

Analysing Media Bias in Health Reporting: Scraping News Articles for Language Patterns in Health Misinformation

Submitted by:

Kesar Tripathi (Roll No. 242807011)

M Aravind (Roll No. 242807019)

Course: M.Sc. Data Science (2024–2026)

Subject: Deep Learning and Text Mining (DDS5205)

Date of Submission: 12th May 2025

Table of Contents

Acknowledgement	4
1. Introduction	5
2. Literature Review.....	6
2.1 Traditional Approaches to Bias and Misinformation Detection	6
2.2 Emergence of Deep Learning and Transformer Models	6
2.3 Framing and Emotion in Health News	7
2.4 Misinformation Classification and Multitask Learning	7
2.5 Justification for the Current Approach.....	8
3. Methodology	8
3.1 Data Collection and Preprocessing.....	9
3.2 Topic Modelling (Unsupervised NLP).....	9
3.3 Bias Detection (Emotion Analysis)	9
3.4 Misinformation Classification (Supervised DL)	9
3.5 Visualization and Interpretation.....	10
4. Dataset Description.....	10
4.1 Acquisition Source and Creation Methodology	10
4.2 Dataset Characteristics	11
4.3 Corpus Statistics and Visualization	11
5. Experimental Setup	14
5.1 Data Collection and Preprocessing.....	14
5.2 Topic Modelling Using BERTopic	16
5.3 Bias Classification Using RoBERTa	17
5.4 Label Encoding and Dataset Transformation	17
5.5 Emotion Classification Using DistilBERT	18
5.6 Hardware Constraints and Environment	18
6. Results and Discussion.....	19
6.1. Quantitative Results	19
6.2. Baseline Model (Logistic Regression + TF-IDF)	19
6.3. DistilBERT Fine-Tuned Model.....	20
6.4. Emotion Distribution Across Bias Labels (roBERTa).....	21
6.5. Observations	23

6.6. Strengths / Improvements Over Baseline	23
6.7. Limitations	24
7. Scope for Future Work	25
8. Conclusion	27
9. References	28

Acknowledgement

We would like to express our sincere gratitude to **Dr. Kavitha Karimbi Mahesh**, Associate Professor, Department of Applied Statistics and Data Science, Prasanna School of Public Health, Manipal Academy of Higher Education, for her valuable guidance, consistent encouragement, and insightful feedback throughout the course of this project. Her expertise and mentorship played a crucial role in shaping the direction and depth of our work.

We also extend our appreciation to the **Department of Applied Statistics and Data Science** for providing the academic framework and resources necessary to carry out this project.

Finally, we would like to acknowledge the collaborative effort of both team members, whose contributions and dedication were essential to the successful completion of this study.

1. Introduction

The rise of digital media has dramatically transformed the way health information is disseminated and consumed. In a populous nation like India, news media play a critical role as a primary source of public health knowledge. However, the increasing prevalence of biased reporting and misinformation—particularly during global health crises like the COVID-19 pandemic—has raised significant concerns about the accuracy and impact of such coverage. Media stories employing **sensationalized language, deceptive framing, or unsubstantiated assertions** can mislead public attitudes and actions, potentially affecting public health on a mass scale.

This project aims to address these challenges by rigorously analysing linguistic patterns in health news articles from major Indian news organizations. By employing a comprehensive pipeline that includes **web scraping, topic modelling, emotion and framing analysis, and deep learning-based misinformation detection**, the project seeks to uncover the ways in which health information is presented—and sometimes misrepresented—in Indian media.

The main contributions of this work include:

- **A specially created dataset** of news articles focused on Indian health.
- **Topic modelling** to capture common themes in the collected articles.
- **Emotion and framing analysis** to identify linguistic bias.
- **A fine-tuned transformer-based classifier** to detect misinformation.

Through the convergence of these approaches, this project not only highlights patterns of media bias but also demonstrates the potential of deep learning in combating misinformation within the public health communication space.

2. Literature Review

2.1 Traditional Approaches to Bias and Misinformation Detection

Early approaches to detecting misinformation relied primarily on traditional supervised algorithms like **Naïve Bayes**, **Support Vector Machines (SVM)**, and **Logistic Regression**. These models used handcrafted features such as **n-grams**, **syntactic structures**, and **sentiment polarity**. While they offered interpretability and computational efficiency, these methods often struggled to generalize across diverse linguistic patterns and capture the contextual nuances inherent in health misinformation.

Sentiment analysis tools like **VADER** and **TextBlob** were commonly employed; however, they lacked the ability to detect domain-specific affect tones, such as **fearmongering and misleading reassurance**. Similarly, topic modeling methods like **Latent Dirichlet Allocation (LDA)**, though effective in identifying broad themes, lacked the coherence and semantic richness needed for complex health discourse, particularly when analyzing long articles.

2.2 Emergence of Deep Learning and Transformer Models

The advent of transformer models such as **BERT**, **RoBERTa**, and **DistilBERT** marked a significant leap in **Natural Language Processing (NLP)**. These models utilize **self-attention mechanisms and contextual embeddings** to capture deep linguistic relationships, significantly outperforming traditional models in tasks like **sentiment analysis, fake news detection, and topic modeling**.

