

Hate Speech Detection using LIME guided Ensemble Method and DistilBERT

N Deepakindresh¹, AviReddy Rohan¹, Aakash Ambalavanan¹ and B. Radhika Selvamani²

¹Vellore Institute of Technology, Chennai, India

²Center for Advanced Data Science, Vellore Institute of Technology, Chennai, India

Abstract

Hate Speech classification has crucial applications in the social media domain. We describe the performance of our classifiers in the Hate Speech and Offensive Content Identification Track (HASOC) of FIRE 2021 conference. The dataset provided is for Indo-European Languages. We chose English tweets and developed two main classifiers as part of HASOC Track 1, which had two Subtasks 1A and 1B. Subtask 1A is a binary Hate Speech identification task, and Subtask 1B is multi grained classification of hate, profane, offensive and neutral content. Our team "Beware Haters" studied Support Vector Machine, Random Forest, Logistic Regression, Bidirectional Long Short Term Memory Model and an Ensemble of the listed models for the Subtask 1A and the highest Macro F1 score we achieved was 0.7722 by our Ensemble model which combined the advantages of SVM, Logistic Regression and Random Forest. We used a model interpretation tool LIME, before integrating the models in a weighted Ensemble approach. For Subtask 1B, we obtained better results using a DistilBERT model that achieved a Macro F1 score of 0.6311. We have compared the performance of the basic DistilBERT Model with a fine tuned version.

Keywords

Hate Speech Identification, TF-IDF, Ensemble, LSTM, SVM, Random Forest, Logistic Regression, LIME, BERT, DistilBERT

1. Introduction

The initial research in hate speech dates back to 1993 and has been legally defined by John T. Nockleby to describe any communication that incites hatred against anyone or any group of people in the name of race, ethnicity, religion, sexual orientation etc [1]. Hate speech identification task caught up with the Machine learning community during the spurt of user created content in social media during the last decade. Huge efforts have been put forth in the task of hate speech identification on the internet by Facebook [2], Twitter and Youtube, to conform to the legal and social responsibilities posed on these sites through governmental policies. Meanwhile, the research community finds the task of hate speech identification to be challenging, due to the diversity of the hate speech statements and the skewed nature of the collected data in most websites. Hate speech also poses challenges concerning the language

Forum for Information Retrieval Evaluation, December 13-17, 2021, India

✉ deepakindresh.n2019@vitstudent.ac.in (N. Deepakindresh); avireddynvsrk.rohan2019@vitstudent.ac.in (A. Rohan); aakash.ambalavanan2019@vitstudent.ac.in (A. Ambalavanan); radhika.selvamani@vit.ac.in (B.R. Selvamani)



© 2021 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

CEUR Workshop Proceedings (CEUR-WS.org)

used and the context in which it originates [3]. These challenges have made the hate speech identification task an interesting topic to be studied in the light of new machine learning algorithms and approaches which have been made available by cloud based libraries. The first HASOC Track of workshops created annotated hate speech in Indo-European languages to enable continued research in this direction. A detailed explanation about the HASOC Track at FIRE 2019 conference and the datasets have been discussed by Mandl et al., 2019 [4]. This paper is a summary of the efforts put forth by our team *Beware Haters* in the HASOC 2021 Track 1 [5]. There were multiple tracks analysing Twitter tweets in different code-mixed languages. We participated in the English hate speech identification subtasks of Track 1 [6]. There were two different subtasks in Track 1. The Subtask 1A requires the participants to identify hate speech in the given tweets and Subtask 1B is about classifying the tweets into multiple classes such as Hate, Profane, Offensive, Neutral etc.

2. Dataset Analysis

The dataset that we chose for analysis consisted only of English tweets. We analysed two different Subtasks 1A and 1B of Track 1. Subtask 1A is a binary classification task on tweets belonging to two distinct categories namely HOF (Hate and Offensive) and NOT (Non Hate-Offensive). HOF consists of hateful, offensive and profane content whereas NOT represents neutral content. Subtask 1B is a fine-grained classification problem, with 4 distinct classes, namely hate, offensive, profane and neutral (usually represented as none). The size of the dataset including training and test data is limited to about 4000 tweets. To overcome the data limitation we collected additional data from other sources. A free and publicly available Twitter hate speech dataset from Kaggle¹ was chosen for data augmentation and solving the class imbalance problem in the dataset provided by HASOC. On manual scrutiny, we found substantial similarity between the HASOC and Kaggle datasets for Subtask 1A. This Kaggle dataset has approximately six thousand tweets with labels comparable to HOF and NOT. The HASOC training dataset contains 65% HOF tweets whereas the Kaggle dataset consists of only 45% HOF tweets [Fig.1]. The combined dataset has 40% of HASOC data and 60% of Kaggle data and was used for the Random Forest model and Ensemble model for Subtask 1A. The dataset provided by HASOC for fine-grained classification of English tweets (Subtask 1B of HASOC 2021 Track 1) has 18% hate, 16% offensive, and 31% profane tweets [Fig.2]. We did not use any additional data for Subtask 1B other than the dataset provided by the HASOC 2021 organizers.

