



Data Science Intern at Data Glacier

Project: Hate Speech Detection using Transformers (Deep Learning)

Week 12: Deliverables

Name: Pranav Boyapati

University: University of Texas at Dallas

Email: pranav13b@gmail.com

Country: United States of America

Specialization: Data Science

Batch Code: LISUM34

Date: July 31, 2024

Submitted to: Data Glacier

Table of Contents:

1. Project Plan	3
2. Problem Statement	3
3. Data Collection	3
4. Data Preprocessing	4
4.1. Text Cleaning.....	4
4.1.1. Lower Case.....	4
4.1.2. Remove Punctuation.....	4
4.1.3. Remove URL.....	4
4.1.4. Remove @tags.....	4
4.1.5. Remove Special Characters.....	4
4.2. Preprocessing Operations	4
4.2.1. Tokenization	4
4.2.2. Remove Stop Words.....	5
4.2.3. Lemmatization.....	5
4.2.4. WorldCloud	5
4.3. Feature Extraction.....	6
4.3.1. TF-IDF Model	6
4.4. Split the Dataset into Train and Test	6
5. Build the Model	6
5.1. CNN with LSTM	6
6. Result Evaluation	7
6.1. Evaluation Criteria	7
7. Application Design	8
7.1. Turning ML Model Into Flask Framework	9
7.2. Running Procedure	10
8. Conclusion	11
Reference	12

1. Project Plan

Weeks	Date	plan
Weeks 07	July 19, 2024	Problem Statement, Data Collection, Data Report
Weeks 08	July 26, 2024	Data Preprocessing (Text Cleaning)
Weeks 09	August 2, 2024	Data Preprocessing (Preprocessing Operation + Feature Extraction)
Weeks 10	August 9, 2024	Building the Model
Weeks 11	August 16, 2024	Model Result Evaluation
Weeks 12	August 23, 2024	Build ML APP Based on Flask
Weeks 13	August 30, 2024	Final Submission (Report + Code + Presentation)

2. Problem Statement

The term hate speech is understood as any type of verbal, written or behavioural communication that attacks or uses derogatory or discriminatory language against a person or group based on what they are, in other words, based on their religion, ethnicity, nationality, race, color, ancestry, sex or another identity factor. In this problem, we will take you through a hate speech detection model with Machine Learning and Python.

Hate Speech Detection is a task of sentiment classification. So, for training, a model that can classify hate speech from a certain piece of text can be achieved by training it on a data that is used to classify sentiments. So, for the task of hate speech detection model, we will use the Twitter tweets to identify tweets containing Hate speech.

3. Data Collection

The Data is about Twitter hate Speech taken from Kaggle [1] which contains the 3 number of features and 31962 number of observations. Dataset using Twitter data, it was used to research hate-speech detection. The text is classified as: hate-speech, offensive language, and neither. Due to the nature of the study, it is important to note that this dataset contains text that can be considered racist, sexist, homophobic, or offensive.

Table 1: Data Information

Total number of observations	31962
Total number of files	1
Total number of features	3
Base format of the file	csv
Size of the data	3.12 MB

4. Data Preprocessing

In part, we explain the data preprocessing approach that we apply in the text data.

4.1 Text Cleaning

First, we clean our text because it was so messy data.

4.1.1 Lowercase

Converting a word to lower case (NLP -> nlp). Words like Racism and racism mean the same but when not converted to the lower case those two are represented as two different words in the vector space model (resulting in more dimensions). Therefore, we convert all text word into lower case letter.

4.1.2 Remove Punctuation

It is important to remove the Punctuation because is not important. Therefore, we remove that Punctuation to do that we use regular expression.

4.1.3 Remove URLs

In this part, we remove URLs because we are working on hate speech application which detect the hate and free speech and to get the output, we need to give only text not URLs therefore, we remove the URLs because we need only clean text input.

4.1.4 Remove @tags

In this part, we remove @tags which used when we mentioned someone So, it's doesn't concern to our application therefore, we remove @tags by using regular expressions.

4.1.5 Remove Special Characters

Remove Special Characters is essentially the following set of symbols [!"#\$%&'()*+,-./:;<=>?@[^_`{|}~] which basically don't have meaning. Therefore, we remove that kind of symbols because we don't need that. To remove we use python isalnum method.

