

For this project, my goal is to develop a machine learning and natural language processing (NLP) solution that can detect and monitor coordinated misinformation campaigns. I believe misinformation poses significant risks during critical events like elections or public health crises. It has the potential to undermine public trust and disrupt decision-making processes. By analyzing patterns in text data, such as fake news articles or social media posts, I aim to identify fabricated content and uncover networks involved in spreading coordinated disinformation.

Misinformation is a growing issue worldwide, fueled by the rapid growth of social media and online communication platforms. I've seen how misinformation campaigns during critical times, such as the COVID-19 pandemic or elections, have far-reaching consequences. These include influencing public opinion, creating confusion, and eroding trust in legitimate information sources. My project will address this issue by building a robust system to detect fake content and monitor coordinated campaigns, giving stakeholders the tools they need to counter disinformation effectively. To achieve this, I will use the WELFake dataset, a labeled collection of real and fake news articles. This dataset offers the comprehensive text data I need to train machine learning models and detect patterns in misinformation campaigns.

I plan to measure the success of this project through several key objectives. First, I will create a classifier capable of accurately identifying fake news articles. Additionally, I will use clustering and other unsupervised techniques to identify networks of coordinated disinformation. Finally, I will ensure the system provides interpretable results, highlighting the features that distinguish fake news from genuine articles and patterns that suggest coordinated efforts. To evaluate my model, I will use metrics like accuracy, precision, recall, and F1 score for the supervised classifier. For network detection, I will rely on clustering coherence and silhouette scores to measure performance.

My approach involves several key steps. I will start by preprocessing the data, which includes cleaning and tokenizing the WELFake dataset, handling missing values, and creating features using text preprocessing techniques such as TF-IDF and word embeddings. Next, I'll conduct exploratory data analysis (EDA) to explore patterns in the dataset, visualize differences between fake and real articles, and identify initial evidence of coordination. I will also engineer advanced NLP features, such as sentiment analysis and linguistic pattern metrics, to enrich the dataset. During the modeling phase, I'll implement supervised models like logistic regression, random forests, and transformer-based models like BERT for fake news classification. I'll also apply clustering algorithms, such as k-means and DBSCAN, to detect coordinated misinformation networks. Finally, I'll validate the model through cross-validation and refine its performance with hyperparameter tuning.

There are challenges I'll need to address throughout the project. Data quality is one of the primary concerns, as incomplete or biased data could impact the accuracy of my models. To overcome this, I'll use preprocessing and balancing techniques. Another challenge is computational resources since training complex models like transformers requires significant processing power. Ethical considerations will also be at the forefront of my work, as I want to ensure the data is used responsibly and that my model does not introduce unintended biases.

The stakeholders I'm focusing on include social media platforms that can use my system to mitigate the spread of misinformation, public sector organizations that need to protect democratic processes and public health systems, and the general public, who rely on accurate and trustworthy information.

The deliverables for this project include a machine learning pipeline that can detect misinformation and coordinated campaigns, a detailed report outlining my methods and findings, and an optional interactive dashboard to visualize results. By tackling this issue, I aim to demonstrate how machine learning and NLP can be used to address one of the most pressing challenges of the digital age while building a meaningful solution that can restore trust and accountability.