3. Feature Engineering

We extensively analysed the tweets from the dataset to decide on the pre-processing step. Tweets have been either removed or have been transformed using pattern matching techniques to deem them fit to the classification models under consideration. We have filtered out non-informative features from the tweets like URLs, white spaces, usernames that start with @ and hashtags. Other features like *emojis* have been filtered out. The tweets have been decontracted. Words

¹<https://www.kaggle.com/vkrahul/twitter-hate-speech>

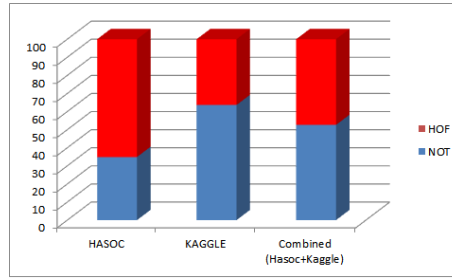


Figure 1: Distribution of HOF and NOT classes within the HASOC and the Kaggle datasets used for the binary classification task of HASOC2021 Subtask 1A.

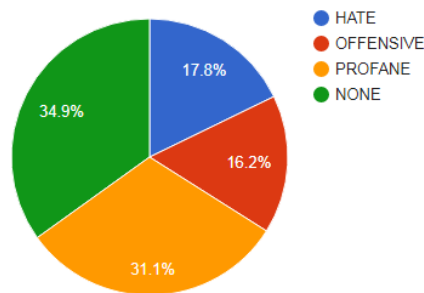


Figure 2: The distribution of classes within the dataset for the multilabel classification Subtask 1B of HASOC 2021.

like *won't*, *don't*, *can't*, *he'll*, *I'll* etc have been converted to their complete forms. Stop words are those that appear very frequently in the tweets but don't help in conveying any meaning. Common stop words like *the*, *not*, *is* and *was* have been removed. In addition, we have performed tokenization and lemmatization of the preprocessed tweets. Table 1 shows some examples of tweets before and after preprocessing.

Table 1

Tweets before and after Preprocessing

Before preprocessing	After preprocessing
@krtoprak_yigit Soldier of Japan Who has dick head	soldier japan dick head
@blueheartedly You'd be better off asking who DOESN'T think he's a sleazy shitbag lmao.	would better ask think sleazy shitbag lmao
@wealth if you made it through this && were not only able to start making money for yourself but sustain living that way all from home, fuck these companies & corporate pigs. power to the people, always. Technically that's still turning back the clock, dick head https://t.co/jbKaPJmpt1	make not able start make money sustain live way home fuck company corporate pig power people always technically still turn back clock dick head

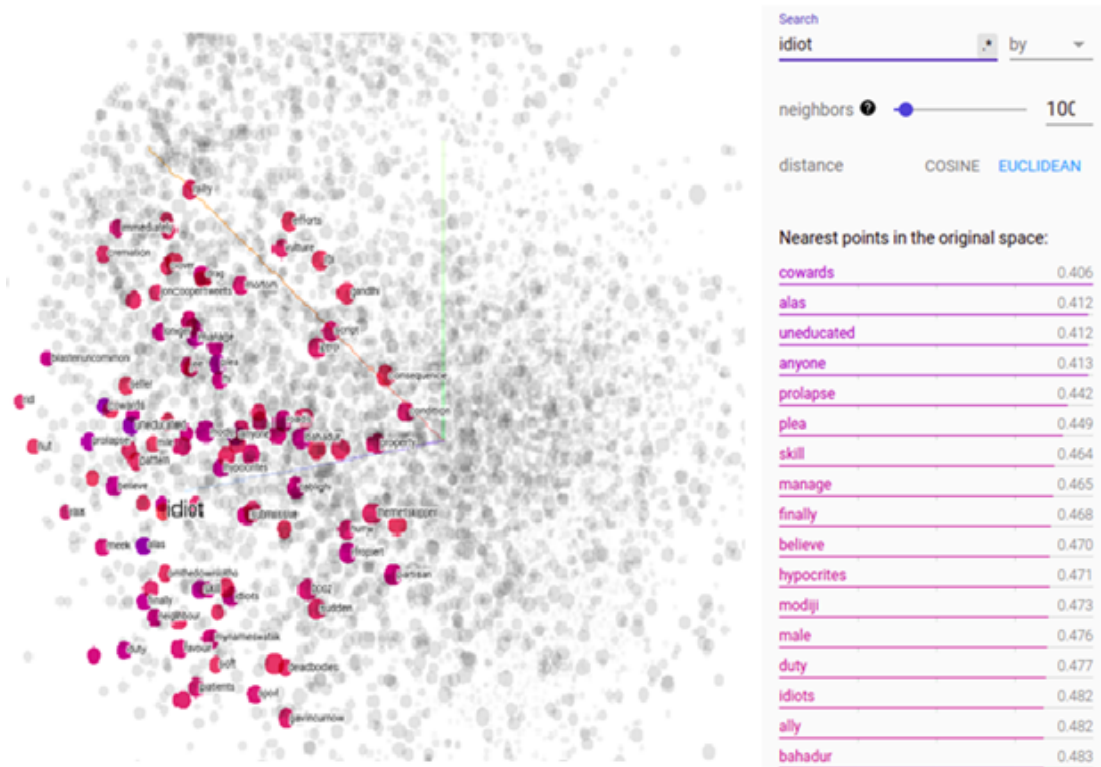


Figure 4: A 3D visualization of the Word2Vec embedding using Principle Component Analysis techniques for dimension reduction.

