# uOttawa

# Project Report
# Group ID : 7

**DTI5125[EG] Data Science Applications [LEC] 20235**

**Submitted By:**

Abdelrahman Ahmed Mansour Badran
Hussien Tarek Ismail Abdelrazik
Mohamed Ahemd Sayed Mohamed
Mohamed Magdy Mahmoud Elasmar

**Instructor:**

Professor: Arya Rahgozar

University of Ottawa, Canada

August 6, 2023

# Contents

# 1 Introduction

## 1.1 Background

Hate speech has become a pervasive concern in modern communication platforms, perpetuating harm and divisiveness. We present a machine learning-based approach for hate speech detection in a text-input chatbot system. Our method aims to promptly identify and mitigate instances of hate speech, fostering a more inclusive online environment.

The significance of our approach lies in real-time hate speech assessment, enabling immediate interventions. Using natural language processing (NLP), we empower the chatbot to classify input text accurately. This proactive identification reduces harmful language spread, enhancing online interactions.

In this report, we detail the design and implementation of our hate speech detection system. We describe data preparation, text feature engineering techniques, and classification algorithms. Additionally, we evaluate system performance, including error analysis and visualizations. We discuss potential adaptation to movie subtitles to redact offensive content.

Throughout this report, we emphasize our contributions, focusing on real-time chatbot assessment to combat digital hate speech.

## 1.2 Problem Formulation

In the digital age, hate speech escalation within online platforms necessitates early detection strategies. Our research addresses this issue:

**Problem Statement:** Develop an efficient chatbot-integrated hate speech detection system, promptly identifying hate speech and offensive language.

The urgency arises from potential harm inflicted, perpetuating prejudice and disrupting online interactions. Our solution employs NLP techniques and machine learning for real-time detection. Our approach combines data preparation, feature engineering, and diverse classification algorithms for accurate, user-friendly chatbot assessment.

Future prospects include adapting the framework to address offensive content within movie subtitles. Recognizing sensitive or hateful language using timestamps can promote inclusive entertainment media.

Subsequent sections detail our methodology, data-driven decisions, classification outcomes, and system performance evaluation. We provide error analysis insights and visualizations, aspiring to curtail digital hate speech's impact and enable broader applications.

## 1.3 Potential Applications

Beyond its primary objective of hate speech detection within chatbots, our developed system holds the potential for versatile applications. Some of the potential applications include:

- **Chatbot Integration:** Our system can seamlessly integrate with various chatbot platforms, enabling real-time hate speech assessment and facilitating healthier online conversations.

- **Content Moderation Automation:** Online platforms often struggle with moderating user-generated content. Our system can automate content moderation by flagging and filtering out hate speech, leading to a safer online environment.

- **Social Media Platforms:** Social media companies can employ our system to automatically detect and prevent hate speech, thereby fostering more constructive discussions and reducing the spread of harmful content.

- **Customer Support:** Customer service chatbots can benefit from hate speech detection to ensure respectful and positive interactions with users.

- **Education and Awareness:** Our system can be used to monitor online educational platforms, ensuring a safe learning environment free from hate speech.

# 2 Data Preparation

The data used in this study is a combination of Fox News data [1] and Twitter data [2], which were integrated to create a comprehensive dataset. The preprocessing steps involved data cleaning and text normalization to ensure the quality and consistency of the textual data.

## 2.1 Text Preprocessing

To prepare the textual data for analysis, several preprocessing steps were performed. These steps included:

- **Lemmatization**:The text data was lemmatized using the spaCy library. Lemmatization reduces words to their base or root form, which helps in standardizing the text.

- **Punctuation and Number Removal**: Punctuation marks and numbers were replaced with spaces using regular expressions. This step eliminates potential noise from the data.

- **Whitespace Removal**: Extra spaces were removed from the text to ensure uniform spacing.

- **Single-Character Word Removal**: Single-character words, which often carry minimal meaning, were filtered out.

- **Stopword Removal**: Common English stopwords, which do not contribute significantly to the content's meaning, were removed from the text.

The above preprocessing steps were applied to both the Fox News and Twitter datasets to ensure consistency in data quality.

## 2.2 Data Cleaning for Classification

The data cleaning process extended to the target labels used for classification. The majority answer in the 'movies_data_cleaned' dataset and the 'label' in both the 'twitter_data_cleaned' and 'fox_news_data_cleaned' datasets were transformed as follows:

- For binary classification, the labels were binary-encoded. Specifically, values labeled as 2 were transformed to 1, while all other values were transformed to 0. This conversion facilitates the classification task.

In summary, the data preparation process involved thorough cleaning and normalization of both the textual content and the target labels, setting the stage for subsequent stages of clustering and classification.