For instance, Zhou et al. (2020) demonstrated the efficacy of **BERT** in detecting COVID-19 disinformation with high precision and recall across multiple datasets. Similarly, Fan et al. (2022) highlighted the effectiveness of emotion-detecting transformer models in identifying emotional biases in health news, emphasizing the crucial role of emotional content in spreading misinformation.

In topic modeling, **BERTopic** has emerged as a superior alternative to LDA by leveraging sentence-transformer embeddings and class-based TF-IDF (c-TF-IDF), providing semantically dense and human-interpretable topics. This approach is particularly effective for uncovering subtle themes in health-related articles.

2.3 Framing and Emotion in Health News

Media framing significantly influences public perception of health issues. According to Entman's Theory of Framing, media frames shape how audiences interpret issues. Recent studies show that health articles frequently use evocative language—such as **fear**, **anger**, and **hope**—to capture attention, often distorting reality and potentially exaggerating misinformation.

Chuang and Tsai (2021) found that emotionally charged news articles not only had higher virality but also a greater tendency to mislead, particularly on topics like vaccines and chronic illnesses. As a result, emotion classification has become a crucial tool in bias detection, with models fine-tuned on datasets like **GoEmotions** to detect tones like **fearfulness**, **sadness**, **optimism**, and **neutrality** in health journalism.

In this project, emotion classification is combined with framing analysis to identify whether articles adopt an **alarmist**, **factual**, or **personal tone**, providing a more comprehensive view of bias.

2.4 Misinformation Classification and Multitask Learning

Detecting misinformation in health media is increasingly viewed as a multifaceted challenge involving **sentiment**, **credibility**, and **topic relevance**. Fine-tuning transformer models like **DistilBERT** for domain-specific classification, coupled with emotion and framing features, has proven effective in improving model accuracy.

While prior studies have often focused on Western datasets like **FakeNewsNet** and **COVID-Lies**, this project bridges the gap by focusing on Indian media, capturing the unique linguistic and narrative styles prevalent in Indian journalism.

2.5 Justification for the Current Approach

To address the unique challenges of Indian health journalism, this project implements the following innovations:

- **A specially developed Indian health news dataset** created through extensive web scraping.
- **BERTopic** for more coherent and semantically rich topic modeling.
- **Emotion and framing detection** using advanced transformer models.
- **Domain-specific misinformation classification** leveraging deep learning.
- **Visualization techniques** to reveal language patterns and model insights.

This integrated approach enables a transparent, in-depth examination of how Indian media shape public discourse on health, both explicitly and implicitly.

3. Methodology

Our methodology involves a multi-phased pipeline combining **data collection**, **preprocessing**, **topic modelling**, **emotion and framing analysis**, and **misinformation detection**. This approach provides a comprehensive understanding of the biases and misinformation present in Indian health news articles.

3.1 Data Collection and Preprocessing

A custom web scraping pipeline was developed using Python libraries like **newspaper3k**, **requests**, and **BeautifulSoup** to extract articles from major Indian news outlets, including **The Hindu**, **NDTV**, **India Today**, and **Times of India**. Each article was stored in a structured **.csv** format with fields for **title**, **content**, **date**, and **source**.

Preprocessing Steps:

- **Text Cleaning:** Removal of HTML tags, special characters, and non-ASCII symbols.
- **Normalization:** Lowercasing, stopword removal, and whitespace normalization.
- **Metadata Extraction:** Extraction and standardization of publication dates and source tags.

3.2 Topic Modelling (Unsupervised NLP)

We used **BERTopic**, a transformer-based topic modelling framework that combines **sentence embeddings** and **class-based TF-IDF** to uncover thematic structures in the dataset. This method produces semantically dense, context-specific topics that capture nuanced themes within the health news articles.

3.3 Bias Detection (Emotion Analysis)

Media framing and emotional tone were analysed using:

- **Emotion Classification:** Utilizing the **GoEmotions** dataset, a fine-tuned **DistilBERT** model was employed to classify text into nuanced emotions like **fear**, **optimism**, **joy**, **anger**, and **sadness**, aiding in the detection of emotionally skewed language.

3.4 Misinformation Classification (Supervised DL)

A **DistilBERT** classifier was trained to distinguish potentially misleading articles from factual health news based on labelled emotion and framing features.

Training Setup:

- **Architecture:** Pre-trained distilbert-base-uncased from HuggingFace Transformers.
- **Data Split:** 80/20 train-test split.
- **Loss Function:** CrossEntropyLoss
- **Evaluation Metrics:** Accuracy, precision, recall, F1-score, with a confusion matrix for performance analysis.

3.5 Visualization and Interpretation

The following visualization techniques were employed to explore and present the findings:

- **Word Clouds:** Displayed dominant vocabulary themes across all articles.
- **Emotion & Framing Charts:** Bar and pie charts illustrating tone and narrative structures.
- **Topic Frequency Graphs:** Highlighted prominent health issues in Indian media.
- **Confusion Matrix:** Illustrated classifier performance.

