

16-08-2025

Fake News Detection using NLP

Project Goals

- Misinformation spreads faster than truth online
- Impacts politics, health, economy, and society
- AI can help flag suspicious news quickly
- Goal: Classify news articles as Fake (0) or Real (1)

NLP

- NLP = AI's way of understanding and processing human language
- Used in:
 - Chatbots
 - Spam filters
 - Language translation
- In this project: process news text so ML can understand it

Dataset Overview

- Files: data.csv (training), validation_data.csv (for predictions)
- Columns:
 - 1.label – 0 (fake), 1 (real), 2 (unknown in validation)
 - 2.title – news headline
 - 3.text – full article
 - 4.subject – topic
 - 5.date – publication date

Data Loading & Preprocessing

Cleaning the Data

- Steps:
 1. Combine title + text
 2. Convert to lowercase
 3. Remove special characters, numbers, and extra spaces
 4. Remove stopwords (common words like "the", "and")
 5. Lemmatize (convert words to base form: "running" → "run")

Turning Text into Numbers

Why? Machines can't understand text directly

- TF-IDF:
 1. Counts word frequency (Term Frequency)
 2. Reduces weight of very common words (Inverse Document Frequency)
 3. Output: Large sparse matrix of numbers representing each article

Model Training

- Tried Linear SVM and Logistic Regression
- SVM performed better:
- Train Accuracy: 99.95%
- Test Accuracy: 99.37%
- Why SVM? Works well with high-dimensional sparse data like text

Predictions

- Used trained SVM + fitted TF-IDF vectorizer
- Applied same preprocessing to validation_data.csv
- Replaced label 2 with predicted 0 or 1
- Saved results as validation_data_predicted.csv

Results & Confusion Matrix

- Accuracy: 99.37%
- High F1-score = balanced precision & recall
- Confusion matrix: shows very few misclassifications
- Strong generalization to unseen data

Flowchart

