



PROJECT REPORT

YENEPLOYA DEEMED TO BE UNIVERSITY, BANGLORE

For the degree of Bachelor of computer application 2023-2026

By

Mohammed Ahsan V A

Reg no: 23BBCDAI101

Industry Project Title	Fake News Detection System Using Social Media Data
Name of the Company	Tata Consultancy Services
Name of the Institute	Yenepoya Deemed to be University

Start Date	End Date	Total Effort (hrs.)	Project Environment	Tools used
13.11.2025	11.02.2026	140 hrs	Jupyter Notebooks, VS Code (Local), GitHub	Python, Pandas, NumPy, Scikit-learn, TensorFlow, NLTK, Streamlit, GitHub

ACKNOWLEDGEMENTS

I would like to express my sincere gratitude to **Tata Consultancy Services (TCS iON)** for providing me with the opportunity to undertake this industry project on **Fake News Detection using Machine Learning**. This project offered valuable exposure to real-world applications of artificial intelligence and natural language processing, enabling me to bridge the gap between theoretical concepts and practical implementation.

I would like to thank my academic mentors and faculty members for their continuous guidance, encouragement, and constructive feedback throughout the duration of this project. Their support helped me understand key concepts such as data preprocessing, feature extraction, machine learning model development, and performance evaluation, and guided me in structuring the project in a systematic and professional manner.

I am also grateful to the developers and contributors of open-source tools and libraries such as **Python, Pandas, Scikit-learn, TensorFlow, and Streamlit**, which played a crucial role in the successful development of this system. The availability of high-quality documentation and learning resources greatly supported my learning process and technical implementation.

Finally, I would like to acknowledge my peers and well-wishers for their motivation and support during the course of this project. This internship has been a valuable learning experience and has strengthened my interest in artificial intelligence, data science, and real-world problem solving using machine learning.

OBJECTIVE AND SCOPE

Objective

The primary objective of this project is to design and develop an automated **Fake News Detection System** using machine learning and deep learning techniques. The project aims to analyze textual news content and classify it as **Real** or **Fake** based on learned linguistic and contextual patterns. By leveraging natural language processing and supervised learning models, the system seeks to provide an efficient and scalable solution for identifying misinformation in digital content.

Another key objective of the project is to compare the performance of traditional machine learning algorithms with deep learning models for fake news classification. The project also focuses on developing a user-friendly web-based interface that allows users to input news content and receive real-time prediction results along with confidence scores. This helps improve accessibility and practical usability of the system.

Through this project, the objective is not only to achieve high classification accuracy but also to gain hands-on experience in data preprocessing, feature extraction, model training, evaluation, and deployment in a real-world AI application.

Scope

The scope of this project is limited to the detection of fake news based on **textual data**. The system analyzes English-language news headlines, articles, and social media content and performs **binary classification** to determine whether the given input is real or fake. The project focuses on supervised learning techniques using labeled datasets and does not involve real-time fact verification from external sources.

Within this scope, the project includes data collection from publicly available datasets, preprocessing of textual data, feature engineering, model training, and performance evaluation. Both traditional machine learning models and deep learning models are implemented and compared to assess effectiveness. The project also covers the deployment of the trained model using a web application framework to allow interactive user testing.

The scope of this project does not include multimedia-based fake news detection such as image, video, or audio analysis. Additionally, aspects such as author credibility analysis, social network analysis, or real-time web scraping are beyond the scope of the current implementation. However, the system is designed in a modular manner, allowing future enhancements and extensions.

PROBLEM STATEMENT

In the modern digital era, online platforms and social media have become primary sources of information for millions of users worldwide. While this has enabled rapid information sharing and improved access to news, it has also led to the widespread circulation of **fake news and misinformation**. False or misleading information can spread quickly across digital platforms, often reaching a large audience before it can be verified or corrected.

Fake news poses serious challenges in various domains, including public health, politics, and social stability. Misinformation related to health issues, government policies, or financial markets can create panic, influence public opinion, and lead to harmful real-world consequences. Due to the massive volume of content generated daily, manual verification of news articles and social media posts is not feasible, making traditional fact-checking methods insufficient.

