

Multi Platform-based Hate Speech Detection

Shane Cooke¹, Damien Graux²^a and Soumyabrata Dev¹^b

¹ADAPT SFI Research Centre, School of Computer Science, University College Dublin, Ireland

²ADAPT SFI Research Centre, Trinity College Dublin, Ireland

shane.cooke@ucdconnect.ie, damien.graux@adaptcentre.ie, soumyabrata.dev@ucd.ie

Keywords: Hate Speech Detection, Multi-Platform, Combining Embeddings and Classifiers.

Abstract: A major issue faced by social media platforms today is the detection, and handling of hateful speech. The intricacies and imperfections of online communication make this a difficult task, and the rapidly changing use of both non-hateful, and hateful language in the online sphere means that researchers must constantly update and modify their hate speech detection methodologies. In this study, we propose an accurate and versatile multi-platform model for the detection of hate speech, using first-hand data scraped from some of the most popular social media platforms, that we share to the community. We explore and optimise 50 different model approaches, and evaluate their performances using several evaluation metrics. Overall, we successfully build a hate speech detection model, pairing the USE word embeddings with the SVC machine learning classifier, to obtain an average accuracy of 95.65% and achieved a maximum accuracy of 96.89%. We also develop and share an application allowing users to test sentences against a collection of the most accurate hate speech detection models. Our application then returns a aggregated hate speech classification, together with a confidence level, and a breakdown of the methodologies used to produce the final classification for explainability.


1 Introduction


The definition of hate speech is a topic of great discussion within society today. Philosophers, researchers and law-makers all have their own variations of the definition, however there are a set of facts upon which most parties agree on, the first being that the message is directed at an individual or group, and the second that based on that message the group is viewed as negative, unwelcome or undesirable which warrants hostility towards them (Rudnicki and Steiger, 2020). In EU law, hate speech is defined as “the public incitement to violence or hatred on the basis of certain characteristics, including race, colour, religion, descent and national or ethnic origin” (Jourová, 2016).

Online hate speech however is a special case of hate speech that occurs in the online environment, making the perpetrators more anonymous, which in turn may make them feel less accountable, and as a result potentially more ruthless. To effectively fight online hate speech, non-government organisations aim to be more flexible than the justice system allow, in particular it is increasingly common to define hate speech much more broadly and include messages that

do not explicitly incite violence only, but instead spread prejudice, stereotypes, biases and a general sense of ostracism. Hate speech online has been a major problem since the spread of the Web, however in light of the rapid rise in popularity of social media sites, this problem has since increased in size exponentially. For instance, (Hawdon et al., 2015) found that approximately 53% of American, 48% of Finnish, 39% of British and 31% of German survey-respondents had been exposed to online hate material. A study conducted by an AI-technology company, found that Twitter hate speech against China and the Chinese had increased 900% in early 2020, and found that a 70% increase in hate between kids and teens in online chatrooms had occurred in the same timeframe (L1GHT, 2020). TikTok removed 380 000 videos in August 2020 alone (Han, 2020) and Facebook reported a record 25 million instances of hate speech in Q1 of 2021.

In this study, we propose an accurate and versatile multi-platform model for the detection of hate speech, using first-hand data scraped from some of the most popular social media platforms. Our contributions are threefold: First, we annotated manually a corpus of 3 000 comments from three social media

^a <https://orcid.org/0000-0003-3392-3162>

^b <https://orcid.org/0000-0002-0153-1095>

platforms and share it to the community¹. Second, we explore and optimise 50 different model approaches, and evaluate their performances using several evaluation metrics. Third, in addition, we also develop and share² an application allowing users to test sentences against a collection of the most accurate hate speech detection models to give the possibility to have finer results made from a combination of several models.

The rest of the article is structured as follows: in Section 2, we briefly remind the state-of-the-art in hate-speech detection. Then, we describe the data acquisition process in Section 3. Section 4 gives the details of the approach followed and Section 5 discusses the obtained results and performances. Section 6 presents the final application we developed. And finally, Sections 7 & 8 respectively mark the limitations of our method and draw our conclusions.

2 Related Work

Detecting hate speech online amongst millions of posts every day is a hard task and carries many associated challenges with it. (Kovács et al., 2020) outlined some of these challenges and reviewed over fifty works on hate detection online such as (Davidson et al., 2017) or (Bhattacharya and Weber, 2019). Some of the main challenges identified in the form of key-word based search approaches. Another huge challenge in hate speech detection which spans all forms of search is the detection of context and nuance. (Röttger et al., 2021), found that many hate speech detection models struggled with “reclaimed slurs” and often mislabelled them as hateful. In parallel, (Sap et al., 2019) outlines bias consideration challenges.

