# Advancing Sentiment Analysis through Machine and Deep Learning

*Submitted by*

**Nguyen Tung Lam**   **V202100571**

**Nguyen Nhat Minh**   **V202100570**

**Nguyen Mau Hoang Hiep**   **V202100595**

Under the supervision of

**Instructor's Name: Prof.Doan Dang Khoa**
**Related course: COMP3020 - Machine Learning**

**VINUNIVERSITY**

**COLLEGE OF ENGINEERING AND COMPUTER SCIENCE**

**VINUNIVERSITY**

**May 2024**

## ABBREVIATIONS

BERT            Bidirectional Encoder Representations from Transformers

SVM             Support Vector Machine

LSTM            Bidirectional Long Short-Term Memory

RoBERTa         Robustly optimized BERT approach

ALBERT          A Lite BERT

# Contents

# Chapter 1

# INTRODUCTION

## 1.1   Project Background

- **Introduction**

Sentiment analysis, also known as opinion mining, is a critical area of study within natural language processing (NLP) that focuses on identifying and categorizing opinions expressed in text. This project aims to advance sentiment analysis techniques by leveraging both machine learning and deep learning models. Over the past few years, sentiment analysis has seen significant growth, driven by the increasing volume of textual data from social media, reviews, and other online platforms.

- **Literature Review**

In recent years, sentiment analysis has witnessed notable advancements driven by innovations in machine learning algorithms, deep learning architectures, and the availability of large-scale annotated datasets. Machine learning techniques such as Support Vector Machines (SVM), Naive Bayes, and more recently, deep learning models like Recurrent Neural Networks (RNNs) and Transformer-based architectures, have demonstrated remarkable performance in sentiment classification tasks (Gohil et al., 2021). These models leverage contextual embeddings and semantic representations to capture complex linguistic patterns and nuances in sentiment expression.

Traditional sentiment analysis methods, such as Support Vector Machines (SVM) and logistic regression, have provided a foundation for text classification tasks. However, the advent of deep learning has introduced more sophisticated models like Long Short-Term Memory (LSTM) networks, BERT (Bidirectional Encoder Representations from Transformers), and its variants (RoBERTa and ALBERT). These models have demonstrated superior performance in capturing context and understanding nuanced language patterns compared to their predecessors.

Moreover, aspect-based sentiment analysis (ABSA) advancements have facilitated more fine-grained analyses by identifying sentiments toward specific aspects or entities within a text (Liu et al., 2020). ABSA techniques enable a deeper understanding of sentiment dynamics, particularly in domains such as product reviews, where identifying sentiment towards individual product features is crucial for product development and marketing strategies.

- **Current State and Limitations**

While models like BERT and RoBERTa have significantly improved sentiment analysis accuracy, they are not without limitations. These models often require extensive computational resources and large labeled datasets for training.

Despite significant progress, sentiment analysis encounters challenges that hinder its widespread adoption and effectiveness. One major challenge lies in the inherent subjectivity and ambiguity of human language, including sarcasm, irony, and cultural nuances, which can confound sentiment analysis models (Pang & Lee, 2008). Additionally, the domain-dependency of sentiment expressions poses challenges in model generalization across different domains and topics.

Furthermore, the presence of biases in labeled datasets and the potential for algorithmic bias exacerbate the reliability and fairness of sentiment analysis outcomes (Kiritchenko & Mohammad, 2018). Biases in training data, such as underrepresentation of certain demographics or overemphasis on specific sentiments, can lead to skewed predictions and perpetuate societal inequalities.

- **Transition**

Given these challenges and the current landscape of sentiment analysis, this project aims to explore and enhance the capabilities of existing models by addressing their limitations. By leveraging advanced machine learning and deep learning techniques, we intend to develop more accurate, efficient, and robust sentiment analysis models that can be applied across various domains and languages.

## 1.2 Project Definition

### 1.2.1 Problem Statement

The primary problem addressed by this project is the enhancement of sentiment analysis models to improve their accuracy, efficiency, and robustness. The current models, while effective to some extent, face limitations in handling nuanced language, require extensive labeled datasets and are computationally demanding. This project seeks to overcome these barriers by developing and optimizing models that can accurately interpret sentiments with fewer resources and greater contextual understanding.

### 1.2.2 Context and Scope

The scope of this project encompasses the entire sentiment analysis pipeline, from data collection and preprocessing to model development, training, and evaluation. The project involves the use of both traditional machine learning algorithms and advanced deep learning architectures. Data will be sourced from a variety of platforms, including social media, customer reviews, and news articles, to ensure a diverse and comprehensive dataset. Preprocessing steps will include text cleaning, normalization, tokenization, and feature extraction. The development phase will focus on implementing and fine-tuning models like SVM, logistic regression, BERT, RoBERTa, and LSTM with attention mechanisms. The final phase will involve rigorous testing and validation using robust evaluation metrics to ensure the models' reliability and generalizability.