Existing approaches to identifying fake news often rely on manual analysis or simple rule-based techniques, which lack scalability and fail to capture complex linguistic patterns. As fake news content becomes more sophisticated and well-written, these methods struggle to differentiate between genuine and misleading information.

To address these challenges, there is a need for an **automated, data-driven solution** that can analyze textual content and accurately classify news as real or fake. By leveraging machine learning and natural language processing techniques, such a system can assist in early detection of misinformation and act as a supportive tool for users, organizations, and content moderation platforms.

EXISTING APPROACHES

Fake news detection has traditionally been addressed using manual and rule-based approaches. One common method involves **manual fact-checking** performed by organizations and experts who verify the authenticity of news articles before publication. While this approach is highly accurate, it is time-consuming and cannot scale to handle the vast amount of content generated daily on digital platforms.

Another widely used approach is **rule-based filtering**, where news articles are flagged based on predefined keywords, source domains, or writing patterns. Although this method is simple to implement, it lacks flexibility and often fails when fake news content avoids obvious keywords or originates from newly created sources. Such approaches require frequent updates and manual intervention to remain effective.

Traditional machine learning techniques such as **Naive Bayes, Support Vector Machines (SVM), and Random Forest classifiers** have also been applied to fake news detection. These models typically rely on handcrafted features such as word frequency, n-grams, and TF-IDF representations. While these methods offer better scalability than manual techniques, they struggle to understand contextual meaning, sarcasm, and sentence structure, which are critical in distinguishing sophisticated fake news from real content.

Due to the limitations of manual verification, rule-based systems, and traditional machine learning approaches, there is a growing need for advanced techniques that can automatically learn contextual and semantic patterns from text. Deep learning models, particularly sequence-based models such as LSTM networks, provide a more effective solution by capturing long-term dependencies in textual data and improving classification accuracy.

APPROACH / METHODOLOGY – TOOLS AND TECHNOLOGIES USED

The methodology adopted for this project follows a structured and systematic approach to ensure accurate detection of fake news using machine learning and deep learning techniques. The overall methodology consists of multiple stages, beginning with data collection and preprocessing and ending with model evaluation and deployment through a web-based interface.

The first stage of the methodology involves **data collection and preprocessing**. Publicly available labeled datasets containing real and fake news articles were collected and combined to form a unified dataset. The raw textual data contained noise such as special characters, URLs, punctuation, and inconsistent formatting. These issues were addressed through text cleaning techniques including lowercasing, removal of stop words, tokenization, and lemmatization to improve data quality.

Once the data was cleaned, **feature extraction and representation** were performed. Traditional machine learning models utilized techniques such as TF-IDF vectorization to convert text into numerical form. For deep learning models, tokenization and sequence padding were applied to prepare the text input for neural network training. These steps ensured that textual data could be effectively processed by the learning algorithms.

The core analytical component of the project involves **model training and evaluation**. Both traditional machine learning algorithms and deep learning models were implemented to compare performance. The dataset was split into training and testing sets, and evaluation metrics such as accuracy, precision, recall, and F1-score were used to assess model effectiveness.

Finally, the trained model was deployed using a **web application framework**, allowing users to input news content and receive real-time predictions. This end-to-end methodology ensures a reliable, scalable, and user-friendly solution for fake news detection.

Tools and Technologies Used

The following tools and technologies were used in the implementation of this project:

- **Python**: Used for data preprocessing, model development, and evaluation
- **Pandas and NumPy**: Used for data manipulation and numerical computations
- **Scikit-learn**: Used for feature extraction, traditional machine learning models, and evaluation
- **TensorFlow / Keras**: Used for implementing deep learning models

- **NLTK**: Used for natural language processing tasks
- **Streamlit**: Used for developing the web-based user interface
- **GitHub**: Used for version control and project sharing.

