

BULLYING DETECTION FROM TWEETS WITH LSTM AND BERT TRANSFORMER

BY

Robin Mia

ID: 0242220004213014

This Report Presented in Partial Fulfillment of the Requirements for the
Degree of Master of Science in Management Information System (MIS)

Supervised By

Mr. Abdus Sattar

Assistant Professor & Coordinator M. Sc.

Department of CSE

Daffodil International University



DAFFODIL INTERNATIONAL UNIVERSITY

DHAKA, BANGLADESH

January 2025

APPROVAL

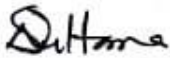
The project, "Bullying Detection from Tweets with LSTM and BERT Transformer," was turned in by Robin Mia to the Daffodil International University Department of Management Information System (MIS). It has been approved in terms of style and content and is considered satisfactory for partially fulfilling the requirements for the M.S in Management Information System (MIS) degree. The date of 11 January 2025 presentation took place.



Dr. S. M. Aminul Haque
Professor & Associate Head
Department of CSE
Faculty of Science & Information Technology
Daffodil International University

Chairman

BOARD OF EXAMINERS



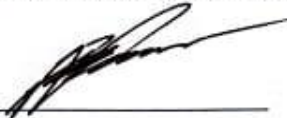
Dr. Naznin Sultana
Associate Professor
Department of CSE
Faculty of Science & Information Technology
Daffodil International University

Internal Examiner



Mr. Md. Sadekur Rahman
Assistant Professor
Department of CSE
Faculty of Science & Information Technology
Daffodil International University

Internal Examiner



Mr. Nazibur Rahman
Technical Lead - Database Administrator,
Wipro Bangladesh Telenor - Grameen Phone

External Examiner

DECLARATION

We hereby certify that we completed this study under the guidance of **Mr. Abdus Sattar**, Assistant professor in the CSE department of Daffodil International University. Furthermore, we affirm that neither this project nor any portion of it has been submitted for consideration for a degree or certificate elsewhere.

Supervised by:



Mr. Abdus Sattar
Assistant Professor & Coordinator M. Sc.
Department of CSE
Daffodil International University

Submitted by:



Robin Mia
ID: **0242220004213014**
Department of Management Information System (MIS)
Daffodil International University

ACKNOWLEDGEMENT

First and foremost, we give the Almighty God our sincere gratitude and appreciation for His divine favor, which enables us to successfully finish the final year Thesis.

We sincerely thank **Mr. Abdus Sattar**, an Assistant professor, Department of CSE, Daffodil International University, Dhaka, and express our sincere gratitude. To complete this project, our supervisor must have deep knowledge of and a strong interest in the subject of "natural language processing." The completion of this project has been made possible by his unending patience, academic guidance, ongoing encouragement, persistent and vigorous supervision, constructive criticism, insightful counsel, and his reading of several subpar drafts and correction of them at every level.

We would like to extend our sincere appreciation to **Dr. Arif Mahmud**, Associate Professor & Program Director MIS & **Dr. Sheak Rashed Haider Noori**, Professor and Head of the CSE Department, as well as to other academic members and staff of Daffodil International University's CSE department, for their kind assistance in completing our project.

We express our gratitude to all of our Daffodil International University classmates who participated in this discussion while finishing their coursework.

Lastly, we have to respectfully thank our parents for their unwavering support and patience.

ABSTRACT

In the internet age, cyberbullying has grown to be a serious concern, particularly for English-speaking nations. This work focuses on identifying cyberbullying in English with the use of deep learning techniques. A specialized English dataset, comprising instances of both cyberbullying and non-cyberbullying text, is utilized for training a deep learning model. Tokenization, preprocessing, and sequence transformation are applied to the dataset so that it may be fed into Random Forest, Naïve Bayes, and BERT classifiers using LSTM cells. The novel LSTM-based deep learning model was used for the dataset and the dropout and word embedding technique were used to improve the model's performance. The best model was evaluated with confusion matrix. Research is being done on a number of approaches, including language-specific preprocessing and data augmentation, to address the particular problems with cyberbullying detection in English. The results demonstrate how well deep learning works to identify cyberbullying in English-speaking contexts and show how promising the technology is for addressing this issue. The study reveals that the BERT achieved an accuracy of 87%, demonstrating its superior performance. Additionally, an alternative approach using LSTM yielded the accuracy 84%. Ensemble models, including Naïve Bayes (NB), and Random Forest, were also employed, with hyperparameter tuning optimizing their performance. Notably, the LSTM and BERT outperformed other models, attaining the highest accuracy rate of 87% in cyberbullying detection, as confirmed by recent experimental inquiries evaluating these findings.

Keywords: BLSTM, LSTM, Deep Learning, Algorithms, Ensemble Model, Dropout, Embedding.

TABLE OF CONTENTS

CONTENTS	PAGE
Board of examiners	i
Declaration	ii
Acknowledgements	iii
Abstract	iv
CHAPTER	PAGE
CHAPTER 1: INTRODUCTION	1-5
1.1 Introduction	1
1.2 Motivation	2
1.3 Rationale of the Study	3
1.4 Research Questions	3
1.5 Expected Output	3
1.6 Project Management and Finance	4
1.7 Report Layout	4
CHAPTER 2: BACKGROUND	5-9
2.1 Preliminaries	5
2.2 Related Works	5
2.3 Comparative Analysis and Summary	7
2.4 Scope of the Problem	8
2.5 Challenges	8

CHAPTER 3: RESEARCH METHODOLOGY	10-19
3.1 Research Subject and Instrumentation	10
3.2 Data Collection Procedure	12
3.3 Statistical Analysis	15
3.4 Proposed Methodology	16
3.5 Implementation Requirements	19
CHAPTER 4: EXPERIMENTAL RESULTS AND DISCUSSION	20-27
4.1 Experimental Setup	20
4.2 Experimental Results & Analysis	20
4.3 Discussion	27
CHAPTER 5: IMPACT ON SOCIETY, ENVIRONMENT AND SUSTAINABILITY	28-29
5.1 Impact on Society	28
5.2 Impact on Environment	28
5.3 Ethical Aspects	28
5.4 Sustainability Plan	29
CHAPTER 6: SUMMARY, CONCLUSION, RECOMMENDATION AND IMPLICATION FOR FUTURE RESEARCH	30-31
6.1 Summary of the Study	30
6.2 Conclusions	30
6.3 Implication for Further Study	31
REFERENCES	32-33

LIST OF TABLES

TABLES	PAGE NO
Table 3.1: Explanation of the Dataset	13

LIST OF FIGURES

FIGURES	PAGE NO
Figure 3.1: The Flow Chart of Methodology	10
Figure 3.2: Count plot of Target Column	13
Figure 3.3: Tweets less than 10 words	14
Figure 3.4: Tweets with high number of words	14
Figure 3.5: Tweets with top 20 most common of words	15
Figure 4.1: Experimental Results of Random Forest	22
Figure 4.2: Confusion Matrix of Random Forest	23
Figure 4.3: Experimental Results of Naïve Bayes	24
Figure 4.4: Confusion Matrix of Naïve Bayes	25
Figure 4.5: Experimental Results of LSTM	25
Figure 4.6: Confusion Matrix of LSTM	26
Figure 4.7: Experimental Results of BERT	26
Figure 4.8: Confusion Matrix of BERT	27

CHAPTER 1

INTRODUCTION

1.1 Introduction

Cyberbullying has proliferated in the digital age and now affects people of all ages and socioeconomic statuses. Social media and other online platforms have opened up new avenues for bad behavior, which is detrimental to the mental, emotional, and social wellbeing of victims. The English-speaking community encounters unique challenges related to cyberbullying, shaped by linguistic nuances and cultural factors influencing the nature and manifestation of harmful online behaviors. For people who interact in English, establishing a safe and welcoming online space requires recognizing and responding to cyberbullying in the language. This research aims to give an approach based on deep learning for detecting cyberbullying in English. Deep learning has promise for the efficient identification and categorization of cyberbullying incidents since it has shown encouraging results in a variety of natural language processing problems. Our objective is to improve methods for detecting cyberbullying that are especially designed for the English language by utilizing deep learning models. The primary objective involves designing and training a deep learning model capable of accurately distinguishing among neutral and bullying words in English. To achieve this, we utilize an English dataset that includes text examples of both cyberbullying and non-cyberbullying. The dataset is painstakingly preprocessed, tokenized, and transformed into sequences that can be fed into an LSTM-cell BERT. In addition to training the model, we investigate other approaches to tackle the difficulties in detecting cyberbullying in English, such as language-specific preprocessing techniques, data augmentation strategies, and the inclusion of English-specific cultural and contextual elements. The purpose of this study's findings is to provide light on how well deep learning techniques work to identify cyberbullying in English. Furthermore, this research contributes to the development of proactive measures to combat cyberbullying and create a safer online environment. To begin tackling this pressing issue in English-speaking Communities, the study clarifies the unique benefits and difficulties related to identifying cyberbullying in English.

