

Detecting threats of violence in online discussions

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 big 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 ...

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.

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.

¹In 2013, YouTube changed its commenting system from using unique usernames, to using "real names", like Facebook

	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

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 the YouTube threat 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.

Figure 1 contains examples of some 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).

and other sites (YouTube, 2013). Some accounts, however, did not provide real names, so they continue to only be identified by their usernames.

096	timpa666	144
097	1 and i will kill every fucking muslim and arab!	145
098		146
099	NimsXdimensions	147
100	0 Need a solution?	148
101	1 Drop one good ol' nuke on that black toilet in Mecca.	149
102		150
103	Ammar Alozaibi	151
104	1 Funny, We will conquer you all in just few years, U will be	152
105	my slave and your women will be my Sex Toy in Bed.	153
106		154
107	LegitZombieSlayer1	155
108	0 As long as i'm alive you'll have no victory.	156
109	1 I'll kill all you cunts	157
110		158
111	Figure 1: Examples of comments from the data set.	159
112		160
113		161
114	4 Experiments	162
115	4.1 Experimental setup	163
116	4.2 Results	164
117		165
118	5 Discussion	166
119		167
120	References	168
121	Hugo Lewi Hammer. 2014. Detecting threats of violence	169
122	in online discussion using bigrams of important words.	170
123	YouTube. 2013. We hear you: Better com-	171
124	menting coming to YouTube. youtube-	172
125	global.blogspot.no/2013/09/youtube-new-	173
126	comments.html/. [Online; accessed May 12th	174
127	2015].	175
128		176
129		177
130		178
131		179
132		180
133		181
134		182
135		183
136		184
137		185
138		186
139		187
140		188
141		189
142		190
143		191