

Classification Methods for Hate Speech Diffusion: Detecting the Spread of Hate Speech on Twitter

Matthew Beatty¹

¹Department of Computer Science, Harvard University
33 Oxford Street, Cambridge, MA
mbeatty@fas.harvard.edu

Abstract - In this paper, we investigate predictive models to detect the spread of hate speech on Twitter based on diffusion patterns. We experiment with a dataset of 10,000 tweets manually labelled as hate speech or not and show that classification based solely on the sharing graph yields strong F1 scores for our task and high hate speech detection precision. We also highlight the vulnerability of existing textual hate speech detection methods to adversarial attacks and demonstrate that while our methods do not outperform state-of-the-art text models, graph-based models provide robust detection mechanisms and are able to detect instances of hate speech that missed by text classifiers. We find that graph convolutional networks produce the strongest hate speech F1 score of 0.58 and that kernel methods offer strong predictive potential. Finally, we also consider the effects of automated bots in the diffusion of hate speech content and conclude that their sharing behavior plays an insignificant role in our experiments.

Keywords: Graph classification, Graph mining, Graph kernels, Hate speech, Twitter

1. Introduction

Online social networks have connected people across the world and made sharing information easier than ever, but rising reports of active hate groups and online hate speech content trouble observers. The anonymity, immediacy, and global reach of online social networks make the platforms potentially attractive broadcast mediums for hate groups. Several large social networking platforms have identified the problem, but their responses have not solved the issue. In 2018, Twitter representatives suggested that the company's automated detection systems used machine learning techniques to identify over 40% of content that requires moderation but admitted that finding all hate speech using automated detection mechanisms is prohibitively difficult [1], [2].

The classification performance of automated hate speech detection systems has improved dramatically in recent years, but three issues still plague existing methods. Hate speech content accounts for a tiny fraction of all shared content, leading to a highly imbalanced classification problem. The second issue is that the subjectivity of defining hate speech makes labelling data for an automated detection system even more difficult. There is no consensus definition of hate speech, and the context becomes critical in many scenarios. Finally, current text-based detection methods can be easily fooled by adversarial attacks. Previous research has shown that trivial manipulations of messages, such as adding positive sentiment words or small typos in hateful terms, can trick systems built to detect toxic language [3], [4]. Current systems also have poor detection of coded language in content, and the ability of simple modifications to impair the abilities of current detection systems significantly undermines the recent performance gains of state-of-the-art systems. Our research is motivated by a desire to uncover new methods of hate speech detection to address these existing shortcomings. Previous research has analyzed online users who frequently post hate speech content, but few studies to our knowledge analyze the diffusion of hate speech online and whether diffusion patterns can offer predictive detection features.

Our contributions are: (1) we investigate models to detect the spread of hate speech tweets on Twitter and show that graph convolutional networks yield robust F1 scores for our imbalanced classification task and high hate speech detection precision. (2) We reproduce existing results showing how adversarial attacks can weaken text detection models for hate speech detection and demonstrate how our methods, while not outperforming state-of-the-art models, are robust against adversarial attacks and hard to manipulate. (3) We consider the role of automated accounts or "bots" in the spread of hate speech and find that automated accounts play an insignificant role in the spread of hate speech in our experimental dataset.

2. Related Work

Hate speech detection systems have relied primarily on natural language processing, and recent deep learning research has produced the current state-of-the-art detection capabilities. Earlier research demonstrated the effectiveness of n-gram and bag-of-words techniques. Specifically, Davidson, et al. show strong classification performance with an F1 score of 0.80 on a multi-class hate speech classification task using word-level n-gram features [5]. Deep learning models, like recurrent neural networks, have been shown to detect hate speech with high levels of precision, and Badjatiya et al. produce a classifier with 0.93 precision and 0.93 recall overall on a dataset of 16,000 labelled tweets using long short-term memory (LSTM) networks with random embeddings and gradient-boosted decision trees [6]. But other researchers have undermined the stated gains of text-based classifiers by demonstrating that simple adversarial attacks can fool them. Grondahl, et al. show that “ appending the word “love” to a comment caused hate speech detection precision to drop by 50% for the detection models from Davidson et al., Badjatiya et al., and Google’s Perspective API which is used to detect toxic language online [7].

