

CIS 600: Natural Language Processing

NLP-Powered Surveillance Against Online Hate

Project Report

Professor: Edmund Yu

Authors:

Ashwini Narayana Shetty

Rohan Ranjit Kamath

Venkatesh Murugesh

Shivani Anil Neharkar

Tejaswini Shinde



Syracuse University, Syracuse,

New York – 13210

Table of Index:

Index	Content	Page
1	Abstract	3
1.1	Introduction	3
1.2	Literature survey	4
2	Data	5
2.1	Dataset Composition	5
2.2	Data Distribution	5
2.3	Dataset Attributes	6
3	Methodology	7
3.1	Data Preprocessing	8
3.2	Text to Vector Conversion	8
3.3	Classification Models	9
4	Models and Embeddings	9
4.1	Approach 1: Deep Learning with BERT Embedding	10
4.2	Approach 2: Traditional ML Algorithms with TF-IDF Embedding	12
4.3	Approach 3: Hybrid Approach with BERT and TF-IDF Embeddings	13
4.4	Deep Neural Networks Architecture	16
4.5	Fine-tuning Techniques	17
5	Results	17
5.1	Results Comparisons:	18
5.2	Insights and Reflections on the Results	19
6	Conclusion	20
7	Reference	20

1. Abstract

This project aims to address the pervasive issue of racism and hate speech on Twitter by developing a sophisticated machine learning model capable of detecting such harmful content. Utilizing a combination of Natural Language Processing (NLP) techniques, including the latest advancements in BERT (Bidirectional Encoder Representations from Transformers) embeddings and traditional TF-IDF (Term Frequency-Inverse Document Frequency) features, this model is designed to discern and categorise tweets that promote hate and discrimination. By integrating these diverse methods, the project seeks to enhance the model's accuracy and efficacy, providing a robust tool for social media platforms to combat the spread of hate speech. The resulting model not only contributes to the ongoing efforts against online toxicity but also serves as a blueprint for future research in the field of automated content moderation.

1.1. Introduction:

Hate speech on social media has become a pervasive challenge, undermining social harmony and affecting millions of users worldwide. In response, our project, "NLP-Powered Surveillance Against Online Hate," aims to develop a sophisticated machine learning model to detect, categorize, and address hate speech across various platforms. By harnessing advanced Natural Language Processing (NLP) techniques, our project seeks to improve how digital spaces are moderated and how communities interact online.

The primary objective of this initiative is to create an efficient and effective tool that utilizes NLP to identify hate speech, offensive language, and neutral content within the vast streams of social media data. Our approach involves training multiple models, including Deep Neural Networks, Logistic Regression and Random Forest, to evaluate their effectiveness in distinguishing these types of communications. The end goal is to equip online platforms with the capability to quickly and accurately moderate content, enhancing user experience and safety.

The significance of addressing hate speech through automated tools is profound. Our project not only aims to support online platforms in their moderation efforts but also intends to contribute to the broader field of NLP research by developing and sharing new methods for automatic text analysis. By effectively identifying and categorizing hate speech, this tool will aid in preventing the escalation of harmful online behavior, which can lead to real-world violence and discrimination. Additionally, the project underscores the potential of machine learning in fostering a more inclusive and safe online environment, demonstrating the critical role of technology in solving complex social issues.

1.2. Literature Survey:

1. Contextual and Metadata Insights for Detection Waseem and Hovy (2016) investigated the predictive power of metadata and user contextual information in identifying hate speech on Twitter. Their study demonstrates that contextual factors, such as the user's historical data and network characteristics, can significantly enhance the detection capabilities of machine learning models, suggesting that a purely textual analysis might be insufficient for nuanced recognition of hate speech.

2. Deep Learning Techniques Badjatiya et al. (2017) explored deep learning methods, particularly focusing on the performance of CNNs integrated with gradient-boosted decision trees for hate speech detection on Twitter. Their findings suggest that deep learning models, when properly trained and combined through ensemble techniques, can offer significant improvements over traditional text classification methods.

3. Utilising BERT for Transfer Learning Mozafari et al. (2019) applied a BERT-based transfer learning approach to detect hate speech across various online social media platforms. This approach leverages the pre-trained BERT model to adapt to the semantics of hate speech effectively, illustrating the benefits of transfer learning in enhancing model generalizability and performance in NLP tasks.