WORKFLOW

The workflow of the Fake News Detection System is designed in a structured and sequential manner to ensure accuracy, efficiency, and clarity at each stage of implementation. The workflow begins with data ingestion and progresses through preprocessing, model training, evaluation, and deployment, resulting in a complete end-to-end solution.

The first stage of the workflow involves **data ingestion**, where labeled datasets containing real and fake news articles are loaded into the Python environment. These datasets include textual information such as news headlines and article content. At this stage, the data is examined to understand its structure, size, and class distribution.

The second stage is **data cleaning and preprocessing**. In this step, unnecessary elements such as URLs, special characters, punctuation, and stop words are removed from the text. All text is converted to lowercase to maintain consistency. Tokenization and lemmatization are applied to standardize words and reduce variations, ensuring that the textual data is suitable for further analysis.

The third stage focuses on **feature extraction and text representation**. For traditional machine learning models, techniques such as TF-IDF vectorization are used to convert textual data into numerical feature vectors. For deep learning models, tokenization and sequence padding are applied to represent text as fixed-length numerical sequences.

Following feature extraction, the next stage is **model training**. The dataset is divided into training and testing sets. Machine learning and deep learning models are trained on the processed data to learn patterns that distinguish real news from fake news. Hyperparameters are adjusted to improve model performance and reduce overfitting.

The fifth stage involves **model evaluation**. The trained models are evaluated using performance metrics such as accuracy, precision, recall, and F1-score. This stage helps compare different models and identify the most effective approach for fake news detection.

The final stage of the workflow is **deployment and testing**. The selected model is integrated into a web-based application using Streamlit, allowing users to input news content and receive real-time prediction results. This workflow ensures a smooth transition from raw data to an interactive and usable fake news detection system.

ASSUMPTIONS

During the development and implementation of the Fake News Detection System, certain assumptions were made to define the scope of analysis and ensure consistency throughout the project. These assumptions are common in machine learning-based systems and help establish clear operational boundaries.

One of the primary assumptions of this project is that the input news content contains sufficient textual information to allow meaningful classification. Extremely short or ambiguous inputs may not provide enough context for accurate prediction and are therefore assumed to be outside the effective operating range of the system.

It is also assumed that the labeled datasets used for training and evaluation are accurate and reliable. The system relies on the correctness of the provided labels to learn patterns distinguishing real news from fake news. Any inconsistencies or mislabeling in the dataset may affect the model's performance.

Another important assumption is that linguistic patterns associated with fake news remain relatively stable over time. The model is trained on historical data and assumes that future fake news content follows similar writing styles and structures. Significant changes in writing patterns or the use of advanced content generation techniques may require retraining the model.

The project further assumes that text-based analysis alone is sufficient for the current scope of fake news detection. External factors such as source credibility, author reputation, and real-time fact verification are not considered in this implementation.

These assumptions help define the limitations of the system while ensuring that the proposed solution remains practical, focused, and suitable for academic evaluation.

IMPLEMENTATION – DATA COLLECTION AND PROCESSING STEPS

The implementation phase of the Fake News Detection System begins with the collection of textual data, followed by systematic preprocessing and transformation to make the data suitable for machine learning and deep learning models. This stage plays a critical role in ensuring the accuracy and reliability of the classification results.

Data Collection

The datasets used in this project were obtained from publicly available and well-known fake news datasets. These datasets contain labeled news articles and social media text categorized as real or fake. The collected data includes news headlines and article content written in English, covering a variety of topics such as politics, health, science, and current events.

The datasets were loaded into the Python environment using the Pandas library. An initial exploratory analysis was conducted to understand the dataset size, structure, class distribution, and presence of missing or inconsistent values. This preliminary analysis helped identify necessary preprocessing steps before model training.

Data Cleaning and Preprocessing

Once the data was loaded, several preprocessing techniques were applied to improve data quality. All text was converted to lowercase to maintain uniformity. Unnecessary elements such as URLs, special characters, punctuation marks, and numerical values were removed using text cleaning techniques.

Stop words that do not contribute meaningful information were eliminated, and tokenization was performed to split text into individual words. Lemmatization was then applied to reduce words to their base form, helping minimize redundancy and improve feature consistency. These preprocessing steps ensured that the textual data was clean, standardized, and suitable for feature extraction.