## 2.3 Code Snippets

The following code snippets illustrate the implementation of the text preprocessing steps and label transformations:

```python
def lemmatize_text(text):
    doc = nlp(text)
    lemmatized_words = [token.lemma_ for token in doc]
    return ' '.join(lemmatized_words)

def remove_punctuation_numbers_stopwords(text):
    cleaned_text = re.sub(r'[0-9]+', ' ', text)
    cleaned_text = re.sub(r'[{}]'.format(re.escape(string.punctuation)), ' ', cleaned_text)
    cleaned_text = re.sub(r'\s+', ' ', cleaned_text).strip()
    cleaned_text = ' '.join(word for word in cleaned_text.split() if len(word) > 1)
    cleaned_text = ' '.join(word for word in cleaned_text.split() if word not in nltk.corpus.
        stopwords.words('english'))
    return cleaned_text
```

Listing 1: Text Preprocessing Functions

```python
# Applying Text Preprocessing to Datasets
movies_data_cleaned = movies_data.copy()
movies_data_cleaned['text'] = movies_data_cleaned['text'].apply(lemmatize_text)
movies_data_cleaned['text'] = movies_data_cleaned['text'].apply(
    remove_punctuation_numbers_stopwords)

fox_news_data_cleaned = fox_news_data.copy()
```

```
7  fox_news_data_cleaned['comment'] = fox_news_data_cleaned['comment'].apply(lemmatize_text)
8  fox_news_data_cleaned['comment'] = fox_news_data_cleaned['comment'].apply(
       remove_punctuation_numbers_stopwords)
9
10 twitter_data_cleaned = twitter_data.copy()
11 twitter_data_cleaned['tweet'] = twitter_data_cleaned['tweet'].apply(lemmatize_text)
12 twitter_data_cleaned['tweet'] = twitter_data_cleaned['tweet'].apply(
       remove_punctuation_numbers_stopwords)
```

Listing 2: Applying Text Preprocessing to Datasets

```
1  # Transforming Labels
2  movies_data_cleaned['majority_answer'] = movies_data_cleaned['majority_answer'].apply(lambda x: 1
       if x == 2 else 0)
3  twitter_data_cleaned['label'] = twitter_data_cleaned['label'].apply(lambda x: 1 if x == 2 else 0)
4  fox_news_data_cleaned['label'] = fox_news_data_cleaned['label'].apply(lambda x: 1 if x == 2 else
       0)
```

Listing 3: Transforming Labels

The provided code snippets demonstrate the implementation of text preprocessing functions and label transformations for the respective datasets. These steps collectively contribute to the preparation of clean and structured data for subsequent analysis.

## 2.4 Preprocessed Datasets

The resulting preprocessed datasets are as follows:

- 'fox_news_data_cleaned': Cleaned and preprocessed Fox News data.
- 'twitter_data_cleaned': Cleaned and preprocessed Twitter data.
- 'movies_data_cleaned': Processed movie-related data, specifically reserved for testing purposes. [3]

These preprocessed datasets will be utilized for subsequent clustering and classification tasks, ensuring reliable and consistent results.

# 3 Text Feature Engineering

Text feature engineering plays a pivotal role in enhancing the performance of hate speech detection models. In this section, we outline the techniques employed to transform the preprocessed textual data into meaningful numerical representations that can be effectively utilized by machine learning algorithms.

## 3.1 Vectorization and TF-IDF

To represent textual data in a numerical format, we employed two fundamental techniques: vectorization and Term Frequency-Inverse Document Frequency (TF-IDF) encoding.

**Vectorization:** Textual content was transformed into numerical vectors using Count Vectorization. This technique creates a matrix where each row corresponds to a document, and each column represents a unique word's count within the document.

**TF-IDF Encoding:** TF-IDF is a more advanced approach that assigns weights to words based on their frequency within a document relative to their frequency across the entire corpus. This technique enhances the significance of words that are informative to a specific document while downplaying common words.

## 3.2 Principal Component Analysis (PCA)

Principal Component Analysis (PCA) was employed to reduce the dimensionality of the feature space while retaining as much variance as possible. This aids in improving model efficiency and mitigating the curse of dimensionality. The top principal components were selected as the new feature representation.

## 3.3 t-Distributed Stochastic Neighbor Embedding (t-SNE)

Additionally, we explored the application of t-Distributed Stochastic Neighbor Embedding (t-SNE) on the combined training data from Twitter and Fox News. t-SNE is a dimensionality reduction technique specifically designed for visualizing high-dimensional data in a lower-dimensional space while preserving the local structure. A placeholder for the t-SNE visualization is presented in Figure 1.



Figure 1: t-SNE Visualization

The combination of vectorization, TF-IDF encoding, PCA, and t-SNE collectively results in a refined and informative feature representation, facilitating the subsequent classification tasks. The following code snippets illustrate the application of these techniques:

```python
# Vectorization and TF-IDF
from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer

vectorizer = CountVectorizer()  # Or TfidfVectorizer()
X_train_vectorized = vectorizer.fit_transform(X_train)
X_test_vectorized = vectorizer.transform(X_test)

# Principal Component Analysis (PCA)
from sklearn.decomposition import PCA

pca = PCA(n_components=10)  # Choose the desired number of components
X_train_pca = pca.fit_transform(X_train_vectorized.toarray())
X_test_pca = pca.transform(X_test_vectorized.toarray())

import numpy as np
```

```
16  import matplotlib.pyplot as plt
17  from sklearn.manifold import TSNE
18  from sklearn.preprocessing import StandardScaler
19
20  # Step 2: Load your data and labels
21  data = x_train_pca
22  labels = np.hstack((y2, y3))
23
24  # Step 3: Standardize your data (optional but recommended)
25  scaler = StandardScaler()
26  data_scaled = scaler.fit_transform(data)
27
28  # Step 4: Apply t-SNE to your data
29  tsne = TSNE(n_components=2, random_state=42)
30  data_tsne = tsne.fit_transform(data_scaled)
31
32  # Step 5: Plot the t-SNE results with labels
33  plt.figure(figsize=(8, 6))
34  # Loop through each class and plot data points with the corresponding label
35  for label in np.unique(labels):
36      plt.scatter(data_tsne[labels == label, 0], data_tsne[labels == label, 1], label=f'Class {label
        }', s=10)
37
38  plt.xlabel('t-SNE Component 1')
39  plt.ylabel('t-SNE Component 2')
40  plt.title('t-SNE Visualization with Labels')
41  plt.legend()
42  plt.show()
```
Listing 4: Text Feature Engineering Code Snippets