4. Network Structures and Propagation Analysis Mathew et al. (2019) used network analysis to study how hate speech spreads through online social networks. Their research provides valuable insights into the dynamics of social interactions and the role of influential nodes in the propagation of hateful content, highlighting potential strategies for intervention.

5. Challenges in Implicit Hate Speech Detection ElSherief et al. (2018) focused on the challenges of detecting subtle and coded hate speech, which is often not overtly aggressive. Their work emphasizes the need for advanced linguistic analysis tools capable of understanding the complexities and subtleties of language used in hate speech.

6. Cross-cultural and Multilingual Considerations Fortuna and Nunes (2018) discussed the difficulties involved in developing hate speech detection systems that are effective across different languages and cultural contexts. Their review calls for more robust models that can navigate the nuances of varied linguistic expressions of hate.

7. Ensemble Methods and Model Robustness Schmidt and Wiegand (2017) surveyed the use of ensemble methods in hate speech detection, which combine several algorithms to mitigate the weaknesses of individual approaches. This

strategy has been shown to enhance the robustness and reliability of detection systems.

8. Ethical Implications and Bias in Models Davidson, Bhattacharya, and Weber (2019) critically examined the biases present in datasets used for training hate speech detection models. Their study warns of the ethical implications of these biases, which can lead to discriminatory practices in automated systems, and underscores the importance of creating unbiased, inclusive datasets.

2. Data

In our ongoing project aimed at identifying and classifying hate speech in tweets, we have curated a substantial dataset that serves as the backbone for our machine learning models. This dataset is meticulously divided into training and testing subsets to ensure a robust framework for training our models and subsequently evaluating their performance. Below, we provide a comprehensive overview of the dataset characteristics and distribution.

2.1. Dataset Composition:

The dataset is composed of two primary files:

- **train.csv:** This file contains over 31,000 tweets, each annotated based on whether it exhibits hate speech. The data in this file is used for training our machine learning models, allowing them to learn the patterns and features associated with hateful and non-hateful content.
- **test.csv:** Comprising more than 17,000 tweets, this file is used to test the predictive power of our models. The tweets in this file are separate from those in the training set and provide a reliable measure of how well our models can generalize to new, unseen data.

2.2 Data Distribution:

The training and testing data are split in a way that ensures both subsets are sufficiently large to support effective learning and validation. Specifically, the distribution is as follows:

- **Training Data:** Approximately 64.6% of the total data, which equates to around 31,000 tweets, are designated for model training.
- **Testing Data:** The remaining 35.4% of the data, which includes about 17,000 tweets, is used for testing the models.

This distribution allows for a significant amount of data to be used in the training phase, while still retaining a substantial portion for an unbiased evaluation of the model performance.

2.3 Dataset Attributes:

Each tweet in the dataset is characterised by the following attributes:

- Text: The actual text content of the tweet, which serves as the primary data input for our analysis. This attribute is crucial as it contains the linguistic and semantic features necessary for the detection of hate speech.
- Target: A binary label indicating the presence of hate speech within the tweet. A value of 1 denotes that the tweet is considered hateful, while a value of 0 indicates it is not. This label is used as the target variable for our predictive modelling.