4.2 Preprocessing Operation

In this part, we implement the preprocessing operation

4.2.1 Tokenization

Tokenization is breaking the raw text into small chunks. Our text data is into paragraph so to convert into work tokenize we use nltk word_tokenize library. These tokens help in

understanding the context or developing the model for the NLP. The tokenization helps in interpreting the meaning of the text by analyzing the sequence of the words.

4.2.2 Removing StopWords

StopWords is basically ‘a,’ ‘is,’ ‘the,’ ‘are’ etc. If we see our dictionary, then these words do not have meaning and don’t need that to build Hate speech detection application. To remove stop words from a sentence, we divide text into words which we did above in tokenization and then remove the word if it exists in the list of stop words provided by NLTK. To do that, we first import the StopWords collection from the nltk.

4.2.3 Lemmatization

Lemmatization is the process of grouping together the different inflected forms of a word so they can be analyzed as a single item. Lemmatization is like stemming but it brings context to the words. So, it links words with similar meanings to one word. Like the word Intelligently, intelligence, convert into root form intelligent.

4.2.4 WordCloud

A Wordcloud is a visual representation of text data, which is often used to depict keyword metadata on websites, or to visualize free form text. Tags are usually single words, and the importance of each tag is shown with font size or color.

Below you see the hate speech and free speech WorldCloud:

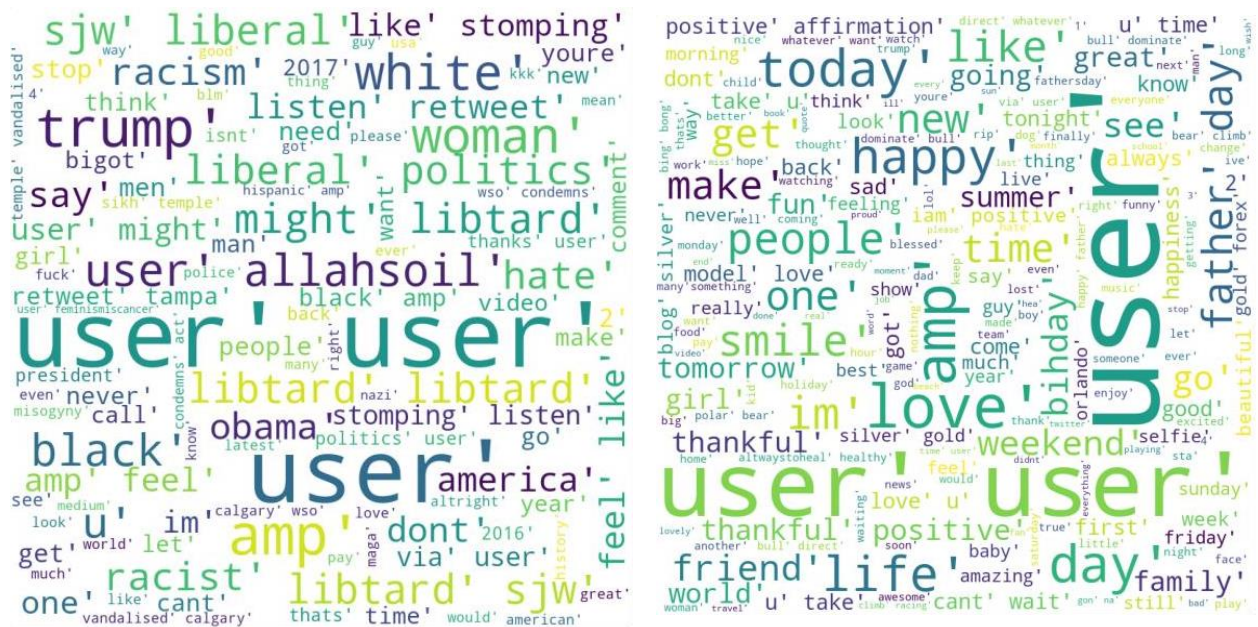


Figure 1: Hate Speech vs FREE Speech WordCloud

4.3 Feature Extraction

4.3.1 TF-IDF Model

Once the dictionary is ready, we apply Term Frequency-Inverse Document Frequency (TF-IDF) model, and we take 2000 most frequent words from dictionaries for each Hate/Free Speech of the whole dataset. Each word count vector contains the frequency of 2000 words in the whole dataset file.

4.4 Split the Data into Train into Test

In this part, we split the data into Train. And we split 80% for training and 20% for test. Data splitting is when data is divided into two or more subsets. Typically, with a two-part split, one part is used to evaluate or test the data and the other to train the model. Data splitting is an important aspect of data science, particularly for creating models based on data.

