

Using Sentiment Analysis to Identify Bullying Using Twitter

Hoanh Nguyen

University Massachusetts Lowell
1 University Ave
Lowell, MA 01854, USA
soujirobot@gmail.com

Michael Meding

University Massachusetts Lowell
1 University Ave
Lowell, MA 01854, USA
mikeymeding@gmail.com

Abstract

Twitter is a social network where users can communicate publicly with short text statements called tweets. However, it should come as no surprise that not all of these tweets have good intention. With the rise of social media usage in youth, Twitter has gone from a calm social environment to connect with others to a hostile place. This program identifies users who have a high likelihood of being a bully through the use of sentiment analysis and machine learning. We did this by taking data from Twitter and analysing the aggressiveness of a tweets and hashtags using both the SentiWordNet database, Harvard Inquirer, and Ark-TweetNLP. This data is then run through a Machine Learning algorithm to determine if a user can be classified as a bully based on the sentiment of the words used and based on previously seen tweets.

1 Introduction

What is bullying? Bullying is when a person is constantly being exposed to negative actions by another individual or group of individuals and can involve both physical and mental abuse. Some common forms of mental bullying include name calling, making threats, and spreading rumours. Since the creation of social media sites, bullying has become a problem outside of the school environment and is now considered a serious national health issue among adolescents [0]. Bullies are using the internet as an outlet to post negative comments about someone, and unlike in a school environment, they

rarely see backlash for what they post on social media sites.

This program identifies Twitter users who are bullying those around them. To accomplish this task, the program trains a binary classifier to decide if a single tweet is bullying or not. However, further research has shown that bullying cannot fully be defined by a simple binary "yes" or "no". Bullying, in a broad sense, can be classified into many different categories beyond that of just the bully and the victim. For example, when observing a bullying in action on Twitter, you will notice that there are a lot of supporting users who do not have anything to do directly with a given bullying situation. This leads to a large amount of ambiguity when trying to decide whether someone is a bully, defender, victim, or accuser.

For a tweet to be cleanly identified by a computer as being bullying it must have both an aggressive statement and a subject or direction towards that anger. One of the major flaws with this is that often times the subject of an aggressive tweet is implied, thus making it very difficult for a machine to identify a subject. Another problem is sarcasm. Computers cannot accurately identify sarcasm using only text and unfortunately bullies often use sarcasm. Our solution to this problem was to use an assisted machine learning algorithm. The data sets were manually annotated by a human who could only use the information contained within a given tweet to decide if that user is a bully or not. By using a brute force approach, we achieved a high accuracy when tweets directly mention a subject. Given the small data set,

the annotated results were surprising.

2 Related Work

WSIC Two years ago, the University of Wisconsin completed a project on the study of bullying in Twitter that heralded lots of media coverage, including several articles from *The Huffington Post* and *Time Magazine*. The University of Wisconsin’s code used a set of static words as search terms for tweets, then classified the tweets using a very small pre-existing bullying model. This model did yield some results but the efficacy left much to be desired. The first version of the University’s code was made public and can easily be seen that it is quite primitive by comparison to our project. Unfortunately, the university didn’t release the data sets that they used for the projects so the reported results could not be confirmed.

Psychology This project required looking at bullying through many different lenses. To be able to identify bullies you need to be able to understand them from a psychological standpoint. The paper *Sticks and Stones Can Break My Bones, But How Can Pixels Hurt Me?* [0] suggests that youth today have adopted a new type of bullying through social media. This new form of electronic communication makes it much easier for young users to say things that they would not normally say in person. Often times users will have a completely different persona online than in person.

3 Methodology

3.1 Machine Learning Methods

3.2 Data Sets

Overview The data for this project was gathered by a data mining algorithm. The algorithm uses a lexicon of Twitter hashtags and search terms that have been previously tagged with a sentiment score. Our program goes through the lexicon to find the words which have medium to high negative sentiment score. These words are then searched for individually and the top results are then stored for analysis. The results are then filtered for spam and relevancy before being saved. Based on the sentiment score of the search term, we are able to control the size of the output file because individual search

Emoticon	Associated Emotion
:~))	positive
;(negative
>:(negative
:-O	neutral

Table 1: Emoticon lexicon example.

Term	Senti Score	P-Count	N-Count
#disgusting	-4.801	3	589
#inferior	-4.839	1	204
#unworthy	-4.915	3	660
#ill	-4.928	35	7802
#fabulous	7.526	2301	2
#excellent	7.247	2612	3
#superb	7.199	1660	2
#perfection	7.099	3004	4

Table 2: Seed word hashtag sentiment example.

terms are limited to 40 tweets per query. Adjusting the threshold sentiment score allows us to control the number of seed words which result from those values thus limiting the size of the corpus. From there, the tweets are manually tagged bullying or not_bullying to create a standard for future data.

Lexicons This project included several lexicons which were used to develop some features. The unigram hashtag sentiment lexicon was used for the generation of aggressive search terms in order to achieve the best potential for bullying tweets. The lexicon was organized into four different categories: term, sentiment score, positive count, and negative count shown in Table 2. Positive and negative counts were calculated from the number of times that the term co-occurred with another matching positive or negative marker. The lexicon was pre-compiled over a large amount of twitter data to give good sentiment scores. Table 2 represents the structure of the data as it exists in the file. This data was supplied by external resources and was compiled for use with Twitter related Sentiment projects. From this data we only cared about the terms which related to negative sentiment or those with a sentiment score below -4. This threshold value was adjusted to control the resulting seed word array as the lower the threshold sentiment score the less and less words would fit into

that category. Another useful lexicon for this project was the emoticon lexicon. This included emoticons such as :D and :(and their related sentiment in 3 categories: positive, negative, and neutral. Each of these categories was used as a feature in the identification of a bully. Unfortunately, this did not result in any significant change in accuracy as tweets that were manually classified as bullying would often include emoticons which related to positive sentiment skewing the results.