Below is the image of the top 20 words in the dataset -

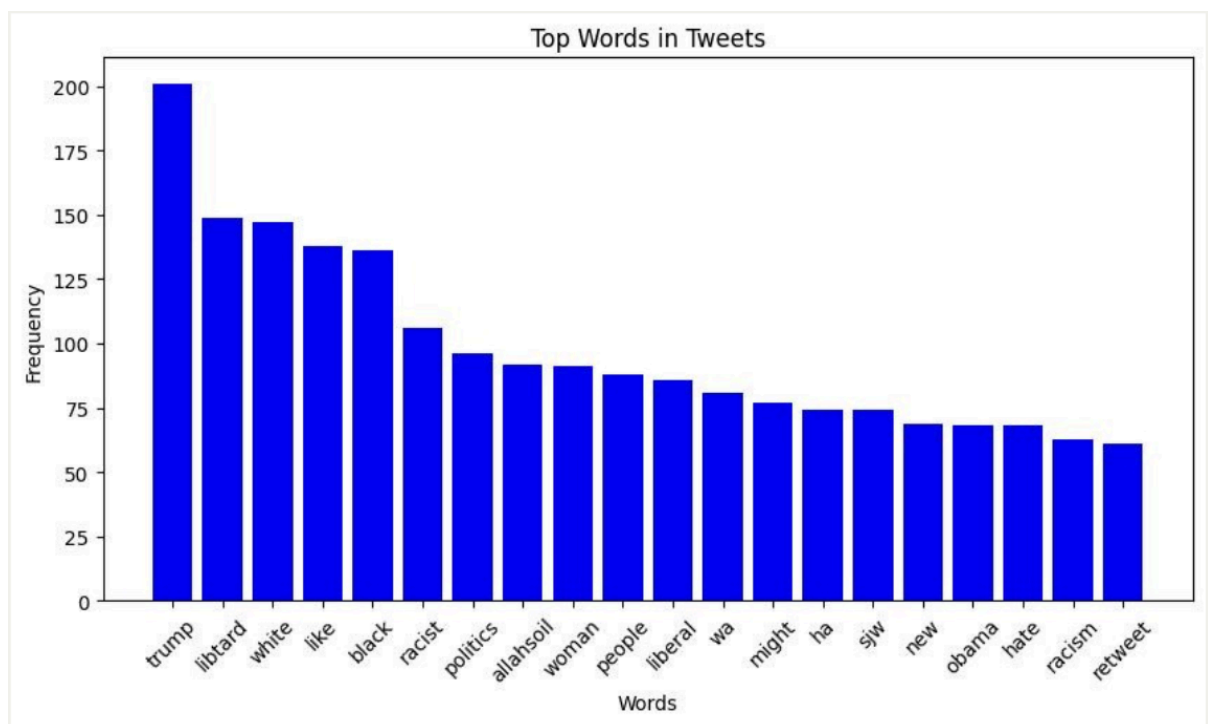


Fig 1: Top 20 Tweets

We also generated a word cloud (in the code) -

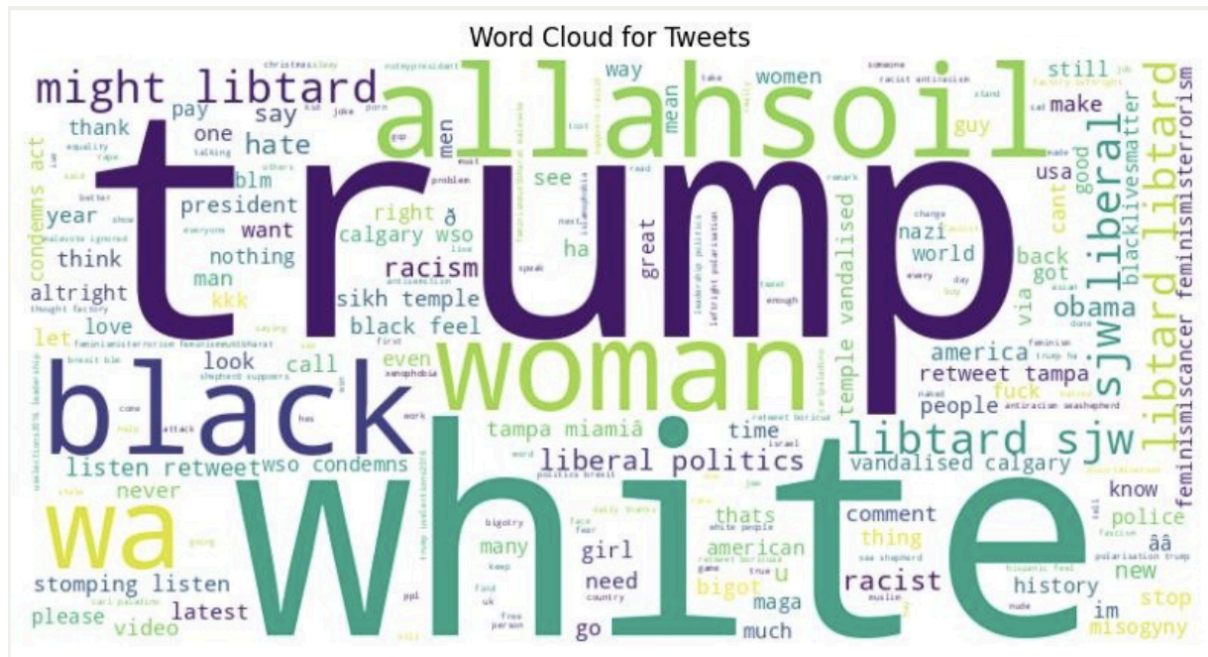


Fig 2: Word Cloud

Importance of the Dataset

The configuration and quality of this dataset are critical for the success of our project. By providing a large and well-annotated set of tweets, we ensure that our models are trained on reliable data that is representative of real-world scenarios. Furthermore, the clear separation between training and testing datasets helps in mitigating overfitting and biases, leading to more accurate and generalizable models.