1.2 Motivation

The motivation behind conducting this study on deep learning for English language cyberbullying detection is based on numerous crucial factors. Foremost among these is the recognition that cyberbullying has evolved into a substantial societal concern, impacting individuals across all age groups, including children, adolescents, and adults. The well-documented adverse effects of cyberbullying on mental health, self-esteem, and overall well-being underscore the urgency of addressing this issue. However, existing research on cyberbullying detection has predominantly centered on the English language, leaving a void in understanding and confronting this problem in English contexts. The English-speaking community, akin to many other English-speaking communities, encounters distinctive challenges in the realm of cyberbullying. The idiosyncrasies of the language, cultural influences, and prevalent online behaviors within the English-speaking population necessitate tailored approaches to effectively identify and counteract cyberbullying in this specific context. This project attempts to bridge the knowledge gap and offer significant new insights on how to address cyberbullying in English-speaking situations through the development of a deep learning model. Our proposed model can be impacted in social and economic life of users. We had tried to enhance the model performances with several different technologies of natural language processing. It can easily detect the cyberbullying words and texts according to model's dependency. Our social and economic life can be motivated to reduce using negative words and texts which means bullying to another person's. Moreover, by directing attention to the English language, this research strives to champion inclusivity and ensure that individuals communicating in English languages receive equitable attention and protection from cyberbullying. This study underscores the significance of linguistic diversity and cultural sensitivity in confronting online abuse, laying the groundwork for future research and initiatives that cater to the needs of diverse language communities. The ultimate motivation for this research is the desire to foster an online community that is safer and more welcoming for people who speak English. By harnessing the capabilities of deep learning techniques and addressing language-specific challenges, this research endeavors to make meaningful contributions toward the overarching goal of combatting cyberbullying and fostering a digital culture characterized by respect.

1.3 Rationale of the Study

The study is motivated by a multifaceted rationale, primarily aimed at addressing critical gaps in cyberbullying detection research specific to English contexts. The study seeks to contribute to the understanding of cyberbullying by delving into the cultural and linguistic intricacies inherent in the English language. By doing so, it aspires to empower English-speaking communities and facilitate the development of targeted solutions for this linguistic group. The exploration of deep learning within the realm of natural language processing serves as a pivotal aspect of this research. The project intends to assess these state-of-the-art technologies' efficacy in the difficult English language context using deep learning techniques. This exploration not only extends the applicability of deep learning but also contributes insights into its efficacy for addressing language-specific challenges associated with cyberbullying. By focusing on who communicate in the English language, the study advocates for the creation of protective measures that cater to linguistic diversity. This emphasis on linguistic and cultural nuances is fundamental to fostering a digital landscape that is secure and inclusive for all, irrespective of the language in which they communicate. In essence, the research's rationale is rooted in bridging gaps, understanding cultural nuances, empowering communities, exploring advanced technologies, and ultimately contributing to the creation of a safer online space for individuals engaging in communication within the English language.

1.4 Research Question

1. Which models work best for detecting cyberbullying in English datasets?
2. Which method is improving the overall performances for detecting cyber-bullying?
3. How does this research address cyberbullying ethically and socially?

1.5 Expected output

This research endeavors to address the existing gap in cyberbullying detection research for English contexts, with a specific focus on the English language. Through investigating the problem of cyberbullying among English-speaking individuals, the study aims to offer tailored solutions that effectively address the unique challenges faced by this linguistic group. To effectively prevent online harassment, certain techniques and interventions are

required due to the unique cultural and linguistic peculiarities inherent in royal languages. The expected outcome of our study is to motivate people about reduce the use of harmful and negative words. In our study, we have collected relevant data for best model training. We have reduced the extra columns and duplicate data for the best outcome. We have proposed the best ultimate solution to detect cyberbullying in social media and social platforms. The system can save the society and young aged nation from negative and harmful society. After successful implement it in real world, the counting of bullying and negative word will be reduced in significant amount. The system will be acceptable among worldwide users. We want to use our system in positive and secure data policy. The ultimate goal is to foster an online community that is more welcoming and safe for those who use the English language for communication.

1.6 Project Management and Finance

The suggested technology is efficient in an affordable solution appropriate for daily usage and has great potential as a useful tool in the field of cyberbullying detection in our nation. While high-configuration tools can improve performance and produce remarkable results, the actual real life, it will be efficient in social media comments and texts. It will automatically detect the negative words from social media platforms. Nevertheless, even with more straightforward tools, our model remains effective, ensuring seamless operation and contributing significantly to the ongoing efforts against cyberbullying. The main finance management of the project was in data collection process and in implementing the study in evaluation process.

1.7 Report Layout

The important components of our suggested study report are logically organized, starting with a full examination of the background research and going into the background and pertinent studies about cyberbullying. A thorough description of the methods, instruments, and strategies used to carry out the investigation and create the suggested model is given in the research methodology section. The study results, conclusions, and a thorough explanation of the findings are then presented in the experimental results and discussion section. A summary that draws conclusions and provides ideas for further research completes the study. The literature and sources that were used to bolster the study's conclusions are included in the references section.

CHAPTER 2

BACKGROUND STUDY

2.1 Preliminaries

The examination of cyberbullying patterns involves a precise analysis leveraging deep learning techniques. Within this segment, our focus is on exploring research related to the evaluation of reports on social media comments. Numerous computational models, such as BERT, LSTM, Random Forest Classifier (RF), and Naïve Bayes Classifier (NB), are employed for this investigation. Deep learning models play a central role in this section's research endeavors. The section also highlights a number of researchers that have used a range of models in their individual investigations, which has improved our knowledge of the dynamics of cyberbullying.

2.2 Related Works

Our developed deep learning techniques prove highly effective. A comprehensive exploration of research studies and methodologies reveals a diverse array of approaches employed to comprehensively address the intricate challenges associated with cyberbullying. These methods employ a variety of natural language processing (NLP) and deep learning algorithms, each of which contributes special advantages and insights to the main goal. Above all, deep learning models become particularly useful tools when trying to find instances of cyberbullying because they use algorithms based on decision trees, or "tree structures," which facilitate efficient decision-making processes that are essential for precise identification. One significant endeavor by Sambhagadi et al. [3] focuses on combatting cyberbullying on social media through the application of NLP approaches. Their approach goes beyond mere identification of cyberbullying, delving into the nuanced context in which profanities are used, facilitating the discrimination between offensive and neutral language. To enhance the accuracy of their annotations, a combination of crowdsourcing and in-lab annotations was employed, resulting in promising results with an F1-Score of 0.59. Notably, their approach achieves these results without relying on Personalized data and employs a distinct dataset. Yao et al. [4] address the recurrent element of cyberbullying on social media, which is characterized by a bully's cruel remarks

to their victim. Their unique approach significantly reduces the number of criteria required for classification while maintaining excellent accuracy by using a sequential hypothesis testing architecture. This method aims for three things: scalability, accuracy, and timeliness. The Instagram dataset, which was gathered using snowball sampling and partially annotated by subject matter experts using semi-supervised machine learning techniques, is used to train the model. It's important to recognize that there are a few drawbacks, though, such as the dataset being unique to Instagram, the absence of techniques for verifying the correctness of labels, and the laborious nature of comment-based label collection. In order to detect cyberbullying, Huang et al. [5] provide a unique method that integrates textual data with social network features. Their approach involves closely investigating the social network structure among individuals in order to extract characteristics such as number of friends, network embeddedness, and connection centrality. The study argues for the underutilization of social media capabilities in previous research and emphasizes the importance of considering the social context in which cyberbullying communications take place.

The approach utilizes the Twitter corpus from December 2008 and January 2009, employing synthetic minority oversampling (SMOTE) to create balanced data for classification. Various machine learning techniques are employed for this purpose. Rakib et al. [6] contribute to the landscape by collecting a corpus from the Reedit database, cleaning the data, and developing a word embedding model based on the word2vec skip-gram model for cyberbullying detection. This innovative word embedding model, enriched with domain knowledge, outperforms pre-trained word embedding models and traditional feature extraction methods. Subsequently, the model's characteristics are employed to educate a random forest classifier, effectively categorizing remarks authored by cyberbullies. Silva and colleagues [7] introduce a unique methodology for cyberbullying detection rooted in psychological research. Their approach involves the creation of an app called "Bully Blocker," designed to notify parents if cyberbullying behaviors are detected in their children's social media interactions. The app assesses the user's social media data by analyzing messages and comments, classifying them as indicators of bullying or warning signs. While the app currently relies on outdated

Facebook detection techniques, it holds potential for development through the incorporation of machine learning classification. Raisi et al. bring a distinctive approach to the table with the development of the Participant Vocabulary Consistency (PVC) method. Recognizing challenges in obtaining high-quality labeled data, this relational model operates with minimal supervision, requiring human experts to provide highly suggestive keywords indicative of harassment. The computer then utilizes these annotations to uncover additional potential keywords and identify specific bullying incidents by identifying victimization patterns within unlabeled social interaction networks. Hosseinmardi et al. [9] address cyberbullying incidents on the prominent social media platform Instagram by examining the most popular comments on a user's public postings. Their program efficiently differentiates cyber aggression data from the rest using Naive Bayes classification.