The forthcoming sections will delve into the classification algorithms applied to the engineered text features, and the subsequent evaluation of their performance.

# 4 Clustering

In this section, we explore the application of clustering algorithms to group similar instances of hate speech. Clustering aids in identifying patterns and relationships within the data, which can provide valuable insights into the distribution and structure of hate speech instances.

## 4.1 K-Means

K-Means is a widely used clustering algorithm that partitions data into a specified number of clusters. We applied K-Means clustering to the combined training data from Twitter and Fox News and evaluated its performance using validation techniques.

**K-Means for Combined Training Data:** We performed K-Means clustering on the PCA-transformed features of the combined training data from Twitter and Fox News. The algorithm aimed to group instances based on their extracted principal components, potentially revealing underlying patterns in hate speech instances.

```
1  from sklearn.cluster import KMeans
2  from sklearn.metrics import silhouette_score
3
4  # K-Means for Combined Training Data
5  kmeans_combined = KMeans(n_clusters=10, n_init='auto').fit(X_train_pca)
6  silhouette_score_combined = silhouette_score(X_train_pca, y_train, metric='euclidean')
7  print('Silhouette score for combined training data:', silhouette_score_combined)
```
Listing 5: K-Means Clustering for Combined Training Data

The silhouette score provides insight into the clustering quality, with higher scores indicating better-defined clusters. A 0.06 Score does not indicate good performance but gives us the motive to try out different classification algorithms.

## 4.2 Discussion and Insights

The application of K-Means clustering, coupled with visualization techniques, offers valuable insights into the distribution of hate speech instances within the dataset. The silhouette score and t-SNE visualization collectively contribute to a deeper understanding of the clustering outcomes. These insights pave the way for informed decisions in subsequent classification tasks and further refinement of our hate speech detection system.

# 5 Classification

In this section, we discuss the implementation of three classification algorithms for the hate speech detection task: XGBoost, Logistic Regression, and BERT transfer learning.

## 5.1 XGBoost Classifier

XGBoost is a powerful gradient boosting algorithm commonly used for classification tasks. In this subsection, we present the XGBoost classifier implementation and evaluation.

```python
# Import necessary libraries and modules
from xgboost import XGBClassifier
from sklearn.metrics import classification_report
from sklearn.metrics import confusion_matrix
import seaborn as sns
import matplotlib.pyplot as plt

# Create an XGBoost model
model_xgb = XGBClassifier()

# Train the model on the combined training data
model_xgb.fit(x_train_pca, np.hstack((y2, y3)))

# Make predictions on the movies data
y_pred_xgb = model_xgb.predict(x1_pca)

# Calculate classification report
classification_report_movies_xgb = classification_report(y1, y_pred_xgb)
print(classification_report_movies_xgb)

# Plot Confusion Matrix for Movies Data
confusion_matrix_movies_xgb = confusion_matrix(y1, y_pred_xgb)
sns.heatmap(confusion_matrix_movies_xgb, annot=True, fmt='g')
plt.xlabel('Predicted')
plt.ylabel('True')
plt.title('Confusion Matrix for Movies Data (XGBoost)')
plt.show()
```

Listing 6: XGBoost Classifier for Hate Speech Detection

## 5.2 Logistic Regression Classifier

Logistic Regression is a linear classification algorithm commonly used for binary classification tasks. Here, we demonstrate the implementation and evaluation of the Logistic Regression classifier.

```python
# Import necessary libraries and modules
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import classification_report, confusion_matrix
import seaborn as sns
import matplotlib.pyplot as plt

# Create a Logistic Regression model
model_logreg = LogisticRegression()

# Train the model on the combined training data
model_logreg.fit(x_train_pca, np.hstack((y2, y3)))

# Make predictions on the movies data
y_pred_logreg = model_logreg.predict(x1_pca)
```

```
16  # Calculate classification report
17  classification_report_movies_logreg = classification_report(y1, y_pred_logreg)
18  print(classification_report_movies_logreg)
19
20  # Plot Confusion Matrix for Movies Data
21  confusion_matrix_movies_logreg = confusion_matrix(y1, y_pred_logreg)
22  sns.heatmap(confusion_matrix_movies_logreg, annot=True, fmt='g')
23  plt.xlabel('Predicted')
24  plt.ylabel('True')
25  plt.title('Confusion Matrix for Movies Data (Logistic Regression)')
26  plt.show()
```