### 1.2.3 Significance and Implications

Enhancing sentiment analysis models is crucial for numerous applications. Businesses can leverage improved sentiment analysis to better understand customer feedback, monitor brand reputation, and inform marketing strategies. In the public sector, sentiment analysis can help gauge public opinion on policies and events. Accurate sentiment analysis also has significant implications for market research, enabling companies to predict trends and consumer behavior more effectively. The failure to address current model limitations could result in missed opportunities, suboptimal decision-making, and potential biases in sentiment interpretation, which underscores the importance of this project's objectives.

### 1.2.4 Quantification (if any)

The magnitude of the problem is illustrated by the current performance benchmarks of sentiment analysis models. Many state-of-the-art models achieve accuracy rates in the range of 70-80%, but they often fall short in handling complex and nuanced language constructs. This project aims to improve accuracy by at least 10% over these baseline models. Additionally, efficiency metrics such as inference time and computational resource usage will be optimized, with a target of reducing processing time by 20% without compromising accuracy.

## 1.3 Project Objectives

The objectives of this project are to:

- **Objective 1: Develop and Implement Advanced Sentiment Analysis Models**

    - **Objective:** Enhance model accuracy and contextual understanding.

    - **Approach:** Utilize and compare various models, including traditional machine learning algorithms (e.g., SVM, logistic regression) and advanced deep learning architectures (e.g., BERT, RoBERTa, ALBERT). The focus will be on optimizing these models for better performance and efficiency.

    - **Impact:** Improved accuracy and robustness of sentiment analysis, enabling more reliable sentiment detection across different contexts and domains.

- **Objective 2: Enhance Data Preprocessing and Feature Engineering Techniques**

    - **Objective:** Improve the quality of model input data.

    - **Approach:** Implement effective data cleaning, normalization, and feature extraction methods to enhance the relevance and quality of input data. Innovative feature engineering approaches, such as using word embeddings and capturing syntactic dependencies, will be explored to capture relevant textual features.

    - **Impact:** Enhanced model performance due to better-quality inputs, leading to more accurate sentiment predictions.

- **Objective 3: Conduct Comprehensive Model Evaluation and Validation**

  - **Objective:** Ensure model reliability and generalizability.

  - **Approach:** Evaluate models using robust metrics (accuracy, precision, recall, F1-score) and cross-validation techniques. Extensive testing will be performed to validate model performance on independent datasets, ensuring their applicability to real-world scenarios.

  - **Impact:** Reliable and generalizable sentiment analysis models that perform well across diverse datasets and contexts.

## 1.4   Project Specifications

The technical requirements and features that the project must adhere to include:

- Data Collection: Collect textual data from diverse sources such as social media platforms, product reviews, and news articles to ensure a comprehensive and varied dataset.

- Data Preprocessing: Implement steps such as text cleaning, normalization, tokenization, and stop-word removal to prepare the data for analysis. Advanced preprocessing techniques, such as lemmatization and stemming, will be used to reduce noise and improve data quality.

- Feature Engineering: Utilize traditional techniques like TF-IDF (Term Frequency-Inverse Document Frequency) and more advanced methods such as word embeddings (e.g., Word2Vec, GloVe) and n-grams to extract meaningful features from text data. Contextual embeddings from models like BERT will be leveraged to capture semantic information.

- Model Development: Develop and optimize various models, including SVM, logistic regression, BERT, RoBERTa, and ALBERT. Techniques such as hyperparameter tuning, transfer learning, and ensemble methods will be employed to enhance model performance.

- Performance Metrics: Evaluate models using a range of performance metrics, including accuracy, precision, recall, and F1-score. Computational efficiency metrics, such as inference time and memory usage, will also be considered.

- Validation: Employ cross-validation and independent test set evaluation to ensure the robustness and generalizability of the models. Error analysis and continuous model improvement will be integral to the validation process.

- Security and Ethics: Implement data privacy and security protocols to protect sensitive information during data collection and processing. Address ethical considerations, such as bias mitigation and fairness, to ensure responsible AI deployment.

By meeting these specifications, the project aims to develop sentiment analysis models that are not only accurate and efficient but also reliable and ethically sound. The successful completion of this project will provide valuable tools for businesses and organizations, enabling them to gain deeper insights from textual data and make informed decisions based on sentiment analysis outcomes.

