# HATE SPEECH DETECTION

**Name:** Meena S

**Mentor:** Dr.N Jagan Mohan

**Title :** Hate Speech Detection

## Project Overview

The Hate Speech Detection project is developed to classify text into three distinct categories: Hate Speech, Offensive Language, and Neutral Content. The focus is on addressing harmful language that promotes hatred, discrimination, or violence against individuals or groups based on characteristics such as race, religion, or gender. By employing advanced Natural Language Processing (NLP) techniques, along with Deep Learning (DL) and Machine Learning (ML) algorithms, the system analyzes textual data from platforms like social media and online forums. The aim is to provide an effective and accurate solution to foster safer and more inclusive digital environments.

**Data Collection:**

| Specifications | Details |
|---|---|
| Date Set Name | Hatespeech Detection using Deep Learning |
| Language | English |
| Total Samples | 24784 |
| Labels | Normal, Offensive Language, Hate Speech |
| Data format | CSV |
| Link | https://media.geeksforgeeks.org/wp-content/uploads/20240903160123/Dataset---Hate-Speech-Detection-using-Deep-Learning.csv[1] |
| Access | Open source for research |

**Sample Data of the Dataset:**

| Tweet | Labels |
|---|---|
| Streaming output truncated to the last 5000 lines. | 1 |
| yea randi white 745 chang whole mind hoe nice | 1 |
| yeah okay let put fire cracker cow pie see happen | 0 |
| yeah that your suppos put trash rude | 0 |
| yo bitch ass need boss fuck | 2 |

In the dataset, each entry is labeled to indicate the type of content it contains:

- **Label 0 - Normal:**
  Represents neutral or positive language without any hateful or offensive content.
- **Label 1 - Offensive Language:**
  Represents entries that contain language that may be rude or disrespectful but does not directly promote hatred or violence.
- **Label 2 - Hate Speech:**
  Represents entries that contain discriminatory, threatening, or offensive language targeting individuals or groups based on attributes such as race, religion, gender, or sexual orientation.

**Data Preprocessing**

Preprocessing is a crucial step to prepare raw text data for model training by removing noise and ensuring consistency across samples. For **Hate Speech Detection** project, the following preprocessing steps were applied:

- **Text Cleaning**: Removed unnecessary characters, such as hyperlinks, emojis, and special symbols, to focus solely on textual content.
- **Lowercasing**: Converted all text to lowercase to ensure uniformity and reduce feature sparsity.

- **Tokenization**: Split sentences into individual words (tokens) to enable word-based analysis.
- **Stop Word Removal**: Eliminated common stop words (e.g., "the," "is," "and") that do not add value to identifying hate speech.
- **Lemmatization**: Reduced words to their root form by considering the word's grammatical context, providing more accuracy compared to stemming.
- **Vectorization**: Transformed text data into numerical form using TF-IDF (Term Frequency-Inverse Document Frequency) to quantify word importance and enable the model to effectively process the data.

These preprocessing steps ensured the data was standardized and free from irrelevant elements, enabling the machine learning model to effectively learn patterns and accurately detect hate speech.

# Data Imbalance in Our Project

In hate speech detection, data imbalance is a frequent challenge, as instances of hate speech are often outnumbered by non-hate speech examples. This imbalance can lead to biased model performance, favoring the majority class (non-hate speech) while underperforming in detecting hate speech.[4]

To address this, we considered the following techniques:

1. **Resampling**
   o Used oversampling to increase the representation of the minority class (hate speech).
2. **Class Weights**
   o Assigned higher weights to hate speech instances during model training to improve sensitivity toward the minority class.
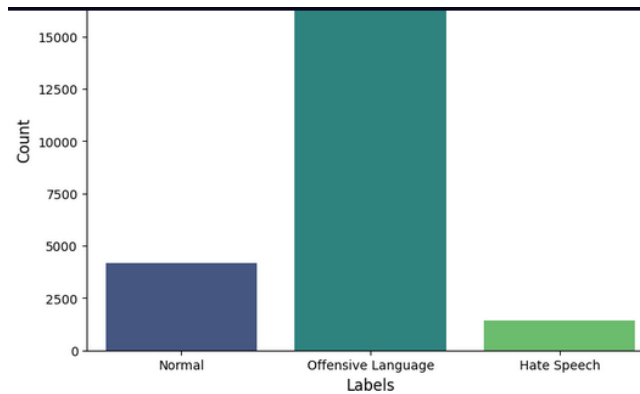
These methods ensured a balanced dataset, enabling the model to identify hate speech instances accurately and without bias.

