# Fake News Detection System
## Project Report

# Contents

# 1. Introduction

Fake news has become a serious issue in the digital era, influencing public opinion and spreading misinformation at a rapid scale. This project focuses on building a machine learning–based Fake News Detection System capable of classifying news articles as Real or Fake using Natural Language Processing (NLP) techniques.

# 2. Objectives

- To preprocess and clean textual news data using NLP techniques.

- To extract meaningful features from text using TF-IDF vectorization.

- To train a supervised machine learning model for fake news classification.

- To deploy the trained model using a Streamlit web application.
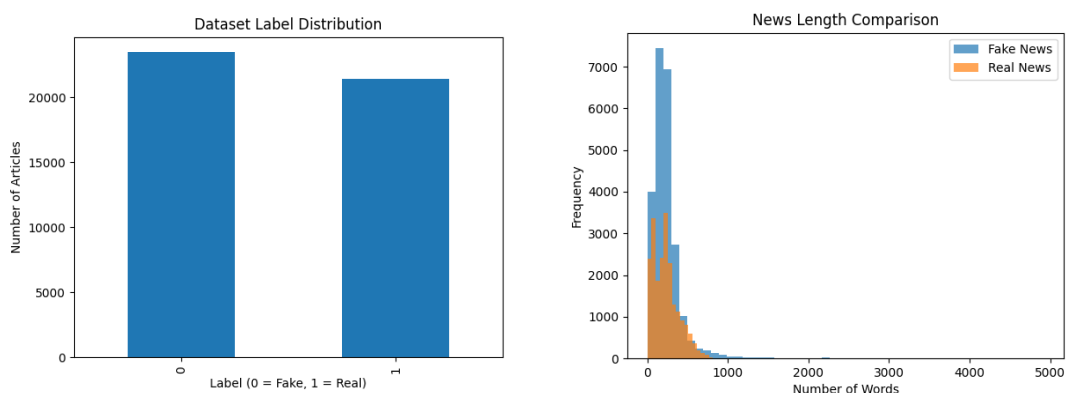
# 3.    System Architecture

The system architecture of the Fake News Detection System is designed in a modular and sequential manner. The process begins with the user providing raw news text either through a dataset or via the Streamlit web interface. The input text then passes through a preprocessing module where it is cleaned by converting to lowercase, removing punctuation, and eliminating stopwords. After preprocessing, the cleaned text is transformed into numerical feature vectors using the TF-IDF vectorizer. These vectors are then passed to the trained Logistic Regression model, which performs classification. Finally, the output is displayed to the user as a prediction along with probability scores for both real and fake classes.
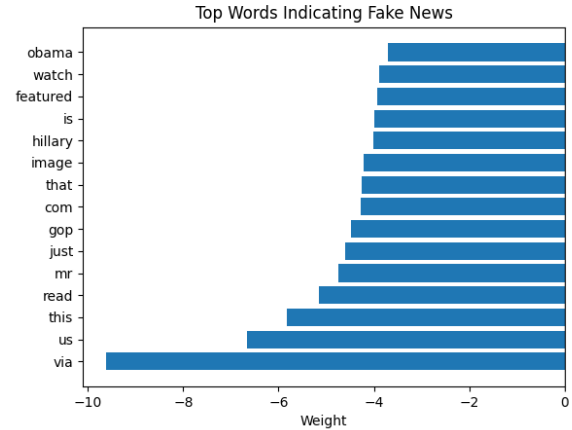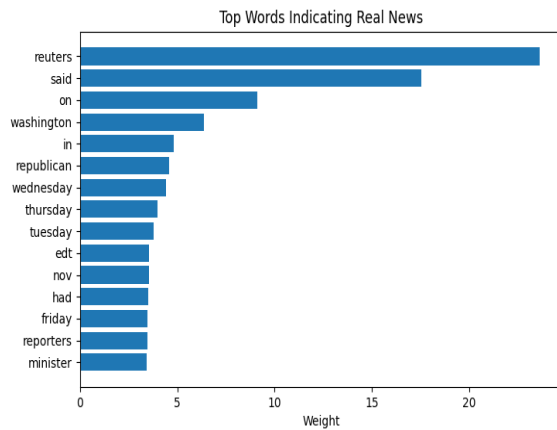
The architecture can be logically represented as:

Input News Text → Text Preprocessing → TF-IDF Feature Extraction → Logistic Regression Model → Prediction Output

# 4. Dataset Description

The dataset consists of two CSV files: Fake.csv and True.csv. Fake.csv contains news articles labeled as fake, while True.csv contains genuine news articles. Each record includes text-based news content. The fake news articles are labeled as 0, and real news articles are labeled as 1.

