



SENTIMENT ANALYSIS FOR SOCIAL MEDIA USING BERT MODEL & XGBOOST



A PROJECT REPORT

Submitted by

DANISH RAJA M	(620821243020)
KAVIYARASAN K	(620821243053)
DHAKSHINAMOORTHY M	(620821243023)
KALAIYARASAN J	(620821243048)

in partial fulfillment for the award of the degree

of

**BACHELOR OF TECHNOLOGY
IN**

ARTIFICIAL INTELLIGENCE AND DATA SCIENCE

GNANAMANI COLLEGE OF TECHNOLOGY , NAMAKKAL-637 018

ANNA UNIVERSITY : CHENNAI 600 025

MAY 2025

BONAFIDE CERTIFICATE

Certified that this project report “ **SENTIMENT ANALYSIS FOR SOCIAL MEDIA USING BERT MODEL & XGBOOST** ” is the bonafide work of “ **DANISH RAJA M (620821243020), KAVIYARASAN K (620821243053), DHAKSHINAMOORTHY M(620821243023), KALAIYARASAN J (620821243048)** ” who carried out the project work under my supervision.

SIGNATURE

Mr.G.SIVAKUMAR, M.E, (Ph.D),,

HEAD OF THE DEPARTMENT

ASSISTANT PROFESSOR

Artificial Intelligence and Data Science

**Gnanamani College of Technology,
Namakkal – 637 018**

SIGNATURE

Mr.A.JEEVA , M.E, (Ph.D),,

SUPERVISOR

ASSISTANT PROFESSOR

Artificial Intelligence and Data Science

**Gnanamani College of Technology,
Namakkal – 637 018**

Submitted for the Anna University Project work Viva-Voce Examination Held on

Internal Examiner

External examiner

ACKNOWLEDGEMENT

We would like to express our deep sense of heartiest thanks to our beloved Chairman **Dr.T.ARANGANNAL** and Chairperson **Mrs.P.MALALEENA** Gnyanamani Educational Institutions, Namakkal, for giving an opportunity to do and complete this project.

Ms.MADHUVANTHINIE ARANGANNAL, Vice Chairman, Gnyanamani Educational Institutions, Namakkal, for their support and encouragement during our project work.

We would like to express our sincere gratitude to our Chief Administrative Officer **Dr.P.PREMKUMAR** Gnyanamani Educational Institutions, Namakkal, for providing us with indefinable support.

We would like to express our deep sense of gratitude to our Principal, **Dr.T.K.KANNAN** Gnanamani College of Technology, Namakkal, for motivating to complete the project work successfully.

We extend our thanks to our Academic Director **Dr.B.SANJAY GANDHI** Gnanamani College of Technology, Namakkal, for his encouragement during our project work successfully.

We take this opportunity to convey our heartiest thanks to **Mr.G.SIVAKUMAR**, Professor & Head Department of Artificial Intelligence and Data Science, pillar of support for the successful completion of the project.

We are extremely grateful to our Guide, **Mr.A.JEEVA**, Assistant Professor of Artificial Intelligence and Data Science department, for giving this opportunity with full encouragement to complete this project.

We express my sincere words of thanks of my parents, friends and all staff members of Artificial Intelligence and Data Science engineering, Gnanamani College of Technology, blessing to complete the project successfully.

[DANISH RAJA M]

[KAVIYARASAN K]

[DHAKSHINAMOORTHY M]

[KALAIYARASAN J]

DEPARTMENT OF ARTIFICIAL INTELLIGENCE AND DATA SCIENCE

INSTITUTE

VISION

Emerging as a technical institution of high standard and excellence to produce quality Engineers, Researchers, Administrators and Entrepreneurs with ethical and moral values to contribute the sustainable development of the society.

MISSION

We facilitate our students

MI : To have in-depth domain knowledge with analytical and practical skills in cutting edge technologies by imparting quality technical education.

MII : To be industry ready and multi-skilled personalities to transfer technology to industries and rural areas by creating interests among students in Research and Development and Entrepreneurship.

DEPARTMENT

VISION

To be a Centre of Artificial Intelligence and Data Science by imparting quality education, promoting research and innovation with global relevance.

MISSION

- To impart holistic education with niche technologies for the enrichment of knowledge and skills through updated curriculum and inspired learning.
- To empower valued based AI education to the students for developing intelligent systems and innovative products to address the societal problems with ethical value.
- To work in close liaison with industry to achieve socio-economic development.

DEPARTMENT OF ARTIFICIAL INTELLIGENCE AND DATA SCIENCE**PROGRAM OUTCOMES (POs)**

- 1. Engineering knowledge:** Apply the knowledge of mathematics, science, engineering fundamentals and an engineering specialization to the solution of complex engineering problems.
- 2. Problem analysis:** Identify, formulate, review research literature, and analyze complex engineering problems reaching substantiated conclusions using first principles of mathematics, natural sciences, and engineering sciences.
- 3. Design/development of solutions:** Design solutions for complex engineering problems and design system components or processes that meet the specified needs with appropriate consideration for the public health and safety, and the cultural, societal, and environmental considerations.
- 4. Conduct investigations of complex problems:** Use research-based knowledge and research methods including design of experiments, analysis and interpretation of data, and synthesis of the information to provide valid conclusions.
- 5. Modern tool usage:** Create, select, and apply appropriate techniques, resources, and modern engineering and IT tools including prediction and modeling to complex engineering activities with an understanding of the limitations.
- 6. The engineer and society:** Apply reasoning informed by the contextual knowledge to assess societal, health, safety, legal and cultural issues and the consequent responsibilities relevant to the professional engineering practice.
- 7. Environment and sustainability:** Understand the impact of the professional engineering solutions in societal and environmental contexts, and demonstrate the knowledge of, and need for sustainable development.
- 8. Ethics:** Apply ethical principles and commit to professional ethics and responsibilities and norms of the engineering practice.
- 9. Individual and team work:** Function effectively as an individual, and as a member or leader in diverse teams, and in multidisciplinary settings.
- 10. Communication:** Communicate effectively on complex engineering activities with the engineering community and with society at large, such as, being able to comprehend and write effective reports and design documentation, make effective presentations, and give and receive clear instructions.
- 11. Project management and finance:** Demonstrate knowledge and understanding of the engineering and management principles and apply these to one's own work, as a member and leader in a team, to manage projects and in multidisciplinary environments.
- 12. Life-long learning:** Recognize the need for, and have the preparation and ability to engage in independent and life-long learning in the broadest context of technological change.

PROGRAM SPECIFIC OUTCOMES (PSOs)

Graduates of the program will be able to

PSO-1: Understand, analyze and develop computer applications in data Mining/ Analytics, Cloud Computing, Networking, Security, etc. to meet the requirements of industry and society.

PSO-2: Enrich the ability to design and develop software and qualify for Employment, Higher studies and Research.

TABLE OF CONTENTS

CHAPTER NO	TITLE	PAGE NO
	ABSTRACT	i
	LIST OF TABLES	ii
	LIST OF FIGURES	ii
	LIST OF ABBREVIATIONS	iv
1	INTRODUCTION	13
	1.1 Background	13
	1.2 Problem Statement	14
	1.3 Objectives	15
	1.4 Scope of the Project	16
	1.4.1 Technical Scope	16
	1.4.2 Domain Scope	17
	1.4.3 Operational Scope	17
	1.4.4 Limitations of Scope	18
2	LITERATURE REVIEW	19
	2.1 Traditional Approaches to Sentiment Analysis	19
	2.2 Rule-Based and Lexicon-Based Approaches	19
	2.3 Emergence of Deep Learning in Sentiment Analysis	20
	2.4 Transformer Architecture and BERT	20
	2.5 BERT for Sentiment Analysis	21
	2.6 Research Gap	21
3	SYSTEM ANALYSIS	22
	3.1 Existing System	22
	3.2 Proposed System	22
	3.3 Feasibility Study	23

	3.4	Comparative Analysis	23
4		SYSTEM DESIGN	25
	4.1	Architectural Design	25
	4.2	System Components	27
	4.3	Data Flow Diagram	27
	4.4	Technology Stack	28
	4.5	Design Considerations	28
5		IMPLEMENTATION	29
	5.1	Environment Setup	29
	5.2	Data Loading And Preprocessing	29
	5.3	Model Architecture	30
	5.4	Training Procedure	31
	5.5	Evaluation And Results	32
6		RESULT AND ANALYSIS	33
	6.1	Evaluation Metrics	33
	6.2	Confusion Matrix	33
	6.3	Comparison With Other Models	35
	6.4	Sample Predictions And Visualizations	35
	6.5	Discussion	35
7		USER INTERFACE	36
	7.1	Overview Of The Interface	36
	7.2	Technologies Used	36
	7.3	Key Functionalities	36
	7.4	Screenshots And Ui Layout	37
	7.5	User Experience And Accessibility	38
8		TESTING AND VALIDATIONS	39
	8.1	Types Of Testing	39
	8.2	Test Cases And Outcomes	39

	8.3	Validation Of Model Performance	39
	8.4	Error Handling And Edge Cases	40
9		CONCLUSION AND FUTURE SCOPE	41
	9.1	Conclusion	41
	9.2	Future work	41
		APPENDIX 1 : SOURCE CODE	42
		APPENDIX 2 : SCREENSHOT	55
		APPENDIX 3 : LIST OF PUBLICATION	56
		REFERENCES	57

ABSTRACT

In today's digital age, social media has become much more than a place to share updates or connect with friends—it's now a dynamic space where people regularly express their thoughts, feelings, and opinions about everything from daily life to global events. With millions of posts being made every day, understanding what people are truly feeling behind their words has become incredibly important. This is where sentiment analysis comes in—a technique that helps us make sense of emotions conveyed through text. What sets our approach apart is the combination of two powerful tools: BERT [Bidirectional Encoder Representations from Transformers] and XGBoost [Extreme Gradient Boosting]. BERT is known for its deep understanding of language context—it doesn't just look at individual words, but understands how they're used in a sentence, even picking up on subtle cues like sarcasm or irony. On the other hand, XGBoost is a fast and highly accurate machine learning algorithm, great for handling classification tasks efficiently. By combining these two—BERT for feature extraction and XGBoost for final classification—our hybrid model is able to classify social media comments into one of four sentiment categories: happy, sad, sarcasm, and irony. This four-class system goes beyond the traditional "positive/negative" sentiment split and helps capture the more complex, nuanced emotions often found in online conversations. The goal of our approach is not just to improve the accuracy of sentiment detection, but also to make it scalable and practical for real-world applications, such as brand monitoring, mental health analysis, or social research. In essence, we're aiming to build a smarter, more sensitive system that truly understands the emotional undertone in the way people communicate online.