4. The Hate Speech Identification Task

Hate speech identification has been perceived as a binary classification problem to determine whether a twitter content is hateful-offensive or not. We compared the performance of various binary text classifiers on the training data. We used an n-gram based tf-idf vector embedding to prepare the training data used for learning the models. Logistic Regression [8] is a well-known simple regression model that serves as a basic model for binary classification. Random Forest is a widely used meta estimator that fits a number of decision tree classifiers on various sub-samples of the dataset and uses averaging to improve the predictive accuracy and control over-fitting [9]. The number of estimators and the sub-sample size for each estimator are some of the parameters used to fine tune the approach. We used 2000 estimators and a minimum sample-size of 2 per estimator as per the default settings of sklearn's Random Forest Classifier. The Support Vector Machine [10] is a state-of-the-art machine learning model with proven performance in countless machine learning applications with sparse high dimensional data. It uses different kernels, namely Linear, polynomial, Radial basis function, and sigmoid to transform the data to a lower dimension, which enables application of maximum margin classifier for obtaining the decision plane.

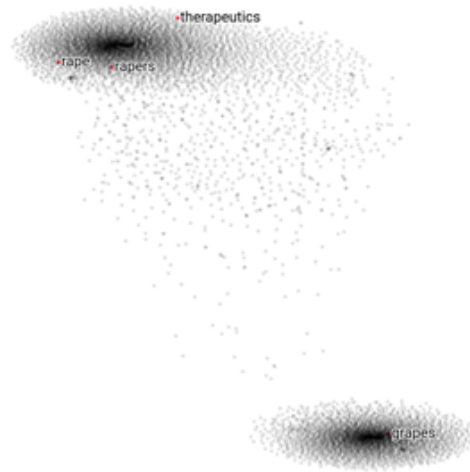


Figure 5: T-SNE word embedding for the rhyming words *rape* and *grape*.

4.1. Sequence Classifiers

In addition to the simple classifiers listed above, we have explored sequence classifiers using Long Short-Term Memory Models (LSTM) [11]. LSTM is a kind of Recurrent Neural Network (RNN) model, which has the added benefit of encoding contextual meaning among the words for a longer time step. Bidirectional LSTMs overcome the directional bias associated with traditional LSTM. Bidirectional LSTMs train two instead of one LSTMs on the input sequence. The first on the input sequence as-is and the second on a reversed copy of the input sequence. This can provide additional context to the network and result in faster and even fuller learning of the concepts. The same training data used for the simple classifiers has been used to train the LSTM models. A Bidirectional LSTM layer with 64 unit follows an input layer of 32 units. The output is then connected to two dense layers activated by two types of functions: relu and sigmoid (for binary prediction) with 64 and 1 unit respectively. In this pipeline we used binary cross-entropy as the loss function and the Adam optimizer [12] to optimize the model parameters. We identified the right stopping epoch by analyzing the validation accuracy and thus finalizing the model parameters.

4.2. Ensemble Methods

Ensemble learning is a process where multiple diverse models are integrated in a way to obtain better predictive performance than what could be obtained by each constituent model independently. We created an Ensemble classifier of Random Forest, Logistic Regression and SVM with the soft voting method [13]. The models used for Ensemble include a Random Forest classifier with 2000 estimators, a Logistic Regression model and a Support Vector Machine using a radial bias function. We used a Local Agnostic Model Interpretation approach to understand the performance of the models before building the Ensemble.

