

Text Classification For Emotion Tracking

Olumuyiwa Obamakin
University of Essex
Computer Science and Electrical Engineering Department
oo21696@essex.ac.uk

Abstract: The goal of this project was to develop a text classification model capable of tracking emotions in Tweets.

I. INTRODUCTION

Sentiment analysis is a common natural language processing (NLP) task where the objective is to classify text into predefined categories or labels based on its content. In this case, the dataset consisted of Tweets annotated with six fundamental emotions: anger, fear, joy, love, sadness, and surprise. The primary objective was to train a model to predict these emotions based on the given text data.

II. METHADODOLOGY

Model Selection

I decided to use a Logistic Regression model with a Term Frequency-Inverse Document Frequency (TF-IDF) vectorizer. Logistic Regression was chosen as it is a widely used and interpretable model for text classification tasks. Its simplicity and efficiency make it suitable for problems with relatively clean well-structured data such as the text.csv I was operating with.

Data Preprocessing

The first step was to preprocess the text data to ensure that the model could learn meaningful patterns from it. Several preprocessing steps were taken, including:

1. **Cleaning the CSV file:** The dataset was loaded and unnecessary columns such as the 'Unnamed: 0' column were dropped. This allowed me to focus on the relevant features (text and the emotion labels).
2. **Text Preprocessing:** Each text entry underwent several key preprocessing steps to prepare the data for training:
 - **Removal of URLs and special characters:** URLs, hashtags, and special characters in tweets add noise and don't

contribute significantly to understanding the sentiment or emotion.

- **Lowercasing all words:** Lowercasing standardizes the text, preventing the model from learning redundant features.
- **Tokenization:** Tokenizing splits the text into individual words, enabling the model to process each word separately and understand its context in the sentence.
- **Stopword removal:** Stopwords are common words which do not carry much semantic value. Removing them ensures the model focuses on more informative words.
- **Lemmatization:** ensures that different word forms are treated as the same word, which helps in capturing the core meaning.

These preprocessing steps were crucial given the nature of Tweets. By applying these preprocessing techniques, the text data was cleaned and standardized. Allowing the model to focus on the core ideas and emotions expressed in each tweet. This was vital for improving the accuracy and effectiveness of my classification model, as it helped reduce the influence of irrelevant or redundant information that could be misleading.

III. TRAINING MY MODEL

Using Logistic Regression/TF-IDF

After preprocessing the text data, it was converted into numerical features using the TF-IDF. TF-IDF calculates the importance of each word in a document by considering both its frequency in the document and its rarity across the entire dataset. This transformation allowed the raw text to be represented as numerical vectors, which could be fed into the Logistic Regression model. The logistic Regression model was then

trained on 90% of the data to learn the relationship between the tweets and their associated emotions. The model evaluated on the remaining 10% of the data to assess its performance in predicting the emotions.

Results

The performance of the Logistic Regression model was evaluated using several metrics:

1. **Accuracy:** The model achieved an accuracy of 90%, which is a good result considering the complexity of emotion classification from text.
2. **Classification Report:** The model's precision, recall, and F1-score were calculated for each emotion, providing a detailed assessment of its performance. The model performed well in detecting emotions like joy and sadness but struggled with emotions such as fear and surprise. With the lowest F1 score of 0.73, and love, which showed similarly poor results. This could be attributed to the difficulty in distinguishing subtle emotional cues in text, as well as the fact that love and surprise had the least amount of support (tweets labeled with these emotions) in the dataset. The limited number of examples for these emotions may have hindered the model's ability to learn their characteristics effectively.
3. **Precision Recall and ROC Curves:** The precision recall curve and ROC curve showed that the model performed well for 'Joy' and 'Sadness', which exhibited higher precision and recall scores. This could be due to the more distinct and recognizable linguistic cues associated with these emotions in the text, making them easier for the model to identify. In contrast, 'Love' and 'Surprise' showed lower performance, likely because these emotions can be more subtle or context-dependent in text, making them harder for the model to detect accurately. Despite this, the model still demonstrated an

overall strong performance.

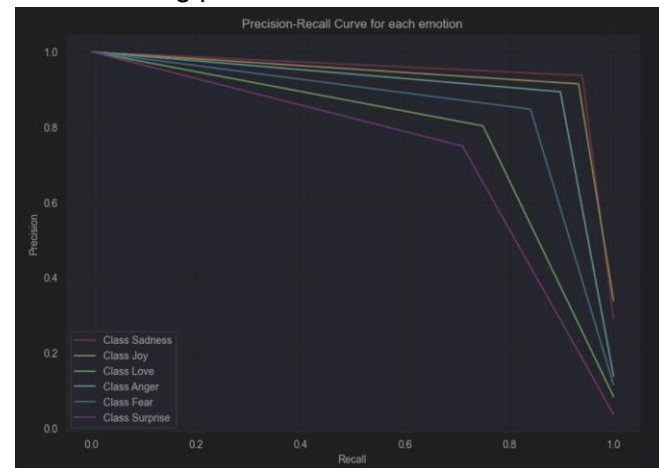


Figure 1: Precision Recall Curve.

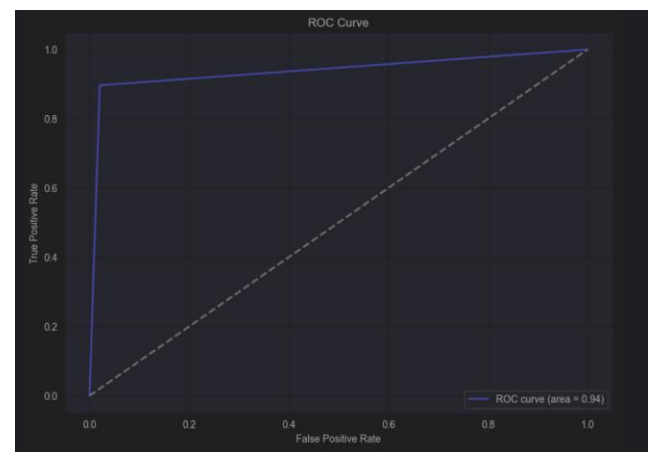


Figure 2: ROC Curve.

Potential Improvements

Hyperparameter Tuning: Fine-tuning the hyperparameters of the Logistic Regression model, such as adjusting the regularization strength. Could help the model achieve better results, especially in improving accuracy on challenging emotions like 'Surprise'.

Advanced Models: More advanced models such as Support Vector Machines, or neural networks like Recurrent Neural Networks might yield better results. Especially for differentiating between subtle emotions like 'Fear' and 'Surprise'.

IV. CONCLUSION

To conclude I was successful in creating a text classification model that can be used to track emotions in tweets. Using Logistic Regression with TF-IDF provided a solid foundation for the task, but there are opportunities to improve performance by exploring more complex models and optimizing hyperparameters. The project

also highlights the importance of data preprocessing in ensuring that the model can effectively learn from the text data. Future work could focus on improving the model's ability to distinguish between more subtle emotions and further investigating deep learning approaches to enhance classification accuracy.