4. Dataset Description

4.1 Acquisition Source and Creation Methodology

The dataset was curated through a systematic web scraping process targeting major Indian news outlets, including **The Times of India**, **NDTV**, **The Hindu**, and **India Today**. Custom Python scripts utilizing libraries such as **requests**, **BeautifulSoup**, and **newspaper3k** were developed to automate the extraction of health-related articles. The scraping focused on sections and keywords pertinent to health topics, ensuring the relevance of the collected data.

Each article was parsed to extract essential metadata, including the **title**, **publication date**, **source**, and **full textual content**. The collected data was structured and stored in **CSV** format, facilitating ease of access and analysis.

4.2 Dataset Characteristics

- **Size:** The final dataset comprises **50,000+** unique health-related news articles.
- **Language:** All articles are in English, reflecting the language of publication of the selected news outlets.
- **Structure:** Each record in the dataset includes the following fields:
 - **Title:** The headline of the article.
 - **Content:** The main body text of the article.
 - **Source:** The news outlet from which the article was obtained.
 - **Date:** The publication date of the article.
 - **URL:** The original link to the article.

4.3 Corpus Statistics and Visualization

To gain a better understanding of the corpus, key statistics were computed on the textual content:

- **Vocabulary Size:** The dataset contains approximately **64,509** unique words after standard preprocessing (lowercasing, punctuation removal, and stopwords filtering).
- **Top 20 Most Frequent Words:**

The 20 most commonly used words in the dataset include **health (11,526)**, **also (8,539)**, **like (5,503)**, **body (5,349)**, **may (5,257)**, **one (4,860)**, **people (4,854)**, **said (4,638)**, **blood (4,257)**, **help (4,213)**, **cancer (4,059)**, **risk (3,985)**, **time (3,707)**, **weight (3,614)**, **india (3,531)**, **disease (3,483)**, **heart (3,393)**, **years (3,336)**, **study (3,269)**, **even (3,221)**. These terms reflect the

dominant themes of the news articles, including health practices, disease mentions, and institutional references.