Feature Preparation

After preprocessing, the cleaned text was transformed into numerical representations required for model training. For traditional machine learning models, TF-IDF vectorization was used to convert text into weighted feature vectors. For deep learning models, tokenization and sequence padding were applied to generate fixed-length numerical input sequences.

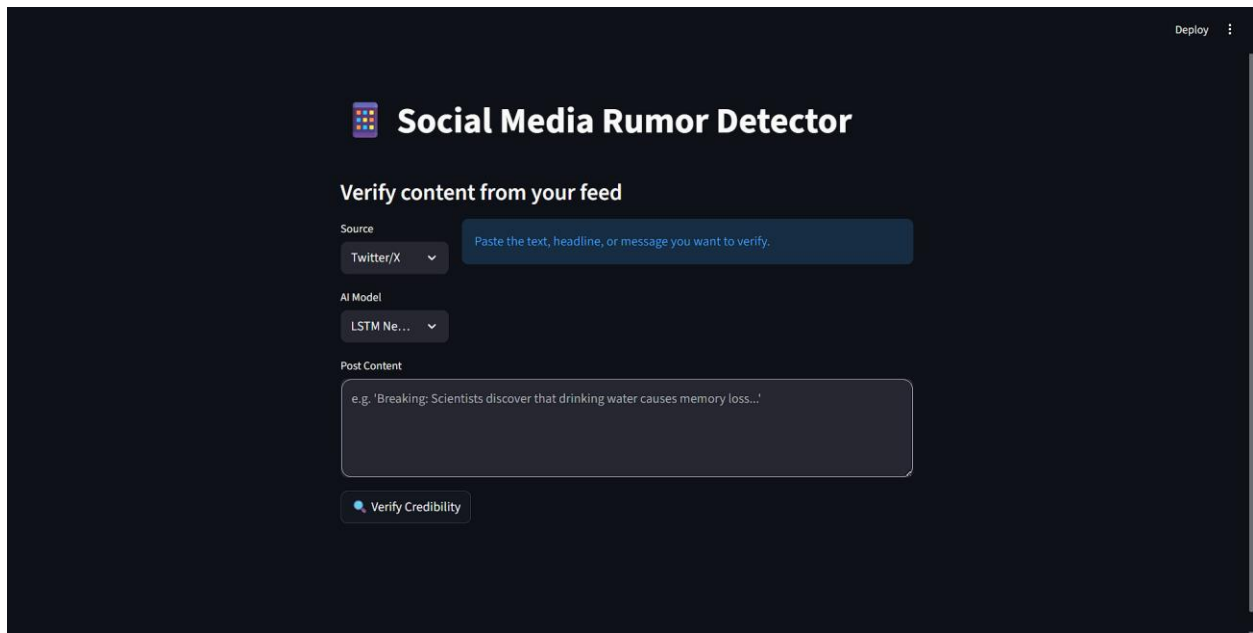
Before training, the dataset was divided into training and testing sets to enable proper evaluation of model performance. This structured data collection and processing pipeline ensured that the input data was accurate, consistent, and representative of real-world news content.

IMPLEMENTATION – DIAGRAMS, CHARTS, AND TABLES

Visual representation of system functionality and output plays an important role in demonstrating the practical implementation of the Fake News Detection System. In this project, screenshots of the web application interface and prediction results were used to illustrate how the system operates and how users interact with it in real time.

Web Application Interface Visualization


To demonstrate the user-facing component of the system, screenshots of the web-based interface developed using Streamlit were included. The interface provides a simple text input field where users can enter news content for verification. This visualization highlights the accessibility and ease of use of the application, making it suitable for non-technical users.



Prediction Output Visualization

Screenshots displaying the prediction results generated by the system were captured to demonstrate the working of the model. After the user submits news content, the system processes the input text and classifies it as either Real News or Fake News. The output is displayed instantly on the interface, allowing users to verify the authenticity of news content in a clear and understandable manner.

