Automating Online Hate Speech Detection: A Survey of Deep Learning Approaches

Ashe Magalhaes

Master of Science School of Informatics University of Edinburgh 2019

Abstract

The problem of automating the detection of online hate speech is made more critical by the moral and legal consequences of its propagation in a social media network. This research surveys the current landscape of online hate speech detection, deep learning model architectures, embedding choices, and the inclusion of user behavior metrics, in order to improve performance on detecting the hateful class in a Twitter dataset. We find that Google's pretrained BERT embeddings outperform the inclusion of user behavior metrics across many model choices and that a single-layer multiple input Convolutional Neural Network with many filters performs best. The contribution of this paper is a range of experiments that inform machine learning researchers on promising future paths and policymakers on the capabilities and limitations of automated systems on the task of detecting hate speech.

Keywords: *Hate Speech Detection, Deep Learning, Natural Language Processing, Twitter*

Acknowledgements

I would like to thank my advisor, Walid Magdy, for his guidance, flexibility, and feedback. It was consistently fun to brainstorm with Walid and this brought a valuable level of excitement to the lengthy process of writing this dissertation. Thank you to Antreas Antoniou, who kept me engaged with this topic and challenged me to grow into a more robust and creative machine learning engineer. Thank you to Abeer AL-Dayel, who offered valuable insight around model performances. Finally, I would like to thank my family, mentors, and friends for their support and encouragement–you are all instrumental in my ability to work towards lofty goals.

Statement of Originality

I confirm that I have read and understood the University's plagiarism guidelines. All work submitted is my own unless referenced, in both the text and the bibliography.

Table of Contents

1	Intr	oduction	1
	1.1	Motivation	1
	1.2	Outline	3
2	Bac	kground	4
	2.1	History of Automating Hate Speech Detection	4
	2.2	Hate Speech Datasets	5
	2.3	State of the Art Approaches	6
		2.3.1 Data Processing	6
		2.3.2 Embedding Choices	7
		2.3.3 Model Architectures	7
		2.3.4 Evaluation Metrics	8
	2.4	Metrics of User Behavior	8
	2.5	Challenges	9
3	Met	hodology	10
	3.1	Data Collection	10
	3.2	Data Processing	11
	3.3	Environment Configurations	13
	3.4	Model Architectures	13
4	Exp	eriment Design	15
	4.1	Phase 1: Tweet Embeddings	15
	4.2	Phase 2: Tweet + Reply Embeddings	15
		4.2.1 Reply Pairings	16
		4.2.2 Reply Network Metrics as Embedding Coefficients	16
		4.2.3 Multiple Input Model Architecture	17
	4.3	Phase 3: Tweet + User Topic Embeddings	17

		4.3.1 LDA Topic Modeling Embeddings	17
	4.4	Hyperparameter Tuning & Evaluation	17
	4.5	Limitations	18
5	Resu	ılts	19
	5.1	Phase 1 Results	19
	5.2	Phase 2 Results	19
	5.3	Phase 3 Results	22
	5.4	Hyperparameter Tuning	24
	5.5	Tuned Model Results Across Experiments	25
	5.6	Results Summary	26
	5.7	Experiment Difficulties	27
6	Con	clusion	31
	6.1	Technical Contributions	31
	6.2	Policy Recommendations	32
Bi	bliogi	aphy	34
A	Exp	eriments	40
	A.1	Code	40
	A.2	Machine Configuration	40
	A.3	Additional Results	40
B	Glos	sary	43
С	Wor	kshop Paper	46

Chapter 1

Introduction

1.1 Motivation

In recent years social media companies like Facebook, Twitter, and Alphabet Inc.'s YouTube have been widely criticized for failing to detect and remove hate speech from their platforms [1]. Terrorist groups like al-Qaeda, ISIS, the Taliban, and Boko Haram have repeatedly been unhindered in their releasing of hateful content, used for the purposes of recruiting and radicalizing followers [2]. In March 2019, a white supremacist posted racist manifestos–which spread through Twitter–and live-streamed the shooting of over 51 people on Facebook and YouTube. While both companies leverage machine learning algorithms to automatically detect and remove such content, they failed to take down the content quickly enough. Facebook's head of public policy defended the platform's slow response time:

"The video was a new type that our machine learning system hadn't seen before. It was a first person shooter with a GoPro on his head...This is unfortunately an adversarial space. Those sharing the video were deliberately splicing and cutting it and using filters to subvert automation. There is still progress to be made with machine learning [3]."

Progress needs to be made sooner rather than later. Hateful content on social media contributes to real-world violence [4], recruitment to and propaganda for terrorist individuals/groups [5], makes other users feel less safe and secure on social platforms [6], and triggers increased levels of toxicity in the network [7].

Machine learning is a powerful tool for natural language processing (NLP) tasks such as hate speech detection. Statistical learning methods learn weights for humandesigned representations of feature data. In contrast, deep learning, a subfield of ML,

Chapter 1. Introduction

shifts the burden of time-intensive feature design from linguistic experts to the learning system. Supervised deep learning methods have achieved high levels of performance in NLP tasks such as sentiment analysis, language generation, and information re-trieval [8]. For the hate speech detection task, we have a corpus of posts, tweets, or videos labeled as hate speech or not, with automated systems learning the nuanced interpretation required to quickly detect hateful content at scale.

Governments are enforcing the need to find reasonable solutions, fast. The Sri Lanken government temporarily banned social media networks three separate times in the wake of the Easter 2019 suicide bombings that killed hundreds of people, in order to prevent "social unrest via hate messages" [9]. Germany and the UK have strict legislation against hate speech [10]. French leader Emmanuel Macron has made fighting online hate speech a priority, taking meetings with Mark Zuckerberg and other high-level Facebook representatives; as a result Facebook has agreed to hand over data on French users suspected of spewing hate speech online. This is an international first, and according to a counsel at law firm Linklaters, "a strong signal in terms of regulation" [11]. EU countries have threatened to fine social networks up to 50 million euro per year if they continue to fail to act fast enough [12].

This paper aims to navigate the challenge of creating systems that distinguish hateful content from abusive, normal, and spam content when there are very few hateful samples to learn from. There are 500 million tweets produced daily [13], and the vast majority of those tweets are not considered hate speech [14]. We echo the concerns of Facebook's head of public policy by asking:

How can we improve the performance of automated systems on identifying hate speech when they must learn from very few hateful samples?

Of particular relevance to this paper is modeling user behavior on Twitter. User timelines and user interactions provide information on each user's surrounding community and offer a more detailed picture of identifying 'hateful users', or those who are prone to disseminating hateful content. Hateful users have been shown to tweet more often and favorite more tweets than 'normal' users [15]. The median hateful user tends to have higher network centrality than normal users according to metrics like betweeness and out-degree. Augmenting inputs to deep learning classifiers with metrics of user behavior could help moderation teams quickly identify hateful content through profiles.

Social media has contributed to a more open and connected world [16]. It has promoted Western liberal values through its effects on protest mobilization [17], community building [18], and accountability in governments and institutions [19]. It is critical that we do not lose the benefits of these platforms as they operate under increased regulation and scrutiny. So, while automated systems can and should enforce uniform policies for user behavior that fosters welcoming online environments, this should not come at the cost of dissenting and fringe views. This research surveys the capabilities and limitations of state-of-the-art (SOTA) deep learning classifiers. We aim to inform policy and decision makers, who must reconcile the benefits of social media platforms with the harms they threaten when hate speech is allowed to propagate.

1.2 Outline

We begin by exploring the history of online hate speech as well as the challenges of collecting and annotating a hate speech dataset from social media platforms. We then describe related research and the model architectures which have shaped the current landscape of automated hate speech detection. We explain how this research informs our methodology and go on to describe our data, how it is collected and processed, the details of our model architectures, and metrics for evaluation. We hypothesize that the inclusion of user behavior features in feature embeddings to deep learning model architectures will improve performance of the hate speech detection task. Finally, we interpret the results from both a technical and policy perspective in order to a) enumerate and qualify the contributions of our experiments in relation to the existing technical body of related research and b) advise policymakers on how such results may be acted on. We find that pairing sophisticated pretrained embeddings with user behavior metrics in a multiple input CNN architecture achieves competitive performance given the scarcity of hateful samples, but does not perform at a level acceptable for complete automation of the task.

While this paper describes technical experiments, its intention is to offer an accessible discussion on this topic to readers from all disciplines. A glossary is offered in Appendix B to define technical jargon. The hope is that computer scientists, linguists, policymakers, lawyers, and others can form interdisciplinary coalitions that are better equipped to tackle the complex and important task of hate speech detection.

Results from this research are compiled in a workshop paper, targeted at Fall 2019 conference submissions. See Appendix C for the full draft.

Chapter 2

Background

2.1 History of Automating Hate Speech Detection

Early work on online hate speech detection focused on the necessity of preserving free speech despite hate speech [20, 21]. The question of whether online forms of expression translate into negative offline actions has been hotly contested. Stormfront, considered the first "hate website", was launched in 1995 by a Ku Klux Klan leader as a forum for discussing ideas related to white separatism and white nationalism. Members of the site claimed to be covered under their American First Amendment rights for expression, arguing that they did not act on the forum's sentiment. A subsequent study conducted in 2014 found users of the site indirectly responsible for "the murders of nearly 100 people in the preceding five years" [22].

In 2014 the UN Human Rights Council observed the proliferation of hate speech online with a focus on ethnicity, nationality, and religion [23]. Subsequent research demonstrated that online hate speech incites real-world violence against refugees in Germany [4] and that in online environments toxicity from one source is likely to foment both reciprocal and sympathetic toxicity from adjacent users [7]. Public criticisms against Facebook, Twitter, and YouTube's slow responses to combating the spread of hate speech have impelled them to hire thousands of people to work on the problem [24]. In the meantime, strict online policies against offensive language or sentiment that falls short of hate speech has resulted in the rapid rise of new platforms that cater to fringe online communities. The social network Gab was created as an alternative to Twitter, welcoming users banned or suspended from other platforms. Researchers have found that its posts contain hate words at a rate 2.4 times higher than Twitter [25]. In spite of the growing phenomenon of hate speech, prior to 2016, NLP research on hate speech was limited [6, 26]. Since then, statistical machine learning techniques like logistic regression and naive Bayes classifiers have been applied to this area with success; logistic regression in particular is widely considered a baseline model for experiments today [26, 27, 28]. Recently, the aforementioned pressures, coupled with concurrent advances in NLP and deep learning, has sparked a resurgence in interest in the field. The success of deep neural networks on NLP tasks like sentiment analysis [29], named entity recognition [30], and part-of-speech tagging [31] have made it a natural choice to apply to the task of classifying tweets or posts as hateful.

2.2 Hate Speech Datasets

It is a nontrivial task to create datasets for training machine learning classifiers on hate speech detection. The common practice is to deploy a collection of tweets or posts to a crowdsourcing site, where a number of human annotators identify an appropriate classification (i.e. hateful, normal, abusive, spam). This is problematic for a few reasons and highlights the ways in which human-annotated data reflects social biases. First, the subjective nature of identifying hate speech makes it difficult to extract ground truth labels for tweets or posts [14]. There are many definitions of hate speech, which vary across different languages, cultures, and governments. Many of the definitions involve the hateful targeting of members of a group, like this one:

"Language used to express hatred towards a targeted individual or group, or is intended to be derogatory, to humiliate, or to insult the members of the group, on the basis of attributes such as race, religion, ethnic origin, sexual orientation, disability, or gender." [32]

Similarly, the UN's definition is:

"The term hate speech is understood as any kind of communication in speech, writing or behaviour, that attacks or uses pejorative or discriminatory language with reference to a person or a group on the basis of who they are, in other words, based on their religion, ethnicity, nationality, race, colour, descent, gender or other identity factor. This is often rooted in, and generates intolerance and hatred and, in certain contexts, can be demeaning and divisive." [33]

Second, because there is no universally agreed upon definition of hate speech [18], it is often difficult for human annotators to separate hate speech from offensive language. Given the legal consequences of hate speech, it is important that annotators and machine learning classifiers distinguish between hateful language and offensive language [14, 32]. It has been found that showing users a definition of hate speech does not improve annotation reliability, hinting that the presence of hate speech is not a binary decision and annotators likely need more context [28], such as sample post-label pairings.

Third, it is likely that annotators skim through tweets or posts too quickly and fail to pick up on the context of the language. Past research demonstrates that annotators will incorrectly label normal text as hate speech if it contains hate or curse words [32]. Annotators are likely to miss out on hate speech without a slur as well as interpret racist and homophobic slurs as hateful but sexist slurs as 'just offensive' [26, 32].

Creating robust hate speech datasets remains an active area of research with persistent challenges. While this research focuses on the implementation of deep learning model architectures, it is clear that progress made in reconciling the social biases in the collection of labeled hate speech data will reduce algorithmic bias.

2.3 State of the Art Approaches

The application of deep neural networks to the task of hate speech detection is fairly recent. While we define SOTA approaches as those specifically cited in the hate speech detection literature, other advances in NLP (i.e. Open AI's recent transformer language model [34]) offer novel contributions to the task. We categorize implementations based on how data is processed, model architecture choices, and the feature embedding types used to convert text into inputs suitable for machine learning. Of particular relevance to this work is research that identifies hateful or abusive language on Twitter.

2.3.1 Data Processing

Hate speech detection research focuses less on data "cleaning" (i.e. removing stop words, punctuation, capitalization, etc. from tweets or posts) than other natural language tasks, likely because these language features capture more information than in more formal text. For example, some users respond to platform restrictions on speech by creating posts that replace letters with symbols (e.g. a\$\$) in order to convey their otherwise 'illegal' message undetected [35]. A standard procedure is for researchers to lowercase and tokenize tweets, as well as remove URLs. Researchers may also remove emojis, although keeping emojis will likely boost NLP-task performances [36].

Another form of data processing is bootstrapping; this method identifies hateful user accounts from annotated tweets and obtains all tweets from their timeline, labeling them hateful as well, in order to increase the number of samples and better balance the dataset in cases where its imbalanced [6].

2.3.2 Embedding Choices

A tweet or post must be made into a feature embedding in order to be fed to our model as a vector of real numbers. An embedding is a representation of model inputs that captures some of the semantics of the input in a relatively low-dimensional space. Pooling layers or trimming and padding of tweets [37] can be used to convert tweets into a fixed length for standardized inputs to our neural networks. N-grams, at both the word-level and character-level, are widely used for converting a sequence of text into a sequence of embeddings [38, 39]. Character n-grams have been used for natural language processing tasks such as authorship identification [40], native language identification [41] and machine translation [42]. Word-level and character-level embeddings have been concatenated and made into feature embedding with success, as demonstrated by [38]. Word2Vec and GloVe (see appendix) embeddings are staples of NLP research and can be compared to randomly initialized embeddings as a baseline. While GloVe is said to better account for global contextual information, in practice researchers have had success with both, as they each vary in efficacy on different datasets. Word-based frequency vectorization methods, such as term frequency-inverse document frequency (TF-IDF), are also commonly used embedding choices and offer the advantage of being language agnostic [37]. Pretrained embeddings, in contrast, are dependent on the language of the corpus they are trained on.