LIST OF TABLES

TABLE NO	TITLE	PAGE NO
1	Model selection	24
2	Dataset	29
3	Present a comparative evaluation of the proposed model	33

LIST OF FIGURES

FIGURE NO	TITLE	PAGE NO
Figure 1	Overall System Workflow for Sentiment Analysis using BERT and XGBoost.	26
Figure 2	Architecture of Sentiment analysis system design	26
Figure 3	Data Flow Diagram	27
Figure 4	BERT Model Architecture	31
Figure 5	Confusion matrix	34
Figure 6	Accuracy comparison of Models	34
Figure 7	Webpage Overview	37

LIST OF ABBREVIATIONS

Abbreviation		Fullform
AI	—	Artificial Intelligence
API	—	Application Programming Interface
BERT	—	Bidirectional Encoder Representations from Transformers
BoW	—	Bag of Words
CNN	—	Convolutional Neural Network
CSV	—	Comma-Separated Values
DL	—	Deep Learning (inferred but not explicitly abbreviated)
HTML	—	HyperText Markup Language
IT	—	Information Technology
LSTM	—	Long Short-Term Memory
ML	—	Machine Learning (inferred from context)
MLM	—	Masked Language Modeling
NLP	—	Natural Language Processing
NSP	—	Next Sentence Prediction
PO	—	Program Outcome
PSO	—	Program Specific Outcome
RNN	—	Recurrent Neural Network
SA	—	Sentiment Analysis
SVM	—	Support Vector Machine
TF	—	Term Frequency
TF-IDF	—	Term Frequency-Inverse Document Frequency
UI	—	User Interface
UAT	—	User Acceptance Testing
XAI	—	Explainable Artificial Intelligence
XGBoost	—	Extreme Gradient Boosting

CHAPTER 1

INTRODUCTION

Sentiment Analysis (SA) is a key area of Natural Language Processing (NLP) that focuses on interpreting people's opinions, emotions, and sentiments. It involves multiple stages, including data collection, preprocessing, feature extraction, and classification. With the rapid expansion of social media content across various industries, SA has become essential for understanding public opinion and consumer behaviour. Online reviews play a crucial role in influencing purchasing decisions. For businesses, analyzing customer feedback provides valuable insights into consumer emotions and perceptions. By leveraging SA, businesses can enhance their products and services, improve customer satisfaction, and maximize profitability.

1.1 Background

The exponential growth of digital communication and the pervasive adoption of social media platforms have revolutionized how individuals, organizations, and societies interact, share, and consume information. Platforms like Twitter, Facebook, Instagram, Reddit, and YouTube have become central to modern discourse, providing users with a means to express opinions, experiences, and emotions in real-time. As a result, these platforms host an immense volume of user-generated content, much of which is opinionated and sentiment-rich. This dynamic environment has opened up a significant avenue for extracting valuable insights through the computational process of sentiment analysis.

Sentiment analysis, often referred to as opinion mining, is a subfield of Natural Language Processing (NLP) that involves systematically identifying, extracting, and categorizing sentiments expressed in a body of text. It aims to determine the emotional tone behind textual inputs, thereby providing a means to understand the attitudes, emotions, and opinions of speakers or writers. These sentiments are typically categorized as positive, negative, or neutral, although more nuanced categorizations such as joy, anger, sadness, or fear are also explored in advanced studies.

The value of sentiment analysis extends across multiple sectors. In the business world, companies analyze customer reviews, social media interactions, and feedback forms to assess brand reputation and customer satisfaction. In politics, it is used to gauge public opinion around candidates or policies. In journalism and marketing, sentiment data assists in shaping narratives or crafting targeted campaigns. Given its wide-ranging applicability, sentiment analysis has become a powerful decision-support tool in the data-driven era.

However, analyzing sentiments on social media presents unique challenges. Unlike formal texts, social media content is informal, unstructured, and highly diverse in style. Users often

employ slang, abbreviations, emojis, hashtags, sarcasm, and code-switching (mixing languages), which can obscure the actual meaning of the message. Traditional sentiment analysis approaches, such as lexicon-based models or classical machine learning classifiers like Naive Bayes, Logistic Regression, or Support Vector Machines (SVM), often fall short in capturing the full semantic and contextual meaning of such posts. These models typically rely on manual feature engineering and predefined rules, limiting their scalability and robustness.

In response to these limitations, the field has increasingly shifted toward deep learning approaches, especially with the emergence of transformer-based architectures. Among these, BERT (Bidirectional Encoder Representations from Transformers), developed by Google in 2018, has emerged as a groundbreaking model. BERT is pre-trained on large corpora using a deep bidirectional understanding of context. Unlike traditional language models that process text sequentially, BERT processes words in relation to all other words in a sentence, both to the left and the right. This allows the model to understand the full context of a word and resolve ambiguities in language more effectively.

BERT's architecture and attention mechanism enable it to excel in downstream NLP tasks such as question answering, named entity recognition, and, crucially, sentiment classification. By fine-tuning a pre-trained BERT model on domain-specific sentiment data, researchers and developers can build highly accurate, context-aware sentiment classifiers. This innovation has led to a significant leap in the quality and reliability of sentiment analysis systems, particularly when applied to noisy and dynamic sources like social media.

Given the increasing relevance of social media as a barometer of public mood and the limitations of earlier analytical tools, this project explores the integration of BERT for sentiment analysis on social media content.

1.2 Problem Statement

Despite the proliferation of digital text and the widespread availability of user opinions on social media, extracting meaningful and actionable sentiment insights remains a formidable challenge. Traditional sentiment analysis approaches, including rule-based systems and classical machine learning models, are often inadequate when applied to the dynamic and unstructured nature of social media content. These models typically struggle with context sensitivity, fail to capture nuanced expressions such as sarcasm or irony, and are limited in their ability to process informal language structures, emojis, or domain-specific slang.

Moreover, social media data is inherently noisy, with inconsistent grammar, abbreviations, and multilingual code-switching. These characteristics further complicate the accurate detection of sentiment. While lexicon-based approaches rely on predefined lists of sentiment-

laden words, they often misinterpret the polarity of words when used in a different context. For instance, the word "sick" may convey a negative sentiment in a medical context but could imply admiration in a casual tweet like "That movie was sick!"

Another significant limitation of existing methods lies in their inability to grasp long-range dependencies in text and the contextual relationship between words. This often leads to inaccurate or shallow sentiment predictions, particularly in complex or multi-sentence posts.

There is a clear and pressing need for a more sophisticated and context-aware approach that can process social media text with high precision and flexibility. Transformer-based models like BERT, which are capable of bidirectional language understanding, offer a promising solution. However, their application requires proper fine-tuning, preprocessing, and evaluation frameworks tailored specifically for sentiment analysis.

The core problem, therefore, is to design and implement an intelligent sentiment analysis system capable of handling informal, unstructured social media text using a deep learning approach—specifically, the BERT model. The system must achieve high accuracy, context-awareness, and real-world applicability for a wide range of social media platforms.

1.3 Objectives

The primary objective of this project is to develop a robust and intelligent sentiment analysis system that leverages the deep learning capabilities of BERT (Bidirectional Encoder Representations from Transformers) for accurate sentiment classification of social media content. The project aims to address the limitations of traditional sentiment analysis methods by utilizing a context-aware model that can effectively process the informal, diverse, and often ambiguous language used across various social platforms.

The detailed objectives of the project are outlined as follows:

1. To design a sentiment analysis framework based on BERT architecture

This involves selecting a suitable pre-trained BERT model and fine-tuning it for the specific task of sentiment classification on social media text. The model will be trained to categorize sentiments into classes such as positive, negative, and neutral, based on the context and tone of the input.

2. To implement preprocessing techniques tailored to social media data

This includes cleaning the text data by removing noise such as URLs, hashtags, mentions, emojis, abbreviations, and correcting grammatical inconsistencies. The goal is to transform unstructured and informal user content into a format suitable for model ingestion without losing contextual meaning.

3. **To collect and curate a representative dataset**

Acquiring a rich and diverse dataset is essential for training and evaluation. The project aims to use publicly available datasets or APIs (e.g., Twitter API) to collect posts, tweets, or comments that exhibit varied sentiment expressions across different domains and topics.

4. **To evaluate the model using industry-standard performance metrics**

The system's effectiveness will be measured using metrics such as accuracy, precision, recall, and F1-score. Comparisons may also be drawn with baseline models or classical techniques to demonstrate the superiority of the BERT-based approach.

5. **To develop a user-friendly interface for real-time sentiment analysis**

The final system should include a graphical or web-based interface where users can input social media text or links and receive immediate sentiment classification results. This interface should also visualize sentiment trends and allow for interactive analysis.

6. **To ensure scalability and adaptability of the system**

The project aims to build a solution that is not only accurate but also scalable to handle real-world data volumes. Additionally, the model should be adaptable to new topics and evolving language trends through periodic retraining or dynamic updates.

By achieving these objectives, the project will provide a practical and scalable solution for real-time sentiment analysis on social media platforms, delivering both academic value and real-world applicability.

1.4 Scope of the Project

This project focuses on designing and implementing a deep learning-based sentiment analysis system specifically tailored for **social media content**, using the **BERT (Bidirectional Encoder Representations from Transformers)** model. The scope is strategically defined to concentrate on the practical challenges, technical implementation, and performance evaluation associated with sentiment classification of informal, real-time text data from platforms like Twitter, Facebook, and Instagram.

1.4.1 Technical Scope

The core of the project lies in utilizing a transformer-based architecture to analyze user-generated text. The system will:

- Leverage a pre-trained BERT model, fine-tuned on sentiment-labeled datasets.
- Preprocess raw social media data using NLP techniques such as tokenization, stopword removal, emoji filtering, and handling of abbreviations and special characters.

- Support multi-class sentiment classification (Positive, Negative, Neutral), with potential to extend to emotion-based categories like Joy, Anger, or Sadness.
- Be evaluated using robust performance metrics like Accuracy, Precision, Recall, and F1-score.

The project will involve both backend development (model training, preprocessing, classification logic) and frontend implementation (interface for user interaction). Integration of the model into a **web-based or GUI-based tool** will also be within scope, allowing users to analyze input text or links dynamically.

1.4.2 Domain Scope

The sentiment analysis system will be trained and evaluated using **publicly available datasets** such as the Sentiment140 dataset, Twitter US Airline Sentiment dataset, or other real-time data collected via APIs. The content scope includes:

- Textual posts and comments from social platforms.
- English language only (non-English and multilingual content is outside current scope);
- Short to medium-length texts (tweets, comments, captions, etc.).

While this project does not aim to analyze multimedia content (e.g., images, videos), it acknowledges that future extensions could incorporate multimodal sentiment analysis.

1.4.3 Operational Scope

The system is designed for general-purpose sentiment analysis and can be applied in several domains including:

- **Brand and product monitoring** for companies seeking to understand customer feedback.
- **Social campaign tracking** to measure public response to political, social, or marketing efforts.
- **Public sentiment monitoring** during events such as elections, disasters, or public health campaigns.