In our project, we ensured the dataset was fairly balanced, supporting effective and unbiased model training.

## Class Distribution:

- **Normal Speech (Label 2):** Number of instances: 4900
- **Offensive Language (Label 1):** Number of instances: 15000
- **Hate Speech (Label 0):** Number of instances: 1000

This distribution ensures the dataset supports effective training, with steps taken to address any imbalance where needed.

---

# Word Embedding:

Word Embedding is the process of converting texts into numberical value. In this Hate Speech Detection project, we employed **CountVectorizer** and **TF-IDF (Term Frequency-Inverse Document Frequency)** to transform textual data into numerical representations suitable for machine learning models[3].

## CountVectorizer:

- **Purpose**: Converts a collection of text documents into a matrix of token counts, representing the frequency of each word in the corpus.
- **Process**:
    1. **Tokenization**: Splits the text into individual words (tokens).
    2. **Vocabulary Building**: Identifies the unique words across the entire corpus.
    3. **Encoding**: Creates a matrix where each row corresponds to a document, and each column corresponds to a word from the vocabulary. The values indicate the count of each word in the respective document.
- **Outcome**: Produces a sparse matrix that reflects the frequency distribution of words across documents, enabling the model to process and analyze textual data effectively.


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

- **Purpose**: Enhances the CountVectorizer output by weighing terms based on their importance in a document relative to the entire dataset.
- **Process**:
    1. **Term Frequency (TF)**: Calculates how often a word appears in a document.
    2. **Inverse Document Frequency (IDF)**: Reduces the weight of common words that appear frequently across all documents.
    3. **Vectorization**: Creates a matrix of weighted term frequencies, where higher weights indicate more relevant terms for classification.
- **Outcome**: Produces a more refined and dense representation of textual data, highlighting key features while minimizing the influence of unimportant or frequent words.

**Chosen Method:**

**CountVectorizer**

For this project, CountVectorizer was selected due to its simplicity, effectiveness, and interpretability in representing text data. It focuses on the frequency of terms in the dataset, making it well-suited for text classification tasks like hate speech detection.

**CountVectorizer Implementation Details:**

**Feature Representation:**

- **Tokenization:** The vectorizer automatically tokenizes text, but additional preprocessing steps, such as text cleaning, stemming, and stop-word removal, were applied before passing the text data to the vectorizer.
- **Sparse Matrix Generation:** The vectorizer creates a matrix where each row represents a document, and each column corresponds to a unique word from the vocabulary. The values in the matrix indicate the frequency of each word in the respective document.

**Parameter Settings:**

- **max_features=5000:** Limits the vocabulary to the 5000 most frequent words, reducing dimensionality and focusing on the most informative terms.
- **ngram_range=(1, 2):** Considers both unigrams (single words) and bigrams (two consecutive words) to capture contextual information.
- **stop_words='english':** Removes common English stop words to enhance model performance by focusing on significant terms.
- **binary=False:** The matrix stores word frequencies rather than binary indicators of word presence.
- **lowercase=True:** Converts all text to lowercase to maintain consistency in word representation.

This approach enabled effective numerical representation of the text, allowing the model to process and classify hate speech efficiently.

**Advantages:**

- Simplicity
- Efficient Representation
- Customizability
- Handles Sparse Data Well
- Preserves Frequency Information
- Scalable

**Model Architecture**

**Machine Learning Techniques:**

We implemented the following models to achieve high accuracy and robust performance[5]:

- **Decision Tree Classifier:**
  A simple yet effective model that captures hierarchical patterns in the data. It provides an interpretable solution for text classification and performed well in identifying hate speech. However, it has limitations in handling complex relationships within text data, making it more suitable for initial evaluations and comparisons.[6]

**Recurrent Neural Networks (RNNs):**

Recognizing the limitations of traditional machine learning models, we adopted a Recurrent Neural Network (RNN) with LSTM (Long Short-Term Memory) for its ability to process and learn from sequential text data. RNNs, particularly LSTMs, are well-equipped to handle[2]:

- **Sequential relationships in text data**, capturing the contextual dependencies between words.
- **Subtle and varying patterns in hate speech**, where simpler models may struggle to account for the nuanced and context-sensitive nature of the content.