Deploy

 **Social Media Rumor Detector**

Verify content from your feed

Source

Twitter/X

Paste the text, headline, or message you want to verify.

AI Model

LSTM Ne...

Post Content

Headline: Scientists Confirm Drinking Hot Water Cures Cancer
Article:
Researchers from an unnamed foreign university have proven that drinking hot water three times a day completely cures cancer. Doctors say pharmaceutical companies are hiding this discovery to protect profits.

Verify Credibility

▲


LIKELY FAKE / MISLEADING

Our AI is 100.0% confident this content is fabricated.

Multiple Test Case Results

To further validate the effectiveness of the system, screenshots of predictions for different input news samples were included. These examples demonstrate the system's ability to handle varied inputs and consistently produce classification results. Presenting multiple test cases helps showcase the reliability and practical applicability of the fake news detection model.

Deploy

 **Social Media Rumor Detector**

Verify content from your feed

Source

Twitter/X

Paste the text, headline, or message you want to verify.

AI Model

LSTM Ne...

Post Content

Headline: WhatsApp Will Start Charging ₹499 Per Month
Article:
WhatsApp has officially announced a ₹499 monthly subscription fee for Indian users starting next month. Users who do not pay will permanently lose access to all chats and contacts.


Verify Credibility

▲

LIKELY FAKE / MISLEADING

Our AI is 95.1% confident this content is fabricated.

Deploy

 **Social Media Rumor Detector**

Verify content from your feed

Source

Twitter/X

Paste the text, headline, or message you want to verify.


AI Model

LSTM Ne...

Post Content

Headline: ISRO Successfully Launches Weather Satellite
Article:
The Indian Space Research Organisation (ISRO) successfully launched a new weather satellite from Sriharikota on Tuesday. The satellite will help improve cyclone prediction and climate monitoring across the Indian Ocean region.

Verify Credibility

 **LIKELY REAL**

Our AI is 100.0% confident this content is authentic.

Summary

The use of interface screenshots and prediction outputs helped clearly demonstrate the working of the Fake News Detection System. These visual elements focus on system usability and real-time performance rather than internal model complexity, making the implementation easy to understand and evaluate. Overall, the visuals support effective presentation of the system's functionality and real-world application.

SOLUTION DESIGN

The solution designed in this project follows a modular and layered architecture, where each component performs a specific function within the overall Fake News Detection System. This design approach ensures clarity, scalability, and ease of maintenance, while allowing future enhancements to be incorporated with minimal changes.

At the core of the solution is the **data processing layer**, which is responsible for handling raw textual data. This layer includes data loading, text cleaning, and preprocessing operations such as tokenization, stop-word removal, and lemmatization. The purpose of this layer is to convert unstructured text into a clean and standardized format suitable for analysis.

The next layer is the **feature representation layer**. In this layer, preprocessed text is transformed into numerical representations that machine learning models can process. Traditional models use TF-IDF vectorization to capture word importance, while deep learning models use tokenized and padded sequences to preserve contextual information within the text.

The **modeling layer** forms the core analytical component of the system. This layer implements machine learning and deep learning models trained to classify news as real or fake. The models learn patterns and relationships within the textual data that help distinguish between genuine and misleading information. Model evaluation is also handled within this layer using standard performance metrics.

Built on top of the modeling layer is the **application layer**, which provides user interaction through a web-based interface. The trained model is integrated into a Streamlit application that allows users to input news content and receive prediction results in real time. This layered solution design ensures a smooth flow from raw input data to actionable output while maintaining system reliability and usability.

CHALLENGES AND OPPORTUNITIES

Challenges

During the development of the Fake News Detection System, several challenges were encountered at different stages of the project. One of the primary challenges was handling large volumes of unstructured textual data. News articles and social media content often contain noise such as irrelevant symbols, informal language, and inconsistent formatting, which required extensive preprocessing to ensure data quality.

Another significant challenge was dealing with class imbalance in the dataset. In some cases, fake news samples outnumbered real news samples, which could lead to biased model predictions. Addressing this issue required careful dataset splitting and evaluation to ensure fair performance across both classes.