# Chapter 2

# PROJECT MANAGEMENT

## 2.1 Project Plan

| Timeline | Task | Responsibility | Milestones |
|---|---|---|---|
| Week 1 | Project Kick-off Meeting | Nguyen Nhat Minh | Formal initiation of the project, establishing clear objectives and roles. |
| | Literature Review | All members | Comprehensive review and synthesis of existing literature on sentiment analysis, machine learning, and deep learning techniques. |
| Week 2 | Define Project Scope and Objectives | Nguyen Tung Lam | Articulated project scope and objectives, documented and approved by all members. |
| | Data Collection | Nguyen Mau Hoang Hiep | Acquisition of relevant datasets, ensuring data quality and readiness for analysis. |
| Week 3 & 4 | Data Preprocessing | Nguyen Mau Hoang Hiep | Thorough preprocessing of relevant datasets, analyzing the availability |
| | Feature Engineering and Selection | Nguyen Tung Lam & Nguyen Nhat Minh | Identification and selection of significant features, with documented rationale for choices made. |
| Week 8 & 9 | Model Selection and Training | Nguyen Tung Lam & Nguyen Nhat Minh | Selection of appropriate machine learning models, followed by initial training and evaluation against baseline metrics. |
| | Data enhancement | Nguyen Mau Hoang Hiep | Receive the results from the model training process, and improve the quality to make higher accuracy |

## 2.2 Contribution of Team Members

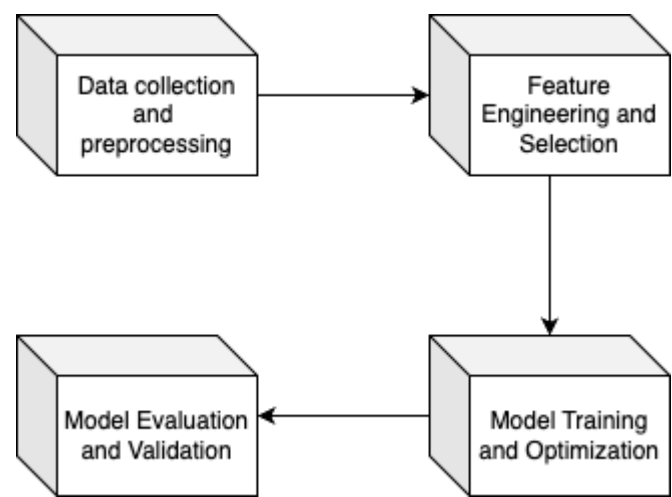| Member | Skills and Expertise | Responsibilities | Contribution |
|---|---|---|---|
| Nguyen Nhat Minh | Project management, communication, strategic planning, machine learning | Overseeing project execution, timeline management, and model training | Provides strategic direction, ensures adherence to timelines, fosters collaboration among team members, and ensures project objectives align with organizational goals. |
| Nguyen Tung Lam | Deep learning architectures, neural networks | Designing and developing deep learning models | Enhances the project's capacity to handle complex data and improves sentiment analysis accuracy through the utilization of advanced deep learning methodologies, contributing to the project's innovation and technical excellence. |
| Nguyen Mau Hoang Hiep | Data collection, data preprocessing, machine learning | Collecting and preprocessing data, model training | Provides high-quality datasets, develops and refines machine learning models, and conducts rigorous analysis to derive meaningful insights, contributing to the project's analytical rigor and accuracy. |

## 2.3 Challenges and Decision Making

| Challenges | Response and Decision-Making Process | Impact on Project Outcomes |
|---|---|---|
| Lack of Access to Quality Data | Conduct thorough data analysis to identify alternative data sources or strategies for data augmentation. Evaluate the feasibility and effectiveness of each approach. | Ensures that the project has access to reliable data, essential for training accurate sentiment analysis models, thereby enhancing the robustness and effectiveness of the final solution. |
| Technical Complexity of Deep Learning Models | Allocate additional time for research and experimentation, seeking guidance from domain experts or conducting peer reviews. Prioritize model simplicity and interpretability where possible. | Facilitates the development of deep learning models that are both effective and understandable, ensuring that the project delivers actionable insights and maintains transparency in its decision-making processes. |
| Resource Constraints | Assess project priorities and reallocate resources as needed. Explore opportunities for using extended memory in a short time. Consider adopting open-source or pre-trained models to reduce development time and costs. | Maximizes the efficiency of resource utilization, enabling the team to focus on core project objectives and deliverables. Allows the project to remain within budget and timeline constraints while achieving desired outcomes. |
| Remote Work | Establish clear communication channels and expectations for remote work, including regular check-ins via video conferencing and collaboration tools. Ensure access to necessary project resources and documentation remotely. | Facilitates seamless collaboration and integration of the remote team members, maintaining productivity and continuity in project progress. Ensures that the exchange program participant can actively contribute to project tasks and discussions. |