(Salminen et al., 2020) presents a multi-platform machine learning approach to online hate detection is proposed. They observed that most studies (e.g. (Kansara and Shekokar, 2015; Ramampiaro, 2018; Lee and Lee, 2018)) tend to focus on one platform, which they saw as problematic because there are no guarantees it generalizes well across platforms.

There are many ways to evaluate the performance of hate speech detection models. (Mozafari, 2020) remark that “classifiers with higher precision and recall scores are preferred in classification tasks. However, due to the imbalanced classes in the hate speech detection datasets, we tend to make a trade-off between these two measures”. For this reason they decided to use macro averaged F1-measures to summarize the performance of their models. On the other hand, (Al-shalan and Al-Khalifa, 2020) decided to evaluate their

models using precision, recall, F1-score, accuracy, hate class recall, and AUROC. (Vigna et al., 2017) proposed that the best evaluation metrics to use are accuracy, precision, recall and F-score.

3 Data Acquisition

In order to create a versatile and well-rounded hate speech detection system, we decided to collect comment and post data from three different sources: Reddit, Twitter and 4Chan. The use of language, both hateful and non-hateful, can vary extremely between platforms, for this reason we believe the use of a multi-platform approach should help hate speech detection. Each of the three platforms boast varying levels and methods of moderation: with Reddit having community-based moderation, Twitter having automatic or employee-based moderation, and 4Chan having virtually zero moderation. Due to these highly differing methods of platform moderation, it is easy to pinpoint the subtleties of the hateful language used on each platform. Due to Reddit’s community-based moderation, the hateful speech exhibited is often very subtle and very few slurs are used, while the automatic and employee-driven moderation used by Twitter promotes “leetspeak” and disguised slurs. These are both in sharp contrast to the language used in the unmoderated 4Chan forums, where extreme slurs are used regularly, and hateful speech is not only tolerated, but encouraged by some. The choice of three data sources was ultimately made to ensure that the hate speech dataset curated for this project would be heterogeneous, and would feature a wide array of different forms of both non-hateful and hateful language.

Overall, we decided that a procured dataset of 3 000 posts and comments would be the best solution. The dataset exhibits an equal split of 1 000 posts from each of the three social media platforms, and each post is classified and labelled as either non-hateful (‘0’) or hateful (‘1’). The posts and comments are split into classifications of 2 400 non-hateful posts (80%), and 600 hateful posts (20%).

4 General Approach

4.1 Word Embeddings

Word Embeddings are a class of techniques in which individual strings are mapped to vector or numerical representations. The chosen form of representation varies widely depending on the word embedding method being employed, however every method

¹Annotated corpus ☞

²Github repository ☞

follows the same core principle of mapping a single string to a single defined value. In order to efficiently and accurately analyse and model the posts contained in our database, we used a variety of **five** different word embedding methods which all employ very different embedding methodologies.

First, we used *TFIDF* (“Term Frequency-Inverse Document Frequency”). It is a machine learning algorithm based on a statistical measure of finding the relevance of words in a text. The “TF”, is calculated by dividing the number of occurrences of words by the total number of words in the text base. The “IDF” is calculated by dividing the total number of comments by the number of comments containing the word. The overall embedding is equal to $(TF) \times (IDF)$. Second, we consider *Doc2Vec* (Le and Mikolov, 2014) which is an NLP tool for representing documents as a vector, and is a generalisation of the “Word2Vec” model. *Doc2Vec* vectorises words to their representative format, and includes a paragraph numerical representation tied to these word vectors. Third, we used a *Hashing Vectorizer* algorithm which converts a text into a matrix of token occurrences, where each token directly maps to a column position in a matrix where its size is predefined. The hash function used is *Murmurhash3*. Fourth, we exploited Google’s *Universal Sentence Encoder* (Cer et al., 2018). It captures the most informative features of a given sentence and discard noise. Finally, we included also *BERT* (Devlin et al., 2018) which uses a Transformer learning the contextual relations between words in a text.

