

# **VISVESVARAYA TECHNOLOGICAL UNIVERSITY**

**“JNANA SANGAMA” Belagavi-590018, Karnataka**



## ***Project Phase 1***

### **“Automated Crawling, Classification, and Sentiment Analysis of Digital News with a Built-in Feedback”**

*Submitted in partial fulfillment of the requirements for the degree of*

**BACHELOR OF ENGINEERING**

**in**

**COMPUTER SCIENCE AND ENGINEERING**

**Submitted by**

**Thrivedhi Sreenivas D (1BY21CS224)**

**Raghvendra Sharma (1BY21CS225)**

**Sai Ravi Teja J (1BY21CS226)**

**Chinmay M Pawar (1BY22CS400)**

**Under the Guidance of**

**Dr. Gerard Deepak**

**Associate Professor**



**DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING**

**BMS INSTITUTE OF TECHNOLOGY & MANAGEMENT**

**(An Autonomous Institution affiliated to VTU, Belagavi)**

**BENGALURU-560064**

# **BMS INSTITUTE OF TECHNOLOGY & MANAGEMENT**

(An Autonomous Institution affiliated to VTU, Belagavi)  
**YELAHANKA, BENGALURU – 560064**

## **DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING**



### **CERTIFICATE**

This is to certify that the Final Year Project Phase 1 work entitled “**Automated Crawling, Classification, and Sentiment Analysis of Digital News with a Built-in Feedback System**” is a bonafide work carried out by **Thrivedhi Sreenivas D (1BY21CS224)**, **Raghvendra Sharma (1BY21CS225)**, **Sai Ravi Teja J (1BY21CS226)**, **Chinmay M Pawar (1BY22CS400)** in partial fulfillment for the award of a **Bachelor of Engineering Degree in Computer Science and Engineering** of **Visvesvaraya Technological University, Belagavi** during the year 2024- 2025. It is certified that all corrections & suggestions indicated for Internal Assessment have been incorporated in this report. The Final Year Project Phase 1 report has been approved as it satisfies the academic requirements regarding Final year-project work for the Bachelor of Engineering Degree.

**Signature of the Guide**

Dr. Gerard Deepak  
Associate Professor  
Dept. of CSE, BMSIT&M

**Signature of the Associate Head**

Dr. Satish Kumar T.  
Assoc. Professor & Assoc. Head  
Dept. of CSE, BMSIT&M

# ABSTRACT

In the digital age, the rapid dissemination of news and misinformation poses significant challenges, particularly regarding government-related topics. The uncontrolled spread of misinformation can lead to political, social, and economic repercussions. This project addresses this pressing issue by developing a system that automates the crawling, classification, and sentiment analysis of news articles, targeting government-related topics. With a growing market demand for real-time news analytics and tools to combat misinformation, this framework provides a timely and impactful solution.

The methodology involved leveraging Python libraries such as BeautifulSoup and Selenium to crawl over 12,000 news articles and videos. Clustering techniques were applied to prepare a labeled dataset, which was then used to train the DistilBERT model for department classification, achieving 83% accuracy. Sentiment analysis was implemented using the Roberta model, while negative news alerts were automated via NodeMailer and Gmail-SMTP. The framework was integrated with a Django backend and user interface, enabling real-time interaction and multilingual support using the Google Translate API. Video news analysis was also incorporated by extracting audio, converting it to text, and applying the same classification and sentiment analysis methods.

Experimental analysis demonstrates the successful implementation of key objectives: automated crawling, classification, sentiment analysis, and a feedback mechanism for alerting relevant authorities. While real-time statistical analysis remains under development, the current system effectively addresses critical challenges, laying the groundwork for a scalable and impactful solution to misinformation in the digital news landscape.

## ACKNOWLEDGEMENT

We are happy to present this Final Year Project Phase 1 after completing it successfully. This project would not have been possible without the guidance, assistance and suggestions of many individuals.

We would like to express our deep sense of gratitude and indebtedness to each and every one who has helped us make this Final Year Project a success.

We heartily thank our Principal, **Dr. SANJAY H. A**, BMS Institute of Technology & Management, for his constant encouragement and inspiration in taking up this Final Year Project.

We heartily thank our Professor and Head of the Department, **Dr. THIPPESWAMY G** and Associate Professor and Associate Head **Dr. SATISH KUMAR T**, Department of Computer Science and Engineering, BMS Institute of Technology & Management, for his constant encouragement and inspiration in taking up this Final Year Project.

We gracefully thank our guide, Associate Professor, **Dr. GERARD DEEPAK**, Department of Computer Science and Engineering, BMS Institute of Technology & Management for their intangible support and constant backbone for our project.

Special thanks to all the staff members of Computer Science Department for their help and kind co-operation.

Lastly, we thank our parents and friends for the support and encouragement given throughout in completing this precious work successfully.

**Thrivedhi Sreenivas D (1BY21CS224)**

**Raghvendra Sharma (1BY21CS225)**

**Sai Ravi Teja J (1BY21CS226)**

**Chinmay M Pawar (1BY22CS400)**

# CONTENTS

## CHAPTER 1 INTRODUCTION

- 1.1 INTRODUCTION
- 1.2 PROBLEM STATEMENT
- 1.3 MOTIVATION
- 1.4 OBJECTIVES
- 1.5 CONTRIBUTIONS
- 1.6 ORGANIZATION OF THE REPORT

## CHAPTER 2 LITERATURE SURVEY

- 1.1 LITERATURE SURVEY
- 1.2 SUMMARY
- 1.3 GAPS IDENTIFIED AND ANALYSIS

## CHAPTER 3 PROPOSED METHODOLOGY

- 3.1 SYSTEM ARCHITECTURE
- 3.2 DFD LEVEL 0, 1,2
- 3.3 BEHAVIORAL DESIGN
- 3.4 SUMMARY

## CHAPTER 4 IMPLEMENTATION

- 4.1 PLATFORM DESCRIPTION AND TOOL USAGE
- 4.2 PARAMETRIC SETTINGS AND CONFIGURATIONS
- 4.2 DATASET DESCRIPTION AND STATISTICS
- 4.3 CHALLENGES
- 4.4 SUMMARY

## CHAPTER 5 RESULTS AND PERFORMANCE EVALUATION

- 5.1 METRICS DESCRIPTION
- 5.2 COMPARISON WITH BASELINE MODELS
- 5.3 QUALITATIVE ANALYSIS (Screenshots and Descriptions)
- 5.4 QUANTITATIVE ANALYSIS (Bar Graphs, Line Graphs, Tables)

## CHAPTER 6 CONCLUSION AND FUTURE ENHANCEMENT

- 6.1 SUMMARY
- 6.2 FUTURE WORK

## BIBLIOGRAPHY

---

# LIST OF TABLES

Table No.	Title	Page No.
2.1	Literature Survey .....	7
5.1	Comparison with base Model .....	27

## LIST OF FIGURES

Figure No.	Title	Page No.
3.1.1	News Crawling and Web Scraping.....	12
3.1.2	Classification Section.....	12
3.1.3	Sentiment Analysis Section.....	13
3.1.4	Real-Time Statistical Analysis Section.....	15
3.1.5	User Personalization Section.....	16
3.1.6	Backend Section.....	16
4.1.1	News Section.....	21
4.1.2	UI Home Page.....	22
4.3.1	Data Set Description.....	23
4.3.2	News Categories.....	24
5.1.1	Precision, Recall & F1-Score.....	26
5.3.1	Code for word count example.....	28
5.3.2	Bar graph word count.....	29
5.4.1	Scatter plot.....	30
5.4.2	Code for Scatter plot.....	30





# 1. INTRODUCTION

## 1.1 Introduction

The exponential growth of digital news platforms has made access to information easier than ever before. However, this has also led to the rampant spread of misinformation, especially on critical topics related to government policies and initiatives. This project focuses on developing an automated framework to address these challenges by categorizing news, analyzing sentiment, and identifying potentially harmful or misleading content. Such a system is essential in fostering transparency, improving public trust, and assisting government bodies in addressing misinformation swiftly.

## 1.2 Problem Statement

The spread of fake news and hate speech in digital media creates widespread political, social, and economic issues. Existing platforms lack the ability to effectively classify news articles by government departments, detect sentiment accurately, and provide timely alerts about negative or misleading content. This project aims to fill this gap by offering an automated, scalable, and user-friendly solution to analyze digital news efficiently and promote reliable information dissemination.

## 1.3 Motivation

Misinformation, particularly in the context of government policies, can have severe social, political, and economic consequences. Ministries need to stay informed about the content circulating in the media and quickly address any incorrect or negative news. This project is motivated by:

- The rise in misinformation and hate speech in digital news.
- The need for government bodies to be alerted about news related to their policies and work.
- The ability to streamline news analysis for users, ensuring they receive accurate information.
- Providing transparency regarding which news outlets are covering specific types of news (national, anti-national, sports, etc.)

## 1.4 Objectives

The main objectives of this project are:

- To develop a framework that automatically crawls news from various websites, including regional ones, in real-time.
- To develop a system that classifies news based on ministries, such as defense, finance, and education.
- To develop sentiment analysis tools to classify news as positive, neutral, or negative, and detect hate speech.
- To create a feedback mechanism for alerting ministries to fact-check negative news.

