (none found in this dataset)
- Combined title and text columns to create a comprehensive content field

#### 2. Text Normalization

The following preprocessing steps were applied:

- Removed non-alphabetic characters using regular expressions
- Converted all text to lowercase for consistency
- Applied Porter Stemmer to reduce words to their root forms
- Removed common English stopwords using NLTK library

#### 3. Feature Engineering

- Stemming: Applied Porter Stemmer to reduce words to their root forms
- Stopword Removal: Eliminated common English stopwords using NLTK
- Text Vectorization: Used CountVectorizer with:
  - Maximum features: 5,000
  - N-gram range: (1, 3) to capture unigrams, bigrams, and trigrams
  - This approach captures both individual words and word combinations

#### 4. Data Splitting

- Training set: 80% of the data (19,482 samples)
- Test set: 20% of the data (4,871 samples)
- Random state: 2 (for reproducibility)

## **D. Data Visualization**

### **Model Performance Visualization**

The project includes comprehensive visualizations to analyze model performance:

#### **1. Confusion Matrices: Generated for each model to visualize:**

- True Positives (correctly identified fake news)
- True Negatives (correctly identified real news)
- False Positives (real news misclassified as fake)
- False Negatives (fake news misclassified as real)

#### **2. Performance Comparison Charts:**

- Bar charts comparing accuracy across different algorithms
- Detailed classification reports showing precision, recall, and F1-scores
- Visual representation of model strengths and weaknesses

#### **3. Model Evaluation Metrics:**

- Accuracy scores for each algorithm
- Precision and recall for both classes
- F1-scores for balanced performance assessment



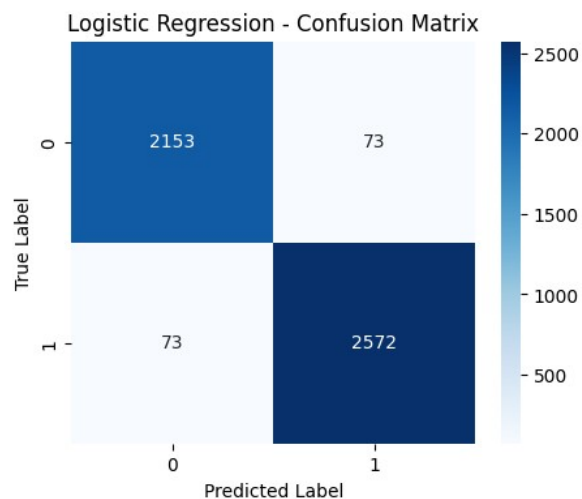
## E. Model Development

### Algorithm Selection and Implementation

Four different machine learning algorithms were implemented and compared using a systematic evaluation approach:

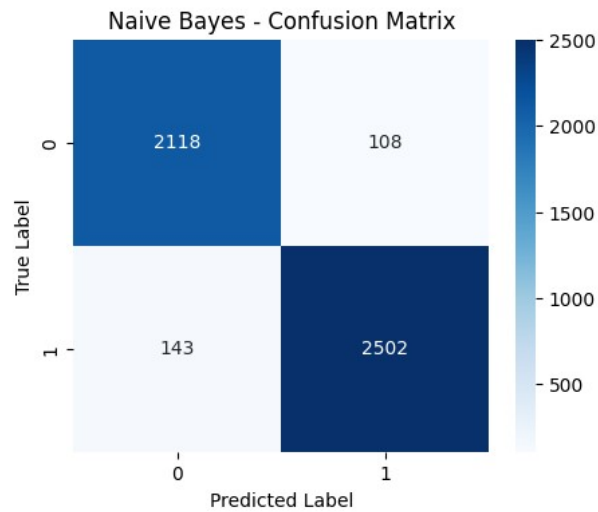
#### 1. Logistic Regression

- Rationale: Simple, interpretable, and effective for binary classification
- Configuration: Maximum iterations set to 1000 for convergence
- Performance: 97.00% accuracy