Listing 7: Logistic Regression Classifier for Hate Speech Detection

## 5.3 BERT Transfer Learning

BERT (Bidirectional Encoder Representations from Transformers) is a pre-trained language model that can be fine-tuned for various NLP tasks, including hate speech detection. Here, we showcase the implementation and evaluation of BERT transfer learning.

```
1   # Import necessary libraries and modules
2   ! pip install transformers
3   from transformers import BertTokenizer, TFBertForSequenceClassification
4   from transformers import InputExample, InputFeatures
5
6   # Load the pre-trained BERT model
7   model_bert = TFBertForSequenceClassification.from_pretrained("bert-base-uncased",
8                                                   trainable=True,
9                                                   num_labels=2)
10  tokenizer = BertTokenizer.from_pretrained("bert-base-uncased")
11
12  # Fine-tune the BERT model using training data
13  # Convert data to BERT-compatible format
14  train_InputExamples, validation_InputExamples = convert_data_to_examples(train, test, DATA_COLUMN,
        LABEL_COLUMN)
15
16  # Convert examples to TensorFlow datasets
17  train_data = convert_examples_to_tf_dataset(list(train_InputExamples), tokenizer)
18  train_data = train_data.batch(32)
19
20  validation_data = convert_examples_to_tf_dataset(list(validation_InputExamples), tokenizer)
21  validation_data = validation_data.batch(32)
22
23  # Compile and fit the BERT model
24  model_bert.compile(optimizer=tf.keras.optimizers.Adam(learning_rate=3e-6, epsilon=1e-08, clipnorm
        =1.0),
25                   loss=tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True),
26                   metrics=[tf.keras.metrics.SparseCategoricalAccuracy('accuracy')])
27
28  hist_bert = model_bert.fit(train_data, epochs=4, validation_data=validation_data, callbacks=[
        cp_callback], use_multiprocessing=True)
29
30  # Make predictions on validation data
31  preds_bert = model_bert.predict(validation_data)
32
33  # Calculate classification report
34  classification_report_movies_bert = classification_report(test['LABEL_COLUMN'], np.argmax(
        preds_bert[0], axis=1), output_dict=True)
35  print(pd.DataFrame(classification_report_movies_bert).transpose())
36
37  # Plot Confusion Matrix for Movies Data (BERT)
38  confusion_matrix_movies_bert = confusion_matrix(test['LABEL_COLUMN'], np.argmax(preds_bert[0],
        axis=1))
39  sns.heatmap(confusion_matrix_movies_bert, annot=True, fmt='g')
40  plt.xlabel('Predicted')
41  plt.ylabel('True')
42  plt.title('Confusion Matrix for Movies Data (BERT)')
43  plt.show()
44
45  # Evaluate the BERT model on movies data
```

```
46 movie_InputExamples = convert_data_to_examples_valid(df_movies, DATA_COLUMN, LABEL_COLUMN)
47 movie_data = convert_examples_to_tf_dataset(list(movie_InputExamples), tokenizer)
48 movie_data = movie_data.batch(32)
49
50 preds_movie_bert = model_bert.predict(movie_data)
51
52 # Calculate classification report for movies data
53 classification_report_movies_bert = classification_report(df_movies['LABEL_COLUMN'], np.argmax(
       preds_movie_bert[0], axis=1), output_dict=True)
54 print(pd.DataFrame(classification_report_movies_bert).transpose())
55
56 # Plot Confusion Matrix for Movies Data (BERT)
57 confusion_matrix_movies_bert = confusion_matrix(df_movies['LABEL_COLUMN'], np.argmax(
       preds_movie_bert[0], axis=1))
58 sns.heatmap(confusion_matrix_movies_bert, annot=True, fmt='g')
59 plt.xlabel('Predicted')
60 plt.ylabel('True')
61 plt.title('Confusion Matrix for Movies Data (BERT)')
62 plt.show()
63
64 # Save the trained BERT model
65 model_bert.save('bert_hate_speech_model.h5')
```

Listing 8: BERT Transfer Learning for Hate Speech Detection

In the subsequent sections, we analyze and discuss the results obtained from each classification algorithm, providing insights into their strengths and limitations in hate speech detection.

# 6 Visualization and Evaluation of ML Results & Error Analysis

In this section, we present the evaluation results for the three classification algorithms: XGBoost, Logistic Regression, and BERT Transfer Learning. We provide a detailed analysis of the precision, recall, and F1-score metrics for each algorithm, along with confusion matrices illustrating the model's performance.