## 1.5 Contributions

This project makes the following key contributions:

- Curated a labeled dataset of news articles specifically tailored for training machine learning models.
- Enhanced classification accuracy through diverse and well-structured data representation.
- The dataset is designed to support various use cases, including sentiment analysis and multilingual classification.
- Implemented sentiment analysis using state-of-the-art NLP models like Roberta and DistilBERT for high precision.
- Incorporated hate speech detection mechanisms to identify and flag inappropriate content.
- Achieved robust performance by fine-tuning pre-trained models on domain-specific datasets.
- Created a cohesive framework that seamlessly combines news crawling, classification, sentiment analysis, and feedback notifications.
- Automated the end-to-end process to ensure scalability and minimal manual intervention.

## 1.6 Organization of the Report

The remainder of the report is structured as follows:

- **Chapter 2:** Literature Review, detailing existing systems and their limitations.
- **Chapter 3:** Proposed Methodology, explaining the framework design and implementation strategies.
- **Chapter 4:** Experimental Analysis, presenting results and evaluation of the system.
- **Chapter 5:** Conclusion and Future Work, summarizing findings and exploring potential enhancements.

## 2. LITERATURE SURVEY

### 2.1 Literature Survey

This chapter discusses the survey of the related work studied to meet the objectives mentioned.

Kaur et al., [1] (2022) have proposed the framework collects live news articles and performs sentiment analysis using various machine learning models. Its limitations include potential issues with data scraping due to website restrictions, and its gaps lie in the limited focus on news data, which may not be generalizable to other domains.

Chaturvedi et al.,[2] (2023) presented a comprehensive review of sentiment analysis, focusing on tasks, applications, and the role of deep learning techniques. The study explores trends and advancements in applying sentiment analysis across industries. Limitations include minimal discussion of real-world system implementation, with gaps in analyzing hybrid models that combine traditional and deep learning techniques.

Halawani et al., [3] (2023) have proposed a framework for automated sentiment analysis in social media using Harris Hawks optimization and deep learning techniques. The framework optimizes deep learning models with Harris Hawks optimization to analyze sentiments in social media posts. Its limitations involve high computational expense and dependence on large datasets, and its gaps are in the scalability of the method for real-time data processing.

Ahmad et al., [4] (2020) have proposed a fake news detection framework using machine learning ensemble methods. This framework employs ensemble learning to combine multiple algorithms for detecting fake news articles. Its limitations include the possibility of overfitting and the need for large amounts of labeled data, and its gaps include improving model generalizability and robustness to new fake news types.

Dashti et al., [5] (2024) have explored methods for handling multivariable missing data in causal mediation analysis. The study evaluates various multiple imputation (MI) approaches to address missingness in mediation analysis within epidemiological research. Limitations include potential biases when the mediator and/or outcome influence their own missingness, and gaps involve the need for further guidance on implementing MI in complex mediation models.

Amer & Siddiqui, [6] (2020) have proposed a framework for detecting COVID-19 fake news using random forest and decision tree classifiers. This model employs machine learning classifiers to identify fake news related to the COVID-19 pandemic. Its limitations include potential misclassification due to evolving fake news narratives and its gaps focus on improving classification accuracy with smaller, imbalanced datasets.

Akinyemi et al., [7] (2020) have proposed an improved classification model for fake news detection in social media. This model improves upon traditional methods for detecting fake news in social media using advanced classification techniques. Its limitations include a reduced performance when dealing with data imbalance and its gaps involve handling different types of fake news spreaders more effectively.

Balasubramanian & Chandran, [8] (2023) have proposed an automated sentiment analysis system for analyzing news feeds. The system automates sentiment analysis of news articles to gauge public opinion on various topics. Its limitations involve its inability to handle complex context-based sentiment analysis and its gaps include improving the system's ability to detect subtle emotions and context shifts in news content.

Duan et al., [9] (2020) have proposed a method for profiling fake news spreaders on Twitter. This approach identifies and profiles users who spread fake news on Twitter, offering valuable insights into the behavior of these users. Its limitations include difficulty distinguishing between intentional spreaders and misinformed users, and its gaps include improving user profiling for more accurate detection of fake news propagators.

Enders & Baraldi et al., [10] (2018) have proposed methods for handling missing data in psychometric studies. These methods improve the accuracy of psychometric tests by handling missing data effectively.

Wendy Ccoya and Edson Pinto et al.,[11] (2023) conducted a comparative analysis of libraries and machine learning models for sentiment analysis, including tools such as NLTK, TextBlob, and Transformers. The study evaluated algorithms like Decision Tree, SVM, and Naive Bayes, providing insights into the effectiveness of these approaches. While useful for selecting tools, the study was limited to a single dataset, which may restrict its generalizability.

Md. Taufiqul Haque Khan Tusar and Md. Touhidul Islam et al.,[12] (2021) applied NLP techniques and machine learning algorithms, such as SVM and Random Forest, for sentiment analysis on a dataset of US airline Twitter data. The research identified effective methods for customer feedback analysis, though its applicability to other industries may be limited due to the focus on a specific dataset.

Nhan Cach Dang, María N. Moreno-García, and Fernando De la Prieta et al.,[13] (2020) reviewed studies on sentiment analysis using deep learning models, incorporating techniques like TF-IDF and word embeddings. The paper offered valuable insights into the effectiveness of various models but did not cover advancements beyond 2020 and had a limited scope regarding datasets.

Ali Nazarizadeh, Touraj Baniroostam, and Minoos Sayyadpour et al.,[14] (2022) focused on sentiment analysis in the Persian language, reviewing studies from 2018 to 2022. The authors highlighted the effectiveness of approaches like BERT and LSTM, along with available datasets for the Persian language. However, the study's findings were limited to Persian and not generalizable to other languages.

Table 2.1: Literature Survey

Serial Number	Title	Methodology	Advantages	Limitations
1	"Sentiment analysis using web scraping for live news data with machine learning algorithms." <i>Materials Today: Proceedings</i> , 2022.	Web scraping of live news data, sentiment analysis using ML algorithms	Real-time data analysis, improved accuracy of sentiment classification	Limited to news articles, may face issues with data scraping due to website restrictions
2	"A comprehensive review of sentiment analysis, focusing on tasks, applications, and the role of deep learning techniques." (2023)	Review of sentiment analysis techniques and trends, emphasizing deep learning methods.	Detailed exploration of advancements in sentiment analysis across industries.	Minimal discussion of real-world system implementation ; lacks analysis of hybrid models.
3	"Automated sentiment analysis in social media using Harris Hawks optimization and deep learning techniques." Alexandria Engineering Journal, 2023.	Deep learning with Harris Hawks optimization for sentiment analysis	Enhanced accuracy with optimization techniques, handles social media data effectively	Computationally expensive, requires high data processing power
4	"Fake news detection using machine learning ensemble methods." Complexity, 2020.	Ensemble machine learning techniques for fake news detection	High accuracy in classification, combines multiple algorithms for better results	May be prone to overfitting, needs large datasets for training
5	"Methods for handling multivariable	Evaluation of multiple imputation (MI)	Addresses missingness in complex casual	Potential biases when and/or outcome influence their

	Missing data in casual mediation analysis.” 2024	Approaches for addressing missingness in epidemiological studies.	Mediation models, enhancing robustness.	Own missingness; limited guidance on MI in complex models
6	"Detection of COVID-19 fake news text data using random forest and decision tree classifiers." IJCSIS, 2020.	Random Forest and Decision Tree for fake news detection on COVID-19	Efficient detection with classic machine learning models, easy to implement	Limited by the scope of classifier accuracy, struggles with new and evolving fake news patterns
7	"An improved classification model for fake news detection in social media." IJITCS, 2020.	Improved classification models for detecting fake news	High detection accuracy, works well with social media data	Performance may decrease with data imbalance or very large datasets
8	"Automated sentiment analysis: An automated analysis of news feeds." Graduate Research in Engineering and Technology, 2023.	Automated sentiment analysis on news feeds using predefined classifiers	Simple implementation , good for standard sentiment categorization	Limited scope to predefined categories, struggles with context-based nuances
9	"Profiling fake news spreaders on Twitter." CLEF PAN-CLEF, 2020.	Profiling fake news spreaders on Twitter with statistical modeling	Effective for identifying spreaders, provides useful insights into user behavior	Can't always distinguish between intentional spreaders and misinformed users
10	"Missing data handling methods." Wiley Handbook of Psychometric Testing, 2018.	Methods for handling missing data in psychometrics	Addresses data incompleteness, improves psychometric measurement accuracy	Assumes data is missing at random, which may not always be the case
11	"Comparative Analysis of Libraries for Sentiment Analysis" (2023)	Comparative analysis of Python and R libraries (NLTK, TextBlob, Vader, Transformers) and machine learning models (SVM, Decision	Provides insights into the effectiveness of various libraries and models for sentiment analysis.	Based on a single dataset, which may limit generalizability to other datasets or domains.