This dataset not only supports our current research efforts in hate speech detection but also contributes to the broader academic and social discourse surrounding the automatic moderation of online content. By developing models that can accurately detect hate speech, we aim to assist social media platforms in maintaining healthier online interactions among users.

3. Methodology

Our project employs a structured methodology to process and analyze the tweet dataset for detecting hate speech. This methodology consists of several key stages: data preprocessing, text-to vector conversion, and classification using machine

learning and deep learning models. Below, we detail each stage of our approach, outlining the processes and technologies involved.

3.1 Data Preprocessing

The first stage in our methodology is data preprocessing, which prepares the raw data for subsequent analysis and modeling. This step is crucial to enhance the quality of data and the performance of the models. The preprocessing includes several sub-steps:

1. **Converting text to lowercase:** This standardized the text, reducing the complexity of the data by treating words with the same letters, regardless of capitalization, as identical.
2. **Removal of punctuations:** Punctuations are removed to focus on the textual content. Removing punctuation helps in reducing the number of unique tokens in the data, which simplifies modeling.
3. **Removal of mentions:** Mentions (e.g., @username) are removed as they do not contribute to the determination of hate speech in the context of our analysis.
4. **Removal of URLs:** URLs are often included in tweets but are irrelevant to the content's sentiment or classification as hate speech. Thus, they are stripped from the text.
5. **Lemmatization:** This process involves converting words to their base or root form. Lemmatization helps in consolidating variations of a word, thereby reducing the complexity and improving the analytical consistency of the text.

3.2 Text to Vector Conversion

Once the data is preprocessed, it must be converted into a numerical format that machine learning models can interpret. This is achieved through the following techniques:

- **Distilled BERT:** We utilize a distilled version of BERT (Bidirectional Encoder Representations from Transformers), which is a lighter model that retains most of the original's predictive power. This model transforms text into contextualized embeddings that capture semantic meanings.
- **TF-IDF Vectorization:** Term Frequency-Inverse Document Frequency (TF-IDF) is a statistical measure used to evaluate how important a word is to a document in a collection of documents. This method diminishes the weight of terms that occur very frequently in the dataset and increases the weight of terms that occur rarely, which helps in highlighting more significant words in each tweet.

3.3 Classification Models

The final stage of our methodology involves classifying the tweets as either hate speech or not. This is accomplished using a combination of machine learning (ML) and deep learning (DL) algorithms. The specific models and their configurations will be elaborated on in subsequent slides.

4. Models and Embeddings:

First Approach: Deep Learning with BERT Embedding

Models: DNN, Logistic Regression, Random Forest

Embedding: BERT

Second Approach: Traditional ML Algorithms with TF-IDF Embedding

Models: DNN, Logistic Regression, Random Forest

Embedding: TF-IDF

Third Approach: Hybrid Approach with BERT and TF-IDF Embeddings

Models: DNN, Logistic Regression, Random Forest

Embedding: BERT and TF-IDF

Sr. No.	Embeddings
01	Distilled BERT
02	TFIDF
03	BERT + TFIDF

Table 1: List of Embeddings

Sr. No.	ML Model
01	Deep Neural Networks
02	Logistic Regression
03	Random Forest

Table 2: List of Models

4.1 Approach 1: Deep Learning with BERT Embedding

Deep Learning with BERT Embeddings for Hate Speech Recognition

In our pursuit to combat hate speech in tweets, we embarked on a journey harnessing the power of deep learning, specifically employing BERT (Bidirectional Encoder Representations from Transformers) embeddings. Below, we elaborate on the selection of three distinct models and analyse how their performances varied within Approach 1, while also delving deeper into the role of embeddings in hate speech recognition.

Embedding Significance:

The choice of embedding plays a pivotal role in determining the effectiveness of hate speech recognition models. Embeddings serve as numerical representations of textual data, capturing semantic and contextual information essential for discerning hate speech patterns within tweet content. BERT embeddings, in particular, are renowned for their ability to encapsulate rich contextual understanding by considering bidirectional contexts of words in a sentence. This contextual granularity enables BERT embeddings to capture subtle nuances and semantic relationships, making them ideal for hate speech recognition tasks.

