
Improving Hate Speech Classification on Twitter

Anonymous Author(s)

Affiliation

Address

email

1 Introduction

Hate speech and offensive language have begun to perpetuate online communities that were originally designed to foster community and bring people together. Both anonymity and the ease of spreading content online have made it easier for hateful speech to infiltrate large communities like Twitter. Many instances of hate speech occur in contexts where no explicit hate terms are used. This problem could be helped by a Machine Learning classifier that identifies hate speech. However, until Davidson et al. (2017)'s research, all classifiers were binary, classifying speech as either offensive or not. Hate speech is a separate from category offensive speech because it targets individuals based on nationality, ethnicity, religion, gender, sexual discrimination, disability or class in an especially aggressive or demeaning manner (Tuckwood, 2017). Davidson et al. (2017) were the first to identify this field as needing at least three classes: Hate Speech, Offensive Language, or Neither. Davidson et al. (2017) identify classification of tweets without explicit hate speech as difficult to correctly classify. We attempted to take Davidson et al. (2017)'s research further by adding more robust features to the Logistic Regression model in order to better capture the context surrounding tweets that don't contain explicit hate terms and correctly classify them.

2 Data and Methodology

The Davidson et al. (2017) data was obtained using the Twitter API to obtain 85.4 million tweets from 33,548 users, of which 24,783 tweets were selected to make up the final dataset. Crowdfunder, a crowd-sourcing website was used and annotators were provided with a formal definition of hate speech and asked to label each tweet as hate speech, offensive but not hate speech, or neither offensive nor hate speech. Every tweet was labeled by at least three annotators, and mean inter-annotator agreement was 92%.

2.1 Training

We trained all models using a set of 19K tweets from the dataset. Each model had a feature set concatenated with the baseline features. These models were then run through 5-fold cross-validation grid search on a Logistic Regression model.

2.2 Features / Baseline Features

We implemented a number of hand-built features and utilized Flair embeddings Akbik et al. (2018) as well. These were used in conjunction with the baseline feature set. We utilized Davidson et al. (2017)'s feature set as our baseline. These features included uni/bi/trigrams weighted by TF-IDF, binary and count indicators for hashtags, mentions, retweets, and URLs. To capture syntactic structure information they used NLTK and Penn Part-of-Speech (POS) taggings

33 2.2.1 Hand Built Features

34 We built a lexicon of ethnic and group membership words. We used this lexicon to create a binary
35 feature capturing if a tweet contains a statement targeting a specific group of people, i.e. "all you
36 Asians" or "every Mexican." This could possibly capture the nuance of a statement that isn't explicitly
37 hateful. Similarly, we tried to identify tweets where the name of a group was followed by a modal
38 verb like "should" or "can." We also tried to capture self-reference when an allusion to specific
39 group was made by implementing a feature that looked for first person pronouns followed by a word
40 indicating group membership.

41 We implemented an indicator feature for offensive / hate speech geared towards women. We sourced
42 gendered insults towards women to form a lexicon, where terms were scraped from a crowd-sourced
43 post (sac, 2018). This context-based feature was performed in two-steps: first, identify if a tweet is
44 aimed at a female (as indicated by pronouns). Second, check if the tweet has a gendered insult. This
45 differs from a simple 'contains check' because many female-specific insults like "feisty", "bossy",
46 etc. are offensive only in the context of being aimed at a woman.

47 Slang is an important characterization of tweets, so we wanted to capture the meaning behind slang
48 words instead of ignoring them. To decode slang terms, we mapped common Twitter slang terms to
49 their definitions and replaced any instance of slang with its definition. After replacing the slang, we
50 extracted the sentiment of the tweet. We utilized the (Marcus et al., 1993) PennTreebank to convert
51 tweets to Wordnet tags (Miller, 1995) to get the sentiment of tweets with slang replaced with their
52 definition.

53 We utilized the NRC Emotion Lexicon (Mohammad and Turney, 2013) to count the number of tokens in
54 a tweet referring to a specific emotion. We used the count for each emotion as its own separate
55 feature.

56 In addition to contextual features, we included some lexical features. The (Davidson et al., 2017)
57 model utilizes a porter stemmer when processing the tweets; we included a feature that did not stem
58 the words in tweets when creating tfidf weightings to see if that helped capture sentiment in tweets
59 that may not be explicitly offensive. We also included indicator features such as if a tweet referenced
60 immigrants directly and a feature that searched for a group membership word inside of quotes.

