

# BullyNet: Unmasking Cyberbullies on Social Networks

Aparna Sankaran, Hannah Johnson, Gaby G. Dagher, Min Long

**Abstract**—One of the most harmful consequences of social media is the rise of cyberbullying, which tends to be more sinister than traditional bullying given that online records typically live on the internet for quite a long time and are hard to control. In this paper, we present a three-phase algorithm, called BullyNet, for detecting cyberbullies on Twitter social network. We exploit bullying tendencies by proposing a robust method for constructing a cyberbullying signed network. We analyze tweets to determine their relation to cyberbullying, while considering the context in which the tweets exist in order to optimize their bullying score. We also propose a centrality measure to detect cyberbullies from a cyberbullying signed network, and we show that it outperforms other existing measures. We experiment on a dataset of 5.6 million tweets and our results shows that the proposed approach can detect cyberbullies with high accuracy, while being scalable with respect to the number of tweets.

**Index Terms**—Cyberbullying, Signed Networks, Social Media Mining.

## I. INTRODUCTION

The Internet has created never before seen opportunities for human interaction and socialization. In the past decade, social media, in particular, has had a popularity explosion. From MySpace to Facebook, Twitter, Flickr, and Instagram, people are connecting and interacting in a way that was previously impossible. The widespread usage of social media across people from all ages created a vast amount of data for several research topics, including recommender systems [1], link predictions [2], visualization, and analysis of social networks [3].

While the growth of social media has created an excellent platform for communications and information sharing, it has also created a new platform for malicious activities such as spamming [4], trolling [5], and cyberbullying [6]. According to the Cyberbullying Research Center (CRC) [7], cyberbullying occurs when someone uses the technology to send messages to harass, mistreat or threaten a person or a group. Unlike traditional bullying where aggression is a short and temporary face-to-face occurrence, cyberbullying contains hurtful messages which are present online for a long time. These messages can be accessed worldwide, and are often irrevocable. Laws about cyberbullying and how it is handled differ from one place to another. For example, in the United States, the majority of the states incorporate cyberbullying into their bullying laws, and cyberbullying is considered a criminal offense in most of them [8]. Popular social media platforms such as Facebook and Twitter are very vulnerable to cyberbullying due to the popularity of these social media sites and the anonymity that the internet offers to the perpetrators. Although strict laws exist to punish cyberbullying, there are very less tools available

to effectively combat cyberbullying. Social media platforms provide users with the option to self-report abusive behavior and content in addition to providing tools to deal with bullying. For example, Twitter has features that include locking accounts for a brief period of time or banning the accounts when the behavior becomes unacceptable. The body of work produced by the research community with regards to cyberbullying in social networks also needs to be expanded to get better insights and help develop effective tools and techniques to tackle the issue.

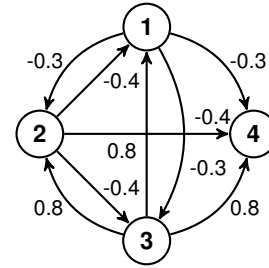


Figure 1: An example of a signed network.

To identify cyberbullies in social media, we first need to understand how social media can be modeled. The common way of modeling relationship in social psychology [9] is to represent it as a signed graph with positive edge corresponds to the good intent and negative edge corresponds to malicious intent between people. Using the signed graph, we model the Twitter social network as a signed network to represent users' behavior [10] where nodes correspond to users, directed edges correspond to communications and/or relations between the users with assigned weight in the range  $[-1, 1]$ , as illustrated in Figure 1.

**Definition 1:** A signed social network (SSN) is a directed, weighted graph  $G = (V, E, W)$ , where  $V$  is the set of users and  $E \subseteq V \times V$  is the set of edges with an edge weight  $w \in W$  in the range of  $[-1, 1]$ .

Mining social media networks to determine cyberbullies imposes several challenges and concerns. First, it is typically hard to accurately interpret user's intentions and meanings in social media based merely on their messages (e.g. posts, tweets, comments), which are typically short, use slang languages, or may include multimedia contents such as pictures and videos. For example, Twitter limits its users' messages to 140 characters, which could be a mix of text, slangs, emojis, and gifs. As a result, it is hard to determine the opinion expressed by a message correctly. For this we utilize sentiment analysis [11], [12] to determine whether the user's attitude towards other users are positive, negative, or neutral. Second, bullying could

be hard to detect if the bully chooses to disguise it through techniques such as sarcasm or passive-aggression. In this situation, a single text (message) cannot determine the user's intention. So, we collect the entire conversation between two or more users to identify the context in which the user attitude exists. Third, the large size and dynamic and complex structure of social media networks makes it challenging to identify cyberbullies. For example, on Twitter, hundreds of millions of tweets are sent every day on the social network platform. In this case we construct the social network as a graph and assign value based on the maliciousness of the user. Because the network analysis reduce the complex relationship between the users to the simple existence of nodes and edges [10] There are several works in the literature concerning detecting malicious users from unsigned networks with positive edge weights, including community detection [13], node classification [14] and link prediction [2]. On the other hand, methods that analyze signed social networks are scarce [15].