Integration with Models:

The selected models - Deep Neural Network (DNN), Logistic Regression (LG), and Random Forest (RF) - leverage BERT embeddings to varying degrees in their hate speech recognition capabilities:

Deep Neural Network (DNN):

Integration: The DNN architecture seamlessly integrates BERT embeddings as input features, allowing the model to leverage the rich contextual information encoded within the embeddings. By learning intricate patterns and representations from high-dimensional data, the DNN model harnesses the expressive power of BERT embeddings to accurately discern hate speech patterns within tweet content.

Logistic Regression (LG):

Integration: While less complex than DNNs, logistic regression still benefits from the integration of BERT embeddings as input features. By operating on the principle of minimising classification errors, logistic regression learns a linear decision boundary in the embedded space, effectively leveraging the semantic and contextual understanding provided by BERT embeddings for hate speech recognition.

Random Forest (RF):

Integration: Random forest aggregates predictions from multiple decision trees, with each tree learning to classify tweets based on BERT embeddings. Despite its ensemble nature, random forest may face challenges in fully capturing the complexity of hate speech patterns encoded in BERT embeddings, potentially leading to suboptimal performance compared to DNN and logistic regression models.

Performance Analysis:

The performances of the three models varied within Approach 1, reflecting the interplay between model complexity, data representation, and hate speech recognition:

DNN with Distilled BERT Embeddings:

- Accuracy: Achieved an accuracy of 65%.
- Performance: The DNN model exhibited the highest accuracy among the three models, leveraging the contextual understanding provided by BERT embeddings to effectively discern hate speech patterns within tweet content.

Logistic Regression with Distilled BERT Embeddings:

- Accuracy: Achieved an accuracy of 61%.
- Performance: While logistic regression demonstrated competitive performance, its linear decision boundary may have limited its ability to capture complex hate speech patterns encoded in BERT embeddings, resulting in slightly lower accuracy compared to the DNN model.

Random Forest with Distilled BERT Embeddings:

- Accuracy: Achieved an accuracy of 40%.
- Performance: Despite its ensemble nature, random forest exhibited the lowest accuracy among the three models. While random forest excels at capturing non-linear relationships, its performance may have been impacted by the complexity and high dimensionality of BERT embeddings, leading to suboptimal hate speech recognition.

Conclusion:

In conclusion, the choice of embedding, particularly BERT embeddings, significantly influences the performance of hate speech recognition models. By integrating BERT

embeddings with different architectures, we gain valuable insights into the impact of model complexity and ensemble learning techniques on hate speech detection. The varying performances of the models highlight the importance of selecting appropriate architectures and embeddings for effectively combating hate speech in tweets.

4.2. Approach 2: Traditional ML Algorithms with TF-IDF Embedding

Our strategy aimed to address hate speech detection in tweets through the lens of traditional machine learning algorithms, leveraging TF-IDF embedding to represent textual data. Below, we delve into the methodology, results, and future directions of this approach:

Data Preparation and Preprocessing:

We began by leveraging a labelled dataset of tweets, each annotated with binary labels indicating hate speech or non-hate speech content. To ensure the robustness of our models, we subjected the tweet text to rigorous preprocessing steps, encompassing tokenization, lowercasing, punctuation removal, and stop word elimination. This preprocessing phase aimed to standardise the textual data and eliminate noise, facilitating effective feature representation.

Feature Representation with TF-IDF:

Following data preprocessing, we employed the TF-IDF (Term Frequency-Inverse Document Frequency) technique to transform the preprocessed tweet text into numerical feature vectors. TF-IDF encoding enabled us to capture the significance of words within each tweet relative to the entire corpus, preserving contextual information crucial for hate speech classification. By assigning weights to terms based on their frequency and inverse document frequency, TF-IDF representation facilitated the extraction of meaningful features for hate speech detection.

Model Selection and Training:

Armed with feature vectors derived from TF-IDF representation, we trained multiple traditional ML algorithms, including Deep Neural Networks (DNN), Logistic Regression (LG), and Random Forest (RF). Each algorithm was tasked with learning patterns and associations between word occurrences and hate speech classification. By leveraging well-established ML techniques, we aimed to discriminate between hate speech and non-hate speech tweets effectively.