Some research studies have analyzed patterns in the types of users associated with hate speech content, but they did not look at the content’s spread. Ribiero, et al. analyze users who create or share hate speech content on Twitter. The authors find that accounts posting hate speech tweets tend to post more frequently than other accounts and also have strong ties in their 1-neighborhood, suggesting strong homophily among that group of accounts [4]. Both Zhong et al. and Waseem and Hovy consider the effect of the number of replies, the replies from followers, and the total number of posts by a user to detect hate speech, but the two studies produce opposite results. Zhong finds no correlation between posting activity and hateful content, while Waseem and Hovy find a strong correlation [8].

To our knowledge, no extensive research has attempted to analyze the spread of hate speech diffusion on Twitter. Previous research has found that information cascades online carrying negative sentiments fade far more rapidly than their positive sentiment counterparts, and Romero, Meeder, and Kleinberg analyze how diffusion cascades on Twitter differed by topic, such as sharing cascades about politics, sports, or movies. [10] [9]. The 2018 research of Vosoughi, Roy, Aral offers a similar investigation of the spread of true and false rumors online. The authors analyze the diffusion cascades of verified true and false news stories on Twitter from 2006 through 2017 and find that false rumors spread faster and more broadly than their truthful counterparts. The authors also find that automated accounts play a limited role in the spread of false news stories [11].

3. Dataset

We use an existing dataset of labelled hate speech tweets and data from Twitter’s public API to collection information about sharing patterns. The tweet dataset contains 99,799 English tweets created between March 30, 2017 and April 9, 2017 and was compiled by Founta, et al. [12]. The authors crowdsource annotations from CrowdFlower to label each tweet as either hateful, abusive, spam, or normal. They use the following definitions for each label:

- Hate Speech (Hateful): 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.
- Abusive Language (Abusive): Any strongly impolite, rude or hurtful language using profanity, that can show a debasement of someone or something, or show intense emotion.
- Spam: Posts consisted of related or unrelated advertising / marketing, selling products of adult nature, linking to malicious websites, phishing attempts and other kinds of unwanted information, usually executed repeatedly.
- Normal: all tweets that do not fall in any of the prior categories. [12]

Each tweet in the dataset has annotations from five labelers and require a majority label for the tweet to be included in the dataset. The original dataset contains 53,790 (54%) normal tweets, 14,024 (14%) spam tweets, 27,037 (27%) abusive tweets, and 4,948 (5%) hateful tweets.

Next, we gather each tweet’s associated diffusion network through Twitter’s public API. We use a process for recreating true retweet paths using the limited information from the Twitter API called time-inferred diffusion, described initially by Goel, et al. and used by Vosoughi, Roy, Aral in their research on the spread of true and false news on Twitter

[13]. The Twitter API does not specify the direct retweet path, meaning all retweets in the API response appear to reshare the original tweet even if they shared the content through a later retweet, so we must use this diffusion inference technique to infer the true sharing patterns using follower relationships and retweet timestamps.

We restrict our gathering to tweets with ten or more retweets, because smaller graphs will not contain enough information for meaningful predictions. At the time we access the Twitter API, 42,872 tweets are unavailable either because users deleted the tweet or Twitter banned the account or tweet. The follower relationships of users who have deleted their accounts are also not available, and a small fraction (2%) of accounts in the dataset are deleted. All of these weaken our ability to infer the true retweet graph, but we are still able to gather information for 86% of all retweets of the tweets in the dataset, according to tweet metadata. Some tweets in the original dataset are retweets, and we remove all retweets so that only original tweets are classified.