61 2.2.2 Flair

62 We wanted to include contextual string embeddings to better capture sentence-level context, since
63 we were interested in capturing the nuanced context of a tweet that contains hateful speech without
64 being explicit. (Akbik et al., 2018) created an embedding library called **Flair** that provides word and
65 sentence-level pre-trained embeddings. One of their word-level embeddings was trained using Twitter.
66 We chose to use this set of embeddings converted to sentence level, which (Akbik et al., 2018) call
67 "Document Pool" embeddings. We also utilized Flair's contextual string embeddings, one of which
68 was BERT embeddings (originally developed by (Devlin et al., 2018)). The second set of embeddings
69 was trained on a 1 billion word corpus from the news. We trained models using Twitter on its own
70 as well as in combination with the news and BERT embeddings. (Akbik et al., 2018) recommend
71 "stacking" word and string embeddings for the best results.

72 3 Results

73 3.1 Model Performance

74 Final results were reported after all models were run on a held-out test set comprised of 5K tweets.
75 We used 5-fold cross-validation grid search on a Logistic Regression model to find the optimal
76 parameters for each model. Our best performing model outperformed the baseline model by 2% in
77 recall of the Hate Speech class and by 3% in macro averaged recall. Interestingly, hand built features
78 did not seem to increase classification recall on the currently labeled data set. However, as will be
79 explored in the error analysis, this does not necessarily mean that the hand-built features are in reality
80 worse at correctly identifying hate speech.

81 **References**

- 82 2018. [Everyday misogyny: 122 subtly sexist words about women \(and what to do about them\)](#).
83 *Sacraparental*.
- 84 Alan Akbik, Duncan Blythe, and Roland Vollgraf. 2018. Contextual string embeddings for sequence
85 labeling. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages
86 1638–1649.
- 87 Thomas Davidson, Dana Warmley, Michael Macy, and Ingmar Weber. 2017. Automated hate speech
88 detection and the problem of offensive language. *International AAAI Conference on Web and
89 Social Media*.
- 90 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of
91 deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- 92 Mitchell Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. 1993. Building a large annotated
93 corpus of english: The penn treebank.
- 94 George A Miller. 1995. Wordnet: a lexical database for english. *Communications of the ACM*,
95 38(11):39–41.
- 96 Saif M Mohammad and Peter D Turney. 2013. Nrc emotion lexicon. *National Research Council,
97 Canada*.
- 98 Christopher Tuckwood. 2017. Hatebase: Online database of hate speech. *The Sentinel Project*.
99 Available at: <https://www.hatebase.org>.

Improving Hate Speech Classification on Twitter

Anonymous Author(s)
Affiliation
Address
email

1 Combined Features

2 We wanted to compare the performance of our hand-built features with the embeddings both separately
3 and combined. We concatenated the embedding features first with the baseline set of features and
4 trained the model on those feature sets. We did the same for the baseline features plus our hand
5 built features. After training separate models for these, we combined the hand-built features with the
6 Twitter embeddings and the combined news, BERT, and Twitter embeddings.

Table 1: [Davidson et al. \(2017\)](#) Logistic Regression Baseline

	Precision	Recall	F1-score
Hate Speech	0.30	0.45	0.36
Offensive	0.94	0.86	0.90
Neither	0.66	0.81	0.73
Micro avg	0.83	0.83	0.83
Macro avg	0.64	0.71	0.66
Weighted avg	0.86	0.83	0.84

Table 2: Twitter Embeddings

	Precision	Recall	F1-score
Hate Speech	0.25	0.35	0.29
Offensive	0.92	0.87	0.90
Neither	0.68	0.75	0.71
Micro avg	0.83	0.83	0.83
Macro avg	0.62	0.66	0.63
Weighted avg	0.84	0.83	0.83

Table 3: News, Twitter, and BERT Embeddings

	Precision	Recall	F1-score
Hate Speech	0.31	0.47	0.38
Offensive	0.95	0.88	0.91
Neither	0.74	0.85	0.79
Micro avg	0.86	0.86	0.86
Macro avg	0.67	0.74	0.69
Weighted avg	0.88	0.86	0.87

Table 4: Hand-Built Features

	Precision	Recall	F1-score
Hate Speech	0.24	0.36	0.29
Offensive	0.91	0.86	0.89
Neither	0.65	0.71	0.68
Micro avg	0.81	0.81	0.81
Macro avg	0.60	0.64	0.62
Weighted avg	0.83	0.81	0.82