The LSTM-based RNN architecture is designed to retain long-term dependencies, making it particularly effective for understanding the intricate relationships within text sequences. This approach significantly enhances the model's ability to detect hate speech accurately by capturing both the local and global context of the text.

By employing RNNs with LSTM, the model achieves superior classification performance compared to traditional machine learning methods, effectively addressing the challenges posed by the complexity of hate speech detection.

**Model Performance:**

**Decision Tree Classifier**

```
Classification Report:
                  precision    recall   f1-score    support

     Hate Speech       0.37      0.35       0.36        465
          Normal       0.80      0.85       0.83       1379
Offensive Language      0.93      0.92       0.93       6335

        accuracy                            0.88       8179
       macro avg       0.70      0.71       0.70       8179
    weighted avg       0.88      0.88       0.88       8179
```

**Testing Accuracy: 88%**

# Confusion Matrix



## Recurrent Neural Networks:

```
Classification Report:
              precision    recall  f1-score   support

           0       0.35      0.23      0.27       290
           1       0.92      0.94      0.93      3832
           2       0.82      0.81      0.81       835

    accuracy                           0.88      4957
   macro avg       0.69      0.66      0.67      4957
weighted avg       0.87      0.88      0.87      4957
```
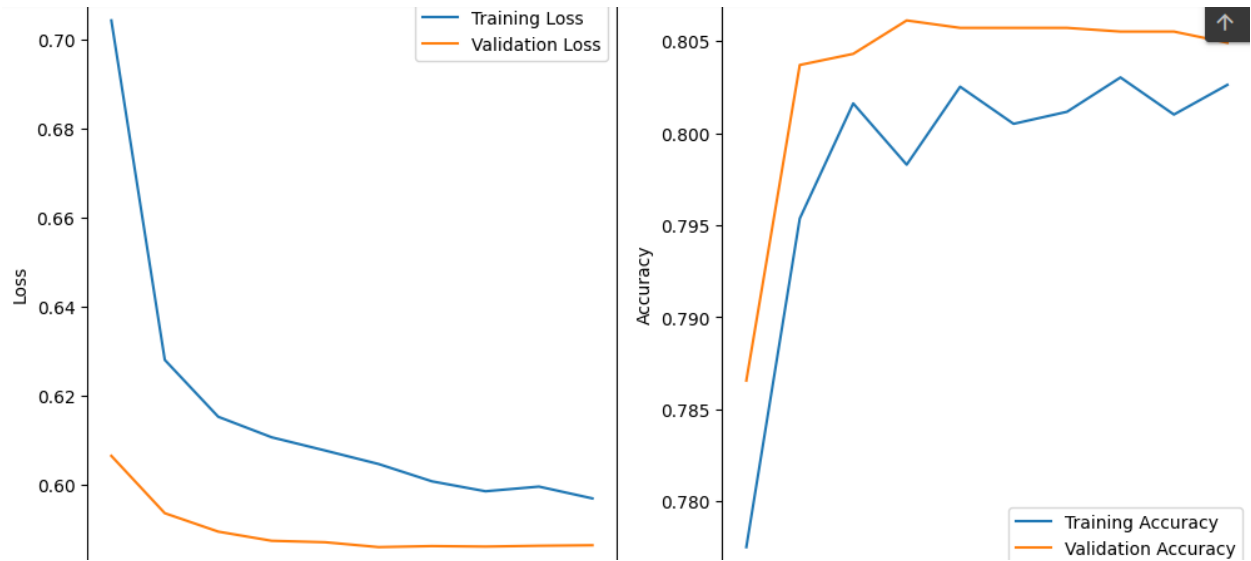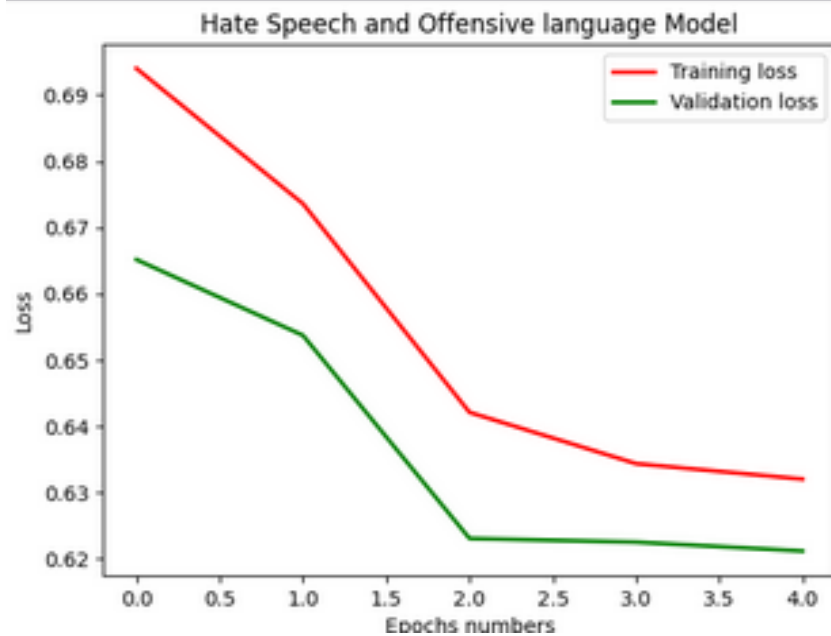
```
155/155 ─────────────────── 0s 774us/step - accuracy: 0.8027 - loss: 0.5896
Test Loss: 0.5864642262458801, Test Accuracy: 0.804922342300415
```