Only 6,924 available tweets from the original dataset meet our minimum retweet requirements, and there are 2,154 “abusive” tweets, 967 “hateful” tweets, and 3,203 “normal” tweets. There are larger proportions of abusive and hateful content than the original dataset, and the cause is a boosted abusive and hateful tweet sample of 20,000 tweets from the original dataset. The authors include 20,000 tweets in the dataset after filtering their sample using sentiment analysis and the presence of hate speech terms to gather tweets, and the tweets from this sample tend to be more accessible and more popular than the rest of the dataset [12]. To counteract this, we resample using the opposite filters as the dataset’s boosted hate speech sample to gather more “normal” tweets to counteract this imbalance. We gather tweets from the same period in April 2017 and filter by tweets that have high positive sentiment scores and that do not contain any of the hate speech terms the authors require in their original boosted sampling process.

In the end, our experimental dataset includes 10,074 tweets. 2,154 (21.4%) tweets are labeled as abusive, 967 (9.6%) are labeled as hateful, and 6,953 (69.2%) are labeled as neither (i.e. normal). For our binary classification experiments, we combine hateful and abusive into a “toxic” class:

Table 1: Dataset Class Labels.

Tweet Class	Number
Normal	6,953
Toxic	3,121
Total	10,074

Given our minimum size constraints, the sizes for most of the eligible sharing graphs cluster towards the minimum size requirement, but several sharing graphs have more than 500 nodes (retweets) and more than 10 layers of depth. We have a higher proportion of normal tweets and slightly higher number of hateful tweets than the original dataset.

4. Experiments

We conduct three sets of hate speech detection experiments on our dataset: classifying retweet networks, classifying tweet texts, and classifying tweet text while applying adversarial attacks. All experiments are run as binary classification tasks between normal tweets as one class and “toxic” tweets, the union of abusive and hateful tweets, as the other class. In each experiment, we optimize for overall F1 score. For each task, we report the precision, recall, and F1 score for each class. The stated results are aggregated over the same ten 90-10 test-train splits, and we use GridSearchCV for model hyperparameter tuning where appropriate.

For our diffusion network classification, we test regression methods, kernel-based methods, and graph convolution network models. We run logistic regression with a set of hand-crafted features based on structural properties of the retweet graphs, such as the graph depth, cascade duration, node degree distribution, and assortativity coefficient, using 50 features in total to make predictions. Next, we include a series of graph kernel-based learning models. We use the open-source GraKel package for implementations of kernel functions. We report results for our SVM classification model using a support vector machine classifier with a precomputed Weisfeiler-Lehman kernel. (We also considered a graph shortest-

paths and random walk kernels, but after their poor performance, we exclude those models from the results.) We test the deep graph kernels proposed in Yanardag, 2015 as well. The deep graph kernel uses a neural language model to learn latent representations of the sub-structures in the graphs, using labeling from the Weisfeiler-Lehman kernel. The last graph-based classification model we test is a graph convolutional network (GCN). GCNs are semi-supervised approach to graph classification, reaching state-of-the-art performance on a number of benchmark graph classification tasks. These are extensions of neural networks applied to capture information specific to graph data structures. We perform training in batches of size 64 for the GCNs with at most 100 epochs.

We compare our results with an implementation of the logistic regression model from Davidson et al. using word-based n-grams and other text-features (readability metrics, tweet length, etc.). We also implement the methods from Malmasi, Zampieri where the authors use SVMs with different word and character-level skip-gram features. We report results from the strongest classifier from their set, specifically using word-level skip-gram features as model input. We also compare our methods against state-of-the-art deep learning models. We investigate two neural network architectures for the task, specifically convolutional neural networks (CNNs) and long short-term memory networks (LSTMs), as described by Badjativa, et al. The authors experiment with both convolution and recurrent neural net architectures for hate speech detection. They find that the recurrent LSTMs have the strongest performance in detection hate speech [6]. As part of our experiments to show the effects of adversarial attacks on existing hate speech detection options, we also implement a few simple adversarial attacks from Grondahl, et al. Specifically we include four adversarial attacks: inserting typos, removing whitespace, and appending the word “love” to the end of hate speech tweets.

5. Results

Table 1 includes results for our graph-based methods for the hate speech classification task. Our diffusion classification methods achieve varying levels of success.

Table 2: Diffusion Graph Classification Results.