## 6.1 XGBoost Evaluation

The XGBoost algorithm was trained and evaluated on the combined training data from Fox News and Twitter. The evaluation results are as follows:

|              | Precision | Recall | F1-Score | Support |
|--------------|-----------|--------|----------|---------|
| Class 0      | 0.98      | 0.99   | 0.99     | 10394   |
| Class 1      | 0.48      | 0.29   | 0.36     | 294     |
| Accuracy     |           |        | 0.97     | 10688   |
| Macro Avg    | 0.73      | 0.64   | 0.67     | 10688   |
| Weighted Avg | 0.97      | 0.97   | 0.97     | 10688   |

Table 1: XGBoost Classification Results



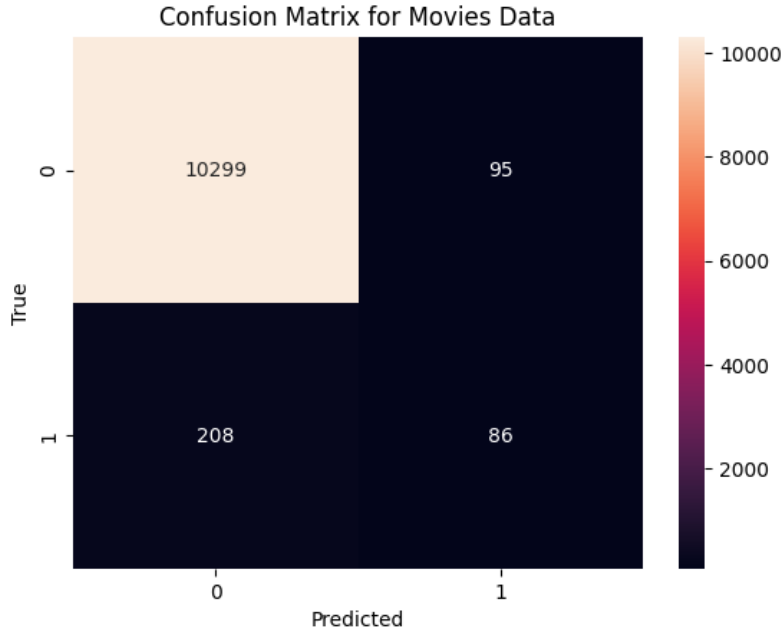Figure 2: Confusion Matrix for XGBoost Model

## 6.2 Logistic Regression Evaluation

The Logistic Regression algorithm was trained and evaluated on the combined training data. The evaluation results are as follows:

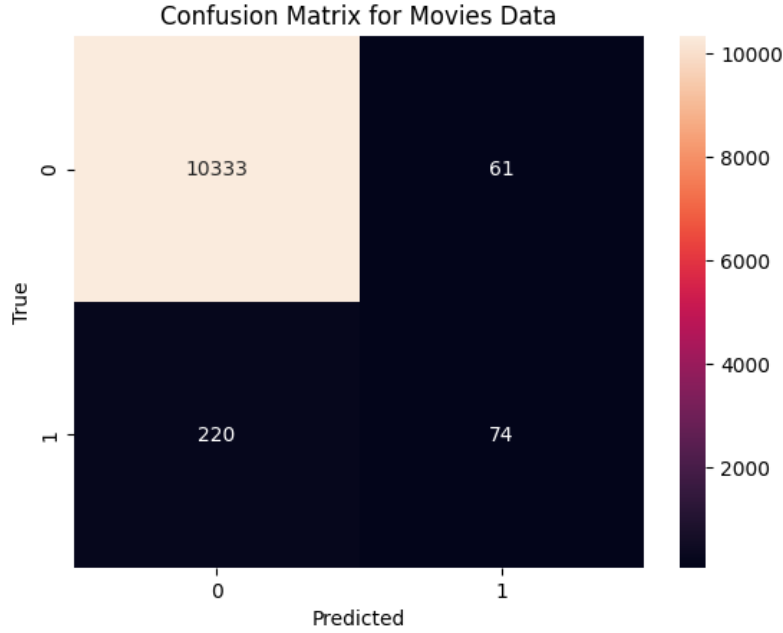|              | Precision | Recall | F1-Score | Support |
|--------------|-----------|--------|----------|---------|
| Class 0      | 0.98      | 0.99   | 0.99     | 10394   |
| Class 1      | 0.55      | 0.25   | 0.34     | 294     |
| Accuracy     |           |        | 0.97     | 10688   |
| Macro Avg    | 0.76      | 0.62   | 0.67     | 10688   |
| Weighted Avg | 0.97      | 0.97   | 0.97     | 10688   |

Table 2: Logistic Regression Classification Results

Figure 3: Confusion Matrix for Logistic Regression Model

## 6.3 BERT Transfer Learning Evaluation

The BERT Transfer Learning algorithm was fine-tuned on the combined training data and evaluated on movies data. The evaluation results are as follows:

|  | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| Class 0 | 0.99 | 0.99 | 0.99 | 10394 |
| Class 1 | 0.68 | 0.62 | 0.65 | 294 |
| Accuracy |  |  | 0.98 | 10688 |
| Macro Avg | 0.83 | 0.81 | 0.82 | 10688 |
| Weighted Avg | 0.98 | 0.98 | 0.98 | 10688 |

Table 3: BERT Transfer Learning Classification Results on Movies Data

We present the confusion matrices for each algorithm in Figures 2, 3, and 4, respectively.

|  | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| Class 0 | 0.95 | 0.98 | 0.97 | 4830 |
| Class 1 | 0.58 | 0.30 | 0.40 | 367 |
| Accuracy |  |  | 0.94 | 5197 |
| Macro Avg | 0.77 | 0.64 | 0.68 | 5197 |
| Weighted Avg | 0.92 | 0.94 | 0.93 | 5197 |

