

**Individual Project Report**  
CPSC 3750:Artificial Intelligence  
OLUWASEENI AJAYI  
JASPREET KAUR, PHD

**Project Title**

Combating Disinformation: A Machine Learning and Natural Language Processing Approach to Detecting Fake News

**Introduction**

In this project, I tackle the increasingly prevalent issue of fake news, leveraging machine learning and natural language processing techniques to differentiate between genuine and fabricated news articles.

**Background Information/Literature Survey**

The spread of fake news has become a global concern, influencing public opinion and potentially undermining the democratic process. Various research efforts have focused on applying machine learning techniques for fake news detection, often using models like Naive Bayes, Support Vector Machines, and Logistic Regression due to their effectiveness in text classification tasks.

**Methodology**

I used a Kaggle dataset of over 40,000 labeled news articles, each marked as 'real' or 'fake.' The data was preprocessed, which included cleaning, tokenization, and vectorization. I trained three machine learning models (Logistic Regression, Multinomial Naive Bayes, and Support Vector Machine) and evaluated their performance.

**Working**

My approach involved three major steps:

1. I preprocessed the data by cleaning, tokenizing, and vectorizing the news articles.
2. I trained Logistic Regression, Multinomial Naive Bayes, and Support Vector Machine models on the preprocessed data.
3. I evaluated the models using various metrics such as accuracy, precision, recall, and F1-score and visualized the performance via a confusion matrix.

**Future Possibilities**

One could consider employing deep learning techniques, such as Recurrent Neural Networks (RNNs), Convolutional Neural Networks (CNNs), and Transformer-based models like BERT, which have proven highly effective in more complex text classification tasks. Additionally, the exploration of different feature extraction methods and the inclusion of additional data sources could also enhance the model's performance.

**Challenges**

Some of the challenges faced during this project included understanding the complexities of natural language processing and text vectorization, selecting suitable machine learning models for the task, and interpreting the evaluation metrics in the context of the task. Further, ensuring the models were not overfitted or biased was also a key challenge.

**Timeline**

The project took approximately three weeks to complete, with the first week dedicated to understanding the problem and data preprocessing, the second week to model selection and training, and the final week to model evaluation and report writing.

### **Results**

The Logistic Regression model performed best, with an accuracy of 99.75%, followed by MultinomialNB with 95.7% accuracy. The SVM model's performance could not be evaluated due to computational limitations.

### **Conclusions**

The project demonstrates the potential of using machine learning and natural language processing to detect fake news. While the chosen models yielded promising results, future work employing more advanced techniques could yield even more accurate and robust fake news detection systems.