2.3.3 Model Architectures

The baseline machine learning model choice, which tends to perform well, is logistic regression [32, 43]. Another baseline model, support vector machines (SVMs)–with linear kernels–perform similarly to logistic regression in practice [44]. Linear models tend not to capture the sparse features in document classification tasks, undermining the detection of relations between words or symbols in sentences. CNNs and RNNs (particularly LSTMs) have achieved SOTA performance due to their ability to capture long-term dependencies in a sentence, tweet, or post through their ability to maintain information in memory for a period of time [27, 37].

2.3.4 Evaluation Metrics

Because of the scarcity of hateful content relative to benign content [14], it is critical we evaluate models based on their classification of actual hateful or abusive language. Otherwise, a model could achieve a high classification accuracy by labeling everything as non-hateful, repeatedly failing our task. F-score, precision, and recall are commonly used as evaluation metrics, because they better capture the class contributions in an imbalanced dataset. Precision is the ratio of the number of tweets correctly classified to a given class divided by the total number of tweets classified to that class; recall is the ratio of tweets correctly classified to a given class. The f-score is the harmonic mean of both. We express them below with true positive (TP) and false positive (FP) rate:

$$Precision = \frac{TP}{TP + FP} \tag{2.1}$$

$$Recall = \frac{TP}{TP + FN}$$
(2.2)

$$F1 = \frac{2 \cdot Precision \cdot Recall}{Precision + Recall}$$
(2.3)

In our task, precision answers: of all the tweets labeled as hateful, how many were actually hateful? Recall answers: of all the hateful tweets, how many did we label? F-score is the weighted average of both precision and recall, providing a measure of if a system can achieve high precision and recall.

Due to the difficulty in acquiring hate speech datasets and relatively small dataset sizes, 10-fold cross validation is a widely used metric of system performance [27, 37]. This is because it allows the model to learn from more observations, as all observations are used for training and validation. Another interesting method is ensemble classification, in which multiple classifiers predict the label of an input and the label with majority consensus is chosen. If all classifiers disagree, the classifier with the strongest confidence in a prediction is chosen [37].

2.4 Metrics of User Behavior

Hateful users tend to make themselves central to the Twitter social network: they target more popular users and they have more statuses and followees per day, although they have younger accounts [6, 15]. While this may increase the user's visibility, the inclusion of hateful or antagonistic content in a tweet reduces the rate of retweet by a factor of 45 [45]. Despite this, hateful users are central to Twitter reply networks as determined by features such as betweenness and eigenvector centrality [15]. Furthermore, the modeling of reply trees—rooted at a tweet and joined by replies—reveals that the in-degree of the tweet, or the number of retweets, plays an important role in the overall shape of the reply tree and a component of the Twitter reply network as a whole [46].

Latent Dirichlet Allocation (LDA) [47] is a method of topic modeling that extracts the main topics in text. It has been used to improve the performance of linear machine learning classifiers on the task of identifying abusive social media content [48]. While LDA-derived topics have been used as annotations [49], to the best of our knowledge it has not been included in feature embeddings in order to improve the performance of deep learning classifiers on the task.

2.5 Challenges

This section highlights relevant background and previous work to our task of hate speech detection; it is also worth summarizing the challenges cited by other researchers in order to inform our methodology. First, extensive data processing may remove valuable information from social media content. Given the limited number of hate speech datasets available, it is important to consider the social biases and fallacies that are inherent to crowdsourced labels. Given the scarcity of hateful content relative to benign content, evaluation metrics that account for imbalanced classes should be used. Because of the legal and moral consequences of identifying a post or user as hateful, it is critical that our classifier distinguish between hateful and abusive language. Finally, the distinction between hateful content and hateful users is worth defining [15], as the task of identifying hateful users is easier but not equivalent to the task of identifying hateful content.

Chapter 3

Methodology

Objective. The overall goal of this research is to conduct a survey of deep learning model architectures, embedding choices, and feature inputs in order to automate the identification of hateful tweets from a dataset with few hateful samples. Unlike previous work, we experiment with user behavior metrics in multiple input model architectures to provide our deep learning classifiers with more context to identify hate speech with.

3.1 Data Collection

The dataset for this research comes from Founta et al.'s work that describes the process of large scale crowdsourcing for annotations of hateful, normal, abusive, and spam tweets [14]. Founta et al. provide status ids for others to obtain tweets through the Twitter API. After learning that a number of tweets were no longer available on Twitter, they provided 100k tweets' status ids and their associated majority labels. Upon using the Twitter API, we find that only 64% of those tweets are still available online; we assume the remainder have been taken down by Twitter or the users. The class distribution is as follows: **4% hateful, 20% abusive, 62% normal, and 14% spam**; there are **64,149 tweets total**. Half the tweets achieved annotation majority with an agreement of more than 3 out of 5 votes, which should give us a degree of confidence in the annotations. Figure 3.1 illustrates a word cloud of each class.

The average retweet count of hateful content is 15% higher than the average retweet count of normal content; the average favorite count of hateful content is 76% higher than the average favorite count of normal content. Finally, on average hateful content is a reply to existing tweets 31% more often than normal content is.



Figure 3.1: Word clouds of tweets from the hateful, abusive, normal, and spam classes.

For each tweet, if it is in reply to a tweet we use Twitter's API to collect the context tweet. 19% of the total tweets are a reply. Of those, 14% are no longer available on Twitter. For each tweet, we crawl the user timeline of the author and collect their 200 most recent tweets, giving us around 12.8 million additional tweets for LDA topic modeling.

3.2 Data Processing

We clean the tweets by tokenizing, lowercasing, and removing punctuation. We keep URLs because they provide important context; tweets with URLs are more likely to get retweeted [50]. We experiment with three types of text embeddings:

- 1. **TF-IDF embeddings.** TF-IDF fits the training set into a weighted vector by normalized frequency of the 10,000 most common words in our vocabulary. Our validation and test set are transformed using the learned weights.
- 2. **Pretrained Twitter embeddings.** These embeddings are from a Word2Vec model trained on 400 million raw English tweets, with an embedding dimension of 400 [51].
- 3. **Pretrained BERT embeddings.** Google's BERT, or Bidirectional Encoder Representations from Transforms, is a novel method of pre-training language representations which obtains SOTA results on a range of NLP tasks described in [52].

	Instigator Tweet	Annotated (Reply) Tweet
ful	want some cheese with that whine?	wow nice retort, putin lover. anything
late		minorities say = whine. kill all f*gs,
		right?
ve	Bombshell report over the weekend	You bunch of ass clowns are as bad
busi	shows a 'very high up' Obama official	as the illegitimate Pres we got stuck
A	unmasking Trump associates for polit-	with. Leaders in bs news. #worstpre-
	ical pu https://t.co/UwNWfa27ea	sever
al	watching the Rays game in NY is frus-	Whoa. Chill. I may have to put this
lorn	trating AF. I have 2 listen 2 the an-	friendship on hold if you bad mouth
	nouncers fave the Yankees & it goes	my Yankees
	against everything I believe in	
ш	This great video from @ will make	Thanks for sharing @! Anyone in-
Spa	you want to hit a bike trail. https:	terested in a MTB road trip to NI,
	//t.co/VvK0KAJh4T	plan your trip at https://t.co/
		mKBUPSDWZc #RideNorth

Table 3.1: Sample tweet reply pairings for hateful, abusive, normal, and spam classes.

Specifically, we use the BERT-12-768-12 model with an embedding dimension of 768.

For our word embeddings, we trim or pad each tweet to 17 words, the average size of a tweet in our dataset. When padding, we add randomly generated embeddings so that each tweet embedding is a n-gram, where n=17. To account for the imbalance of classes in our dataset, for our training experiments we draw samples through a multinomial distribution that accounts for each class weight. Our balanced batches allow the classifier to learn class representations more consistently each epoch. We do not do this for validation and testing because a) we run the risk of tweet samples that are seen multiple times or never in our evaluation and b) it is not a realistic representation of tweets in the wild. Table 3.1 displays a random sample of each of our categories of tweets, with personally identifiable information removed.

3.3 Environment Configurations

Each experiment type is run with three or more different seeds, and the mean and standard deviation of results are reported. This adds robustness to our experiments and higher confidence in our results. In order to ensure reproducibility, we set the random number values and random number generators of our data providers, PyTorch (deep learning platform) backends, and NumPy (scientific computing package) to our experiment seed. This is so that computations are deterministic and replicable.

All experiments use the cross entropy loss function which takes in the classifier's prediction score for each class, applies a softmax layer, and measures its divergence from the true label. For our deep learning experiments, we use the Adam optimizer, an adaptive learning rate optimization algorithm designed for deep neural networks; for our logistic regression experiment we use limited-memory BFGS, a popular optimizer for linear machine learning models. Finally, we use a scheduled learning rate that relies on a cosine annealing schedule to more quickly reach a learning rate that reduces noise in our parameter updates [53].

3.4 Model Architectures

Logistic Regression Our baseline model is multi-class logistic regression which uses the L2 norm penalty as a form of regularization. We do not expect it to perform well on the hateful and abusive classes because this model assumes the data is linearly separable.

Multilayer Perceptron. Known as a "vanilla neural network", this model consists of fully connected layers with the non-linear activation function leaky ReLu applied. We use Leaky ReLus here and in other models because it allows a small, positive gradient when the unit is not active and tends to perform better than nonleaky ReLu in practice.

Convolutional Neural Network (CNN). While this model is traditionally used for image inputs, it has performed quite well on NLP tasks and has had success when applied to hate speech detection [38]. The convolutional layer computes the output of neurons (i.e. mathematical operations) that are connected to local regions in the input. It consists of a set of learnable filters that produce an activation map used to ultimately compute the class scores. For a more thorough explanation of CNNs, see [53]. Our CNN consists of a series of convolutional layers, batch normalization, leaky ReLu,

and dropout. The output is processed through a max pooling layer and then a fully connected layer which computes class probabilities. Our CNN contains a context list that gives every layer access to all the previous layers, allowing for more connectivity and greater chance of the network learning from long-term dependencies in the tweet.

DenseNet. CNNs are more accurate and efficient if they contain shorter connections between layers close to the input and output [54]. DenseNet leverages this observation by connecting every layer to every other layer: for each layer, the feature maps of all preceding layers are used as inputs into all subsequent layers. Our DenseNet implementation has achieved SOTA results for image datasets; we experiment with its application to NLP.

LSTM. LSTMs (and RNNs more broadly) improve on traditional deep neural networks by propagating information in both directions. This loop in the architecture acts as a memory state by which the network makes adjustments in the information flow, allowing the model to selectively remember or forget. This model has performed well on NLP tasks in practice because of its ability to remember dependencies and between sentences; it is applied to the task of hate speech detection by [27] [44] [38]. Our implementation passes our input through a series of LSTM layers to a final fully connected layer.

Chapter 4

Experiment Design

The experiments consist of three phases that demonstrate the effect of user behavior metrics on a combination of models and embedding types. Each of our 4 model choices are run with our 3 embedding types, for 4 rounds of experiments, with 3 different seeds, for a total of **144+ experiments before tuning**.

4.1 Phase 1: Tweet Embeddings

We compare our deep learning models to a baseline logistic regression model to begin with an idea around how much of an improvement they can offer. Here the feature embedding to our deep learning models are simply the tweet embeddings. Table 4.1 enumerates regularization choices, number of layers, and model-specific hyperparameter values. Regularization reduces the number of parameters and amount of computation in the network, which helps to control overfitting. The number of layers and hyperparameter choices are set with values that are commonly used in the literature and with the expectation that tuning will be conducted down the line on the best performing models.

4.2 Phase 2: Tweet + Reply Embeddings

The goal of this round of experiments is to apply some type of context to our embeddings.

MODEL	REGULARIZATION	LAYERS
LR	L2-NORM PENALTY	-
MLP	MP, λ=1ε-4	3 FC-LERELU-FC
CNN	MP, BN, δ=0.5, λ=1ε-4	3 CONV-BN-LERELU-FC
DENSENET	AP, BBN, δ=0.5, λ=1ε-4	4 DENSE-BBN-RELu-CONV-FC
LSTM	δ=0.5, λ=1E-4	3 LSTM-FC

Table 4.1: Baseline model architectures. Terms include Weight decay (λ), Max Pooling (MP), Average Pooling (AP), Dropout (δ), Batch Normalization (BN), and Bottleneck Layer, which includes BN (BBN). Model-specific hyperparameters include CNN filters=8, DENSENET growth-rate=12, LSTM hidden-layers=5.

4.2.1 Reply Pairings

First, because of the statistics collected around the behavior of hateful users and retweeting, we define pairs of tweets: the original tweet and a context tweet. The context tweet is for the case that the tweet was a response to something else. If the tweet is not a response to something else, we add null embeddings to indicate that there is no context for the tweet. Although only 19% of our tweets in the dataset are replies, the extra information around context (when it is there and not there) may help our network learn in some way. See Table 3.1 for an example of reply pairings.

4.2.2 Reply Network Metrics as Embedding Coefficients

Next, because tweet in-degree has been shown to be significant ([46], see Background) we focus on in-degree in terms of number of times someone has retweeted a given tweet and number of times someone has favorited a given tweet. Our maximum **retweet value is 154,565** and our **maximum favorite value is 27,741**. Because these numbers are so large, we log each (and use 0 when the count is 0), so as to not skew the learned weights. We concatenate the logged retweet and favorite counts to our tweet and reply embeddings. We are interested to see if the network can better learn from the retweet and favorite numbers as a measure of context.

4.2.3 Multiple Input Model Architecture

Here and for the remaining experiments, we shift to building a neural network with multiple inputs in order to have the network learn from the annotated tweet in addition to other types of embeddings. We have the annotated tweet embedding and the context tweet embedding as separate inputs; they are processed by different parts of the model architecture. The learned features for each are concatenated and fed into a final fully connected layer. See Figure 5.2 for a visualization of a multiple input network architecture.

4.3 Phase 3: Tweet + User Topic Embeddings

This phase aims to add context in a more sophisticated way.

4.3.1 LDA Topic Modeling Embeddings

For each tweet in our dataset, we crawl the author's user timeline and collect 200 tweets. We then conduct topic modeling through the LDA approach (discussed in Background). We use LDA, as opposed to other topic modeling techniques, because LDA represents documents (or tweets) as random mixtures over topics in the corpus, which reflects what we expect from tweets on a user's timeline [47]. For each of our users, we find the dominant topic in their timeline tweets and create embeddings of a sequence of the top 10 words related to the topic. We choose the same embedding that is used to embed our tweets for a given experiment. We also concatenate the coherence and perplexity scores of the user's timeline to each embedded topic word to add a global measure of topic modeling. Perplexity is a measure of how well a probability model can predict a topic (lower is better) and coherence captures how well topics can be defined (higher is better).

Our multiple input model architecture processes the topic embeddings and tweet embeddings. Again, we concatenate the learned features and use this as input to the final fully connected layer.

4.4 Hyperparameter Tuning & Evaluation

We tune our best performing models-by model type and embedding type-by experimenting with learning rate, regularization, number of layers, and the model specific parameter. We define best performing model as the model with the greatest validation set f-score on the hateful classes. We evaluate our results by overall metrics of f-score, precision, and recall as well as f-score of the hateful and abusive classes. This is because we are concerned with our classifier's ability to correctly identify hateful and abusive samples. To achieve stability in the results, we evaluate the mean and standard deviation of 3 experiments each.