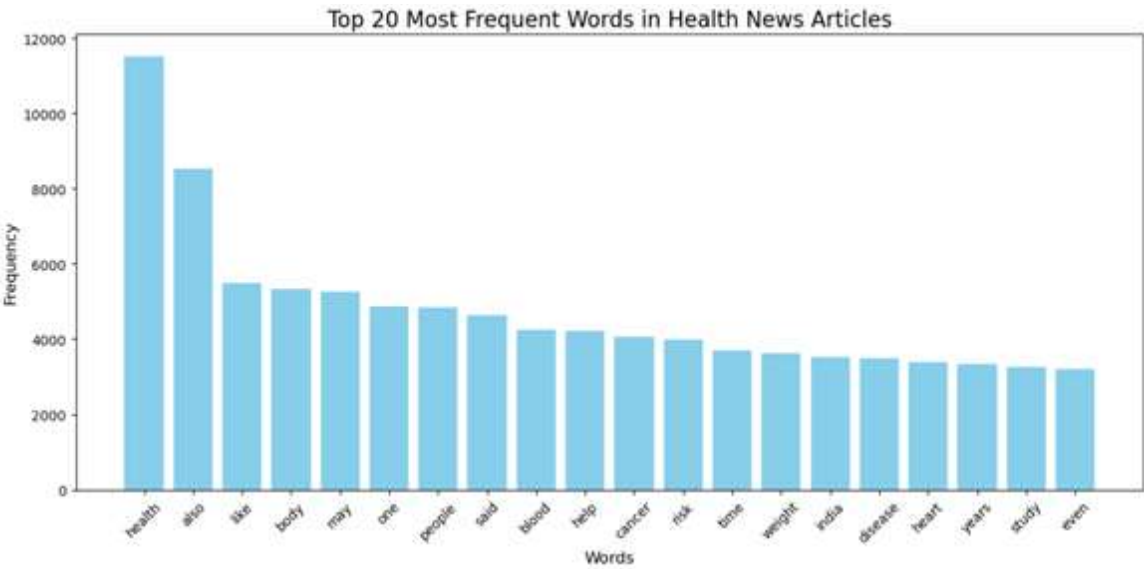


Figure 1: Top 20 Most Frequent Words in Health News Articles

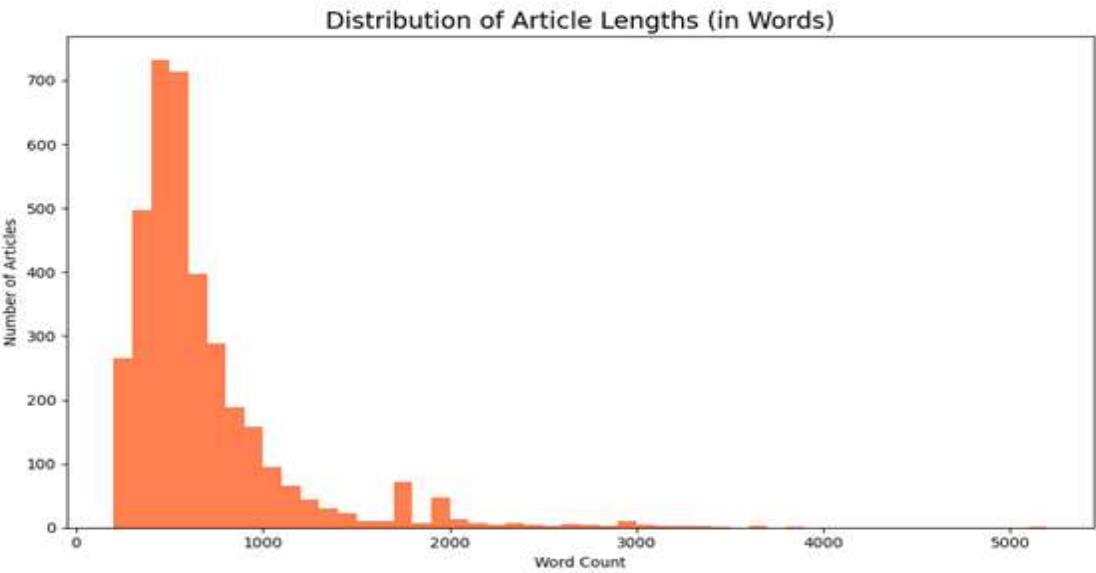


Figure 2: Distribution of Article Lengths (in Words)

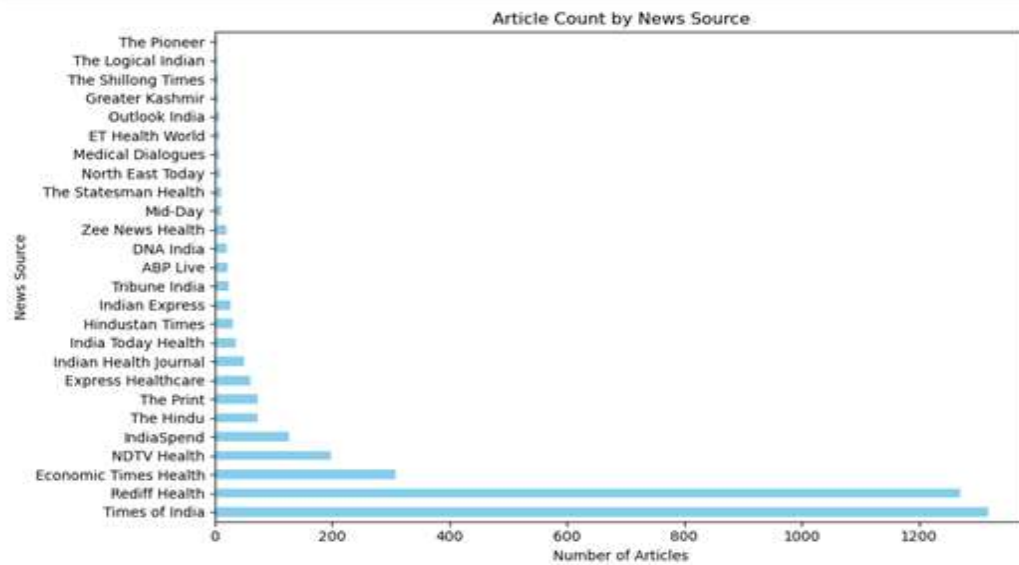


Figure 3: Article Count by News Source



Figure 4: Word Cloud of Most Frequent Terms

These corpus-level insights help contextualize downstream tasks like topic modelling and bias detection by grounding them in the actual vocabulary and linguistic patterns present in the media landscape.

5. Experimental Setup

The primary objective of this project is to investigate emotional framing and potential media bias in health-related news reporting by Indian media houses. This section outlines the multi-stage process used to build a robust pipeline for analysing emotional framing and bias, including data collection, preprocessing, topic modelling, and both emotion and bias classification using transformer models.

5.1 Data Collection and Preprocessing

Health-related news articles were collected from multiple Indian news outlets using a custom-built web scraping pipeline implemented in Python. This pipeline extracted articles covering medical topics, health tips, fitness trends, disease outbreaks, and wellness practices to ensure a comprehensive dataset representing diverse journalistic styles and perspectives in health reporting.

Key Tools and Libraries Used:

- **requests** / **httpx**: For making HTTP calls to news websites.
- **BeautifulSoup**: For parsing HTML and extracting article titles, dates, and full text.
- **pandas**: For managing structured tabular data, saving to CSV/Excel, and cleaning text.
- **re** (Regular Expressions): For removing unwanted characters, HTML tags, and formatting artifacts.

Preprocessing Steps:

- **Cleaning**: Removal of HTML elements, JavaScript snippets, duplicate lines, and noise (ads, links).
- **Normalization**: Converted all text to lowercase, removed special symbols, and normalized punctuation.
- **Metadata Extraction**: Extracted article source, publication date, and category (when available).
- **Keyword-Based Weak Supervision**:
 - **"misinfo"** if articles contained sensational keywords (e.g., "miracle", "cure", "no side effects").
 - **"accurate"** if scientific cues were present (e.g., "study", "trial").
 - **"neutral"** if neither condition was met.

This semi-automated labelling enabled early experimentation without extensive manual annotation.

- **Label Encoding**: Emotion categories (anger, joy, optimism, sadness) and misinformation labels (neutral, accurate, misinfo) were encoded using **LabelEncoder** from **scikit-learn**, ensuring compatibility with machine learning and deep learning models.