This study utilizes a media-based social network for studying cyberbullying, considering both photos and comments for labeling. The inclusion of multi-modal elements from text, photos, and media session metadata significantly increases the accuracy of cyberbullying detection to 87%. Dadvar et al. [10] leverage deep neural networks to identify cases of cyberbullying, applying various models, including CNN, LSTM, BLSTM, and BLSTM with attention, on datasets from Form spring, Wikipedia, and Twitter. Their study demonstrates the potential of deep neural networks in identifying cyberbullying across different platforms. Nadine et al. [11] report a high accuracy of 91% using the MySpace.com dataset and a naive Bayes-based learning model. They employ a small dataset to train a neural network and an SVM classifier, achieving notable success in identifying cyberbullying. To sum up, the many methods and tactics discussed demonstrate how successful the strategies used to identify cyberbullying in Bangla are. Researchers have looked at a wide range of tools and approaches, such as deep learning algorithms, novel word embedding models, and NLP techniques, to solve this urgent problem. Even while every strategy has its own advantages and disadvantages, taken as a whole, they support the continuous efforts to successfully tackle cyberbullying.

2.3 Comparative Analysis and Summary

It took a lot of work to navigate the deep learning model space, especially considering the growing demand for these models. Numerous parallel projects encountered challenges,

witnessing less-than-optimal model outcomes and constrained accuracy. Achieving peak accuracy in dataset detection necessitated the adoption of a unique deep learning model. This undertaking involved deploying cutting-edge hardware infrastructure to ensure the efficient execution of these models. The methodology encompassed conducting meticulous computations to assess categorization rates, with the incorporation of high-end GPUs enabling the execution of intricate models, albeit with a potential extension of runtime.

2.4 Scope of the Problem

Cyberbullying is a worldwide problem in recent worlds. This presents a complex landscape, involving language-specific hurdles, cultural influences, various online platforms, and a diverse demographic. A comprehensive approach is essential to comprehend and tackle these facets effectively. Taking into account the intricacies of language, cultural dynamics, and the diverse array of online platforms allows for the development of holistic solutions. Consequently, this helps to lessen the incidence of cyberbullying and creates a safer online space for those who communicate in English. Achieving peak accuracy in dataset detection necessitated the adoption of a unique deep learning model. This undertaking involved deploying cutting-edge hardware infrastructure to ensure the efficient execution of these models. The methodology encompassed conducting meticulous computations to assess categorization rates, with the incorporation of high-end GPUs enabling the execution of intricate models, albeit with a potential extension of runtime.

2.5 Challenges

The identification of cyberbullying in English is hampered by language complexities, dynamic character of cyberbullying behaviors. Developing language-specific algorithms, taking into account cultural settings, adjusting to platform changes, and staying up to date on new developments in cyberbullying are all necessary to overcome these challenges. These difficulties must be addressed in order to properly identify and prevent cyberbullying in the English language and to enhance English language communicators. The many methods and tactics discussed demonstrate how successful the strategies used to identify

Cyberbullying in Bangla are. Researchers have looked at a wide range of tools and approaches, such as deep learning algorithms, novel word embedding models, and NLP techniques, to solve this urgent problem. Even while every strategy has its own advantages and disadvantages, taken as a whole, they support the continuous efforts to successfully tackle cyberbullying.

CHAPTER 3

RESEARCH METHODOLOGY

Proposed Methodology Flow chart:

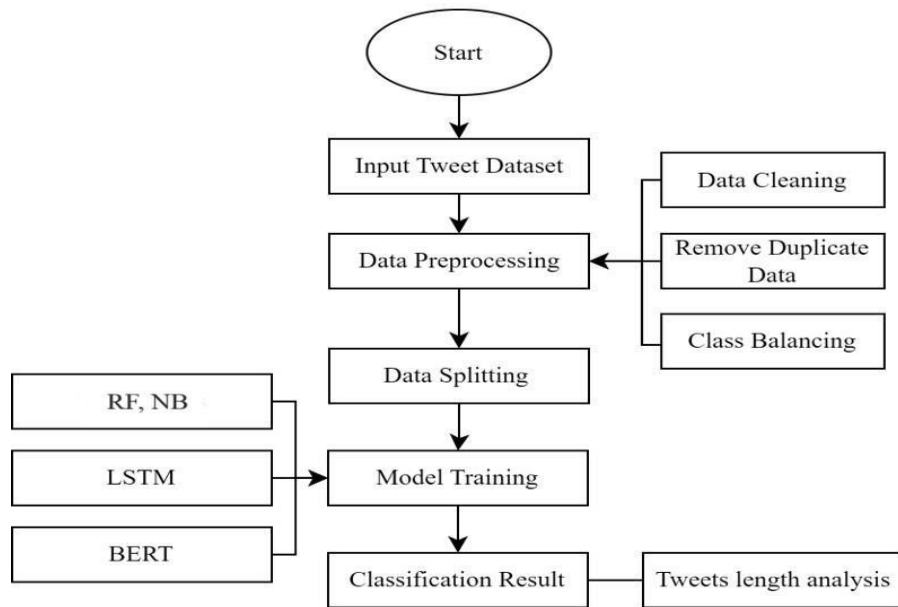


Figure 3.1: The Flow Chart of Methodology

3.1 Research Subject and Instrumentation

Tools like Jupyter Notebook, Google Colab, and Anaconda were essential to our workflow since they were essential to our data processing and model training procedures. Python's vast library ecosystem and adaptability made code creation and execution easy, resulting in productive operations. Furthermore, browser-based coding features improved accessibility and collaborative activities. Examples of these platforms include Jupyter Notebook, Google Colab, and Anaconda. By utilizing this all-inclusive combination of tools and resources, we were able to push the limits of dataset correctness and provide reliable findings in our quest for perfection.

The study makes use of frameworks and programming languages like TensorFlow and Python to create deep learning models and related techniques. Deep learning tasks are quite complicated, and computational infrastructure—which includes GPUs for rapid processing—is essential to successfully addressing them.

Here we will discuss about the tools were required to implement the technique of detecting cyberbullying in social media platforms.

Here's an overview of the materials used in Anaconda, Jupyter Notebook, NumPy, Pandas, Matplotlib, and Seaborn:

Anaconda: Anaconda is a distribution of the Python and R programming languages for scientific computing, aimed at simplifying package management and deployment.

- Anaconda Navigator: A desktop graphical user interface (GUI) that allows users to manage environments, packages, and channels easily.
- Conda: A package manager that installs, runs, and updates packages and their dependencies.
- Python Interpreter: The core Python programming language, including standard libraries and additional scientific computing packages.
- Jupyter Notebook: An interactive computing environment for creating and sharing documents that contain live code, equations, visualizations, and narrative text.

Jupyter Notebook: Jupyter Notebook is an open-source web application that allows you to create and share documents containing live code, equations, visualizations, and narrative text.

- Markdown Cells: For text-based documentation and explanation.
- Code Cells: For writing and executing Python code interactively.
- Output Cells: Display the output of executed code, including text, images, plots, and interactive widgets.
- Kernel: The computational engine that executes the code contained in the notebook.

NumPy: NumPy is a fundamental package for scientific computing in Python. It provides support for arrays, mathematical functions, linear algebra operations, random number generation, and more.

- ndarray: An efficient multidimensional array object that supports mathematical operations.
- Mathematical functions: Functions for array manipulation, mathematical operations, linear algebra, Fourier transforms, and random number generation.

Pandas: Pandas is a powerful data manipulation and analysis library for Python. It provides data structures and functions for efficiently working with structured data.

- DataFrame: A two-dimensional labeled data structure with columns of potentially

different data types, similar to a spreadsheet or SQL table.

- Series: A one-dimensional labeled array capable of holding any data type.
- Data manipulation functions: Functions for reading/writing data, cleaning, filtering, transforming, aggregating, and merging datasets.
- Matplotlib: Matplotlib is a comprehensive library for creating static, animated, and interactive visualizations in Python. It provides a MATLAB-like interface for plotting.
 - Figure and Axes objects: Components of the plot where data and visual elements are rendered.
 - Various types of plots: Including line plots, scatter plots, bar plots, histogram plots, pie charts, etc.

Seaborn: Seaborn is a statistical data visualization library built on top of Matplotlib. It provides a high-level interface for drawing attractive and informative statistical graphics.

- Simplified plotting functions: Functions that allow users to create complex plots with minimal code.
- Statistical visualization functions: Functions for visualizing statistical relationships, distributions, and categorical data.
- Styling and customization options: Options for customizing the appearance of plots, including colors, styles, and themes.