4.5 Limitations

In order to keep our research in the scope of our central questions, we choose experiments that demonstrate the effect of user behavior in feature embeddings on model performance. We acknowledge that there are alternative paths within our methodology that may result in higher performance and are worth exploring in future work. We list a few here. First, we choose to trim or pad tweets instead of processing tweets in batch sizes by tweet length or applying a pooling layer. Second, we do not experiment with character embeddings. Third, we do not add attention to our LSTM architecture despite its SOTA performance in other related NLP tasks. Instead, we depend on our experiments with feature embeddings to add a form of attention. Fourth, we do not conduct cross-validation or ensemble classification, even though it may augment training, because of time considerations. Finally, we collect only 200 tweets from user timelines instead of their full timelines due to rate-limiting restrictions. This research would benefit from collecting all tweets in the timeline in addition to crawling the timelines of all followers of all users in our dataset.

Additionally, aside from rate-limiting, we are limited by the data returned by Twitter's API. The reply tree study conducted by [46] in 2016 is no longer easily replicable, as Twitter only provides a subset of data. For example, when trying to collect replies from a tweet with 78 replies, only 15 were available through the API. According to Twitter:

Please note that Twitter's search service and, by extension, the Search API is not meant to be an exhaustive source of Tweets. Not all Tweets will be indexed or made available via the search interface.

Chapter 5

Results

5.1 Phase 1 Results

Table 5.1 illustrates Phase 1 results across our evaluation metrics. Logistic regression performs surprisingly well, which is in line with the results from [38]. The BERT embeddings perform the best across all models, with CNN-BERT being the best performing combination for both f-score on the hateful class and f-score overall. The MLP does not perform well at all. We hypothesize that this is because of the scarcity of hateful tweets in addition to the fact that our feature embeddings only consist of tweet embeddings; as a result the MLP and LSTM models do not have enough signals or cues to interpret the input. In contrast, we hypothesize that the CNN and DenseNet models perform well because, similar to how they handle images, they are interpreting patterns from an otherwise noisy input.

5.2 Phase 2 Results

Table 5.2 illustrates Phase 2 results across our evaluation metrics. While different model-embedding combinations see different performance boosts/losses from this round of experiments, we achieve our highest score with the CNN-BERT with network metrics as embedding coefficients, but only by a slight margin. We see the most benefit to TF-IDF embeddings for CNN and DenseNet, where the overall f-score is improved with both types of reply metrics. For pretrained Twitter and BERT embeddings, using reply metrics as embedding coefficients by passing the context tweets through a tweet-level version of our model did not help the final model in learning. MLP and LSTM were likely not affected much because similar to our hypothesis from our baselines,

Model	F _H	FA	F	Р	R
LR-TFIDF	0.13 ± 0.05	0.28 ± 0.12	0.52 ± 0.05	0.53 ± 0.07	0.57 ± 0.06
LR-TWIT	0.15 ± 0.01	0.63 ± 0.0	0.65 ± 0.01	0.65 ± 0.01	0.64 ± 0.01
LR-BERT	0.2 ± 0.0	0.74 ± 0.0	0.71 ± 0.0	0.71 ± 0.0	0.72 ± 0.0
MLP-TFIDF	0.0 ± 0.0	0.22 ± 0.16	0.05 ± 0.02	0.03 ± 0.01	0.18 ± 0.03
MLP-TWIT	0.08 ± 0.04	0.42 ± 0.01	0.32 ± 0.13	0.57 ± 0.03	0.32 ± 0.07
MLP-BERT	0.05 ± 0.04	0.24 ± 0.16	0.19 ± 0.18	0.33 ± 0.22	0.23 ± 0.16
CNN-TFIDF	0.1 ± 0.01	0.22 ± 0.09	0.38 ± 0.06	0.49 ± 0.06	0.33 ± 0.06
CNN-TWIT	0.27 ± 0.04	0.78 ± 0.02	0.68 ± 0.01	0.75 ± 0.01	0.65 ± 0.01
CNN-BERT	0.29 ± 0.03	0.75 ± 0.02	0.71 ± 0.03	0.76 ± 0.02	0.69 ± 0.03
LSTM-TFIDF	0.18 ± 0.01	0.62 ± 0.06	0.52 ± 0.03	0.67 ± 0.02	0.49 ± 0.03
LSTM-TWIT	0.22 ± 0.01	0.67 ± 0.02	0.62 ± 0.02	0.72 ± 0.01	0.58 ± 0.02
LSTM-BERT	0.25 ± 0.02	0.74 ± 0.01	0.67 ± 0.01	0.76 ± 0.0	0.64 ± 0.01
DENSE-TFIDF	0.08 ± 0.0	0.14 ± 0.1	0.27 ± 0.19	0.29 ± 0.21	0.26 ± 0.15
DENSE-TWIT	0.22 ± 0.01	0.74 ± 0.01	0.67 ± 0.01	0.7 ± 0.0	0.66 ± 0.01
DENSE-BERT	0.26 ± 0.01	0.73 ± 0.0	0.68 ± 0.01	0.72 ± 0.0	0.67 ± 0.01

Table 5.1: Phase 1 Experiment Results

F is f1 score, F_H indicates hateful class, and F_A indicates abusive class. *P* is precision and *R* is recall. The red highlight indicates best performing model and the blue highlight indicates the TF-IDF embeddings that see the most improvement with the inclusion of user behavior metrics.

Model	F _H	F_A	F	P	R
MLP-TFIDF	0.06 ± 0.04	0.11 ± 0.16	0.03 ± 0.03	0.02 ± 0.02	0.1 ± 0.07
MLP-TWIT	0.12 ± 0.01	0.36 ± 0.06	0.26 ± 0.14	0.44 ± 0.2	0.28 ± 0.09
MLP-BERT	0.07 ± 0.02	0.38 ± 0.01	0.4 ± 0.0	0.52 ± 0.02	0.38 ± 0.02
CNN-TFIDF	0.06 ± 0.04	0.23 ± 0.08	0.26 ± 0.14	0.38 ± 0.1	0.28 ± 0.05
CNN-TWIT	0.25 ± 0.01	0.77 ± 0.03	0.66 ± 0.03	0.75 ± 0.01	0.64 ± 0.03
CNN-BERT	0.29 ± 0.01	0.76 ± 0.02	0.7 ± 0.02	0.75 ± 0.01	0.68 ± 0.01
LSTM-TFIDF	0.19 ± 0.0	0.69 ± 0.0	0.56 ± 0.0	0.7 ± 0.0	0.53 ± 0.0
LSTM-Twit	0.13 ± 0.0	0.41 ± 0.0	0.5 ± 0.0	0.61 ± 0.0	0.46 ± 0.0
LSTM-BERT	0.22 ± 0.0	0.7 ± 0.0	0.65 ± 0.0	0.74 ± 0.0	0.61 ± 0.0
DENSE-TFIDF	0.11 ± 0.01	0.3 ± 0.16	0.48 ± 0.07	0.5 ± 0.06	0.48 ± 0.07
DENSE-TWIT	0.09 ± 0.01	0.31 ± 0.07	0.46 ± 0.03	0.55 ± 0.02	0.43 ± 0.04
DENSE-BERT	0.19 ± 0.01	0.69 ± 0.03	0.64 ± 0.02	0.7 ± 0.0	0.62 ± 0.02
MLP-TFIDF	0.0 ± 0.0	0.22 ± 0.16	0.2 ± 0.19	0.15 ± 0.16	0.34 ± 0.19
MLP-TWIT	0.06 ± 0.02	0.37 ± 0.02	0.14 ± 0.02	0.38 ± 0.21	0.23 ± 0.01
MLP-BERT	0.05 ± 0.0	0.35 ± 0.01	0.13 ± 0.03	0.29 ± 0.23	0.21 ± 0.02
CNN-TFIDF	0.1 ± 0.01	0.25 ± 0.07	0.41 ± 0.08	0.5 ± 0.05	0.39 ± 0.12
CNN-TWIT	0.26 ± 0.03	0.77 ± 0.01	0.66 ± 0.01	0.75 ± 0.01	0.63 ± 0.01
CNN-BERT	0.3 ± 0.01	0.76 ± 0.0	0.69 ± 0.01	0.76 ± 0.0	0.67 ± 0.01
LSTM-TFIDF	0.17 ± 0.01	0.64 ± 0.04	0.52 ± 0.03	0.67 ± 0.02	0.49 ± 0.03
LSTM-Twit	0.21 ± 0.01	0.65 ± 0.04	0.6 ± 0.03	0.71 ± 0.02	0.56 ± 0.02
LSTM-BERT	0.25 ± 0.0	0.73 ± 0.01	0.66 ± 0.01	0.75 ± 0.01	0.63 ± 0.01
DENSE-TFIDF	0.07 ± 0.02	0.29 ± 0.05	0.49 ± 0.06	0.53 ± 0.07	0.5 ± 0.1
DENSE-TWIT	0.22 ± 0.02	0.75 ± 0.01	0.67 ± 0.0	0.71 ± 0.02	0.65 ± 0.0
DENSE-BERT	0.25 ± 0.0	0.73 ± 0.01	0.68 ± 0.0	0.71 ± 0.0	0.67 ± 0.01

Table 5.2: Phase 2: Reply Pairings (Top) & Network Metrics (Bottom) Experiment Results

F is f1 score, F_H indicates hateful class, and F_A indicates abusive class. *P* is precision and *R* is recall. The red highlight indicates best performing model and the blue highlight indicates the TF-IDF embeddings that see the most improvement with the inclusion of user behavior metrics.

MODEL	F_H	F _A	F	Р	R
MLP-TFIDF	0.14 ± 0.02	0.47 ± 0.02	0.39 ± 0.04	0.58 ± 0.01	0.39 ± 0.03
MLP-TWIT	0.08 ± 0.04	0.36 ± 0.1	0.31 ± 0.16	0.41 ± 0.24	0.32 ± 0.11
MLP-BERT	0.06 ± 0.01	0.39 ± 0.02	0.35 ± 0.09	0.53 ± 0.03	0.34 ± 0.07
CNN-TFIDF	0.24 ± 0.01	0.82 ± 0.01	0.64 ± 0.02	0.77 ± 0.01	0.61 ± 0.02
CNN-TWIT	0.19 ± 0.07	0.71 ± 0.08	0.62 ± 0.06	0.71 ± 0.04	0.61 ± 0.05
CNN-BERT	0.3 ± 0.02	0.76 ± 0.01	0.68 ± 0.03	0.76 ± 0.01	0.65 ± 0.03
LSTM-TFIDF	0.2 ± 0.0	0.72 ± 0.01	0.59 ± 0.01	0.73 ± 0.0	0.55 ± 0.01
LSTM-TWIT	0.21 ± 0.01	0.66 ± 0.02	0.61 ± 0.02	0.72 ± 0.01	0.57 ± 0.02
LSTM-BERT	0.25 ± 0.01	0.73 ± 0.01	0.65 ± 0.01	0.75 ± 0.01	0.62 ± 0.01
DENSE-TFIDF	0.21 ± 0.0	0.76 ± 0.01	0.66 ± 0.01	0.72 ± 0.02	0.64 ± 0.01
DENSE-TWIT	0.19 ± 0.01	0.74 ± 0.02	0.69 ± 0.01	0.71 ± 0.0	0.67 ± 0.01
DENSE-BERT	0.24 ± 0.0	0.73 ± 0.0	0.68 ± 0.02	0.72 ± 0.0	0.66 ± 0.02

Table 5.3: Phase 3: User Behavior Metrics Experiment Results

F is f1 score, F_H indicates hateful class, and F_A indicates abusive class. *P* is precision and *R* is recall. The red highlight indicates best performing model and the blue highlight indicates the TF-IDF embeddings that see the most improvement with the inclusion of user behavior metrics.

the sparse amount of hateful samples coupled with the lack of cues or signal around the content of our tweets and reply tweets inhibits the model from learning. This tells us that it is not enough to feed our models more context: we need to find a better way to define that context within our feature embeddings.

5.3 Phase 3 Results

Figure 5.1 illustrates the topic keywords from our LDA model of user timeline tweets from a user that authored a sample hateful tweet. We see that our largest topic, with **16.1% of tweet tokens**, includes terms like 'potus', 'russian', 'nytime', and 'vote', implying that this user frequently tweets about US political issues. Another tweet annotated as hateful listed the dominant topic's top words as 'film', 'good', 'award', 'location', 'series', 'international'; upon further inspection of the tweet and the ac-



Figure 5.1: Topic LDA Modeling of User's Timeline from Sample Hateful Tweet

Sample hateful tweet is: "Muslims immigrants are living like they're used to where their from, squalor, filth and violence are all they know." This illustrates output from the Gensim library for LDA Topic modeling.

count, we found that the tweet was quoting a film and was mislabelled as hateful. In this case, topic modeling highlighted the difference between hateful users and hateful tweets by providing much needed context.

Table 5.3 illustrates phase 3 results across our evaluation metrics. We see that we achieve our best performing model on the hateful class, CNN-BERT, by using topic modeling features, but with only a slight improvement from our previous experiment phases. Figure 5.2 illustrates the model architecture of this best performing model. Notably, adding the user behavior metrics to the TF-IDF embedded models improves performance for the CNN and DenseNet classifiers, more than in phase 2 of our experiments. This tell us that the pretrained Twitter and BERT embeddings add a great deal of semantic value and cannot be much improved upon by user behavior metrics, such as topics from the user timeline. In contrast, TF-IDF embeddings do not offer this, and dominant topic words compensate for that.



Figure 5.2: Best Performing Model Architecture

CNN-BERT with User Behavior Metrics. Illustrates a multiple input model architecture where the BERT embeddings of annotated tweets and user timeline tweets are passed through distinct CNN models. The learned features are concatenated and passed as input to the final fully connected layer, which maps to the output.

5.4 Hyperparameter Tuning

We take our best performing model and embedding combination, CNN-BERT with user behavior metrics, and tune it by experimenting with dropout, number of layers, number of filters, and learning rate. Because the CNN-BERT model converged on average at **epoch 44**, we do not experiment with number of epochs. Figure 5.3 illustrates the 17 rounds of tuning experiments, with runtimes ranging from 5 hours and 20 minutes to 4 minutes (see Appendix A for machine configurations). We use the Bayes algorithm with the adaptive Parzen-Rosenblatt estimator (facilitated by the Comet.ML library's optimization framework [55]), which balances exploring unknown space with exploiting the known hyperparameter values that yield the best results. The following values perform best of the values explored: **num-filters=47**, **num-layers=1**, **dropout=0.8877**, **learning-rate=0.0196**. This leaves us with a single layer CNN with many filters and with a high level of regularization.





This chart illustrates the parameter choices and validation f-score on the hateful class across 17 experiments. Hyperparameter tuning was facilitated by the Comet.ML library's optimizer and visualization tools [55].

5.5 Tuned Model Results Across Experiments

Table 5.4 and Figure 5.5 illustrate the tuned model performances. See Appendix A.1 for additional results. All the tuned models demonstrate higher f-score on the minority class, with the tuned phase 2 reply-pairings model performing best at 0.33 ± 0.01 . This model also performs best on all classes, with an overall f-score of 0.72 ± 0.01 . It achieves the highest recall but not the highest precision of the experiments conducted so far. We hypothesize that the tuned phase 2 reply-pairings model performs better than the other models because the single layer CNN with many filters is overfitting with the BERT and LDA topic embeddings as well as with the reply network coefficients. With only the tweet and reply-pairing BERT embeddings, the network seems to have a more generalizeable notion of context.