Table 5: Hand-Built Features with Twitter Embeddings

	Precision	Recall	F1-score
Hate Speech	0.27	0.38	0.32
Offensive	0.92	0.87	0.90
Neither	0.68	0.75	0.71
Micro avg	0.83	0.83	0.83
Macro avg	0.62	0.67	0.64
Weighted avg	0.85	0.83	0.84

Table 6: Hand-Built Features with News, BERT, and Twitter Embeddings

	Precision	Recall	F1-score
Hate Speech	0.31	0.47	0.38
Offensive	0.95	0.89	0.92
Neither	0.74	0.84	0.79
Micro avg	0.86	0.86	0.86
Macro avg	0.67	0.73	0.69
Weighted avg	0.88	0.86	0.87

7 References

- 8 Thomas Davidson, Dana Warmley, Michael Macy, and Ingmar Weber. 2017. Automated hate speech
9 detection and the problem of offensive language. *International AAAI Conference on Web and*
10 *Social Media*.

Improving Hate Speech Classification on Twitter

Anonymous Author(s)
Affiliation
Address
email

1 0.1 Error Analysis

2 Error analysis with this dataset proved difficult because not all examples of tweets labeled as hate
3 speech are truly hate speech. Some examples include :

4 "I'm not really a phone kinda guy.. I actually hate talking on the phone amp; texting kinda trash to
5 me also."

6 "Whipped out some french in front of some babes at the post office. winning"
7

8 The Hand-Built Features with News, BERT, and Twitter Embedding (Table 6) model had 591
9 misclassified tweets. Of these, 88 were in the hate speech class, 99 in the neither class, 404 in the
10 offensive class.

11 The News, Twitter, and BERT Embeddings ((without our hand built features)) model had 489
12 misclassified tweets. Of these, 32 were in the hate speech class, 82 neither, and 375 offensive. Out
13 of the 489 misclassified tweets we do not agree with the labeling of 12% of the tweets with 36% of
14 those both wrongly classified by our model and wrongly labeled and 64% correctly predicted by our
15 model.

16 In the following sections we take a deep dive comparing the models Hand-Built Features with News,
17 BERT, and Twitter Embedding (Table 6) and News, Twitter, and BERT Embeddings (Table 3).

18 0.2 Hand-Built Features with News, BERT, and Twitter Embedding model Class: Hate 19 Speech

20 In the Hand-Built Features with News, BERT, and Twitter Embedding model we found 29 instances
21 of tweets that we believe were incorrectly labeled. That's 35% of all missed tweets in the hate speech
22 category. Of those, 3% of the tweets were both incorrectly labeled and incorrectly predicted by our
23 model, leaving 32% of tweets in the hate speech class that our model correctly predicted.

24 Amongst the hate speech labels we agree with, the targeted groups were: 25% Female with one
25 instance threatening violence and another suggesting the target commit suicide. 73% of these were
26 predicted as offensive showing a bias towards the offensive class when concerning women and
27 variations of the terms 'hoe' and 'bitches'.

28 20% Gay Community. The diverse variations and spellings of the term 'fag' make it difficult to weight
29 it towards hate speech. One possible feature could be a regex for the term 'fag' or gay in conjunction
30 with a swear word. Targets of the term are as follows: 6/12 males as a means to emasculate; 4/12
31 women; 1/12 males; 1/12 the gay community in general.

32 < 14% Males with six instances including insults meant to emasculate with three of those also
33 including threats of violence. A feature looking at males as a target and the usage of terms 'bitch'
34 and 'pussy' could weigh it from offensive to hate speech.

35 < 12% African American; with most of the tweets having hard to discern context identifying them
36 as hate speech such as a link to an article, usage of the n word in different spellings, and one tweet

37 where the hate speech was in quotes making it difficult to know if it was commentary condemning it
38 or agreeing with it. Targets of the term are as follows: 1/7 African American females; 6/7 African
39 Americans in general.

40 < 12% White people with one threat of violence and two including politics. The term 'white trash'
41 was in almost every instance. When running our features, we ran them over each token so a possible
42 feature is to do a regex for the term and format it as one token.

43 8% Asian; with every instance targeting Chinese people including some variation of the term 'chink',
44 making it hard to understand why our model failed to flag it as hate speech especially since half
45 of these instances were predicted to the class 'neither' by the model. One possibility is the low
46 number of hate speech in the training class making use of the term and the alternate definition of
47 chink meaning "a narrow opening or crack, typically one that admits light." Targets of the term are as
48 follows: 4/5 Chinese; 1/5 Asians and the gay community as a means to emasculate.

49 Remainder: 2/59 Politics, 1/59 general racist, 1/59 Jewish, 1/59 Latinos

50 Many of the tweets included masked swear or hate words surrounded by hash tags, or html tags.
51 Making a regex for these terms could help identify the intent.