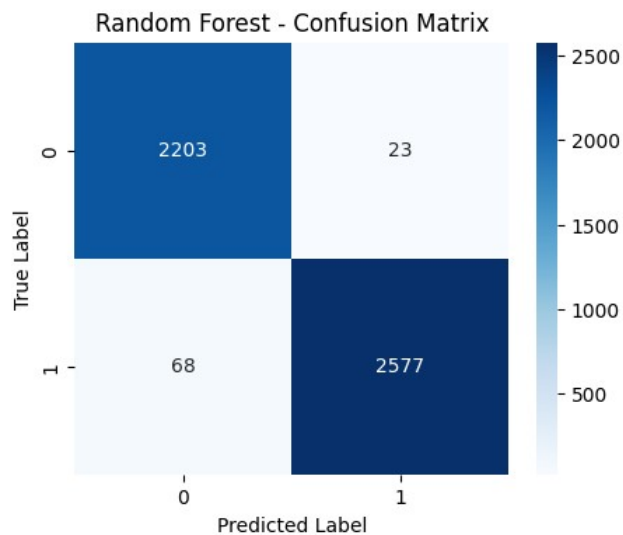
#### 2. Naive Bayes (Multinomial)

- Rationale: Excellent for text classification due to its probabilistic approach
- Configuration: Default parameters with multinomial distribution
- Performance: 94.85% accuracy



### 3. Random Forest

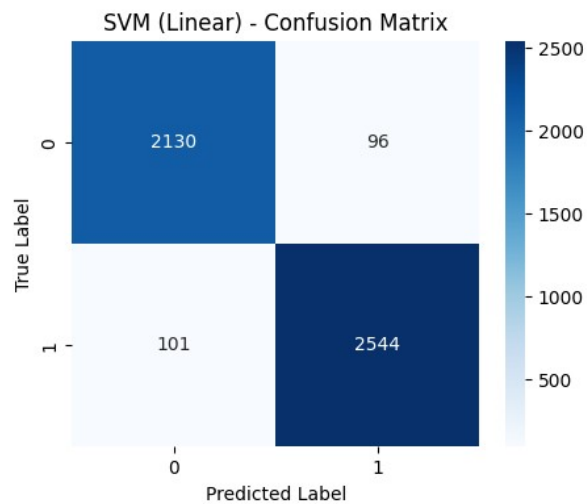
- Rationale: Robust ensemble method that handles overfitting well
- Configuration: 100 estimators, random state 42
- Performance: 97.95% accuracy (best performing)



### 4. Support Vector Machine (SVM)

- Rationale: Effective for high-dimensional text data

- Configuration: Linear kernel for efficiency
- Performance: 95.96% accuracy



## Model Training Implementation

### Model Training Process

1. Feature Extraction: Converted preprocessed text to numerical features using CountVectorizer
2. Model Training: Each algorithm was trained on the training set using systematic evaluation
3. Performance Comparison: All models evaluated using consistent metrics and visualization
4. Model Selection: Random Forest was selected as the best-performing model

## F. Model Evaluation

### Performance Metrics Analysis

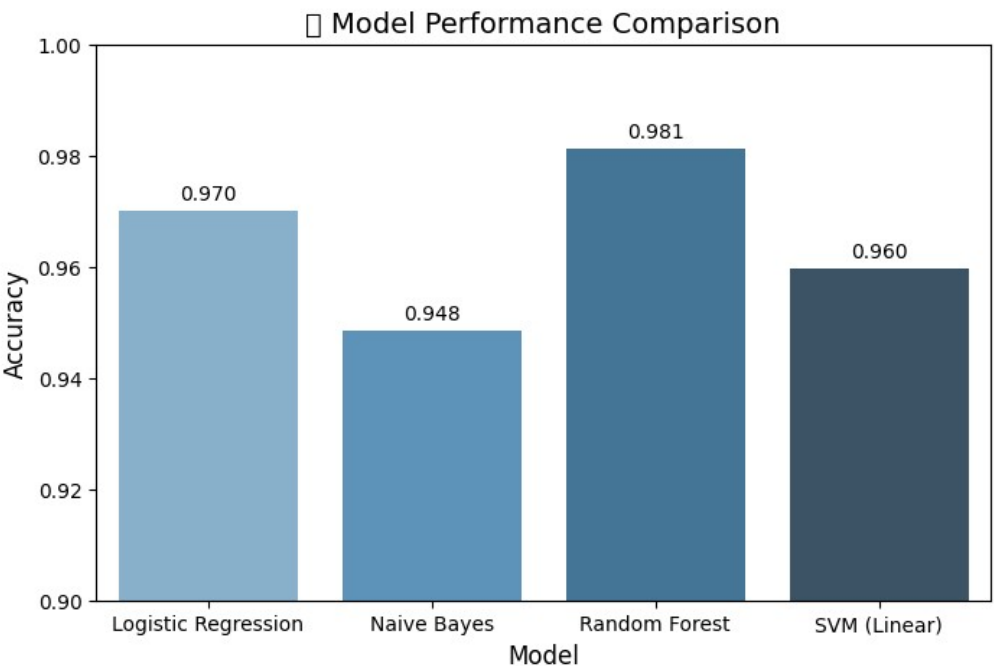
#### Random Forest (Best Model) Results:

- Overall Accuracy: 97.95%
- Precision (Real News): 97%
- Recall (Real News): 99%
- F1-Score (Real News): 98%
- Precision (Fake News): 99%
- Recall (Fake News): 97%
- F1-Score (Fake News): 98%

#### Model Comparison Summary:

1. Random Forest: 97.95% accuracy (Selected)
2. Logistic Regression: 97.00% accuracy
3. SVM (Linear): 95.96% accuracy
4. Naive Bayes: 94.85% accuracy

## Performance Comparison Implementation



## G. Conclusion and Recommendations

### Key Insights and Findings

#### 1. Model Performance:

The Random Forest algorithm achieved the highest accuracy of 97.95%, demonstrating that ensemble methods are particularly effective for fake news detection. The model shows excellent balance between precision and recall for both real and fake news classification.

#### 2. Feature Engineering Impact:

The combination of title and article text, along with comprehensive preprocessing (stemming, stopwords removal, and n-gram features), significantly improved model performance. The 5,000-feature vectorization with 1-3 gram ranges captured important linguistic patterns.

### 3. Algorithm Comparison:

While all algorithms performed well (above 94% accuracy), Random Forest's ensemble approach provided the most robust and reliable classification, making it the optimal choice for this application.

### 4. Evaluation Tools:

The development of comprehensive evaluation functions and visualization tools enhanced the analysis process, providing clear insights into model performance and enabling systematic comparison across different algorithms.

## Practical Applications

### 1. Web Application:

The developed Streamlit application provides an intuitive interface for real-time fake news detection, making the technology accessible to end users.

### 2. Educational Tool:

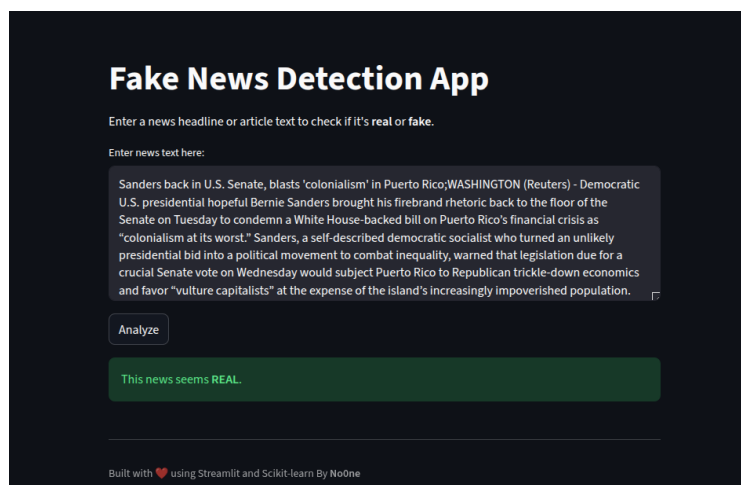
The system can serve as an educational resource to help users understand the characteristics of fake news and improve their media literacy.

### 3. Content Moderation:

The model can be integrated into social media platforms and news websites to automatically flag potentially fake content for human review.

### 4. Research Platform:

The comprehensive evaluation tools and modular design make this system suitable for further research and development in fake news detection.



## Limitations and Future Improvements

### 1. Current Limitations:

- **Language Dependency:** Model trained only on English text
- **Temporal Bias:** Performance may degrade with evolving language patterns
- **Context Understanding:** Limited ability to understand nuanced context and satire
- **Source Verification:** Cannot verify factual accuracy, only linguistic patterns

### 2. Recommended Future Enhancements:

- **Multilingual Support:** Expand to other languages for global applicability
- **Real-time Learning:** Implement online learning to adapt to new patterns
- **Fact-checking Integration:** Combine with external fact-checking databases
- **Ensemble with Other Models:** Integrate with transformer-based models like BERT