In order to better understand our precision and recall scores, we analyze a confusion matrix in table 5.4. We see that the hateful class was correctly classified 40% of the

time and was most commonly confused with the abusive class. The abusive class was correctly classified 82% of the time, which is higher than the model's ability to detect normal tweets or spam tweets. In the below samples, the model predicts a tweet is hateful when it is labeled as hateful, abusive, and normal, respectively. Personally identifiable information is removed.

RHate when folk ask me questions that they already know the answers to. *B*TCH what you fishing for*? (Hateful Label)

That was the most patronising thing I've ever heard about young people and propaganda. What about old people and the mail?! (Abusive Label)

TEACHERS too! Another "trusted" member of society, doing the nasty with teen. (Normal Label)

There are a few problems with this. First, the tweet that is labelled as hateful does not appear to be hate speech from the definitions explored in our background. Second, the tweet that is labelled as normal could be considered abusive towards teachers. All of these tweets have negative sentiment, so it does appear our model has picked up on that. This deeper dive into the actual content that is misclassified hints at a problem with the annotation quality of our dataset.

5.6 Results Summary

The models that used Google's pretrained BERT embeddings performed better than TF-IDF and Twitter pretrained embeddings across most models in our three phases of experiments. CNN-BERT outperformed the logistic regression, MLP, LSTM, and DenseNet models for all three phases of experiments. Before tuning, our best performing model is the CNN multiple input model architecture with tweet and LDA topic words from user timeline tweets, all embedded with Google's pretrained BERT embeddings. After tuning, our best performing model is the CNN multiple input model architecture with tweet and reply pairings embedded using Google's pretrained BERT embeddings. We hypothesize that this is because the parameter choices of a single layer, 47 filter CNN with high dropout still leads to overfitting when additional context is added to the embeddings. This may be because the pretrained BERT embeddings add a lot of semantic information to begin with, which is demonstrated in our results where BERT embeddings are not much improved upon between experiment phases.

Model	F_H	F_A	F	Р	R
PHASE-1	0.32 ± 0.02	0.77 ± 0.01	0.71 ± 0.01	0.76 ± 0.01	0.69 ± 0.01
PHASE-2-1	0.33 ± 0.01	0.78 ± 0.01	0.72 ± 0.01	0.76 ± 0.0	0.7 ± 0.01
PHASE-2-2	0.32 ± 0.0	0.77 ± 0.01	0.7 ± 0.01	0.76 ± 0.0	0.69 ± 0.01
PHASE-3	0.32 ± 0.01	0.78 ± 0.01	0.7 ± 0.02	0.77 ± 0.01	0.68 ± 0.02

Table 5.4: Tuned CNN-BERT Model Performance

F is f1 score, F_H indicates hateful class, and F_A indicates abusive class. *P* is precision and *R* is recall.

Ultimately, we are able to improve on our logistic regression baseline performance on the hateful class by **.13** on a dataset with scarce hateful labels. If we were to randomly annotate a tweet as hateful with 4% probability, we'd achieve around 4% accuracy on the hateful class. Thus, we interpret the f-score of .33 as relatively high.

5.7 Experiment Difficulties

Below are a few experiment difficulties.

Nondeterministic evaluation error in university cluster's ability to load a saved model. This bug was identified when the test metrics had a standard deviation of .3+, which implied a bug in the metric computation, in sampling, in the experiment builder or in the dataset. After verifying the performance of each of those sections, the bug was isolated to the model load functionality. When the model was evaluated within training on the test set, the metrics were normal, but when the best performing model was loaded in using PyTorch, the abnormal metrics occurred in one or two out three experiment seeds. After verifying all parts of the models state dictionary and parameter weights, we guessed it was a hardware stack error, and deployed the experiment on the Google Cloud Platform (GCP) instead of the Machine Learning Practical cluster. The experiment ran a) significantly faster and b) with normal test metrics.

BERT embeddings. In order to expedite experiments, we do not load BERT embeddings from the Python library every time, because that is expensive in terms of time and computation. Instead we collect the BERT embeddings of our tweets (including our reply and user timeline tweets) 1,000 at a time, process the embeddings with

Chapter 5. Results

trimming or padding, and store a dictionary mappings tweet ids to embeddings.

Twitter data that has been taken down. Only 64% of the tweets in the original dataset could be recovered. For phase 2, about 16% of the instigator tweets could not be recovered. While this signals that potentially hateful or abusive content is being taken down rapidly (within a 9 month time window), it limits our ability to conduct research. A recommendation is for future researchers to provide the social media content they collect directly instead of references to it, although the retention of user data may violate some form of General Data Protection Regulation.



Figure 5.4: Confusion Matrix of Final Model.



Figure 5.5: Tuned CNN-BERT Experiment Results

All the tuned models demonstrate higher f-score on the minority class, with the tuned phase 2 model performing best at 0.33.

Chapter 6

Conclusion

6.1 Technical Contributions

Our final model successfully picked up on negative tweet sentiment and identified the abusive class at the highest rate, of 82% accuracy and 0.78 f-score. The model offers a significant improvement on detecting hate speech, as we are able to improve on our logistic regression baseline performance on the hateful class by .13 f-score on a dataset with scarce hateful labels. If we were to randomly annotate a tweet as hateful with 4% probability, we'd achieve around 4% accuracy on the hateful class. Thus, we interpret the final f-score on the hateful class of .33 as relatively high.

The models that used Google's pretrained BERT embeddings performed better than TF-IDF and Twitter pretrained embeddings across most models in our three phases of experiments. CNN-BERT outperformed the logistic regression, MLP, LSTM, and DenseNet models for all three phases of experiments. Before tuning, our best performing model is the CNN multiple input model architecture with tweet and user topic BERT embeddings. After tuning, our best performing model is the CNN multiple input model architecture with tweet and reply BERT embeddings. We hypothesize that this is because the parameter choices of a single layer, 47 filter CNN with high dropout will overfit with BERT and a large measure of user context. The pretrained BERT embeddings add enough semantic information to give us our most competitive models and adding additional metrics of context through aspects of the social network does not improve performance.

The task of automating the detection of hate speech on social media platforms remains a challenge, in part due to the difficulty in obtaining high-quality, large-scale annotated datasets and the scarce hateful samples available for machine learning models to learn from. Our experiments reflect this and suggest that improving the quality and consistency of annotations in our dataset is likely to result in more accurate automated systems. A promising area of future research is to use boosting methods, like oversampling or undersampling, in order to adjust the class distribution of our dataset.

6.2 Policy Recommendations

We now outline the key takeaways for policy and decision makers. First, if the task is completely automated it seems we risk removing posts that are abusive, normal, or spam which would look like a form of censorship on social media platforms. In this environment, the promotion of free speech may warrant keeping a human in the loop and having hateful content online for longer amounts of time. When content is published on the internet, especially controversial hate speech, it will likely be copied very quickly. Therefore, the speed at which a social media platform takes down the content does not guarantee that the content will not be made available elsewhere. It may be preferable to keep the human in the loop, especially if the detection of hate speech carries legal penalties, as is the case for certain European Union countries.

Second, topic modeling and other forms of user behavior metrics can be incorporated into a contagion tool. A contagion tool monitors trends and patterns and how they might spread through a network or change in response to triggering events. Topic modeling can provide helpful backgrounds on users and support the automated system's decision. While the inclusion of this in the feature embedding may not add much information to the model, it can add a level of interpretability to the decision by summarizing how coherent a user's timeline is around a topic and the words that make up a dominant topic. For example, we found that the author of a hate speech tweet in our dataset had a timeline defined by the topics of US politics and Russia. With this type of tool, researchers could track if hate speech from this user's follower Twitter network increases as we approach the US 2020 elections.

Third, given the massive amounts of content produced on social media platforms daily, it is not reasonable today to expect that a company can immediately detect and remove online hate speech. While increased scrutiny has furthered our capabilities, a "good" automated hate speech detector may mean detecting detecting 30% of hate speech and passing merely abusive content to human moderators.

The technical limitations described indicate that this problem will not be solved through automation alone; policy-makers must cooperate with social media companies in order to guide the handling of classified hateful and abusive content. This research is meant to inform interdisciplinary coalitions on the current landscape of hate speech detection. Ultimately, social media has contributed to a more open and connected world. It is critical that we mitigate the negative consequences of hate speech while preserving the benefits of online discussion.

Bibliography

- [1] D. Cave, "Countries want to ban 'weaponized' social media. what would that look like?" Mar 2019. [Online]. Available: https://www.nytimes.com/2019/03/ 31/world/australia/countries-controlling-social-media.html
- [2] S. Alava, D. Frau-Meigs, and G. Hassan, *Youth and violent extremism on social media: mapping the research.* UNESCO Publishing, 2017.
- [3] J. Wakefield, "Hate speech: Facebook, twitter and youtube told off by mps," Apr 2019. [Online]. Available: https://www.bbc.co.uk/news/technology-48037482
- [4] K. Müller and C. Schwarz, "Fanning the flames of hate: Social media and hate crime," *Available at SSRN 3082972*, 2018.
- [5] I. Awan, "Cyber-extremism: Isis and the power of social media," *Society*, vol. 54, no. 2, pp. 138–149, 2017.
- [6] M. ElSherief, S. Nilizadeh, D. Nguyen, G. Vigna, and E. Belding, "Peer to peer hate: Hate speech instigators and their targets," in *Twelfth International AAAI Conference on Web and Social Media*, 2018.
- [7] J. Cheng, M. Bernstein, C. Danescu-Niculescu-Mizil, and J. Leskovec, "Anyone can become a troll: Causes of trolling behavior in online discussions," in *Proceedings of the 2017 ACM conference on computer supported cooperative work and social computing*. ACM, 2017, pp. 1217–1230.
- [8] L. Deng and Y. Liu, *Deep Learning in Natural Language Processing*. Springer, 2018.
- [9] T. Marcin, "Facebook, youtube, whatsapp banned again sri lanka after violence against muslims," 2019. in May [Online]. Available: https://news.vice.com/en_us/article/597dm5/ facebook-youtube-whatsapp-banned-again-in-sri-lanka-after-violence-against-muslims

- [10] E. Stein, "History against free speech: The new german law against the" auschwitz": And other:" lies"," *Michigan Law Review*, vol. 85, no. 2, pp. 277– 324, 1986.
- "Exclusive: first. face-[11] M. Rosemain, In world а give data speech ... " 2019. [Online]. book to on hate Jun https://www.reuters.com/article/us-france-tech-exclusive/ Available: exclusive-in-a-world-first-facebook-to-give-data-on-hate-speech-suspects-to-french-courts-idl
- [12] E. Thomasson, "German cabinet agrees to fine somedia speech," [Online]. cial hate Apr 2017. Availover https://uk.reuters.com/article/uk-germany-hatecrime-facebook/ able: german-cabinet-agrees-to-fine-social-media-over-hate-speech-idUKKBN1771FK
- [13] "Twitter usage statistics." [Online]. Available: https://www.internetlivestats.com/ twitter-statistics/
- [14] A. M. Founta, C. Djouvas, D. Chatzakou, I. Leontiadis, J. Blackburn, G. Stringhini, A. Vakali, M. Sirivianos, and N. Kourtellis, "Large scale crowdsourcing and characterization of twitter abusive behavior," in *Twelfth International AAAI Conference on Web and Social Media*, 2018.
- [15] M. H. Ribeiro, P. H. Calais, Y. A. Santos, V. A. Almeida, and W. Meira Jr, "" like sheep among wolves": Characterizing hateful users on twitter," *arXiv preprint arXiv:1801.00317*, 2017.
- [16] H. Rainie, J. Q. Anderson, and J. Albright, *The future of free speech, trolls, anonymity and fake news online.* Pew Research Center Washington, DC, 2017.
- [17] A. Breuer, T. Landman, and D. Farquhar, "Social media and protest mobilization: Evidence from the tunisian revolution," *Democratization*, vol. 22, no. 4, pp. 764– 792, 2015.
- [18] S. J. Jackson, M. Bailey, and B. Foucault Welles, "# girlslikeus: Trans advocacy and community building online," *New Media & Society*, vol. 20, no. 5, pp. 1868– 1888, 2018.
- [19] R. Enikolopov, M. Petrova, and K. Sonin, "Social media and corruption," American Economic Journal: Applied Economics, vol. 10, no. 1, pp. 150–74, 2018.

- [20] R. Weintraub-Reiter, "Hate speech over the internet: A traditional constitutional analysis or a new cyber constitution," *BU Pub. Int. LJ*, vol. 8, p. 145, 1998.
- [21] L. Leets and H. Giles, "Words as weaponswhen do they wound? investigations of harmful speech," *Human Communication Research*, vol. 24, no. 2, pp. 260–301, 1997.
- [22] R. Cohen-Almagor, "Taking north american white supremacist groups seriously: The scope and the challenge of hate speech on the internet," *International journal of crime, justice, and social democracy*, vol. 7, no. 2, pp. 38–57, 2018.
- [23] E. Keen, "Mapping study on projects against hate speech online," Council of Europe, 2014.
- [24] I. Lunden and I. Lunden, "Facebook to add 3,000 to team reviewing posts with hate speech, crimes, and other harming posts," May 2017. [Online]. Available: https://techcrunch.com/2017/05/03/ facebook-to-hire-3000-to-review-posts-with-hate-speech-crimes-and-other-harming-posts/
- [25] S. Zannettou, B. Bradlyn, E. De Cristofaro, M. Sirivianos, G. Stringhini, H. Kwak, and J. Blackburn, "What is gab? a bastion of free speech or an altright echo chamber?" *arXiv preprint arXiv:1802.05287*, 2018.
- [26] Z. Waseem and D. Hovy, "Hateful symbols or hateful people? predictive features for hate speech detection on twitter," in *Proceedings of the NAACL student research workshop*, 2016, pp. 88–93.
- [27] P. Badjatiya, S. Gupta, M. Gupta, and V. Varma, "Deep learning for hate speech detection in tweets," in *Proceedings of the 26th International Conference on World Wide Web Companion*. International World Wide Web Conferences Steering Committee, 2017, pp. 759–760.
- [28] A. Schmidt and M. Wiegand, "A survey on hate speech detection using natural language processing," in *Proceedings of the Fifth International Workshop on Natural Language Processing for Social Media*, 2017, pp. 1–10.
- [29] D. Tang, F. Wei, B. Qin, T. Liu, and M. Zhou, "Coooolll: A deep learning system for twitter sentiment classification," in *Proceedings of the 8th international workshop on semantic evaluation (SemEval 2014)*, 2014, pp. 208–212.