52 **0.3 News, BERT, and Twitter Embedding model Class: Hate Speech**

53 For the News, Twitter, and BERT Embeddings model, we found 4 instances of tweets that we disagree
54 are hate speech. Three of these were correctly predicted by our model and one tweet was both
55 incorrectly labeled and wrongly predicted by our model.

56 For the News, Twitter, and BERT Embeddings model (without our hand built features), out of the
57 12% of the tweets that were incorrectly labeled, 36% of those were both wrongly classified by our
58 model and wrongly labeled and 64% correctly predicted by our model.

59 In the hate speech class, our model correctly predicted three instances labeled as hate speech as
60 offensive and neither.

61 It also correctly predicted 17 instances as hate speech: 76% were incorrectly labeled as offensive
62 speech; 18% were incorrectly labeled as neither; While some of the following target groups are
63 represented in the same tweet, instances include: 11% African Americans; 28% Female; 28% male
64 with three instances aiming to emasculate men; 22% gay community in general; 5% Asian; and 5 %
65 White people.

66 We note that this model improved our recall for hate speech targeting African Americans and White
67 people compared to the combined usage of the hand built features. It would be useful to comment out
68 certain features to see what is reducing performance.

69 Our model missed 28 instances of hate speech our model incorrectly predicted 39% into the neither
70 class, and 61% of tweets into the offensive class.

71 Of the 17 instances incorrectly predicted as offensive: 35% female with two of those encouraging
72 suicide; 24% gay community; 24% male with three including threats of sexual violence / general
73 violence; 12% targeted African Americans; and the remaining targeting politics and using the term
74 'retarded'.

75 We notice a reduction of overall missed classifications over all groups and a removal of missed hate
76 speech tweets targeting Latinos and Asians. This provides us with a starting point of inspecting the
77 ethnic group feature that focuses on Latinos and Asians (although also African Americans).

78 Of the 11 instances that were classified in the neither class: 36% African American with one including
79 threat of violence; 27% gay community with one including encouragement of suicide and the rest
80 emasculation of males; the remaining targeting Chinese, and White people.

81 We notice that the number of missed instances targeting African Americans, White people, and
82 Asians rises showing a bias of this model to classify as neither. It should be noted that there are no
83 missed instances targeting females.

84 **0.4 Hand-Built Features with News, BERT, and Twitter Embedding model Class: Offensive**

85 Out of the 403 Offensive class using the Hand-Built Features with News, BERT, and Twitter Em-
86 bedding model, we disagree with 45% of the tweets and believe they are incorrectly labeled, 7% of
87 which our model also incorrectly predicted the label. The remaining 28% of the tweets were correctly
88 predicted by our model.

89 Of the 28% tweets where we believe our model correctly predicted the label, 68% of the tweets
90 should be labeled as hate speech. Within these, the targeted communities are: 26% gay community;
91 11 instances where it was used as a means to emasculate men and the remaining targeting the gay
92 community in general

93 22% African American; this continues to be difficult as variations of 'nigga', 'nigguh', 'nig', 'niglet'
94 are used both as an in-group and by other parties as part of an insult. 15% female; three instances
95 including violence / sexual violence. 14% male; 3 emasculating. 8% used the term 'retarded';
96 a possible feature capturing instances of [What / He's][a][retard/ed] to weigh them towards the
97 offensive class could reduce the errors. 8% White people. 4% Latinos; an interesting insight as that
98 all instance in the missed tweets mentioning Latinos were either offensive or hate speech. A feature to
99 capture more of these instances would be to decode masked hate words, e.g. 'buck all the beaners', in
100 fact almost all tweets targeting Latinos included a variation of the term beaner so a feature checking
101 for the term along with swear words would properly flag it as hate. The remaining 3% were generally
102 offensive

103 32% should be labeled as class neither. Our model appropriately captured self-referential statements
104 and usage of terms that were used in a self-affirmative manner, such as:

105 'RT @kaitlyn_jardi: "@17Seniors: so i basically become a fearless bitch when i'm mad"'

106 'RT @G0ldenG0ddess: Turn up about to be real , marriott with my bitches for the weekend , mansion
107 tonight , adult swim tomorrow 128131;'

108 Of the 7% of tweets that were both incorrectly labeled and incorrectly predicted by our model and
109 whose class should be hate speech, there were the following instances targeting: 3 female, 3 African
110 American; 1 Latino; 1 using the term 'retarded'.