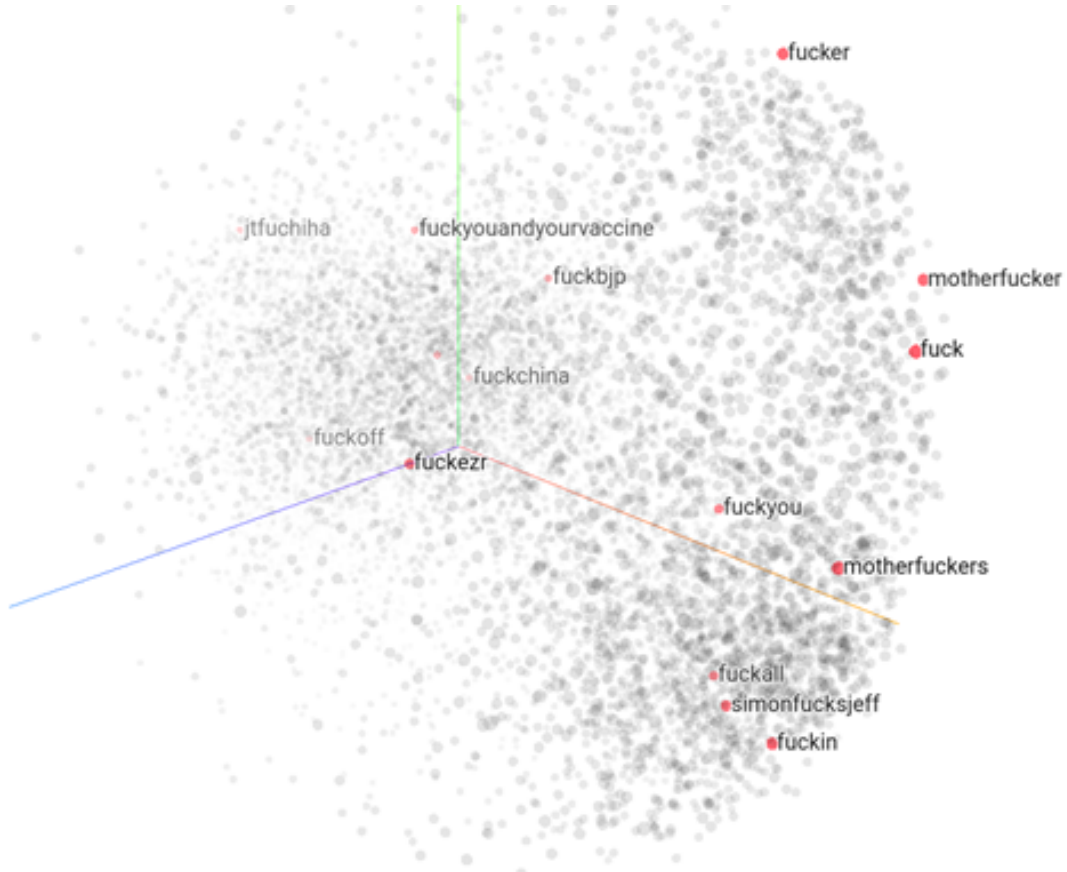


Figure 6: Visualization using Principle Component Analysis for the Word *fuck*.

4.3. Lime Guided Ensemble Approach

To have a better understanding of the models, we decided to use model explanation strategies. LIME is a local agnostic model interpreter [14]. The advantage of using LIME is that it provides uniform explanations across different models since it is model agnostic. It has been recently proposed by Sangani et.al.,2021 [15] to be used to compare the high performing models in the Kaggle platform. LIME supports explanations on both regression and classification models. We used LIME to compare the predictions of our models. The explanations are provided for each model by means of highlighted text and a relevance bar chart [Fig.9]. The words highlighted in orange denote a hate content, whereas those highlighted in blue support neutral decision. They also provide a score for the overall final decision of the interpreter summarizing over the highlighted words. We used LIME to analyse the shortcomings of each model and fine-tuned the weights accordingly for the Ensemble technique. The best performance by the Ensemble models was achieved by assigning a weight of 1,2 and 1 to the Logistic Regression, Random Forest, and



Figure 7: T-SNE visualization for the word *ass*.



Figure 8: T-SNE visualization for the word *bjp*.

SVM respectively. We chose the soft voting method over other methods as it predicts the class label based on the argmax of the sums of the predicted probabilities, which is recommended for an Ensemble of well-calibrated classifiers [13]. The particular LIME implementation we used was time-consuming, hence we could only make a qualitative decision based on the manual analysis of the limited number of explanations obtained for each model. We also had issues trying to implement LIME in other models such as LSTM as a lot of processing was required before and after training and prediction. We couldn't fit LSTM into the LIME pipeline to give explanations for predictions.

The LIME explanations are provided in [Fig.9] and [Fig.10] for predictions made on a hateful content and a neutral content by Random Forest, Logistic Regression, SVM, and Ensemble model respectively. We used LIME for analysing 6 manually sampled instances from the dataset. From [Fig.9] we can infer that SVM has done better than Logistic Regression and Random Forest as it has assigned a higher weight to the word *covidvaccine* and other models have either made a mistake or have assigned lesser weights.

In [Fig.10] Random Forest has done best since it has the highest prediction score for hate speech along with SVM. But SVM misses relevant words. Random forest has correctly classified *mattancock* as a word of Hate Speech, unlike SVM which has mistaken it for Non-Hate Speech. The Ensemble model's prediction has also been explained and it is clearly visible that it has