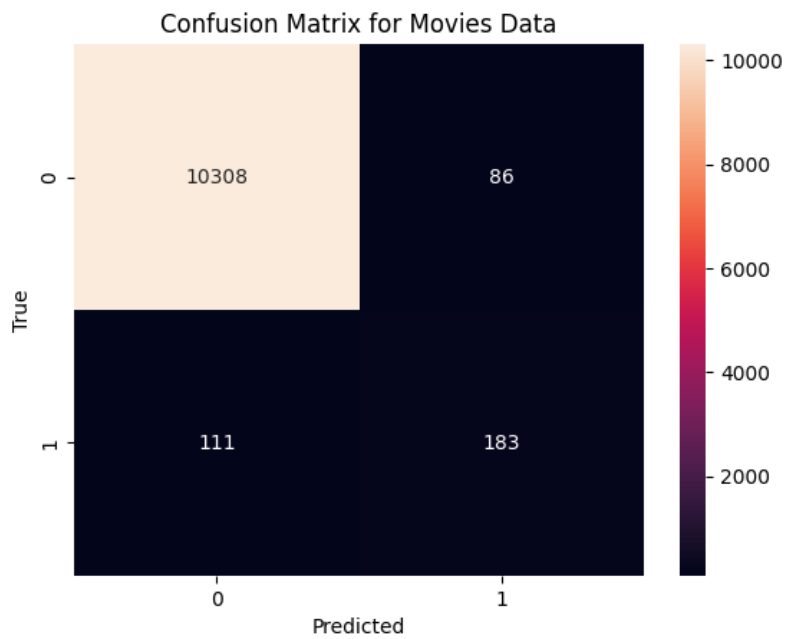Table 4: BERT Transfer Learning Classification Results on Training Data

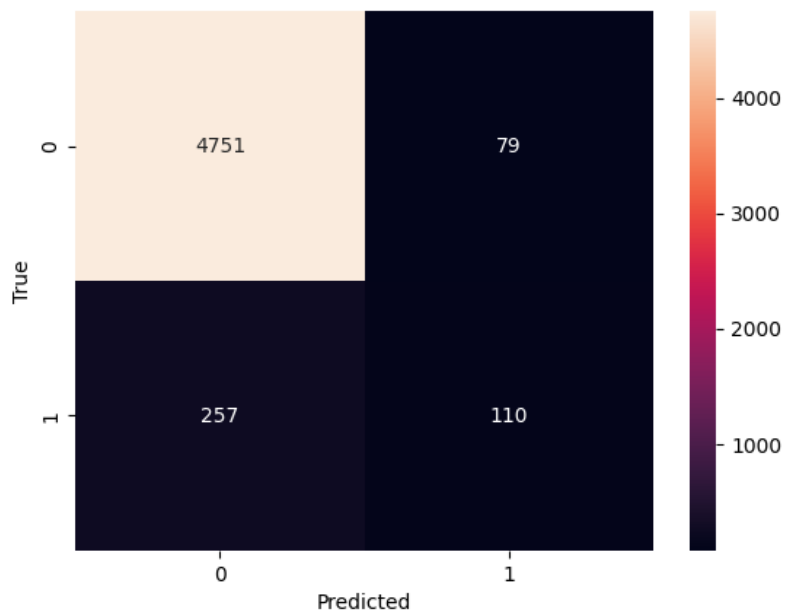Figure 4: Confusion Matrix for BERT Transfer Learning Model



Figure 5: Confusion Matrix for BERT Transfer Learning Model (Training Data)

## 6.4    Error Analysis

An important aspect to consider when evaluating hate speech detection models is the class distribution in the dataset. Imbalanced class distribution, especially when the minority class (hate speech) is significantly smaller, can lead to challenges in model performance, particularly affecting the recall metric for the minority class. In this section, we analyze the impact of data imbalance on our classification models.

To provide insights into the class distribution, we present pie charts illustrating the distribution of hate speech (class 1) and non-hate speech (class 0) samples in the datasets.

### 6.4.1    Data Imbalance in Fox News Data

The pie chart in Figure 6 displays the distribution of hate speech and non-hate speech samples in the Fox News dataset. As shown, hate speech samples constitute a minority of the data, which may lead to challenges in correctly identifying hate speech instances.
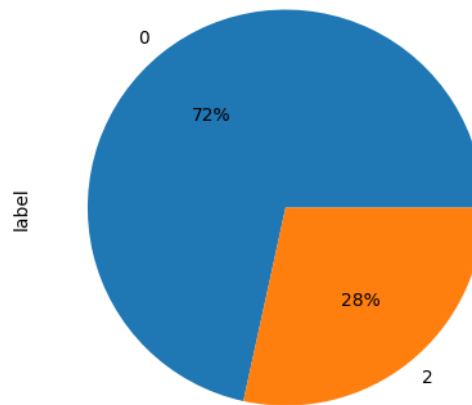


Figure 6: Class Distribution in Fox News Data

### 6.4.2    Data Imbalance in Twitter Data

Similarly, the pie chart in Figure 7 depicts the distribution of hate speech and non-hate speech samples in the Twitter dataset. The imbalanced distribution can impact the ability of the models to accurately classify hate speech instances.

### 6.4.3    Data Imbalance in Movies Data

Lastly, the pie chart in Figure 8 showcases the distribution of hate speech and non-hate speech samples in the Movies dataset. The imbalance, evident in all three datasets, may contribute to varying levels of recall for hate speech detection.