In this paper, we study the problem of cyberbullying in social media in an attempt to answer the following research question: *Can tweet contexts (conversations) help improve the detection of cyberbullying in Twitter?*. Our intuition is that each tweet should be evaluated not only based on its contents, but also based on the context in which it exists. We call such a context a conversation, which is a set of tweets between two or more people exchanging information about a certain subject. Thus, our solution consists of three parts. First, for each conversation, a conversation graph is generated based on the sentiment and bullying words in the tweets. Second, we compute the bullying score for each pair of users in a conversation graph, and then combine all graphs to create an SSN called bullying signed network ( $\mathcal{B}$ ). The inclusion of negative links can bring out information that would otherwise be missed with only positive links [16]. Finally, we propose a centrality measure called attitude & merit ( $A\&M$ ) to detect bullying users from the signed network  $\mathcal{B}$ .

Our main contributions are organized as follows:

- 1) Collected, preprocessed and labelled the Twitter dataset.
- 2) Proposed a novel efficient algorithm for detecting cyberbullies on Twitter.
  - a) Built conversation.
  - b) Constructed Bullying Signed Network.
  - c) Proposed Attitude and Merit Centrality.
- 3) Experimented on 5.6 million tweets collected over 6 months. The results show that our approach can detect cyberbullies with high accuracy, while being scalable with respect to the number of tweets.

## II. RELATED WORKS

In this section, we review the literature on areas related to cyberbullying detection and signed social networks.

### A. Cyberbullying Detection

There isn't a lot of works in the literature that utilizes signed networks to detect cyberbullies. The papers [6], [17], are aimed at detecting trolls in a signed network. Wu et

al. [17] proposed a method for ranking nodes to identify trolls without using a PageRank algorithm. Kumar *et al.* [6] proposed an iterative algorithm involving new decluttering operations and various centrality measures to detect trolls. Unlike the proposed method in this paper, the authors begin their process with an already created signed network.

A significant amount of work has been done over the past decade in the area of cyberbullying detection in general. There have been two broad methods in identifying bullies - one aims to detect bullying messages (e.g. [18]–[21]), while the other approach is to detect the cyberbullies responsible for the messages (e.g. [22]–[25]).

The first method of determining bullying messages was done using a combination of text-based analytics and a mix of text and user features. Zhao *et al.* [18] proposed a text based Embeddings-Enhanced Bag-of-Words (EBoW) model that utilizes a concatenation of bullying features, bag-of-words, and latent semantic features to obtain a final representation, which is then passed through a classifier to identify cyberbullies. Xu *et al.* [21] used textual information to identify emotions in bullying traces, as opposed to determining whether or not a message was bullying. Singh *et al.* [19] proposed a probabilistic socio-textual information fusion for cyberbullying detection. This fusion uses social network features derived from a 1.5 ego network and textual features, such as density of bad words and part-of-speech-tags. Hosseinmardi *et al.* [20] used images and text to detect cyberbullying incidents. The text and image features were gathered from media sessions containing images and the corresponding comments, which was then fed into various classifiers. Chen [25] proposed a novel method in identifying cyberbullies within a multi-modal context. To understand cyberbullying Kao *et al.* [26] proposed a framework by studying social role detection. By using words and comments, temporal characteristics, and social information of a session as well as peer influence Cheng *et al.* [27], [28] proposed frameworks for detecting cyberbullies.

The second method was aimed at identifying the person behind the cyberbullying incidents. Squicciarini *et al.* [22] used MySpace data to create a graph, which integrated user, textual, and network features. This graph was used to detect cyberbullies and predict the spreading of bullying behavior through node classification. Galán-García *et al.* [23] used supervised machine learning to detect the real users behind troll profiles on Twitter, and demonstrated the technique in a real case of cyberbullying. In a recent paper on aggression and bullying in Twitter, Chatzakou *et al.* [24] found cyberbullies and aggressors using user, text, and network-based features.

From the above methods we determined that these approaches focus on how offensive the content of the message is and based on that they identify cyberbullies but does not consider why the message was offensive i.e., the above papers do not analyze the context of the entire conversation just the content of the message. Our approach utilizes the Bag-of-words with the text to identify curse words, use sentiment analysis to determine the emotions or attitude of the sender and finally we analyze the entire context in which the sender and receiver communicates. These overlooked factors could significantly or completely change the results of Cyberbullying

Detection.