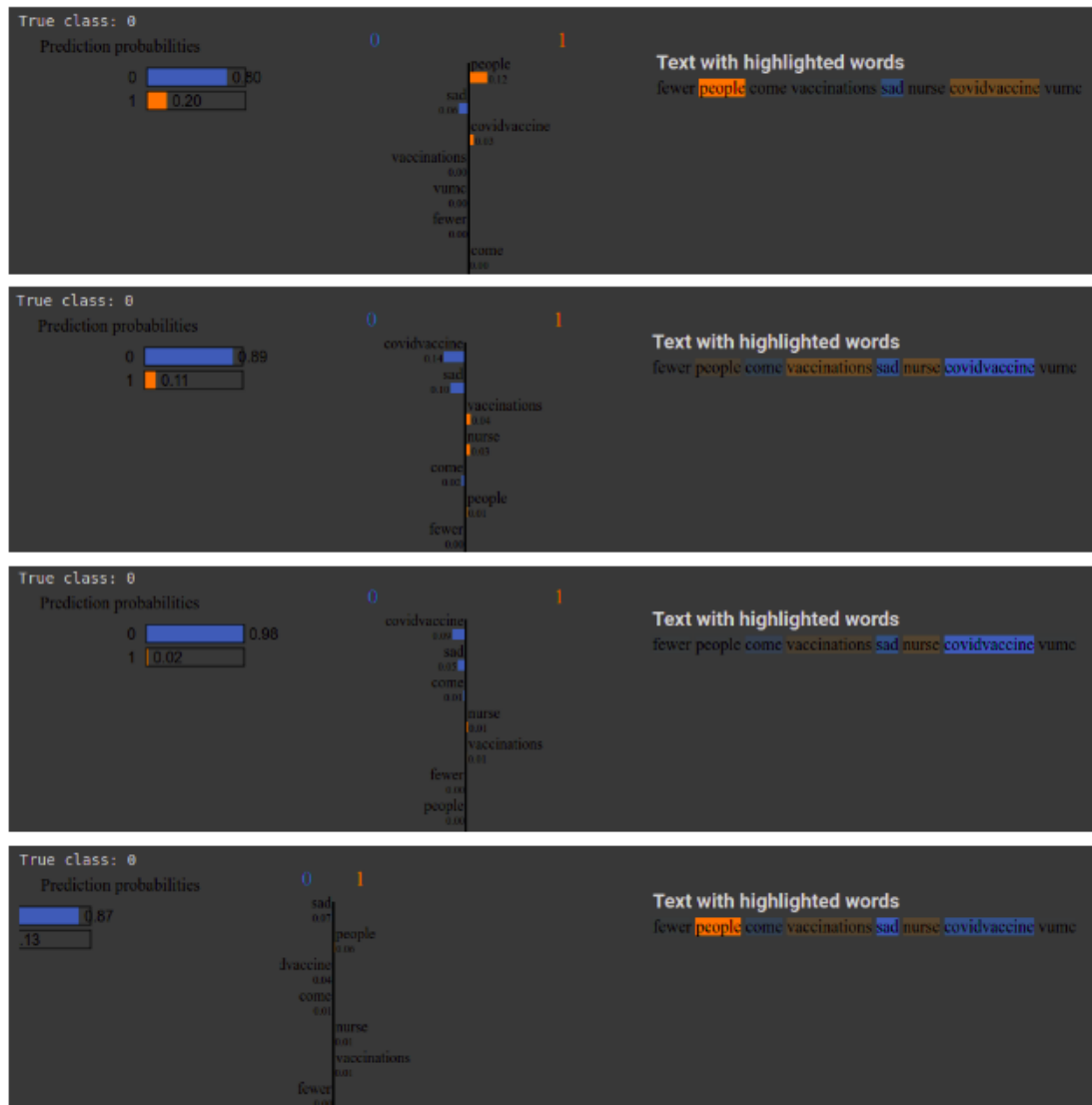


Figure 9: LIME Explanations for Non Hate Sentences for Random Forest, Logistic Regression, SVM and Ensemble model respectively

used a combined result of all three models for its prediction making it unbiased and being able to perform better.

5. Fine Grained Classification for Subtask 1B

For Subtask 1B which requires a multi-label classifier, we turned to DistilBERT [16]. Google provides pre-trained partial BERT [17] language models which may be fine tuned based on our



Figure 10: LIME Explanations for Hate Sentences for Random Forest, Logistic Regression, SVM and Ensemble model respectively

application. BERT is designed to pre-train deep bidirectional representations from the unlabeled text by jointly conditioning on both left and right context words. BERT uses the advantages of transfer learning, a method where a model developed for a task is reused as the starting point for a model on a second task. We have built DistilBERT, a model trained in a self-supervised fashion using the BERT base model as a teacher. It is smaller and faster than BERT. We have trained DistilBERT on the hate speech corpus [18]. The main reason for choosing this model, unlike other sequence classifiers like RNN is that, the DistilBERT model works on the entire sequence at once instead of reading the tokens sequentially. This process can be further accelerated using GPU support provided by Google Colab. We could adapt the pre-trained DistilBERT model by training it further on our relatively smaller dataset, further fine-tuning the parameters for better accuracy without much computational overhead. The tweets in the dataset were of different

lengths. We used padding to normalize all the tweets to have the maximum sequence length. The DistilBERT model we chose was pre-trained on unlabelled Wikipedia corpus. It was then fine-tuned for our corpus by adding the output layer.

The DistilBERT [18] model we chose facilitated in developing a multi label regression classifier. The output vector from the model is a 4-dimensional numeric vector. We used one-hot encoding on the class labels of the existing HASOC dataset to obtain a 4-dimensional boolean output vector for each tweet. The first value of the vector indicates hate class, the second profane, the third offensive and the fourth represents neutral. For eg., a tweet that is profane will be encoded [0 1 0 0]. While classifying a tweet using DistilBERT, on obtaining the 4 dimensional output vector, the tweet is assigned the class corresponding to the vector position which has the highest numerical value.