4.5 Build the Model

In this part, we build the CNN with LSTM Model using Tensorflow.

4.5.1 CNN with LSTM

The CNN LSTM architecture involves using Convolutional Neural Network (CNN) layers for feature extraction on input data combined with LSTMs to support sequence prediction. This architecture was originally referred to as a Long-term Recurrent Convolutional Network or LRCN model, although we will use the more generic name “CNN LSTM” to refer to LSTMs that use a CNN as a front end in this lesson. [2].

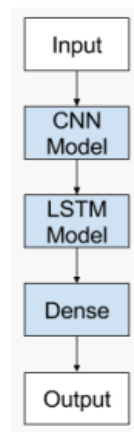


Figure 2: CNN with LSTM Model [3]

5. Result Evaluation and Discussion

In this chapter, we evaluate our Result and also define the evaluation criteria to calculate the performances of our best classification model.

5.1 Evaluation Criteria

The confusion matrix was used to evaluate the classification models throughout the training process. The confusion matrix is a table that compares predicted and actual outcomes. It is frequently used to describe a classification model's performance on a set of test data.

Table 1: Confusion Matrix

Class	Predicted Negative	Predicted Positive
Actual Negative	TN	FP
Actual Positive	FN	TP

Important metrics were constructed from the confusion matrix in order to evaluate the classification models. In addition to the accurate classification rate or accuracy, other metrics for evaluation included True Positive Rate (TPR), True Negative Rate (TNR), False Positive Rate (FPR), False Negative Rate (FNR), Precision, F1 score, and Misclassification rate.

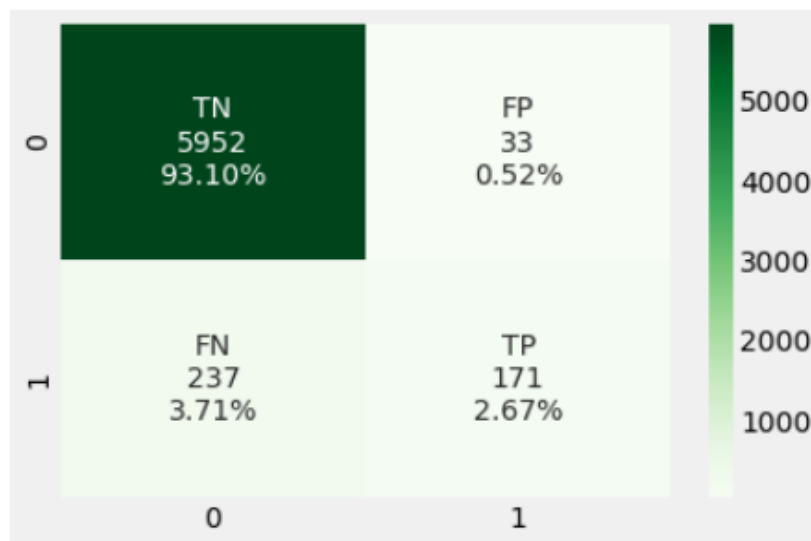


Figure 3: Confusion Matrix

Table 6.2 shows the final result that we evaluate on the basis of confusion matrix result

Table 6.2: Final Results

Classifiers	Accuracy	Precision	TPR	FPR	F1 Score	Error Rate	Specificity
CNN with LSTM	0.9577	0.8382	0.4191	0.0055	0.5588	0.0422	0.9944

Below you can see the visualization result of above table as well.

