

Fake News Detection

Submitted for

CSET211 - Statistical Machine Learning

Submitted by:

E23CSEU1226	AKSHAT RATHI
E23CSEU1238	YOGITA TOMAR
E23CSEU1225	GOVIND SINGH

Submitted to:

PRASHANT KAPIL

July-Dec 2024

SCHOOL OF COMPUTER SCIENCE AND ENGINEERING



BENNETT
UNIVERSITY
THE TIMES GROUP

INDEX

Sr.No	Content	Page No
1	INTRODUCTION	3
2	SURVEY AND DATASET	4
3	DATA PREPROCESSING	5 ,6,7
4	DATA VISUALISATION AND ANALYSIS	8
5	TRAINING AND TESTING	9
6	CONCLUSION AND FUTURE WORK	10
7	REQUIREMENTS	11

ABSTRACT

- . The spread of misinformation and fake news through social media platforms has become a significant societal problem, impacting political discourse, public opinion, and public health.

The project aims to build an efficient fake news detection system using machine learning algorithms, natural language processing (NLP), and deep learning models. The goal is to identify fake news articles and minimize the influence of misinformation.

The paper explores existing fake news detection systems, discusses current methodologies, evaluates performance, and proposes potential future improvements.

2. Literature Survey on Fake News Detection Systems

A Systematic Literature Review

- **Machine Learning Algorithms:**
 - Support Vector Machines (SVM):
Achieved an accuracy of 93.61%. ◦ Random Forest: Popular for classification tasks.
 - Naive Bayes: Efficient for text classification.
 - Deep Learning Models: Applied alongside traditional methods for improved accuracy. ◻
 - Datasets:
 - Datasets used for training include features such as article text, metadata, and user engagement.
 - Example dataset: 7,796 news articles (50% real, 50% fake).
 - **NLP:** Techniques like TF-IDF, sentiment analysis, and feature extraction play a crucial role in improving detection accuracy.
-

SURVEYS:

- 1) Fake news and the spread of misinformation: A research roundup by [Denise-Marie Ordway](#) . This research roundup explores the impact of fake news and misinformation, examining factors such as cognitive biases, the role of social media, and the effectiveness of debunking. Studies highlight how misinformation spreads, why people believe it, and its influence on politics, trust, and attitudes. Scholars also analyze methods to identify and counter false narratives, including detailed corrections and leveraging source credibility. The findings emphasize the need for robust tools to address misinformation effectively in the digital age. For more details, visit the source <https://journalistsresource.org/politics-and-government/fake-news-conspiracy-theories-journalism-research/>

- 2) Misinformation, Disinformation, and Propaganda: Fake News [[David M. J. Lazer, et al., "The Science of Fake News," Science 09 Mar 2018: Vol. 359, Issue 6380, pp. 1094-1096.](#)]. "Fake news" is "fabricated information that mimics news media content in form but not in organizational process or intent. Fake-news outlets, in turn, lack the news media's editorial norms and processes for ensuring the accuracy and credibility of information. Fake news overlaps with other information disorders, such as misinformation (false or misleading information) and disinformation (false information that is purposely spread to deceive people)." https://guides.library.cornell.edu/evaluate_news/fakenews

- 3) The science of fake news [DAVID M. J. LAZER, MATTHEW A. BAUM, YOCHAI BENKLER, ADAM J. BERINSKY, KELLY M. GREENHILL, FILIPPO MENCZER, MIRIAM J. METZGER, BRENDAN NYHAN, GORDON PENNYCOOK](#) The rise of fake news reveals the weakening of traditional defenses against misinformation, particularly online. Global concerns focus on societal vulnerabilities and the urgent need for new protective systems. Research explores its spread, belief mechanisms, and the challenges of combating politically-driven misinformation effectively. <https://www.science.org/doi/10.1126/science.aao2998>

4)

Beyond Fake News

Words matter. Misinformation, fake news, disinformation, conspiracy theories...so many different ways of describing how we can be deceived, either deliberately or inadvertently with misleading information. The BBC encourages people to think critically about what they read, see and view so that they can spot misleading or bad information, and resist sharing content that might be false or out of context.

<https://www.bbc.com/beyondfakenews/fakenewsdefinitions>

5) Fake News [Bente Kalsnes](#)

Fake news, though not new, gained international attention after the 2016 U.S. presidential election. Modern technologies have revolutionized its creation, spread, and consumption, making misinformation harder to identify. Research focuses on its characterization, production motives, methods of dissemination, and counter-strategies, which include technical solutions, improved media literacy, and fact-checking efforts.

[Fake News | Oxford Research Encyclopedia of Communication](#)

6) The Real Impact of Fake News: The Rise of Political Misinformation—and How We Can Combat Its Influence

Misinformation is unintentional falsehood, while disinformation is deliberately deceptive. A Columbia University panel highlighted social media, AI, and underfunded local journalism as key factors amplifying these issues. Solutions proposed include media literacy programs, regulation, and support for local news, with AI offering both risks and opportunities for combating misinformation.