6. Results

6.1. Metrics For Comparison

The models were compared using precision, recall, F-Measure and accuracy.

- Precision: Precision is also known as the positive predicted value. It is the proportion of predictive positives which are actually positive (true positives).
- Recall: It is the proportion of actual positives which are predicted positive.
- F-Measure: It is the harmonic mean of precision and recall. The standard F-measure (F1) gives equal importance to precision and recall.
- Accuracy: It is the number of correctly classified instances (true positives and true negatives).

6.2. Comparing the Models for Subtask 1A

For Subtask 1A the combined training dataset from HASOC and Kaggle were used for training and the test dataset was randomly sampled from the HASOC training dataset. The split was 90:10 for training and testing. We have experimented with the HASOC dataset as well as the combined dataset. In both cases the test cases were sampled from the HASOC test dataset. The augmented kaggle dataset did play a significant role to improve the accuracy of Ensemble and Random Forest while it decreased the performance for the others such as Logistic Regression, Bidirectional LSTM and SVM. We have tabulated the best results obtained for the different models among the different training sets used.

The performance of Random Forest and Ensemble models trained on the combined HASOC and Kaggle datasets and the other models trained on the HASOC dataset have been provided in Table 2. The data has been ordered based on accuracy. From the plots 11 and 12 we can infer that the Ensemble model performs much better than the independent models as it combines the advantages of all the other models and reduces the ill effects of each model.

We have trained the LSTM model for 100 epochs, at which it gives the best performance. LSTM model does not fare well compared to other models with a Macro F1 of just 0.7198. Yet it provides valuable insights about the hate words through embedding vectors that we projected

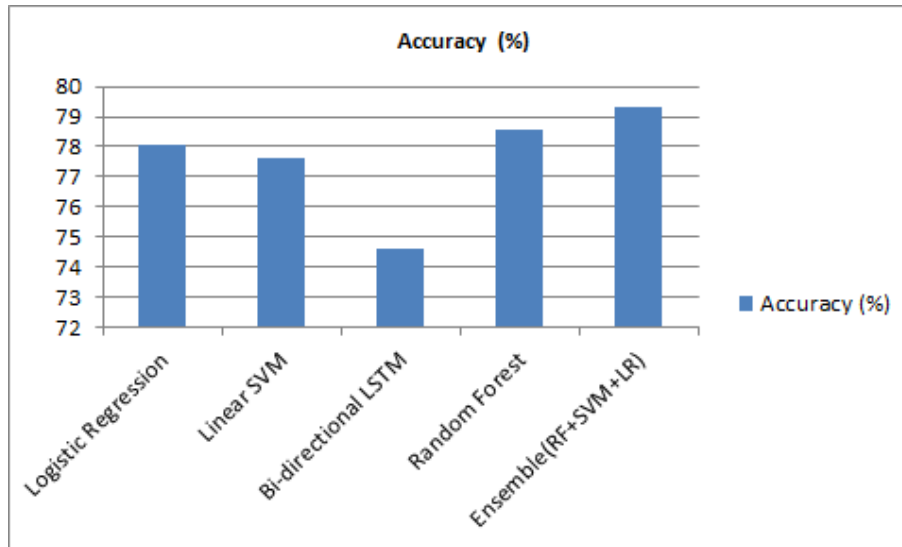


Figure 11: Accuracy of the models built for Task 1A

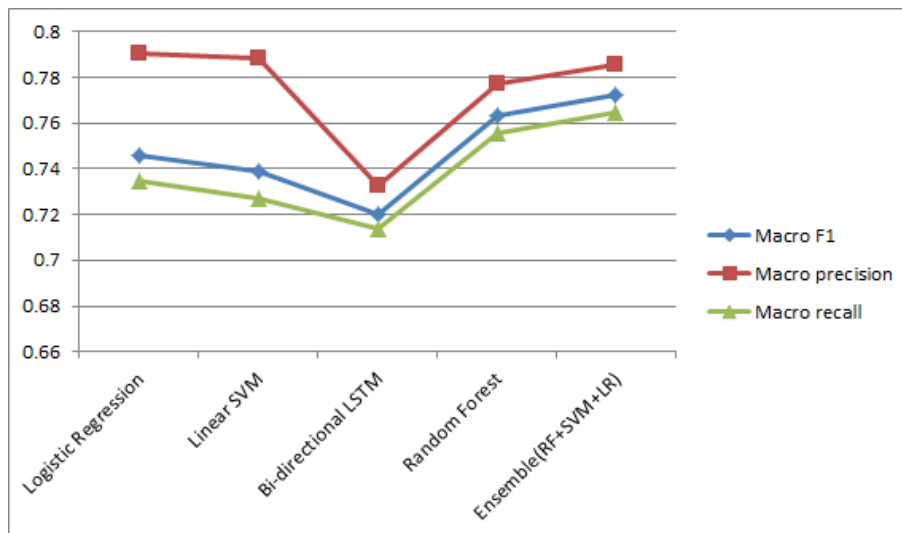


Figure 12: Macro F1, Macro Precision and Macro Recall of all the models built for Task1A

using embedding projectors explained in section 3. With more data the model could have performed better by learning improved hate word embeddings and understanding the context better.