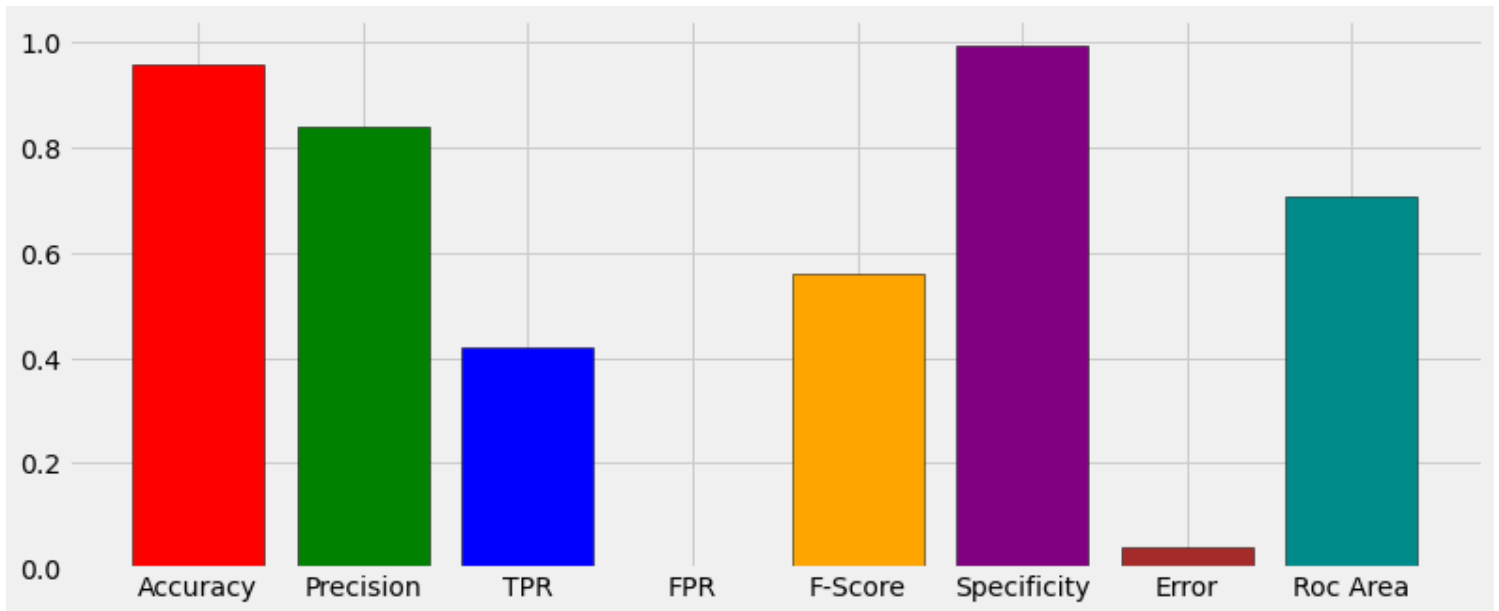


Figure 4: Visualization of Final Result

5. Application Design

In this chapter, we develop a Model & Deploy It with Flask. Our model systems workflow is like this: Train offline -> Make model available as a service -> Predict online.

- A classifier is trained offline with Fake and True news.
- The trained model is deployed as a service to serve users.

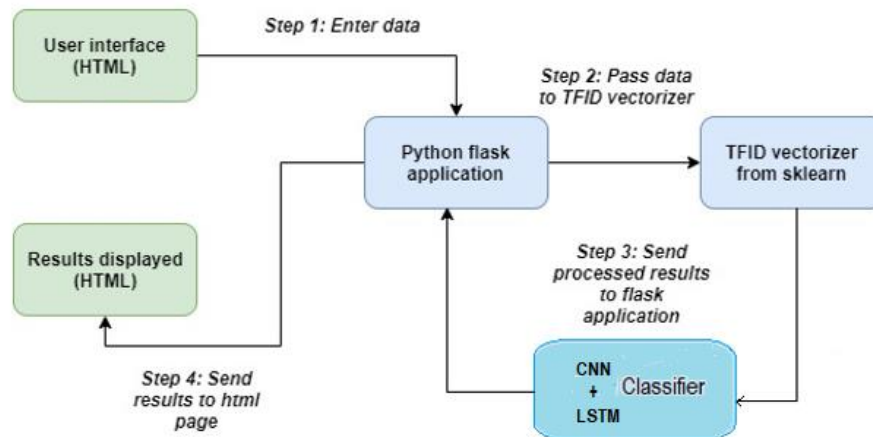
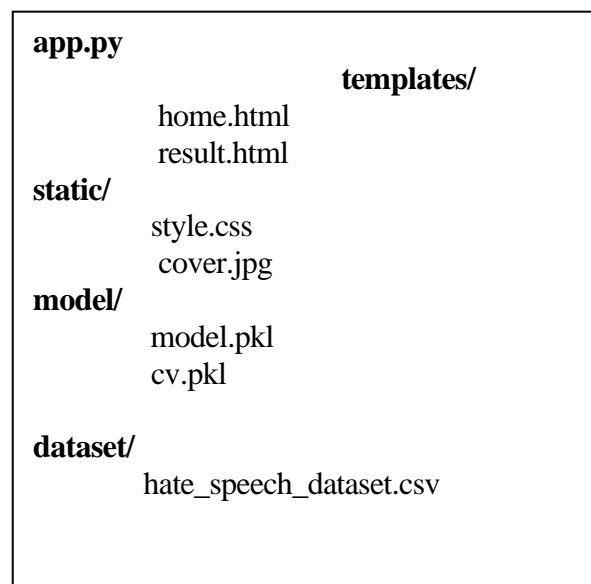


Figure 5. Application Workflow

5.1. Turning Model into a Web Application

First, we use the twitter hate speech dataset to build a prediction model that will accurately classify hate speech. We build the CNN with Model. After training the model, it is desirable to have a way to persist the model for future use without having to retrain. To persist we make a pickle file of our model.