- [30] M. Habibi, L. Weber, M. Neves, D. L. Wiegandt, and U. Leser, "Deep learning with word embeddings improves biomedical named entity recognition," *Bioinformatics*, vol. 33, no. 14, pp. i37–i48, 2017.
- [31] C. D. Santos and B. Zadrozny, "Learning character-level representations for partof-speech tagging," in *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, 2014, pp. 1818–1826.
- [32] T. Davidson, D. Warmsley, M. Macy, and I. Weber, "Automated hate speech detection and the problem of offensive language," in *Eleventh international aaai conference on web and social media*, 2017.
- [33] [Online]. Available: https://www.un.org/en/genocideprevention/documents/ UNStrategyandPlanofActiononHateSpeech18JuneSYNOPSIS.pdf
- [34] A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever, "Improving language understanding by generative pre-training," URL https://s3-us-west-2. amazonaws. com/openai-assets/researchcovers/languageunsupervised/language understanding paper. pdf, 2018.
- [35] Y. Mehdad and J. Tetreault, "Do characters abuse more than words?" in *Proceed-ings of the 17th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, 2016, pp. 299–303.
- [36] I. Rahwan, "What can we learn from emojis?" [Online]. Available: https://www.media.mit.edu/posts/what-can-we-learn-from-emojis/
- [37] G. K. Pitsilis, H. Ramampiaro, and H. Langseth, "Detecting offensive language in tweets using deep learning," *arXiv preprint arXiv:1801.04433*, 2018.
- [38] B. Gambäck and U. K. Sikdar, "Using convolutional neural networks to classify hate-speech," in *Proceedings of the first workshop on abusive language online*, 2017, pp. 85–90.
- [39] S. O. Sood, E. F. Churchill, and J. Antin, "Automatic identification of personal insults on social news sites," *Journal of the American Society for Information Science and Technology*, vol. 63, no. 2, pp. 270–285, 2012.
- [40] U. Sapkota, S. Bethard, M. Montes, and T. Solorio, "Not all character n-grams are created equal: A study in authorship attribution," in *Proceedings of the 2015*

conference of the North American chapter of the association for computational linguistics: Human language technologies, 2015, pp. 93–102.

- [41] C. Nobata, J. Tetreault, A. Thomas, Y. Mehdad, and Y. Chang, "Abusive language detection in online user content," in *Proceedings of the 25th international conference on world wide web*. International World Wide Web Conferences Steering Committee, 2016, pp. 145–153.
- [42] P. Nakov and J. Tiedemann, "Combining word-level and character-level models for machine translation between closely-related languages," in *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Short Papers-Volume 2.* Association for Computational Linguistics, 2012, pp. 301– 305.
- [43] G. Xiang, B. Fan, L. Wang, J. Hong, and C. Rose, "Detecting offensive tweets via topical feature discovery over a large scale twitter corpus," in *Proceedings of the* 21st ACM international conference on Information and knowledge management. ACM, 2012, pp. 1980–1984.
- [44] F. Del Vigna12, A. Cimino23, F. DellOrletta, M. Petrocchi, and M. Tesconi, "Hate me, hate me not: Hate speech detection on facebook," 2017.
- [45] P. Burnap and M. L. Williams, "Cyber hate speech on twitter: An application of machine classification and statistical modeling for policy and decision making," *Policy & Internet*, vol. 7, no. 2, pp. 223–242, 2015.
- [46] R. Nishi, T. Takaguchi, K. Oka, T. Maehara, M. Toyoda, K.-i. Kawarabayashi, and N. Masuda, "Reply trees in twitter: data analysis and branching process models," *Social Network Analysis and Mining*, vol. 6, no. 1, p. 26, 2016.
- [47] D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent dirichlet allocation," *Journal of machine Learning research*, vol. 3, no. Jan, pp. 993–1022, 2003.
- [48] J.-M. Xu, K.-S. Jun, X. Zhu, and A. Bellmore, "Learning from bullying traces in social media," in *Proceedings of the 2012 conference of the North American chapter of the association for computational linguistics: Human language technologies.* Association for Computational Linguistics, 2012, pp. 656–666.

- [49] G. Wiedemann, E. Ruppert, R. Jindal, and C. Biemann, "Transfer learning from lda to bilstm-cnn for offensive language detection in twitter," *arXiv preprint arXiv*:1811.02906, 2018.
- [50] B. Suh, L. Hong, P. Pirolli, and E. H. Chi, "Want to be retweeted? large scale analytics on factors impacting retweet in twitter network," in 2010 IEEE Second International Conference on Social Computing. IEEE, 2010, pp. 177–184.
- [51] F. Godin, B. Vandersmissen, W. De Neve, and R. Van de Walle, "Multimedia lab@ acl wnut ner shared task: Named entity recognition for twitter microposts using distributed word representations," in *Proceedings of the workshop on noisy* user-generated text, 2015, pp. 146–153.
- [52] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.
- [53] A. Karpathy, "Cs231n convolutional neural networks for visual recognition." [Online]. Available: http://cs231n.github.io/neural-networks-3/
- [54] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 4700–4708.
- [55] Comet.ml, "Comet.ml supercharging machine learning." [Online]. Available: https://www.comet.ml/

Appendix A

Experiments

A.1 Code

All code can be found at https://github.com/ashemag/hate_speech_detection

A.2 Machine Configuration

We use the Google Cloud Platform and create an instances with:

- 32 vCPUs, 208 GB memory
- 2 NVIDIA Tesla P4 GPUs
- 100 GB disk of the Deep-Learning-Pytorch Image

A.3 Additional Results

R_N	0.67 ± 0.03	0.68 ± 0.02	0.66 ± 0.02	0.65 ± 0.04	
P_N	0.88 ± 0.01	0.88 ± 0.01	0.88 ± 0.0	0.88 ± 0.01	
F_N	0.76 ± 0.02	0.77 ± 0.01	0.75 ± 0.01	0.75 ± 0.03	
R_A	0.8 ± 0.03	0.8 ± 0.01	0.81 ± 0.0	0.79 ± 0.04	
$P_{ m A}$	0.75 ± 0.04	0.76 ± 0.02	0.73 ± 0.02	0.77 ± 0.05	
F_A	0.77 ± 0.01	0.78 ± 0.01	0.77 ± 0.01	0.78 ± 0.01	
R_H	0.42 ± 0.05	0.42 ± 0.01	0.41 ± 0.03	0.44 ± 0.08	
P_H	0.26 ± 0.02	0.27 ± 0.02	0.27 ± 0.02	0.27 ± 0.03	
F_H	0.32 ± 0.02	0.33 ± 0.01	0.32 ± 0.0	0.32 ± 0.01	
MODEL	PHASE-1	PHASE-2-1	PHASE-2-2	PHASE-3	

Table A.2: Test set results for classes hateful, abusive, and normal.

Results
Experiment
A.1: Tuned
Table ,

Appendix A. Experiments



Figure A.1: Training vs. Validation Performances

Training f-score overall on and the hateful class is more consistent than on the validation set. This may be because our training set consists of balanced batches; even though our model is learning from the training set its unable to consistently classify the validation. This may be due to poor annotation quality.

Appendix B

Glossary

API: application program interface; allows applications to communicate with one another.

Convolutional layer: a mathematical operation that takes two inputs, such as an image and a filter, and learns features by using smaller parts of the input.

Cross validation: technique that reserves a sample of a dataset that the model is not trained on. Later, this sample is used for testing.

Crowdsourcing: the practice of obtaining information for a task by enlisting the services of a large number of people, typically via the Internet.

Feature embedding: the input fed into a machine learning model. It can consist of embeddings of different text features, like a tweet or topic. Embeddings are mappings of data into numerical vectors.

Fully connected layer: layer that connects every input neuron (i.e. operation) to outputs. We use this layer to map our learned feature representations to our final classes (i.e. hateful, abusive, normal, spam).

General Data Protection Regulation: a regulation in EU law (put into practice in 2018) on data protection and privacy for all individual citizens of the European Union (EU) and the European Economic Area (EEA).

Ground truth labels: the reality, or labels, we want our model to predict.

Hyperparameter: a parameter whose value is set before the learning process begins.

In-degree: in graph theory, this is the number of edges going into a node.

Logistic regression: a classification algorithm used to assign observations to classes. In our case, we use multiclass logistic regression.

Natural language processing (NLP): the process of a computer extracting information from natural language input and producing natural language output.

N-gram: a contiguous sequence of n items from a given sample of text or speech.

NumPy: package for scientific computing with Python.

Out-degree: in graph theory, this is the number of edges going out of a node.

Overfitting: model behavior where training corresponds too closely to data and thus may fail to fit additional data.

Personally identifiable information: data that could potentially identify an individual.

Pooling layer: combines several values into a single one. Helps the layer generalize because it effectively decreases the chance of overfitting.

PyTorch: a machine learning library based on the Torch library, for Python, used for applications such as deep learning and natural language processing.

Regularization: a way of penalizing model complexity to prevent overfitting.

Softmax layer: used to map the output of a network to a probability distribution over predicted classes. From this, we select label samples with the class with the high-

est probability.

Appendix C

Workshop Paper

The Benefits of BERT: A Survey of Deep Learning Approaches to Online Hate Speech Detection

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

Ashe Magalhaes

University of Edinburgh Edinburgh, Scotland a.m.magalhaes@ed.ac.uk

Abstract

The problem of automating the detection of online hate 1 speech is made more critical by the moral and legal con-2 3 sequences of its propagation in a social media network. This research surveys the current landscape of online 4 hate speech detection, deep learning model architec-5 tures, embedding choices, and the inclusion of user be-6 havior metrics, in order to improve performance on de-7 8 tecting the hateful class in a Twitter dataset. We find that Google's pretrained BERT embeddings outperform the 9 10 inclusion of user behavior metrics across many model choices and that a single-layer multiple input Convolu-11 tional Neural Network with many filters performs best. 12 The contribution of this paper is a range of experiments 13 that inform machine learning researchers on promising 14 15 future paths and policymakers on the capabilities and limitations of automated systems on the task of detect-16 ing hate speech. 17

Introduction

In recent years social media companies like Facebook, Twit-19 ter, and Alphabet Inc.'s YouTube have been widely criti-20 cized for failing to detect and remove hate speech from their 21 platforms (Cave 2019). In March 2019, a white supremacist 22 posted racist manifestos-which spread through Twitter-and 23 live-streamed the shooting of over 51 people on Face-24 book and YouTube. While both companies leverage machine 25 learning algorithms to automatically detect and remove such 26 content, they failed to take down the content quickly enough. 27 Facebook's head of public policy defended the platform's 28 slow response time: 29

- "The video was a new type that our machine learning
 system hadn't seen before. It was a first person shooter
 with a GoPro on his head...This is unfortunately an ad-
- versarial space. Those sharing the video were deliber-
- ately splicing and cutting it and using filters to subvert
- automation. There is still progress to be made with ma-
- chine learning (Wakefield 2019)."

18

Progress needs to be made sooner rather than later. Hateful content on social media contributes to real-world violence (Müller and Schwarz 2018), recruitment to and propaganda for terrorist individuals/groups (Awan 2017), makes

Copyright © 2019, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

Walid Magdy

University of Edinburgh Edinburgh, Scotland wmagdy@inf.ed.ac.uk

other users feel less safe and secure on social platforms
(ElSherief et al. 2018), and triggers increased levels of toxicity in the network (Cheng et al. 2017).

Governments are enforcing the need to find reasonable 44 solutions, fast. The Sri Lanken government temporarily 45 banned social media networks three separate times in the 46 wake of the Easter 2019 suicide bombings that killed hun-47 dreds of people, in order to prevent "social unrest via hate 48 messages" (Marcin 2019). French leader Emmanuel Macron 49 has made fighting online hate speech a priority, taking meet-50 ings with Mark Zuckerberg and other high-level Facebook 51 representatives; as a result Facebook has agreed to hand over 52 53 data on French users suspected of spewing hate speech on-54 line. This is an international first, and according to a counsel at law firm Linklaters, "a strong signal in terms of regula-55 tion" (Rosemain 2019). EU countries have threatened to fine 56 social networks up to 50 million euro per year if they con-57 tinue to fail to act faster (Thomasson 2017). 58

Social media has contributed to a more open and connected world (Rainie, Anderson, and Albright 2017). It has promoted Western liberal values through its effects on protest mobilization (Breuer, Landman, and Farquhar 2015), community building (Jackson, Bailey, and Foucault Welles 2018), and accountability in governments and institutions (Enikolopov, Petrova, and Sonin 2018). It is critical that we do not lose the benefits of these platforms as they operate under increased regulation and scrutiny. So, while automated systems can and should enforce uniform policies for user behavior that fosters welcoming online environments, this should not come at the cost of dissenting and fringe views.

This paper aims to navigate the challenge of creating systems that distinguish hateful content from abusive, normal, and spam content. We echo the concerns of Facebook's head of public policy by asking:

How can we improve the performance of automated systems on identifying hate speech when they must learn from few hateful samples?

This research surveys the capabilities and limitations of state-of-the-art (SOTA) deep learning model architectures with an emphasis on comparing the inclusion of user behavior features to using advanced pretrained word embeddings. We begin by describing related research and the model architectures which have shaped the current landscape of automated hate speech detection. We explain how this research

informs our methodology and go on to describe our data, 141 85 how it is collected and processed, the details of our model 86 142 architectures, and metrics for evaluation. 87

143 We find that Google's pretrained BERT embeddings add 88 144 enough semantic information to our model that the inclusion 89 145 of user behavior metrics does not offer further improve per-90 146 formance. 91 147

Related Work

148

149

150

151

152

153

154

155

156

161

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

Using Deep Learning Approaches 93

92

In spite of the growing phenomenon of hate speech, prior to 94 2016, NLP research on hate speech was limited (ElSherief 95 et al. 2018) (Waseem and Hovy 2016). Since then, statisti-96 cal machine learning techniques like logistic regression and 97 naive Bayes classifiers have been applied to this area with 98 success; logistic regression in particular is widely consid-99 ered a baseline model for experiments today (Waseem and 157 100 Hovy 2016) (Badjatiya et al. 2017) (Schmidt and Wiegand 158 101 2017). Recently, the aforementioned pressures, coupled with 102 159 concurrent advances in NLP and deep learning, has sparked 160 103 a resurgence in interest in the field. The success of deep neu-104 ral networks on NLP tasks like sentiment analysis (Tang et 105 162 al. 2014), named entity recognition (Habibi et al. 2017), and 106 107 part-of-speech tagging (Santos and Zadrozny 2014) have made it a natural choice to apply to the task of classifying 108 tweets or posts as hateful. 109

The baseline machine learning model choice, which tends 110 to perform well, is logistic regression (Davidson et al. 2017) 111 (Xiang et al. 2012). Another baseline model, support vec-112 tor machines (SVMs)-with linear kernels-perform similarly 113 to logistic regression in practice (Del Vignal et al. 2017). 114 Linear models tend not to capture the sparse features in doc-115 ument classification tasks, undermining the detection of re-116 lations between words or symbols in sentences. CNNs and 117 RNNs (particularly LSTMs) have achieved SOTA perfor-118 mance due to their ability to capture long-term dependencies 119 in a sentence, tweet, or post through their ability to maintain 120 information in memory for a period of time (Badjatiya et al. 121 2017) (Pitsilis, Ramampiaro, and Langseth 2018). 122

The Difficulty with Hate Speech Datasets 123

A connected area of research is the task of creating robust 124 hate speech datasets. The common practice is to deploy a 125 collection of tweets or posts to a crowdsourcing site, where 126 a number of human annotators identify an appropriate clas-127 sification (i.e. hateful, normal, abusive, spam). This is prob-128 lematic for a few reasons and highlights the ways in which 129 human-annotated data reflects social biases. First, the sub-130 jective nature of identifying hate speech makes it difficult 131 to extract ground truth labels for tweets or posts (Founta et 132 al. 2018). There are many definitions of hate speech, which 133 vary across different languages, cultures, and governments. 134 Many of the definitions involve the hateful targeting of mem-135 bers of a group, like this one: 136

"Language used to express hatred towards a targeted 137 individual or group, or is intended to be derogatory, to 138

humiliate, or to insult the members of the group, on the 139 basis of attributes such as race, religion, ethnic origin, 140

sexual orientation, disability, or gender." (Davidson et al. 2017)

Second, because there is no universally agreed upon definition of hate speech [18], it is often difficult for human annotators to separate hate speech from offensive language. Given the legal consequences of hate speech, it is important that annotators and machine learning classifiers distinguish between hateful language and offensive language (Founta et al. 2018) (Davidson et al. 2017). It has been found that showing users a definition of hate speech does not improve annotation reliability, hinting that the presence of hate speech is not a binary decision and annotators likely need more context (Schmidt and Wiegand 2017), such as sample post-label pairings.