3.2 Data Collection Procedure

After being downloaded from Kaggle [1], the dataset was nearly ready for usage, but two essential elements were still lacking: the textual data and the labels that went with it. Six different categories made up this supervised dataset, which offered a thorough representation of the material. We have assumed the 80:20 ratio of train and test the dataset for best evaluation. This division designated 80% of the data for in-depth analysis and training, and 20% of the data for the test portion.

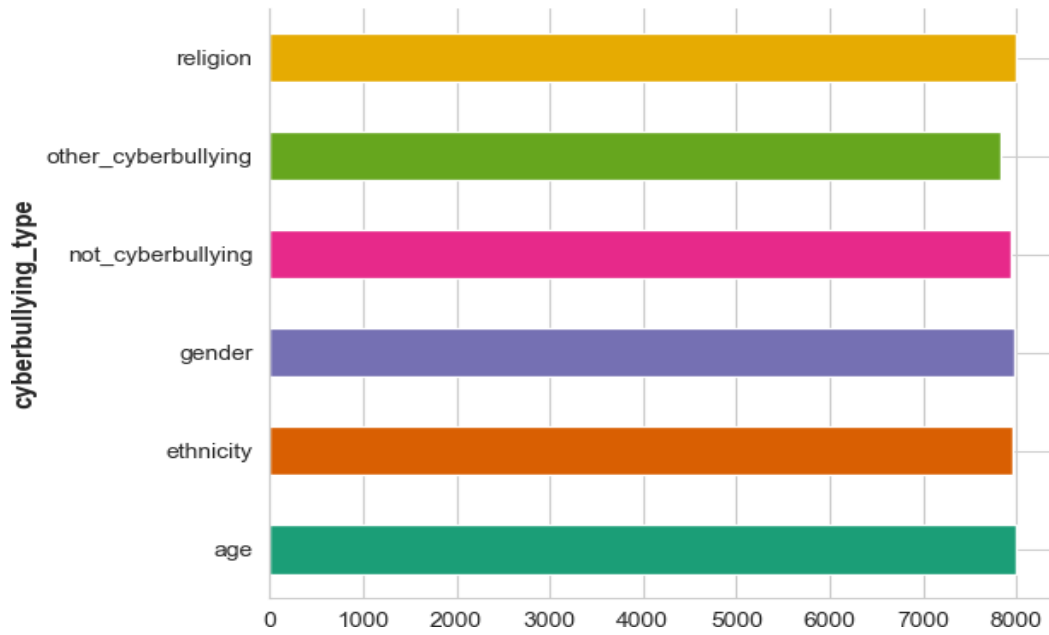


Figure 3.2: Count plot of Target Column

We have checked the dataset for missing value removal and duplicate rows. We have removed the missing values and duplicate values from the dataset.

Columns	Description	Count
text	Contains text values	47692
sentiments	Contains categorical values	47692

Table 3.1: Explanation of the Dataset

Count of tweets with less than 10 words

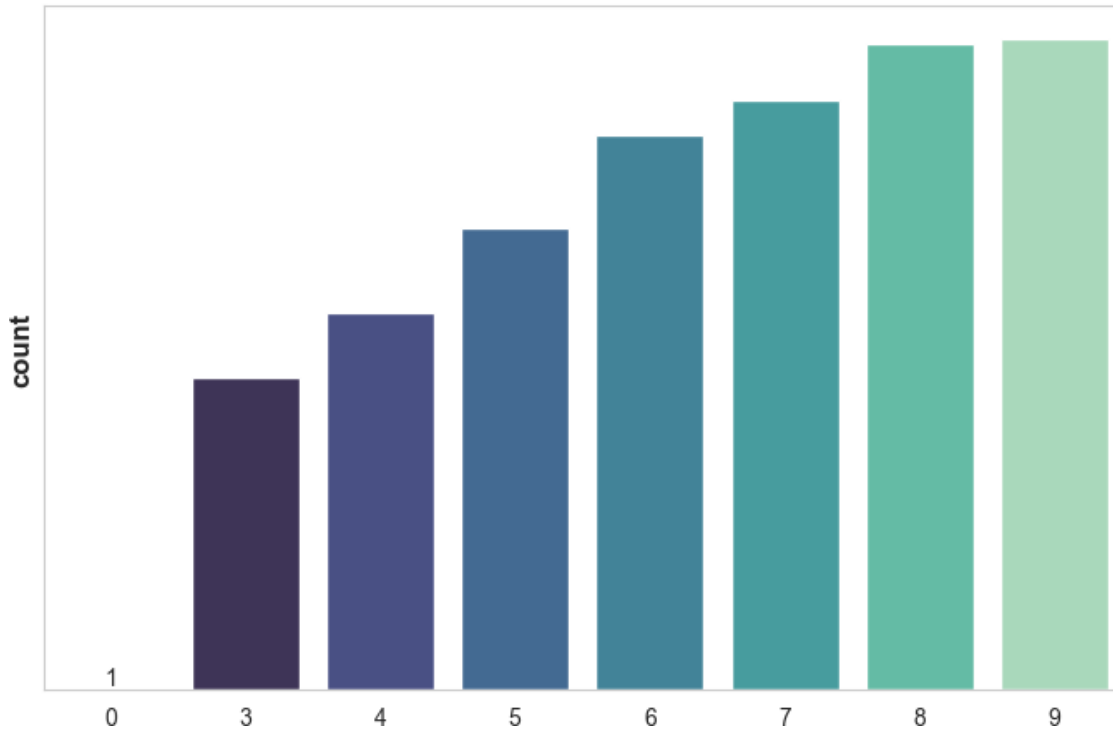


Figure 3.3: Tweets less than 10 words

Count of tweets with high number of words

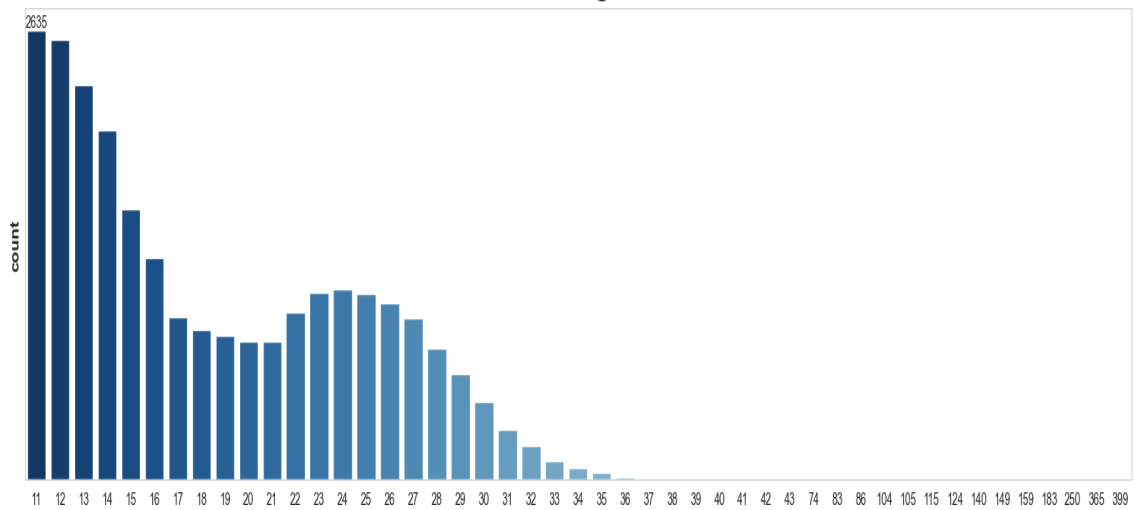


Figure 3.4: Tweets with high number of words

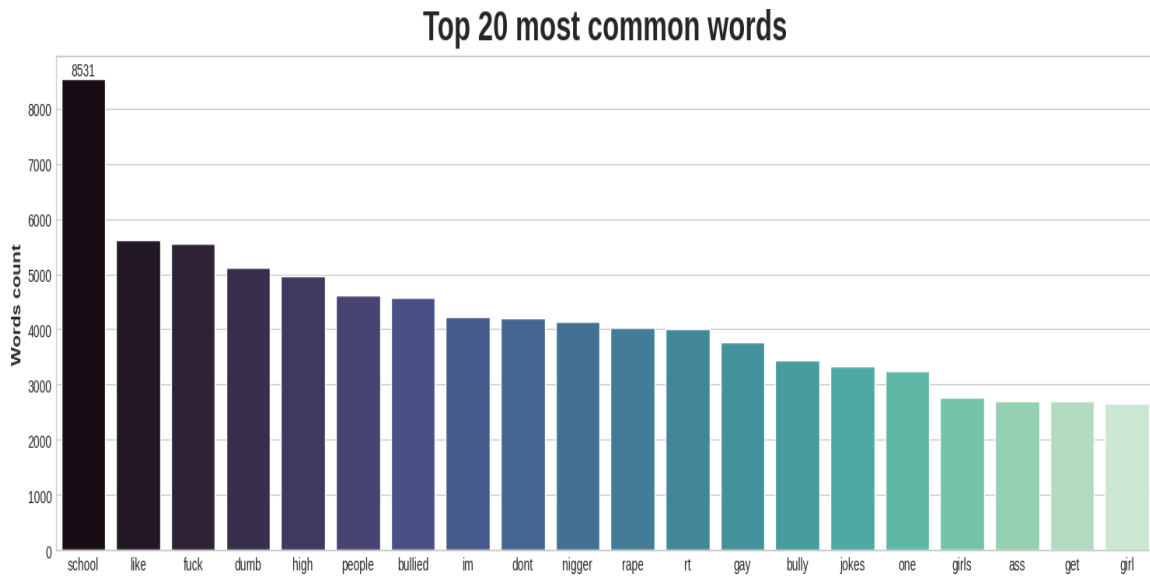


Figure 3.5: Tweets with top 20 most common of words

3.2.1 Categorical Data Encoding

We have used categorical data to numerical data for model's best results. We also used label encoder to do the task easily. Given that, machine learning algorithms demand numerical data as input and output by nature, this system was critical to our study.

