

FAKE NEWS DETECTION

RAUNAK SHAHI - 102203054 DEV MEHTA-102203073



PROBLEM STATEMENT



2013 Muzaffarnagar riots (started due to misinformation online)

In today's digital age, the rapid spread of misinformation and fake news poses significant threats to public trust, social harmony, and informed decision-making. With the rise of social media and online platforms, distinguishing credible news from fabricated content has become increasingly challenging.



2018 Myanmar Rohingya Crisis (lot of fake news was spread online on platforms like faceook)

OUR SOLUTION

To combat the spread of fake news, we developed a machine learning-based detection system using Logistic Regression. By transforming textual data into numerical features (e.g., TF-IDF), logistic regression classifies news articles as real or fake based on patterns in the data. This solution is lightweight, interpretable, and computationally efficient for datasets of moderate size.

Transformer-based models like DistilBERT, which excel at understanding the contextual and semantic nuances of text, can also be used offering significantly better accuracy and generalization. Memory limitations during training constrained our ability to implement this solution fully.





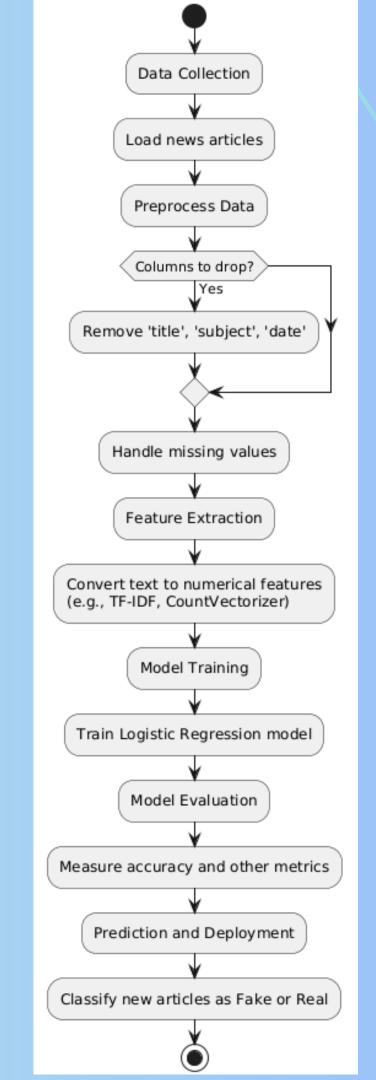
Logistic regression is a statistical method used in machine learning for binary classification problems, where the goal is to predict one of two possible outcomes (e.g., yes/no, true/false, 0/1). It predicts the probability that a given input belongs to a specific category.

TRANSFORMERS

Transformers are a type of neural network architecture designed to handle sequential data, such as text, speech, or time-series, while addressing the limitations of traditional sequence models like RNNs and LSTMs.