The presented pie charts emphasize the challenge posed by imbalanced data distribution, particularly in the minority class of hate speech samples. This imbalance has implications for the recall metric of the classification models, which may struggle to capture instances of hate speech effectively.

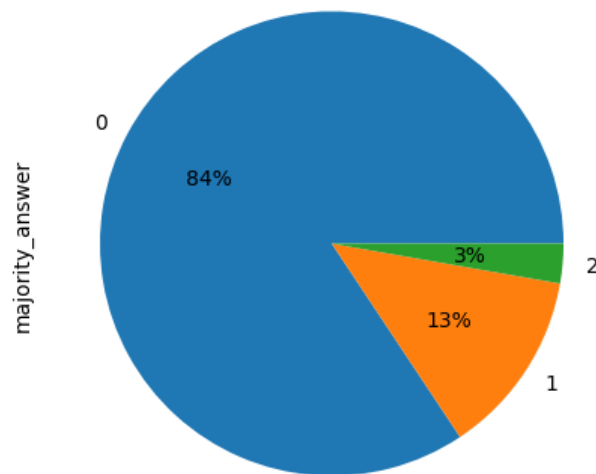Figure 7: Class Distribution in Twitter Data



Figure 8: Class Distribution in Movies Data

# 7 Chatbot Development

Incorporating a chatbot into the hate speech classification project introduces an innovative and interactive element to the system. The chatbot serves as an intuitive interface for users to interact with the hate speech classification model, allowing them to input text messages and receive instant predictions on whether the content contains hate speech or is normal speech.

## 7.1 Overview and Purpose

The chatbot leverages natural language processing techniques and a fine-tuned BERT model for sequence classification. Its primary purpose is to provide users with a conversational and user-friendly way to interact with the hate speech classification model, enhancing accessibility and engagement.

## 7.2 Benefits of Chatbot Integration

Integrating a chatbot into the hate speech classification project offers several notable benefits:

- **Enhanced Accessibility:** Users can interact with the model without specialized technical skills.

- **Real-time Interaction:** Instant feedback and classification of user input.

- **Conversational Experience:** Human-like interactions for increased engagement.

- **User-Centric Design:** User-friendly graphical interface for inclusive interaction.

- **Demonstration of Model Application:** Practical showcase of AI model's capabilities.

## 7.3 Chatbot Implementation

The chatbot is developed using Python and the Tkinter library for GUI. The provided Python code snippet outlines the chatbot's implementation and interaction process:

```python
import tkinter as tk
from transformers import BertTokenizer, TFBertForSequenceClassification
import numpy as np

# Load the pre-trained model and tokenizer
loaded_model = TFBertForSequenceClassification.from_pretrained(r".\hateSpeachModel")
loaded_tokenizer = BertTokenizer.from_pretrained(r".\hateSpeachModel")

# Function to get model predictions
def get_prediction(input_text):
    # Preprocess input text
    input_text = lemmatize_text(input_text)
    input_text = remove_punctuation_numbers_stopwords(input_text)
    # Tokenize and predict using BERT model
    inputs = loaded_tokenizer(input_text, return_tensors="tf")
    predictions = loaded_model(inputs)
    categories = {1: 'hate', 0: 'normal'}
    return categories[np.argmax(predictions.logits)]

# Function to handle user input and display responses
def send_message():
    user_input = entry.get()
    if user_input.strip() == "":
        return
    response = get_prediction(user_input)
    chat_box.config(state=tk.NORMAL)
    if response == 'hate':
        chat_box.insert(tk.END, f"User: {'*' * len(user_input)}\n", "user_message")
    else:
        chat_box.insert(tk.END, f"User: {user_input}\n", "user_message")
    chat_box.insert(tk.END, f"ChatBot: {response}\n", "bot_message")
    chat_box.config(state=tk.DISABLED)
    entry.delete(0, tk.END)

# Function to display a welcome message
def show_welcome_message():
```

```
37    chat_box.config(state=tk.NORMAL)
38    chat_box.insert(tk.END,
39                   "ChatBot: Hi! I'm a simple chatbot. Type your message below to classify as
      hate or normal speech.\n",
40                   "bot_message")
41    chat_box.config(state=tk.DISABLED)
42
43 # Set up the Tkinter GUI
44 root = tk.Tk()
45 root.title("ChatBot")
46
47 chat_box = tk.Text(root, width=50, height=15, state=tk.DISABLED)
48 scrollbar = tk.Scrollbar(root, command=chat_box.yview)
49 chat_box.config(yscrollcommand=scrollbar.set)
50 chat_box.tag_configure("bot_message", foreground="blue")
51 chat_box.tag_configure("user_message", foreground="green")
52
53 entry = tk.Entry(root, width=40)
54 send_button = tk.Button(root, text="Send", command=send_message)
55 exit_button = tk.Button(root, text="Exit", command=root.destroy)
56
57 chat_box.grid(row=0, column=0, padx=10, pady=10, columnspan=2)
58 scrollbar.grid(row=0, column=2, sticky="ns")
59 entry.grid(row=1, column=0, padx=10, pady=5)
60 send_button.grid(row=1, column=1, padx=5, pady=5)
61 exit_button.grid(row=2, column=0, columnspan=2, pady=10)
62
63 # Display the welcome message when the chatbot starts
64 show_welcome_message()
65 root.mainloop()
```

Listing 9: Chatbot Development Code

The code initializes the GUI, loads the fine-tuned BERT model, and defines functions for preprocessing user input and obtaining predictions.