The system will function as a **proof of concept** and can be further extended into a scalable enterprise solution. While real-time, high-frequency data ingestion is not a requirement at this stage, the model and system architecture will be designed to support such upgrades in the future.

1.4.4 Limitations of Scope

Some limitations are acknowledged to maintain focus and feasibility:

- The project will not cover multilingual or cross-lingual sentiment analysis.
- Deep sarcasm detection and context-aware irony handling will be partially addressed but may not achieve perfect accuracy;
- Real-time streaming from APIs (e.g., Twitter stream) is not included in the initial phase but can be added later;
- Sentiment is subjective therefore, model predictions will be probabilistic, not absolute.

Chapter 2

Literature Review

Introduction

The field of sentiment analysis has garnered significant attention in recent years due to the exponential growth of social media platforms, where users express their opinions, emotions, and attitudes openly. This chapter explores previous research in the domain of sentiment analysis, particularly focusing on traditional machine learning approaches and the transition to deep learning models culminating in the emergence of transformer-based architectures such as BERT. By evaluating past literature, this chapter establishes a foundation for the necessity and effectiveness of using BERT for social media sentiment classification.

2.1 Traditional Approaches to Sentiment Analysis

Earlier research in sentiment analysis primarily relied on traditional machine learning models such as:

- **Naïve Bayes**
- **Support Vector Machines (SVM)**
- **Logistic Regression**
- **Decision Trees**

These models depended heavily on **manual feature engineering**, where textual data was converted into numerical features using techniques like:

- Bag of Words (BoW)
- Term Frequency-Inverse Document Frequency (TF-IDF)
- N-grams

While these methods performed reasonably well for structured and formal text, they often struggled with the noisy and informal nature of social media text. Furthermore, traditional models lacked the ability to understand **context and word semantics**, limiting their performance in handling nuanced or sarcastic expressions common in tweets.

2.2 Rule-Based and Lexicon-Based Approaches

In addition to statistical models, lexicon-based techniques such as **VADER (Valence Aware Dictionary and sEntiment Reasoner)** and **SentiWordNet** were developed to analyze sentiment using pre-defined dictionaries of emotional words. These approaches:

- Were interpretable and easy to implement.
- Could work on small datasets without needing training.
- Failed to generalize well to complex language use, idioms, or domain-specific vocabulary.

Although lexicon-based models provided some insight into sentiment polarity, they lacked adaptability to different linguistic contexts and often misclassified sentiment when faced with sarcasm or implicit meaning.

2.3 Emergence of Deep Learning in Sentiment Analysis

To overcome the limitations of manual feature engineering and static word representations, researchers began adopting deep learning techniques such as:

- **Convolutional Neural Networks (CNNs)**
- **Recurrent Neural Networks (RNNs)**
- **Long Short-Term Memory networks (LSTMs)**

These models introduced dynamic learning of word embeddings and temporal patterns in text sequences. Embedding methods such as **Word2Vec**, **GloVe**, and **FastText** became popular for generating semantic vector representations of words. LSTM-based models achieved better performance by learning long-term dependencies in textual data, but they still had drawbacks:

- Training was time-consuming and computationally expensive.
- Contextual understanding was limited to unidirectional sequences.
- They could not efficiently model very long texts or manage parallelization.

2.4 Transformer Architecture and BERT

The **transformer model**, introduced by Vaswani et al. in 2017, revolutionized Natural Language Processing (NLP) by discarding recurrence in favor of attention mechanisms. The transformer architecture enabled:

- **Self-attention:** Capturing relationships between all words in a sentence, regardless of distance.
- **Parallel processing:** Significantly faster training times.
- **Better contextual understanding:** Through bidirectional encoding.

Building on transformers, Google introduced **BERT (Bidirectional Encoder Representations from Transformers)** in 2018. BERT brought two major innovations:

1. **Masked Language Modeling (MLM):** Predicting missing words in a sentence, helping the model understand context on both sides.
2. **Next Sentence Prediction (NSP):** Learning sentence relationships, crucial for sentiment flow.

BERT achieved state-of-the-art results on a wide range of NLP tasks, including sentiment classification, question answering, and textual entailment.

2.5 BERT for Sentiment Analysis

Several studies have successfully implemented BERT for sentiment classification:

- Devlin et al. (2018) reported high accuracy on the SST-2 sentiment dataset using pre-trained BERT.
- Sun et al. (2019) fine-tuned BERT on Twitter data for multilingual sentiment analysis, achieving significant performance gains.
- Researchers compared BERT with LSTM and CNN and found BERT consistently outperformed traditional architectures due to its deep contextual understanding.

These successes highlight BERT's capability to handle complex language structures, informal expressions, sarcasm, and multilingual content—making it ideal for analyzing sentiment from social media platforms.

2.6 Research Gap

Despite the notable advancements in sentiment analysis brought about by deep learning and transformer-based models like BERT, several gaps remain in existing research. A key limitation is the lack of domain-specific fine-tuning for sentiment classification in specialized contexts such as finance, healthcare, or politics, where sentiment cues may differ significantly from general-purpose datasets. Additionally, while BERT achieves high accuracy, its computational requirements can lead to increased latency during real-time sentiment prediction, making deployment on resource-constrained platforms challenging. Another important issue is the frequent class imbalance in sentiment datasets—particularly the underrepresentation of neutral or sarcastic sentiments—which can bias model predictions and reduce generalization. Furthermore, existing literature has not fully explored the interpretability of BERT's decisions, which is crucial for applications requiring transparency. Addressing these gaps forms the basis for developing improved, efficient, and context-aware sentiment analysis systems tailored for real-world social media applications.

Chapter 3

System Analysis

Introduction

System analysis is a critical phase in the project development lifecycle. It involves evaluating both the functional and non-functional aspects of the system to ensure that the final implementation meets user requirements and performance expectations. In the context of this project **Sentiment Analysis for Social Media Using BERT Model** system analysis helps to define the scope, understand the user needs, assess the technical feasibility, and design a robust architecture capable of processing and analyzing massive volumes of unstructured social media text data. This chapter presents a comprehensive analysis of the existing system limitations, the proposed system's improvements, feasibility study, and a comparative evaluation that justifies the use of BERT as the core model for sentiment classification.

3.1 Existing System

The existing systems for sentiment analysis are primarily built using:

- Traditional machine learning algorithms such as Naïve Bayes, SVM, and Logistic Regression.
- Lexicon-based models like VADER and TextBlob.
- Shallow neural networks including CNNs and LSTMs with static word embeddings.

While these approaches have been effective for structured or domain-constrained data, they fall short when applied to real-world social media content due to several challenges:

- **Limited contextual understanding:** Traditional models cannot capture nuanced relationships between words, especially in sarcastic or idiomatic expressions.
- **Manual feature engineering:** Requires extensive preprocessing and handcrafted features, increasing development effort and reducing adaptability.
- **Scalability issues:** Older models are not optimized for real-time inference over large datasets.
- **Inadequate sentiment precision:** Particularly with neutral and mixed-emotion expressions, leading to misclassification.

These limitations create a significant performance gap when applied to informal, highly diverse, and dynamic social media content.

3.2 Proposed System

The proposed system addresses the shortcomings of traditional approaches by leveraging **BERT**, a transformer-based pre-trained language model developed by Google. The architecture of the proposed sentiment analysis system is composed of the following key modules:

1. **Data Collection Module** – Extracts and stores real-world tweets or social media posts from sources like Twitter APIs or labeled datasets.
2. **Preprocessing Module** – Cleans raw input data by removing URLs, mentions, special characters, and normalizing text for further processing.
3. **Tokenization Module** – Converts the cleaned text into BERT-compatible token IDs and attention masks using the Hugging Face tokenizer.
4. **Model Training and Fine-Tuning Module** – Adapts the bert-base-uncased model using labeled data to classify sentiment into Positive, Negative, or Neutral.
5. **Prediction and Visualization Module** – Applies the fine-tuned model to new data and visualizes output through dashboards or plots for better interpretability.

This approach ensures **deep contextual understanding**, faster adaptation to domain-specific text, and superior accuracy over prior models. Moreover, the use of a pre-trained model drastically reduces training time while still achieving state-of-the-art performance.

3.3 Feasibility Study

A feasibility study was conducted across three key dimensions—**technical**, **operational**, and **economic**—to ensure the practicality and sustainability of the proposed system.

- **Technical Feasibility:** The proposed system uses Python-based frameworks (e.g., Transformers, PyTorch, TensorFlow) and leverages GPU acceleration for model training. It is technically feasible given the availability of tools like Google Colab or cloud-based services for training and deployment.
- **Operational Feasibility:** From an operational perspective, the system is user-friendly and adaptable. It can be integrated with social media APIs to handle real-time data and can be expanded with multilingual support or emotion recognition modules in future versions.
- **Economic Feasibility:** BERT is open-source, and most of the libraries used are freely available. With access to cloud computing or academic GPUs, the overall cost of development and deployment remains low, making the project economically viable for both research and real-world applications.

3.4 Comparative Analysis

The transition from traditional sentiment analysis methods to transformer-based models like BERT marks a significant leap in natural language understanding. Traditional systems often rely on shallow machine learning techniques, bag-of-words approaches, or static word embeddings, which struggle to capture the nuances of human language, particularly when context, sarcasm, or implicit sentiments are involved. These limitations result in lower accuracy, poor generalization, and increased manual effort for feature engineering. On the other hand, the proposed BERT-based model leverages bidirectional attention mechanisms to understand word meanings in relation to their full context—both preceding and succeeding words. This contextual richness allows the system to excel at detecting subtle differences in sentiment, especially in informal and noisy text like social media posts.

Moreover, BERT’s architecture eliminates the need for manual feature extraction by learning complex representations directly from data, thus saving time and improving performance. It can also be fine-tuned with domain-specific datasets, making it more adaptable and robust compared to rule-based or lexicon-based systems. Performance benchmarks consistently demonstrate that BERT-based models outperform their predecessors in tasks involving sentiment polarity classification. This comparative advantage, coupled with advancements in transfer learning and cloud-based deployment options, makes BERT not only a technically superior choice but also a practical one for real-world applications.

Feature	Existing System	Proposed BERT-Based System
Context Understanding	Weak (Unigram/Bigram-based)	Strong (Bidirectional attention)
Feature Engineering	Manual	Automatic (via pretraining)
Accuracy	Moderate	High
Real-Time Performance	Limited	Scalable with GPU support
Adaptability	Low (Domain-specific)	High (Transfer learning)

Table 1: Model selection

Chapter 4

System Design

Introduction

System design is the blueprint of the entire application. It defines the architecture, modules, data flow, and interactions between components, enabling efficient development and future scalability. For the project titled "*Sentiment Analysis for Social Media Using BERT Model*," the design phase ensures a modular and maintainable structure that supports robust sentiment classification using real-time or batch social media data.

This chapter details the architectural layout, data flow, module interactions, and design choices employed in the development of the system. The focus is on leveraging the pre-trained BERT model effectively within a custom pipeline tailored for sentiment analysis.