3.2.2 Missing Value Imputation

Imputed values are utilized in this process to fill in any missing or incomplete data that was found during the examination of data from other datasets. It is important to note that this kind of imputation was not required because we have reduced the blank values from the dataset.

3.2.3 Handling Imbalanced Data

By methodically adding new instances, this technique may be used to achieve data balance and modify the class distribution within a dataset. Improving minority data representation is the main objective, with the complete dataset serving as input.

3.3 Statistical Analysis

An essential part of every research study is the analysis segment, which focuses on creating and assessing the algorithms used. Since we decided to deal with an Excel file format in

this instance, a few prior actions needed guarantee for the collected dataset usefulness. These steps included gathering optimum datasets and carefully handle it, both were necessary for our research to be carried out successfully. Four distinct kinds of algorithms—LSTM, BERT, RF, and NB—were used in this investigation. Achieving an astounding 87% accuracy, the BERT was tested at RF's 85% accuracy.

3.4 Proposed Methodology

For the best result of cyberbullying detection, we have used many types of algorithms with machine learning and deep learning encompassing BERT, LSTM, Random Forest Classifier (RF) and Naïve Bayes Classifier (NB) algorithms. Our extensive algorithm evaluation and data analysis relied heavily on these classifiers.

3.4.1 Long Short-Term Memory (LSTM)

The Long Short-Term Memory (LSTM) model is a type of recurrent neural network (RNN) designed to effectively learn and retain long-term dependencies in sequential data. Unlike traditional RNNs, which struggle with vanishing and exploding gradient problems, LSTMs address these issues through a unique cell structure featuring input, output, and forget gates. These gates regulate the flow of information, allowing the network to maintain, update, or forget specific pieces of data over time. This gating mechanism enables LSTMs to capture complex temporal patterns and long-range dependencies, making them particularly well-suited for tasks such as time series forecasting, natural language processing, and speech recognition [19-20].

3.4.2 Random Forest (RF)

As a powerful and versatile machine learning algorithm, the Random Forest classifier receives a lot of praise for its ability to perform well in functions building a collection of decision trees, each built using a distinct subset of the feature set and training data. By adding variation, this method reduces over fitting and improves generalization. In regression, the Random Forest calculates the average prediction, but in classification, it combines decision tree results by majority voting. The technique reduces over fitting when compared to single decision trees since it uses bagging and feature bagging. It is useful because it can be used to complex and noisy datasets and has a low sensitivity to

hyper parameter adjustment. Moreover, the Random Forest finds significant traits and provides information about how they affect predictive capability. It is the preferred choice in image analysis, finance, and healthcare due to its scalability, resilience, and adaptability; nevertheless, real-time or resource-constrained applications may be impacted by its increased computing complexity and cost. In spite of this, the Random Forest continues to be a trustworthy workhorse in machine learning, providing insightful insights and precise predictions in a variety of problem-solving settings [16, 17]. The algorithm's effectiveness is demonstrated by finding the average of two decision tree techniques.

3.4.3 BERT

In 2018, Google researchers developed the groundbreaking BERT natural language processing (NLP) technology. It is a significant advancement in natural language processing, particularly in terms of contextualized word representation pre-training. One of BERT's main achievements is its bidirectional way for understanding context inside a given text. Using a transformer design, BERT may continuously evaluate words that come before and after, in contrast to previous language models that read text either from left to right or from left to left. By capturing a deeper understanding of the environment in which words occur, BERT's bidirectional processing allows it to more precisely represent sentence structures and word meanings. Pre-training on large text data corpora and subsequent fine-tuning on specific downstream NLP tasks are the reasons for BERT's remarkable performance. During pre-training, BERT learns to predict missing words in a sentence by utilizing the surrounding context. By using this process, BERT can produce intricate, contextually-aware word embeddings that faithfully capture the nuances of language use. One of the unique features of BERT is its ability to capture contextual information at different levels of detail. By combining attention processes, BERT is able to dynamically identify the relative worth of words in a sentence depending on their contextual importance. This allows BERT to create contextually aware word representations that are very informative. BERT has demonstrated remarkable performance in a range of natural language processing (NLP) tasks, including text classification, named entity recognition, sentiment analysis, question answering, and language translation. Due to its versatility and effectiveness, it is a cornerstone of many state-of-the-art NLP systems and applications. Furthermore, BERT's pre-trained models—like BERT-base and BERT

large—are now publically accessible thanks to Google, enabling researchers and developers to use it for a range of natural language processing applications. Moreover, BERT has undergone several extensions and alterations that have facilitated the advancement of contextualized word embeddings and NLP research.

3.4.4 Naïve Bayes

The core of it is the Bayes theorem, a basic idea in probability theory that determines the likelihood of a hypothesis based on observable data. Based on the observed attributes, Naive Bayes determines the likelihood that a given data point belongs to a specific class in the context of classification. The "naive" assumption of feature independence, which maintains that each feature independently contributes to the likelihood of the data point belonging to a given class, sets Naive Bayes apart from other classifiers and makes probability computation simpler. Despite being simple, this assumption frequently holds up well in real-world situations, particularly in jobs involving text categorization and other areas where feature correlations are not statistically significant. The simplicity of Naive Bayes stems from its straightforward approach to modeling. Rather than learning complex relationships between features, Naive Bayes directly estimates probabilities from the training data. This simplicity makes Naive Bayes particularly well-suited for situations with limited training data, where more complex models may struggle due to overfitting. Moreover, Naive Bayes classifiers are computationally efficient, making them ideal for large datasets and high-dimensional feature spaces. The algorithm's efficiency arises from its ability to compute probabilities independently for each feature, resulting in a linear time complexity with respect to the number of features. There are several variations of naive Bayes classifiers, each designed to handle distinct kinds of data. For classification tasks involving continuous-valued data, the Gaussian Naive Bayes version is appropriate since it makes the assumption that continuous features have a Gaussian (normal) distribution. This variant is commonly used in applications such as medical diagnosis, where features like patient age or blood pressure may follow a normal distribution. On the other hand, the Multinomial Naive Bayes variant is designed for text classification tasks, where features are words or keywords that often occur in texts. Because it assumes that features have a multinomial distribution, this technique performs well in applications such as document categorization, sentiment analysis, and spam detection. Moreover, the binary feature

vectors that represent whether or not words are present in texts are the only application for which the Bernoulli Naive Bayes variant is intended. It is comparable to the version that is multinomial. This version is often used in binary classification applications such as spamfiltering and emotion analysis. The cyberbullying detection methodology applied to the English dataset encompasses several crucial steps. The first step is to compile a dataset of English social media comments. After that, preliminary tasks including removing pointless columns, handling missing data, and cleaning the text are completed. To help with model training, the target column is numerically mapped after the dataset is split into training and testing sets at an 80-20 ratio [12]. Categorical data are converted to numerical values so that they can be further analyzed. To prepare text data for entry into different models, such as BERT, LSTM, Random Forest, and Naïve Bayes, it is tokenized and padded. The prepared data is used to train these models. Evaluation metrics, which include accuracy, precision, recall, and F1 score for models utilizing probability estimates, are used to assess how well the models work. Choosing the top-performing model is the last stage, which is determined by the evaluation's findings [13]. For researchers looking to create efficient cyberbullying detection algorithms, this thorough technique gives a solid framework for identifying cyberbullying in the English dataset [14].

3.5 Implementation Requirements

To ensure the effectiveness of our proposed model for cyberbullying detection, our methodology began with the acquisition of necessary datasets. A stringent data cleaning process was implemented, employing various filtering methods to guarantee the dataset's integrity. Following this, crucial data preprocessing steps were executed, including the application of the Scaler Transform to convert categorical data into numerical values. Then for the best output, we have separated the dataset into train and test part, giving a strong basis for evaluating algorithms. The chosen algorithms were implemented, and their results underwent a comprehensive assessment. To improve the detection accuracy, ensemble methods were used, and the results of these ensemble algorithms were thoroughly verified by hyperparameter tweaking, which required a careful analysis of the underlying mathematical models. Then came the data analysis stage, which included creating a personalized detection approach and learning a mode

CHAPTER 4

EXPERIMENTAL RESULTS AND DISCUSSION

4.1 Experimental Setup

The methodology used in this work is based on supervised learning and consists of distinct training and assessment phases. The testing dataset was used to extensively assess the model's performance, while the training dataset was utilized to begin developing the deep learning model. The following parts will provide a brief but thorough explanation of the deep learning method used in this study.

