

Project Report: Fake News Prediction (Veritas)

Date: 08 September 2025

1. Introduction

The rapid spread of digital information has led to a surge in "fake news," which can mislead public opinion and erode trust. This project tackles this issue by building a machine learning model to automatically classify news articles as real or fake, aiming to provide a tool that helps curb the spread of misinformation.

2. Abstract

This report outlines the creation of a fake news detector using Python and scikit-learn. The model was trained on a balanced dataset derived from Fake.csv, True.csv, and BBC News articles. The process involved comprehensive text preprocessing (stemming, stopword removal), feature extraction using TF-IDF, and classification with a Logistic Regression algorithm. The final model achieved a high accuracy of 98.37%, proving its effectiveness in distinguishing between real and fake news.

3. Tools Used

- **Data Handling:** Pandas, NumPy
 - **NLP:** NLTK, Re (Regular Expressions)
 - **Machine Learning:** Scikit-learn (TfidfVectorizer, LogisticRegression, train_test_split)
 - **Model Persistence:** Pickle
-

4. Project Workflow

Step 1: Data Preparation

- Three datasets (Fake.csv, True.csv, BBC News Train.csv) were loaded and merged.
- The dataset was cleaned by removing duplicates and null values.
- To prevent bias, the data was balanced by down-sampling the majority class ("real news") to match the count of "fake news" articles.

Step 2: Text Preprocessing

- The article text was standardized through normalization (lowercasing, removing punctuation) and by removing common English stopwords.
- Words were reduced to their root form using Porter Stemming to group related terms.

Step 3: Feature Extraction & Model Training

- The cleaned text was converted into numerical features using the TfidfVectorizer.
- The data was split into training (80%) and testing (20%) sets.

- A **Logistic Regression** model was trained on the training data.

Step 4: Evaluation and Saving the Model

- The model's performance was measured on the unseen test set, achieving **98.37% accuracy**.
 - The trained model and the TF-IDF vectorizer were saved as .pkl files for future use.
-

5. Conclusion

This project successfully produced a highly accurate fake news detection model. By combining a balanced dataset, effective text preprocessing, and a robust classification algorithm, the model can reliably identify fake news with high precision. This tool serves as a practical solution that can be integrated into digital platforms to help combat misinformation.