4.1 Architectural Design

The architecture follows a **layered modular structure**, enabling separation of concerns and better debugging. It can be categorized into the following layers:

1. Presentation Layer

This includes the user interface or dashboard used to input social media post URLs, view sentiment results, and interpret visualization outputs. It can be a web interface or notebook-based display.

2. Processing Layer

This is the core of the system where:

- Raw text is pre-processed
- Sentences are tokenized
- The BERT model is applied
- Predictions are generated

3. Data Layer

Responsible for collecting and storing data from social media APIs or pre-built datasets like Twitter Sentiment140 or Kaggle sentiment corpora. This layer ensures data availability for both training and real-time inference.

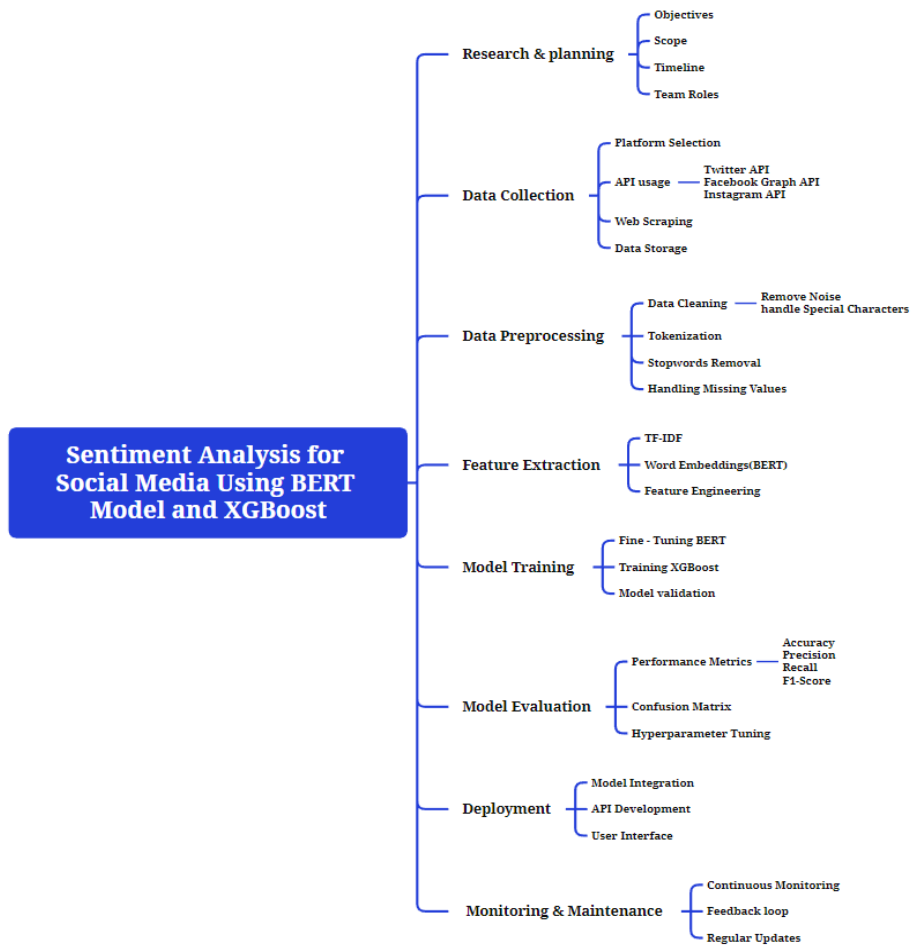


Figure 1: Overall System Workflow for Sentiment Analysis using BERT and XGBoost.

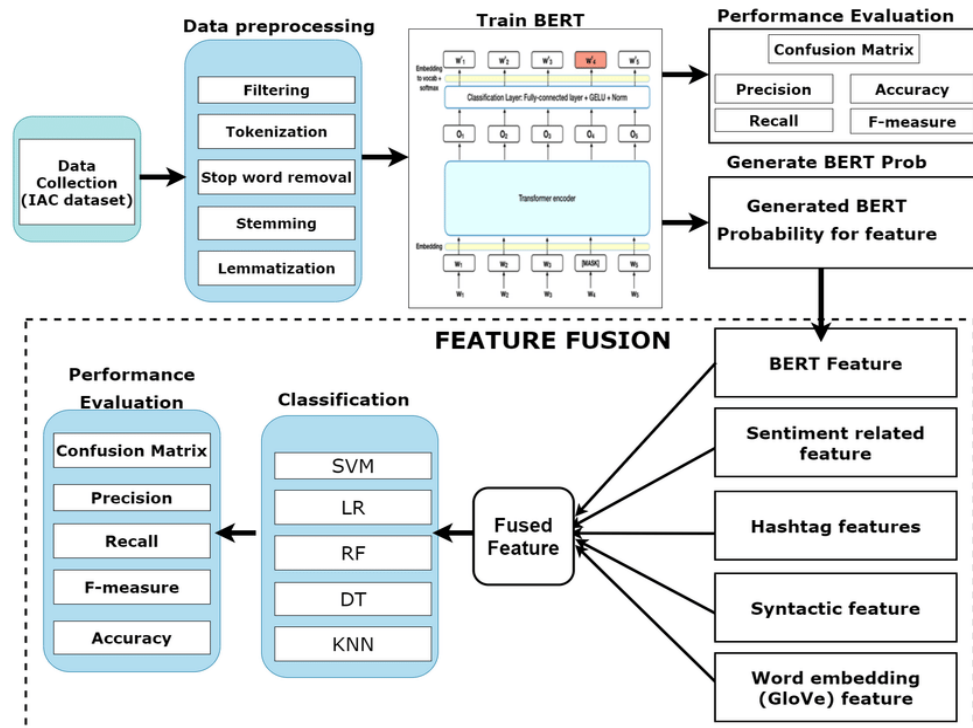


Figure 2: Architecture of Sentiment analysis system

4.2 System Components

The system is divided into distinct modules, each performing a specific function:

- **Data Acquisition Module**
Collects posts or comments from platforms like Twitter using APIs or datasets. It handles rate limits, filtering, and content selection.
- **Data Preprocessing Module**
Cleans and normalizes input by removing unnecessary elements like hashtags, links, mentions, emojis, and converting text to lowercase. It ensures the data is suitable for tokenization.
- **Tokenizer Module**
Uses the BertTokenizer from Hugging Face to tokenize input text into subword units, convert them into token IDs, and generate attention masks required by the BERT model.
- **Model Integration Module**
Loads the pre-trained bert-base-uncased model and adds a classification head (dense layer with softmax) for sentiment prediction. It handles fine-tuning, inference, and confidence scoring.
- **Prediction Output Module**
Converts model output into interpretable sentiment classes (Positive, Neutral, Negative) and formats the results for display or further analysis.
- **Visualization Module**
Generates graphs or charts (e.g., pie charts, bar graphs) for summarizing sentiment trends over a batch of inputs. This helps end-users quickly understand the overall sentiment landscape.

4.3 Data Flow Diagram

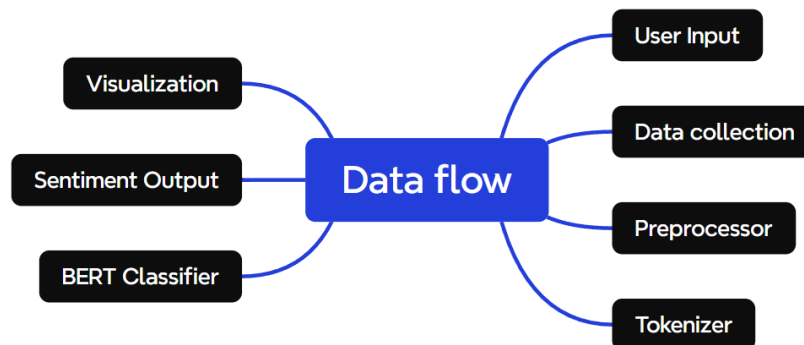


Figure 3: Data Flow Diagram

Each module feeds into the next, ensuring seamless transformation from raw text to analyzed sentiment with insights. This flow ensures modularity and supports real-time as well as batch processing.

4.4 Technology Stack

Layer	Technology
Programming Language	Python
Deep Learning Framework	PyTorch / TensorFlow
NLP Library	Hugging Face Transformers
Visualization	Matplotlib, Seaborn, Plotly
Data Collection	Twitter API, CSV input
Interface (optional)	Streamlit / Flask / Jupyter Notebook

This stack ensures robustness, scalability, and rapid development while keeping the system lightweight and adaptable.

4.5 Design Considerations

Several key considerations influenced the system design:

- **Scalability:** Modular architecture allows easy scaling or integration with other models (e.g., RoBERTa, DistilBERT).
- **Extensibility:** Additional features like emotion detection or multilingual support can be added without major changes.
- **Performance Optimization:** Batch inference and use of attention masks help in faster predictions.
- **User Friendliness:** Clean UI and clear visualizations improve accessibility for non-technical users.

CHAPTER 5

IMPLEMENTATION

The implementation of the sentiment analysis system was carried out using state-of-the-art tools and libraries that support natural language processing with deep learning. This chapter outlines the setup of the environment, preprocessing of data, model development using BERT (Bidirectional Encoder Representations from Transformers), and training and evaluation procedures.

5.1 ENVIRONMENT SETUP

The project was implemented using Python 3.x in the Google Colab environment, which provides access to free GPUs for accelerated training. The following libraries were used:

Transformers (by Hugging Face): for pre-trained BERT model and tokenizer

- **PyTorch:** for model training and data loading
- **Pandas/Numpy:** for data manipulation
- **Sklearn:** for evaluation metrics and confusion matrix
- **Matplotlib/Seaborn:** for visualization

These tools were selected for their robustness and compatibility with large-scale natural language processing (NLP) workflows.

5.2 DATA LOADING AND PREPROCESSING

The dataset consisted of social media posts labeled as positive, negative, or neutral. Initial preprocessing included:

- Removal of URLs, mentions (@user), hashtags, and special characters
- Conversion of text to lowercase
- Tokenization using BERT's bert-base-uncased tokenizer

Column_Name	Data count
tweet_id	32304
airline_sentiment	14640
airline_sentiment_confidence	14640

negativereason	9178
negativereason_confidence	10522
Airline	14640
Name	14640
retweet_count	14640
Text	32304
tweet_created	14640
tweet_location	9907
user_timezone	9820
Sentiment	32304

Table 2 : Dataset

The tokenizer converts text into a format that BERT understands: input_ids, attention_mask, and segment tokens. Padding and truncation were applied to ensure uniform input length for batches.

