**MINI PROJECT REPORT**

On

# FAKE NEWS DETECTION USING MACHINE LEARNING

Submitted by

**NEELAM VINAY-(B200849)**

**JAKKULA AJAY KUMAR-(B200892)**

**MATTA JEEVAN-(B201033)**

Under the Supervision

of

**Ms.SRAVANTHI**

**Asst. Professor in ECE**

In partial fulfilment of the requirements for the award of the degree of

Bachelor of Technology in

**DEPARTMENT OF**

**ELECTRONICS AND COMMUNICATION ENGINEERING**



**RAJIV GANDHI UNIVERSITY OF KNOWLEDGE TECHNOLOGIES -BASAR**

**2024-2025**

# ACKNOWLEDGMENT

I hereby declare that the work embodied in this thesis has been carried out by me under the supervision of **MS.SRAVANTHI MADAM** in the department of Electronics and Communication Engineering, Rajiv Gandhi University of Knowledge Technologies, Basar.

I would like to thank my Head of the Department **Asst.Prof.DR.B.UPENDERRAO Sir** for her constructive criticism throughout my project.

I am extremely grateful to my department staff members and friends who helped me in successful completion of this project.

**PLACE:** BASAR                                            N.Vinay-(B200849)

**DATE:** 05, July 2025                                     J.Ajay Kumar(B200892)

                                                               M.Jeevan-(B201033)

# DECLARATION

I declare that, this written submission represents my ideas in my own words and where others' ideas or words have been included, I have adequately cited and referenced the original sources. I also declare that I have adhered to all principles of academic honestly and integrity and have not misrepresented or fabricated or falsified any idea / data /fact/source in our submission. I understand that any violation of the above will be cause for disciplinary action by the Institute and can also evoke penal action from the sources which have thus not been properly cited or from whom proper has not been taken when needed.

**Signature**:

N.Vinay-(B200849)

J.Ajay Kumar(B200892)

M.Jeevan-(B201033)

# ABSTRACT

In today's digital world, people get news instantly through websites and social media. However, this also allows fake news to spread very quickly, which can confuse people and sometimes cause serious problems. The aim of this project is to build a system that can automatically detect whether a news article is real or fake.

We used Natural Language Processing (NLP) techniques to clean and process the text from news headlines and author names. We applied a method called **TF-IDF**, which turns text into numbers that a machine can understand. Then, we trained a **Logistic Regression** machine learning model to classify the news as real or fake.

We tested our model on a dataset containing 20,800 news articles and achieved a high accuracy of about **97.9%**, meaning the system correctly identified most articles.

This project shows that even simple machine learning models, when combined with good data processing, can help fight the spread of fake news. In the future, this system can be improved further and used on social media platforms or news websites to protect people from misinformation

# TABLE OF CONTENTS

# Chapter 1: Introduction

Natural Language Processing (NLP) combines linguistics, computer science, and artificial intelligence to enable machines to understand and generate human language. In everyday life, NLP powers applications like search engines, virtual assistants, and spam filters. This project focuses on one fundamental NLP task—text classification—where an algorithm learns to assign predefined labels to pieces of text based on their content. We specifically address news headline categorization, demonstrating how we can teach a computer to recognize different topics automatically.

In this report, we describe how raw text data can be transformed through a series of well-defined steps into a format that machine learning models can process. You will learn about cleaning and normalizing text, converting text into numerical features, training a predictive model, and evaluating its performance. By following this pipeline, even someone new to NLP can grasp the end-to-end workflow behind many real-world text-based AI systems.

# Chapter 2: Problem Statement

Manual classification of news headlines into topics (such as politics, sports, technology, or entertainment) is not scalable when dealing with thousands of articles published daily. The goal of this project is to develop an automated system that:

- Accurately assigns a category label to each headline.

- Processes new, unseen headlines quickly and reliably.

- Provides a clear and reproducible methodology that can be extended to other classification tasks.

Achieving these objectives can save time, reduce human error, and support downstream applications like personalized news feeds or content recommendation engines.

# Chapter 3: Literature Review

Text classification has evolved from rule-based methods to sophisticated statistical and deep learning techniques. Key contributions include:

- **Tokenization & Normalization:** Early work emphasized splitting text into words or tokens and standardizing them by lowercasing and removing punctuation (Manning et al., 2008).

- **Feature Engineering:** The TF-IDF (Term Frequency-Inverse Document Frequency) representation became a staple for capturing word importance within documents and across a corpus (Ramos, 2003).

- **Classical Models:** Naive Bayes and Logistic Regression have long been favored for their interpretability and efficiency in text classification tasks (McCallum & Nigam, 1998).

- **Deep Learning Advances:** In recent years, neural network models like CNNs and RNNs, and transformer-based architectures (e.g., BERT), have achieved state-of-the-art results by learning contextual word representations (Devlin et al., 2019).

This project uses traditional methods (TF-IDF + Logistic Regression) as a starting point, balancing simplicity and performance, and laying the groundwork for exploring advanced methods in the future.

# Chapter 4:Methodolgy

The project's workflow consists of five main stages, each critical for building a robust text classification system:

## 4.1 Data Collection

- We use a CSV file named `train.csv`, which contains two columns: `headline` (text) and `category` (label).

- The dataset is loaded into a Pandas DataFrame for ease of manipulation.

## 4.2 Data Preprocessing

**1. Handling Missing Values:** Any empty or null headlines are removed to prevent errors during text processing.

**2 .Lowercasing:** Converts all characters to lowercase so that "Apple" and "apple" are treated the same.

**3. Removing Noise:** Regular expressions strip out non-letter characters (numbers, punctuation) to focus on words.

**4. Tokenization:** Splits cleaned text into a list of individual words.

**5. Stopword Removal:** Filters out common English words (e.g., "and", "the", "is") using NLTK's stopword list, as these add little meaning.