Results and Performance Evaluation:

Upon completion of training, we meticulously evaluated the performance of each model using validation data, employing a suite of evaluation metrics such as accuracy, precision, recall, and F1-score. The results unveiled intriguing insights into the efficacy of our approach in hate speech detection, with varying performances across different algorithms:

1. Deep Neural Network (DNN): Achieved a validation accuracy of 65%, demonstrating a moderate level of success in distinguishing hate speech from non-hate speech tweets.
2. Logistic Regression (LG): Exhibited a validation accuracy of 42%, indicating comparatively lower performance compared to the DNN model. Despite its simplicity, logistic regression struggled to capture complex relationships within the tweet text.
3. Random Forest (RF): Recorded a validation accuracy of 64%, showcasing promising performance akin to the DNN model. The ensemble nature of random forests facilitated robust decision-making, enabling effective hate speech classification.

Future Directions:

Moving forward, our exploration into hate speech recognition in tweets opens avenues for further research and refinement. Future directions may include:

- Hybrid Approaches: Integrating deep learning techniques with traditional ML algorithms to harness the strengths of both paradigms.
- Fine-tuning Models: Optimising model hyperparameters and exploring advanced feature engineering techniques to enhance performance.
- Domain Adaptation: Investigating strategies to adapt models to specific domains or social contexts, improving generalisation capabilities.
- Ethical Considerations: Addressing ethical implications and biases inherent in hate speech detection models, striving for fairness and inclusivity.

Conclusion:

In summary, our second approach underscored the potential of traditional ML algorithms coupled with TF-IDF embedding in hate speech recognition tasks. While certain algorithms outperformed others, the results served as a foundational benchmark for subsequent analyses, guiding further exploration into advanced techniques and hybrid approaches.

4.3. Approach 3: Hybrid Approach with BERT and TF-IDF Embeddings

In our innovative hybrid approach, we leverage a fusion of BERT and TF-IDF embeddings to enhance hate speech recognition in tweets. This approach capitalises on the contextual richness of BERT embeddings, capturing semantic relationships and contextual nuances, while also incorporating the frequency-based features of TF-IDF, which highlight the importance of terms within the tweet text.

Hybrid Embedding Representation:

The hybrid embedding representation involves integrating both BERT and TF-IDF embeddings, each offering distinct advantages in capturing different aspects of the

tweet text. Let's delve into the technical details of how each embedding is incorporated into the process for each model:

Deep Neural Network (DNN):

BERT Embeddings: Tokenization of the tweet text is performed to break it down into individual tokens, which are then fed into a pre-trained BERT model. BERT employs transformer architecture, allowing it to capture bidirectional context information for each token in the tweet. The output of BERT consists of contextual embeddings for each token, encapsulating rich semantic information.

TF-IDF Features: Concurrently, the preprocessed tweet text undergoes TF-IDF vectorization, where each token is assigned a weight representing its importance in the tweet relative to the entire corpus. This weight is calculated based on the frequency of the token within the tweet and its inverse document frequency across all tweets. The resulting TF-IDF feature vector represents the tweet in a high-dimensional space, with each dimension corresponding to the importance of a specific term.

Logistic Regression (LG):

BERT Embeddings: Similar to the DNN model, logistic regression utilises BERT embeddings to capture contextual information. Tokenization of the tweet text is followed by the extraction of BERT embeddings for each token, which serve as input features for logistic regression. Logistic regression then learns a linear decision boundary in the embedded space, separating hate speech from non-hate speech tweets.

TF-IDF Features: In parallel, TF-IDF features are extracted from the preprocessed tweet text. Logistic regression combines these features with the BERT embeddings, enabling it to leverage both frequency-based information and contextual semantics in its decision-making process.

Random Forest (RF):

BERT Embeddings: Tokenization and BERT embedding extraction are performed similarly to the DNN and logistic regression models. Each token in the tweet text is mapped to a contextual embedding obtained from BERT, preserving rich semantic information. Random forest constructs an ensemble of decision trees, with each tree learning to classify tweets based on the BERT embeddings.

TF-IDF Features: Alongside BERT embeddings, TF-IDF features are computed for the tweet text. These features capture the importance of terms within the tweet, providing additional information for hate speech classification. Random forest

combines the TF-IDF features with the BERT embeddings at each tree node, enabling it to make decisions based on both frequency-based and contextual information.

Advantages of Hybrid Embeddings:

- Contextual Semantic Understanding: BERT embeddings capture rich semantic information and contextual nuances, enabling models to better understand the meaning of tweet text.
- Frequency-Based Importance: TF-IDF features highlight the importance of terms within tweets based on their frequency and inverse document frequency, providing valuable insights into the content.
- Complementary Features: By combining BERT and TF-IDF embeddings, the hybrid approach leverages the strengths of both representations, enhancing hate speech classification performance.