The final pre-processed dataset included both framing bias labels and emotion labels, enabling a deeper study of the intersection between misinformation framing and emotional tone in health reporting.

5.2 Topic Modelling Using BERTopic

To uncover the most reported themes in health journalism, we employed **BERTopic**, a state-of-the-art topic modelling tool that combines BERT embeddings with TF-IDF-based clustering.

Key Tools:

- **BERTopic:** Uses transformer-based embeddings for dense document vectors and clusters them into interpretable topics.
- **sentence-transformers:** Provides pretrained models (e.g., paraphrase-MiniLM-L6-v2) for generating sentence-level embeddings.
- **UMAP + HDBSCAN:** Used within BERTopic for dimensionality reduction and topic clustering.

Why BERTopic?

- Traditional models like LDA assume a bag-of-words structure and ignore word order.
- BERTopic uses semantic similarity via embeddings, making it robust to synonyms and nuanced language.

The output revealed 10 dominant topics, including themes like:

- Fitness and body image
- Cancer treatment
- Mental health
- Diet and weight loss
- Skin care and UV protection

These topics provided insights into the types of health content that might carry emotional or biased tones.

5.3 Bias Classification Using RoBERTa

After topic modelling, we employed the **RoBERTa (Robustly Optimized BERT)** model for framing bias detection. This phase focused on identifying editorial slants such as fear appeals, sensationalism, or framing through authority figures.

Key Tools and Approach:

- **transformers**: HuggingFace's library for fine-tuning transformer models.
- **Model**: roberta-base, known for high performance on classification tasks due to its larger pretraining corpus and longer context window.

Training Arguments:

- **Epochs**: 4
- **Batch size**: 8
- **Learning rate**: 2e-5

The model was trained on texts labelled with framing styles, learning linguistic cues that indicate subtle bias in health narratives. This step was completed before emotion classification to ensure training data was pre-analysed for editorial bias.

5.4 Label Encoding and Dataset Transformation

Before model training, both emotion and bias labels were converted from categorical text to numerical format using **LabelEncoder** from **sklearn.preprocessing**, essential for PyTorch-based models that require numerical class indices for classification tasks.

Additionally:

- The dataset was transformed into **HuggingFace Dataset** objects for efficient tokenization and training.
- **Tokenization**: Applied using model-specific tokenizers (RobertaTokenizer and DistilBertTokenizer) with padding and truncation for consistent input.

These steps optimized memory usage and ensured reproducibility during training.

5.5 Emotion Classification Using DistilBERT

The final classification task involved predicting the emotional tone of each article using **DistilBERT**, a lighter, faster version of BERT that maintains ~97% of its performance.

Training Configuration:

- **Model:** distilbert-base-uncased
- **Training/Testing Split:** 80/20 (stratified)
- **Epochs:** 5
- **Batch size:** 8
- **Evaluation Metrics:** Accuracy, Precision, Recall, F1-score
- **Loss Function:** CrossEntropyLoss

Why DistilBERT?

- Requires fewer computational resources (suitable for mid-range laptops).
- Fine-tunes faster than full BERT.
- Offers a good balance between accuracy and training efficiency.

5.6 Hardware Constraints and Environment

All training and inference tasks were performed locally on a personal machine with the following specifications:

- **CPU:** Intel i5-8265U
- **RAM:** 16 GB
- **GPU:** NVIDIA GeForce MX250
- **Platform:** Jupyter Notebook on VS Code (Windows 11)

Despite limited GPU power, using optimized models like DistilBERT and RoBERTa enabled efficient model training and evaluation within reasonable timeframes.

6. Results and Discussion

6.1. Quantitative Results

We evaluated two models for the emotion classification task:

- **Baseline:** Logistic Regression with TF-IDF features.
- **Fine-tuned Model:** DistilBERT fine-tuned on the labelled emotion dataset.

6.2. Baseline Model (Logistic Regression + TF-IDF)

- **Test Accuracy:** 80.89%
- **Macro-average F1-Score:** 0.51
- **Weighted-average F1-Score:** 0.79

Emotion	Precision	Recall	F1-Score	Support
Anger	1.00	0.06	0.12	16
Joy	0.74	0.62	0.67	185

Emotion	Precision	Recall	F1-Score	Support
Optimism	0.69	0.26	0.38	42
Sadness	0.83	0.96	0.89	479