Feature representation also posed a challenge, as converting textual data into meaningful numerical features is a complex task. Traditional feature extraction techniques sometimes failed to capture contextual meaning, while deep learning models required careful tuning to avoid overfitting and excessive computational cost.

From a deployment perspective, integrating the trained model into a web-based application required additional effort to ensure smooth performance, real-time prediction, and user-friendly interaction.

.

Opportunities

Despite these challenges, the project revealed several opportunities for improvement and future expansion. One key opportunity lies in adopting more advanced deep learning models such as transformer-based architectures to further enhance classification accuracy.

There is also potential to improve system reliability by incorporating additional features such as source credibility analysis and real-time fact verification. Expanding the dataset to include multiple languages can help make the system more versatile and widely applicable.

Overall, the challenges faced during the project contributed to valuable learning experiences, while the opportunities identified highlight the potential for extending the system into a more comprehensive fake news detection solution.

REFLECTIONS ON THE PROJECT

This project served as a valuable learning experience that significantly enhanced my understanding of machine learning, natural language processing, and real-world data analysis. Working on the Fake News Detection System allowed me to apply theoretical concepts to a practical problem, helping me gain hands-on experience across the complete machine learning workflow.

One of the most important lessons learned during this project was the importance of data preprocessing and text cleaning. I realized that the quality of input data plays a critical role in determining model performance. Handling unstructured textual data, removing noise, and standardizing text required careful attention and significantly influenced the effectiveness of the classification models.

The project also helped me understand the strengths and limitations of different machine learning approaches. Comparing traditional machine learning models with deep learning models provided valuable insights into how contextual understanding impacts performance in text classification tasks. This experience strengthened my understanding of model selection, evaluation metrics, and performance trade-offs.

Additionally, deploying the trained model through a web-based interface improved my practical skills in integrating machine learning models into real applications. Overall, this project enhanced my technical skills, analytical thinking, and problem-solving ability, and strengthened my interest in artificial intelligence and data science.

RECOMMENDATIONS

Based on the analysis and outcomes of the Fake News Detection System, several recommendations can be made to enhance the effectiveness and real-world applicability of the solution. The system can be further improved by incorporating advanced natural language processing techniques that better capture contextual and semantic relationships within text.

Regular retraining of the model using updated datasets is recommended to ensure that the system remains effective against evolving patterns of fake news. As misinformation techniques change over time, continuous learning will help maintain classification accuracy.

Integrating additional features such as source credibility analysis and metadata-based evaluation can further strengthen the system's reliability. Providing users with explainable outputs, such as highlighting suspicious phrases or keywords, can also improve transparency and trust in the system.

Overall, adopting a continuous improvement approach and expanding the system's analytical capabilities can help transform the project into a more robust and scalable fake news detection solution.

OUTCOME / CONCLUSION

The Fake News Detection System developed in this project demonstrates the effective application of machine learning and natural language processing techniques to address the growing problem of misinformation in digital media. By systematically preprocessing textual data and applying supervised learning models, the system is capable of classifying news content as real or fake with a high level of accuracy.

Both traditional machine learning models and deep learning approaches were implemented and evaluated to assess their effectiveness in fake news classification. The comparative analysis helped identify the strengths of context-aware models in capturing complex linguistic patterns. The evaluation results indicate that the developed system performs reliably on unseen data and provides consistent prediction outcomes.

The deployment of the trained model through a web-based interface further enhances the practical usability of the system. Users can easily input news content and receive real-time predictions, making the solution accessible and interactive. Overall, the project successfully met its objectives and highlights the importance of AI-driven solutions in supporting misinformation detection and informed decision-making.

ENHANCEMENT SCOPE

While the current implementation of the Fake News Detection System successfully meets the project objectives, there are several opportunities for future enhancement that can further improve accuracy, scalability, and real-world applicability.

One important enhancement is the adoption of advanced deep learning models such as **transformer-based architectures**. Models like BERT or RoBERTa can capture deeper contextual relationships within text and may significantly improve performance on complex and subtle fake news content. Although these models require higher computational resources, they offer improved semantic understanding.