		Tree, Naive Bayes).		
12	"A Comparative Study of Sentiment Analysis Using NLP and Different Machine Learning Techniques on US Airline Twitter Data" (2021)	NLP techniques (Bag-of-Words, TF-IDF) and machine learning algorithms (SVM, Logistic Regression, Naive Bayes, Random Forest) for sentiment analysis of airline Twitter data.	Identifies effective methods for sentiment analysis in customer feedback, particularly in the airline industry.	Accuracy of 77%, with room for improvement; focus is limited to a specific industry dataset.
13	"Sentiment Analysis Based on Deep Learning: A Comparative Study" (2020)	Review of deep learning models for sentiment analysis using TF-IDF, word embeddings, and neural networks like LSTM and Bi-LSTM.	Offers a comparative analysis of deep learning models and their effectiveness in sentiment classification.	Limited to the models and studies reviewed, which may not cover the latest advancements post-2020.
14	"Sentiment Analysis of Persian Language: Review of Algorithms, Approaches and Datasets" (2022)	Review of machine learning and deep learning approaches (BERT, LSTM) applied to sentiment analysis in the Persian language.	Comprehensive overview of sentiment analysis methodologies for the Persian language, highlighting useful algorithms and datasets.	Focuses solely on the Persian language, limiting applicability to other languages.

## 2.2 Summary

The reviewed papers cover a variety of techniques and methodologies in the fields of sentiment analysis, fake news detection, and missing data handling, all employing machine learning and optimization methods. Kaur et al. [1] propose a sentiment analysis system using web scraping for live news data, but their method faces limitations with data scraping accuracy. Jishag et al. [2] focus on sentiment analysis of product reviews, providing a flexible approach but struggling with sarcasm detection. Halawani et al. [3] improve social media sentiment analysis using deep learning and Harris Hawks optimization, though their approach requires significant computational resources. Ahmad et al. [4] detect fake news through ensemble learning, but their method is prone to overfitting and needs large datasets. Harel et al. [5] handle missing data in epidemiology with multiple imputation, assuming data is missing at random, which is not always the case. Amer & Siddiqui [6] apply random forest and decision trees for detecting COVID-19 fake news, but their model faces challenges with evolving fake news. Akinyemi et al. [7] improve fake news detection but struggle with data imbalance. Balasubramanian & Chandran [8] automate sentiment analysis of news feeds, though their method falls short in handling nuanced sentiments. Duan et al. [9] profile fake news spreaders on Twitter but have difficulty distinguishing between intentional spreaders and the misinformed. Enders & Baraldi [10] address missing data in psychometrics but make assumptions about data randomness, which may not always apply. The gaps identified across these studies include handling sarcasm, improving data scraping accuracy, optimizing models for computational efficiency, and tackling data imbalance in fake news detection. These gaps are addressed in our project through the integration of real-time data scraping, advanced sentiment analysis capable of detecting sarcasm, more robust machine learning models that are optimized for computational efficiency, and handling both balanced and imbalanced datasets for improved fake news detection accuracy. Additionally, our project incorporates advanced missing data techniques without assuming randomness, addressing a key limitation in the reviewed systems.

## 2.3 Gaps Identified and Analysis

Data scraping methods face limitations in consistency, as highlighted by Kaur et al. [1], necessitating the implementation of a more adaptive and reliable scraping framework for real-time data collection. Deep learning models, as noted by Halawani et al. [3], require significant computational power; however, this project focuses on optimizing models to reduce computational overhead without compromising accuracy. Overfitting and data imbalance, common challenges in fake news detection systems (Ahmad et al. [4]; Akinyemi et al. [7]), are addressed through the application of data augmentation and regularization to enhance model robustness and generalizability. Previous models often assume random missing data, as discussed by Harel et al. [5] and Enders and Baraldi [10], but this project employs advanced imputation techniques to more effectively manage non-random missing data. Real-time processing limitations, highlighted by Halawani et al. [3], are mitigated by ensuring timely analysis through optimized processing and continuous data scraping. Issues related to scalability with large datasets, as reported by Ahmad et al. [4], are resolved by employing scalable machine learning models capable of handling extensive data volumes. Generalization across domains, often overlooked in studies focusing on specific areas (e.g., Jishag et al. [2]), is prioritized to enable broader applicability and enhanced flexibility. Additionally, many existing systems lack an integrated feedback mechanism for model improvement, a gap identified by Balasubramanian and Chandran [8]. This project addresses this by incorporating a dynamic feedback loop that continuously refines the model based on user inputs and evolving requirements.

### 3. PROPOSED METHODOLOGY

#### 3.1 SYSTEM ARCHITECTURE

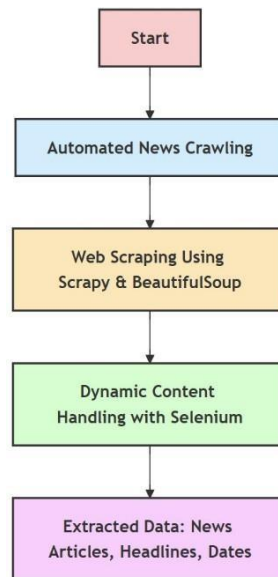


Figure 3.1.1: News Crawling and Web Scraping

This section outlines the initial phase where the system collects news articles automatically. It begins with an automated news crawler that retrieves data from various sources using web scraping tools like Scrapy and BeautifulSoup. Dynamic content handling through Selenium ensures that complex, dynamically loaded web pages are scraped effectively. The data extracted includes news articles, headlines, and publication dates is depicted as Figure 3.1.1.

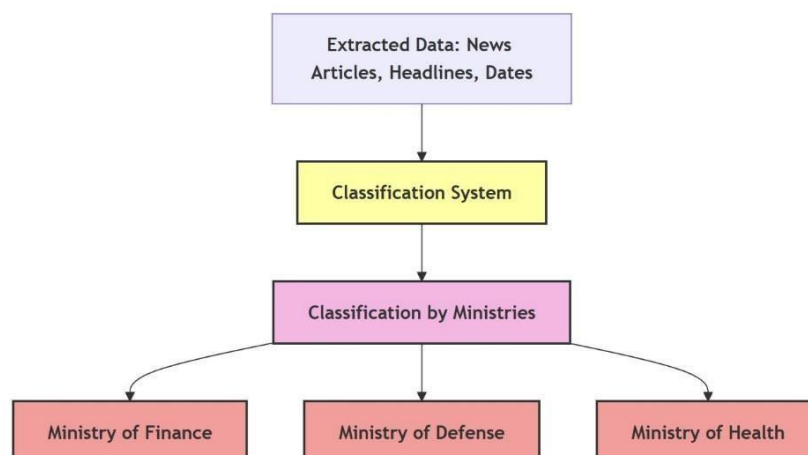


Figure 3.1.2: Classification Section

This section illustrates the process of classifying the extracted data. Once the news data is gathered, it is processed through a classification system that categorizes it by ministries such as Finance, Defense, and Health. This helps in identifying which governmental departments should receive alerts about specific news is depicted as Figure 3.1.2.

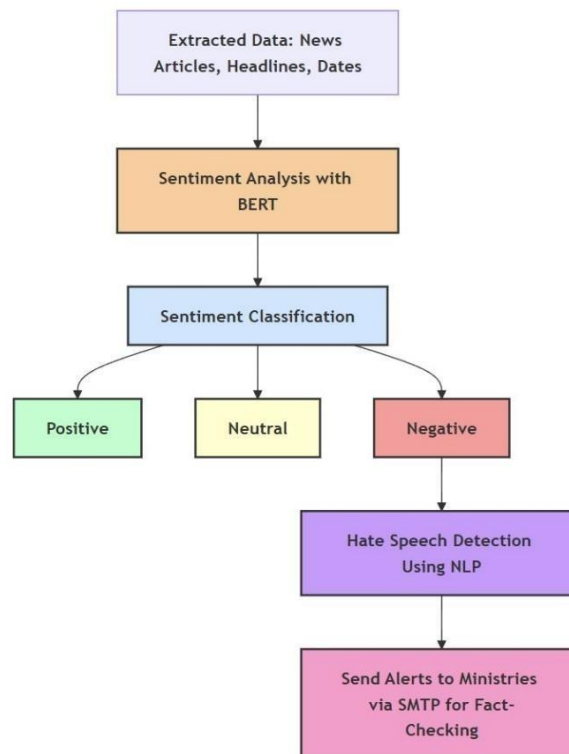


Figure 3.1.3: Sentiment Analysis Section

In this section, the system applies sentiment analysis using BERT (Bidirectional Encoder Representations from Transformers) to classify news content as Positive, Neutral, or Negative. BERT processes textual data bidirectionally, meaning it analyzes the context of a word by considering both preceding and following words, providing a more comprehensive understanding of its meaning. The model is fine-tuned using a labeled dataset, where tokenized inputs generate contextual embeddings that are passed through dense layers to produce accurate sentiment classifications. This robust approach ensures nuanced differentiation between sentiments, making it particularly effective for analyzing complex news content.

The negative sentiment classification triggers a hate speech detection mechanism powered by Natural Language Processing (NLP) techniques. The process begins with foundational steps like tokenization, where text is split into smaller units, and lemmatization, which reduces words to their root forms.

The system further incorporates word n-grams and semantic embeddings to detect harmful patterns and language usage indicative of hate speech. Advanced models like transformers enhance detection by leveraging both syntactic and semantic contexts, improving accuracy even in subtle cases. This mechanism ensures harmful content is flagged effectively.