Third, it is likely that annotators skim through tweets or posts too quickly and fail to pick up on the context of the language. Past research demonstrates that annotators will incorrectly label normal text as hate speech if it contains hate or curse words (Davidson et al. 2017). Annotators are likely to miss out on hate speech without a slur as well as interpret racist and homophobic slurs as hateful but sexist slurs as 'just offensive' (Waseem and Hovy 2016) (Davidson et al. 2017).

The Inclusion Of User Behavior Metrics

Hateful users tend to make themselves central to the Twitter social network: they target more popular and popular users and they have more statuses and followees per day, although they have younger accounts (ElSherief et al. 2018) (Ribeiro et al. 2017). While this may increase the user's visibility, the inclusion of hateful or antagonistic content in a tweet reduces the rate of retweet by a factor of 45 (Burnap and Williams 2015). Despite this, hateful users are central to Twitter reply networks as determined by features such as betweenness and eigenvector centrality (Ribeiro et al. 2017). Furthermore, the modeling of reply trees-rooted at a tweet and joined by replies-reveals that the in-degree of the tweet, or the number of retweets, plays an important role in the overall shape of the reply tree and a component of the Twitter reply network as a whole (Nishi et al. 2016).

Latent Dirichlet Allocation (LDA) (Blei, Ng, and Jordan 2003) is a method of topic modeling that extracts the main topics in text. It has been used to improve the performance of linear machine learning classifiers on the task of identifying abusive social media content (Xu et al. 2012). While LDA-derived topics have been used as annotations (Wiedemann et al. 2018), to the best of our knowledge it has not been included in feature embeddings in order to improve the performance of deep learning classifiers on the task.

Methodology

The overall goal of this research is to conduct a survey of deep learning model architectures, embedding choices, and feature inputs in order to automate the identification of hateful tweets from a dataset with few hateful samples. Unlike previous work, we experiment with user behavior metrics in multiple input model architectures to provide our deep



Figure 1: Word clouds of tweets from the hateful, abusive, normal, and spam classes.

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

277

278

279

280

281

282

283

284

285

286

287

288

289

learning classifiers with more context to identify hate speechwith.

198 Dataset

The dataset for this research comes from Founta et al.'s work 199 that describes the process of large scale crowdsourcing for 200 annotations of hateful, normal, abusive, and spam tweets 201 (Founta et al. 2018). Founta et al. provide status ids for oth-202 ers to obtain tweets through the Twitter API. After learning 203 that a number of tweets were no longer available on Twitter, 204 they provided 100k tweets' status ids and their associated 205 majority labels. Upon using the Twitter API, we find that 206 only 64% of those tweets are still available online; we as-207 sume the remainder have been taken down by Twitter or the 208 users. The class distribution is as follows: 4% hateful, 20% 209 abusive, 62% normal, and 14% spam; there are 64,149 210 tweets total. Half the tweets achieved annotation majority 211 with an agreement of more than 3 out of 5 votes, which 212 should give us some confidence in the annotations. Figure 213 1 illustrates a word cloud of each class. 214

The average retweet count of hateful content is 15% higher than the average retweet count of normal content; the average favorite count of hateful content is 76% higher than the average favorite count of normal content. Finally, on average hateful content is a reply to existing tweets 31% more often than normal content is.

For each tweet, if it is in reply to a tweet we use Twitter's API to collect the context tweet. 19% of the total tweets are a reply. Of those, 14% are no longer available on Twitter. For each tweet, we crawl the user timeline of the author and collect their 200 most recent tweets, giving us around 12.8 million additional tweets for LDA topic modeling.

227 Data Processing

We clean the tweets by tokenizing, lowercasing, and removing punctuation. We keep URLs because they provide important context; tweets with URLs are more likely to get retweeted (Suh et al. 2010). We experiment with three types of text embeddings:

- 1. Term Frequency-Inverse Document Frequency (TF-IDF) embeddings. TF-IDF fits the training set into a weighted vector by normalized frequency of the 10,000 most common words in our vocabulary. Our validation and test set are transformed using the learned weights.
- 2. **Pretrained Twitter embeddings.** These embeddings are from a Word2Vec model trained on 400 million raw English tweets, with an embedding dimension of 400 (Godin et al. 2015).
- 3. **Pretrained BERT embeddings.** Google's BERT, or Bidirectional Encoder Representations from Transforms, is a novel method of pre-training language representations which obtains SOTA results on a range of NLP tasks devlin2018BERT. Specifically, we use the BERT-12-768-12 model with an embedding dimension of 768.

For our word embeddings, we trim or pad each tweet to 17 words, the average size of a tweet in our dataset. When padding, we add randomly generated embeddings so that each tweet embedding is a n-gram, where n=17. To account for the imbalance of classes in our dataset, for our training experiments we draw samples through a multinomial distribution that accounts for each class weight. Our balanced batches allow the classifier to learn class representations more consistently each epoch. We do not do this for validation and testing because a) we run the risk of tweet samples that are seen multiple times or never in our evaluation and b) it is not a realistic representation of tweets in the wild.

Model Architectures

Logistic Regression Our baseline model is multi-class logistic regression which uses the L2 norm penalty as a form of regularization.

Multilayer Perceptron. Known as a "vanilla neural network", this model consists of fully connected layers with the non-linear activation function leaky ReLu applied. We use Leaky ReLus here and in other models because it allows a small, positive gradient when the unit is not active and as a result, tends to perform better than nonleaky ReLu in practice.

Convolutional Neural Network (CNN). While this model is traditionally used for image inputs, it has performed well on NLP tasks and has had success when applied to hate speech detection (Gambäck and Sikdar 2017). The convolutional layer computes the output of neurons (i.e. mathematical operations) that are connected to local regions in the input. It consists of a set of learnable filters that produce an activation map used to ultimately compute the class scores. Our CNN consists of a series of convolutional layers, batch normalization, leaky ReLu, and dropout. The output is processed through a max-pooling layer and then a fully connected layer which computes class probabilities. Our CNN contains a context list that gives every layer access to all the previous layers, allowing for more connectivity and greater chance of the network learning from long-term dependencies in the tweet.

DenseNet. CNNs are more accurate and efficient if they contain shorter connections between layers close to the input and output (Huang et al. 2017). DenseNet leverages this

observation by connecting every layer to every other layer: 290 for each layer, the feature maps of all preceding layers are 291 used as inputs into all subsequent layers. Our DenseNet im-292 plementation has achieved SOTA results for image datasets; 293 we experiment with its application to NLP. 294

LSTM. LSTMs (and RNNs more broadly) improve on 295 traditional deep neural networks by propagating informa-296 tion in both directions. This loop in the architecture acts as 297 a memory state by which the network makes adjustments in 298 the information flow, allowing the model to selectively re-299 member or forget. This model has performed well on NLP 300 tasks in practice because of its ability to remember depen-301 302 dencies and between sentences; it is applied to the task of hate speech detection by (Badjatiya et al. 2017) (Del Vignal 303 et al. 2017) (Gambäck and Sikdar 2017). Our implementa-304 tion passes our input through a series of LSTM layers to a 305 final fully connected layer. 306

Experiment

345

346

347

348

349

357

358

359

360

361

362

363

364

365

366

367

368

369

376

378

Phase 1: Tweet Embedding 308

307

We compare our deep learning models to a baseline logistic 309 350 regression model to begin with an idea around how much of 310 an improvement they can offer. Here the feature embedding 311 to our deep learning models are simply the tweet embed-351 312 dings. Table 1 enumerates regularization choices, number 313 of layers, and model-specific hyperparameter values. Reg-352 314 ularization reduces the number of parameters and amount of 353 315 computation in the network, which helps to control overfit-354 316 ting. 355 317 356

Phase 2: Tweet + Reply Embeddings 318

The goal of this round of experiments is to apply some type 319 of context to our embeddings. First, because of the statistics 320 collected around the behavior of hateful users and retweet-321 ing, we define pairs of tweets: the original tweet and a con-322 text tweet. The context tweet is for the case that the tweet 323 was a response to something else. If the tweet is not a re-324 sponse to something else, we add null embeddings to indi-325 cate that there is no context for the tweet. Although only 326 19% of our tweets in the dataset are replies, the extra infor-327 mation around context (when it is there and not there) may 328 help our network learn in some way. 329

Next, because tweet in-degree has been shown to be sig-330 nificant ((Nishi et al. 2016), see Background) we focus 331 on in-degree in terms of number of times someone has 332 370 retweeted a given tweet and number of times someone has 333 favorited a given tweet. Our maximum retweet value is 371 334 154,565 and our maximum favorite value is 27,741. Be-335 372 cause these numbers are so large, we log each (and use 0 336 373 when the count is 0), so as to not skew the learned weights. 337 374 We concatenate the logged retweet and favorite counts to 375 338 our tweet and reply embeddings. We are interested to see if 339 the network can better learn from the retweet and favorite 340 377 numbers as a global measure of context. 341

Here and for the remaining experiments, we shift to build-342 379 ing a neural network with multiple inputs in order to have the 343 380 network learn from the annotated tweet in addition to other 381 344



Figure 2: Multiple Input CNN Architecture

types of embeddings. We have the annotated tweet embedding and the context tweet embedding as separate inputs; they are processed by different parts of the model architecture. The learned features for each are concatenated and fed into a final fully connected layer. See Figure 2 for a visualization of a multiple input network architecture.

Phase 3: Tweet + Dominant Topic Embeddings

This phase aims to add context in a more sophisticated way. For each tweet in our dataset, we crawl the author's user timeline and collect 200 tweets. We then conduct topic modeling through the LDA approach (discussed in Background). For each of our users, we find the dominant topic in their timeline tweets and create embeddings of a sequence of the top 10 words related to the topic. We choose the same embedding that is used to embed our tweets for a given experiment. We also concatenate the coherence and perplexity scores of the user's timeline to each embedded topic word in order to add a global measure of topic modeling. Perplexity is a measure of how well a probability model can predict a topic (lower is better) and coherence captures how well topics can be defined (higher is better).

Our multiple input model architecture processes the topic embeddings and tweet embeddings. Again, we concatenate the learned features and use this as input to the final fully connected layer.

Hyperparameter Tuning & Evaluation

We tune our best performing models-by model type and embedding type-by experimenting with learning rate, regularization, number of layers, and the model specific parameter. We define best performing model as the model with the greatest validation f-score on the hateful classes. We evaluate our results by overall metrics of f-score, precision, and recall as well as f-score of the hateful and abusive classes. This is because we are concerned with our classifier's ability to correctly identify hateful and abusive samples. To achieve stability in the results, we evaluate the mean and standard deviation of 3 experiments each.

MODEL	REGULARIZATION	LAYERS
LR	L2-NORM PENALTY	-
MLP	MP, $\lambda = 1E-4$	3 FC-LERELU-FC
CNN	MP, BN, δ =0.5, λ =1E-4	3 CONV-BN-LERELU-FC
DENSENET	AP, BBN, δ =0.5, λ =1E-4	4 DENSE-BBN-RELU-CONV-FC
LSTM	δ =0.5, λ =1E-4	3 LSTM-FC

Table 1: Baseline model architectures. Terms include Weight decay (λ), Max Pooling (MP), Average Pooling (AP), Dropout (δ), Batch Normalization (BN), and Bottleneck Layer, which includes BN (BBN). Model-specific hyperparameters include CNN filters=8, DENSENET growth-rate=12, LSTM hidden-layers=5.

426

427

428

429

430

431

432

433

434

435

436

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

Limitations 382

398

399

417

We acknowledge that there are alternative experiment design 383 384 choices within our methodology that are worth exploring in future work. First, we choose to trim or pad tweets instead 385 of processing tweets in batch sizes by tweet length or by ap-386 plying a pooling layer. Second, we do not experiment with 387 character embeddings. Third, we are limited by the data re-388 turned by Twitter's publicly available API. The reply tree 389 study conducted by (Nishi et al. 2016) in 2016 is no longer 390 easily replicable, as Twitter only provides a subset of data. 391 For example, when trying to collect replies from a tweet with 392 78 replies, only 15 were available through the API. Accord-393 ing to Twitter: *Please note that Twitter's search service and*, 394 437 by extension, the Search API is not meant to be an exhaus-395 tive source of Tweets. Not all Tweets will be indexed or made 396

available via the search interface. 397

Results

Phase 1 Results

Table 2 illustrates Phase 1 results across our evaluation met-400 rics. Randomly classifying a sample as hate speech with the 401 distribution found in our dataset results in 4% classifica-402 tion accuracy. Relative to this, logistic regression is a strong 403 baseline with 0.20 f-score on the hateful class. This is in 404 line with the results from (Gambäck and Sikdar 2017). The 405 BERT embeddings perform the best across all models, with 406 CNN-BERT being the best performing combination for both 407 f-score on the hateful class (0.29) and f-score overall (0.71). 408 The MLP does not perform well; we hypothesize that this is 409 because of the scarcity of hateful tweets in addition to the 410 fact that our feature embeddings only consist of tweet em-411 bedddings. As a result, the MLP and LSTM models do not 412 have enough signals or cues to interpret the input. In con-413 trast, we hypothesize that the CNN and DenseNet models 414 perform well because, similar to how they handle images, 415 they are interpreting patterns from an otherwise noisy input. 416

Phase 2 Results

Table 3 illustrates Phase 2 results across our evaluation met-418 rics. While different model-embedding combinations see 419 different performance boosts/losses from this round of ex-420 periments, we achieve our highest score with the CNN-421 BERT with network metrics as embedding coefficients, but 422 only by a slight margin. We see the most benefit to TF-423 IDF embeddings for CNN and DenseNet, where the over-424 all f-score is improved with both types of reply metrics. For 425

pretrained Twitter and BERT embeddings, using reply metrics as embedding coefficients by passing the context tweets through a tweet-level version of our model did not help the final model in learning. MLP and LSTM were likely not affected much because similar to our hypothesis from phase 1, the sparse amount of hateful samples coupled with the lack of cues or signal around the content of our tweets and reply tweets inhibits the model from learning. This tells us that it is not enough to feed our models more context: we need to find a better way to define that context within our feature embeddings.