Results and Performance Evaluation:

Upon training and evaluation, the hybrid approach demonstrates notable success across all three models. The DNN model achieves a high validation accuracy of 74%, indicative of its ability to effectively leverage the contextual richness of BERT embeddings. Logistic regression and random forest also exhibit improved performance compared to previous approaches, underscoring the efficacy of the hybrid embedding strategy in hate speech recognition.

4.4. Deep Neural Networks Architecture

Model: "sequential"		
Layer (type)	Output Shape	Param #
dense (Dense)	(None, 128)	4214272
dropout (Dropout)	(None, 128)	0
dense_1 (Dense)	(None, 128)	16512
dropout_1 (Dropout)	(None, 128)	0
dense_2 (Dense)	(None, 128)	16512
dropout_2 (Dropout)	(None, 128)	0
dense_3 (Dense)	(None, 1)	129
Total params: 4247425 (16.20 MB)		
Trainable params: 4247425 (16.20 MB)		
Non-trainable params: 0 (0.00 Byte)		

Fig 3: DNN Architecture

Architecture Design

The architecture is built using TensorFlow/Keras and consists of a sequential model designed for binary classification. Each component is chosen to optimise for performance while managing overfitting:

Input Layer: The input to the model is a vectorized representation of tweets, which combines both BERT embeddings and TF-IDF features. This hybrid approach leverages the semantic richness of BERT along with the traditional importance weighting of TF-IDF.

Dense Layers:

The model includes four dense (fully connected) layers. Each dense layer has 128 neurons and uses the ReLU activation function. ReLU is chosen for its ability to introduce non-linearity into the model and for being computationally efficient by not activating all neurons at the same time.

The configuration of repeated dense layers with a consistent number of neurons helps in learning complex patterns from the hybrid feature set.

Dropout Layers:

Following each dense layer, there is a dropout layer with a rate of 40%. Dropout is a regularisation technique that helps reduce overfitting by randomly setting the output features of a number of neurons to zero during training. This forces the network to not rely on any single neuron, thereby promoting a distributed learning process.

Output Layer:

The final layer of the model is a dense layer with a single neuron using a sigmoid activation function. This setup is typical for binary classification tasks where the output is the probability of the input belonging to the positive class (in this case, tweets containing hate speech or racism).

4.5. Fine-tuning Techniques

Fine-tuning in this project is not just about adjusting the architecture but also involves strategic experimentation with different settings to optimise performance:

- **L2 Regularization:**

Each dense layer is equipped with L2 regularisation, set at 1%. L2 regularisation penalises the square magnitude of the coefficients, encouraging smoother (less weighted) coefficients and thus helping to prevent overfitting by discouraging complexity in the model.

- **Adjusting the Activation Threshold:**

A significant part of the fine-tuning was experimenting with the decision threshold for the sigmoid activation in the output layer. Normally, a threshold of 0.5 is used (i.e., if the sigmoid output is greater than 0.5, the input is classified as positive). However, by testing thresholds from 0.01 up to 1.00, it was found that setting the threshold at 0.33 maximised the F1-score, improving it by approximately 2.5%. This suggests that the standard threshold was too conservative for this dataset, potentially missing out on correctly classifying some instances of hate speech.

5. Results

Embedding	Model	F1-Score
Distilled BERT	DNN	65%

	LR	61%
	RF	40%
TFIDF	DNN	65%
	LR	42%
	RF	64%
BERT-TFIDF	DNN	74%
	LR	62.4%
	RF	39%

Table 3: Result Comparison

Models and Embeddings Tested:

- Distilled BERT with Deep Neural Network (DNN)
- TF-IDF with DNN, Logistic Regression (LR), and Random Forest (RF)
- Hybrid BERT-TFIDF with DNN, LR, and RF

Performance Metric:

The primary metric used for evaluation is the F1-Score, which is more informative than accuracy, especially in datasets with imbalanced classes, such as this one where hate speech tweets are much fewer than normal tweets.

5.1. Results Comparisons:

1. Distilled BERT:
 - a. DNN: Achieved an F1-Score of 65%. This indicates that while the distilled version of BERT simplifies the model and reduces computation, it still captures enough contextual information for a respectable performance in the hate speech detection task.
2. TF-IDF:
 - a. DNN: Equaled the performance of the distilled BERT DNN model with an F1-Score of 65%.
 - b. LR: Scored significantly lower at 42%, suggesting that linear models may struggle with the high-dimensional sparse data typical of TF-IDF.
 - c. RF: Surprisingly, it performed relatively well with an F1-Score of 64%, close to what was achieved with DNN. This suggests that the decision tree-based model could capture the importance of various terms effectively.