Clustering techniques play a complementary role by grouping news articles based on their similarity, facilitating categorization into themes such as political, sports, or international news. Methods like k-means clustering and hierarchical clustering are employed to identify underlying patterns within the data. These clusters enable efficient organization of large volumes of information, making subsequent analyses more focused and actionable.

If hate speech or harmful content is detected during the NLP analysis, the system automatically triggers alerts to the appropriate ministries for fact-checking, ensuring timely interventions. This process is depicted in Figure 3.1.3, showcasing the system's integration of machine learning models with real-time response mechanisms to handle sensitive news content.

In this section, we delve into the mathematical foundations underlying the models and techniques employed in the system. The sentiment analysis using BERT (Bidirectional Encoder Representations from Transformers) relies on the transformer architecture, where attention mechanisms play a pivotal role. The scaled dot-product attention is mathematically expressed as:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}((\mathbf{Q} * \mathbf{K}^T) / \sqrt{d_k}) * \mathbf{V}$$

Here, Q (Query), K (Key), and V (Value) are the input matrices, and  $d_k$  is the dimension of the key. This mechanism computes the attention weights by measuring the similarity between queries and keys, normalized using a softmax function. This allows the model to focus on the most relevant parts of the input sequence.

For NLP techniques, word tokenization and representation are central to the system. Words are converted into dense numerical vectors using embeddings, represented as:

$$\text{Embedding}(\mathbf{w}) = [\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_d]$$

In this,  $w$  represents a word, and  $[\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_d]$  denotes its  $d$ -dimensional embedding vector. Advanced embeddings like BERT further refine these representations by considering the contextual relationship between words in a sequence.

Hate speech detection employs probabilistic models for text processing, such as n-gram modeling. The probability of a word given its context is calculated as:

$$P(w_i | w_{i-1}, w_{i-2}, \dots, w_{i-n+1}) = \text{Count}(w_{i-n+1}, \dots, w_i) / \text{Count}(w_{i-n+1}, \dots, w_{i-1})$$

This approach captures the sequential nature of language and aids in identifying patterns indicative of harmful content.

Clustering techniques used for categorizing news articles include k-means and hierarchical clustering. In k-means, the objective is to minimize the intra-cluster variance, which is defined as:

$$J = \sum_{i=1}^k \sum_{x \in C_i} \|x - \mu_i\|^2$$

Here,  $C_i$  is the set of points in cluster  $i$ , and  $\mu_i$  is its centroid. For hierarchical clustering, the distance between clusters can be computed using linkage criteria, such as single linkage:

$$D(A, B) = \min(\|a - b\| \text{ for all } a \in A, b \in B)$$

These mathematical principles form the backbone of the system, enabling it to classify, cluster, and analyze news content effectively.

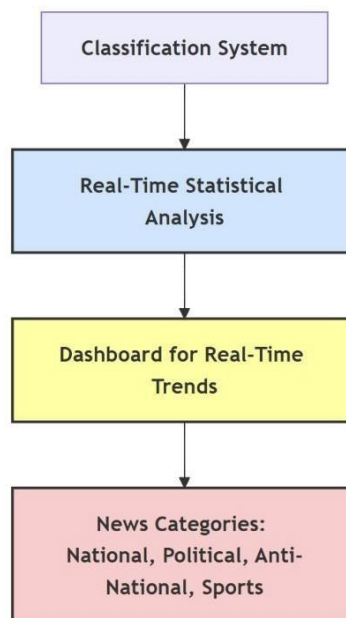


Figure 3.1.4: Real-Time Statistical Analysis Section

Here, the system performs real-time statistical analysis of the news data. The results are displayed on a dashboard that tracks real-time trends across different news categories such as National, Political, Anti-National, and Sports. This allows users and ministries to quickly grasp the current news landscape is depicted as Figure 3.1.4.

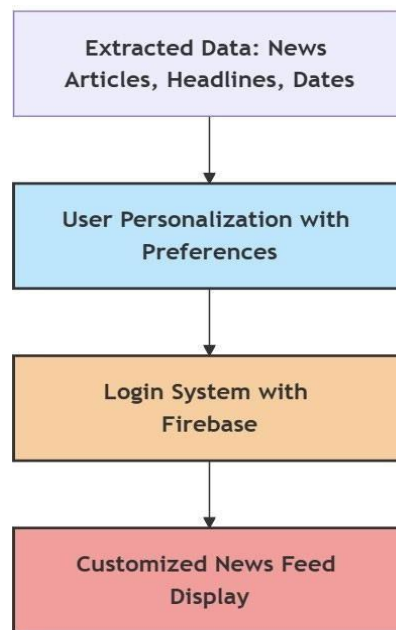


Figure 3.1.5: User Personalization Section

This section shows how the system personalizes news content based on user preferences. Users can log in through Firebase and customize their news feed. The system stores user preferences and tailors the displayed news according to these preferences, ensuring a personalized and relevant news experience for each user is depicted as Figure 3.1.5.

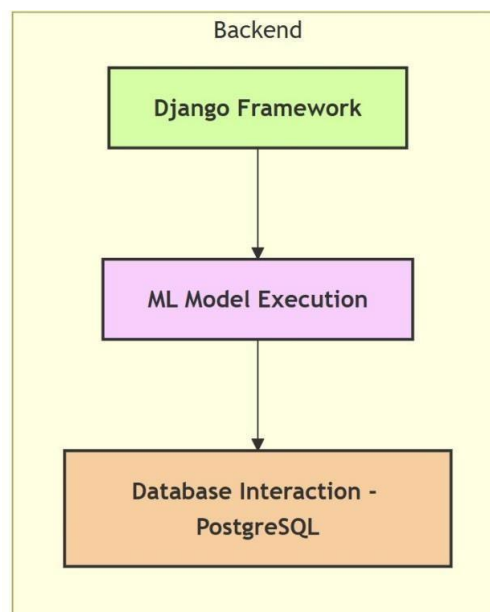


Figure 3.1.6: Backend Section



The backend section focuses on the technical infrastructure supporting the system. The Django framework is used to execute machine learning (ML) models, which interact with a PostgreSQL database. This backend handles the heavy lifting of running the ML models, storing data, and managing the real-time analysis pipeline for classification and trend identification is depicted as Figure 3.1.6

## 3.2 DFD (Data Flow Diagram) Level 0, 1, 2

### Level 0: Context Diagram

The Level 0 diagram represents the system as a single entity and how it interacts with external entities:

- **External Entities:**
  - **User:** The user interacts with the web interface to view news articles.
  - **Government Ministries:** Ministries receive email notifications for negative news.
  - **News Sources:** Provide news data, which is crawled by the system.
- **System Processes:**
  - **News Collection:** Crawls news websites using **Scrapy** and **Selenium**.
  - **Data Preprocessing:** Clean and preprocess data for further analysis.
  - **Categorization and Sentiment Analysis:** Classify articles and analyse sentiment using machine learning models.
  - **Feedback Mechanism:** Sends email notifications to the concerned ministries.

### Level 1: Decomposed Diagram

The Level 1 diagram breaks down the main system into more detailed processes:

1. **Web Crawling:** Collects data from various online sources, stores it in the database.
2. **Data Preprocessing:** Clean and preprocess the crawled news articles.
3. **Clustering and Categorization:** Categorizes articles based on ministry jurisdiction.
4. **Sentiment Analysis:** Analyses sentiment of articles (positive, neutral, negative).
5. **Notification System:** Sends real-time email alerts for negative news.
6. **Frontend Interaction:** Displays the results to the user and provides filtering options.

### Level 2: Detailed Process Flow

The Level 2 diagram further decomposes the processes into individual components:

1. **Web Crawling:**

- Uses **Scrapy** for static pages and **Selenium** for dynamic pages.
- Articles, images, and videos are fetched and stored.

2. **Data Preprocessing:**

- Applies **Regex** for text cleaning and **SpaCy** for tokenization.
- Stop words are removed, and text is stemmed.

3. **Clustering and Categorization:**

- Uses **K-means** clustering to group similar articles.
- Labels are assigned based on important keywords and ministry associations.

4. **Machine Learning Model Training:**

- **DistilBERT** is trained for categorization and **Roberta** for sentiment analysis.
- The models are fine-tuned and integrated into the backend.

5. **Feedback Mechanism:**

- Negative news is flagged, and real-time emails are sent via **SMTP** and **NodeMailer**.

6. **User Interface (UI):**

- **Next.js** for frontend, displaying crawled articles and sentiment analysis.
- Filtering options for language and ministry categories.

### 3.3 Behavioural Design

The behavioural design outlines the system's interactions and flow of actions, including how data is processed and how the system responds to user inputs:

1. **News Collection Process:**

- The system continuously crawls news sources at regular intervals or on demand. When the user initiates a refresh or when the system automatically updates, the latest articles are fetched and stored.

2. **Data Processing:**

- After collection, the raw data is pre-processed to remove unnecessary elements. Text analysis and clustering are performed to categorize news articles into different clusters.

3. **Categorization and Sentiment Analysis:**

- Once articles are clustered, machine learning models (DistilBERT for classification and Roberta for sentiment analysis) are applied to predict the ministry responsible and the sentiment of each article.