111 Of the missed predictions to true class offensive, the tweets were not targeted and were said as a
112 statement as opposed to an attack. The terms 'hoe', 'pussy', and variations of 'nigga' were common.
113 Making use of our targeted features could help to capture these tweets by toggling the targeted to off.

114 **0.5 News, BERT, and Twitter Embedding model Class: Offensive**

115 Of the 331 missed tweets in the offensive class, we only disagree with the labeling of 8 tweets. Five
116 targeted males, showing a higher bias towards labeling male offensive speech as hate speech. 3
117 targeted the gay community and 1 was generally racist.

118 48% of the tweets were incorrectly predicted as neither; about 90% of the missed tweets referenced
119 offensive speech towards women. This provides a case for our hand built feature that checks if a
120 tweet is offensive to women and we look forward to doing additional manipulation of combining our
121 features to improve performance.

122 52% predicted as hate speech; the overwhelming majority were offensive to women and then African
123 Americans showing a bias our model has towards labeling offensive speech targeting these groups as
124 hate speech. The tweets also included heavy usage of slang, hinting that our slang decoder does help
125 in contextualizing the tweet to capture more information.

126 **0.6 Hand-Built Features with News, BERT, and Twitter Embedding model Class: Neither**

127 26% of the 99 missed tweets in the class neither should have alternate labels. Our model correctly
128 predicted the label for 21% of the tweets with the remaining 5% both wrongly labeled and incorrectly
129 predicted by our model.

130 Our model incorrectly classified 41% of tweets as offensive. The upside is that not many included
131 slang outside as terms 'nig' and variations thereof, meaning that our slang decoder seemed to help
132 minimizing previous biases that labeled tweets with slang as offensive. The majority of the tweets

133 included pronouns, which may have triggered our pronouns checker feature to label these as offensive.
134 We could fine tune that feature by noting if the tweet is a question, that it may not be offensive.
135 Checking to see if their tweet is commentary that flips the negative sentiment could be helpful in
136 correctly labeling tweets as neither, e.g. 'RT @MobJoe: Word. And it don't make u a hoe RT
137 @100granHman: It's okay to have sex on first date long as the feeling is mutual'.

138 More concerning is that our model labeled 32% of the missed tweets in the neither class as hate
139 speech. Many of the tweets combined the term 'trash' with a noun. Capturing the term 'white trash'
140 could down weigh and other instances of the term trash just by itself. The use of the term 'Jihadis'
141 created a strong bias towards labeling the tweet offensive, even when in the context of reporting news.
142 We could create a feature where it searched for the term 'Jihad' along with profanity to differentiate
143 it from the term 'Jihad' just by itself.

144 **0.7 News, BERT, and Twitter Embedding model Class: Neither**

145 We disagree with the labeling of 47% of the missed tweets in the neither class. And of these we believe
146 our model correctly predicts 65% of these tweets and 35% of these tweets were both incorrectly
147 labeled and incorrectly predicted by our model.

148 Our model correctly identified 19 instances of hate speech targeting: 26% female; 26% male; 21%
149 gay community; 11% African Americans; and the remaining 16% evenly distributed targeting Asians,
150 general racist, and White people.

151 Of the instances where both the labeling and predictions were incorrect, 7 were hate speech targeting:
152 36% African Americans; 36% female; 28% gay community

153 Of our model's incorrect predictions of neither into the offensive class, the model failed to pick up
154 on self-referential and in-group tweets highlighting the need for these hand built features. For the
155 tweets that were incorrectly predicted as hate speech, the NRC emotions feature may help capture
156 more nuanced information about the tone of the tweet.

157 **0.8 Summary**

158 Overall we hope to have highlighted how difficult it is to gauge model performance when dealing
159 with a dataset that is almost 2/3 offensive and with which we feel a great disagreement in the labeling
160 process. We look forward to continuing our work in this space as we have just received academic
161 research API access from Twitter and plan to work on creating new labeled multi-class datasets
162 available to everyone and to test our current and future models. (Davidson et al., 2017) used a Logistic
163 Regression model, as did we. This task would benefit from other models (LSTMs, NNs) being run on
164 larger, more accurately labeled datasets, as Wang (2018) began to explore.

165 **References**

166 Thomas Davidson, Dana Warmusley, Michael Macy, and Ingmar Weber. 2017. Automated hate speech
167 detection and the problem of offensive language. *International AAAI Conference on Web and
168 Social Media*.

169 Cindy Wang. 2018. [Lexicon integrated deep neural networks for fine-grained hate speech detection
170 on twitter](#).