Method	Class	Precision	Recall	F1
Logistic Regression	Normal	0.91	0.69	0.78
	Hateful	0.37	0.53	0.33
SVM, WL Kernel	Normal	0.80	1.00	0.89
	Hateful	0.33	0.12	0.20
SVM, Deep WL Kernel	Normal	0.82	0.18	0.26
	Hateful	0.39	0.18	0.26
Graph Convolutional Network	Normal	0.86	0.93	0.90
	Hateful	0.55	0.34	0.42

The graph convolution network produces the best F1 score of 0.40 for the hate speech class and the highest weighted accuracy of our diffusion graph classification methods. Our kernel-based methods reliably detect normal content but have severe difficulty detection hate speech, producing very low recall scores. Our baseline logistic regression produces the highest hate speech recall at 0.53, but it has a precision of only 0.37. In general, the detection methods can classify the majority class of normal tweets well but struggle to separate the retweet networks for abusive and hateful tweets. Another finding is that graph classification accuracy varies with graph size. Our models like the graph convolution network are able to classify graphs with more nodes and those with higher max depths much more accurately. Larger graphs offer more information to make a prediction, and we also find that larger networks in our dataset are disproportionately those of normal content. The diffusion networks of hate speech tweets tend to propagate widely and quickly at first but do not spread as deeply as other content.

Next, we compare existing textual methods on the dataset. Table 2 shows the results of existing text classification implementations on tweets from the dataset.

Table 3: Source Tweet Text Classification Results.

Method	Class	Precision	Recall	F1
Logistic Regression	Normal	0.83	0.97	0.89
	Hateful	0.68	0.46	0.56
SVM, n-grams	Normal	0.83	0.93	0.87
	Hateful	0.60	0.54	0.56
LSTM CNN	Normal	0.88	0.98	0.93
	Hateful	0.74	0.62	0.67
Bi-directional LSTM CNN	Normal	0.87	0.97	0.92
	Hateful	0.81	0.54	0.65

The LSTM classifier has the best performance of our implemented detection methods. It produces the highest F1 scores for all classes and highest overall accuracy, and the bi-directional LSTM CNN model follows closely in performance. The classifiers using n-gram features do not match the performance of the deep learning architectures but still produce higher F1 scores overall and for the hate speech class than our diffusion network classifiers. While these models outperform our sharing network detection techniques, they still exhibit performance significantly below the stated performance on their original datasets. We attempt to tune hyperparameters to optimize performance, specifically F1 score, but the results suggest either there are peculiarities in our dataset that make classification especially difficult or the initial results are over-trained on their original datasets. In total, there are 116 hateful tweets (12% of total) that were detected by our graph convolution network model that are not detected by any of the text-based methods. One major benefit of our diffusion-based approach is that we are able to pick out instances of hate speech that do not use obviously derogatory terms unlike the text-based classifiers.

The convolutional neural networks provide the best classification performances of the implemented text classifiers. They exhibit slight improvements in classifying normal tweets but show major improvements classifying the text of toxic messages. Some of the text classifiers exhibit F1 scores below our graph convolution network model after applying adversarial attacks. The more advanced classifiers using pretrained embeddings have their performance suffer significantly.

Table 4: Toxic Class F1 Scores with Adversarial Attacks.

Classifier	Original F1	Typos F1	Whitespace F1	Love F1
Logistic Regression	0.56	0.35	0.36	0.39
SVM, n-grams	0.56	0.33	0.34	0.43
LSTM	0.67	0.47	0.52	0.43
Bi-directional LSTM CNN	0.65	0.51	0.49	0.39

The simple addition of the word “love” at the end of each tweet is the most effective adversarial attack for those models, and the n-grams model has the largest performance loss when typos are added or whitespace is removed. Certain examples of toxic language were misclassified by all text classifiers after adversarial attacks by adding “love” including the following:

- “you, or a relative is a pig. That’s why I’m telling you to delete your account. You’re retarded **love**”
- “I ain’t never had soo much anger for just one person. Like I just wanna beat yo ass bro **love**”
- “Imagine giving an American cop a Pepsi at a demonstration! He would shoot you in the face and tip foaming soda over your twitching black corpse **love**”