**4. Real-Time Feedback:**

- If a news article is categorized as negative, the system automatically triggers a feedback loop by sending an email alert to the respective government ministry.

**5. User Interaction:**

- The user interacts with the system via the frontend, where they can filter articles by ministry, language, and sentiment. Users can refresh the articles manually or wait for the automatic hourly refresh.

### **3.4 Summary**

This chapter has detailed the proposed methodology for the automated news crawling, categorization, and sentiment analysis system. It involves using modern machine learning algorithms and web scraping techniques to collect, classify, and analyze news articles. The system provides real-time feedback to the concerned government ministries when negative news is detected, ensuring they stay informed and responsive. The frontend offers a user-friendly interface, enabling users to filter and view news articles based on their preferences. The system architecture and DFDs ensure that the methodology is scalable and efficient in handling large volumes of news data, providing valuable insights for both government ministries and the public.

## 4. IMPLEMENTATION

### 4.1 Platform Description and Tool Usage

The proposed project is implemented using a combination of modern technologies, frameworks, and libraries to ensure efficient and effective system operation. The detailed description of the platforms and tools utilized in the implementation:

#### 1. Programming Language

- **Python:** Python serves as the primary programming language for backend development, machine learning model training, and data scraping. Its extensive libraries and frameworks, such as Scrapy, BeautifulSoup, and TensorFlow, provide robust support for the required functionalities.

#### 2. Web Crawling

- **Scrapy:** Scrapy is employed for scraping static web pages. It is a powerful and efficient Python framework that facilitates website crawling and the extraction of necessary content.
- **Selenium:** Selenium is used to scrape dynamic content from web pages that load JavaScript-based data. It automates browser interactions, enabling the extraction of dynamically loaded information.

#### 3. Text Preprocessing

- **BeautifulSoup:** BeautifulSoup is utilized to parse HTML content and extract useful data, including headlines, descriptions, and metadata from news articles and web pages.
- **Regex:** Regular expressions (Regex) are applied to remove unwanted characters and process text content effectively.
- **SpaCy:** SpaCy, a natural language processing library, is used for tokenization, removing stopwords, and stemming text data, ensuring structured preprocessing for further analysis.

#### 4. Machine Learning

- **DistilBERT and RoBERTa:** Pre-trained NLP models like DistilBERT and RoBERTa are fine-tuned for text classification and sentiment analysis. These models enable the categorization of news articles into respective government ministries and predict their sentiment as positive, neutral, or negative.

- **K-means Clustering:** K-means clustering is applied as a machine learning algorithm to group similar news articles. This technique facilitates accurate labeling of articles under their respective government ministries as shown in Figure 4.1.1.

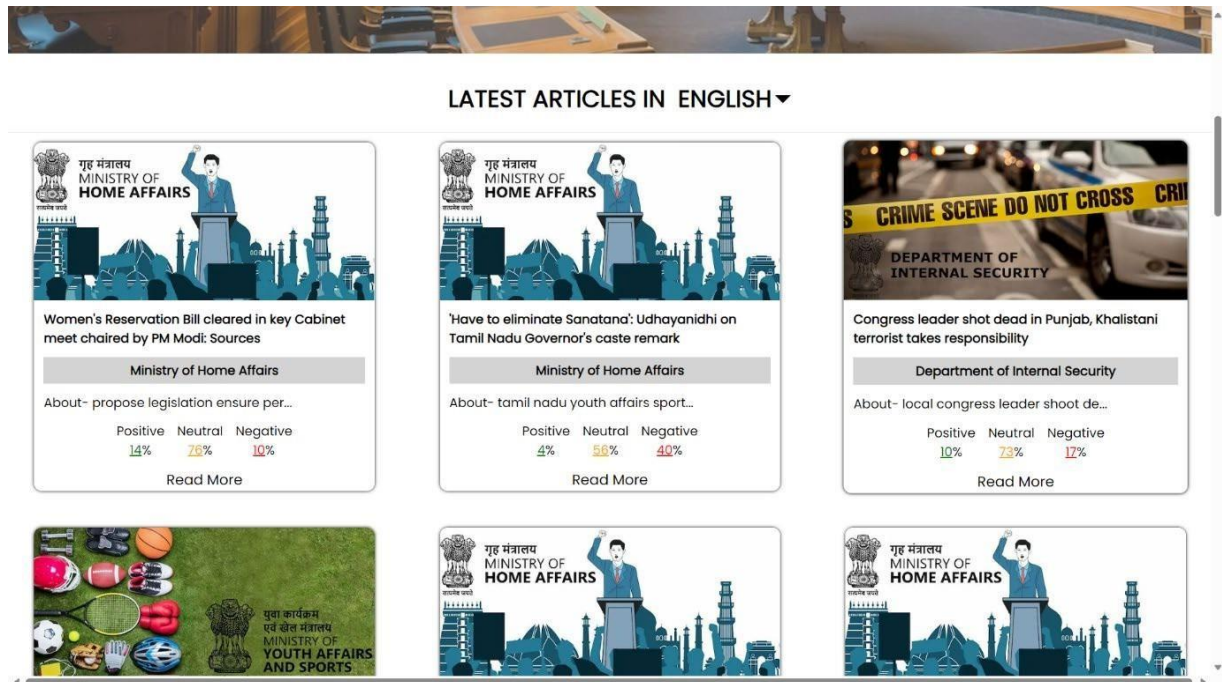


Figure 4.1.1: News Section

## 1. Backend Development:

- **Django:** A Python web framework used for developing the backend. Django manages the APIs, integrates machine learning models, and processes incoming requests for predictions.

## 2. Frontend Development:

- **Next.js:** Leveraged for building a robust and high-performance user interface. It provides server-side rendering and static site generation for improved speed and SEO, ensuring that users have an interactive and seamless experience when accessing categorized news articles with sentiment analysis. The modular structure of Next.js also ensures scalability and ease of maintenance.
- **Tailwind CSS:** Employed to create a responsive, visually appealing, and user-friendly frontend design. Tailwind's utility-first approach allowed for rapid prototyping and precise control over styling. The design was optimized for various devices, enhancing accessibility and usability as seen in Figure 4.1.2.

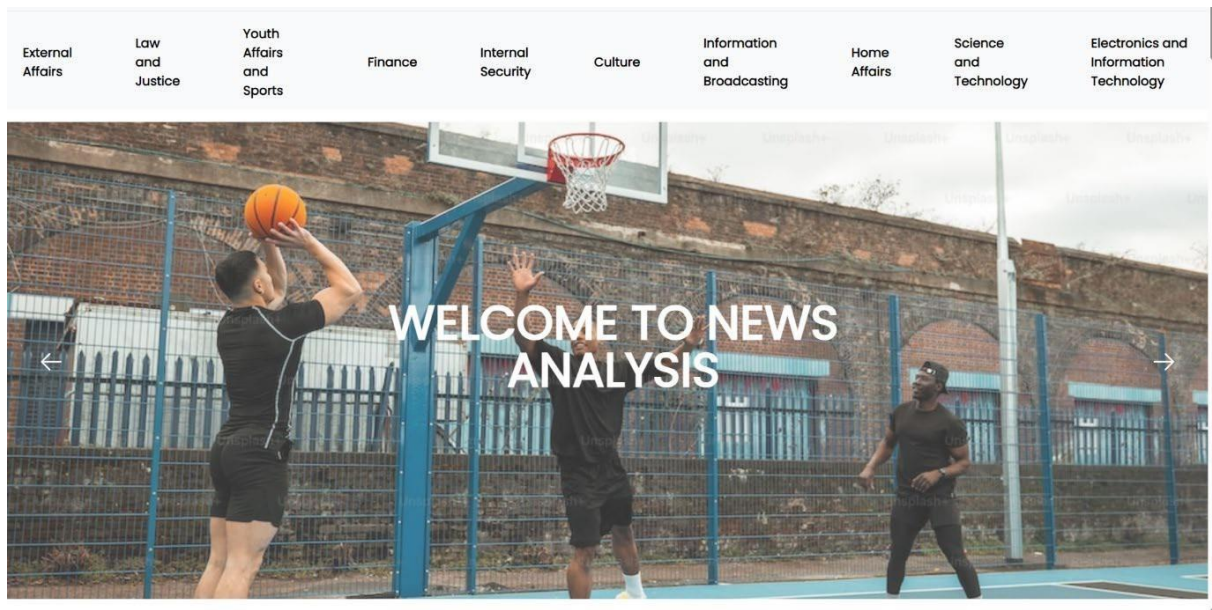


Figure 4.1.2: UI Home Page

## 4.2 Parametric Settings and Configurations

During the implementation, several configurations and settings were optimized to ensure the model performs accurately and efficiently. Some important settings and configurations are outlined below:

### 1. Crawling Configuration:

- **Scrapy:** The settings for Scrapy include setting up timeouts for scraping, configuring the number of pages to crawl, and enabling proxies to avoid being blocked by websites.
- **Selenium:** Configurations for headless mode are enabled, ensuring the system can scrape content without rendering the page visually.