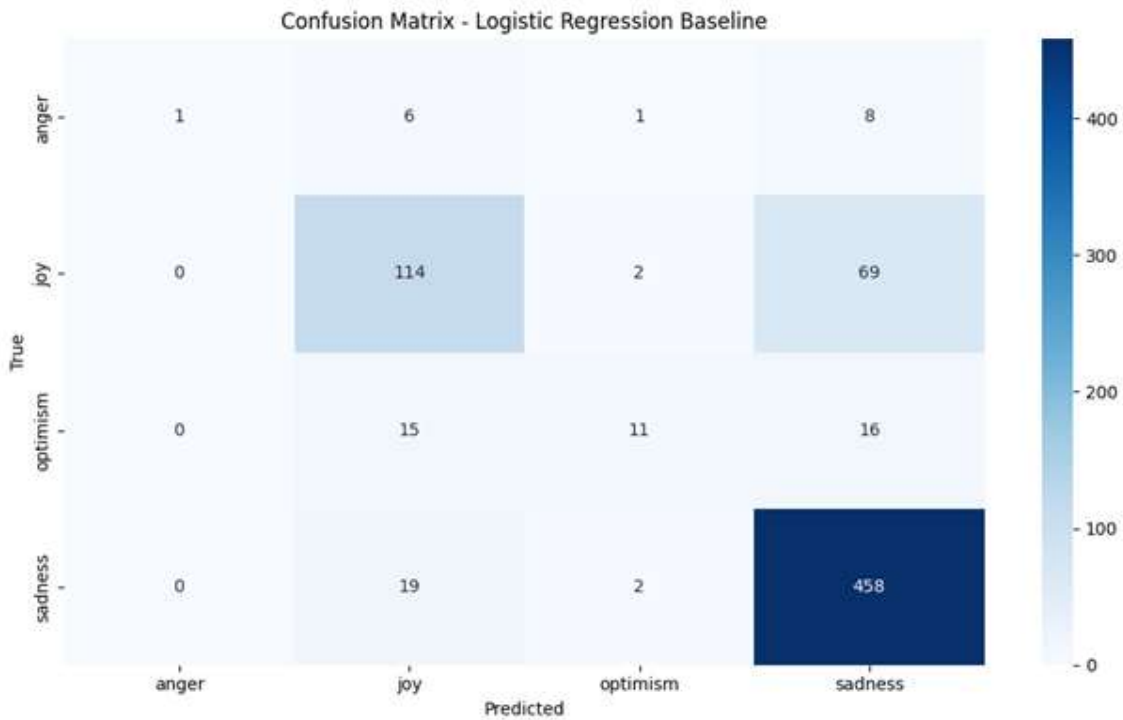


Figure 5.1: Confusion matrix for the Logistic Regression baseline model.

6.3. DistilBERT Fine-Tuned Model

- **Test Accuracy:** 86.29%
- **Macro-average F1-Score:** 0.7431
- **Weighted-average F1-Score:** 0.8622

Emotion	Precision	Recall	F1-Score	Support
Anger	0.6111	0.6875	0.6471	16

Emotion	Precision	Recall	F1-Score	Support
Joy	0.8000	0.7784	0.7890	185
Optimism	0.6410	0.5952	0.6173	42
Sadness	0.9134	0.9248	0.9191	479

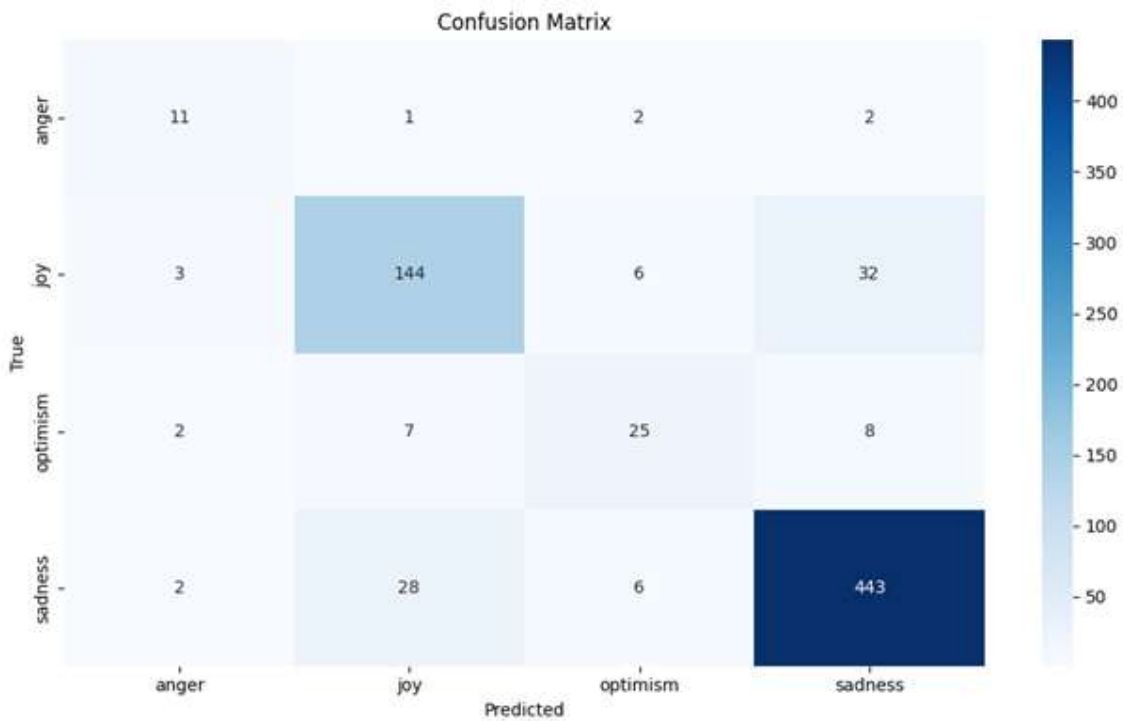


Figure 5.2: Confusion matrix for the DistilBERT fine-tuned model.