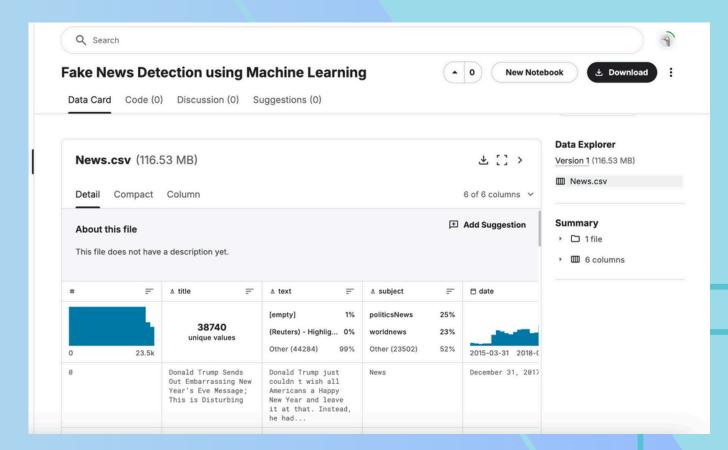


DATA FLOW DIAGRAM

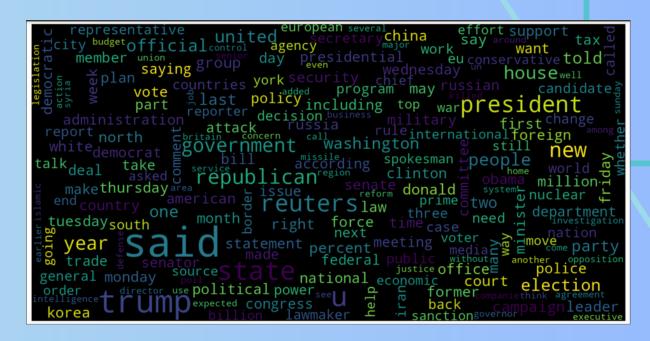


DATASET

- The dataset used for this project is a comprehensive news dataset consisting of 44,919 entries.
- Title: Provides the headline of the news article.
- Text: Represents the main body of the article, crucial for analyzing the substance and identifying patterns.
- Subject: Categorizing the news into topics which can influence the style and nature of the writing.
- Class: Serves as the label for fake news detection, where
 O indicates real news and 1 indicates fake news.



The dataset used, on kaggle



Word Cloud of dataset used





• Preprocessing Steps:

- The textual data is processed using tokenization, each article split into individual tokens (words or phrases)
- Stopwords were removed to focus on meaningful content.
- The cleaned data was then transformed into numerical features using TF-IDF (Term Frequency-Inverse Document Frequency) capturing the significance of words in the dataset.

Model Used:

0000

0000

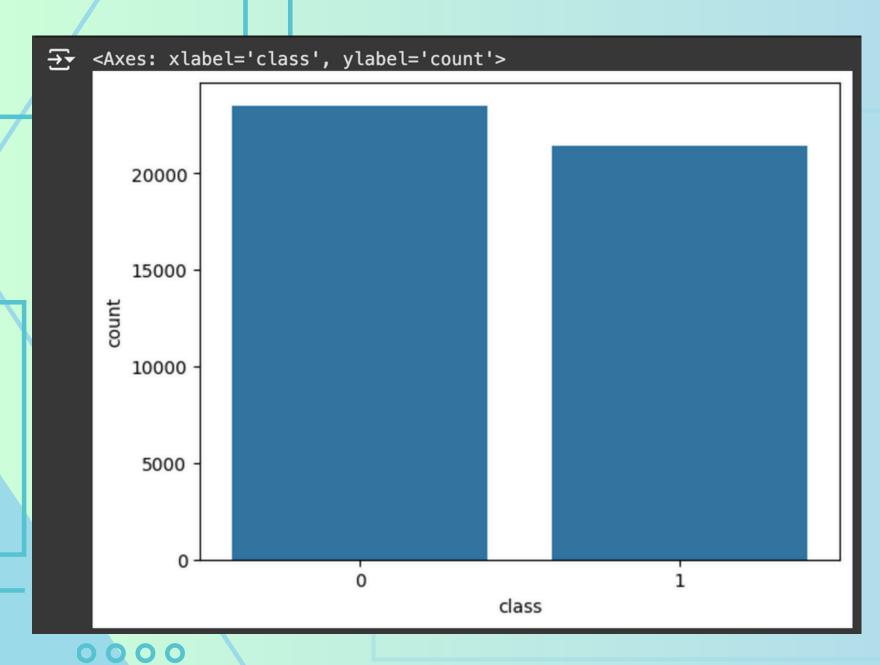
0000

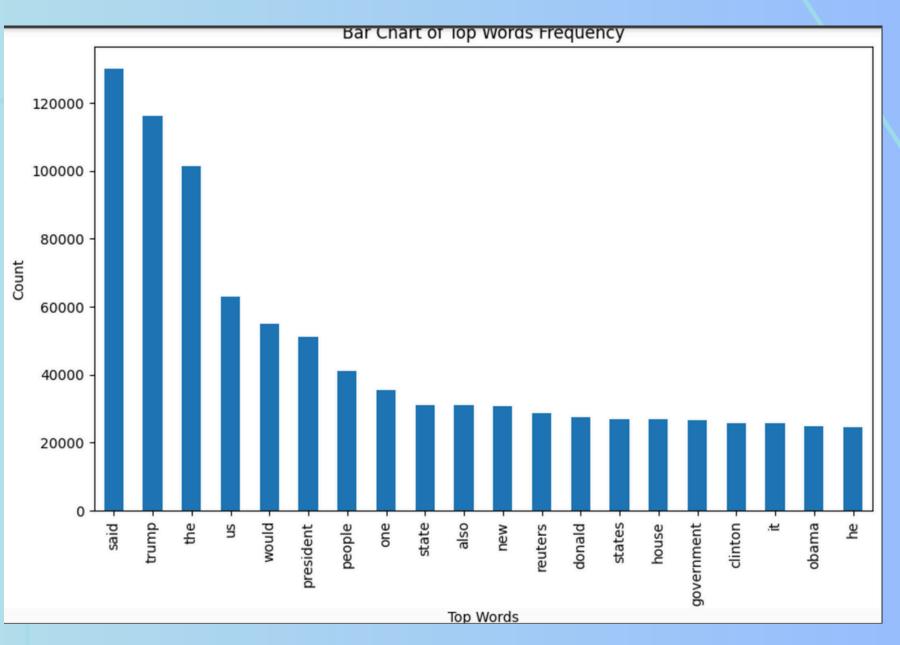
0000

- Logistic Regression: A lightweight, binary classification algorithm suitable for moderately sized datasets which maps input features to probabilities, enabling effective classification of news articles as real (0) or fake (1).
- Alternate Solution: Transformer-based models, such as DistilBERT, which excel at understanding contextual and semantic nuances in text, can also be used.

MODEL GRAPHS







X Axis: 'class', Y Axis: 'count'

0000

0000

0000

0000

Bar Chart of Top words frequency
X Axis: words
Y Axis: Word Count

RESULTS



0000

0000

0000

0000

0000

0000

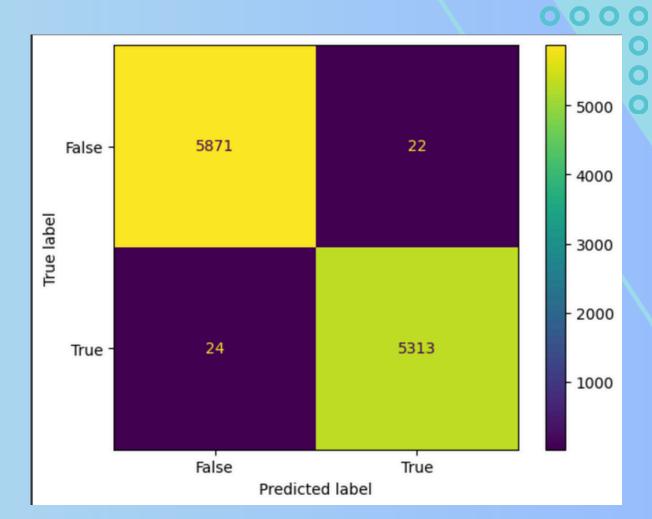
0.9999703167205913
0.9967052537845058

Accuracy

- 1. Training dataset
 - 2. Test dataset

_ Classifica		oort: cision	recall	f1-score	support	
	0 1	1.00 1.00	1.00 1.00	1.00 1.00	5809 5421	
accura macro a weighted a	vg	1.00 1.00	1.00 1.00	1.00 1.00 1.00	11230 11230 11230	

Classification Report



0000

0000

Confustion Matrix

			[4976/6738
Epoch	Training Loss	Validation Loss	Accuracy
1	0.000000	0.003259	0.999443
2	0.000000	0.001305	0.999889

Accuracy with Transformer