[The Real Impact of Fake News: The Rise of Political Misinformation—and How We Can Combat Its Influence | Columbia University School of Professional Studies](#)

7) Fake news, social media, and "The Death of Truth" *by Dustin Stephens. Editor: Ed Givnish.*

The rise of alternate facts, social media echo chambers, and AI-driven misinformation undermines shared realities and trust in democracy. Section 230 shields online platforms from liability, enabling unchecked proliferation of fake news and disinformation. Efforts like NewsGuard aim to address this, but increasing skepticism, fake local news, and deepfakes raise concerns about fair elections and societal stability.

[Fake news, social media, and "The Death of Truth" - CBS News](#)

8) The Problem of Fake News in India: Issues, Concerns and Regulation

Fake news, amplified by social media and bots, spreads misinformation, causing unrest and political manipulation. In India, it has led to violence, while countries like Russia and China use it for political gain. Tackling the issue requires education, media literacy, and effective tech regulation to combat misinformation without restricting free speech.

[Problem of Fake News in India: Issues, Concerns and Regulation](#)

9)

Finland is winning the war on fake news. What it's learned may be crucial to Western democracy By [Eliza Mackintosh](#), CNNVideo by Edward Kiernan, CNN

Finland has been fighting disinformation since 2014 by teaching media literacy and critical thinking in schools and through government initiatives. The country has trained citizens to spot fake news, especially from Russian sources. Finland's approach has inspired other nations, but officials warn that the battle against misinformation is ongoing and will only grow more complex.

[Finland is winning the war on fake news. Other nations want the blueprint](#)

Data Preprocessing for Fake News Detection □ 1. Data

Collection:

Collect data from reliable sources such as news websites, social media, and datasets like LIAR, FakeNewsNet.

Ensure data is in structured formats like CSV, JSON.

Data Cleaning:

Remove duplicate articles, handle missing values. ◦ Normalize the text (convert to lowercase) for consistency.

Text Preprocessing:

Tokenization: Break text into words or tokens.

Removing Stop Words: Filter out common nonmeaningful words (e.g., "and," "the").

Stemming/Lemmatization: Reduce words to their base form.

Remove Punctuation/Special Characters: Clean the text for further analysis.

Feature Extraction:

Bag of Words (BoW): Represent text as word frequency counts.

TF-IDF: Weigh words based on their importance.

Word Embeddings: Use pre-trained embeddings like Word2Vec, GloVe, or BERT.

Data Splitting:

Split data into training and testing sets (e.g., 80/20 split).

Stratified Sampling: Ensure balanced representation of fake and true news articles.

3. Data Visualization and Exploratory Data Analysis (EDA)

Libraries:

`pandas`: Data manipulation and analysis. ◦

`matplotlib`, `seaborn`: Visualization of data distributions.

`nltk`, `wordcloud`: Text analysis and visualization. ◦

`gensim`, `plotly.express`: For advanced text analysis and interactive plotting. ◦ `sklearn`: Used for model building, including feature extraction and splitting datasets.

Techniques:

Visualizing the distribution of fake vs true news articles.

Analyzing the most frequent words or phrases in the dataset. ◦ Visualizing relationships between features like article length, sentiment, and authenticity.

4. Model Building and Creation

Train-Test Split: Divide the data into training and testing sets to evaluate the model's performance.

Algorithms:

Decision Tree: Simple and interpretable model for classification.

Logistic Regression: Models the probability of a binary outcome (fake/true).

Gradient Boosting Classifier: Builds an ensemble of models sequentially, often providing high accuracy.

Random Forest: Combines multiple decision trees to enhance prediction accuracy.

5. Performance Evaluation

Metrics:

Precision, Recall, F1 Score: Metrics used to evaluate model performance, focusing on false positives/negatives and balancing precision and recall.

Accuracy: Overall correctness of the model.

Confusion Matrix: Visual representation of true positives, false positives, true negatives, and false negatives.

6. Conclusions and Future Work

Summary of Achievements:

Built a fake news detection system using machine learning algorithms.

Achieved promising results with traditional and deep learning models. ◦ Successfully preprocessed and analyzed the dataset using NLP techniques.

Challenges:

Imbalanced datasets: Fake news is often less frequent than true news.

Ethical concerns: Bias in the models can arise due to biased training data. ◦ Real-time detection: Ensuring the model performs well with new, unseen data.

Future Work:

Explore more advanced models such as transformers (e.g., BERT, GPT) for higher accuracy. ◦ Focus on reducing model bias and ensuring fairness in predictions.

Implement multi-modal detection, combining text, image, and video analysis for a more robust system.

7. Hardware and Software Requirements □

Hardware: ◦ Modern CPU or GPU (for deep learning models).

Sufficient RAM (minimum 8GB recommended).

Software:

Python: The primary programming language used.

Libraries: `pandas`, `sklearn`, `tensorflow`, `pyTorch`, `nltk`, `spaCy`, `matplotlib`, etc.

IDE: Jupyter Notebook or VSCode for model development.