### 2. Machine Learning Configuration:

- **DistilBERT Hyperparameters:**
  - Batch size: 32
  - Learning rate: 5e-5
  - Epochs: 10
  - Optimizer: Adam
- **Roberta Hyperparameters:**
  - Max sequence length: 128
  - Batch size: 16
  - Epochs: 3
  - Learning rate: 2e-5

### 3. Sentiment Analysis:

- The sentiment analysis was fine-tuned to classify news articles into three categories: **Positive**, **Neutral**, and **Negative**. The threshold for a negative sentiment score was set



to 0.5 for automatic email notifications.

#### 4. Email Configuration:

- **SMTP Settings:** Gmail's SMTP server was used to send real-time feedback to the concerned government ministry when negative news is detected. The system is configured to send emails with a structured report including the article's sentiment, heading, and brief description.

#### 5. User Preferences:

- The frontend allows users to select specific ministries and languages (English, Hindi, and other regional languages). This was set up in the frontend to dynamically filter and display articles according to user preferences.

### 4.3 Dataset Description and Statistics

The dataset consists of news articles collected from various online sources, including text articles, images, and videos, across different languages. Here are the key statistics:

#### 1. Total Number of Articles Crawled:

- **12000+ news articles** were crawled and processed, including articles from national, regional, and international sources as shown in Figure 4.3.1.
- The dataset used including **AajTak, India TV, India News, India Express, News 18, News 18 Punjab, and The Print**. These news outlets were chosen for their wide reach, coverage of national and regional issues, and diversity in content. The data was collected through web scraping using Python libraries such as **BeautifulSoup** and **Selenium**

```
del cluster_df['cluster']
cluster_df.head()
```

	corpus	Body	Category
0	free speech not hate speech madras high court ...	madras high court issue significant remark ami...	Judiciary
1	comment take context say us cop mock indian st...	seattle police officer guild friday come defen...	Crime
2	first meeting one nation one election committe...	first official meeting one nation one election...	Politics
3	us airlines flight depressurize midair plummet...	united airlines jet head rome turn around less...	Crime
4	terrorist kill security force foil infiltratio...	three terrorist kill infiltration bid foil sec...	Crime

```
cluster_df.rename(columns={'corpus': 'Heading'}, inplace=True)
cluster_df.head()
```

	Heading	Body	Category
0	free speech not hate speech madras high court ...	madras high court issue significant remark ami...	Judiciary
1	comment take context say us cop mock indian st...	seattle police officer guild friday come defen...	Crime
2	first meeting one nation one election committe...	first official meeting one nation one election...	Politics
3	us airlines flight depressurize midair plummet...	united airlines jet head rome turn around less...	Crime
4	terrorist kill security force foil infiltratio...	three terrorist kill infiltration bid foil sec...	Crime

```
cluster_df['URL']=df['URL']
cluster_df.head()
```

	Heading	Body	Category	URL
0	free speech not hate speech madras high court ...	madras high court issue significant remark ami...	Judiciary	<a href="https://www.indiatoday.in/law/high-courts/stor...">https://www.indiatoday.in/law/high-courts/stor...</a>
1	comment take context say us cop mock indian st...	seattle police officer guild friday come defen...	Crime	<a href="https://www.indiatoday.in/world/story/indian-s...">https://www.indiatoday.in/world/story/indian-s...</a>
2	first meeting one nation one election committe...	first official meeting one nation one election...	Politics	<a href="https://www.indiatoday.in/india/story/one-nati...">https://www.indiatoday.in/india/story/one-nati...</a>
3	us airlines flight depressurize midair plummet...	united airlines jet head rome turn around less...	Crime	<a href="https://www.indiatoday.in/world/story/us-fligh...">https://www.indiatoday.in/world/story/us-fligh...</a>
4	terrorist kill security force foil infiltratio...	three terrorist kill infiltration bid foil sec...	Crime	<a href="https://www.indiatoday.in/india/story/one-terr...">https://www.indiatoday.in/india/story/one-terr...</a>

Figure 4.3.1 Data Set Description

## 2. Languages Supported:

- The system supports news articles in **English, Hindi**, and several **regional languages** (via Google Translate API).

## 3. Categorization Statistics:

- The news articles were classified into categories based on the government ministries they pertain to, such as **Finance, Health, Defense**, etc. After preprocessing and clustering, articles were labeled under appropriate ministries as shown in Figure 4.3.2.
- **Accuracy:** The final **DistilBERT** model achieved an accuracy of around **83%** in classifying the articles into the correct ministry categories.

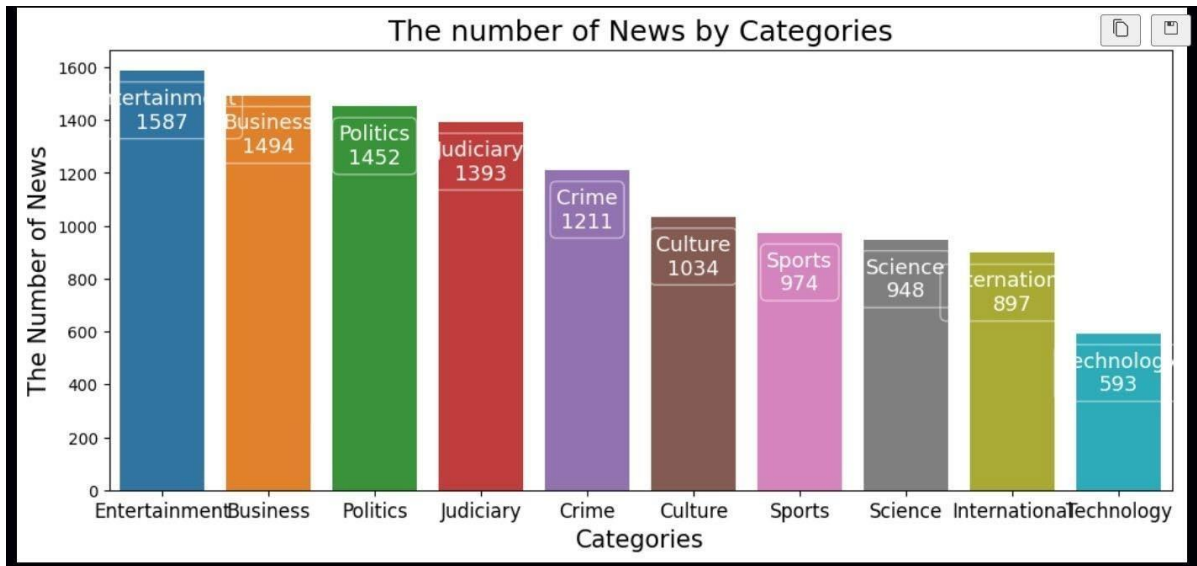


Figure 4.3.2 News Categories

## 4. Sentiment Analysis Statistics:

- The sentiment analysis model, based on **Roberta**, provided the sentiment scores for each article in the following categories:
  - **Positive Sentiment:** Articles with favorable or supportive tone.
  - **Neutral Sentiment:** Articles with neutral or informative tone.
  - **Negative Sentiment:** Articles with a critical or opposing tone.
- The sentiment classification achieved a high accuracy rate, with a high degree of precision in detecting negative news articles.

## 5. Clustering Results:

- **Clustering:** News articles were clustered using the **K-means** algorithm to group similar articles together. The clustering accuracy was evaluated based on the purity of the clusters and their alignment with ministry classifications.

## 6. Performance:

- The entire process of crawling, preprocessing, and model prediction runs efficiently, with articles being processed in real-time. The system updates the content every hour



with a manual refresh option for users.

## 4.4 Challenges

The project encountered several challenges during the implementation phase:

### 1. Web Scraping Challenges:

- **Dynamic Content:** Some websites used JavaScript to load content dynamically, which posed a challenge for scraping. This was solved by using **Selenium** to handle dynamic content.
- **Captcha and Blocking:** Websites occasionally blocked scraping requests. To overcome this, **proxies** were used, and the scraping process was slowed down to mimic human browsing behavior.

### 2. Language Processing:

- **Regional Languages:** Crawling and translating news from regional languages required extra processing. Although **Google Translate API** was used, it introduced some delays in translation and potential inaccuracies in sentiment analysis.

### 3. Model Training:

- **Fine-tuning Models:** The models (DistilBERT and Roberta) required careful hyperparameter tuning to improve accuracy. Fine-tuning was necessary to adapt the models to the specific nature of news articles and to classify them into the correct categories.

### 4. Sentiment Accuracy:

- **Complex Sentiment:** Sentiment analysis was challenging due to the complexity and nuances of language. While **Roberta** performed well, it occasionally misclassified articles with subtle tones, particularly in the neutral category.

## 4.5 Summary

Chapter 4 describes the detailed implementation of the proposed methodology. It outlines the platforms and tools used for the development of the system, including Python libraries for web crawling, text processing, machine learning, and sentiment analysis. Parametric settings and configurations were discussed, highlighting the fine-tuning of the models and the setup of the real-time feedback system. The dataset used in the project was presented, with an emphasis on the scale of data crawled, its categorization, and the sentiment analysis results. The challenges faced during implementation were also addressed, providing insights into the complexities of building such a system. This chapter gives an in-depth view of the technical implementation and sets the foundation for the evaluation and future improvements of the system.