Another enhancement involves expanding the system to support **multilingual fake news detection**. Currently, the system processes only English-language text. Extending the model to handle multiple languages can increase its usability across different regions and platforms.

The system can also be enhanced by integrating **source credibility analysis and metadata features**, such as publisher reliability, publication date, and historical accuracy. Combining textual analysis with metadata-based evaluation can improve detection reliability.

Additionally, incorporating **explainable AI techniques** can help users understand why a particular piece of news is classified as fake or real. Highlighting influential words or phrases can improve transparency and user trust.

These enhancements demonstrate that the current system provides a strong foundation for developing a more advanced, scalable, and intelligent fake news detection solution in the future.

.

LINK TO CODE AND EXECUTABLE FILE

As part of the project deliverables, the complete source code, trained models, and application files for the Fake News Detection System have been organized and shared through a public repository. This ensures transparency, reproducibility, and ease of evaluation.

The repository contains Python scripts used for data preprocessing, feature extraction, model training, evaluation, and deployment. It also includes the trained model files and supporting configuration files required to run the application.

The web-based application developed using Streamlit allows users to input news text and obtain real-time prediction results. The project can be executed locally by installing the required dependencies and running the application script.

Project Repository Link:

<https://github.com/AhsanVA/fake-news-detection>

Execution Instructions:

- Install dependencies: `pip install -r requirements.txt`
- Run application: `streamlit run app/app.py`

RESEARCH QUESTIONS AND RESPONSES

This section presents the key research questions addressed during the project along with concise responses based on the implementation and evaluation results.

Research Question 1

Can machine learning models effectively detect fake news based on textual data?

Response:

Yes. The results of this project demonstrate that machine learning and deep learning models can effectively identify fake news by learning patterns from textual data. Proper preprocessing and feature representation significantly contribute to classification accuracy.

Research Question 2

How does deep learning perform compared to traditional machine learning models in fake news detection?

Response:

Deep learning models showed improved performance due to their ability to capture contextual and sequential information in text. Compared to traditional models, deep learning approaches provided better generalization on complex news content.

Research Question 3

What role does text preprocessing play in fake news detection?

Response:

Text preprocessing plays a crucial role in improving model performance. Cleaning, tokenization, and lemmatization help reduce noise and ensure consistent input representation, which directly impacts classification accuracy.

Research Question 4

How does deploying the model as a web application add value?

Response:

Deploying the model through a web interface improves accessibility and usability. It allows users to interact with the system easily and receive real-time predictions, making the solution practical for real-world use.

Summary

The research questions highlight the effectiveness of machine learning-based approaches for fake news detection and demonstrate the importance of preprocessing, model selection, and deployment in building a complete and usable system.

REFERENCES

1. Wang, W. Y. (2017). *“Liar, Liar Pants on Fire”: A New Benchmark Dataset for Fake News Detection*.
Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (ACL).
<https://aclanthology.org/P17-2067/>
2. Ahmed, H., Traore, I., & Saad, S. (2017). *Detection of Online Fake News Using N-Gram Analysis and Machine Learning Techniques*.
Springer – Intelligent, Secure, and Dependable Systems in Distributed and Cloud Environments.
https://link.springer.com/chapter/10.1007/978-3-319-69155-8_9
3. Shu, K., Sliva, A., Wang, S., Tang, J., & Liu, H. (2017). *Fake News Detection on Social Media: A Data Mining Perspective*.
ACM SIGKDD Explorations Newsletter, 19(1).
<https://dl.acm.org/doi/10.1145/3137597.3137600>
4. Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*.
Proceedings of NAACL-HLT.
<https://aclanthology.org/N19-1423/>
5. Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press.
<https://www.deeplearningbook.org/>
6. Pedregosa, F., et al. (2011). *Scikit-learn: Machine Learning in Python*.
Journal of Machine Learning Research, 12.
<https://jmlr.org/papers/v12/pedregosa11a.html>