### B. Signed Social Networks

This section reviews the previous work done on signed networks (e.g. [6], [10], [15], [17], [29]). The idea of signed networks is not new but its application and analysis of them were only developed in recent years. We extended its application to establish node classification in our model. Previously, in 2010, Leskovec *et al.* [10] reviewed the balance and status theory and their relation to social media and proposed a modified status theory that better reflects patterns found in signed networks in social media. Tang *et al.* have done a broad survey of signed networks in social media [15] and proposed a new framework for node classification in signed social networks [29]. The authors incorporated negative links in the signed network and proposed an approach to mathematically model both independent and dependent information from the links.

Over the last few years, a number of methods have been designed for signed network analysis with both positive and negative links [30]–[33]. Most of these methods are based on simple modifications of the PageRank or Eigenvector centrality that accounts for negative weights on the links. However, some of these measures do not consider, how the incoming edges of a node depends upon the outgoing edges from the same node and vice versa i.e., interactions between incoming and outgoing links in a signed networks. Mishra *et al.* [34] employ this scenario and proposed bias and deserve measures. The deserve of a node depends on the opinions of other nodes whereas the trustworthiness of a node depends on how a node gives a correct opinion about other nodes. From the experiment in the Section VI-D, we can find the bias and deserve measures is not effective for identifying bullies in the network.

### III. PROBLEM FORMULATION

In this section, the Twitter social network is represented as a directed, weighted graph  $G = (U, E)$  with  $U$  being the set of users (represented as nodes) and  $E$  being the set of tweets  $T$  sent between the users (represented as edges). Each user  $u \in U$  has a set of features including an ID, the number of followers, the number of friends, and the number of the tweets that they sent.

Each tweet  $t \in T$  is associated with certain features: source ID ( $SID$ ), destination ID ( $DID$ ), the date of creation, a user ID ( $UID$ ), a reply ID ( $RID$ ), and mentions ( $MID$ ). If the tweet includes mentions (i.e., if a given @username is included in a tweet anywhere else but at the very start), then Twitter interprets this as a mention and the user gets a notification that someone has mentioned them.

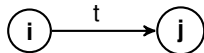


Figure 2: An example of a tweet.

As shown in the figure 2, the notation  $e_{ij}$  represents a tweet  $t$  directed edge from node (user)  $i$  to node (user)  $j$ . The existence of an edge  $e_{ij}$  denotes an interaction from node

$i$  to node  $j$  which is  $t$ . Each tweet has a set of features, as shown in Table II. ( $SID$ ) is assigned when a new tweet is created, in this case it is 101. ( $DID$ ) is an ID to which this current tweet is in response to where the destination ID is 3001. The ( $UID$ ), ( $RID$ ) and ( $MID$ ) corresponds to IDs of that particular users/nodes. Finally, the text is the content of the tweet sent from the node  $i$  to node  $j$ .

From the above Twitter data, we extract conversations and build a directed weighted graph for each conversation  $C = \{c_1, c_2, \dots, c_{|C|}\}$ . In our model, each  $c_i$  is a set of two or more tweets between two or more users:

**Definition 2:** A conversation  $c$  is a set of time-ordered tweets  $c = \{t_1, t_2, \dots, t_{|c|}\}$  such that:

- 1) The first tweet  $t_1$  is the *initiator* tweet that starts the conversation, and can be one of the two following types:
  - $DID(t_1) = \text{NULL}$ , and either  $MID(t_1)$  or  $RID(t_1)$  is not null.
  - $DID(t_1) \neq \text{NULL}$ , and  $\forall t \subseteq T : SID(t) \neq DID(t_1)$ .
- 2) All tweets in  $c$  satisfy the following:  
 $SID(t_i) = DID(t_{i+1}) : 1 \leq i \leq |c| - 1$ .

Our model will analyze the nodes and conversions and will output a list as results for detecting cyberbullies on the Twitter social network:

$$L = \{(u_1, s_1), (u_2, s_2), \dots, (u_{|L|}, s_{|L|})\},$$

where  $u_i$  is a user (node) and  $s_i$  is a confidence value for the likelihood of user  $u_i$  being a bully.

### IV. OUR SOLUTION: BULLYNET ALGORITHM

In this section, we first present an overview of the proposed three-phase bully finding algorithm and elaborate the steps in each phase.

The objective of our solution is to identify the bullies from raw Twitter data based on the context as well as the contents in which the tweets exist. Given a set of tweets  $T$  containing Twitter features such as user ID, reply ID etc, the proposed approach consists of three algorithms: (i) Conversation Graph Generation Algorithm, (ii) Bullying Signed Network Generation Algorithm, (iii) Bully Finding Algorithm. The first algorithm constructs a directed weighted conversation graph  $G_c$  by efficiently reconstructing the conversations from raw Twitter data while enabling a more accurate model of human interactions. The second algorithm constructs a bullying signed network  $\mathcal{B}$  to analyze the behaviour of users in social media. The third algorithm consists of our proposed attitude and merit centrality measures to identify bullies from  $\mathcal{B}$ . Figure 3 shows the process flow of BullyNet where the raw data is extracted from Twitter using Twitter API from which the conversation graph is constructed for each conversation using algorithm 1. Then from the conversation graphs, a bullying signed network is generated using algorithm 2. Finally, the bullies from Twitter are identified by applying algorithm 3.