4.2 Classifiers

We trial a diverse collection of **ten** machine learning classifiers. We ensure that within this group of classifiers there are both classical machine learning algorithms such as the Decision Tree classifier, and more modern, task-specific algorithms such as the XGBoost classifier. The selected classifiers are:

1. Random Forest Classifier (Pal, 2005);
2. Decision Tree (Safavian and Landgrebe, 1991);
3. Naive Bayes (Rish et al., 2001);
4. SVC (Vapnik, 1998);
5. AdaBoost (Freund and Schapire, 1997);
6. Gaussian Process (Gibbs, 1998);
7. K-Neighbours (Guo et al., 2003);
8. Multi-layer Perceptron (Hornik et al., 1989);
9. XGBoost (Chen and Guestrin, 2016);
10. Linear Discrimination (Izenman, 2013).

Once the word embedding methods are chosen and implemented, we test all possible combinations of

word embedding, machine learning classifier pairs. First, the classifiers are trained using the training data vectors produced by the word embedding process. Each one of these vectors has a corresponding “Hateful” value of either ‘0’ or ‘1’, which is the ‘target’ variable. We run each classifier twenty times with a new train and test data split for each iteration, and take an average of each of the evaluation metrics across the twenty iterations and achieved a set of final results.

4.3 Optimisation Strategies...

4.3.1 ...for Classifiers

In order to optimise these results, the parameters or configuration variables of each classifier had to be tested and refined in pursuit of the highest possible results. While some classifiers do not take parameters such as the ‘GaussianNB’, ‘GaussianProcess’ and ‘XGBoost’ classifiers, the other classifiers can take upwards of eight parameters³. We thus implement exhaustive searches over specified parameter values, and implements fitting and scoring methods to evaluate each combination of parameters, see Figure 1.

Results The parameter optimisation of the machine learning classifiers had a majorly positive effect on the results produced by the classifiers across all evaluation metrics. While some algorithms do not take parameters such as ‘GaussianNB’, ‘GaussianProcess’ and ‘XGBoost’, the majority of algorithms do take parameter variables, and the overall optimisation process was extremely effective. For instance, Figure 2 shows an example of the results produced by the parameter optimisation for the SVC classifier. The highest performing kernel parameter (‘rbf’) has more than 1% higher accuracy than the lowest one (‘poly’).

4.3.2 ...for Word Embeddings

Similarly, we optimise the parameters of the word embedding methods. While some of the word embedding methods such as USE and BERT do not take parameters due to the fact that they are pre-trained algorithms, the *TFIDF*, *Doc2Vec* and *Hashing Vectorizer* methods take parameters. In pursuit of the most efficient word embedding parameter optimisation process possible, we created multiple different word embedding instances from the same word embedding method, each initialised with different parameters.

Results The parameter optimisation of the word embedding methods was also successful, and had a

³We used *GridSearchCV* from the *sci-kit learn* library.

```

randomForest_grid = {'n_estimators':[200,400,600,800,1000], 'criterion':['gini', 'entropy']}
decisionTree_grid = {'max_depth':[2,4,6,8], 'splitter':['best', 'random'], 'criterion':['gini', 'entropy']}
svc_grid = {'kernel':['linear', 'poly', 'rbf', 'sigmoid']}
adaboost_grid = {'n_estimators':[50,100,150,200], 'algorithm':['SAMME', 'SAMME.R']}
mlp_grid = {'max_iter':[500,1000,1500], 'activation':['identity', 'logistic', 'tanh', 'relu']}
linearDis_grid = {'solver':['svd', 'lsqr', 'eigen']}

```

Figure 1: Parameter dictionaries for the GridSearchCV algorithm.

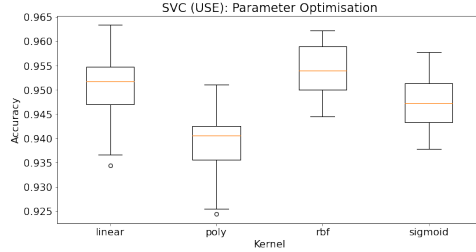


Figure 2: Parameter Optim. of the SVC classifier.

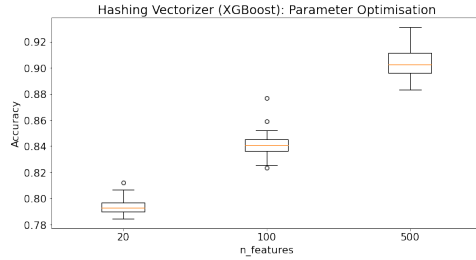


Figure 3: Parameter Optim. of the HV embeddings.

much more pronounced and noticeable effect on the results produced by the models, across all evaluation metrics. While the pre-trained word embedding algorithms such as ‘USE’ and ‘BERT’ do not take parameters, the other ones do, and the overall optimisation process was extremely effective. For example, Figure 3 shows some examples of the results achieved by the parameter optimisation for the Hashing Vectorizer. When 500 ‘n_features’ are selected as opposed to 20, there is more than a 10% increase in the accuracy produced by the word embedding method.