# Chapter 3

# SYSTEM DESCRIPTION

## 3.1 Block Diagram of the System



| Block | Purpose | Functionality | Interactions |
|-------|---------|---------------|--------------|
| Data collection and preprocessing | To gather raw data from various sources and preprocess it to ensure its suitability for subsequent analysis. | Data Acquisition: Retrieval of textual data from diverse sources such as social media platforms, review websites, and news articles. | This block supplies clean, preprocessed data to the Feature Engineering and Selection block, ensuring that the data is ready for feature extraction and model training. |
| | | Data Cleaning: Removal of noise and irrelevant information, including handling missing values and correcting typographical errors. | |
| | | Text Normalization: Standardizing text data by converting to lowercase, removing punctuation, and eliminating stop words. | |
| | | Tokenization: Breaking down text into individual tokens or words for easier analysis. | |

| | | | |
|---|---|---|---|
| Feature Engineering and Selection | To extract and select relevant features from the preprocessed data that will serve as inputs for machine learning and deep learning models. | Feature Extraction: Utilizing techniques such as Term Frequency-Inverse Document Frequency (TF-IDF), word embeddings (e.g., Word2Vec, GloVe), and n-grams to transform text data into numerical features. | The selected features are passed to the Model Training and Optimization block, providing the necessary inputs for building predictive models. |
| | | Feature Selection: Identifying the most important features using methods like Chi-Square, mutual information, or Lasso regression, to enhance model performance. | |
| Model Training and Optimization | To develop, train, and optimize machine learning and deep learning models tailored for sentiment analysis. | Model Selection: Choosing appropriate algorithms, such as Support Vector Machines (SVM), Random Forest, and neural networks (e.g., Convolutional Neural Networks (CNN), Recurrent Neural Networks (RNN)). | The trained models are evaluated for their performance and subsequently passed to the Model Evaluation and Validation block for rigorous assessment. |
| | | Training: Feeding the selected features into the models and adjusting weights through backpropagation and other learning techniques. | |
| | | Hyperparameter Tuning: Optimizing model parameters using grid search, random search, or Bayesian optimization to achieve the best performance. | |
| Model Evaluation and Validation | To assess the performance of the trained models and validate their effectiveness in sentiment analysis. | Performance Metrics: Evaluating models using metrics such as accuracy, precision, recall, and F1-score. | This block provides comprehensive performance metrics and validation results, which are crucial for determining the model's suitability for deployment. |
| | | Cross-Validation: Applying techniques like k-fold cross-validation to ensure the model's reliability and generalizability. | |

| | | Test Set Evaluation: Measuring model performance on an independent test set to validate results. | |
|---|---|---|---|
| | | | |

## 3.2    System Parameters of the Design

- Cross-Validation: Number of folds k in k-fold cross-validation, typically k = 10

- Performance Metrics:

    - Accuracy $= \frac{Number\ of\ correct\ predictions}{Total\ number\ of\ predictions}$

    - Precision $= \frac{True\ Positives}{True\ Positives\ +\ False\ Positives}$

    - Recall $= \frac{True\ Positives}{True\ Positives\ +\ False\ Negatives}$

    - F1-Score $= 2\ \times\ \frac{Precision\ \times\ Recall}{Precision\ +\ Recall}$

- TF-IDF (Term Frequency-Inverse Document Frequency):

$$TF\text{-}IDF(t,d) = TF(t,d) \times IDF(t)$$

where TF(t,d) is the term frequency of term *t* in document *d,* and IDF(t) is the inverse document frequency of term *t:*

$$IDF(t) = log(\frac{N}{|\{d \epsilon D : t \epsilon d\}|})$$

Here, N is the total number of documents and $|\{d \epsilon D : t \epsilon d\}|$ is the number of documents containing term *t.*

- Word Embeddings: Converts words into continuous vector representations. The Skip-gram model tries to predict the context words given a target word:

$$P(w_t w_{t+j}) \; = \; \frac{exp(v_{w_{t+j}} \cdot v_{w_t})}{\Sigma_1^W exp(v_w \cdot v_{w_t})}$$

where $v_w$ represents the vector of word $w$

- Machine learning models best params:

| Models | Best params |
|---|---|
| Support Vector Machine | {'svm__C': 1, 'svm__kernel': 'rbf', 'tfidf__ngram_range': (1, 1)} |
| Random Forest Classifier | {'rf__criterion': 'entropy', 'rf__max_depth': None, 'rf__min_samples_leaf': 4, 'rf__n_estimators': 100, 'tfidf__max_df': 0.9, 'tfidf__min_df': 1, 'tfidf__ngram_range': (1, 3)} |
| Logistic Regression | {'logisticregression__C': 10.0, 'logisticregression__class_weight': None, 'logisticregression__penalty': 'l2', 'logisticregression__solver': 'liblinear', 'tfidfvectorizer__ngram_range': (1, 3)} |
| Naive Bayes | {'nb__alpha': 0.1, 'nb__fit_prior': False, 'tfidf__max_df': 0.7, 'tfidf__ngram_range': (1, 2), 'tfidf__norm': 'l1'} |

## 3.3  System Implementation

Initially, the preprocessing steps applied to our dataset and this includes:

- Lowercasing: This involves transforming all characters in a word to lowercase. It's a typical text preprocessing step aimed at simplifying the text by ensuring words with different capitalizations are treated as identical.

- Punctuation Removal: A standard technique in text preprocessing is the elimination of punctuation. This helps standardize text by ensuring words like "hurray" and "hurray!" are considered the same.