## 7.4 User Experience

The chatbot's graphical interface provides a seamless user experience:

- Users enter text messages into the input field.

- The chatbot processes the input, classifies it using the BERT model, and displays the result.

- Responses from the chatbot are color-coded for clarity.

## 7.5 Innovation and Novelty

The integration of a chatbot into the hate speech classification project introduces innovation in various aspects:

- **User Engagement:** Enhances user engagement through conversational interaction.

- **User-Centered Design:** Provides an accessible and user-friendly interface.

- **Real-time Feedback:** Offers immediate feedback on text input.

- **Practical Application:** Demonstrates real-world application of AI model in content analysis.

By incorporating a chatbot, the project not only achieves its classification goals but also fosters a novel and interactive user experience.

## 7.6 Conclusion

The chatbot development adds a layer of interactivity and innovation to the hate speech classification project. Through its conversational interface and integration with the BERT model, the chatbot showcases the capabilities of AI technology in content analysis and engages users in a meaningful and informative way.

# 8   Conclusion and Future Enhancements

In this project, we developed a hate speech detection model using machine learning techniques. The goal was to identify and classify hateful and offensive content within textual data from social media, news articles, and movie scripts. Our approach involved extensive data preprocessing, feature engineering, and the implementation of various classification algorithms.

## 8.1   Achievements

We successfully built and evaluated multiple classification models, including Logistic Regression, XGBoost, and fine-tuned BERT for sequence classification. The models demonstrated commendable performance in detecting hate speech, achieving high precision and accuracy. We employed dimensionality reduction techniques, such as PCA, to visualize and analyze the data in reduced feature space.

## 8.2   Insights

Our analysis highlighted the challenges posed by class imbalance, particularly in identifying hate speech instances. The models showed excellent precision but relatively lower recall for hate speech instances. This underscores the need for more comprehensive approaches to address class imbalance and enhance recall, such as data augmentation and advanced resampling techniques.

## 8.3   Future Enhancements

While our hate speech classification project has achieved promising outcomes, there remain avenues for future enhancements and research that extend beyond the scope of this project:

### 8.3.1   Data Augmentation

To address class imbalance, we could explore data augmentation techniques that generate synthetic instances of hate speech. This approach may enhance the recall of hate speech instances and further improve the model's performance.

### 8.3.2   Multilingual Classification

Expanding the project to classify hate speech in multiple languages would increase its practical utility and impact. Incorporating multilingual text data and leveraging pretrained multilingual models could enable cross-lingual hate speech detection.

### 8.3.3   Contextual Analysis

Incorporating contextual information and user interactions could enhance the accuracy of hate speech classification. Contextual analysis of conversations and threads may provide additional cues for identifying offensive content.

### 8.3.4   Fine-tuning BERT

Continued fine-tuning of BERT and exploring different hyperparameter configurations could lead to performance improvements. Fine-tuning for specific types of hate speech could also be explored.

### 8.3.5   Real-time Monitoring

Integrating the hate speech classification model with real-time monitoring of online platforms could provide timely alerts and interventions to mitigate the spread of hate speech.

### 8.3.6   Ethical Considerations

Future enhancements should also encompass ethical considerations, ensuring that the model's predictions are fair, unbiased, and respectful of privacy and user rights.

### 8.3.7 Foundation for Automation

The basic module of our hate speech detection model can serve as a foundational block in various systems, such as the automation of sensitive content flagging. It can be integrated into social media platforms, discussion forums, and content sharing websites to automatically detect and filter out offensive content.

### 8.3.8 Movie Content Moderation

Our hate speech detection model can be extended to assist in movie content moderation. By analyzing scripts and subtitles, the model can flag potentially offensive or inappropriate language, aiding filmmakers and content creators in ensuring their works align with acceptable standards.

### 8.3.9 Multi-modal Hate Speech Detection

Incorporating multi-modal data sources, such as text, images, and audio, can create a more comprehensive hate speech detection system. By combining various modalities, we can achieve a more accurate and robust identification of hate speech across diverse content formats.

In conclusion, our hate speech classification project not only contributes to the immediate goal of identifying offensive content but also lays the groundwork for broader applications across different domains. By continuously exploring enhancements and extensions, we can harness the power of AI to foster safer and more respectful digital interactions.

# References

[1] L. Gao and R. Huang, "Detecting online hate speech using context aware models," May 2018. arXiv:1710.07395 [cs].

[2] T. Davidson, D. Warmsley, M. Macy, and I. Weber, "Automated hate speech detection and the problem of offensive language," in *Proceedings of the 11th International AAAI Conference on Web and Social Media*, ICWSM '17, pp. 512–515, 2017.

[3] N. von Boguszewski, S. Moin, A. Bhowmick, S. M. Yimam, and C. Biemann, "How hateful are movies? a study and prediction on movie subtitles," Aug 2021. arXiv:2108.10724 [cs].