Unlike the parameter optimisation of the machine learning classifiers where changes in accuracy were often subtle, the word embedding parameter optimisation exhibited major improvements to the accuracy of the models. Increases in accuracy due to word embedding parameter optimisations were non-uniform and varied widely, however an increase in the range of +0.25% (TFIDF) and +10.5% (Hashing Vectorizer) was exhibited across all embeddings.

5 Experimental Validation

To evaluate and analyse the performance of the hate detection models, we relied on four main evaluation metrics: accuracy, precision, recall and F1-Score.

The proportion of positive identifications that were actually correct (Precision) and the proportion of actual positives that were identified correctly (Recall) are major factors in determining the overall performance of a hate speech detection system, and F1-Score (a harmonic mean of both precision and recall) is also an extremely valuable metric when judging overall performance. Accuracy is used to determine the ability of the model to accurately identify patterns or relationships between the data in a dataset based on the training data that it has received.

5.1 Single Platform

The overall goal is to produce an efficient and effective multi-platform hate speech detection model, however it is important to analyse and evaluate the performance of each of the three individual single-platform models. To carry out this evaluation, we first split the full 3000 comment dataset into three datasets of 1000 comments, with each dataset only containing data from one specific platform. This resulted in each platform having its own dataset with 800 (80%) of comments being non-hateful, and 200 (20%) of those being hateful. Each dataset was then initialised used the same training to testing data split ratio of 0.3, and tested using the exact same methodology, classifiers and word embeddings. Each platforms data was then used to create and test fifty models (number of word embeddings times the number of classifiers). Once this testing had been completed and all evaluation metrics had been noted, we singled out the top performing machine learning classifiers for each of the five word embedding methods, for each of the three platform datasets. The results of this process are shown in Figures 4, 5 & 6. We notice that there is a wide variety in the highest achieving machine learning classifier and word embedding pairs depending on what data the model had been trained and tested on. The Reddit datasets highest performing model was the USE word embeddings paired with the Naïve Bayes classifier, for the Twitter datasets it was the TFIDF paired with the Decision Tree classifier, and with 4Chan datasets it was the USE paired with the Multi-Layer Perceptron. There is also diversity in the highest average accuracy achieved by each dataset, with Reddit achieving maximum of 94.73%, Twitter 98.78% and 4Chan 96.02%.

Reddit Dataset					
Word Embedding / ML Model	Accuracy	Precision	Recall	F1 Score	
TFIDF / AdaBoost	0.912333	0: 0.92 1: 0.88	0: 0.93 1: 0.63	0: 0.95 1: 0.74	
Doc2Vec / Linear Discrimination	0.809667	0: 0.84 1: 0.54	0: 0.94 1: 0.27	0: 0.89 1: 0.36	
Hashing / Multi-Layer Perceptron	0.878000	0: 0.88 1: 0.87	0: 0.98 1: 0.46	0: 0.97 1: 0.86	
USE / Naive Bayes	0.947333	0: 0.96 1: 0.90	0: 0.98 1: 0.83	0: 0.97 1: 0.86	
BERT / Multi-Layer Perceptron	0.921667	0: 0.94 1: 0.84	0: 0.96 1: 0.75	0: 0.95 1: 0.79	

Figure 4: Top classifiers for each embedding with Reddit.

Twitter Dataset					
Word Embedding / ML Model	Accuracy	Precision	Recall	F1 Score	
TFIDF / Decision Tree	0.987833	0: 0.99 1: 0.99	0: 1.00 1: 0.95	0: 0.99 1: 0.97	
Doc2Vec / Linear Discrimination	0.865167	0: 0.88 1: 0.75	0: 0.96 1: 0.51	0: 0.92 1: 0.60	
Hashing / XGBoost	0.947500	0: 0.96 1: 0.90	0: 0.98 1: 0.84	0: 0.97 1: 0.86	
USE / SVC	0.960167	0: 0.95 1: 0.99	0: 1.00 1: 0.81	0: 0.98 1: 0.88	
BERT / SVC	0.934333	0: 0.94 1: 0.93	0: 0.99 1: 0.73	0: 0.96 1: 0.81	

Figure 5: Top classifier for each embedding with Twitter.