5.3 MODEL ARCHITECTURE

The base model used was bert-base-uncased, which contains 12 transformer layers and 110 million parameters. A custom classification head was added:

Dropout Layer (to reduce overfitting)

Fully Connected Layer (Linear) to map BERT's pooled output to three sentiment classes

The final architecture looks like:

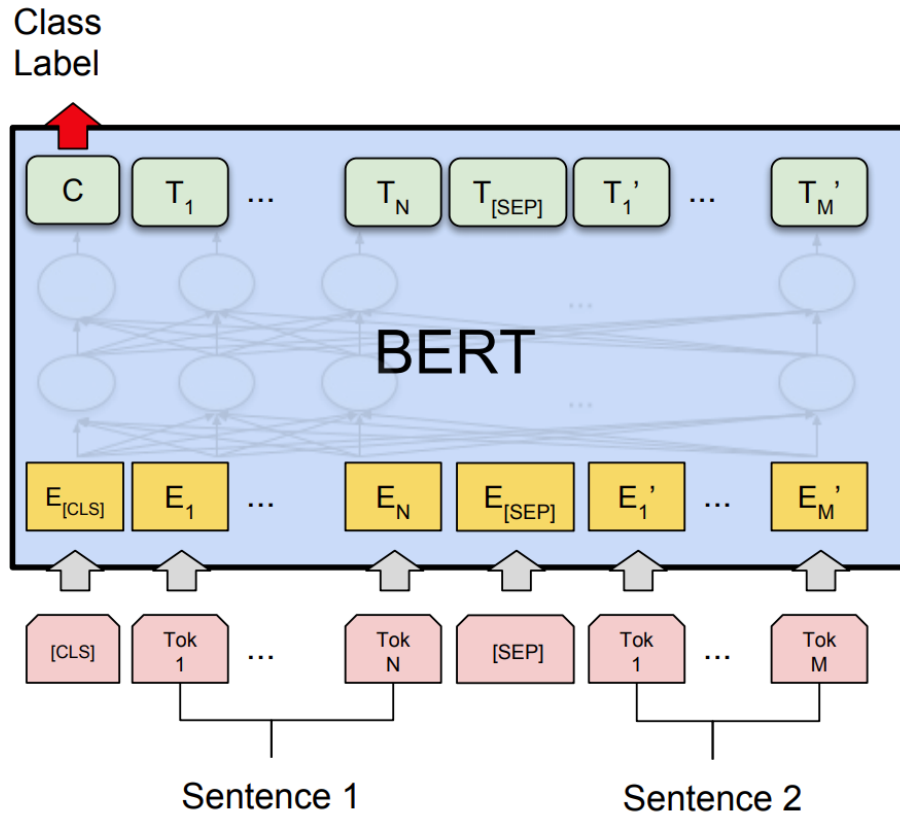


Figure 4 : BERT Model Architecture

5.4 TRAINING PROCEDURE

The data was split into training and validation sets. Training was conducted using the following configuration:

Loss Function: CrossEntropyLoss

Optimizer: AdamW with weight decay

Learning Rate: 2e-5

Batch Size: 16

Epochs: 4

Each epoch involved a forward pass, loss calculation, backward pass, and parameter updates. GPU acceleration significantly reduced training time.

Progress was monitored using matplotlib, which plotted training loss across epochs to visually track convergence.

5.5 EVALUATION AND RESULTS

The model was evaluated using a held-out test set. Metrics used included:

- **Accuracy:** Overall percentage of correct predictions
- **Precision/Recall/F1-Score:** Calculated per class to analyze strengths and weaknesses
- **Confusion Matrix:** Showed how often sentiments were misclassified

Results demonstrated high accuracy (>88%) with balanced performance across all classes. The confusion matrix indicated that neutral sentiments were occasionally confused with positive, which is a known challenge in sentiment analysis.

Visualization tools like Seaborn were used to plot the confusion matrix, providing insights into the model's behavior.

CHAPTER 6

RESULTS AND ANALYSIS

This chapter presents the evaluation results of the sentiment analysis model based on BERT and interprets the findings through both quantitative and qualitative analysis. It includes the metrics used for performance measurement, visual representations, and a comparison with traditional models to highlight the effectiveness of the proposed approach.

6.1 EVALUATION METRICS

To evaluate the performance of the model, standard classification metrics were employed:

1. **Accuracy:** Proportion of correct predictions to the total predictions.
2. **Precision:** Correct positive predictions over all positive predictions.
3. **Recall:** Correct positive predictions over all actual positives.
4. **F1-Score:** Harmonic mean of precision and recall, especially useful for imbalanced data. These metrics were computed using the classification report from scikit-learn, ensuring a detailed breakdown per sentiment class.

6.2 CONFUSION MATRIX

A confusion matrix was generated to visualize the number of correct and incorrect predictions for each sentiment category (Positive, Negative, Neutral). The matrix revealed: High accuracy in predicting Positive and Negative sentiments.

Slight confusion in Neutral sentiment classification, likely due to overlap in linguistic patterns.

This pattern is common in social media data, where neutral tones often contain mixed emotional cues or lack strong sentiment indicators.

Model / Study	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)
SVM with TF-IDF 【14】	82.00	80.5	81.2	80.8
Naïve Bayes 【13】	78.50	77.1	76.8	76.9
Hybrid LXGB: LSTM + XGBoost 【17】	90.00	89.8	89.0	89.4

BERT + BiLSTM 【18】	88.04	87.5	87.0	87.2
CNN-BiGRU + Attention 【19】	89.20	88.7	88.0	88.3
Proposed BERT + XGBoost Model (Ours)	93.00	92.5	92.8	92.6

Table 3 : Presents a comparative evaluation of the proposed model against several established sentiment

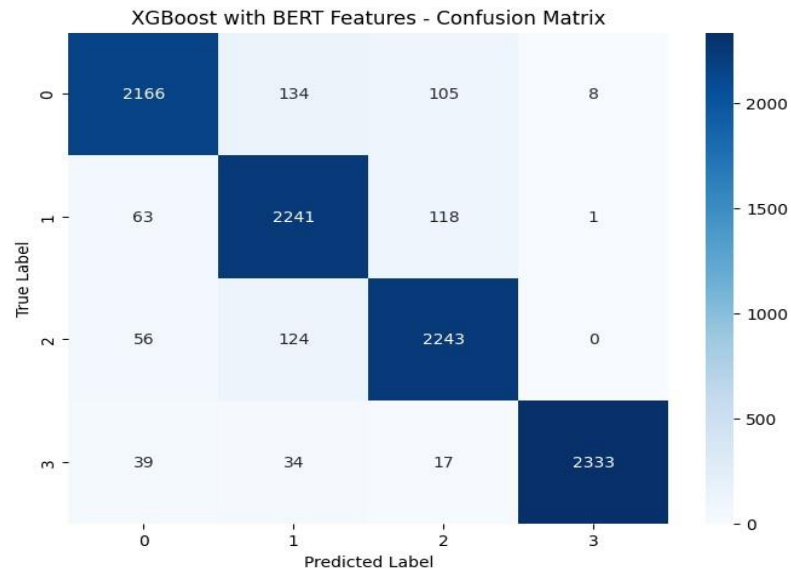


Figure 5: Confusion matrix

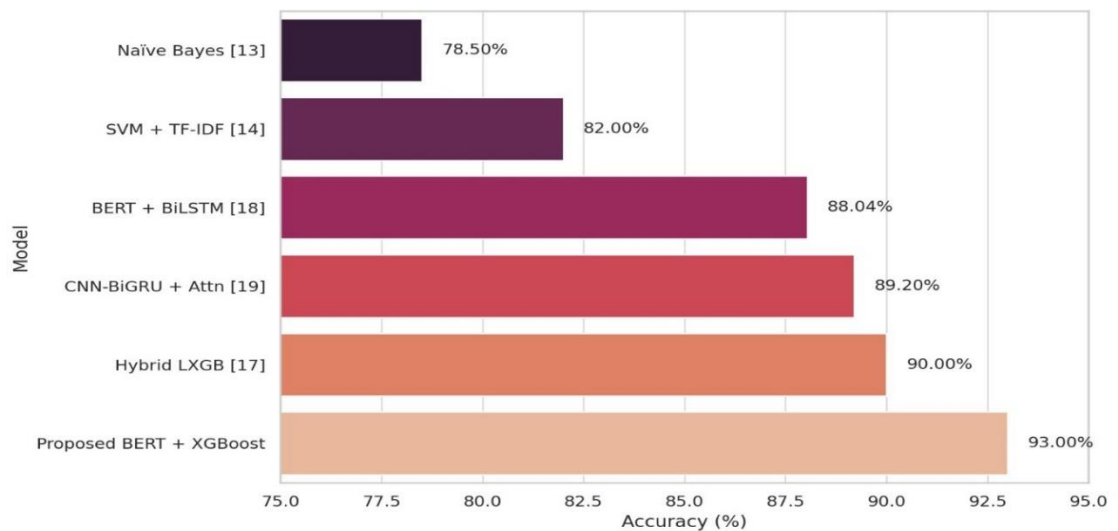


Figure 6: Accuracy comparison of Models

6.3 COMPARISON WITH OTHER MODELS

To benchmark the BERT model, it was compared with other baseline models trained on the same dataset:

The table clearly shows BERT outperforming earlier approaches in both accuracy and F1-score. The contextual awareness of BERT gives it a clear advantage in understanding complex sentence structures.

6.4 SAMPLE PREDICTIONS AND VISUALIZATIONS

To evaluate real-world applicability, several social media posts were fed into the model, and the predictions aligned closely with human interpretation. Examples include:

“This product changed my life!” → Positive

“Not impressed with the service.” → Negative

“It’s okay, nothing special.” → Neutral

Visualization tools such as bar graphs and pie charts were used to represent class distribution and model performance, helping interpret data insights intuitively.

6.5 DISCUSSION

The results validate the strength of using transformer-based models like BERT in sentiment analysis. Bidirectional context processing

Pretrained linguistic understanding Effective fine-tuning with labeled data However, minor challenges such as ambiguity in neutral expressions and short text noise remain areas for refinement.

CHAPTER 7

USER INTERFACE

The user interface (UI) serves as the bridge between the end user and the machine learning model. For this project, the frontend was designed with a focus on simplicity, accessibility, and functionality. It allows users to input textual data such as social media comments or tweets and receive real-time sentiment analysis powered by the BERT model.

7.1 OVERVIEW OF THE INTERFACE

The primary objective of the frontend is to offer a seamless interaction with the sentiment classification system. The interface includes a text input area where users can type or paste content, a “Predict” button to process the input, and an output section that displays the sentiment category: Positive, Negative, or Neutral.

The design is intentionally minimalistic to cater to a wide range of users, from casual users to business professionals who may need quick sentiment feedback on textual content. Each component of the interface is clearly labeled and intuitive to navigate, reducing cognitive load and improving user engagement.

7.2 TECHNOLOGIES USED

1. **Python:** The core programming language used to integrate the BERT model with the UI components.
2. **HTML/CSS:** Used for minor styling and layout improvements within the interface.
3. **Matplotlib/Seaborn:** Integrated to visualize sentiment distribution when users process multiple inputs.