Although from the given LIME examples the significance of the Logistic Regression model in the Ensemble might not be visible, we could get the best Ensemble performance only when Logistic Regression was included.

6.3. Comparison between Fine tuned and Untuned DistilBERT Models for Subtask 1B

For Subtask 1B we tested two versions of DistilBERT model on a 90:10 training and test data split of the HASOC dataset. The untuned initial version has 4 hidden layers each with 128, 64, 32 and 4 units. We removed the hidden layer with 128 units while fine tuning the DistilBERT model. To prevent underfitting we reduced dropout from 0.1 to 0.05. We trained the fine tuned DistilBERT model for 6 epochs, unlike the untuned version which had only 5 epochs. The major changes that led to a colossal increase in accuracy were basically with respect to simplifying the hidden layer, reducing the drop out and increasing the training epochs. The results are found in Table 3.

Table 2

Submissions for Subtask 1A

Classifier	Macro F1	Macro prec	Macro recall	Accuracy (%)
Logistic Regression	0.7462	0.7902	0.7344	78.064
Linear SVM	0.7386	0.7882	0.7270	77.596
Bi-directional LSTM	0.7198	0.7326	0.7138	74.629
Random Forest	0.7631	0.7775	0.7558	78.532
Ensemble(RF+SVM+LR)	0.7722	0.7859	0.7649	79.313

Table 3

Submissions for Subtask 1B

Classifier	Macro F1	Macro Prec	Macro Recall	Accuracy (%)
Untuned DistilBERT	0.6106	0.6223	0.6098	66.667
Finetuned DistilBERT	0.6311	0.6364	0.6303	67.681

7. Conclusion

We classified the Hate Speech dataset provided by the HASOC Track 1 of FIRE 2021. Though Subtask 1A is a simple binary classifier we could not achieve the expected accuracy and stood 25th out of the total 56 submissions [Fig.13]. The reason is that Subtask 1A had a limited data and our decision to augment it with more data from Kaggle probably did not work. Meanwhile, we observed that an Ensemble method improved our accuracy over other approaches used. Work is under progress in this line to quantify the interpretations provided by LIME. The multigrain classification problem was the toughest with the benchmark performance as low as 0.5 (Precision and F1 score). With respect to the fine-grained classifier developed for Subtask 1B, DistilBERT gave superior performance compared to other approaches and scored the 8th rank among the total 37 submissions [Fig.14]. It is interesting to note that both Subtask 1A and Subtask 1B required entirely different approaches though they were all Hate Speech classifiers.

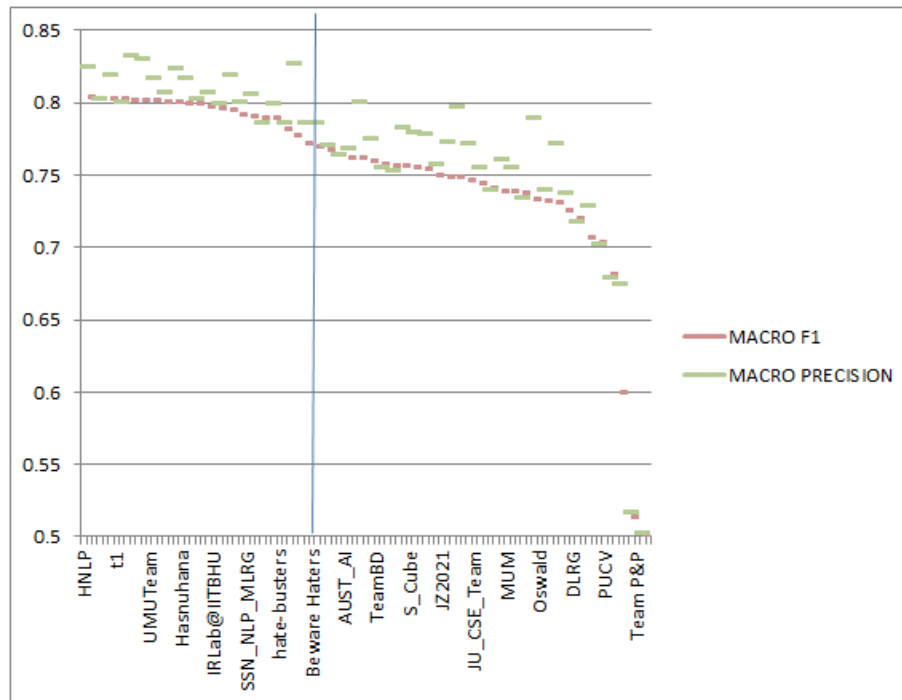


Figure 13: Performance of the Ensemble Model based Hate Speech Classifier (Beware Haters) in HASOC 2021 Track 1 Subtask A