4Chan Dataset					
Word Embedding / ML Model	Accuracy	Precision	Recall	F1 Score	
TFIDF / Random Forest	0.921167	0: 0.91 1: 0.98	0: 1.00 1: 0.63	0: 0.95 1: 0.77	
Doc2Vec / Random Forest	0.795833	0: 0.80 1: 0.45	0: 0.99 1: 0.04	0: 0.89 1: 0.06	
Hashing / XGBoost	0.894833	0: 0.90 1: 0.84	0: 0.97 1: 0.59	0: 0.94 1: 0.69	
USE / Multi-Layer Perceptron	0.956833	0: 0.96 1: 0.97	0: 0.99 1: 0.80	0: 0.97 1: 0.88	
BERT / SVC	0.934500	0: 0.93 1: 0.93	0: 0.99 1: 0.72	0: 0.96 1: 0.81	

Figure 6: Top classifier for each embedding with 4Chan.

Due to the fact that the exact same methodologies, embeddings and classifiers were employed on each of the three datasets, this diversity in embedding and classifier pairs, and in the results achieved by these pairs can be explained by the differences in the collected data. As stated before, each of the three social media platforms has specific moderation methods. We believe that the diversity in results achieved by each single-platform model is largely due to the relative difficulty to detect the specific forms of hateful speech exhibited on that platform. Reddit's community based and strong moderation leads to "slurless" and subtle hateful language (e.g. "Get them all out of our country") which is difficult to detect and classify, whereas the automated and somewhat in-

adequate moderation on Twitter promotes the common use of typical hateful slurs (e.g. n*gger, f*ggot) which is much easier to detect and classify. 4Chan's no moderation policy leads to a diverse array of hateful speech, language and slurs (e.g. k*ke, n*gger, f*ggot, towelhead, mudskin), which ultimately makes it easier to detect and classify than the subtle Reddit hate speech, but harder to detect and classify than the repetitive Twitter hate speech.

For this reason it is extremely important in this day and age to produce multi-platform, versatile hate speech detection models that don't rely on the specific language used on a single social media platform.

5.2 Multi Platform

The goal of our study is to produce a high-achieving hate speech detection model that could span multiple social media platforms and produce reliable and replicable results. To carry out a multi platform analysis, we use the full 3 000 comment dataset of combined platform data. The dataset was split into training and testing data in a 0.3 ratio, and tested against the fifty word embedding and machine learning classifier pairs. The results of this process are in Figure 7.

The highest performing multi-platform hate speech detection model produced in this project was a combination of the Universal Sentence Encoder word embeddings paired with the Support Vector Machine (SVC) machine learning classifier, which achieved a peak average accuracy of 95.85%. The USE word embeddings achieving the highest accuracy result is relatively unsurprising due to the analysis carried out on the single-platform models, in which USE was identified as an extremely consistent and versatile word embedding method, regardless of the platform data. The box plot shown in Figure 8 shows that the pair (USE,SVC) exhibits the highest upper bound accuracy of all combinations at a value of 96.89%, and also exhibits a lower variation in accuracy results when compared to all other embedding methods.

The USE embedding combined with SVC exhibited a maximum average precision of 0.96, recall of 0.82 and F1 of 0.88 when classifying a comment as hateful. Each one of these individual values were the highest evaluation metric results achieved by any model trained using the multi-platform dataset (Figure 9). It also exhibited a maximum average precision of 0.96, recall of 0.99 and F1 of 0.97 when classifying comments as non-hateful, which apart from recall where some models equalled the highest result, were also the highest evaluation metrics achieved by any model trained on the multi-platform dataset. Results for non-hateful are shown in Figure 10.

Word Embedding / ML Model	Accuracy	Precision	Recall	F1 Score
TFIDF / Random Forest	0.949444	0: 0.95 1: 0.95	0: 0.99 1: 0.79	0: 0.97 1: 0.86
Doc2Vec / Linear Discrimination	0.900778	0: 0.90 1: 0.90	0: 0.99 1: 0.56	0: 0.94 1: 0.69
Hashing / XGBoost	0.902378	0: 0.92 1: 0.84	0: 0.97 1: 0.66	0: 0.94 1: 0.73
USE / SVC	0.956500	0: 0.96 1: 0.96	0: 0.99 1: 0.82	0: 0.97 1: 0.88
BERT / SVC	0.933944	0: 0.94 1: 0.91	0: 0.98 1: 0.74	0: 0.96 1: 0.81

Figure 7: Best ML classifier for each of the word embedding methods using the multi-platform dataset.