Phase 3 Results

Figure 3 illustrates the topic keywords from our LDA model of user timeline tweets from a user that authored a sample hateful tweet. We see that our largest topic, with 16.1% of tweet tokens, includes terms like 'potus', 'russian', 'nytime', and 'vote', implying that this user frequently tweets about US political issues. Another tweet annotated as hateful listed the dominant topic's top words as 'film', 'good', 'award', 'location', 'series', 'international'; upon further inspection of the tweet and the account, we found that the tweet was quoting a film and was mislabelled as hateful. In this case, topic modeling highlighted the difference between hateful users and hateful tweets by providing much needed context.

Table 4 illustrates phase 3 results across our evaluation metrics. We see that we achieve our best performing model on the hateful class, CNN-BERT, by using topic modeling features, but with only a slight improvement from our previous experiment phases. Figure 2 illustrates this multiple input model architecture. Notably, adding the user behavior metrics to the TF-IDF embedded models improves performance for the CNN and DenseNet classifiers, more than in phase 2 of our experiments. This tell us that the pretrained Twitter and BERT embeddings add a great deal of semantic value and cannot be much improved upon by user behavior metrics, such as topics from the user timeline. In contrast, TF-IDF embeddings do not offer this, and dominant topic words compensate for that.

Hyperparameter Tuning

We take our best performing model and embedding combination, CNN-BERT with user behavior metrics, and tune it by experimenting with dropout, number of layers, number of filters, and learning rate. Because the CNN-BERT

MODEL	F_H	F_A	F	P	R
L R TEIDE	0.13 ± 0.05	0.28 ± 0.12	0.52 ± 0.05	0.53 ± 0.07	0.57 ± 0.06
LR-ITIDI LD Tuur	0.15 ± 0.05	0.20 ± 0.12	0.52 ± 0.05	0.55 ± 0.07	0.37 ± 0.00
LR-IWIT	0.15 ± 0.01	0.63 ± 0.0	0.65 ± 0.01	0.65 ± 0.01	0.64 ± 0.01
LR-BERT	0.2 ± 0.0	0.74 ± 0.0	0.71 ± 0.0	0.71 ± 0.0	0.72 ± 0.0
MLP-TFIDF	0.0 ± 0.0	0.22 ± 0.16	0.05 ± 0.02	0.03 ± 0.01	0.18 ± 0.03
MLP-TWIT	0.08 ± 0.04	0.42 ± 0.01	0.32 ± 0.13	0.57 ± 0.03	0.32 ± 0.07
MI P-BERT	0.05 ± 0.04	0.24 ± 0.16	0.19 ± 0.18	0.33 ± 0.22	0.23 ± 0.16
	0.05 ± 0.04	0.24 ± 0.10	0.17 ± 0.10	0.55 ± 0.22	0.25 ± 0.10
CNN-TFIDF	0.1 ± 0.01	0.22 ± 0.09	0.38 ± 0.06	0.49 ± 0.06	0.33 ± 0.06
CNN-TWIT	0.27 ± 0.04	0.78 ± 0.02	0.68 ± 0.01	0.75 ± 0.01	0.65 ± 0.01
CNN-BERT	0.29 ± 0.03	0.75 ± 0.02	0.71 ± 0.03	0.76 ± 0.02	0.69 ± 0.03
LSTM-TFIDF	0.18 ± 0.01	0.62 ± 0.06	0.52 ± 0.03	0.67 ± 0.02	0.49 ± 0.03
LSTM-TWIT	0.22 ± 0.01	0.67 ± 0.02	0.62 ± 0.02	0.72 ± 0.01	0.58 ± 0.02
LSTM-BERT	0.25 ± 0.02	0.74 ± 0.01	0.67 ± 0.01	0.76 ± 0.0	0.64 ± 0.01
					0.06 1.0.15
DENSE-TFIDF	0.08 ± 0.0	0.14 ± 0.1	0.27 ± 0.19	0.29 ± 0.21	0.26 ± 0.15
DENSE-TWIT	0.22 ± 0.01	0.74 ± 0.01	0.67 ± 0.01	0.7 ± 0.0	0.66 ± 0.01
DENSE-BERT	0.26 ± 0.01	0.73 ± 0.0	0.68 ± 0.01	0.72 ± 0.0	0.67 ± 0.01

Table 2: Phase 1 Experiment Results, Tweet Embeddings. F is f1 score, F_H indicates hateful class, and F_A indicates abusive class. P is precision and R is recall. The red highlight indicates best performing model and the blue highlight indicates the TF-IDF embeddings that see the most improvement with the inclusion of user behavior metrics.

MODEL	F_H	F_A	F	P	R
MLP-TFIDF MLP-Twit MLP-BERT	$ \begin{vmatrix} 0.06 \pm 0.04 \\ 0.12 \pm 0.01 \\ 0.07 \pm 0.02 \end{vmatrix} $		$\begin{array}{c} 0.03 \pm 0.03 \\ 0.26 \pm 0.14 \\ 0.4 \pm 0.0 \end{array}$	$\begin{array}{c} 0.02 \pm 0.02 \\ 0.44 \pm 0.2 \\ 0.52 \pm 0.02 \end{array}$	$ \begin{vmatrix} 0.1 \pm 0.07 \\ 0.28 \pm 0.09 \\ 0.38 \pm 0.02 \end{vmatrix} $
CNN-TFIDF CNN-Twit CNN-BERT	$\begin{array}{c} 0.06 \pm 0.04 \\ 0.25 \pm 0.01 \\ 0.29 \pm 0.01 \end{array}$	$\begin{array}{c} 0.23 \pm 0.08 \\ 0.77 \pm 0.03 \\ 0.76 \pm 0.02 \end{array}$	$\begin{array}{c} 0.26 \pm 0.14 \\ 0.66 \pm 0.03 \\ 0.7 \pm 0.02 \end{array}$	$\begin{array}{c} 0.38 \pm 0.1 \\ 0.75 \pm 0.01 \\ 0.75 \pm 0.01 \end{array}$	$\begin{array}{c} 0.28 \pm 0.05 \\ 0.64 \pm 0.03 \\ 0.68 \pm 0.01 \end{array}$
LSTM-TFIDF LSTM-Twit LSTM-BERT	$\begin{array}{c} 0.19 \pm 0.0 \\ 0.13 \pm 0.0 \\ 0.22 \pm 0.0 \end{array}$	$\begin{array}{c} 0.69 \pm 0.0 \\ 0.41 \pm 0.0 \\ 0.7 \pm 0.0 \end{array}$	$\begin{array}{c} 0.56 \pm 0.0 \\ 0.5 \pm 0.0 \\ 0.65 \pm 0.0 \end{array}$	$\begin{array}{c} 0.7 \pm 0.0 \\ 0.61 \pm 0.0 \\ 0.74 \pm 0.0 \end{array}$	$ \begin{vmatrix} 0.53 \pm 0.0 \\ 0.46 \pm 0.0 \\ 0.61 \pm 0.0 \end{vmatrix} $
DENSE-TFIDF DENSE-Twit DENSE-BERT	$\begin{array}{c} 0.11 \pm 0.01 \\ 0.09 \pm 0.01 \\ 0.19 \pm 0.01 \end{array}$	$\begin{array}{c} 0.3 \pm 0.16 \\ 0.31 \pm 0.07 \\ 0.69 \pm 0.03 \end{array}$	$\begin{array}{c} 0.48 \pm 0.07 \\ 0.46 \pm 0.03 \\ 0.64 \pm 0.02 \end{array}$	$\begin{array}{c} 0.5 \pm 0.06 \\ 0.55 \pm 0.02 \\ 0.7 \pm 0.0 \end{array}$	$\begin{array}{c} 0.48 \pm 0.07 \\ 0.43 \pm 0.04 \\ 0.62 \pm 0.02 \end{array}$
MLP-TFIDF MLP-Twit MLP-BERT	$\begin{array}{c} 0.0 \pm 0.0 \\ 0.06 \pm 0.02 \\ 0.05 \pm 0.0 \end{array}$	$\begin{array}{c} 0.22 \pm 0.16 \\ 0.37 \pm 0.02 \\ 0.35 \pm 0.01 \end{array}$	$\begin{array}{c} 0.2 \pm 0.19 \\ 0.14 \pm 0.02 \\ 0.13 \pm 0.03 \end{array}$	$\begin{array}{c} 0.15 \pm 0.16 \\ 0.38 \pm 0.21 \\ 0.29 \pm 0.23 \end{array}$	$\begin{array}{c} 0.34 \pm 0.19 \\ 0.23 \pm 0.01 \\ 0.21 \pm 0.02 \end{array}$
CNN-TFIDF CNN-Twit CNN-BERT	$\begin{array}{c} 0.1 \pm 0.01 \\ 0.26 \pm 0.03 \\ 0.3 \pm 0.01 \end{array}$	$\begin{array}{c} 0.25 \pm 0.07 \\ 0.77 \pm 0.01 \\ 0.76 \pm 0.0 \end{array}$	$\begin{array}{c} 0.41 \pm 0.08 \\ 0.66 \pm 0.01 \\ 0.69 \pm 0.01 \end{array}$	$\begin{array}{c} 0.5 \pm 0.05 \\ 0.75 \pm 0.01 \\ 0.76 \pm 0.0 \end{array}$	$\begin{array}{c} 0.39 \pm 0.12 \\ 0.63 \pm 0.01 \\ 0.67 \pm 0.01 \end{array}$
LSTM-TFIDF LSTM-Twit LSTM-BERT			$\begin{array}{c} 0.52 \pm 0.03 \\ 0.6 \pm 0.03 \\ 0.66 \pm 0.01 \end{array}$		$ \begin{vmatrix} 0.49 \pm 0.03 \\ 0.56 \pm 0.02 \\ 0.63 \pm 0.01 \end{vmatrix} $
DENSE-TFIDF DENSE-Twit DENSE-BERT	$\begin{array}{c} 0.07 \pm 0.02 \\ 0.22 \pm 0.02 \\ 0.25 \pm 0.0 \end{array}$	$\begin{array}{c} 0.29 \pm 0.05 \\ 0.75 \pm 0.01 \\ 0.73 \pm 0.01 \end{array}$	$\begin{array}{c} 0.49 \pm 0.06 \\ 0.67 \pm 0.0 \\ 0.68 \pm 0.0 \end{array}$	$\begin{array}{c} 0.53 \pm 0.07 \\ 0.71 \pm 0.02 \\ 0.71 \pm 0.0 \end{array}$	$\begin{array}{c} 0.5 \pm 0.1 \\ 0.65 \pm 0.0 \\ 0.67 \pm 0.01 \end{array}$

Table 3: Phase 2 Experiment Results, Tweet + Reply (Top) & Network Metrics (Bottom) Embeddings. F is f1 score, F_H indicates hateful class, and F_A indicates abusive class. P is precision and R is recall. The red highlight indicates best performing model and the blue highlight indicates the TF-IDF embeddings that see the most improvement with the inclusion of user behavior metrics.

MODEL	F_H	F_A	F	<i>P</i>	R
MLP-TFIDF MLP-Twit MLP-BERT		$\begin{array}{c} 0.47 \pm 0.02 \\ 0.36 \pm 0.1 \\ 0.39 \pm 0.02 \end{array}$	$ \begin{vmatrix} 0.39 \pm 0.04 \\ 0.31 \pm 0.16 \\ 0.35 \pm 0.09 \end{vmatrix} $	$ \begin{vmatrix} 0.58 \pm 0.01 \\ 0.41 \pm 0.24 \\ 0.53 \pm 0.03 \end{vmatrix} $	
CNN-TFIDF CNN-Twit CNN-BERT	$\begin{array}{c} 0.24 \pm 0.01 \\ 0.19 \pm 0.07 \\ 0.3 \pm 0.02 \end{array}$	$\begin{array}{c} 0.82 \pm 0.01 \\ 0.71 \pm 0.08 \\ 0.76 \pm 0.01 \end{array}$	$\begin{array}{c} 0.64 \pm 0.02 \\ 0.62 \pm 0.06 \\ 0.68 \pm 0.03 \end{array}$	$\begin{array}{c} 0.77 \pm 0.01 \\ 0.71 \pm 0.04 \\ 0.76 \pm 0.01 \end{array}$	$\begin{array}{c} 0.61 \pm 0.02 \\ 0.61 \pm 0.05 \\ 0.65 \pm 0.03 \end{array}$
LSTM-TFIDF LSTM-Twit LSTM-BERT	$\begin{array}{c} 0.2 \pm 0.0 \\ 0.21 \pm 0.01 \\ 0.25 \pm 0.01 \end{array}$	$\begin{array}{c} 0.72 \pm 0.01 \\ 0.66 \pm 0.02 \\ 0.73 \pm 0.01 \end{array}$	$ \begin{vmatrix} 0.59 \pm 0.01 \\ 0.61 \pm 0.02 \\ 0.65 \pm 0.01 \end{vmatrix} $	$\begin{array}{c} 0.73 \pm 0.0 \\ 0.72 \pm 0.01 \\ 0.75 \pm 0.01 \end{array}$	$\begin{array}{c} 0.55 \pm 0.01 \\ 0.57 \pm 0.02 \\ 0.62 \pm 0.01 \end{array}$
DENSE-TFIDF DENSE-Twit DENSE-BERT	$\begin{array}{c} 0.21 \pm 0.0 \\ 0.19 \pm 0.01 \\ 0.24 \pm 0.0 \end{array}$	$\begin{array}{c} 0.76 \pm 0.01 \\ 0.74 \pm 0.02 \\ 0.73 \pm 0.0 \end{array}$	$\begin{array}{c} 0.66 \pm 0.01 \\ 0.69 \pm 0.01 \\ 0.68 \pm 0.02 \end{array}$	$\begin{array}{c} 0.72 \pm 0.02 \\ 0.71 \pm 0.0 \\ 0.72 \pm 0.0 \end{array}$	$\begin{array}{c} 0.64 \pm 0.01 \\ 0.67 \pm 0.01 \\ 0.66 \pm 0.02 \end{array}$

Table 4: Phase 3 Experiment Results, Tweet + User Topic Embeddings. F is f1 score, F_H indicates hateful class, and F_A indicates abusive class. P is precision and R is recall. The red highlight indicates best performing model and the blue highlight indicates the TF-IDF embeddings that see the most improvement with the inclusion of user behavior metrics.



Figure 3: **Topic LDA Modeling of User's Timeline from Sample Hateful Tweet.** Sample hateful tweet is: "Muslims immigrants are living like they're used to where their from, squalor, filth and violence are all they know." This illustrates output from the Gensim library for LDA Topic modeling.