6.4. Emotion Distribution Across Bias Labels (roBERTa)

To complement the emotion classification model evaluation, we further explored the distribution of emotional tones across different bias categories using the **cardiffnlp/twitter-roberta-base-emotion model**. This model outputs softmax probabilities across four emotions: anger, joy, optimism, and sadness. Each article in the dataset was passed through

the model, and the resulting emotion scores were aggregated by bias label (neutral, possibly biased, possibly misinformative).

The resulting boxplot (see Figure 5.3) visualizes the emotional intensity across these three bias categories.

Interpretation:

- Articles labelled as possibly misinformative tend to show relatively higher probabilities for joy and sadness, suggesting a strategy of emotional persuasion.
- Neutral articles show more balanced and less extreme emotional tone, with sadness being dominant but less intense on average.
- The possibly biased group displays a moderate rise in joy and sadness, pointing to emotionally loaded but not overtly deceptive content.

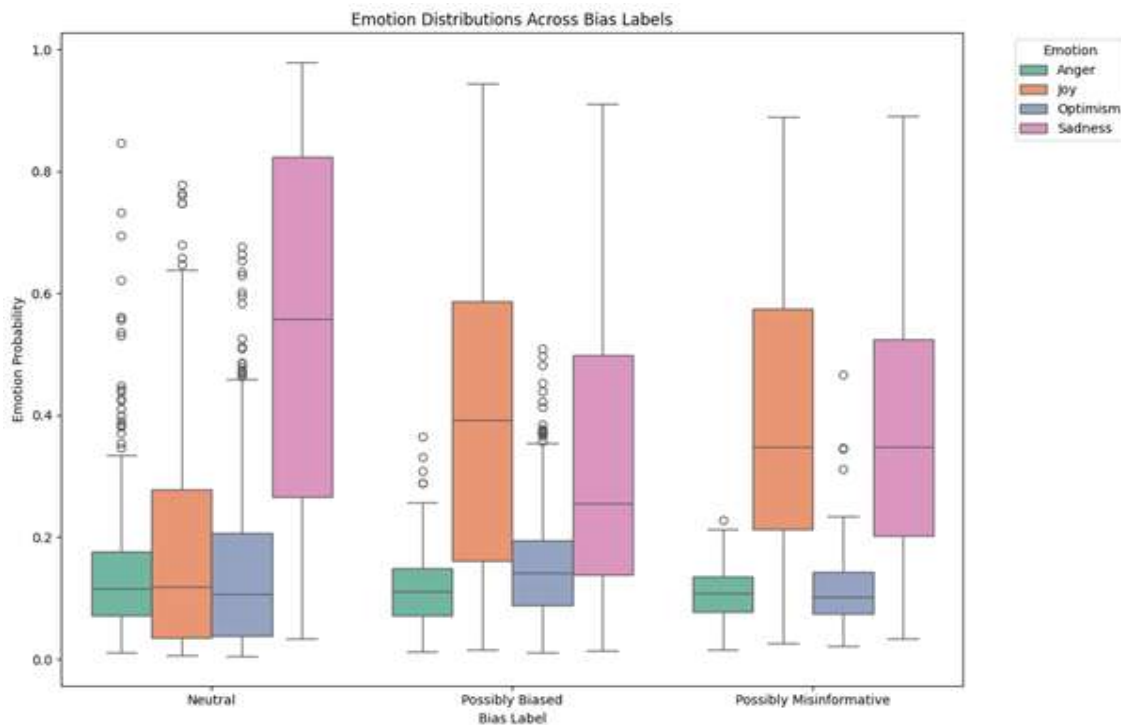


Figure 5.3: Boxplot of emotion distributions across bias labels using RoBERTa.

6.5. Observations

Performance by Class

- Sadness was consistently the most accurately predicted emotion across both models, attributed to its large class support.
- Anger and Optimism suffered in the baseline model with very low recall values (0.06 and 0.26), making them nearly invisible to the classifier.
- DistilBERT significantly improved recall for underrepresented classes such as "anger" (from 0.06 to 0.69) and "optimism" (from 0.26 to 0.60), making it more balanced in handling minority classes.
- A separate emotion detection model (RoBERTa fine-tuned on tweets) was used to analyse emotion distributions across bias categories.
- Articles labelled as possibly biased and possibly misinformative exhibited stronger emotional cues, particularly joy and sadness.
- This observation aligns with existing media research, which highlights emotional framing as a tool for persuasion and deception in health reporting.
- In contrast, neutral articles showed a more balanced and less emotionally polarized distribution.
- These insights indicate that emotion-aware misinformation detection could be improved by jointly modelling emotional intensity and bias.

6.6. Strengths / Improvements Over Baseline

- **Substantial F1-Score Gains:** The macro-average F1-score rose from 0.51 (baseline) to 0.74 (DistilBERT), a strong indicator of better performance across all classes.

- **Balanced Predictions:** While both models performed well on "sadness", DistilBERT was able to generalize across all classes with much more balanced recall and precision.
- **Efficient Transformer Fine-Tuning:** Despite being lightweight, DistilBERT showed strong performance and was feasible to train locally on a modest GPU (NVIDIA MX250).
- **Contextual Understanding:** Unlike the bag-of-words approach of TF-IDF, DistilBERT could understand semantic relationships in health-related texts, critical for emotion recognition.
- **Topic Modelling & Framing Insights:** BERTopic and RoBERTa-based framing models were used in parallel to segment articles by topic and framing bias, helping in future explainability.