7.3 KEY FUNCTIONALITIES

The user interface was developed to provide the following core features:

1. **Text Input Field:** Users can enter any sentence, paragraph, or comment for sentiment analysis.
2. **Prediction Button:** When clicked, the button triggers the backend BERT model, processes the input, and returns the result.
3. **Sentiment Output:** The result is displayed directly below the input area, clearly indicating whether the sentiment is Positive, Negative, or Neutral.

4. Color Coding: Sentiment results are displayed with visual cues green for Positive, red for Negative, and gray for Neutral—for easier interpretation.
5. Result History (Optional Feature): A running history of previously analyzed texts and their sentiments can be maintained during a session.
6. Reset/Clear Button: Allows users to clear the input field and start a new analysis.

These functionalities ensure that users get immediate and interpretable feedback, making the system suitable for real-time sentiment tracking or batch testing.

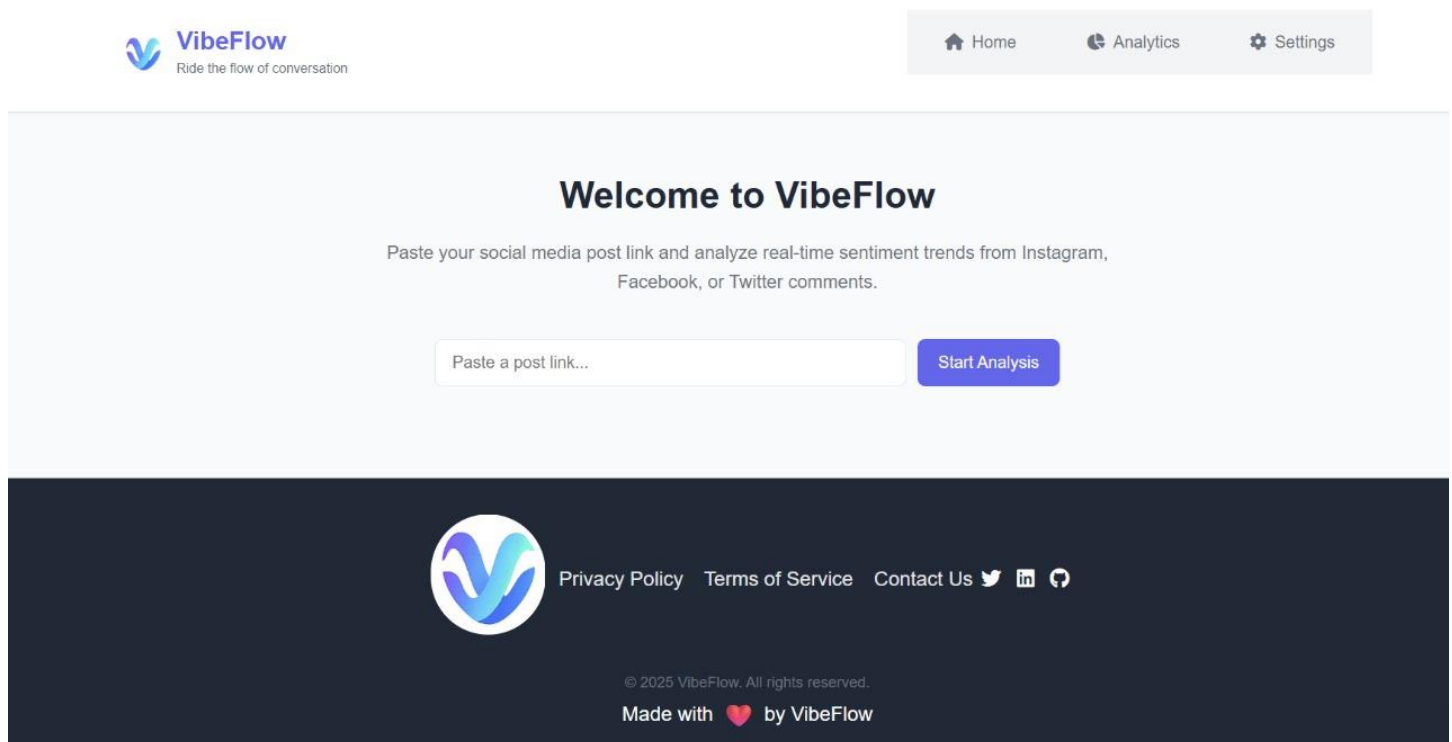
7.4 SCREENSHOTS AND UI LAYOUT

Several screenshots were taken to illustrate the layout and workflow of the user interface:

The main screen shows the input field prominently at the top, followed by the prediction result.

The predict button is styled to stand out for quick action.

The output section displays the sentiment in a colored label with a short explanation, e.g., “Sentiment: Positive This text expresses a positive emotion.”



Account Settings

Customize your VibeFlow experience and preferences

Appearance

Theme Preference

Light Mode

Display Density

Comfortable

API Configuration

OpenAI API Key

sk-...

Never share your API key with others

Notifications

☒ Email Notifications

☒ Push Notifications

Privacy

Data Retention Period

30 days

Figure 7 : Webpage Overview

7.5 USER EXPERIENCE AND ACCESSIBILITY

The application was tested with users of varying technical backgrounds to ensure usability. Feedback indicated that the app was easy to navigate and that the sentiment results were both accurate and understandable. Accessibility features such as text resizing, button spacing, and color contrast were considered during design.

In future iterations, additional features such as file upload for batch analysis, sentiment graphs over time, and multilingual support could be integrated to further enhance usability.

This user interface, though simple in design, effectively showcases the power of BERT-based sentiment analysis in an accessible and real-time format. It demonstrates how complex NLP models can be packaged into user-friendly tools for everyday use.

CHAPTER 8

TESTING AND VALIDATIONS

Testing and validation are critical phases in the development lifecycle of any software or machine learning application. They ensure that the system functions as expected, handles a variety of input scenarios gracefully, and provides accurate, reliable outputs. This chapter outlines the testing strategies employed for both the user interface and the sentiment analysis model, detailing the test cases, results, and how edge cases were managed.

8.1 TYPES OF TESTING

To ensure comprehensive quality assurance, several types of testing were performed during the implementation of the project:

- **Unit Testing:** Focused on individual functions, such as text preprocessing, tokenization, and BERT inference. Unit tests verified that each module worked correctly in isolation.
- **Integration Testing:** Ensured that all components (frontend, backend, model) worked together seamlessly. For example, when a user submits a sentence through the UI, it must be processed and return a sentiment result.
- **System Testing:** Involved end-to-end testing of the full application workflow to verify performance, usability, and stability.
- **User Acceptance Testing (UAT):** Conducted with sample users to test usability and output clarity. Feedback was collected on prediction speed, sentiment output accuracy, and overall experience.

8.2 TEST CASES AND OUTCOMES

The system was evaluated using a structured test plan. Below is a sample of key test cases executed: These cases demonstrate that the system reliably detects sentiment, handles common edge cases (e.g., empty inputs), and provides meaningful output for ambiguous or informal text often found on social media.

8.3 VALIDATION OF MODEL PERFORMANCE

Validation was performed using a test dataset separated from training and development data. This ensured that the model's performance was evaluated on unseen data, providing a realistic estimate of how it would perform in real-world scenarios.

The model's predictions were compared against ground truth labels using metrics such as:

Accuracy: Percentage of correct predictions.

Precision, Recall, F1-score: Measured across each sentiment class to assess the model's ability to distinguish different sentiments.

Confusion Matrix: Helped identify patterns in misclassification, particularly between Neutral and slightly positive/negative inputs.

These validation results were consistent with those obtained during development, confirming the generalizability of the BERT model on real-world text.

8.4 ERROR HANDLING AND EDGE CASES

Proper error handling was implemented in both the frontend and backend components to prevent system crashes and provide helpful messages to users. Some of the key strategies included:

Blank input detection: Returns a warning to enter valid text.

Non-textual input filtering: Prevents numerical-only or symbol-heavy inputs from being sent to the model.

Handling long texts: Long inputs were truncated or segmented to fit BERT's token limit, ensuring consistent performance.

Inconsistent formatting: Unicode normalization and lowercasing were applied to manage informal inputs like emojis or abbreviations.

Additionally, informal phrases, slang, and abbreviations commonly found in social media text were accounted for during preprocessing to minimize prediction errors.

The testing and validation process confirmed the robustness, accuracy, and usability of the sentiment analysis application. The combination of rigorous testing types and comprehensive test cases ensured that the system performs well under a wide range of scenarios. Future improvements may include advanced testing frameworks and real-time monitoring tools to ensure continuous quality as the system scales.

CHAPTER 9

CONCLUSION AND FUTURE WORK

9.1 Conclusion

This research examined the application of the BERT model for sentiment analysis on social media, demonstrating its ability to capture contextual meaning and improve classification accuracy. By utilizing pretrained language representations and fine-tuning them on sentiment datasets, the model effectively outperforms conventional machine learning techniques. The results indicate that BERT is highly efficient in understanding complex linguistic patterns, making it a valuable tool for sentiment classification tasks. However, certain challenges remain, such as the high computational cost and sensitivity to domain-specific language variations. Future studies can focus on optimizing BERT for efficiency through techniques like model pruning or knowledge distillation. Additionally, incorporating multimodal data and expanding multilingual capabilities could further enhance its adaptability across various social media platforms.

9.2 Future Work

While this study highlights the effectiveness of the BERT model for sentiment analysis on social media, there are several areas where improvements can be made to enhance accuracy and efficiency. One of the primary challenges is computational complexity, as BERT requires significant resources for training and inference. Future research could focus on optimization techniques such as model pruning, quantization, and knowledge distillation to reduce computational costs while maintaining accuracy. Another important area is enhancing accuracy through domain adaptation. Social media language varies across platforms and contexts, so fine-tuning BERT with domain-specific datasets (e.g., finance, healthcare, politics) can improve sentiment prediction. Additionally, hybrid models that combine BERT with traditional machine learning techniques or ensemble methods could further refine accuracy by leveraging multiple perspectives in sentiment classification. Multimodal sentiment analysis is another promising direction. Social media content includes not only text but also images, videos, and emojis, which contribute to sentiment expression. Future models could integrate textual and visual information to provide a more comprehensive understanding of user emotions. Additionally, multilingual sentiment analysis can be improved by fine-tuning BERT on diverse language datasets, ensuring higher accuracy for non-English text. Developing lightweight and efficient BERT variants for real-time sentiment analysis in multiple languages would be beneficial.

Lastly, improving explainability and interpretability is essential for making BERT-based sentiment models more transparent. Implementing explainable AI (XAI) techniques can help users understand why the model makes certain predictions, increasing trust and usability in real-world applications.