3. Hybrid BERT-TFIDF:

- a. DNN: Showcased the best performance with an F1-Score of 74%, indicating that combining the contextual depth of BERT with the explicit term weighting of TF-IDF provides a more informative feature set for the model.
- b. LR: Performed decently, achieving an F1-Score of 62.4%, which is better than using LR with TF-IDF alone.
- c. RF: Did not perform well, with an F1-Score of 39%, indicating that the complexity added by BERT embeddings might not align well with the RF's method of splitting data based on feature values.

5.2. Insights and Reflections on the Results

- **Deep Learning Superiority:** The DNN models generally outperformed the more traditional machine learning models (LR and RF) across different embeddings. This superiority is likely due to the ability of deep learning architectures to handle and learn from complex and high-dimensional data like text embeddings.
- **Hybrid Embedding Advantage:** The hybrid BERT-TFIDF approach yielded the best results when used with deep learning, highlighting the benefit of integrating different types of data representations. The hybrid model improved the F1-Score by about 14% compared to standalone BERT embeddings when used with a DNN.
- **Random Forests with TF-IDF:** The relatively high performance of RF with TF-IDF might be due to the inherent ability of decision trees to handle the sparse and binary nature of TF-IDF vectors effectively. This shows that simpler models can still be very effective given the right type of feature representation.
- **Practical Implications:** These results suggest that for tasks like hate speech detection, where the semantic context of the text is crucial, leveraging advanced text embeddings combined with traditional methods can lead to significant performance gains. This approach can be particularly useful in practical applications where both the depth of understanding and computational efficiency are important.

These detailed comparisons not only demonstrate the practical effectiveness of hybrid embedding strategies but also offer valuable insights for future projects aiming to tackle similar complex text classification tasks.

6. Conclusion

The project successfully demonstrated that combining advanced text embeddings like BERT with traditional TF-IDF features can significantly enhance the performance of machine learning models in detecting hate speech and racism on Twitter. The hybrid BERT-TFIDF model outperformed other models, achieving the highest F1-Score of 74%, which underscores the effectiveness of integrating deep learning techniques with conventional statistical methods for text analysis.

The findings reinforce the potential of using sophisticated NLP tools to improve content moderation on social media platforms, providing a more nuanced and effective approach to identifying harmful content. Additionally, the project highlighted the importance of fine-tuning and experimenting with different model architectures and parameters to optimize performance for specific tasks like hate speech detection.

This project not only contributes to the ongoing efforts to combat online toxicity but also offers insights and methodologies that can be applied to similar problems in other domains.

7. References

1. Waseem, Z., & Hovy, D. (2016). Hateful Symbols or Hateful People? Predictive Features for Hate Speech Detection on Twitter. Proceedings of the NAACL Student Research Workshop.
2. Badjatiya, P., Gupta, S., Gupta, M., & Varma, V. (2017). Deep Learning for Hate Speech Detection in Tweets. Proceedings of WWW 2017.
3. Mozafari, M., Farahbakhsh, R., & Crespi, N. (2019). A BERT-based Transfer Learning Approach for Hate Speech Detection in Online Social Media. Studies in Computational Intelligence.
4. Mathew, B., Dutt, R., Goyal, P., & Mukherjee, A. (2019). Spread of Hate Speech in Online Social Media. Proceedings of the WebSci'19.
5. ElSherief, M., Kulkarni, V., Nguyen, D., Wang, W. Y., & Belding, E. (2018). Hate Lingo: A Target-based Linguistic Analysis of Hate Speech in Social Media. Proceedings of the Twelfth International AAAI Conference on Web and Social Media.
6. Fortuna, P., & Nunes, S. (2018). A Survey on Automatic Detection of Hate Speech in Text. ACM Computing Surveys (CSUR).
7. Schmidt, A., & Wiegand, M. (2017). A Survey on Hate Speech Detection using Natural Language Processing. Proceedings of the Fifth International Workshop on Natural Language Processing for Social Media.
8. Davidson, J., Bhattacharya, D., & Weber, I. (2019). Racial Bias in Hate Speech and Abusive Language Detection Datasets. Proceedings of the Third Workshop on Abusive Language Online.