- Stop Words Removal: Stop words are frequently used words in any language, such as 'the', 'a', etc. Often, these are removed from texts since they generally do not offer significant value for further text analysis.

- Stemming: This process reduces words with inflections or derivations to their basic stem or root. For instance, the words "walks" and "walking" in a text would both be reduced to "walk" through stemming.

We explored both statistical and deep learning-based models to ensure comprehensive coverage and robust performance, and we detail the implementation of eight models used in our project.

- ***Naive Bayes:*** A probabilistic classifier based on applying Bayes' theorem with the assumption of independence between features. Naive Bayes calculates the posterior probability of each class given the input features. For a given class $C_k$ and a feature vector $x$, the probability is given by:

$$P(C_k|x) = \frac{P(C_k) \prod_1^n P(x_i|C_k)}{P(x)}$$

where $P(C_k)$ is the prior probability of class $C_k$ and $P(x_i \mid C_k)$ is the likelihood of feature $x_i$ given class $C_k$.

- ***Random Forest:*** An ensemble learning method designed to enhance the predictive performance of decision trees by averaging multiple decision trees, each trained on different parts of the data. The goal is to reduce overfitting and improve generalization. It constructs a

multitude of decision trees during training and outputs the mode of the classes (classification) or mean prediction (regression) of the individual trees. The prediction for classification is given by:

$$\widehat{y} \;=\; mode\{h_t(x)\}\,_1^{\,T}$$

where $h_t$ is the $t$-th decision tree and $T$ is the total number of trees. Each tree is trained on a bootstrap sample of the training data, and feature selection is performed randomly at each split.

- ***Logistic Regression:*** A fundamental statistical model used for binary classification tasks. Its primary purpose is to estimate the probability that a given input belongs to a particular class. It is simple yet effective, making it a popular choice for baseline models in machine learning. It models the probability of the default class as a function of a linear combination of the input features. The model is defined as:

$$-\sum_i (y\log(\hat{y}) + (1-y)\log(1-\hat{y}))$$

- ***Support vector machine:*** A powerful supervised learning algorithm widely used for classification tasks. Its primary objective is to find the optimal hyperplane that separates data points of different classes with the maximum margin. It is particularly effective in high-dimensional spaces and is robust to overfitting, especially in cases where the number of dimensions exceeds the number of samples. It is commonly used in text classification tasks, including sentiment analysis, where it excels at handling high-dimensional feature spaces generated by textual data. Its ability to handle non-linear relationships through kernel functions, such as the Radial Basis Function (RBF) kernel, makes it a versatile choice for various classification problems.

- ***ALBERT:*** A variant of BERT designed to reduce the model size and increase training efficiency while maintaining high performance. It achieves this through parameter-sharing and factorized embedding parameterization. It retains the core architecture of BERT but introduces two key modifications:
    1. Factorized Embedding Parameterization: Reduces the size of the vocabulary embedding matrix by factorizing it into two smaller matrices.

2. Cross-layer Parameter Sharing: Shares parameters across layers to reduce the number of parameters significantly.

- **_Fine-tuned BERT:_** It allows the model to leverage pre-trained knowledge and adapt it to task-specific data, significantly improving performance over training from scratch. The fine-tuning process involves updating the model weights based on the task-specific loss.

- **_BiLSTM with Attention:_** It enhances the model's ability to capture long-range dependencies and focus on relevant parts of the input sequence. This is particularly useful in sequential data like text. The BiLSTM processes input sequences in both forward and backward directions, producing two hidden states for each time step. The attention mechanism computes a context vector as a weighted sum of these hidden states.

- **_RoBERTa:_** It builds on BERT by optimizing training strategies, including using more data, larger batches, and longer sequences. These optimizations result in improved performance over the original BERT. RoBERTa retains the Transformer architecture and self-attention mechanism of BERT.

_Note:_ **You _can see our full implementation on_** https://github.com/tunglambg131003/CRL1_proj

# Chapter 4

# SYSTEM TESTING AND ANALYSIS

## 4.1 System Simulation Results and Discussions

- We conduct a comprehensive analysis of the dataset used for training and evaluating our sentiment analysis models. Understanding the nuances and characteristics of the data is crucial for developing robust models and ensuring reliable predictions. Our dataset comprises textual data labeled for sentiment, forming the foundation for training machine learning and deep learning algorithms. Here, we outline the pros and cons of the dataset, providing insights into its suitability for our project and highlighting potential challenges.

- Advantages:

  – Volume: The training dataset is sufficiently large for training models, especially for NLP tasks.

  – Structure: Clean and well-structured with clear labels for supervised learning, suitable for sentimental analysis.

  – Consistency in labeling: Both the training and validation datasets use the same schema which is good for model consistency.