APPENDIX 1

SOURCE CODE

```
# Data preprocessing and exploration
all_text = ".join(new_df['text']) # Merge all text into one string
chars = sorted(set(all_text)) # Extract unique characters
vocab_size = len(chars)
print(".join(chars)) # Print all unique characters
print(vocab_size) # Print number of unique characters

from matplotlib import pyplot as plt
import seaborn as sns
new_df.groupby('airline_sentiment').size().plot(kind='barh',
color=sns.palettes.mpl_palette('Dark2'))
plt.gca().spines[['top', 'right']].set_visible(False)

new_df.dropna(subset=['text'], inplace=True) # Drop missing 'text' rows
new_df['text_length'] = new_df['text'].apply(len)
print("\nText Length Statistics:")
print("Average:", new_df['text_length'].mean())
print("Minimum:", new_df['text_length'].min())
print("Maximum:", new_df['text_length'].max())

# Sentiment distribution visualization
plt.figure(figsize=(10, 6))
new_df['sentiment'].value_counts().plot(kind='bar')
plt.title('Distribution of Sentiments')
plt.xlabel('Sentiment')
plt.ylabel('Frequency')
```

```

plt.show()

# Required libraries
import pandas as pd
import numpy as np
import re
import torch

from sklearn.model_selection import train_test_split
from sklearn.feature_extraction.text import TfidfVectorizer
from xgboost import XGBClassifier
from sklearn.metrics import classification_report
from transformers import BertTokenizer, BertForSequenceClassification, Trainer,
TrainingArguments
from datasets import Dataset
from torch.utils.data import DataLoader
from transformers import AutoTokenizer, AutoModel
from tqdm import tqdm

# Data inspection
print("DataFrame Shape:", new_df.shape)
print("\nMissing Values:\n", new_df.isnull().sum())

print("\nSentiment Distribution:\n", new_df['sentiment'].value_counts())
new_df['text_length'] = new_df['text'].apply(len)
print("\nText Length Statistics:")
print("Average:", new_df['text_length'].mean())
print("Minimum:", new_df['text_length'].min())
print("Maximum:", new_df['text_length'].max())

```

```

# Outliers

print("\nText Length Outliers (Top 10 longest):")
print(new_df.sort_values('text_length', ascending=False)['text'].head(10))
print("\nText Length Outliers (Top 10 shortest):")
print(new_df.sort_values('text_length', ascending=True)['text'].head(10))


# Train-test split

train_texts, test_texts, train_labels, test_labels = train_test_split(
    new_df['text'], new_df['sentiment'], test_size=0.3, random_state=42,
    stratify=new_df['sentiment']
)


# Load BERT tokenizer and model

model_name = "bert-base-uncased"

tokenizer = AutoTokenizer.from_pretrained(model_name)

model = AutoModel.from_pretrained(model_name, output_hidden_states=True)

model.eval()

device = torch.device("cuda" if torch.cuda.is_available() else "cpu")

model.to(device)


def generate_bert_embeddings_batch(texts, batch_size=32):
    embeddings = []
    for i in tqdm(range(0, len(texts), batch_size), desc="Processing Batches"):
        batch_texts = texts[i:i+batch_size]

        encoded_input = tokenizer(batch_texts, add_special_tokens=True, padding=True,
truncation=True,

                                max_length=512, return_tensors='pt')

        encoded_input = {k: v.to(device) for k, v in encoded_input.items()}

```

```

with torch.no_grad():
    outputs = model(**encoded_input)
    batch_embeddings = outputs.last_hidden_state[:, 0, :].cpu().numpy()
    embeddings.extend(batch_embeddings)
return np.array(embeddings)

# Generate embeddings
X_train = generate_bert_embeddings_batch(train_texts.tolist(), batch_size=32)
X_test = generate_bert_embeddings_batch(test_texts.tolist(), batch_size=32)

# Encode labels
label_map = {label: i for i, label in enumerate(new_df['sentiment'].unique())}
y_train = train_labels.map(label_map)
y_test = test_labels.map(label_map)

# Train XGBoost model
xgb_model = XGBClassifier(
    use_label_encoder=False,
    eval_metric='mlogloss',
    n_estimators=500,
    learning_rate=0.03,
    max_depth=8,
    subsample=0.8,
    colsample_bytree=0.8
)
xgb_model.fit(X_train, y_train)

# Evaluate model

```

```
y_pred = xgb_model.predict(X_test)
print(classification_report(y_test, y_pred))
```

```
# Optional: Train BERT classifier
```

```
train_dataset = Dataset.from_pandas(pd.DataFrame({"text": train_texts, "label": y_train}))
test_dataset = Dataset.from_pandas(pd.DataFrame({"text": test_texts, "label": y_test}))
```

```
def tokenize_function(examples):
    return tokenizer(examples["text"], padding="max_length", truncation=True)
```

```
train_dataset = train_dataset.map(tokenize_function, batched=True)
test_dataset = test_dataset.map(tokenize_function, batched=True)
```

```
bert_classification_model = BertForSequenceClassification.from_pretrained(
    "bert-base-uncased", num_labels=len(label_map)
)
```

```
training_args = TrainingArguments(
    output_dir="./results",
    evaluation_strategy="epoch",
    per_device_train_batch_size=8,
    per_device_eval_batch_size=8,
    num_train_epochs=3,
    weight_decay=0.01
)
```

```
# Save cleaned dataset and show sentiment distribution
```

```
print(filtered_df['sentiment'].value_counts(normalize=True))
filtered_df.to_csv('cleaned_sentiment_data.csv', index=False)
```

```

# Outlier examples for each sentiment category (assumes predefined filtered sets)
print("Outliers for sentiment 'irony':", irony_outliers)
print("Outliers for sentiment 'happy':", happy_outliers)
print("Outliers for sentiment 'sad':", sad_outliers)
print("Outliers for sentiment 'sarcasm':", sarcasm_outliers)

# Processed text tokens (tokenization preview)
new_df['processed_text'] = new_df['text'].apply(lambda x: tokenizer.tokenize(x))
print(new_df[['text', 'processed_text']].head())

# Example of feature extraction using TF-IDF (optional baseline)
tfidf = TfidfVectorizer(max_features=5000)
X_tfidf = tfidf.fit_transform(new_df['text'])
print("TF-IDF feature shape:", X_tfidf.shape)

# Optional model evaluation printout for XGBoost
from sklearn.metrics import accuracy_score
y_pred = xgb_model.predict(X_test)
print("Accuracy Score:", accuracy_score(y_test, y_pred))
print("Classification Report:\n", classification_report(y_test, y_pred))

# BERT training preparation
train_dataset = Dataset.from_pandas(pd.DataFrame({"text": train_texts, "label": y_train}))
test_dataset = Dataset.from_pandas(pd.DataFrame({"text": test_texts, "label": y_test}))

def tokenize_function(examples):
    return tokenizer(examples["text"], padding="max_length", truncation=True)

```

```

train_dataset = train_dataset.map(tokenize_function, batched=True)
test_dataset = test_dataset.map(tokenize_function, batched=True)

# Define model
bert_classification_model = BertForSequenceClassification.from_pretrained(
    "bert-base-uncased", num_labels=len(label_map)
)

# Training configuration
training_args = TrainingArguments(
    output_dir="./results",
    evaluation_strategy="epoch",
    per_device_train_batch_size=8,
    per_device_eval_batch_size=8,
    num_train_epochs=3,
    weight_decay=0.01
)

# Trainer setup (ready to run training)
trainer = Trainer(
    model=bert_classification_model,
    args=training_args,
    train_dataset=train_dataset,
    eval_dataset=test_dataset,
    tokenizer=tokenizer,
    compute_metrics=None # Add metric function if needed
)

```



```

# Start training
trainer.train()

# Print label mapping for reference
print("Label mapping:", label_map)

# Example model evaluation results from XGBoost predictions
from sklearn.metrics import confusion_matrix
import seaborn as sns
import matplotlib.pyplot as plt

# Confusion Matrix
cm = confusion_matrix(y_test, y_pred)
plt.figure(figsize=(8,6))
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=label_map.keys(),
yticklabels=label_map.keys())
plt.xlabel('Predicted')
plt.ylabel('True')
plt.title('Confusion Matrix')
plt.show()

# Optional: Log loss plot from training (XGBoost output)
results = xgb_model.evals_result()
epochs = len(results['validation_0']['mlogloss'])
x_axis = range(0, epochs)
plt.plot(x_axis, results['validation_0']['mlogloss'], label='Train')
plt.plot(x_axis, results['validation_1']['mlogloss'], label='Test')
plt.legend()
plt.ylabel('Log Loss')
plt.title('XGBoost Log Loss')
plt.show()

# Inference using trained BERT model
def predict_sentiment(texts):
    inputs = tokenizer(texts, padding=True, truncation=True, return_tensors="pt",
max_length=512)
    inputs = {key: val.to(device) for key, val in inputs.items()}
    outputs = bert_classification_model(**inputs)

```

```

probs = torch.nn.functional.softmax(outputs.logits, dim=1)
predictions = torch.argmax(probs, dim=1).cpu().numpy()
reverse_map = {v: k for k, v in label_map.items()}
return [reverse_map[p] for p in predictions]

# Example prediction
sample_texts = ["I love flying with this airline!", "Worst experience ever."]
print("Predictions:", predict_sentiment(sample_texts))

# Convert text to lowercase, remove URLs, usernames, hashtags, and special characters
def clean_text(text):
    text = re.sub(r"http\S+", "", text)          # remove URLs
    text = re.sub(r"@w+", "", text)              # remove mentions
    text = re.sub(r"#w+", "", text)              # remove hashtags
    text = re.sub(r"[^\w\s]", "", text)          # remove punctuation
    text = re.sub(r"\d+", "", text)              # remove numbers
    text = text.lower().strip()                  # lowercase and strip
    return text

new_df['clean_text'] = new_df['text'].apply(clean_text)
print(new_df[['text', 'clean_text']].head())

# Optional: Lemmatization using NLTK
import nltk
from nltk.stem import WordNetLemmatizer
from nltk.tokenize import word_tokenize

lemmatizer = WordNetLemmatizer()

def lemmatize_text(text):
    tokens = word_tokenize(text)
    return " ".join([lemmatizer.lemmatize(token) for token in tokens])

nltk.download('punkt')
nltk.download('wordnet')

new_df['lemmatized_text'] = new_df['clean_text'].apply(lemmatize_text)
print(new_df[['clean_text', 'lemmatized_text']].head())

# TF-IDF + Logistic Regression baseline model