#### A. Algorithm 1 - Conversation Graph Generation

The conversation graph generation algorithm 1, is constructed from a set of tweets  $T$  to generate directed weighted

Table I: Comparative evaluation of main features in related approaches including our proposed approach

Approach	Detect			Attributes based on				Signed Network		Dataset			
	Cyberbullying Message	User	Other	Content	Context	User	Network	Yes	No	Twitter	YouTube	Slashdot	Instagram
Zhao <i>et al.</i> [18]	●			●					●	●			
Xu <i>et al.</i> [21]	●			●					●	●			
Hosseinmardi <i>et al.</i> , [20]	●			●					●				●
Dadvar <i>et al.</i> , [35]	●			●					●		●		
Dinakar <i>et al.</i> , [36]	●			●		●			●		●		
Squicciarini <i>et al.</i> [22]		●		●	●	●	●		●	●			
Chen <i>et al.</i> [37]		●		●					●		●		
Galán-García <i>et al.</i> [23]		●		●		●			●	●			
Chatzakou <i>et al.</i> [24]		●		●		●	●		●	●			
Mishra & Bhattacharya [34]			●				●	●				●	
Kumar <i>et al.</i> [6]			●				●	●				●	
Wu <i>et al.</i> [17]			●				●	●				●	
Ortega <i>et al.</i> [38]			●				●	●				●	
Our proposed protocol		●		●	●		●	●		●			

Table II: Tweet features

<i>STD</i>	<i>DID</i>	<i>UID</i>	<i>RID</i>	<i>MID</i>	Text
101	3001	User1	User1	UserX, UserY	@UserX @UserY Lets meet at the central park

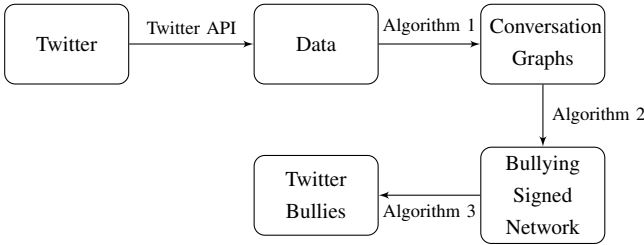
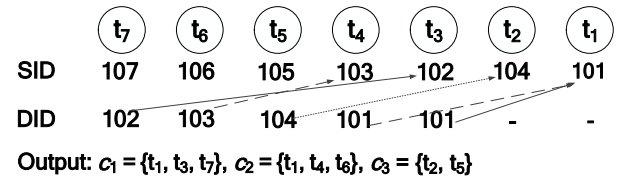


Figure 3: Protocol Flowchart of BullyNet.

conversation graphs  $G_c$  for each conversation. The weights between the nodes or users are determined by analyzing the sentiment behind the text of a tweet and examining for curse words. We then provide a score based on the expression the text represents. For each tweet  $t_i$  in  $T$ , the conversations are built by doing a binary search  $DID(t_i)$  with the  $STD$  of the remaining tweets. If a match is found as  $t'$ , then it is appended with  $t_i$  to form a new conversation. If a binary search match is found with an already existing tweet in a conversation  $c_i$  then,  $t_i$  is appended to tweets in  $c_i$ . The graphs are represented as  $G_c = (V, E, I)$  where  $V$  is the set of users involved in the conversation,  $E$  is the set of edges representing the tweets in the conversation, and each edge is assigned a bullying indicator value  $I$  as the edge weight which is in the range of  $[-1, +1]$ . When  $I_{ij} = -1$ , it indicates a negative interaction by  $i$  towards  $j$ , and when  $I_{ij} = 1$ , indicating a positive interaction. The bullying indicator for each tweet is computed as  $I = \beta * SA + \gamma * CS$ , based on sentiment analysis  $SA$  (VADER) and cosine similarities  $CS$  with a list of commonly used insulting words. The factors of  $\beta$  and  $\gamma$  are 0.9 and 0.1 respectively, which are determined by the experiment (see section VI-C).

**Example 1:** Figure 4 illustrates the conversation extracted from the set of tweets  $T = \{t_1, \dots, t_7\}$ . First, the tweets are sorted in descending order i.e.,  $t_7, t_6, \dots, t_1$ . Next,  $DID(t_7)$  is searched with the  $STD$  of the remaining tweets ( $t_6$  through  $t_1$ ). A match is found in  $t_3$  and conversation  $c_1$  is formed. This

Figure 4: Matching tweets based on  $DID$  and  $STD$  to construct conversations. Given tweets  $\{t_1, \dots, t_7\}$ , the output is three conversations:  $c_1$ ,  $c_2$ , and  $c_3$ .

process is repeated for each tweet. The conversations  $c_2$  and  $c_3$  are created with tweets  $\{t_6, t_4\}$  and  $\{t_5, t_2\}$  respectively. Since  $DID(t_4)$  and  $DID(t_3)$  match with the  $STD(t_1)$ , the tweet  $t_1$  is appended with  $t_4$  and  $t_3$ . So, the final conversations are  $c_1 = \{t_7, t_3, t_1\}$ ,  $c_2 = \{t_6, t_4, t_1\}$  and  $c_3 = \{t_5, t_2\}$ . This process can be represented by Algorithm 1.