Corpus After going through the lexicon of negative sentiment words we end up with an array of negative words. This array now represents our search terms. Using the Twitter API, the program queries the top 40 English tweets one at a time for each word in the array of negative words. We are limited to a maximum of 40 tweets per request and 180 requests every 15 minutes. After each request has finished, the received tweets are filtered for relevancy and appended into an untagged corpus text file which is then ready for the final stage of corpus construction.

Annotation The final stage of building the data set for this project is manual tagging. After exhausting the array of negative search words, we end up with a large file of tweets which have a high potential of being tagged bullying due to the large amount of negative words. As this program uses an assisted machine learning algorithm, it requires a basis of comparison to verify the tagging efficacy. Therefore, it is required for someone to manually go through each tweet and tag each as either bullying or not bullying. This, as you can imagine, is exceptionally time consuming and has required the combined efforts of Alison Rose, Uwe Meding, and Michael Meding.

3.3 Evaluation

4 Results

Machine Learning

Corpus Constructing the corpus was a multi step process which had varying results at each step. First, seed word collection from the existing lexicon yielded just under 2500 words using a threshold sentiment score of -4. Examples of this lexicon can be seen in Table 2. After manual review some interesting words came back such as ipad2 which had a negative sentiment score of -4.999. After feeding

User	Tweet	Tag
@AprilH1...	@track_maniac17 I hate Brynn. She's like 50 and has like 10 baby daddies #golddigger	bullying
@emilybr...	And the award for "Tackiest Couple" goes to.. (link)	bullying
@Rachael...	@KTHopkins you really are something else. #vile	bullying
@MrsDiva...	Ur worse than screen doors on a submarine. #worthless	bullying
@chuckpa...	@stlcountypd You are a disgrace. Thanks for proving that tonight. Unreal. #immature	bullying

Table 3: Bullying tweet examples. (usernames shortened for clarity)

the 2500 words into the twitter search algorithm we ended up with just under 25000 tweets to be analysed. As of this writing we have managed to manually tag almost 4000 tweets giving us a tagging accuracy of around 30% on unseen testing data. An example of what this data looks like is shown in Table 3. This table also showcases what this program will identify as bullying tweets. This is an improvement from training on 700 tweets and seeing less than 20% accuracy. With more tagged tweets we expect to see additional improvement. The difference in results can be seen comparatively in Table 4 and Table 5. These two tables show the confusion matrix for two separate runs. One run using only 700 tweets as training data the other used 3500 tweets as training data. Resultant values are shown as counts over a smaller test corpus.

5 Future Improvements

Non Binary Classifier This project only trains a binary classifier for bullying. Unfortunately, bullying is a complex issue and has many more parts than

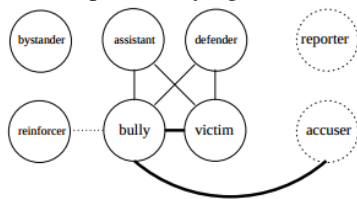
	Bullying	Not_Bullying
Bullying	5.0	11.0
Not_Bullying	12.0	103.0

Table 4: 700 tweets confusion matrix.

	Bullying	Not_Bullying
Bullying	1.0	27.0
Not_Bullying	18.0	502.0

Table 5: 3500 tweets confusion matrix.

Figure 1: Graph of bullying and its relations.



just bullying and not bullying. Often, someone who may appear to be a bully is simply reporting that someone else is a bully. Additionally, users who are bullies often have others who support them but who they themselves are not actually bullies. For this reason it is logical to improve the classifier to include tags for supporters, victims, bullies, and reporters. An approach which was taken by a similar project was to attempt to construct a graphical bullying map with included bullies and victims and all supporting roles. The more complete that the graph was the higher the likelihood that bullying is occurring. A visual for this can be seen in Figure 5 [0].

Confirmation One advancement that can be made for this project is to add another layer of evaluation known as confirmation. This next layer would take the list of tweets which had been tagged as bullying by the initial evaluation and run them through another algorithm. This algorithm will then search for only that particular users tweets and run those through the analyser as well. If the second analysis comes back with several more bullying tweets then we can confirm that this person is a bully. This would improve on the fact that occasionally people have a bad day and post negative things online. However, if someone is doing it all the time it gives us proof that this person is a bully.

Figure 2: The difference between emoticon and emoji.



Emojis Twitter now has support for iPhone Emojis. Emojis are similar to emoticons but they are far more expressive than simple smile and frown faces as shown in Figure 5. They use a format similar to html symbols using a unique unicode table representing faces and actions. Having a feature which could take into account the sentiment of Emojis would be a significant improvement to this project.

Seed Words Although the list of potential bullying words that are in use for this project is large it is still static. As bullying progresses in social media so does the associated vernacular. A feature which would improve the searching algorithm is one that will identify new seed words for future searches in the tweets which have been classified as bullying.

Acknowledgments

We thank Alison Rose, and Uwe Meding for their contributions in annotating our data set for this project. Without their efforts this project would not have been possible. Additionally, we would like to thank Willam Boag for providing us with lexicons from a similar project which proved to be immensely helpful.

References

- Jun-Ming Xu, Xiaojin Zhu, Amy Bellmore. Fast learning for sentiment analysis on bullying. pages.cs.wisc.edu/~jerryzhu/pub/wisdom12.pdf.
- Wanda Cassidy, Margaret Jackson, and Karen N. Brown. Sticks and stones can break my bones, but how can pixels hurt me?