We develop a web application that consists of a simple web page with a form field that lets us enter a message. After submitting the message to the web application, it will render it on a new page which gives us the result of hate speech/free speech.



The sub-directory templates are the directory in which Flask will look for static HTML files for rendering in the web browser, in our case, we have two HTML files: home.html and result.html.

5.2. Running Procedure

Once we have done all of the above, we can start running the API by either double click `app.py`, or executing the command from the Terminal:

```
/Library/Frameworks/Python.framework/Versions/3.12/bin/python3 "/Users/pranav13b/Documents/Data Glacier Internship/Week 12/Application/app.py"
○ pranav13b@Pranavs-MBP Week 12 % /Library/Frameworks/Python.framework/Versions/3.12/bin/python3 "/Users/pranav13b/Documents/Data Glacier Internship/Week 12/Application/app.py"
* Serving Flask app 'app'
* Debug mode: on
WARNING: This is a development server. Do not use it in a production deployment. Use a production WSGI server instead.
* Running on http://127.0.0.1:5000
Press CTRL+C to quit
* Restarting with stat
* Debugger is active!
* Debugger PIN: 124-107-806
```

Figure 6: Command Execution

Now we could open a web browser and navigate to <http://127.0.0.1:5000/>, we should see a simple website with the content like so

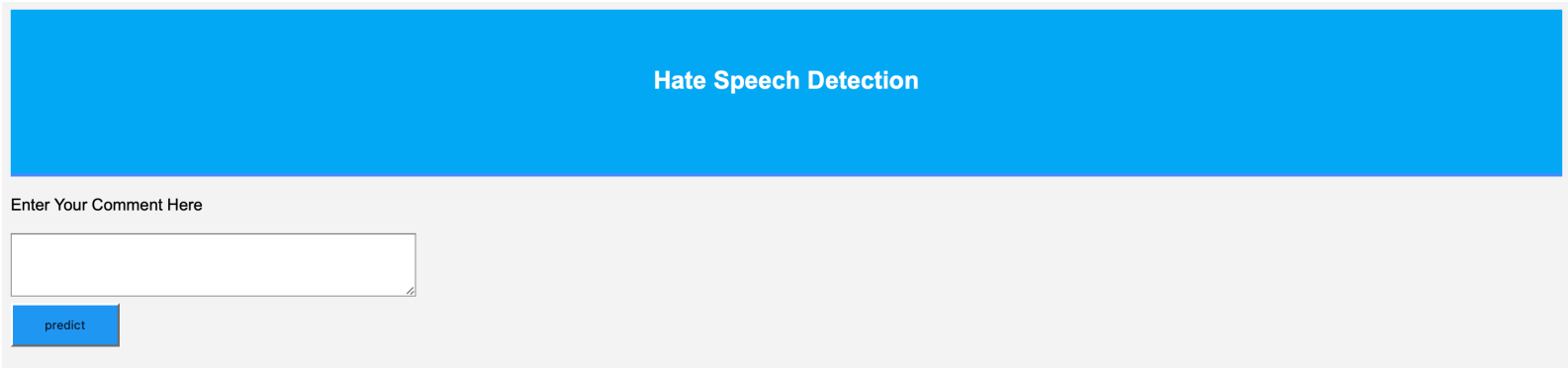
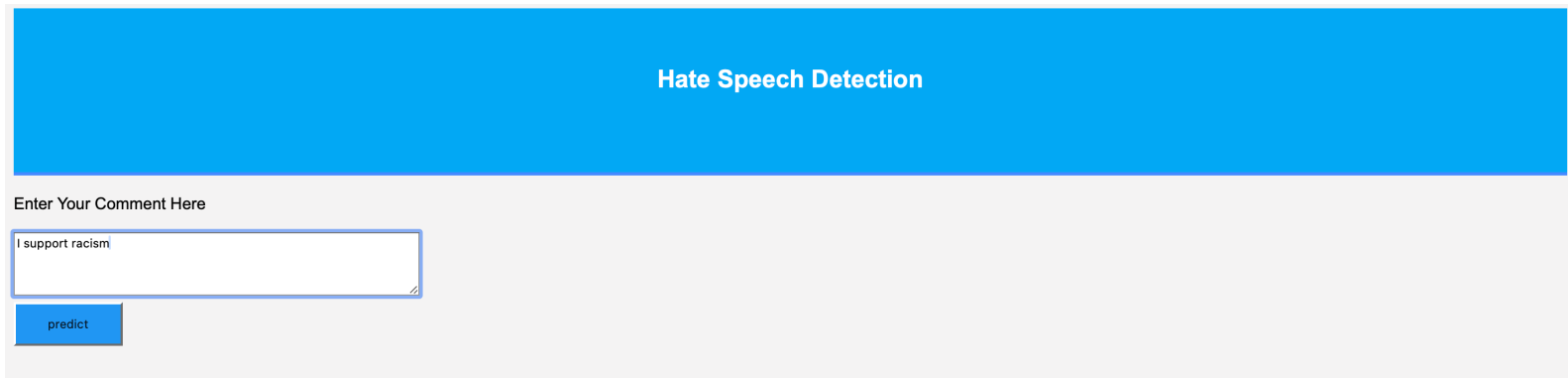
The screenshot shows a web application titled "Hate Speech Detection" in a blue header. Below the header, there is a light gray section with the text "Enter Your Comment Here" above a white text input field. To the right of the input field is a blue button with the word "predict" in white text.