Figure 4: Hyperparameter Tuning Parallel Coorodinates Chart. This chart illustrates the parameter choices and validation f-score on the hateful class across 17 experiments.

model converged on average at epoch 44, we do not ex-470 periment with number of epochs. Figure 4 illustrates the 17 471 rounds of tuning experiments, with runtimes ranging from 472 5 hours and 20 minutes to 4 minutes (see Appendix X for 473 machine configurations). We use the Bayes algorithm with 474 the adaptive Parzen-Rosenblatt estimator (facilitated by the 475 Comet.ML library's optimization framework (Comet.ml)), 476 which balances exploring unknown space with exploiting 477 the known hyperparameter values that yield the best re-478 sults. The following values perform best of the values ex-479 plored: num-filters=47, num-layers=1, dropout=0.8877, 480 learning-rate=0.0196. This leaves us with a single layer 481 CNN with many filters and with a large regularization value 482 for dropout. 483

Tuned Model Results Across Experiments 484

Table and Figure 5 illustrate the tuned model performances. 485 All the tuned models demonstrate higher f-score on the mi-486 nority class, with the tuned phase 2 reply-pairings model 487 performing best at 0.33 ± 0.01 . This model also performs 507 488 best on all classes, with an overall f-score of 0.72 ± 0.01 . 489 It achieves the highest recall but not the highest precision of 509 490 the experiments conducted so far. We hypothesize that the 510 491 tuned phase 2 reply-pairings model performs better than the 511 492 other models because the single layer CNN with many fil-512 493 ters is overfitting with the tweet and topic embeddings. With 513 494 only the tweet and reply-pairing BERT embeddings, the net- 514 495 work seems to have a more generalizeable notion of context. 515 496

In order to better understand our precision and recall 516 497 scores, we analyze a confusion matrix in table 6. We see that 517 498 the hateful class was correctly classified 40% of the time 518 499 and was most commonly confused with the abusive class. 500 519 The abusive class was correctly classified 82% of the time, 501 which is higher than the model's ability to detect normal 502 tweets or spam tweets. In the below samples, the model pre-503 dicts a tweet is hateful when it is labeled as hateful, abusive, 504 523 and normal, respectively. Personally identifiable information 505 524 is removed. 506



Figure 5: Tuned CNN-BERT Experiment Results. All the tuned models demonstrate higher f-score on the minority class, with the tuned phase 2 model performing best at 0.33.



Figure 6: Confusion Matrix of Final Model. The model can detect the abusive class most accurately.

RHate when folk ask me questions that they already know the answers to . B*TCH what you fishing for ? (Hateful Label)

508

520

521

522

525

That was the most patronising thing I've ever heard about young people and propaganda. What about old people and the mail?! (Abusive Label)

TEACHERS too! Another "trusted" member of society, *doing the nasty with teen.* (Normal Label)

There are a few problems with this. First, the tweet that is labelled as hateful does not appear to be hate speech from the definitions explored in our background. Second, the tweet that is labelled as normal could be considered abusive towards teachers. All of these tweets have negative sentiment, so it does appear our model has picked up on that. This deeper dive into the actual content that is misclassified hints at a problem with the annotation quality of our dataset.

MODEL	F_H	F_A	F	<i>P</i>	R
phase-1	0.32 ± 0.02	0.77 ± 0.01	0.71 ± 0.01	0.76 ± 0.01	0.69 ± 0.01
phase-2-1	0.33 ± 0.01	0.78 ± 0.01	0.72 ± 0.01	0.76 ± 0.0	0.7 ± 0.01
phase-2-2	0.32 ± 0.0	0.77 ± 0.01	0.7 ± 0.01	0.76 ± 0.0	0.69 ± 0.01
phase-3	0.32 ± 0.01	0.78 ± 0.01	0.7 ± 0.02	0.77 ± 0.01	0.68 ± 0.02

Table 5: **Tuned CNN-BERT Model Performance.** F is f1 score, F_H indicates hateful class, and F_A indicates abusive class. P is precision and R is recall.

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

Conclusion

526

Our final model successfully picked up on negative tweet 527 528 sentiment and identified the abusive class at the highest 529 rate, of 82% accuracy and 0.78 f-score. The models that used Google's pretrained BERT embeddings performed bet-530 ter than TF-IDF and Twitter pretrained embeddings across 531 most models in our three phases of experiments. CNN-532 BERT outperformed the logistic regression, MLP, LSTM, 533 and DenseNet models for all three phases of experiments. 534 Before tuning, our best performing model is the CNN multi-535 ple input model architecture with tweet and user topic BERT 536 embeddings. After tuning, our best performing model is the 537 CNN multiple input model architecture with tweet and reply 538 BERT embeddings. We hypothesize that this is because the 539 parameter choices of a single layer, 47 filter CNN with high 540 dropout will overfit with BERT and a large measure of user 541 context. The pretrained BERT embeddings add enough se-542 mantic information to give us our most competitive models 543 544 and adding additional metrics of context through aspects of the social network actually hurts performance. 545

The task of automating the detection of hate speech on 546 social media platforms remains a challenge, in part due to 547 the difficulty in obtaining high-quality, large-scale annotated 548 datasets and the scarce hateful samples available for ma-549 chine learning models to learn from. Our experiments reflect 550 this and suggest that improving the quality and consistency 551 of annotations in our dataset is likely to result in more accu-552 rate automated systems. 553

Even so, our final model offers a significant improvement 554 on detecting hate speech, as we are able to improve on our 555 logistic regression baseline performance on the hateful class 556 by .13 f-score on a dataset with scarce hateful labels. If we 557 were to randomly annotate a tweet as hateful with 4% proba-558 bility, we'd achieve around 4% accuracy on the hateful class. 559 Thus, we interpret the final f-score on the hateful class of .33 560 as relatively high. 561

The benefits of using BERT in our experiments suggest 562 that further work on pretrained embeddings, from larger or 563 more targeted datasets, holds more promise than adding in-564 formation specific to the social media platform. However, 565 LDA topic modeling is a useful tool for gauging the quality 566 of the dataset annotations, as it allowed us to quickly see a 567 hateful tweet was mislabeled by clarifying that the user is a 568 film reviewer. Policy and decision-makers may benefit from 569 using topic modeling and other forms of user behavior met-570 rics in a contagion tool, i.e. a tool that monitors trends and 615 571 patterns of hate speech and how they might spread through 616 572

a network or change in response to triggering events.

This research hopes to inform interdisciplinary coalitions on the current landscape of hate speech detection. While our model's were able to classify abusive text with a high accuracy, they struggled with the task of distinguishing hate speech from normal and abusive content. A "good" automated hate speech detector may mean detecting detecting 30% of hate speech and passing merely abusive content to human moderators. The technical limitations described indicate that policy-makers must cooperate with social media companies in order to guide the handling of classified hateful and abusive content. Ultimately, social media has contributed to a more open and connected world. It is critical that we mitigate the negative consequences of hate speech while preserving the benefits of online discussion.

References

- 589 [Awan 2017] Awan, I. 2017. Cyber-extremism: Isis and the590 power of social media. *Society* 54(2):138–149.
- [Badjatiya et al. 2017] Badjatiya, P.; Gupta, S.; Gupta, M.;
 and Varma, V. 2017. Deep learning for hate speech detection in tweets. In *Proceedings of the 26th International Conference on World Wide Web Companion*, 759–760. International World Wide Web Conferences Steering Committee.
- [Blei, Ng, and Jordan 2003] Blei, D. M.; Ng, A. Y.; and Jordan, M. I. 2003. Latent dirichlet allocation. *Journal of machine Learning research* 3(Jan):993–1022.
- [Breuer, Landman, and Farquhar 2015] Breuer, A.; Landman, T.; and Farquhar, D. 2015. Social media and protest
 mobilization: Evidence from the tunisian revolution. *De- mocratization* 22(4):764–792.
- [Burnap and Williams 2015] Burnap, P., and Williams, M. L.
 2015. Cyber hate speech on twitter: An application of machine classification and statistical modeling for policy and
 decision making. *Policy & Internet* 7(2):223–242.
- ⁶⁰⁷ [Cave 2019] Cave, D. 2019. Countries want to ban ⁶⁰⁸ 'weaponized' social media. what would that look like?
- ⁶⁰⁹ [Cheng et al. 2017] Cheng, J.; Bernstein, M.; Danescu⁶¹⁰ Niculescu-Mizil, C.; and Leskovec, J. 2017. Anyone can
 ⁶¹¹ become a troll: Causes of trolling behavior in online dis⁶¹² cussions. In *Proceedings of the 2017 ACM conference on*⁶¹³ computer supported cooperative work and social comput⁶¹⁴ ing, 1217–1230. ACM.
 - [Comet.ml] Comet.ml. Comet.ml supercharging machine learning.

- 617 [Davidson et al. 2017] Davidson, T.; Warmsley, D.; Macy, 672
- 618 M.; and Weber, I. 2017. Automated hate speech detection
- and the problem of offensive language. In *Eleventh interna*-
- tional aaai conference on web and social media.
- 621 [Del Vignal et al. 2017] Del Vignal, F.; Cimino, A.; Del-
- ⁶²² lâĂŹOrletta, F.; Petrocchi, M.; and Tesconi, M. 2017. Hate
- me, hate me not: Hate speech detection on facebook.
- 624 [ElSherief et al. 2018] ElSherief, M.; Nilizadeh, S.; Nguyen,
- D.; Vigna, G.; and Belding, E. 2018. Peer to peer hate: Hate
- speech instigators and their targets. In Twelfth International
- 627 AAAI Conference on Web and Social Media.
- 628 [Enikolopov, Petrova, and Sonin 2018] Enikolopov, R.;
- 629 Petrova, M.; and Sonin, K. 2018. Social media and cor-
- ruption. American Economic Journal: Applied Economics 10(1):150–74.
- 632 [Founta et al. 2018] Founta, A. M.; Djouvas, C.; Chatzakou,
- D.; Leontiadis, I.; Blackburn, J.; Stringhini, G.; Vakali, A.;
- 634 Sirivianos, M.; and Kourtellis, N. 2018. Large scale crowd-
- sourcing and characterization of twitter abusive behavior. In
- 636 Twelfth International AAAI Conference on Web and Social
- 637 Media.
- 638 [Gambäck and Sikdar 2017] Gambäck, B., and Sikdar, U. K.
- 639 2017. Using convolutional neural networks to classify hate-
- speech. In *Proceedings of the first workshop on abusive lan-*
- 641 guage online, 85–90.
- 642 [Godin et al. 2015] Godin, F.; Vandersmissen, B.; De Neve,
- 643 W.; and Van de Walle, R. 2015. Multimedia lab@ acl wnut
- ner shared task: Named entity recognition for twitter micro-
- 645 posts using distributed word representations. In *Proceedings*
- of the workshop on noisy user-generated text, 146–153.
- 647 [Habibi et al. 2017] Habibi, M.; Weber, L.; Neves, M.; Wie-
- gandt, D. L.; and Leser, U. 2017. Deep learning with word
- embeddings improves biomedical named entity recognition. *Bioinformatics* 33(14):i37–i48.
- 651 [Huang et al. 2017] Huang, G.; Liu, Z.; Van Der Maaten, L.;
- and Weinberger, K. Q. 2017. Densely connected convolu-
- tional networks. In Proceedings of the IEEE conference on
- 654 *computer vision and pattern recognition*, 4700–4708.
- 655 [Jackson, Bailey, and Foucault Welles 2018] Jackson, S. J.:
- 656 Bailey, M.; and Foucault Welles, B. 2018. # girlslikeus:
- 657 Trans advocacy and community building online. New Me-
- 658 *dia & Society* 20(5):1868–1888.
- [Marcin 2019] Marcin, T. 2019. Facebook, youtube, whatsapp banned again in sri lanka after violence against muslims.
- [Müller and Schwarz 2018] Müller, K., and Schwarz, C.
 2018. Fanning the flames of hate: Social media and hate
 crime. *Available at SSRN 3082972*.
- ⁶⁶⁴ crime. Available at SSRN 3082972.
 ⁶⁶⁵ [Nishi et al. 2016] Nishi, R.; Takaguchi, T.; Oka, K.; Mae-
- hara, T.; Toyoda, M.; Kawarabayashi, K.-i.; and Masuda,
- 667 N. 2016. Reply trees in twitter: data analysis and branch-
- ing process models. Social Network Analysis and Mining
- 669 6(1):26.
- 670 [Pitsilis, Ramampiaro, and Langseth 2018] Pitsilis, G. K.;
- Ramampiaro, H.; and Langseth, H. 2018. Detecting offen-

- sive language in tweets using deep learning. *arXiv preprint arXiv:1801.04433*.
- 674 [Rainie, Anderson, and Albright 2017] Rainie, H.; Ander-

673

- son, J. Q.; and Albright, J. 2017. *The future of free speech, trolls, anonymity and fake news online*. Pew Research Cen-
- 677 ter Washington, DC.
- 678 [Ribeiro et al. 2017] Ribeiro, M. H.; Calais, P. H.; Santos,
- Y. A.; Almeida, V. A.; and Meira Jr, W. 2017. " like sheep among wolves": Characterizing hateful users on twitter. *arXiv preprint arXiv:1801.00317*.
- [Rosemain 2019] Rosemain, M. 2019. Exclusive: In a world
 first, facebook to give data on hate speech...
- ⁶⁸⁴ [Santos and Zadrozny 2014] Santos, C. D., and Zadrozny, B.
- ⁶⁸⁵ 2014. Learning character-level representations for part-of-⁶⁸⁶ speech tagging. In *Proceedings of the 31st International*
- 687 Conference on Machine Learning (ICML-14), 1818–1826.
- [Schmidt and Wiegand 2017] Schmidt, A., and Wiegand, M.
 2017. A survey on hate speech detection using natural language processing. In *Proceedings of the Fifth International*Workshop on Natural Language Processing for Social Media, 1–10.
- [Suh et al. 2010] Suh, B.; Hong, L.; Pirolli, P.; and Chi, E. H.
 2010. Want to be retweeted? large scale analytics on factors
 impacting retweet in twitter network. In 2010 IEEE Second International Conference on Social Computing, 177–
 184. IEEE.
- [Tang et al. 2014] Tang, D.; Wei, F.; Qin, B.; Liu, T.; and
 Zhou, M. 2014. Coooolll: A deep learning system for twitter sentiment classification. In *Proceedings of the 8th inter- national workshop on semantic evaluation (SemEval 2014)*,
- 702 208–212.
- Thomasson 2017] Thomasson, E. 2017. German cabinetagrees to fine social media over hate speech.
- [Wakefield 2019] Wakefield, J. 2019. Hate speech: Face-book, twitter and youtube told off by mps.
- ⁷⁰⁷ [Waseem and Hovy 2016] Waseem, Z., and Hovy, D. 2016.
 ⁷⁰⁸ Hateful symbols or hateful people? predictive features for
 ⁷⁰⁹ hate speech detection on twitter. In *Proceedings of the*⁷¹⁰ NAACL student research workshop, 88–93.
- 711 [Wiedemann et al. 2018] Wiedemann, G.; Ruppert, E.; Jin-
- dal, R.; and Biemann, C. 2018. Transfer learning from lda to
 bilstm-cnn for offensive language detection in twitter. *arXiv*
- 714 *preprint arXiv:1811.02906*.
- 715 [Xiang et al. 2012] Xiang, G.; Fan, B.; Wang, L.; Hong, J.;
 716 and Rose, C. 2012. Detecting offensive tweets via topical
- ⁷¹⁷ feature discovery over a large scale twitter corpus. In *Pro*-
- 718 ceedings of the 21st ACM international conference on Infor-
- 719 mation and knowledge management, 1980–1984. ACM.
- [Xu et al. 2012] Xu, J.-M.; Jun, K.-S.; Zhu, X.; and Bellmore, A. 2012. Learning from bullying traces in social
 media. In *Proceedings of the 2012 conference of the North American chapter of the association for computational lin-*guistics: Human language technologies, 656–666. Association for Computational Linguistics.