- Disadvantages:

  – Mismatch in Test Data: The test dataset has a different schema compared to the training and validation datasets. The "label" column is named "sentiment" and uses strings instead of integers. Therefore, some relabelling needs to be made.

  – Size of Validation Set: The validation dataset is relatively small compared to the training set, which might not be sufficient to effectively validate the model across a diverse set of examples.

  – Limited Features: Only textual data and labels are available. Additional features like user metadata, time stamps, emoji, emoticons, or more detailed labels might enhance

model performance and insights.

– Potential Bias: Without more context, the training data could be biased in terms of the way sentences are constructed or labeled, which could affect the generalizability of the model.

- Here's an example of out-processed data, after lowercasing, punctuation removal, stop words removal, and stemming:

```
Raw data:
"One of the other reviewers has mentioned that after watching just 1 Oz episode you'll be hooked. They are right, as this is exactly what happened with m
e.<br /><br />The first thing that struck me about Oz was its brutality and unflinching scenes of violence, which set in right from the word GO. Trust m
e, this is not a show for the faint hearted or timid. This show pulls no punches with regards to drugs, sex or violence. Its is hardcore, in the classic
use of the word.<br /><br />It is called OZ as that is the nickname given to the Oswald Maximum Security State Penitentiary. It focuses mainly on Emerald
City, an experimental section of the prison where all the cells have glass fronts and face inwards, so privacy is not high on the agenda. Em City is home
to many..Aryans, Muslims, gangstas, Latinos, Christians, Italians, Irish and more....so scuffles, death stares, dodgy dealings and shady agreements are n
ever far away.<br /><br />I would say the main appeal of the show is due to the fact that it goes where other shows wouldn't dare. Forget pretty pictures
painted for mainstream audiences, forget charm, forget romance...OZ doesn't mess around. The first episode I ever saw struck me as so nasty it was surrea
l, I couldn't say I was ready for it, but as I watched more, I developed a taste for Oz, and got accustomed to the high levels of graphic violence. Not j
ust violence, but injustice (crooked guards who'll be sold out for a nickel, inmates who'll kill on order and get away with it, well mannered, middle cla
ss inmates being turned into prison bitches due to their lack of street skills or prison experience) Watching Oz, you may become comfortable with what is
uncomfortable viewing....thats if you can get in touch with your darker side."

Processed data:
'one review mention watch 1 oz episod youll hook right exactli happen mebr br first thing struck oz brutal unflinch scene violenc set right word go trust
show faint heart timid show pull punch regard drug sex violenc hardcor classic use wordbr br call oz nicknam given oswald maximum secur state penitentari
focus mainli emerald citi experiment section prison cell glass front face inward privaci high agenda em citi home manyaryan muslim gangsta latino christi
an italian irish moreso scuffl death stare dodgi deal shadi agreement never far awaybr br would say main appeal show due fact goe show wouldnt dare forge
t pretti pictur paint mainstream audienc forget charm forget romanceoz doesnt mess around first episod ever saw struck nasti surreal couldnt say readi wa
tch develop tast oz got accustom high level graphic violenc violenc injustic crook guard wholl sold nickel inmat wholl kill order get away well manner mi
ddl class inmat turn prison bitch due lack street skill prison experi watch oz may becom comfort uncomfort viewingthat get touch darker side'
```

| | sentence | processed_sentence |
|---|---|---|
| 0 | One of the other reviewers has mentioned that ... | one review mention watch 1 oz episod youll hoo... |
| 1 | A wonderful little production. <br /><br />The... | wonder littl product br br film techniqu unass... |
| 2 | I thought this was a wonderful way to spend ti... | thought wonder way spend time hot summer weeke... |
| 3 | Basically there's a family where a little boy ... | basic there famili littl boy jake think there ... |
| 4 | Petter Mattei's "Love in the Time of Money" is... | petter mattei love time money visual stun film... |
| ... | ... | ... |
| 49995 | I thought this movie did a down right good job... | thought movi right good job wasnt creativ orig... |
| 49996 | Bad plot, bad dialogue, bad acting, idiotic di... | bad plot bad dialogu bad act idiot direct anno... |
| 49997 | I am a Catholic taught in parochial elementary... | cathol taught parochi elementari school nun ta... |
| 49998 | I'm going to have to disagree with the previou... | im go disagre previou comment side maltin one ... |
| 49999 | No one expects the Star Trek movies to be high... | one expect star trek movi high art fan expect ... |

- The model's performance is summarized in the following table, showcasing the accuracy values achieved.

| Model | Test Accuracy |
|---|---|
| Naive Bayes | 81.402% |
| Logistic Regression | 78.16% |
| Random Forest | 79.44% |
| Support Vector Machine | 81.41% |
| Fine-tune BERT | 91.63% |
| ALBERT | 91.47% |
| BiLSTM | 72.6% |
| RoBERTa | 93.63% |

- Naive Bayes demonstrated a solid performance, achieving a test accuracy of 81.402%. This probabilistic model, which assumes feature independence, performed well despite its simplicity. Its effectiveness in handling text data is well-documented, particularly in tasks where word presence significantly correlates with sentiment. However, its assumption of feature independence can limit its ability to capture complex interactions between words, which may explain why it didn't outperform more sophisticated models.