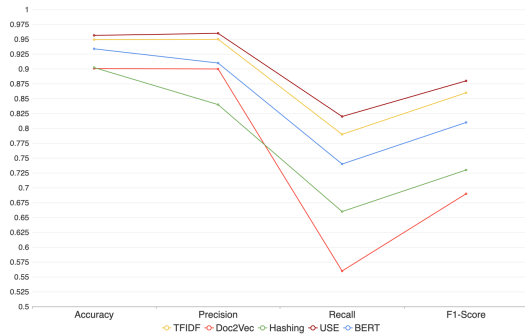


Figure 9: Evaluation metrics achieved by each of the best word embedding and classifier pairs when classifying data as hateful (X-Axis starts at 0.5).

	TFIDF	Doc2Vec	Hashing	USE	BERT	Average
Random Forest	15.5070	6.9077	5.9888	18.0517	27.1622	14.7235
Decision Tree	0.4148	0.0747	0.0620	0.6194	0.9643	0.4270
Naive Bayes	0.3136	0.0101	0.0574	0.0570	0.0933	0.1063
SVC	9.6417	0.1189	1.1423	0.7407	0.9411	2.5170
AdaBoost	2.0077	0.4864	0.3037	3.8224	5.8252	2.4891
Gaussian Process	10.4660	4.1397	5.1421	5.7909	2.4921	5.6062
K Neighbours	0.5193	0.0549	0.1422	0.1472	0.1317	0.1991
Multi-Layer Perceptron	14.2263	0.9683	2.7452	3.4606	2.8556	4.8512
XGBoost	4.2010	0.4059	0.7215	1.3560	2.3949	1.8159
Linear Discrimination	5.6602	0.1172	0.4563	0.4749	0.4895	1.4396
Average	6.2958	1.3284	1.6762	3.4521	4.3350	

Figure 11: Comparative average run times of each model tested on the multi-platform dataset.

5.3 Run Time & Efficiency

In the course of this study we measured both the single run time and the twenty run time average of the word embeddings and machine learning classifiers. The twenty runtime average for all classifiers and embeddings came out to exactly 20x the single run time, so we decided to only focus on the single runtime metric. In order to fairly evaluate this metric, we calculated both the average runtime of each machine

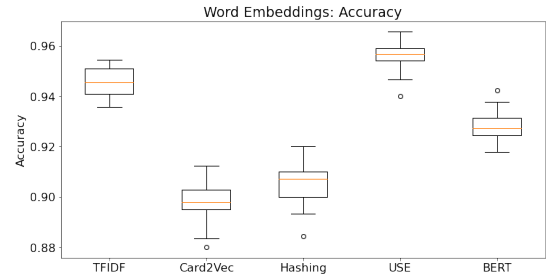


Figure 8: Accuracies achieved by the best machine learning classifier for each of the word embedding methods for the multi-platform dataset.

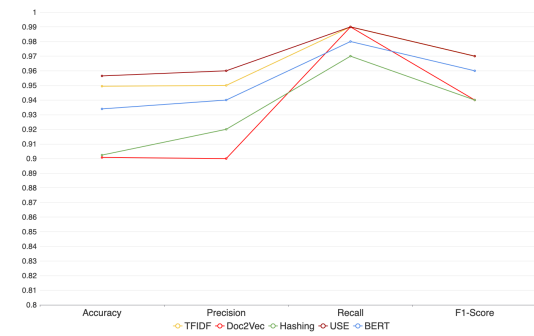


Figure 10: Evaluation metrics achieved by each of the best word embedding and classifier pairs when classifying data as non-hateful (X-Axis starts at 0.8).

learning classifier across all word embeddings, and the average run time of each word embedding method across all classifiers and noted the results.

The classifier with the highest average run time by quite a large margin was the Random Forest Classifier (14.7235s), and the classifier with the lowest average run time was the Naïve Bayes classifier (0.1063s). The Gaussian Process and Multi-Layer Perceptron classifiers also had notably high average run times (5.6062s and 4.8512s respectively), while the K-Neighbours and Decision Tree classifiers had notably low average run times (0.1991s and 0.4270s respectively). Regarding the word embedding with the highest average run-time, it was TFIDF (6.2958s), with BERT having the second highest average run-time of 4.3350s. Doc2Vec was the fastest performing word embedding method with an average run-time of 1.3284s and the Hashing word embeddings also performed well with a 1.6762s average run time.

