

جامعة أم القرى
UMM AL-QURA UNIVERSITY



Report on Fake News Detection

Student name	Student ID	Group
Sara Abed Alzulfie	443002749	2
Raghad Adel Alzulfie	444003159	1

Instructor name: Umaima Falath

Date: 25-10-2024



Contents

3.....	Introduction
4.....	Data Exploration
4.....	1. Data Loading and Preprocessing:
4.....	2. Text Preprocessing and Feature Extraction:
4.....	Tokenization:
4.....	Stop Word Removal:
4.....	Lemmatization:
5.....	Model Training and Evaluation:
5.....	Naive Bayes
6.....	Support Vector Machine (SVM):
7.....	Logistic Regression
8.....	Conclusion:
9.....	References

Introduction

Fake News and Real News: The dataset is designed for distinguish between fake and real news articles. Which means it aims to binary classification using Text-Based. The primary features used for classification are text-based, derived from the title, author, and text content of the news articles. Develop a machine learning program to identify when an article might be fake news. Run by the UTK Machine Learning Club. This project aims to detect fake news using machine learning techniques, specifically focusing on Natural Language Processing (NLP). It leverages libraries like Pandas, Scikit-learn, and NLTK for data manipulation, model training, and text preprocessing, respectively.

Dataset Characteristics:

- **train[1].csv:** This file likely contains the primary training data for the fake news detection model. It's expected to have the following columns:
 - **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:** A binary label indicating whether the article is fake (1) or real (0).
- **test[1].csv:** it contains the data used for evaluating the trained model. It should have the same columns as train except for the 'label' column, which the model needs to predict.
- **submit.csv:** This file probably contains the predictions made by the model on a separate dataset. It's used to submit results to a competition or for further analysis. It's assumed to have 'id' and 'label' columns, where 'id' matches the identifiers in the test dataset, and 'label' contains the predicted labels.

Data link; <https://www.kaggle.com/competitions/fake-news>



Data Exploration

1. Data Loading and Preprocessing:

The code begins by loading training, testing, and submission datasets. It handles missing data by filling them with spaces. Text features are combined for comprehensive analysis. Wordclouds are generated to visualize frequent words in real and fake news articles.

2. Text Preprocessing and Feature Extraction:

The code performs crucial preprocessing steps:

Tokenization: Dividing text into individual words.

Stop Word Removal: Eliminating common words (like "the", "a", "is") that don't carry much meaning.

Lemmatization: Reducing words to their base form (e.g., "running" to "run").

TF-IDF (Term Frequency-Inverse Document Frequency) vectorization converts text data into numerical format for model training, highlighting important words in the documents.

Model Training and Evaluation:

The code employs three machine learning models:

Naive Bayes: A probabilistic classifier based on Bayes' theorem.

The Naive Bayes model performs better on identifying fake news (higher recall for class 1), while its performance in classifying real news is slightly lower. The precision and recall are balanced, with a solid F1-score for both classes, meaning the model is fairly robust for this task.

The model correctly classified about **79%** of the news as either fake or real. This suggests the model performs reasonably well but may still misclassify about 21% of the data.

2. Confusion Matrix Interpretation

Class 0 (Real News): Precision = 0.80, Recall = 0.76, F1-Score = 0.78

Precision: Out of all news classified as real, 80% were correct.

Recall: The model correctly identified 76% of real news instances.

F1-Score: A harmonic mean of precision and recall at 0.78.

Class 1 (Fake News): Precision = 0.79, Recall = 0.83, F1-Score = 0.81

Precision: Out of all news classified as fake, 79% were correct.

Recall: The model identified 83% of fake news correctly.

F1-Score: Indicates strong performance in classifying fake news (0.81).

3. Macro and Weighted Averages

Macro Average: Indicates equal weighting for each class, averaging around 79%.

Weighted Average: Accounts for class imbalance and provides an overall performance score of 79%.



Support Vector Machine (SVM):

Support Vector Machines (SVMs) are powerful supervised learning algorithms used for classification and regression tasks. SVMs aim to find a hyperplane in an N-dimensional space that best separates the classes. The goal is to maximize the margin between the closest points (support vectors) of different classes, ensuring optimal classification. SVMs can handle both linear and non-linear classification using kernels.

SVM Model Results for Fake News Detection

Accuracy: 86.99%

The SVM model correctly classifies about 87% of the news as either fake or real, showcasing strong overall performance.

Confusion Matrix Breakdown:

Class 0 (Real News): **Precision** = 0.85, Recall = 0.86, F1-Score = 0.86

Precision: 85% of news classified as real are correctly labeled.

Recall: The model captures 86% of the actual real news instances.

F1-Score: Combines precision and recall, showing balanced performance.

Class 1 (Fake News): Precision = 0.88, Recall = 0.88, F1-Score = 0.88

Precision: 88% of news classified as fake are correctly identified.

Recall: The model catches 88% of actual fake news cases.

Overall Metrics:

Macro Average: Averages precision, recall, and F1 across both classes, each scoring around 87%.

Weighted Average: Incorporates the class distribution and provides an overall accuracy of 87%.

Summary:

The SVM model performs very well in detecting both real and fake news, with slightly better performance in identifying fake news (higher precision and recall for class 1). The model's high accuracy and balanced F1-scores for both classes suggest it is robust for this classification task, outperforming Naive Bayes in this context.



Logistic Regression: Logistic regression is a supervised machine learning algorithm used for classification tasks where the goal is to predict the probability that an instance belongs to a given class or not. Logistic regression is a statistical algorithm which analyze the relationship between two data factors.

Precision, Recall, and F1-Score:

These metrics provide a more detailed insight into your model's performance for each class:

Precision: Measures the accuracy of positive predictions.

Precision for class 0: 83%

Precision for class 1: 90%

Recall: Measures the ability of the model to identify all actual positive instances.

Recall for class 0: 88%

Recall for class 1: 86%

F1-score: A harmonic mean of precision and recall, providing a balanced measure.

F1-score for class 0: 0.85

F1-score for class 1: 0.88

Support:

This refers to the number of actual instances in each class in the test set.

Support for class 0: 672

Support for class 1: 888

Macro Average and Weighted Average:

These are overall averages of precision, recall, and F1-score:

Macro Average: Calculates the average without considering class imbalance.

Weighted Average: Calculates the average, giving more weight to classes with more instances.

In essence, your model is performing quite well, with high accuracy, precision, recall, and F1-scores for both classes. However, there's always room for improvement. You might explore techniques like hyperparameter tuning or different algorithms to further enhance its performance.



Conclusion:

The provided code demonstrates a comprehensive approach to fake news detection by combining text preprocessing techniques with machine learning models. It focuses on preparing data for analysis, training classifiers, and assessing their performance using appropriate evaluation metrics. This analysis helps determine the effectiveness of each model in identifying fake news articles.



References

- William Lifferth. Fake News.
<https://kaggle.com/competitions/fake-news>, 2018. Kaggle.
- Logistic Regression in Machine Learning.GeeksforGeeks Improve & GeeksforGeeks
<https://www.geeksforgeeks.org/understanding-logistic-regression/>
- ADA Fake news.swetha2096
<https://www.kaggle.com/code/swetha2096/ada-fake-news/>









جامعة أم القرى
UMM AL-QURA UNIVERSITY