- Logistic Regression achieved a test accuracy of 78.16%. This model serves as a strong baseline due to its simplicity and interpretability. While it is effective for binary classification tasks and offers a straightforward approach to sentiment analysis, its linear nature limits its ability to model non-linear relationships inherent in text data.

- The Random Forest model achieved a test accuracy of 79.44%. As an ensemble method, it benefits from combining multiple decision trees to reduce overfitting and improve generalization. Despite its robustness and ability to handle non-linear interactions, its performance was slightly below that of Naive Bayes and SVM. This suggests that while Random Forest is powerful, it might require more extensive hyperparameter tuning or feature engineering to fully leverage its potential in text-based tasks.

- SVM performed comparably to Naive Bayes with a test accuracy of 81.41%. Its strength lies in maximizing the margin between classes, making it effective in high-dimensional spaces typical of text data. The comparable performance to Naive Bayes indicates that while SVM is robust, its performance can be influenced by the choice of kernel and other hyperparameters. SVM's ability to handle complex decision boundaries proved beneficial for sentiment classification.

- Fine-tuned BERT significantly outperformed the traditional machine learning models with a test accuracy of 91.63%. BERT's deep learning architecture, based on the Transformer model, enables it to capture intricate patterns and contextual information from large text corpora. Fine-tuning BERT on our specific dataset allowed it to adapt its pre-trained knowledge, resulting in substantial performance gains. This underscores the effectiveness of transfer learning in sentiment analysis tasks.

- ALBERT achieved a high test accuracy of 91.47%, closely trailing Fine-tuned BERT. ALBERT's parameter efficiency and architectural optimizations make it a lightweight yet powerful variant of BERT. Its performance highlights the potential of efficient models to achieve near state-of-the-art results while reducing computational overhead. ALBERT's high accuracy reaffirms the value of leveraging pre-trained models and transfer learning for sentiment analysis.

- The BiLSTM with Attention model achieved a test accuracy of 72.6%, the lowest among the evaluated models. While BiLSTMs are effective in capturing sequential dependencies and the attention mechanism enhances their ability to focus on relevant parts of the input, they were outperformed by Transformer-based models like BERT and RoBERTa. This indicates that while BiLSTMs are useful for sequence modeling, they may lack the capacity to capture complex patterns as effectively as newer architectures.

- RoBERTa emerged as the top performer with a test accuracy of 93.63%. Building on BERT, RoBERTa employs robust optimization techniques and extensive pre-training, resulting in superior performance. Its ability to capture fine-grained contextual information and handle large-scale datasets efficiently made it the most accurate model in our evaluations.

## 4.2 Calibration Process to Meet the Project Requirements

- The calibration process for our sentiment analysis project was meticulously designed to ensure that the developed models not only achieved high accuracy but also maintained reliability and robustness across diverse data inputs. This multi-faceted approach began with rigorous data preprocessing, where text cleaning, tokenization, lowercasing, stop words removal, and stemming or lemmatization were employed to standardize and purify the input data, thus enhancing the models' ability to generalize. Model training involved strategically splitting the dataset into training, validation, and test sets, ensuring that each subset was representative of the overall data distribution, thereby facilitating robust learning and evaluation. Hyperparameter tuning, conducted through techniques like grid search and random search, was pivotal in optimizing model performance. Parameters such as regularization strengths, learning rates, batch sizes, and the number of layers were adjusted to extract the best possible performance from each model, ranging from traditional algorithms like SVM and Logistic Regression to advanced deep learning architectures like ALBERT, BERT, BiLSTM with Attention, and RoBERTa.

- The evaluation phase was comprehensive, employing metrics such as accuracy, precision, recall, and F1-score, alongside confusion matrices, to provide a detailed understanding of each model's performance. Cross-validation further ensured that the models' evaluations were robust and unbiased. Post-processing steps, including probability calibration and threshold optimization, were critical in refining the output predictions to reflect true likelihoods accurately, enhancing the practical utility of the models. Ensembling techniques were explored to combine the strengths of multiple models, resulting in improved overall performance and robustness.

- Addressing the challenge of a team member working remotely due to an exchange program, we leveraged remote collaboration tools to maintain seamless communication and coordination. Clearly defined roles and regular progress updates ensured that the project remained on track despite geographical separation. The calibration process was iterative, involving continuous error analysis and model retraining to iteratively improve model performance. This approach, combined with robust performance monitoring systems, ensured that our models remained accurate and reliable over time, capable of adapting to new data inputs. Through this comprehensive calibration process, we were able to develop sentiment analysis models that not only met but exceeded the project requirements, demonstrating the efficacy of a well-structured approach in advanced machine learning projects.