With the run-time of each model calculated and noted, we determine which models exhibited the highest levels of efficiency. Efficiency refers to the ability of a machine learning model to produce accurate results, while also exhibiting a very short relative run-time. Figure 11 contains all of the run-times of the machine learning models tested during

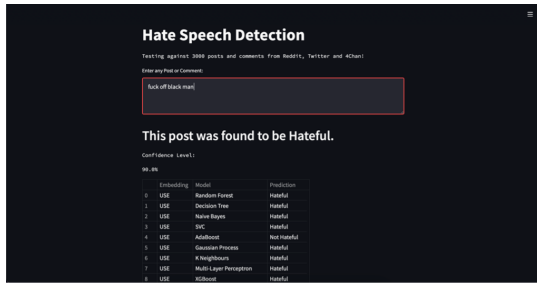


Figure 12: Output screenshot of our HateChecker.

this project. We then used this run-time table along with the accuracies table produced by each model to come to a conclusion as to the most efficient models.

Ultimately, we determined that the three models circled in red in the above table exhibited the highest levels of efficiency out of all tested models. The Doc2Vec and Naïve Bayes model, the Doc2Vec and K-Neighbours model, and the USE and Naïve Bayes model all exhibited an average run time of less than one second (0.0101s, 0.0549s and 0.0570s respectively), and also all exhibited a notably high degree of accuracy when compared to other models. The Doc2Vec and Naïve Bayes model achieved an overall accuracy of 83.02%, the Doc2Vec and K-Neighbours model achieved an overall accuracy of 88.55%, and the USE and Naïve Bayes model achieved an overall accuracy of 93.61%. All three of these models exhibited fast run-times in comparison to other models, and also achieved an high accuracy in comparison to other models, which makes them the most highly efficient and economical models.

6 HateChecker Application

HateChecker is the application we developed using the “streamlit” python library for the purpose of testing some of the most accurate hate speech detection models against a wide variety of different user inputted comments and posts. The HateChecker application takes input in the form of a comment or post like sequence of strings. Each of the twenty individual models selected will then use its own methodology to classify the user inputted comment as non-hateful (‘0’) or hateful (‘1’), and an aggregation of the classifications produced by these twenty models is then calculated and returned as an overall classification with a confidence level percentage included. Practically, the models that we selected for use in the HateChecker application employ each of the ten classifiers paired with both the USE and BERT word embeddings. We selected USE and BERT word embeddings for the HateChecker application over the other word embed-

dings, because the models produced using these word embeddings exhibited an absolute minimum average of 81.00%, while other word embedding methods exhibited accuracies as low as 43.78%. To ensure that the results produced by the HateChecker application where to a sufficiently high standard, we decided to rule out the other word embedding methods.

The HateChecker application was ultimately created and designed so that individual models could be tested by carefully designed user inputted data which may have not occurred in either the training or testing data which the model was built on. An example of this would be using the HateChecker application to analyse the models ability to classify sequences of strings filled with punctuation, numbers and special characters. By calculating an overall classification, confidence level percentage, and displaying which individual models made which classifications, we could begin to test our models on a wide array of different sequences of strings allowing us to evaluate and analyse weaknesses and strengths within our models.

7 Limitations

The presented findings of this study rely on the first (manual-)step of annotating comments, thereby the consideration of additional platforms and more comments could refine our results, indeed, each time we added more comments to the database, the results achieved by the models across the board would increase by 0.5% - 1%. Similarly, considering a larger set of embeddings techniques and classifiers could lead to finer results and different explanations at the end of the process when HateChecker returns its confidence score. Regarding the application, technical strategies could be deployed to improve the performance of HateChecker as many intermediate results are not yet saved, leading to long loading times when opening the software. On the hate detection performances, advance annotations could be explored to detect subtle hate-sentences. Finally, the quality of the presented and shared set of annotations is bound in time as languages and usages evolve across time, with haters often finding new ways to convey their ideas.

8 Conclusion

In this study, we explore a strategy to detect hate-speech. We based our approach on considering messages from several online social-media platforms at once, betting that their different internal moderation policies would provide a larger set of haters’ meth-

ods. In additions to sharing our annotated set with the community, we also develop an application building up our strategy of combining/comparing multiple pairs of word embeddings and classifiers. Overall, we successfully build a hate speech detection model, pairing USE and SVC, to obtain an average accuracy of 95.65% and achieved a maximum accuracy of 96.89%. Moreover, our application allows to define an aggregating strategy by *e.g.* choosing which pairs should be taken more into account. Therefore, we hope that this two-side strategy of involving several platforms and combining multiple pairs of embeddings and classifiers, will inspire the community to improve our results and refine our performance score.