4.2 Experimental Results & Analysis

We turned our focus to examining the performance of our proposed model in the area of cyberbullying detection and comparing it with other models that are now in use. We methodically used a variety of performance evaluation techniques using never-before-seen data to assess overall performance. An analytical analysis based on our experimental results—more particularly, deep learning models customized for cyberbullying datasets—is summarized in this section. Applying our chosen dataset was the initial phase in our procedure, during which we carefully managed any missing or erroneous values to preserve the integrity of the dataset. After that, a large range of algorithms were shown, and their effectiveness was carefully investigated. To thoroughly evaluate the effectiveness of these algorithms, key metrics such as the F-1 Score, Accuracy, Precision, Recall, and Confusion Matrices were employed. Both the conventional techniques and our suggested models were evaluated. Specifically designed for our dataset, the algorithms used in our research were Random Forest Classifier (RF), Naïve Bayes Classifier (NB), LSTM, and BERT. Several ensemble strategies were investigated to improve our evaluation even further, and Confusion Matrices offered insights for both datasets. The goal of this thorough approach was to provide a full knowledge of the efficacy and performance of the various models being considered in the context of cyberbullying detection.

4.2.1 Accuracy

The concept of accuracy is examined in this section, with particular attention paid proved to be precise or accurate. By comparing the model's predictions to actual measurements made in the real world, accuracy shows the prediction correctness. It is one of the simplest and most popular model assessment methods since it concentrates on just one variable and mostly deals with deliberate mistakes. A critical component of model validation and performance evaluation is ensuring the correctness.

$$Accuracy = \frac{(TruePositive+TrueNegative)}{(TruePositive+FalsePositive+TrueNegative+FalseNegative)} \dots\dots\dots(i)$$

4.2.2 Precision

Precision is the measure of the percentage of favorably predicted observations that really happened, and it is covered in this section. The real positive rate is reflected in precision, which shows the actual percentage of cases in which the model accurately predicted genuine positive outcomes. It's crucial to remember that, even though many models prefer a high recall, this characteristic can occasionally be deceptive if accuracy and other performance measures aren't taken into account.

$$Precision = \frac{TruePositive}{FalsePositive+TruePositive} \dots\dots\dots(ii)$$

4.2.3 Recall

Recall, also known as sensitivity or true positive rate, is a metric used to evaluate the performance of a classification model. It measures the proportion of actual positive instances that are correctly identified by the model. High recall indicates that the model successfully captures most of the positive cases, making it particularly useful in contexts where identifying positives is crucial, such as medical diagnosis or fraud detection. However, recall alone does not account for the number of false positives, so it is often considered alongside precision to provide a more comprehensive evaluation of a model's effectiveness, commonly in the form of the F1 score.

$$Recall = \frac{TruePositive}{TruePositive+FalseNegative} \dots\dots\dots(iii)$$

4.2.4 F-1 Score

The assessment metrics of accuracy and recall are covered in this part, with a focus on their use in evaluating a model's performance. Since they provide light on the model's overall accuracy and correctness in detecting noteworthy occurrences, the recall and accuracy ratios are crucial metrics to take into account. It is important to keep in mind that a relatively low evaluation might indicate that more adjustments are needed in order for the model to perform as intended.

$$F - 1 \text{ Score} = 2 * \frac{(Recall*Precision)}{(Recall+Precision)} \dots\dots\dots(iv)$$

```

Classification Report for Random Forest:
              precision    recall  f1-score   support

 religion      0.94      0.96      0.95      1574
    age        0.96      0.98      0.97      1566
 ethnicity    0.98      0.98      0.98      1537
    gender    0.93      0.82      0.87      1481
other_cyberbullying 0.62      0.63      0.63      1401
 not bullying 0.57      0.63      0.60      1051

 accuracy                    0.85      8610
 macro avg      0.84      0.83      0.83      8610
 weighted avg   0.85      0.85      0.85      8610

```

Figure 4.1: Experimental Results of Random Forest

The Random Forest model's categorization report offers a thorough assessment of its performance in several classifications. With scores over 0.94 for age, race, and religion, the model shows excellent precision for the majority of classes when it comes to precision, which gauges the accuracy of positive predictions. It does, however, demonstrate worse predictive accuracy for gender, non-bullying, and other types of cyberbullying, indicating a higher frequency of false positives in these areas. The model's recall refers to its ability to identify instances of each class. The main distinctions between the categories are in terms of gender, other forms of cyberbullying, and non-bullying. The model does well in identifying age, race, and religion; but, it struggles in identifying gender, exhibiting a very poor recall rate. In contrast, the F1-score, which accounts for both accuracy and recall, illustrates how the overall recall and precision of each class are balanced. The model often has high F1-scores for age, ethnicity, and religion, demonstrating its effectiveness in these areas. However, it shows lower F1-scores for gender, other cyberbullying, and not bullying, suggesting areas for potential improvement. Notwithstanding these fluctuations, the model exhibits robust performance in most classes, with an accuracy rate of 85% overall. In light of the class distribution, the weighted-average and macro-average scores offer more information about the model's overall performance. In conclusion, even though the Random Forest model shows promise in predicting certain classes—like age, ethnicity, and religion—it still needs to be improved, especially in terms of gender classification and the ability to discern between instances of bullying and other types of cyberbullying.

**Random Forest Sentiment Analysis
Confusion Matrix**

Test	religion	1512	1	3	4	43	11
	age	1	1527	0	4	19	15
	ethnicity	2	4	1504	5	3	19
	gender	5	5	3	1216	156	96
	other_cyberbullying	71	37	12	37	887	357
	not_bullying	12	12	6	44	318	659
		religion	age	ethnicity	gender	other_cyberbullying	not_bullying
		Predicted					

Figure 4.2: Confusion Matrix of Random Forest

Classification Report for Naive Bayes:				
	precision	recall	f1-score	support
religion	0.81	0.97	0.88	1574
age	0.75	0.98	0.85	1566
ethnicity	0.86	0.93	0.89	1537
gender	0.83	0.83	0.83	1481
other_cyberbullying	0.73	0.38	0.50	1401
not bullying	0.59	0.45	0.51	1051
accuracy			0.78	8610
macro avg	0.76	0.76	0.74	8610
weighted avg	0.77	0.78	0.76	8610

Figure 4.3: Experimental Results of Naïve Bayes

The classification report comprehensively assesses the effectiveness of the Naive Bayes classifier for six different classes: age, gender, religion, ethnicity, other cyberbullying, and not bullying. The F1-score, which has scores ranging from 0 to 1, is a single statistic that combines the accuracy and recall metrics. While recall assesses the proportion of correctly identified genuine positives, precision examines the accuracy of positive predictions. The precision scores indicate that the classifier performs well in predicting religion, ethnicity, and gender, with values ranging from 0.81 to 0.86. However, it struggles more with other cyberbullying and not bullying classes, achieving lower precision scores of 0.73 and 0.59, respectively. Recall scores demonstrate high performance across most classes, particularly for age, ethnicity, and religion, with values ranging from 0.93 to 0.98. However, the classifier exhibits lower recall for other cyberbullying and not bullying classes, indicating difficulty in correctly identifying instances of these classes. The F1-scores, which balance precision and recall, highlight the classifier's overall effectiveness in each class. The weighted average F1-score of 0.76 and accuracy of 0.78 suggest that the classifier performs reasonably well across all classes, albeit with some variation in performance depending on the class characteristics. In summary, while the Naive Bayes classifier demonstrates strong performance in certain classes, particularly those with clear distinctions in features, it shows room for improvement in accurately classifying instances of cyberbullying-related classes.

