

Detecting threats of violence in YouTube comments

Anonymous NAACL submission

Abstract

Here comes an abstract for this article

1 Introduction

Threats of violence is an increasingly common occurrence in online discussions. It disproportionately affects women and minorities, often to the point of effectively eliminating them from taking part in discussions online. Moderators of social networks operate on such a large scale that manually reading all posts is an insurmountable task. Methods for automatically detecting threats could therefore potentially be very helpful, both to moderators of social networks, and to the members of those networks.

In this article, we evaluate different types of features for the task of detecting threats of violence in YouTube comments. We show that ...

2 Previous work

There is little previous work specifically devoted to the detection of threats of violence in text, however, there is previous work which examines other types of closely related phenomena, such as cyberbullying and hate-speech.

Dinakar et al. (2011) proposes a method for the detection of cyberbullying by targeting combinations of profane or negative words, and words related to several predetermined sensitive topics. Their data set consists of over 50,000 YouTube comments taken from videos about controversial topics. The topics included sexuality, race, culture and intelligence. The comments were grouped by topic, and then 12 % of these were manually annotated to

check that they were placed in the right category. The first stage of the detection was the same across all categories. It consisted of using a lexicon of negative words and a list of profane words, as well as part-of-speech tags from the training data that were correlated with bullying. The second stage was category-specific, and used commonly observed uni- and bigrams from each category as features. The experiments reported accuracies from 63 % to 80 %, but did not report precision or recall.

There has been quite a bit of work focused on the detection of threats in a data set of Dutch tweets (Oostdijk and van Halteren, 2013a; Oostdijk and van Halteren, 2013b). The data set used for these experiments consisted a collection of 5000 threatening tweets collected by a website over a period of about two years. In addition, a large number of random tweets were collected for development and testing. A set of 2.3 million random tweets was used for development, and a set of 1 million was used for testing. This set was not annotated in any way prior to being used in the experiments. The system relies on manually constructed recognition patterns, in the form of n-grams (uni-, bi- and trigrams, as well as skip bi- and trigrams), but do not go into detail about the methods used to construct these patterns, stating that the researchers relied on their (linguistic) intuition as speakers of Dutch (Oostdijk and van Halteren, 2013a). In Oostdijk and van Halteren (2013b), a manually crafted shallow parser is added to the system. This improves results to a precision of 39% and a recall of 59%.

Warner and Hirschberg (2012) present a method for detecting unwanted or illegal comments in user-

generated web text from the internet, which relies on machine learning in combination with template-based features. The data set used in the research came from two sources. The first consists of posts from Yahoo news groups, the second were web-pages collected by the American Jewish Congress that had been identified as offensive. The data set was manually annotated, and then the hate speech was assigned to a category, such as antisemitic, anti-woman, anti-asian, etc. The research then focused on the antisemitic category. The research approaches the problem as a word-sense disambiguation task, since the same words can be used in both hateful and non-hateful contexts. The features used in the classification were combinations of uni-, bi- and trigrams, part-of-speech-tags and Brown clusters. The best results of the classifications were obtained using only unigrams as features, with a precision of 67 % and a recall of 60 %. The other feature sets garnered much lower results, and the authors suggest that deeper parsing could reveal significant phrase patterns.

3 The YouTube threat data set

The YouTube threat data set is comprised of user-written comments from eight different YouTube videos (Hammer, 2014). A comment consists of a set of sentences, each of them manually annotated to be either a threat of violence (or support for a threat of violence) or not. The data set furthermore records the username¹ of the user that posted the comment. The eight videos that the comments were posted to cover religious and political topics like halal slaughter, immigration, Anders Behring Breivik, Jihad, etc. (Hammer, 2014).

The data set consists of 9,845 comments, comprised of 28,643 sentences, see table 1. In total there are 402,673 tokens in the sentences in the data set. There are 1,285 comments containing threats, and 1,384 sentences containing threats, as seen in table 1. (Hammer, 2014) report inter annotator agreement on this data set to be 98 %, as calculated on 120 of the comments, doubly annotated for evaluation.

¹In 2013, YouTube changed its commenting system from using unique usernames, to using "real names", like Facebook and other sites (YouTube, 2013). Some accounts, however, did not provide real names, so they continue to only be identified by their usernames.

	Comments	Sentences	Users posting
Total	9,845	28,643	5,483
Threats	1,285	1,384	992

Table 1: Number of comments, sentences and users in the YouTube threat data set

Figure 1 provides some examples of comments containing threats of violence taken from the data set. The first line is the username or name, and the subsequent lines are the sentences of the comment. An empty line indicates the end of a comment. The sentences are annotated with a number indicating whether they contain a threat of violence (1), or not (0).

4 Experiments

4.1 Experimental setup

4.2 Results

5 Discussion

References

- Karthik Dinakar, Roi Reichart, and Henry Lieberman. 2011. Modeling the detection of textual cyberbullying. In *Proceedings of The Social Mobile Web*.
- Hugo Lewi Hammer. 2014. Detecting threats of violence in online discussion using bigrams of important words.
- Nelleke Oostdijk and Hans van Halteren. 2013a. N-gram-based recognition of threatening tweets. In *Proceedings of Computational Linguistics and Intelligent Text Processing*, pages 183–196. Springer.
- Nelleke Oostdijk and Hans van Halteren. 2013b. Shallow parsing for recognizing threats in Dutch tweets. In *Proceedings of IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining*, Niagara, Canada. ACM.
- William Warner and Julia Hirschberg. 2012. Detecting hate speech on the world wide web. In *Proceedings of the Second Workshop on Language in Social Media*, pages 19–26. Association for Computational Linguistics.
- YouTube. 2013. We hear you: Better commenting coming to YouTube. youtube-global.blogspot.no/2013/09/youtube-new-comments.html/. [Online; accessed May 12th 2015].

192		240
193		241
194		242
195		243
196		244
197		245
198		246
199		247
200		248
201		249
202		250
203		251
204		252
205		253
206		254
207		255
208		256
209		257
210	timpa666	258
211	1 and i will kill every fucking muslim and arab!	259
212		260
213	NimsXdimensions	261
214	0 Need a solution?	262
215	1 Drop one good ol' nuke on that black toilet in Mecca.	263
216		264
217	Ammar Alozaibi	265
218	1 Funny, We will conquer you all in just few years, U will be	266
219	my slave and your women will be my Sex Toy in Bed.	267
220		268
221	Figure 1: Examples of comments from the data set.	269
222		270
223		271
224		272
225		273
226		274
227		275
228		276
229		277
230		278
231		279
232		280
233		281
234		282
235		283
236		284
237		285
238		286
239		287