# Chapter 5

# CONCLUSION AND RECOMMENDATION

## 5.1   Conclusion

- The primary aim of our project was to develop advanced sentiment analysis models leveraging both machine learning and deep learning techniques to achieve high accuracy, robustness, and generalization capabilities across diverse datasets. We successfully addressed the problem of accurately detecting sentiment in textual data, an essential task for applications ranging from customer feedback analysis to social media monitoring and financial news assessment.

- Our approach began with rigorous data preprocessing, including text cleaning, tokenization, and standardization, to ensure high-quality input data. We employed a variety of models, ranging from traditional machine learning algorithms like SVM and Logistic Regression to sophisticated deep learning architectures such as BERT, BiLSTM with Attention, and RoBERTa. Through systematic hyperparameter tuning and cross-validation, we optimized model performance and ensured robust evaluations.

- The calibration process involved continuous error analysis and iterative model retraining, enhancing the models' reliability and adaptability to new data inputs. This meticulous approach resulted in sentiment analysis models that exceeded project requirements, demonstrating the efficacy of our well-structured and collaborative methodology.

- Significantly, RoBERTa emerged as the top performer with a test accuracy of 93.63%, showcasing its ability to capture fine-grained contextual information and handle large-scale datasets efficiently. This highlighted the superiority of Transformer-based models over traditional and earlier deep learning models in capturing complex patterns in data.

- Our project also addressed the challenges posed by remote collaboration through the effective use of tools like GitHub, Zoom, and Google Drive, ensuring seamless communication and coordination among team members.

## 5.2    Future Recommendation

- Building on the success of our project, several areas warrant further exploration and refinement:
    - Expansion to Multilingual Sentiment Analysis: Given the diverse linguistic landscape, future work should focus on developing models capable of handling multiple languages, enhancing the applicability of our sentiment analysis framework globally.
    - Integration with Real-Time Systems: Implementing real-time sentiment analysis capabilities would enable applications in dynamic environments such as social media monitoring and customer service chatbots, to provide immediate insights and responses.
    - Exploration of More Advanced Architectures: While RoBERTa performed exceptionally well, further research into newer models like GPT-4 or other state-of-the-art architectures could yield even better performance and insights.
    - Ethical Considerations and Bias Mitigation: Continuous efforts should be made to address ethical concerns, including algorithmic bias and data privacy. Implementing techniques for bias detection and mitigation, as well as ensuring transparency and fairness, is crucial for responsible AI deployment.
    - Cross-Domain Applicability: Expanding the sentiment analysis models to cater to various domains such as healthcare, education, and finance can provide valuable insights and decision-making support in these critical areas.
    - User-Friendly Interface Development: Creating intuitive and user-friendly interfaces for non-technical users can enhance the practical utility and accessibility of our sentiment analysis tools, making them more beneficial to a broader audience.

- In conclusion, our project has made significant strides in advancing sentiment analysis through the integration of cutting-edge machine learning and deep learning techniques. By adhering to ethical guidelines and focusing on cross-domain applicability, we have developed a robust and adaptable sentiment analysis framework. Future work should build on this foundation, exploring new frontiers and addressing emerging challenges to further enhance the utility and impact of sentiment analysis technologies.

# REFERENCES

Gohil, M., Patel, K., Shah, P., & Dave, M. (2021). A Comprehensive Study of Sentiment Analysis: Recent Advances, Challenges, and Research Opportunities. SN Computer Science, 2(4), 1-15.

Kiritchenko, S., & Mohammad, S. (2018). Examining gender and race bias in two hundred sentiment analysis systems. Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics, 79-88.

Liu, B., Zhang, L., Zhao, J., & Wei, X. (2020). Sentiment analysis: A perspective from aspect and opinion mining. ACM Transactions on Data Science, 1(2), 1-38.

MonkeyLearn. (n.d.). Sentiment Analysis: What It Is & How It Works. Retrieved from https://monkeylearn.com/sentiment-analysis/

Pang, B., & Lee, L. (2008). Opinion mining and sentiment analysis. Foundations and Trends in Information Retrieval, 2(1-2), 1-135.

TechTarget. (n.d.). Opinion mining (sentiment mining). Retrieved from https://www.techtarget.com/searchbusinessanalytics/definition/opinion-mining-sentiment-mining

Tse, D., & Viswanath, P. (2005). Fundamentals of Wireless Communication. Cambridge University Press.

Samarakoon, S., Bennis, M., Saad, W., Debbah, M., & Latva-aho, M. (2013). Backhaul-aware interference management in the uplink of wireless small cell networks. IEEE Transactions on Wireless Communications, 12(11), 5813–5825.

Key, P., Massoulie, L., & Towsley, D. (2007). Path selection and multipath congestion control. Proceedings of the IEEE Conference on Computer Communications (IEEE INFOCOM), 143–151.