**Naive Bayes Sentiment Analysis
Confusion Matrix**

Test	Predicted					
	religion	age	ethnicity	gender	other_cyberbullying	not_bullying
religion	1529	10	10	10	10	5
age	8	1540	5	4	3	6
ethnicity	36	48	1434	7	4	8
gender	44	35	38	1225	85	54
other_cyberbullying	180	245	76	121	528	251
not_bullying	88	186	108	105	94	470

Figure 4.4: Confusion Matrix of Naïve Bayes

Classification Report for Bi-LSTM :

	precision	recall	f1-score	support
religion	0.96	0.94	0.95	1574
age	0.97	0.97	0.97	1565
ethnicity	0.99	0.95	0.97	1537
gender	0.89	0.86	0.87	1481
other_cyberbullying	0.65	0.55	0.59	1400
not bullying	0.54	0.74	0.63	1051
accuracy			0.84	8608
macro avg	0.83	0.83	0.83	8608
weighted avg	0.85	0.84	0.85	8608

Figure 4.5: Experimental Results of LSTM

An extensive study of the LSTM model's performance across several classes can be found in the classification report. Measures like as precision, recall, and F1-score are employed to assess how well the model can categorize examples of each class. With accuracy ranging from 0.54 to 0.99 and recall ranging from 0.55 to 0.97 across many classes, the model performs admirably overall. Specifically, the model achieves good accuracy, recall, and F1-scores when it comes to predicting classes such as 'age' and 'ethnicity'. However, its performance in classifying cases of 'other cyberbullying' and 'not bullying' is low, suggesting challenges with accurately recognizing these classifications, particularly with respect to accuracy. In spite of this, the model obtains an overall accuracy of 0.84, meaning

that most occurrences in all classes can be accurately classified by it. The model's strong performance across all classes is further supported by the macro and weighted average F1-scores, both of which are continuously above 0.83. Although the model performs well overall in most classes, there may be space for improvement, especially in correctly identifying cases of "other cyberbullying" and "not bullying".

**PyTorch Bi-LSTM Sentiment Analysis
Confusion Matrix**

Test	religion	1472	2	2	13	78	7
	age	2	1524	2	3	20	14
	ethnicity	10	7	1459	22	6	33
	gender	2	5	0	1277	95	102
	other_cyberbullying	45	27	3	66	763	496
	not bullying	3	12	1	57	203	775
		religion	age	ethnicity	gender	other_cyberbullying	not bullying
		Predicted					

Figure 4.6: Confusion Matrix of LSTM

Classification Report for BERT :

	precision	recall	f1-score	support
religion	0.97	0.97	0.97	1574
age	0.98	0.98	0.98	1566
ethnicity	0.99	0.99	0.99	1538
gender	0.89	0.88	0.89	1481
other_cyberbullying	0.69	0.65	0.67	1398
not bullying	0.63	0.69	0.66	1048
accuracy			0.87	8605
macro avg	0.86	0.86	0.86	8605
weighted avg	0.87	0.87	0.87	8605

Figure 4.7: Experimental Results of BERT

The classification report for the BERT model provides a comprehensive overview of its performance across multiple classes. It shows high precision, recall, and F1-score values for most classes, indicating strong predictive capabilities. Specifically, the model demonstrates excellent performance in classifying instances related to religion, age,

ethnicity, and gender, with precision, recall, and F1-score values ranging from 0.97 to 0.99. However, it exhibits slightly lower performance in distinguishing instances of other cyberbullying and not bullying, with precision, recall, and F1-score values around 0.69 and 0.65, respectively. Despite these minor discrepancies, the overall accuracy of the model is impressive at 0.87, suggesting that it effectively classifies the majority of instances across all classes. The macro and weighted average scores further reinforce the model's robustness and generalizability, with both averaging above 0.86. Overall, the evaluation indicates that the BERT model performs well across various classes, with particularly strong performance in identifying instances related to religion, age, ethnicity, and gender, while maintaining respectable accuracy and generalization capabilities across the entire dataset.

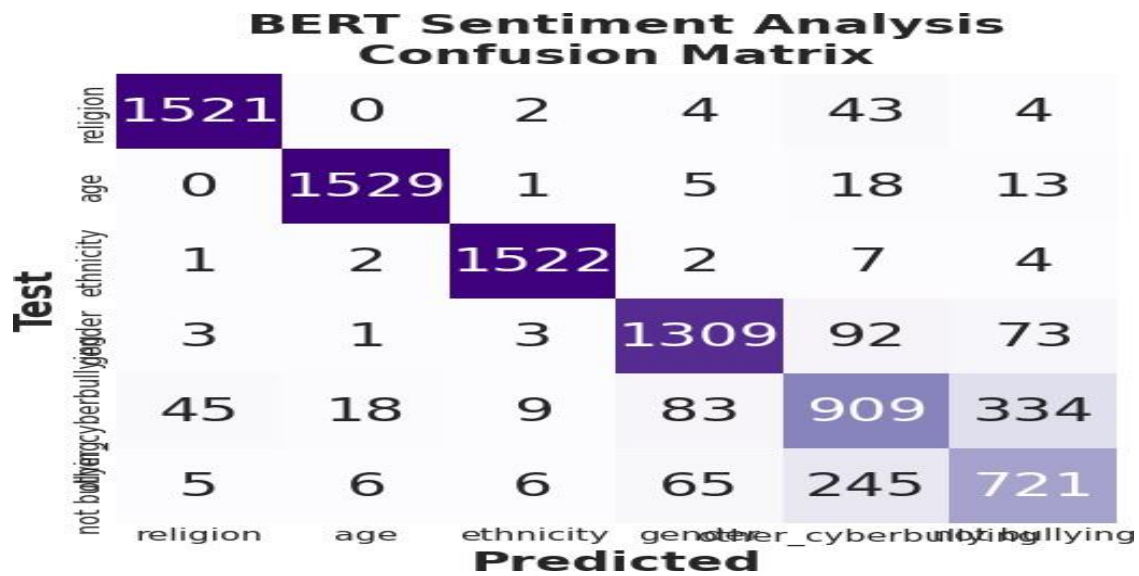


Figure 4.8: Confusion Matrix of BERT

4.3 Discussion

We shall clarify the evaluation framework for our suggested model at this point. We used accuracy, precision, recall, and the F-1 score as our evaluation criterion [22]. The study shows that the BERT performed better than the others, with an accuracy rate of 87%. Furthermore, another strategy that made use of LSTM produced an accuracy of 84%. Ensemble models, including Naïve Bayes (NB), and Random Forest, were also employed, with hyper-parameter tuning optimizing their performance. Notably, the LSTM and BERT outperformed other models, attaining the highest accuracy rate of 87% in cyberbullying detection, as confirmed by recent experimental inquiries evaluating these findings.

CHAPTER 5

IMPACT ON SOCIETY, ENVIRONMENT AND SUSTAINABILITY

5.1 Impact on Society

The ramifications of cyberbullying within the English-speaking society are far-reaching, affecting diverse aspects of individuals' lives. The psychological well-being of victims is profoundly impacted, leading to heightened levels of stress, anxiety, and depression. Academic pursuits and educational achievements also bear the brunt of cyberbullying, with victims often grappling with challenges such as difficulty concentrating, diminished motivation, and increased absenteeism. Social relationships undergo strain, as individuals subjected to cyberbullying may opt for isolation due to fear or shame, resulting in feelings of loneliness and social exclusion. Furthermore, the reputation and image of victims are at risk, influencing both personal and professional spheres of their lives. The pervasive nature of cyberbullying in the English language underscores the urgency of implementing proactive measures to mitigate its impact and foster a safer online environment.

5.2 Impact on Environment

It is crucial to emphasize that cyberbullying predominantly manifests in the digital sphere rather than the physical environment. However, its repercussions on the online landscape are profound. Cyberbullying establishes an atmosphere characterized by fear, hostility, and negativity, discouraging individuals from expressing themselves freely and participating in online discussions. The prevalence of cyberbullying undermines the cultivation of a supportive and inclusive online community, resulting in the erosion of trust, empathy, and mutual respect. Initiatives aimed at tackling cyberbullying and cultivating a positive online environment are imperative to facilitate healthy digital interactions and enhance the well-being of individuals navigating the digital space.

5.3 Ethical Aspects

The identification of cyberbullying in English poses significant moral questions. Protecting sensitive data and personal information requires maintaining research participants' privacy and confidentiality. In order to respect the rights and permission of the people whose data is being used, social media comments must be collected and ethically assessed. Fairness and impartiality must be maintained in order to stop stigmatization and prejudice against

particular groups. The implementation of algorithms and models must be transparent in order to allow for accountability and inspection. Preserving the rights and well-being of those affected by cyberbullying necessitates the application of ethics and conscientious conduct. This increases the English-speaking community's comprehension of the problem and preventative measures.

5.4 Sustainability Plan

In order to ensure long-lasting effectiveness and impact, long-term methods and procedures are needed to provide a sustainable strategy to cyberbullying detection in the English language. In order to carry out awareness campaigns, policies, and support systems, this entails forming partnerships with important stakeholders including social media platforms, educational institutions, and community groups. In order to adjust to changing cyberbullying strategies, it is imperative that detection methods and algorithms be continuously monitored and evaluated. Over time, the accuracy and relevance of the detection model may be improved with regular updates and enhancements that take into account user input and technological breakthroughs. Additionally, by encouraging, it may be able to make the internet world safer and healthier for English-speaking individuals.