Sensitive Content Warning Due to the nature of this study, there are places in this article where hateful language and terms are used. While we did try and keep the use of these terms and phrases to a minimum, and while we obviously do not approve these message, it was vital to provide the reader with a proper understanding of the context and methodologies used in the process of completing this project.

REFERENCES

- Alshalan, R. and Al-Khalifa, H. (2020). A deep learning approach for automatic hate speech detection in the saudi twittersphere. *MDPI*.
- Bhattacharya, D. and Weber, I. (2019). Racial bias in hate speech and abusive language detection datasets. *ACL Anthology*.
- Cer, D., Yang, Y., Kong, S.-y., Hua, N., Limtiaco, N., John, R. S., Constant, N., Guajardo-Cespedes, M., Yuan, S., Tar, C., et al. (2018). Universal sentence encoder. *arXiv preprint arXiv:1803.11175*.
- Chen, T. and Guestrin, C. (2016). Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pages 785–794.
- Davidson, T., Warmusley, D., Macy, M., and Weber, I. (2017). Automated hate speech detection and the problem of offensive language.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv:1810.04805*.
- Freund, Y. and Schapire, R. E. (1997). A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of computer and system sciences*, 55(1):119–139.
- Gibbs, M. N. (1998). *Bayesian Gaussian processes for regression and classification*. PhD thesis, Citeseer.
- Guo, G., Wang, H., Bell, D., Bi, Y., and Greer, K. (2003). Knn model-based approach in classification. In *OTM Confederated Int. Conf.*, pages 986–996. Springer.
- Han, E. (2020). Countering hate on tiktok.
- Hawdon, J., Oksanen, A., and Räsänen, P. (2015). Online extremism and online hate exposure among adolescents and young adults in four nations.
- Hornik, K., Stinchcombe, M., and White, H. (1989). Multi-layer feedforward networks are universal approximators. *Neural networks*, 2(5):359–366.
- Izenman, A. J. (2013). Linear discriminant analysis. In *Modern multivariate statistical techniques*, pages 237–280. Springer.
- Jourová, V. (2016). Code of conduct - illegal online hate speech questions and answers.
- Kansara, K. and Shekokar, N. (2015). A framework for cyberbullying detection in social network. *Semantic Scholar*.
- Kovács, G., Alonso, P., and Saini, R. (2020). Challenges of hate speech detection in social media. *Springer Nature*.
- LIGHT (2020). Rising levels of hate speech & online toxicity during this time of crisis.
- Le, Q. and Mikolov, T. (2014). Distributed representations of sentences and documents. In *International conference on machine learning*, pages 1188–1196. PMLR.
- Lee, H. S. and Lee, H. R. (2018). An abusive text detection system based on enhanced abusive and non-abusive word lists. *Yonsei University*.
- Mozafari, M. (2020). Hate speech detection and racial bias mitigation in social media based on bert model. *PLOS ONE*.
- Pal, M. (2005). Random forest classifier for remote sensing classification. *International journal of remote sensing*, 26(1):217–222.
- Ramampiaro, H. (2018). Detecting offensive language in tweets using deep learning. *Cornell University*.
- Rish, I. et al. (2001). An empirical study of the naive bayes classifier. In *IJCAI 2001 workshop on empirical methods in artificial intelligence*.
- Rudnicki, K. and Steiger, S. (2020). Online hate speech.
- Röttger, P., Vidgen, B., Nguyen, D., Waseem, Z., Margetts, H., and Pierrehumbert, J. (2021). Hatecheck: Functional tests for hate speech detection models.
- Safavian, S. R. and Landgrebe, D. (1991). A survey of decision tree classifier methodology. *IEEE transactions on systems, man, and cybernetics*, 21(3):660–674.
- Salminen, J., Hopf, M., Chowdhury, S., gyo Jung, S., Almerexhi, H., and Jansen, B. (2020). Developing an online hate classifier for multiple social media platforms. *Human-centric Computing and Info. Sciences*.
- Sap, M., Card, D., Gabriel, S., Choi, Y., and Smith, N. (2019). The risk of racial bias in hate speech detection.
- Vapnik, V. (1998). Statistical learning theory new york. NY: Wiley, 1(2):3.
- Vigna, F. D., Cimino, A., and Petrocchi, M. (2017). Hate me, hate me not: Hate speech detection on facebook. *First Italian Conference on Cybersecurity*.