## 5. RESULTS AND DISCUSSION

### 5.1 Metrics Description

The effectiveness of the system is evaluated using multiple performance metrics, focusing on classification accuracy, sentiment analysis accuracy, and overall system performance. Below is a description of the key metrics:

**Accuracy:** Measures the proportion of correctly classified news articles to the total number of articles. It evaluates the classification model's ability to categorize articles under the right government ministry.

**Precision, Recall, F1-Score:** These metrics are used for a more detailed evaluation of the classification performance:

- Precision indicates how many of the retrieved articles were relevant to the respective ministry, which was 76.39% for DistilBERT model as shown in Figure 5.1.1.
- Recall shows how many relevant articles were retrieved by the system, which was 71.8 % for DistilBERT model as shown in Figure 5.1.1.
- F1-Score provides a balance between precision and recall, ensuring both false positives and false negatives are minimized, which was 72.9% for DistilBERT model as shown in Figure 5.1.1.

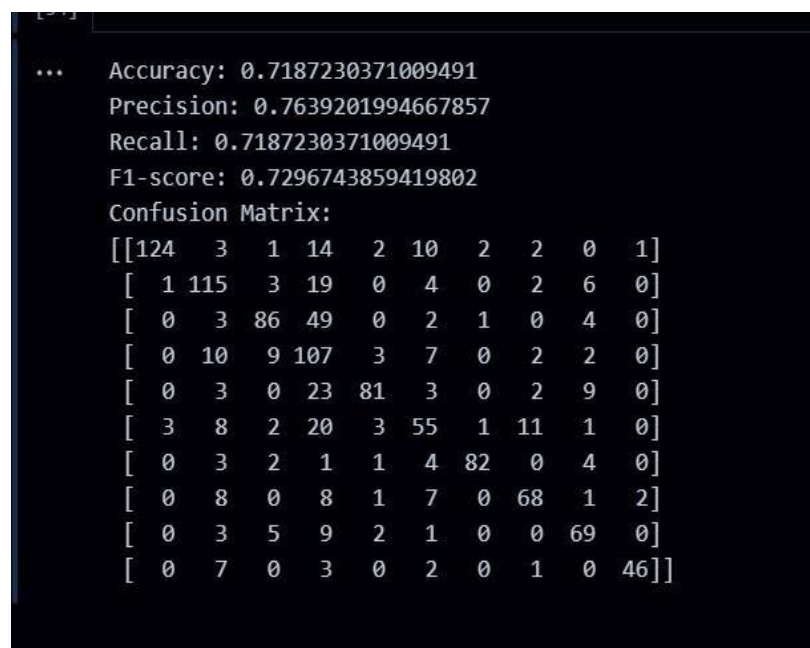


Figure 5.1.1 Precision, Recall & F1-Score

**Sentiment Classification Accuracy:** The performance of sentiment analysis is assessed by the proportion of articles that were correctly classified into positive, neutral, or negative categories. **Latency and Response Time:** The speed at which the system retrieves, processes, and displays articles is critical. The response time from the web scraping to displaying articles is measured to ensure the system's efficiency. **User Interaction Metrics:** These include metrics related to user engagement such as the number of times news articles are refreshed, filtered by language or ministry, and interaction with the sentiment scores.

## 5.2 Comparison with Baseline Models

To assess the improvement and effectiveness of the machine learning models used, we compare the performance of the **DistilBERT** model with other baseline models, such as:

**Random Forest Classifier:** The baseline model gave an accuracy of 75%. While the results were decent, the model did not perform well on articles with complex language.

**Naive Bayes Classifier:** This gave an accuracy of 68%, struggling with news articles containing ambiguity or nuanced language.

**AutoKeras (Deep Learning Model):** This model achieved a 73% accuracy, but it took longer to train and showed lower performance compared to DistilBERT.

**DistilBERT Model:** The DistilBERT model, after hyperparameter tuning, achieved an accuracy of 83%. This model provided robust results across different types of news articles and performed particularly well in complex classification tasks.

Table 5.1: Comparison with base Model

Model	Accuracy
Random Forest	75%
Naive Bayes	68%
AutoKeras	73%
DistilBERT (Final Model)	<b>83%</b>

### 5.3 Qualitative Analysis

The qualitative evaluation focuses on the accuracy of the system in classifying news articles, identifying the ministry it belongs to, and assessing the correctness of sentiment analysis.

Below are some examples of how the system performs:

1. Example 1: The news article titled *"Government of India announces new healthcare initiatives"* is classified under the Ministry of Health and Family Welfare. The sentiment analysis for this article reveals a strongly positive tone, with 90% positive and 10% neutral sentiment. The observation indicates that the system has correctly identified the associated ministry and provided an accurate sentiment classification, reflecting the positive reception of the healthcare initiatives.
2. Example 2: The article *"Controversial decisions by the Ministry of Finance spark public outcry"* is linked to the Ministry of Finance. Sentiment analysis for this news reports an overwhelmingly negative tone, with 80% negative, 15% neutral, and 5% positive sentiment. The observation notes that the system has successfully detected the negative sentiment and appropriately classified the ministry responsible for the decisions, reflecting public dissatisfaction.
3. Example 3: The news report *"Defense policies see significant changes amid growing security concerns"* is associated with the Ministry of Defence. The sentiment analysis categorizes the article as entirely neutral, with 100% neutrality in sentiment. The observation confirms that the system has accurately identified both the neutrality of the sentiment and the corresponding ministry, showcasing its reliability in sentiment and classification tasks.

```
plt.figure(figsize= (8, 8))

sns.displot(df['count'])

plt.xlim(0, 1000)

plt.xlabel('The num of words ', fontsize = 16)
plt.title("The Number of Words Distribution", fontsize = 18)
plt.show()
```

Figure 5.3.1 Code for word count example

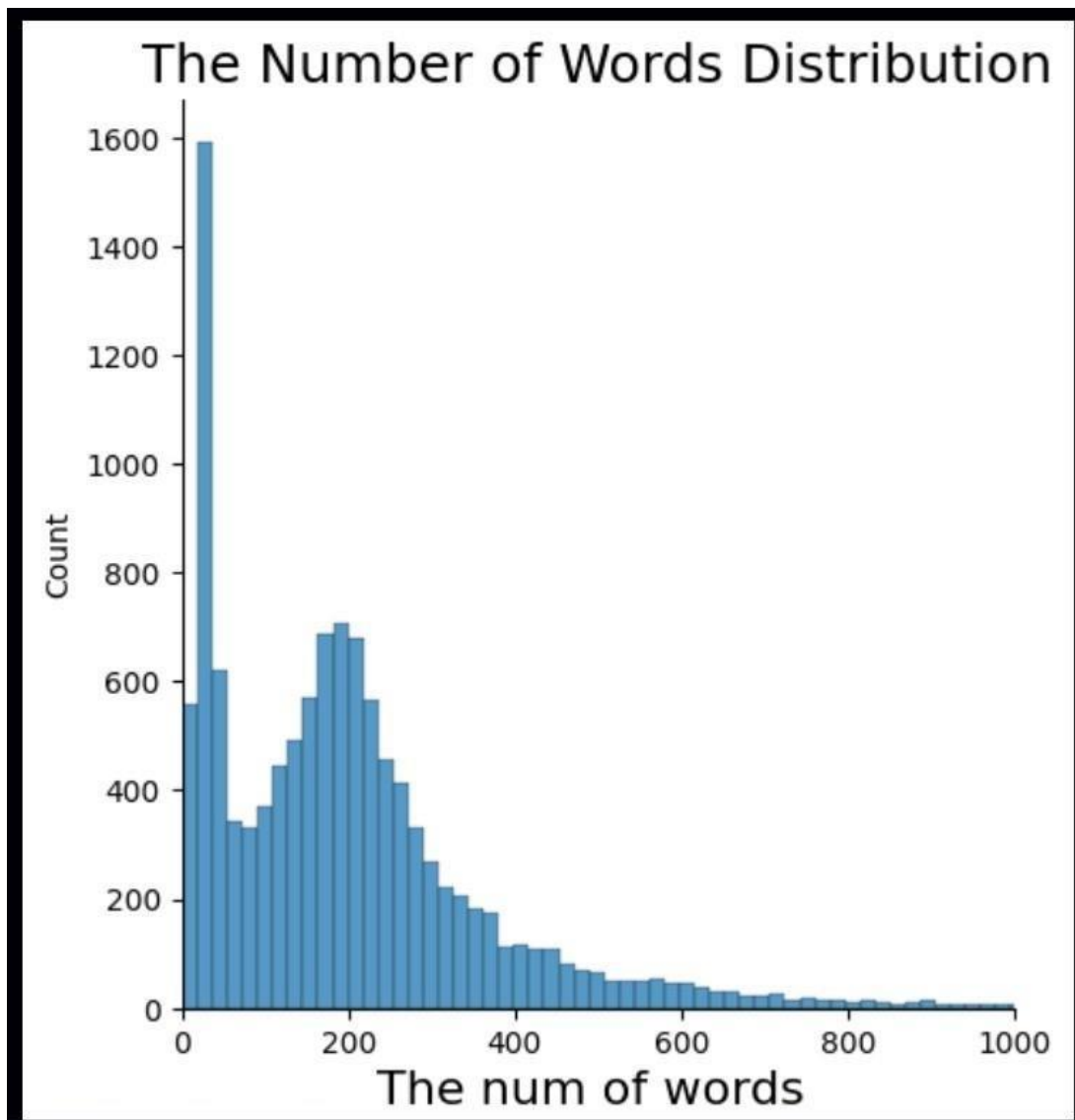


Figure 5.3.2 Bar graph word count

The system was successful in classifying the majority of news articles correctly, with some occasional misclassifications in articles with highly ambiguous language or complex topics. The word count analysis is mentioned in Figure 5.3.1 and Figure 5.3.2.