6.7. Limitations

- **Class Imbalance:** Despite improvements, "anger" and "optimism" still lagged in performance due to limited samples in the training data.
- **Resource Constraints:** While the model performed well on local hardware, training time and inference latency were higher compared to the logistic regression baseline.
- **Lack of Additional Emotional Labels:** Only four emotions were considered. Expanding this spectrum could improve the real-world utility of the classifier.
- **No Temporal or Source-based Analysis:** The model treats all samples equally, without accounting for changes over time or source reliability—critical in health misinformation studies.
- **Framing Bias Evaluation:** While a RoBERTa-based classifier was trained for framing detection, it lacked quantitative evaluation due to the absence of gold-standard framing labels.

7. Scope for Future Work

While the current project demonstrates promising results in detecting emotional framing and bias in health-related news articles, several avenues for improvement and further exploration remain:

- **Joint Multi-task Learning:** Instead of treating emotion classification and bias detection as separate tasks, future work can adopt a multi-task learning framework where both are learned jointly. This could allow the model to capture shared representations and contextual dependencies between emotion and bias more effectively.
- **Modelling Context Beyond Sentences:** Current models primarily focus on individual article texts without considering broader context such as publication source, author background, or article history. Integrating metadata and longitudinal patterns may enhance bias detection accuracy.
- **Incorporation of Expert-Annotated Labels:** While weak supervision (keyword-based labelling) was efficient for initial bias tagging, it lacks nuance. Future iterations could use expert-annotated datasets for higher-quality training and evaluation, especially for subjective categories like framing bias.
- **Domain Adaptation for Health-Specific Language:** Although DistilBERT and RoBERTa are strong general-purpose models, health journalism often uses domain-specific terminology. Fine-tuning on a corpus of biomedical or health-related texts may yield better emotion and bias understanding.
- **Inclusion of Multilingual News Sources:** Given India's linguistic diversity, expanding the dataset to include regional language articles would make the system more representative and generalizable across different reader demographics.
- **Visual and Multimedia News Analysis:** Many online health articles include images, videos, or infographics, which also carry emotional framing. Future extensions could use multimodal models to incorporate visual sentiment analysis alongside text classification.

- **Explainability and Interpretability Tools:** As media bias detection is sensitive and impactful, integrating explainability frameworks such as LIME or SHAP could help interpret predictions and improve transparency in model outputs.
- **Real-Time Misinformation Monitoring System:** Building an interactive dashboard or pipeline to detect and flag emotionally manipulative or biased health content in real time could turn this research into a practical media monitoring tool for journalists or public health authorities.

8. Conclusion

Summary of Task

This project aimed to build an emotion classification model for health-related news articles, focusing on four core emotions—anger, joy, optimism, and sadness.

Approach and Relevance

We explored both a baseline ML pipeline (logistic regression + TF-IDF) and a fine-tuned transformer-based architecture (DistilBERT). The transformer model's superior ability to capture contextual language patterns proved highly effective in emotion classification tasks.

Obtained Results

The fine-tuned DistilBERT achieved an accuracy of 86.29% and macro-average F1-score of 0.7431, outperforming the baseline on nearly every metric and emotion class.

Reflection on Relevance

The ability to detect emotion in health news has wide implications for understanding misinformation, public anxiety, and information framing. Our approach demonstrates how even modest computational setups can yield powerful models using transformer architectures.

Application in Larger Context

This work has potential applications in automated journalism analysis, public health sentiment monitoring, and real-time emotional trend analysis during health crises such as pandemics.

9. References

1. Ayoub, J., Yang, X. J., & Zhou, F. (2021). Combat COVID-19 infodemic using explainable natural language processing models. *arXiv preprint arXiv:2103.00747*. <https://arxiv.org/abs/2103.00747>
2. Chen, B., Gao, D., Chen, Q., Huo, C., Meng, X., Ren, W., & Zhou, Y. (2021). Transformer-based language model fine-tuning methods for COVID-19 fake news detection. *arXiv preprint arXiv:2101.05509*. <https://arxiv.org/abs/2101.05509>
3. Chuang, Y., & Tsai, M. (2021). Cultural differences in emotions and opinions on painless delivery: A comparative study of Twitter and Weibo. *International Communication Studies*, 3(2), 45–60.
4. Fan, Y., Wang, Y., & Li, Q. (2022). Emotion classification utilizing transformer models with ECG signals. *International Journal of Modern Education and Computer Science*, 16(6), 24–31. <https://doi.org/10.5815/ijmeecs.2022.06.03>
5. Grootendorst, M. (2022). BERTopic: Neural topic modeling with a class-based TF-IDF procedure. *arXiv preprint arXiv:2203.05794*. <https://arxiv.org/abs/2203.05794>
6. Zhou, X., Yang, J., & Feng, Z. (2021). Combat COVID-19 infodemic using explainable natural language processing models. *arXiv preprint arXiv:2103.00747*. <https://arxiv.org/abs/2103.00747>