In the step 3 of Algorithm 1, a directed, weighted graph  $g_{c_i} = (V, E)$  is constructed for every conversation  $c_i$  where nodes  $V$ , represented as the users, and the edges  $E$ , represented as the tweets, are directed from one user to another in a conversation. For every edge  $e$ , an edge weight is calculated as  $I = \beta * SA + \gamma * CS$ . This is known as the Bullying Indicator which is in range of  $[-1, +1]$ . The sentiment analysis ( $SA$ ) and cosine similarity ( $CS$ ) is computed on the tweet (edge) to evaluate the emotion and behaviour of the user. The  $\beta$  and  $\gamma$  are constants, which will be determined by the experiment (See section VI-C). In Step 4, the algorithm outputs the conversation graphs  $G_c$ .

Sentiment analysis ( $SA$ ) is the process of analyzing the sentiment of a message based on the user's opinion, attitude, and emotion towards an individual. Depending on the analysis, the polarity of the text is classified into positive, negative or neutral. The sentiment reflects feeling or emotion while emotion reflects attitude. There are different libraries or tools available to determine the sentiment of the content which includes sarcasm, emoji, images etc. Some of the them are: VADER, TextBlob, Python NLTK etc. We use VADER (Valence Aware Dictionary and sEntiment Reasoner) [39], which is a lexicon and rule-based sentiment analysis tool that is specifically attuned to sentiments expressed in social media. It performs well with emojis, emoticons, slangs and acronyms in



**Conversation Graph Generation Algorithm****Input:** Set of tweets,  $T = \{t_1, \dots, t_n\}$ **Output:** Conversation graphs  $G_c = \{g_{c_1}, \dots, g_{c_m}\}$ 

- 1) Sort all tweets in  $T$  in reverse-chronological order based on date of creation.
- 2) For each tweet  $t_i$  in  $T$ , where  $1 \leq i \leq |T|$ :
  - a) If  $t_i$  does not belong to a conversation, then create a new conversation  $c \in C$  and associate  $t_i$  with  $c$ .
  - b) If there is a tweet  $t' \in \{t_i, t_{i+1}, \dots, t_{|T|}\}$  where  $DID(t_i) = SID(t')$  then associate  $t'$  with all  $t_i$ 's conversations.
- 3) For each conversation  $c_i \in C$ :
  - a) Construct a conversation graph  $g_{c_i} \in G_c$ , where users are represented as nodes and tweets as edges.
  - b) For each edge  $e = (u, v)$  in  $g_{c_i}$ :
    - i) Compute the sentiment of the tweet (SA).
    - ii) Compute the cosine similarity of the tweet with bullying bag of words (CS).
    - iii) Calculate the bullying indicator  $I_{t_i}$  (weight) of the edge as follows:
 
$$I_{uv} = \beta * SA + \gamma * CS$$
- 4) Return  $G_c$

Algorithm 1: Conversation Graph Generation

a sentence. Cosine Similarity (CS) [40] measures the similarity between two vectors using their inner product. In the Twitter and some tweets may contain curse or insulting words that are reasonable indications of the existence of bullying. Thus, we select a reference list of insulting words commonly used in Twitter and some external linguistic resources for insulting analysis seeds. This list contains words indicating curse or negative emotions such as *nigga*, *bitch*, *slut* etc. and are compared with individual tweets with cosine similarity to compute a score. In this context, each tweet and insulting seeds are represented as vectors, where each vector has the word frequencies.

*Example 2:* Figure 5a shows a sample conversation of tweets. From algorithm 1, the conversation graphs are constructed, as shown in Figure 5b. It contains two conversation graphs shown with dashed blue edges and with solid red edges. The rounded number on the edges indicate the tweet order of that particular conversation. Figure 5c and 5d represent the two conversation graphs as  $g_{c_1}$  and  $g_{c_2}$  with the bullying indicator as the edge weight. With  $\beta$  and  $\gamma$  values as 0.9 and 0.1 respectively which was determined experimentally (see section VI-C), the edge weight  $I_{31}$  i.e., the edge from  $P3$  to  $P1$ , is calculated as -0.23. Similarly, the score of the other edges are calculated as shown in Figure 5c and 5d.