Finally, we also consider the role of automated bots in our experimental results. Previous studies have considered the role of bots on Twitter and their role in spreading divisive content. We pass each source tweet in our dataset through IUNI’s Botometer tool for automated Twitter bot detection. The tool provides a metric called “Complete Automation Probability” which is the probability that the account and its postings are entirely automated [16]. Few accounts in the dataset have high CAP scores, and only 94 have CAP scores over 0.5. This suggests that the vast majority of accounts are authentic users and that bots are playing a very limited role in creating and spreading hate speech online. We also find no significant difference in the presence of bots across labels. These results align with the findings from Vosoughi who found that bots played little role in the spread of fake news online [11].

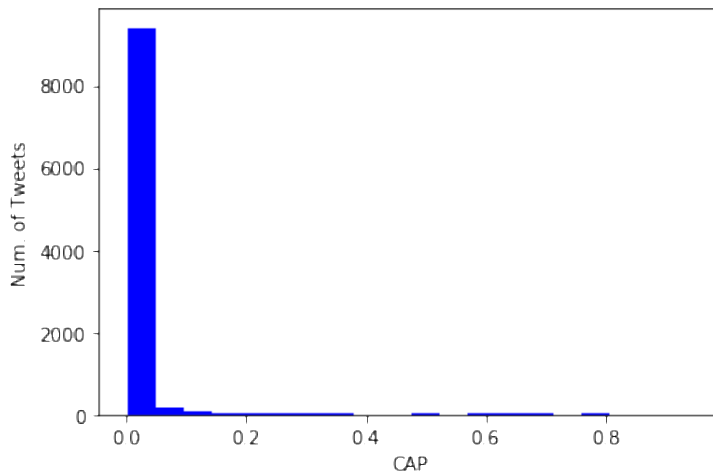


Fig. 1: Histogram of the distribution of Complete Automation Probability (CAP) values for tweets in the dataset. CAP values are in the range [0, 1].

5. Conclusion

We present novel methods for hate speech detection using the diffusion of tweets as features for classification. Our detection models are more robust than text classifiers which are susceptible to simple adversarial attacks, and our methods can detect a subset of hate speech missed by state-of-the-art text classifiers. Our study presents a new avenue of hate speech detection to explore, and there are several interesting implications. First, the dataset for this task is vital. We used an existing dataset to suit our needs but sampling specifically to produce a dataset of retweet networks could avoid the pitfalls introduced by our data collection methods. Future work should also focus on applications of advanced graph-based machine learning techniques. New graph-based classification techniques may provide models with better performance than our experimental methods. Also, in our experiments we assume that the adversarial attacks come from a single actor producing adversarial attacks against our diffusion-based detection difficult. Future robustness experiments should consider the effects of a coordinated group attacks that could change the sharing patterns.

Examining our methods using other datasets could demonstrate further robustness. For instance, our methods are language agnostic and therefore should demonstrate similar results on non-English datasets. Exploring differences in the spread of hate speech contents among different countries and communities could provide for better automated detection systems and behavioral insights into communities. Likewise, we do not analyze the specific accounts associated with hate speech diffusion. Our results show that hate speech tweets tend to be shared rapidly but not propagate broadly, suggesting that perhaps the hateful content remains constrained within an insular community. Investigating community dynamics and communication patterns of hate groups and their followers could provide more features for detection.

Acknowledgements

We would like to first thank Dr. Nir Rosenfeld. We are indebted to Dr. Rosenfeld for his guidance and insights during the research process. We would also like to thank Prof. Yaron Singer for his timely advice and perspective on the research. We further thank Antigoni Founta for her assistance with details on using her dataset and the Indiana University of Network Science Institute, specifically Matthew Hutchinson, Chathuri PeliKankanamalage, and Clayton Davidson, for allowing us to use their tools.