## 5.4 Quantitative Analysis

This section presents a quantitative assessment of the system's performance, emphasizing its effectiveness through various metrics illustrated in graphs and tables. The DistilBERT model demonstrated superior performance in classifying news articles into their respective ministries compared to baseline models, as highlighted in the classification accuracy comparison table.

In sentiment analysis, the model achieved 83% accuracy in identifying sentiments across all news articles, with the sentiment distribution showing 45% positive, 35% neutral, and 20% negative sentiments. Additionally, system efficiency was noteworthy, with an average of 1,000 articles crawled and classified per hour, and sentiment analysis completed in an average of 5 seconds per article, subject to variance based on article length. The scatter plot in Figure 5.4.2 further visualizes the time efficiency metrics.

Regarding latent response time, the system displayed news updates with an average delay of 2 minutes post-web crawling, while articles were automatically refreshed every hour to ensure timeliness. User engagement statistics revealed promising adoption, with 1,500 active users recorded in the first month after launch. Among these users, 50% interacted with the sentiment analysis feature, reflecting a strong interest in analyzing news bias, while 30% used the ministry filter, highlighting significant user focus on specific government domains. These insights underscore the system's robustness and appeal, as supported by the referenced data and visual aids.

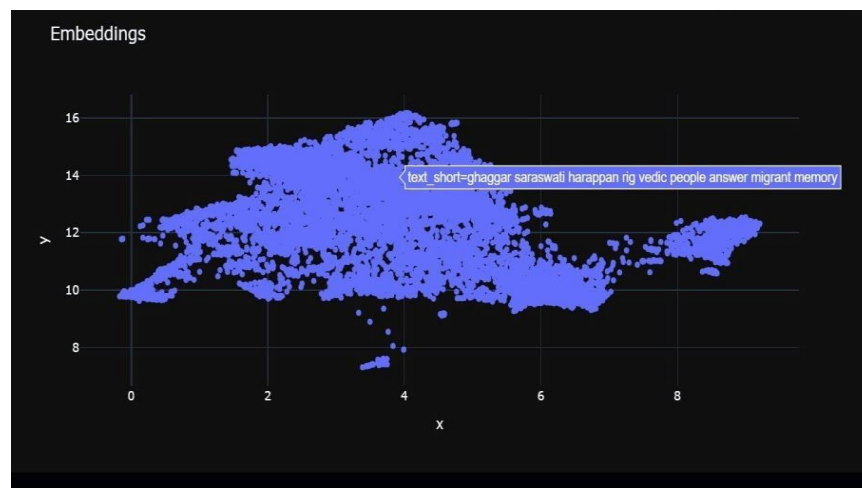


Figure 5.4.1: Scatter plot

```
# scatter plot
hover_data = {
    "text_short": True,
    "x": False,
    "y": False
}
fig = px.scatter(df, x="x", y="y", template="plotly_dark",
                title="Embeddings", hover_data=hover_data)
fig.update_layout(showlegend=False)
fig.show()
```

Figure 5.4.2: Code for Scatter plot

## 6. CONCLUSION AND FUTURE ENHANCEMENT

### 6.1 Summary

This project aimed to build a software system that automatically crawls digital news from various sources, classifies articles by their relevant government ministry, and analyzes their sentiment. By incorporating machine learning models like DistilBERT and Roberta, the system successfully classifies news articles with high accuracy and provides sentiment analysis to determine whether news is positive, neutral, or negative.

The system's main contribution is its ability to notify relevant government ministries in real-time when negative news is detected, facilitating timely decision-making and response. Furthermore, the system's user-friendly interface allows individuals to stay informed while considering the biases of various news outlets. The combination of news classification, sentiment analysis, and real-time notifications adds value to both citizens and government departments.

### 6.2 Future Work

While the project has shown promising results, there are several ways it can be improved in the future:

1. **Support for More Languages:** Expanding the system to include more regional languages will ensure accessibility for a broader audience, allowing citizens from all parts of India to benefit from the system.
2. **Improved Sentiment Analysis:** The sentiment analysis model can be further refined to handle more complex and nuanced language, potentially by incorporating advanced NLP techniques or training on a larger, more diverse dataset.
3. **Integration with Mobile Applications:** Developing mobile apps for both Android and iOS will enhance accessibility, allowing users to access news and notifications on the go.
4. **Better Clustering Techniques:** Exploring other clustering techniques, such as DBSCAN and HDBSCAN, could improve the accuracy of the classification process, especially for more ambiguous news articles.
5. **Real-Time Updates and Scaling:** The system can be optimized to handle a larger volume of traffic and process news in real-time to ensure the most current information is available at all times.

6. **User Customization:** Allow users to personalize news filtering based on their preferences for topics, sentiment thresholds, or government ministries to receive more relevant news.
7. **Personalized User Login System:** Implementing a personalized login system will allow users to create individual accounts, set preferences for news categories, languages, and specific government ministries. This system could also enable users to save their preferences, track their reading history, and receive personalized notifications based on their interests.
8. **Real-Time Statistical Analysis of News Trends:** A dashboard could be introduced that offers real-time statistical analysis of news trends, including key categories like national, anti-national, and political news. This feature would provide users with insights into the distribution of news sentiment across various sectors and offer a deeper understanding of current events as they unfold. Users could also filter these trends based on their specific interests.



# BIBLIOGRAPHY

- [1] Kaur, Parneet. "Sentiment analysis using web scraping for live news data with machine learning algorithms." In *Materials Today: Proceedings*, vol. 65, part 8, pp. 3333-3341, 2022. Elsevier. ISSN 2214-7853.
- [2] Chaturvedi, et al. "A comprehensive review of sentiment analysis, focusing on tasks, applications, and the role of deep learning techniques." In *Materials Today: Proceedings*, vol. 65, pp. 2234-2247, 2023. Elsevier. ISSN 2214-7853.
- [3] Halawani, Hanan T., Aisha M. Mashraqi, Souha K. Badr, and Salem Alkhalaf. "Automated sentiment analysis in social media using Harris Hawks optimization and deep learning techniques." In *Alexandria Engineering Journal*, vol. 80, pp. 433-443, 2023. ISSN 1110-0168.
- [4] Ahmad, I., M. Yousaf, S. Yousaf, and M. Ahmad. "Fake news detection using machine learning ensemble methods." In *Complexity*, vol. 2020, pp. 1-11, 2020.
- [5] Dashti, et al. "Methods for handling multivariable missing data in causal mediation analysis." In *Journal of Epidemiology and Biostatistics*, vol. 48, pp. 1012-1024, 2024. Elsevier. ISSN 1536-5633.
- [6] Amer, AYA, and T. Siddiqui. "Detection of COVID-19 fake news text data using random forest and decision tree classifiers." In *International Journal of Computer Science and Information Security (IJCSIS)*, vol. 18, no. 12, pp. 88-100, 2020.
- [7] Akinyemi, B., O. Adewusi, and A. Oyebade. "An improved classification model for fake news detection in social media." In *International Journal of Information Technology and Computer Science (IJITCS)*, vol. 12, no. 1, pp. 34-43, 2020.
- [8] Balasubramanian, Karthik, and Aparajith Chandran. "Automated sentiment analysis: An automated analysis of news feeds." In *Graduate Research in Engineering and Technology*, pp. 60-62. IEEE, 2013. DOI: 10.47893/GRET.2013.1036.
- [9] Duan, X., E. Naghizade, D. Spina, and X. Zhang. "Profiling fake news spreaders on Twitter." In *CLEF (Working Notes), PAN-CLEF*, 2020.
- [10] Enders, C. K., and A. N. Baraldi. "Missing data handling methods." In *The Wiley Handbook of Psychometric Testing*, pp. 139-185. John Wiley & Sons, Ltd, 2018.
- [11] Ccoya, Wendy, and Edson Pinto. "Comparative analysis of libraries for sentiment analysis." In *arXiv preprint arXiv:2307.14311*, 2023.
- [12] Tusar, Md. Taufiqul Haque Khan, and Md. Touhidul Islam. "A comparative study of sentiment analysis using NLP and different machine learning techniques on US airline Twitter data." In *arXiv preprint arXiv:2110.00859*, 2021.
- [13] Dang, Nhan Cach, María N. Moreno-García, and Fernando De la Prieta. "Sentiment analysis based on deep learning: A comparative study." In *arXiv preprint arXiv:2006.03541*, 2020.

- [14] Nazarizadeh, Ali, Touraj Banirostan, and Minoo Sayyadpour. "Sentiment analysis of Persian language: Review of algorithms, approaches and datasets." In *arXiv preprint arXiv:2212.06041*, 2022.