**B. Algorithm 2 - Bullying Signed Network Generation**

In many real-world social systems, the relation between two nodes can be represented as signed networks with positive and negative links. Since this research focuses on identifying the bullying nodes in the network, the algorithm 2 is designed to

determine the final outgoing edge weight,  $w_{ij}$  for the users in the conversation graphs  $G_c$ .

**Bullying Signed Network Generation Algorithm****Input:** Set of conversation graphs,  $G_c$ **Output:** Bullying Signed Network  $\mathcal{B}$ 

- 1) For each conversation graph  $g_{c_i}$  in  $G_c$ :
  - a) For each set of edges with the same order, sorted ascendingly, compute the bullying score of source node  $u$  toward target node  $v$  for each edge  $e = (u, v)$  as follows:
 
$$S_{uv} = I_{uv} + \alpha(I_{uv} - S_{vu}).$$
 and then determine the average score of node  $u$  for the same set of edges.
  - b) Compute the overall bullying score  $S$  of each node in  $g_{c_i}$  as follows:
    - i) If the node is the root node, then:
 
$$S = \frac{\sum S}{1 + 2.2(n-1)}$$
    - ii) Otherwise:  $S = \frac{\sum S}{2.2(n)}$
- 2) Construct the bullying signed network graph  $\mathcal{B}$  by merging all the conversation graphs together.
- 3) Return  $\mathcal{B}$ .

Algorithm 2: Bullying Signed Network Generation

In Step 1a of the algorithm 2, for every conversation graph  $g_{c_i}$ , a bullying score  $S$  is calculated based on how a node/user interact with other nodes/users in the graph based on the tweet order (sorted in ascending order) i.e., tweets are arranged based on the conversation. For an edge  $e = (u, v)$ , the bullying score  $S_{uv} \equiv I_{uv}$  if the edge towards  $v$  is not a reply from  $u$ . Otherwise, the bullying score  $S_{uv}$  is calculated as  $I_{uv} + \alpha(I_{uv} - S_{vu})$  where  $\alpha$  is a constant determined by the experiment as 0.6. Here,  $\alpha$  is used to calculate how much percent of the difference between the sender and receiver should be taken to determine the bullying score  $S$ .  $I_{uv}$  is the Bullying Indicator between the nodes  $u$  to  $v$  and  $S_{vu}$  is the Bullying Score between the nodes  $v$  to  $u$ . The difference between  $I_{uv}$  and  $S_{uv}$  is that,  $I_{uv}$  computes a score for the content on a tweet based on the sentiment analysis and cosine similarity, whereas  $S_{uv}$  computes a score based on the entire conversation between  $u$  and  $v$  i.e., the context in which opinion of  $u$  towards  $v$ . If there are more than one edge for a user with the same order, an average bullying score is computed for the same set of orders after the bullying score is evaluated.

*Example 3:* Table III shows the bullying score calculation for the conversation graph  $g_{c_1}$  in Figure 5c. In order 1, the bullying score  $S_{21} = I_{21} = 0$  since, the edge from  $P2$  to  $P1$  is not a reply edge. The user in the parenthesis represents to whom the edge responds. In order 2, there are two edges from  $P3$  to  $P1$  and  $P3$  to  $P2$  and the bullying score  $S_{31} = -61$  and  $S_{32} = -61$  is the same as  $I_{31}$  and  $I_{32}$  respectively. The order 3 also has two edges,  $P2$  to  $P3$  and  $P2$  to  $P1$ . Since the edge  $P2$  to  $P3$  is a reply to the edge  $P3$  to  $P2$ , the bullying score is calculated as  $S_{23} = I_{23} + \alpha(I_{23} - S_{32}) = 0.99$  where