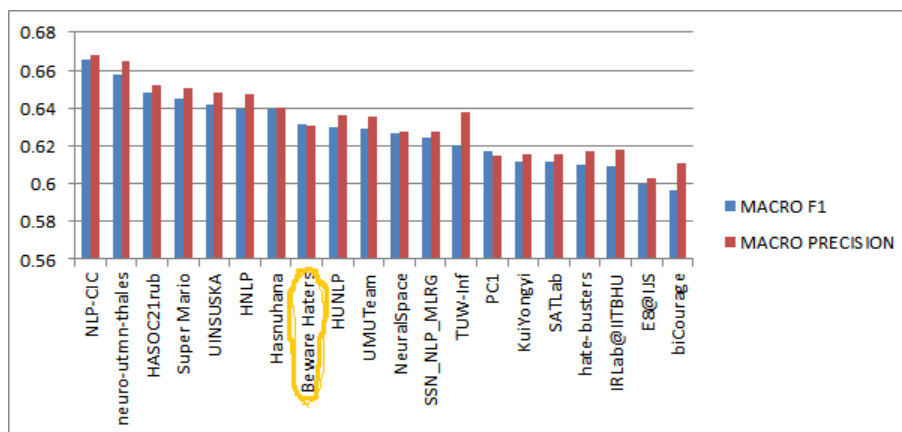


Figure 14: Performance of the DistilBERT Model based Hate Speech Classifier (Beware Haters) in HASOC 2021 Track 1 Subtask B

References

- [1] J. Nockleyby, Hate speech in encyclopedia of the american constitution, Electronic Journal of Academic and Special librarianship (2000).
- [2] T. Guardian, Zuckerberg on refugee crisis: 'hate speech has no place on

facebook', 2016. URL: <https://www.theguardian.com/technology/2016/feb/26/mark-zuckerberg-hate-speech-germany-facebook-refugee-crisis>.

- [3] Z. Zhang, L. Luo, Hate speech detection: A solved problem? the challenging case of long tail on twitter, *Semantic Web* 10 (2019) 925–945.
- [4] T. Mandl, S. Modha, P. Majumder, D. Patel, M. Dave, C. Mandalia, A. Patel, Overview of the HASOC track at FIRE 2019: Hate speech and offensive content identification in indo-european languages, *Proceedings of the 11th Forum for Information Retrieval Evaluation* (2019).
- [5] S. Modha, T. Mandl, G. K. Shahi, H. Madhu, S. Satapara, T. Ranasinghe, M. Zampieri, Overview of the HASOC Subtrack at FIRE 2021: Hate Speech and Offensive Content Identification in English and Indo-Aryan Languages and Conversational Hate Speech, in: *FIRE 2021: Forum for Information Retrieval Evaluation*, Virtual Event, 13th-17th December 2021, ACM, 2021.
- [6] T. Mandl, S. Modha, G. K. Shahi, H. Madhu, S. Satapara, P. Majumder, J. Schäfer, T. Ranasinghe, M. Zampieri, D. Nandini, A. K. Jaiswal, Overview of the HASOC subtrack at FIRE 2021: Hate Speech and Offensive Content Identification in English and Indo-Aryan Languages, in: *Working Notes of FIRE 2021 - Forum for Information Retrieval Evaluation*, CEUR, 2021. URL: <http://ceur-ws.org/>.
- [7] J. E. Ramos, Using tf-idf to determine word relevance in document queries, 2003.
- [8] O. Oriola, E. Kotzé, Evaluating machine learning techniques for detecting offensive and hate speech in south african tweets, *IEEE Access* 8 (2020) 21496–21509.
- [9] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, E. Duchesnay, Scikit-learn: Machine learning in Python, *Journal of Machine Learning Research* 12 (2011) 2825–2830.
- [10] D. Robinson, Z. Zhang, J. Tepper, Hate speech detection on twitter: Feature engineering vs feature selection, in: *European Semantic Web Conference*, Springer, 2018, pp. 46–49.
- [11] M. Sundermeyer, R. Schlüter, H. Ney, Lstm neural networks for language modeling, in: *Thirteenth annual conference of the international speech communication association*, 2012.
- [12] D. Kingma, J. Ba, Adam: A method for stochastic optimization, *International Conference on Learning Representations* (2014).
- [13] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, E. Duchesnay, Scikit-learn: Machine learning in Python, *Journal of Machine Learning Research* 12 (2011) 2825–2830.
- [14] M. T. Ribeiro, S. Singh, C. Guestrin, "why should i trust you?": Explaining the predictions of any classifier, 2016. [arXiv:1602.04938](https://arxiv.org/abs/1602.04938).
- [15] R. B. Sangani, A. Shukla, R. B. Selvamani, Comparing deep sentiment models using quantified local explanations, in: *Accepted for publication in Proceedings of IEEE-Smart Technologies, Communication Robotics 2021 Conference*, IEEE, 2021.
- [16] V. Sanh, L. Debut, J. Chaumond, T. Wolf, Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter, *ArXiv abs/1910.01108* (2019).
- [17] S. Yu, J. Su, D. Luo, Improving bert-based text classification with auxiliary sentence

and domain knowledge, IEEE Access 7 (2019) 176600–176612. doi:10.1109/ACCESS.2019.2953990.

- [18] R. Mutanga, N. Naicker, O. O. Olugbara, Hate speech detection in twitter using transformer methods, International Journal of Advanced Computer Science and Applications 11 (2020).