References

- [1] Cox, J., Koebler, J.: Twitter Won't Treat White Supremacy Like ISIS Because It'd Have to Ban Some GOP Politicians Too, https://www.vice.com/en_us/article/a3xgq5/.
- [2] Benner, K.: Twitter Adds New Ways to Curb Abuse and Hate Speech, <https://www.nytimes.com/2016/11/16/technology/twitter-adds-new-ways-to-curb-abuse-and-hate-speech.html>.
- [3] Sonnad, N.: Alt-right trolls are using these code words for racial slurs online, <https://qz.com/798305>
- [4] Ribeiro, M. H., Calais, P. H., Santos, Y. A., Almeida, V. A., Meira Jr, W.: Characterizing and Detecting Hateful Users on Twitter. In: Proceedings of the Twelfth International AAAI Conference on Web and Social Media, 2018.
- [5] Davidson, T., Warmsley, D., Macy, M., Weber, I.: Automated Hate Speech Detection and the Problem of Offensive Language. In: Eleventh International AAAI Conference on Web and Social Media, 2017.
- [6] Badjatiya, P., Gupta, S., Gupta, M., Varma, V.: Deep Learning for Hate Speech Detection in Tweets. In: Proceedings of the 26th International Conference on World Wide Web Companion, pp. 759-760, 2017.
- [7] Grondahl, T., Pajola, L., Juuti, M., Conti, M., Asokan, N.: All You "Need" is "Love" Evading Hate Speech Detection. In: Proceedings of the 11th ACM Workshop on Artificial Intelligence and Security, pp. 2-12, 2018.
- [8] Waseem, Z., Hovy, D.: Hateful Symbols or Hateful People? Predictive Features for Hate Speech Detection on Twitter. In: Proceedings of the NAACL Student Research Workshop, pp. 88-93, 2016.
- [9] Romero, D. M., Meeder, B., Kleinberg, J.: Differences in the Mechanics of Information Diffusion across Topics: Idioms, Political Hashtags, and Complex Contagion on Twitter. In: Proceedings of the 20th International Conference on World Wide Web, pp. 695-704, 2011.
- [10] Wu, B. and Shen, H.: Analyzing and Predicting News Popularity on Twitter. International Journal of Information Management, vol. 35, no. 6, pp. 702- 711, 2015.
- [11] Vosoughi, S., Roy, D., Aral, S.: The Spread of True and False News Online. Science, vol. 359, no. 6380, pp. 1146-1151, 2018.
- [12] Founta, A. M., Djouvas, C., Chatzakou, D., Leontiadis, I., Blackburn, J., Stringhini, G., Vakali, A., Sirivianos, M., Kourtellis, N.: Large Scale Crowdsourcing and Characterization of Twitter Abusive Behavior. In: Proceedings of the Twelfth International AAAI Conference on Web and Social Media, 2018.
- [13] Goel, S., Watts, D. J., Goldstein, D. G.: The Structure of Online Diffusion Networks. In: Proceedings of the 13th ACM Conference on Electronic Commerce, 2012.
- [14] Yanardag, P., Vishwanathan, S. V. N.: Deep Graph Kernels. In: Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 1365-1374, 2015.
- [15] Zampieri, M., Malmasi, S., Nakov, P., Rosenthal, S., Farra, N., Kumar, R.: Detecting Hate Speech in Social Media. In: Proceedings of Recent Advances in Natural Language Processing, 2017.
- [16] Davis, C. A., Ciampaglia, G. L., Aiello, L. M., Chung, K., Conover, M. D., Ferrara, E., Flammini, A., Fox, G. C., Gao, X., Goncalves, B., Grabowicz, P. A., Hong, K., Hui, P., McCaulay, S., McKelvey, K., Meiss, M. R., Patil, S., Peli Kankanamalage, C., Pentchev, V., Qiu, J., Ratkiewicz, J., Rudnick, A., Serrette, B., Shiralkar, P., Varol, O., Weng, L., Wu, T., Younge, A. J., Menczer, F. OSoMe: the IUNI Observatory on Social Media. PeerJ Computer Science, 2,e87, pp. 2376-5992, 2016. <https://doi.org/10.7717/peerj-cs.87>.