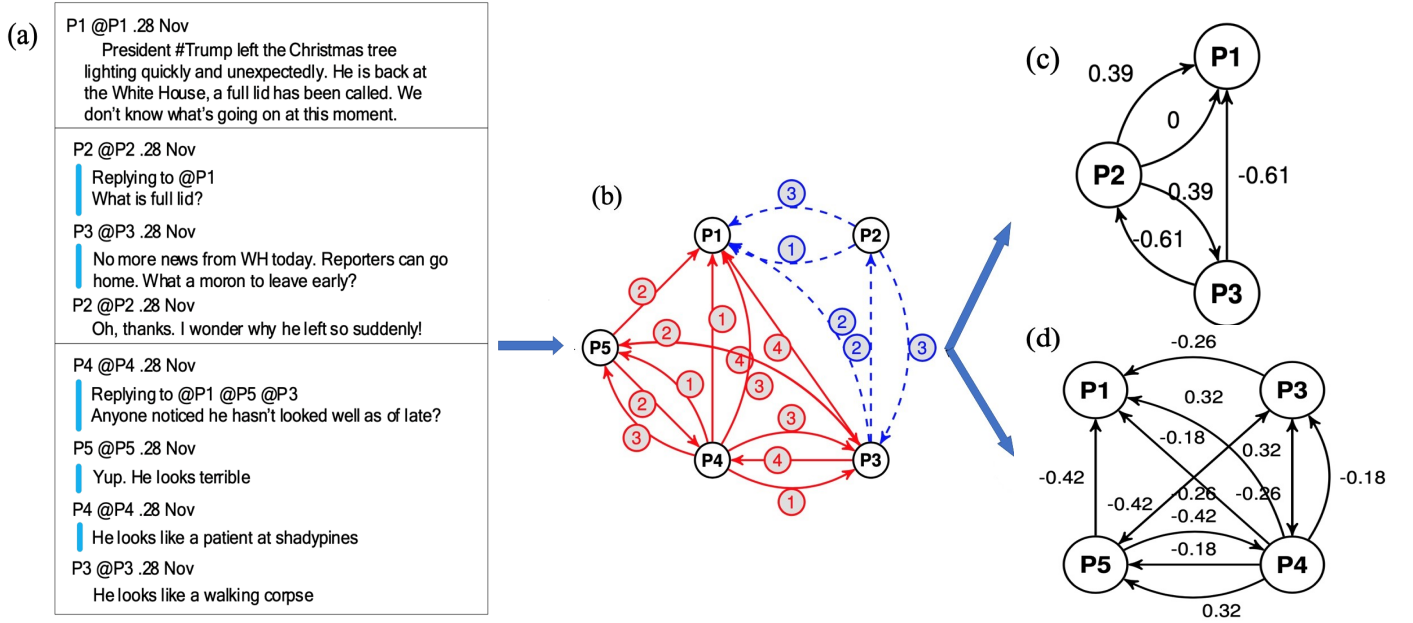


Figure 5: Conversation Graph Generation. (a) Sample conversation of tweets. (b) Conversations graph. blue and red sub-graphs represent 2 different conversations. (c)(d) Bullying indicators are added as edge weights to the conversation graphs.

$\alpha = 0.6$  was determined by the experiment. Next, the average of the score for the same order of the user is computed i.e., order 2 of the user  $P3$  is  $-0.61$  and order 3 of the user  $P2$  is  $0.69$ . Following a similar approach, the bully score  $S$  is calculated in Table IV for the second conversation graph  $g_{c_2}$  in Figure 5d.

Table III: Bullying score table for  $g_{c_1}$

Tweet #	P1	P2	P3
1	-	0 (P1)	
2	-		-0.61 (P1,P2)
3	-	0.99 (P3) 0.39 (P1)	
Total	-	0.69	-0.61

Table IV: Bullying score table for  $g_{c_2}$

Tweet #	P1	P4	P5	P3
1	-	-0.18 (P1,P3,P5)		
2	-		-0.56(P4) -0.42(P1,P3)	
3	-	0.84(P5) 0.32(P1,P3)		
4	-			-0.60(P4) -0.16(P5) -0.26(P1)
Total	-	0.4	-0.49	-0.34

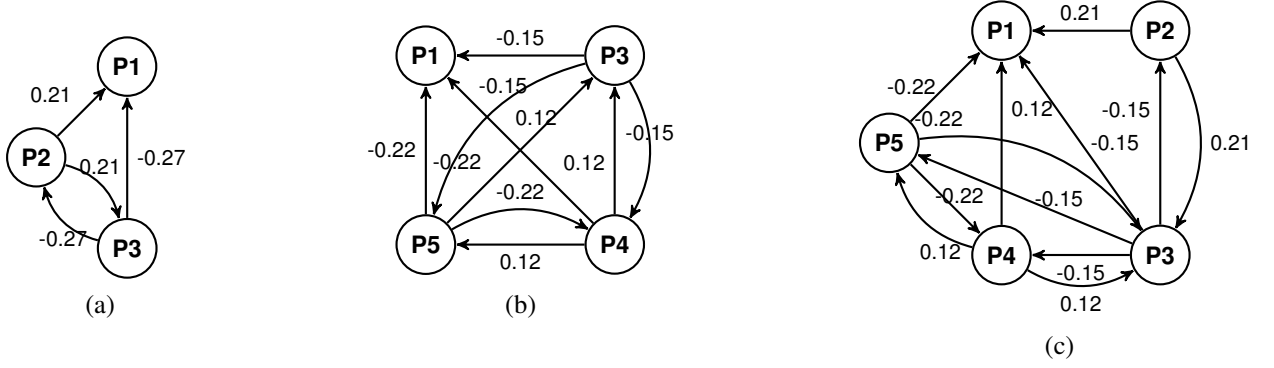
In Step 1b, the bullying score which was computed in the previous step for the users in every conversation graph  $g_{c_i}$  is normalized in  $[-1, 1]$ . The normalization is performed in two ways i.e., for the user that initiated the conversation, known as root nodes, and the users that are involved in the conversation. For the first type of users, the normalization is computed as  $\sum S / (1 + 2.2(n - 1))$  and for the second type of users as  $\sum S / 2.2(n)$  where,  $n$  is the number of times the user occurs in the order and the value 2.2 is computed using  $1 + (Maxdiff)(\alpha)$  in which  $Maxdiff$  is the range i.e., 2. This normalized score of the users becomes the edge weight to the other users in  $g_{c_i}$ .