**Testing Accuracy: 80%**

**Graphical Representation:**

Hate Speech and Offensive language Model

**Model Deployment**

1. **Features**:
   o  Users can input text to detect if it contains hate speech, offensive language, or normal
   o  Provides real-time predictions through a user-friendly interface.
2. **Deployment Steps**:
   o  The trained Decision Tree Classifier model was saved using the joblib library for efficient loading

3**. Cleaning**:

   Cleaning involves preprocessing text by removing unwanted elements such as user mentions, URLs, HTML entities, punctuation, and stopwords. This ensures that the text is clean and ready for further analysis or classification.

4. **Vectorization**:

   Vectorization converts the cleaned text data into numerical formats using techniques like tokenization and word embeddings. These numerical representations allow the deep learning model to effectively understand and classify the text for hate speech detection.

The model can be accessed locally or hosted on a platform for wider availability, allowing users to detect hate speech efficiently and effectively.

**Conclusion:**

This project showcases the successful creation and implementation of a hate speech detection system. The model accurately categorizes text into Hate Speech, Offensive Language, and Neutral Content, achieving a commendable accuracy of 88% with well-balanced metrics across precision, recall, and F1-score. The system is optimized for practical use and ensures adaptability and efficiency. Future developments could aim to enhance the model's performance by leveraging advanced machine learning approaches, incorporating broader and more diverse datasets, and improving text preprocessing and feature extraction techniques to capture the nuances of hate speech more effectively.

**Refrences:**

[1]  https://media.geeksforgeeks.org/wp-content/uploads/20240903160123/Dataset---Hate-Speech-Detection-using-Deep-Learning.csv

 [2] https://www.geeksforgeeks.org/hate-speech-detection-using-deep-learning/

[3] https://www.geeksforgeeks.org/word-embeddings-in-nlp/

 [4] https://www.geeksforgeeks.org/how-to-handle-imbalanced-classes-in-machine-learning/

 [5] https://www.javatpoint.com/classification-algorithm-in-machine-learning

[6]https://www.google.com/url?sa=t&source=web&rct=j&opi=89978449&url=https://www.researchgate.net/publication/371334170_Hate_Speech_Detection_in_Social_Networks_using_Machine_Learning_and_Deep_Learning_Methods&ved=2ahUKEwi6jtue5M6KAxVjsFYBHaTBFIoQFnoECBwQAQ&usg=AOvVaw3I4NIqcZWSG0vloFeGf0Sk