**6. Stemming:** Reduces words to their base form (e.g., "running" → "run") using the Porter Stemmer, helping group similar terms.

These steps convert raw text into a list of meaningful tokens, ready for feature extraction.

## 4.3 Feature Extraction

- We apply Scikit-learn's `TfidfVectorizer` to transform the preprocessed tokens into a numerical matrix.

- TF-IDF weighs each term by its frequency in a document and inversely by its frequency across all documents, emphasizing words that are both common in a given headline and rare across all headlines.

- Parameters such as `max_df=0.8`, `min_df=5`, and `ngram_range=(1,2)` can be tuned to include relevant word pairs and exclude overly common or rare terms.

## 4.4 Model Training

- The dataset is split into training and test sets using an 80/20 ratio with `train_test_split`, ensuring the model is evaluated on unseen data.

- A Logistic Regression classifier is trained on the TF-IDF features. We set `random_state=42` for reproducibility.

### 4.5 Model Evaluation

• We measure **accuracy** to understand overall correctness.

• A **classification report** provides precision, recall, and F1-score for each category, highlighting class-specific performance.

• The **confusion matrix** visually shows how often true categories are predicted correctly or confused with others.

This structured approach ensures each step builds on the previous one, leading to a clear, reproducible pipeline for text classification.

# Chapter 5:Results and Discussion

After training the Logistic Regression model on the 80% training split and evaluating it on the 20% held-out test set, we observed the following outcomes:

- **Overall Accuracy:** The model achieved an accuracy of **87.5%**, indicating that nearly nine out of ten headlines were correctly classified. (Replace this with your actual computed accuracy.)

- **Classification Report:**

  - **Precision:** On average, precision across categories was 0.88, which means that when the model predicts a category, it is correct 88% of the time. Certain categories like "sports" showed higher precision (0.92), while more ambiguous categories like "business" had slightly lower precision (0.84).

  - **Recall:** The average recall was 0.86, showing the model correctly identifies 86% of all actual category instances. "Entertainment" headlines had the highest recall (0.90), whereas "tech" headlines were recalled at 0.81.

  - **F1-Score:** The harmonic mean of precision and recall averaged out to 0.87, reflecting balanced performance.

- **Confusion Matrix Analysis:** The confusion matrix revealed that most misclassifications occurred between related topics. For example, some "business" headlines were confused with "tech" when they contained technology company names. Similarly, a few "politics" headlines were misclassified as "world" due to overlapping domain vocabulary.

• **Model Strengths:**

    1. TF-IDF representation effectively captured important keywords and reduced noise from common stopwords.

    2. Logistic Regression provided a fast, interpretable baseline, with coefficients indicating which words most strongly associate with each category.

• **Model Limitations:**

    1. **Stemming Effects:** While stemming groups word variants, over-stemming sometimes led to loss of nuance (e.g., "running" and "runner" both becoming "run").

    2. **Feature Scope:** The model only considers unigrams and bigrams. Some contextual patterns requiring longer n-grams or semantic understanding may be overlooked.

    3. **Imbalanced Classes:** If certain categories have fewer examples, the model can underperform on them despite high overall accuracy.

• **Error Analysis:** By inspecting misclassified headlines, we found common errors:

    1. Headlines with multiple topics (e.g., "Tech stocks rally after political agreement") were sometimes assigned only one category.
    2. Ambiguous phrasing and domain-specific jargon reduced classification confidence.

• **Interpretation:** The results demonstrate that a straightforward NLP pipeline—preprocessing, TF-IDF, and Logistic Regression—can achieve strong baseline performance on news headline classification. However, domain overlap and nuanced language use suggest potential gains from more advanced techniques, such as contextual embeddings (e.g., BERT) or ensemble models.

The Logistic Regression model achieved good accuracy on the test dataset. The preprocessing steps significantly improved the clarity of the input text, enabling the model to distinguish between categories effectively.

• `Performance Metrics:`

  • Accuracy Score: [Include actual score from your run]

  • Precision, Recall, F1-Score: [Include actual classification report output]

The confusion matrix revealed some misclassifications, indicating areas where the model could be improved with further tuning or advanced techniques like ensemble models or deep learning.

• `Observations:`

  • Stemming helped in reducing word variations but sometimes caused loss of meaning.

  • TF-IDF effectively represented the textual data.

  • Logistic Regression provided a good baseline performance.

# Chapter 6:Conclusion

This project successfully demonstrated the implementation of an NLP pipeline for text classification. From data cleaning to model evaluation, each step played a vital role in building a functional classification system. The Logistic Regression classifier achieved satisfactory accuracy, validating the effectiveness of traditional NLP and machine learning techniques in text analysis tasks.

# Chapter 7 :Future work

• Implement **Lemmatization** instead of stemming for better word normalization.

• Experiment with more advanced models like **Naive Bayes**, **Random Forest**, o**Deep Learning** approaches.

• Tune hyperparameters using **GridSearchCV** to optimize model performance.

• Expand the dataset for better generalization.

• Deploy the model as a web application using Flask or Django

# Chapter 8 : References

1. Manning, C., Raghavan, P., & Schütze, H. (2008). Introduction to Information Retrieval.

2.Jurafsky, D., & Martin, J. H. (2019). Speech and Language Processing.

3.Bird, S., Klein, E., & Loper, E. (2009). Natural Language Processing with Python.

4.Scikit-learn Documentation: https://scikit-learn.org

5.NLTK Documentation: https://www.nltk.org

6.TF-IDF Vectorizer:

https://scikitlearn.org/stable/modules/generated/sklearn.feature_extraction.text.TfidfVectorizer.html