CHAPTER 6

SUMMARY, CONCLUSION, RECOMMENDATION AND IMPLICATION FOR FUTURE RESEARCH

6.1 Summary of the Study

Aim of the study is to use English-specific deep learning algorithms to identify cyberbullying. With a focus on addressing the language and cultural peculiarities common in the English-speaking population, the goal is to close the current research gap regarding cyberbullying detection in English environments. Through the use of deep learning models—BERT and LSTM cells, in particular—the research aims to develop a strong model that can recognize instances of cyberbullying in English text, particularly in comments on social media. The findings of this research should significantly aid in the development of language-specific instruments, guidelines, and interventions, keeping in mind the unique challenges associated with cyberbullying in the English language, such as linguistic complexity and cultural effects. In response, these work to stop and lessen instances of cyberbullying, creating a more secure and welcoming online space for people using English to communicate.

6.2 Conclusions

Cyberbullying detection in the English language, utilizing deep learning techniques, significantly contributes to addressing the critical issue of cyberbullying in English contexts. The development of a tailored deep learning model for English acknowledges and accounts for the linguistic and cultural nuances that shape cyberbullying dynamics within this specific community. In order to close the research gap in cyberbullying detection for English-language learners, the study highlights the need for inclusive strategies that accommodate various linguistic groups. It also emphasizes how successful deep learning techniques are, especially LSTM cell-based recurrent neural networks. By focusing on English, the study hopes to empower those who use the language to communicate, guaranteeing fair treatment and defense against online harassment. The

Study's findings can be used to develop language-specific tools, policies, and initiatives that will reduce and eliminate cyberbullying among English-speaking individuals. The study highlights how crucial it is to keep an eye on, assess, and improve detection algorithms in order to keep up with the latest developments in cyberbullying strategies. It also emphasizes how important it is to use awareness campaigns and educational programs to foster conduct. Among essence, addressing cyberbullying among communities of English speakers. It emphasizes how important it is for academics, legislators, social media companies, academic institutions, and community organizations to work together to create an online space that is safer and more welcoming for everyone, regardless of the language they speak.

6.3 Implication for Further Study

Future research endeavors in the realm of cyberbullying detection in the English language should prioritize the refinement and expansion of datasets to encompass a broader spectrum of cyberbullying instances and contextual factors. Sophisticated deep learning models, particularly transformer-based architectures like BERT or GPT, can increase the efficacy and precision of cyberbullying detection systems. The exploration of model transferability, assessing the effectiveness of models trained on larger, multilingual datasets when applied to English, presents an avenue for investigation. Additionally, longitudinal studies can offer a deeper understanding of the evolving nature of cyberbullying and its enduring impact on individuals over time, thereby furnishing valuable insights for the development of proactive interventions and support systems. These directions underscore the continuous evolution and enhancement of cyberbullying detection strategies tailored.

Reference

- [1] J. Wang, K. Fu, C.T. Lu, "SOSNet: A Graph Convolutional Network Approach to Fine-Grained Cyberbullying Detection," Proceedings of the 2020 IEEE International Conference on Big Data (IEEE BigData 2020), December 10-13, 2020
- [2] L. Yang and A. Shami, "On hyperparameter optimization of machine learning algorithms: theory and practice," Neurocomputing, vol. 415, pp. 295–316, 2020.
- [3] Samghabadi, Niloofar Safi, et al." Detecting nastiness in social media." Proceedings of the First Workshop on Abusive Language Online. 2017.
- [4] Yao, Mengfan, Charalampos Chelmiss, and Daphney? Stavroula Zois." Cyberbullying ends here: Towards robust detection of cyberbullying in social media." The World Wide Web Conference.2019.
- [5] Huang, Qianjia, Vivek Kumar Singh, and Pradeep Kumar Atrey." Cyberbullying detection using social and textual analysis." Proceedings of the 3rd International Workshop on Socially-Aware Multimedia. 2014.
- [6] T. Bin Abdur Rakib, L. K. Soon, in Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics) (Springer Verlag, 2018), vol. 10751 LNAI, pp. 180–189.
- [7] Y. N. Silva, D. L. Hall, C. Rich, BullyBlocker: toward an interdisciplinary approach to identify cyberbullying. Social Network Analysis and Mining. 8 (2018), doi:10.1007/s13278-018-0496-z.
- [8] E. Raisi, B. Huang, Weakly supervised cyberbullying detection with participant-vocabulary consistency. Social Network Analysis and Mining. 8 (2018), doi:10.1007/s13278-018-0517-y.
- [9] Homa Hosseinmardi, Sabrina Arredondo Mattson, Rahat Ibn Rafiq, Richard Han, Qin Lv, Shivakant Mishra. (2015). Detection of Cyberbullying Incidents on the Instagram Social Network."
- [10] Dadvar, Maral Eckert, Kai. (2018). Cyberbullying Detection in Social Networks Using Deep Learning Based Models; A Reproducibility Study. 10.13140/RG.2.2.16187.87846.
- [11] Nandhini, B. Sri, and J. I. Sheeba." Cyberbullying detection and classification using information retrieval algorithm." Proceedings of the 2015 International Conference on Advanced Research in Computer Science Engineering Technology (ICARCSET 2015).
- [12] GeeksforGeeks, StandardScaler, MinMaxScaler and RobustScaler techniques– ML, Accessed: January, 2022, Available: <https://www.geeksforgeeks.org/standardscaler-minmaxscaler-and-robustscaler-techniques-ml/>
- [13] StackExchange, when not to use cross validation? Accessed: January 9, 2022, Available: <https://stats.stackexchange.com/questions/320154/when-not-to-use-cross-validation>
- [14] Han, Jun, and Claudio Moraga. "The influence of the sigmoid function parameters on the speed of ackpropagation learning." In International workshop on artificial neural networks, pp. 195-201. Springer, Berlin, Heidelberg, 1995.

- [15] Zarrin Tasnim, S. Chakraborty, F. M. J. M. Shamrat, A. N. Chowdhury, H. Alam Nuha, A. Karim, m abrina Binte Zahir, and M. Billah. "Deep learning predictive model for colon cancer patient using CNN-based classification." *Int. J. Adv. Comput. Sci. Appl* 12 (2021).
- [16] Russell, Stuart, and Peter Norvig. "Artificial intelligence: a modern approach." (2002).
- [17] zam, Md Shafiul, Aishe Rahman, SM Hasan Sazzad Iqbal, and Md Toukir Ahmed. "Prediction of liver iseases by using few machines learning based approaches." *Aust. J. Eng. Innov. Technol* 2, no. 5 (2020): 85-90
- [18] Hasib, Khan Md, Md Ahsan Habib, Nurul Akter Towhid, and Md Imran HossainShowrov. "A Novel deep Learning based Sentiment Analysis of Twitter Data for US Airline Service." In *2021 International Conference on Information and Communication Technology for Sustainable Development (ICICT4SD)*, pp. 450-455. IEEE, 2021.
- [19] Greff, Klaus, Rupesh K. Srivastava, Jan Koutník, Bas R. Steunebrink, and Jürgen Schmidhuber. "LSTM: A search space odyssey." *IEEE transactions on neural networks and learning systems* 28, no. 10 (2016): 2222-2232.
- [20] Devi, Usha M., and A. Marimuthu. Donor-Recipient for Liver Transplantation Using CNN and LSTM Deep Learning Techniques (No. 4923). EasyChair, 2021.
- [21] Shamima Akter, F. M. Shamrat, Sovon Chakraborty, Asif Karim, and Sami Azam. "COVID-19 detection using deep learning algorithm on chest X-ray images." *Biology* 10, no. 11 (2021): 1174.
- [22] Rakibul Islam, Abhijit Reddy Beeravolu, Md Al Habib Islam, Asif Karim, Sami Azam, and Sanzida Akter Mukti. "A Performance Based Study on Deep Learning Algorithms in the Efficient Prediction of Heart Disease." In *2021 2nd International Informatics and Software Engineering Conference (IISEC)*, pp. 1-6. IEEE, 2021

0242220004213014

[Handwritten Signature]
28/12/24

ORIGINALITY REPORT

14% SIMILARITY INDEX	11% INTERNET SOURCES	8% PUBLICATIONS	7% STUDENT PAPERS
--------------------------------	--------------------------------	---------------------------	-----------------------------

PRIMARY SOURCES

- 1** dspace.daffodilvarsity.edu.bd:8080 Internet Source **2%**
- 2** Submitted to Daffodil International University Student Paper **2%**
- 3** Dinesh Goyal, Bhanu Pratap, Sandeep Gupta, Saurabh Raj, Rekha Rani Agrawal, Indra Kishor. "Recent Advances in Sciences, Engineering, Information Technology & Management - Proceedings of the 6th International Conference "Convergence2024" Recent Advances in Sciences, Engineering, Information Technology & Management, April 24-25, 2024, Jaipur, India", CRC Press, 2025 Publication **1%**
- 4** Submitted to Higher Education Commission Pakistan Student Paper **<1%**
- 5** R. N. V. Jagan Mohan, Vasamsetty Chandra Sekhar, V. M. N. S. S. V. K. R. Gupta. "Algorithms in Advanced Artificial Intelligence", CRC Press, 2024 **<1%**