Figure 6: Twitter Hate Speech Website

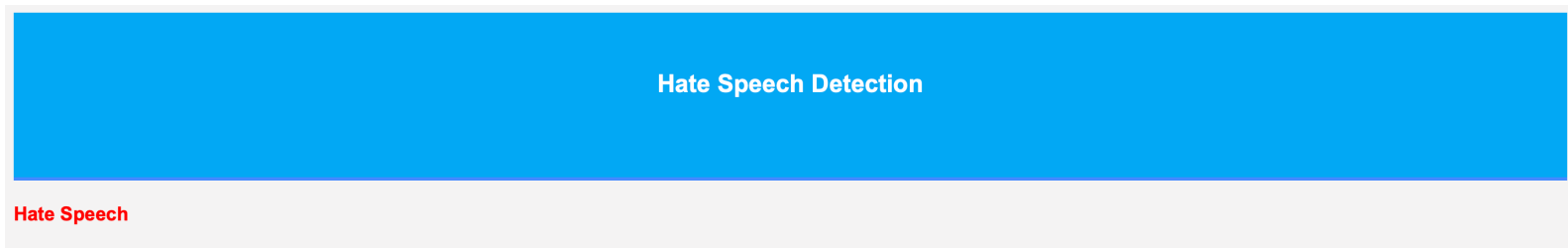
Now we enter input in the comments form



The screenshot shows a web interface for 'Hate Speech Detection'. It features a blue header with the title. Below the header is a light gray area containing a text input field with the placeholder 'Enter Your Comment Here'. The input field contains the text 'I support racism'. Below the input field is a blue button labeled 'predict'.

Figure 7: Input In The Comments Form

After entering the input click the predict button now we can the result of our input.



The screenshot shows the same 'Hate Speech Detection' interface as Figure 7, but with the result displayed. The blue header is still present. Below the header, the text 'Hate Speech' is displayed in red, indicating the model's classification for the input 'I support racism'.

Figure 8: Hate Speech Example

7. Conclusion

The goal of this project was to find capable methods and settings that could be used to help the detection of Hate and Free Speech of twitter. The error rate of the model is not zero, so still, some incorrect can be classified as true by the model. In future we will enhance this work by implementing Temporal Convolutional Network (TCN) and Random Multimodel Deep Learning (RMDL) Techniques.

Reference

- [1] https://www.kaggle.com/datasets/vkrahul/twitter-hate-speech?select=train_E6oV3lV.csv
- [2] <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>
- [3] Donahue, J., Hendricks, L. A., Rohrbach, M., Venugopalan, S., Guadarrama, S., Saenko, K., & Darrell, T. (2016, May 31). Long-term recurrent convolutional networks for visual recognition and description. arXiv.org. Retrieved May 8, 2022, <https://arxiv.org/abs/1411.4389>