## 5. Methodology

### 5.1 Data Preprocessing

The text data is cleaned by converting it to lowercase, removing punctuation, and eliminating English stopwords. This step reduces noise and improves the quality of the extracted features.

### 5.2 Feature Extraction

TF-IDF (Term Frequency–Inverse Document Frequency) vectorization is used to convert textual data into numerical form. A maximum of 5000 features is selected to balance performance and computational efficiency.
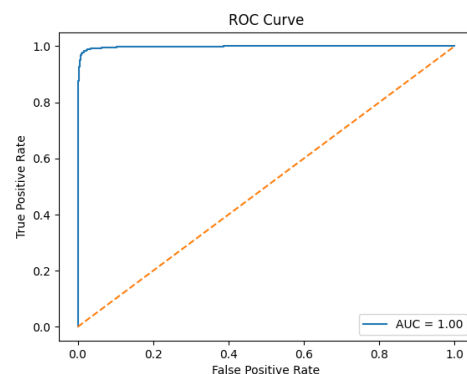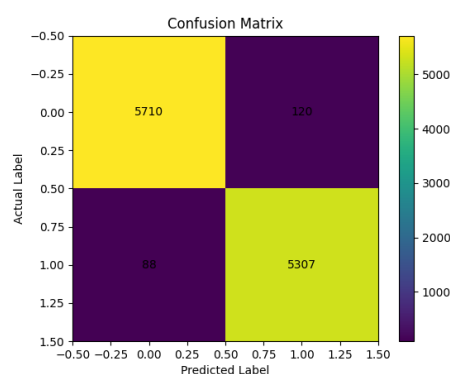
### 5.3 Model Training

Logistic Regression is used as the classification algorithm. The dataset is split into training (75%) and testing (25%) subsets. The model is trained with a maximum of 1000 iterations to ensure convergence.

### 5.4 Model Saving

The trained Logistic Regression model and the TF-IDF vectorizer are saved using pickle. This allows reuse of the trained components during deployment without retraining.

## 6. Model Evaluation

The model performance is evaluated using Accuracy, Confusion Matrix, and Classification Report. Accuracy measures the overall correctness of the model, while precision, recall, and F1-score provide deeper insight into classification performance for each class.

# 7. Application Deployment

A Streamlit-based web application is developed to provide a user-friendly interface. Users can input news text, and the system predicts whether the news is Fake or Real. The application also displays probability scores for both classes, increasing transparency.

# 8.    Pseudo-code

The following pseudo-code describes the overall working of the Fake News Detection System in a simplified manner:

Step 1: Load Fake and Real news datasets

Step 2: Assign labels (0 for Fake, 1 for Real)

Step 3: Merge and shuffle the dataset

Step 4: Preprocess text (lowercase, remove punctuation and stopwords)

Step 5: Convert text into TF-IDF feature vectors

Step 6: Split data into training and testing sets

Step 7: Train Logistic Regression model

Step 8: Evaluate model performance

Step 9: Save trained model and vectorizer

Step 10: Accept user input and predict news authenticity

# 9. Observations

- Text preprocessing significantly improves classification accuracy.

- TF-IDF effectively captures important keywords for fake news detection.

- Logistic Regression provides fast training and reliable results for text classification.

- Probability outputs help users understand prediction confidence.

# 10.    Limitations

- The model relies only on textual content and ignores images or metadata.

- Performance may degrade on news topics not well represented in the training data.

- Sarcasm and highly contextual language remain challenging to classify.

# 11.    Conclusion

This project presents a complete and practical Fake News Detection System using machine learning and NLP. The modular architecture, clear algorithmic steps, and visual representation make the system easy to understand, implement, and explain during examinations or project evaluations.

## 12. Reference And Important links

### 12.1 Streamlit

https://fake-news-detection-5fgselsibthklf4ssacqkg.streamlit.app/

### 12.2 Dataset

https://www.kaggle.com/datasets/clmentbisaillon/fake-and-real-news-dataset