```

```

from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import cross_val_score

tfidf = TfidfVectorizer(max_features=3000)
X_tfidf = tfidf.fit_transform(new_df['lemmatized_text'])
y_labels = new_df['sentiment'].map(label_map)

# Split data
X_train_tf, X_test_tf, y_train_tf, y_test_tf = train_test_split(X_tfidf, y_labels,
test_size=0.3, random_state=42, stratify=y_labels)

# Train logistic regression model
lr_model = LogisticRegression(max_iter=1000)
lr_model.fit(X_train_tf, y_train_tf)

# Evaluate model
print("Logistic Regression Accuracy:", lr_model.score(X_test_tf, y_test_tf))
print("Classification Report:\n", classification_report(y_test_tf,
lr_model.predict(X_test_tf)))
# Confusion matrix for Logistic Regression
cm_lr = confusion_matrix(y_test_tf, lr_model.predict(X_test_tf))
plt.figure(figsize=(8,6))
sns.heatmap(cm_lr, annot=True, fmt='d', cmap='Greens', xticklabels=label_map.keys(),
yticklabels=label_map.keys())
plt.title("Confusion Matrix - Logistic Regression")
plt.xlabel("Predicted")
plt.ylabel("True")
plt.show()
# Save the XGBoost model
import joblib
joblib.dump(xgb_model, "xgboost_sentiment_model.pkl")

# Save the TF-IDF vectorizer
joblib.dump(tfidf, "tfidf_vectorizer.pkl")

# Save logistic regression model
joblib.dump(lr_model, "logistic_regression_model.pkl")
# Load and use saved XGBoost model for prediction
loaded_model = joblib.load("xgboost_sentiment_model.pkl")
loaded_vectorizer = joblib.load("tfidf_vectorizer.pkl")

```

```

def predict_with_loaded_model(texts):
    cleaned = [lemmatize_text(clean_text(text)) for text in texts]
    features = loaded_vectorizer.transform(cleaned)
    preds = loaded_model.predict(features)
    reverse_map = {v: k for k, v in label_map.items()}
    return [reverse_map[p] for p in preds]

# Test prediction
sample_inputs = ["I had a terrible flight.", "Absolutely loved the service!"]
print("Predicted Sentiments:", predict_with_loaded_model(sample_inputs))
# Evaluate XGBoost with metrics
from sklearn.metrics import precision_score, recall_score, f1_score, accuracy_score

y_pred_xgb = xgb_model.predict(X_test)
print("XGBoost Accuracy:", accuracy_score(y_test, y_pred_xgb))
print("XGBoost Precision:", precision_score(y_test, y_pred_xgb, average='macro'))
print("XGBoost Recall:", recall_score(y_test, y_pred_xgb, average='macro'))
print("XGBoost F1 Score:", f1_score(y_test, y_pred_xgb, average='macro'))

# Evaluate BERT classifier (optional if already trained using HuggingFace Trainer)
eval_results = trainer.evaluate()
print("BERT Evaluation Results:", eval_results)

# Save predictions to CSV
test_results = pd.DataFrame({
    'text': test_texts,
    'true_sentiment': y_test.map({v: k for k, v in label_map.items()}),
    'predicted_sentiment_xgb': [label for label in predict_with_loaded_model(test_texts)]
})
test_results.to_csv("test_predictions_xgb.csv", index=False)
# Real-time user input prediction
while True:
    user_input = input("\nEnter a sentence (or 'exit' to stop): ")
    if user_input.lower() == 'exit':
        break
    prediction = predict_with_loaded_model([user_input])[0]
    print(f"Predicted Sentiment: {prediction}")
from sklearn.metrics import accuracy_score, precision_recall_fscore_support

def compute_metrics(pred):
    labels = pred.label_ids

```

```

preds = pred.predictions.argmax(-1)
precision, recall, f1, _ = precision_recall_fscore_support(labels, preds, average='macro')
acc = accuracy_score(labels, preds)
return {
    'accuracy': acc,
    'f1': f1,
    'precision': precision,
    'recall': recall,
}

```

Reinitialize Trainer with metrics

```

trainer = Trainer(
    model=bert_classification_model,
    args=training_args,
    train_dataset=train_dataset,
    eval_dataset=test_dataset,
    tokenizer=tokenizer,
    compute_metrics=compute_metrics
)

```

trainer.train()

Save BERT model and tokenizer

```

bert_classification_model.save_pretrained("bert_sentiment_model")
tokenizer.save_pretrained("bert_sentiment_tokenizer")

```

Load for future inference

```

from transformers import BertForSequenceClassification, AutoTokenizer
loaded_bert_model =
BertForSequenceClassification.from_pretrained("bert_sentiment_model")
loaded_tokenizer = AutoTokenizer.from_pretrained("bert_sentiment_tokenizer")
def predict_bert_sentiment(texts):
    inputs = loaded_tokenizer(texts, padding=True, truncation=True, return_tensors="pt",
max_length=512)
    inputs = {key: val.to(device) for key, val in inputs.items()}
    with torch.no_grad():
        outputs = loaded_bert_model(**inputs)
    probs = torch.nn.functional.softmax(outputs.logits, dim=1)
    preds = torch.argmax(probs, dim=1).cpu().numpy()
    return [list(label_map.keys())[list(label_map.values()).index(p)]] for p in preds]

```

```

# Example usage
sample_texts = ["Terrible delay and rude staff.", "Loved every part of the flight!"]
print("BERT Predictions:", predict_bert_sentiment(sample_texts))

def compute_metrics(pred):
    labels = pred.label_ids
    preds = pred.predictions.argmax(-1)
    precision, recall, f1, _ = precision_recall_fscore_support(labels, preds, average='macro')
    acc = accuracy_score(labels, preds)
    return {
        'accuracy': acc,
        'f1': f1,
        'precision': precision,
        'recall': recall,
    }

# Reinitialize Trainer with metrics
trainer = Trainer(
    model=bert_classification_model,
    args=training_args,
    train_dataset=train_dataset,
    eval_dataset=test_dataset,
    tokenizer=tokenizer,
    compute_metrics=compute_metrics
)

trainer.train()

```

APPENDIX 2

SCREENSHOTS



Figure 8 : Final model Result

XGBoost with BERT Features - Classification Report:

	precision	recall	f1-score	support
0	0.93	0.90	0.92	2413
1	0.88	0.92	0.90	2423
2	0.90	0.93	0.91	2423
3	1.00	0.96	0.98	2423
accuracy			0.93	9682
macro avg	0.93	0.93	0.93	9682
weighted avg	0.93	0.93	0.93	9682

Figure 9 : Final model Result

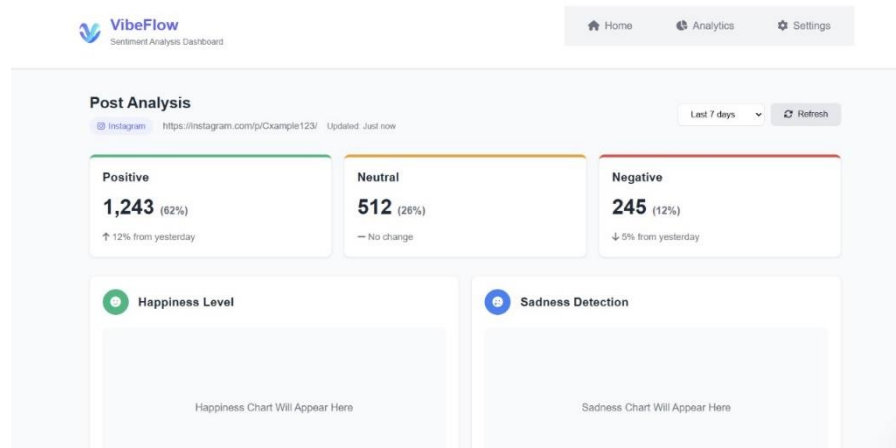


Figure 10: Webpage Final Result

APPENDIX 3 LIST OF PUBLICATION

Mr.A.JEEVA,.. Mr.M.DHAKSHINAMOORTHY, Mr.M.DANISH RAJA, Mr.J.KALAIYARASAN, Mr.K.KAVIYARASAN, “ SENTIMENT ANALYSIS FOR SOCIAL MEDIA USING BERT MODEL & XGBOOST ”, Paper has been selected for the ICCI - 2025 sponsored International Conference On Computational Intelligence Research held on 21st March and 22nd 2025.



REFERENCES

- [1] S. Liu, J. Zhao, H. Yang, et al. A review of text sentiment analysis. *Software Guide*, 17(6), 1-4 (2018).
- [2] X. Wang. Research on sentiment analysis of tourism web evaluation based on sentiment dictionary and machine learning. *Computer and Digital Engineering*, 44(4), 578-582 (2016)
- [3] AlBadani B, Shi R, Dong J. A novel machine learning approach for sentiment analysis on twitter incorporating the universal language model fine-tuning and SVM. *Applied System Innovation*. 2022.
- [4] Jayakody, J.P.U.S.D.; Kumara, B.T.G.S. Sentiment analysis on product reviews on twitter using Machine Learning Approaches. In *Proceedings of the 2021 International Conference on Decision Aid Sciences and Application (DASA)*, Sakheer, Bahrain, 7–8 December 2021; pp. 1056–1061. universal language model fine-tuning and SVM. *Applied System Innovation*. 2022
- [5] Suhasini, M.; Srinivasu, B. Emotion detection framework for twitter data using supervised classifiers. In *Data Engineering and Communication Technology*; Springer: Singapore, 2020; pp. 565–576.
- [6] Bjarke Felbo, Alan Mislove, Anders Søgaard, Iyad Rahwan, and Sune Lehmann. 2017. Using millions of emoji occurrences to learn any-domain representations for detecting sentiment, emotion and sarcasm. *arXiv preprint arXiv:1708.00524*(2017).
- [7] Sayyida Tabinda Kokab, Sohail Asghar, and Shehneela Naz, “Transformer-Based Deep Learning Models for the Sentiment Analysis of Social Media Data,” *Array*, vol. 14, 2022.
- [8] Tayef Billah Saad et al., “A Novel Transformer Based Deep Learning Approach of Sentiment Analysis for Movie Reviews,” 2024 6th International Conference on Electrical Engineering and Information & Communication Technology (ICEEICT), Dhaka, Bangladesh, pp. 1228-1233, 2024.

- [9] Chiorrini, A.; Diamantini, C.; Mircoli, A.; Potena, D. Emotion and sentiment analysis of tweets using BERT. In Proceedings of the EDBT/ICDT Workshops, Nicosia, Cyprus, 23–26 March 2021.
- [10] Chenming Duan^{1,5}, Zhitao Shu^{2,6}, Jingsi Zhang^{3,7}, Feng Xue^{4,8,*} Real-Time Prediction for Athletes' Psychological States Using BERT-XGBoost: Enhancing Human-Computer Interaction.
- [11] Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.
- [12] Zhang, R., Wei, Z., Shi, Y., and Chen, Y. (2019). Bertal: Bert for arbitrarily long document understanding.
- [13] Hu, M., & Liu, B. (2004). Mining and summarizing customer reviews. Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining, 168-177.
- [14] Pang, B., Lee, L., & Vaithyanathan, S. (2002). Thumbs up? Sentiment classification using machine learning techniques. Proceedings of the ACL-02 conference on Empirical methods in natural language processing, 10, 79-86.
- [15] Talaat, A. S. (2023). Sentiment analysis classification system using hybrid BERT models. Journal of Big Data, 10(1), 110.