In Step 2, the bullying signed network graph  $\mathcal{B}$  is constructed by merging all the conversation graphs  $G_c$ . If there is more than one edge i.e.,  $e = (u, v)$  then a single edge weight is calculated by taking the difference between average and standard deviation of all  $w_{uv}$ . Step 4 outputs the bullying signed network graph  $\mathcal{B}$ .

*Example 4:* Figure 6c illustrates the bullying signed network by merging the two normalized conversation graphs in Figure 6a and 6b. From Figure 6, it can be seen that there are two different edges from the user  $P3$  to  $P1$  ( $-0.27$  and  $-0.15$ ). So, the difference between the average and the standard deviation of the two edges are calculated as  $-0.15$  which is the final edge weight of  $P3$  to  $P1$  in the bullying signed network.

### C. Algorithm 3 - Bully Finding

This work, is to identify bullies from  $\mathcal{B}$  using centrality measures. Since this paper is about social networks the importance is defined as the behaviour. Among several centrality

Figure 6: Normalized conversation graphs (a)  $g_{c1}$  (b)  $g_{c2}$  and (c) Bullying signed network.

measures, we consider Bias and Deserve (BAD) by Mishra and Bhattacharya [34] a state of art method, that handles signed network because, their measure is computed on how the outgoing edge from a node/user depends on the incoming edges from other nodes/users. However, BAD is modelled on a trust based network i.e., the users that have a propensity to trust/distrust other users. Also, the edge weight denotes the trust score rather than the bullying score as in this research.

So, we proposed a centrality measure  $A\&M$  Attitude and Merit, similar to that of BAD to identify bullies from our proposed signed network  $\mathcal{B}$ . Merit is a measure of the opinion (good or bad) that the other nodes have towards a particular node and Attitude is a measure of the behaviour of a node towards the other node. However, in a given bullying signed network, the attitude or likes or dislike of a node towards other nodes in the network is not known. Therefore the expressions to compute the Merit and Attitude metrics in a mutually recursive manner.

$$M^{n+1}(j) = \frac{1}{2|in(j)|} \sum_{k \in in(j)} (w_{kj})(A^n(i)) \quad (1)$$

$$A^{n+1}(i) = \frac{1}{2|out(i)|} \sum_{j \in out(i)} (w_{ij} + X_{ij}) \quad (2)$$

$$X_{ij} = \begin{cases} M(j) & \text{if } (w_{ij} \times M(j)) > 0 \\ -M(j) & \text{otherwise} \end{cases}$$

Let  $in(j)$  denote the set of all incoming edges to node  $j$  and  $out(i)$  denote the set of all outgoing edges from node  $i$ . Normalization is done to maintain the value in the range of  $[-1, 1]$ . An auxiliary variable  $X_{ij}$  is introduced to measure the effect of the merit score of a node  $j$  on its incoming edge to node  $i$ . Since merit is about whether the node is considered good or bad, it is calculated to be the sum of all its incoming edges from other nodes. Likewise, since attitude is about the particular node's view of others, it is calculated using the outgoing edges of a node towards others and its corresponding merit score in the network. Although we use two metrics similar to BAD, the calculation of the incoming and the outgoing edges of a node differs. Since Bias in BAD is about how truly it rates other nodes, it is calculated by the difference in the edge weight and the real trust of a node (deserve). The explanation of the proposed metric follows.

From the above expression, it can be seen that if the outgoing edge weight from node  $i$  to node  $j$  has a positive

value and the merit score of node  $j$  is also positive, then the attitude of node  $i$  to  $j$  is calculated by the sum of both values. If the outgoing edge weight from node  $i$  to  $j$  is negative and the merit score of node  $j$  is positive or vice-versa, then the attitude of node  $i$  to  $j$  is calculated by subtracting the merit score from the edge weight, which means if a node has a positive edge weight towards a benign merit node then the attitude score increases. Similarly, holding a negative edge weight towards a benign merit node decreases that node's attitude score. However, if a node has a positive edge weight towards a negative merit node, the attitude of a node decreases.

Table V: Example showing the values of the graph (Figure 6c) after each iteration. A denotes attitude and M denotes merit.

No.	P1		P2		P3		P4		P5	
	M	A	M	A	M	A	M	A	M	A
0	-1	-	-1	-1	-1	-1	-1	-1	-1	-1
1	0.02	-	0.01	0.11	-0.01	-0.13	0.09	0.06	0.01	-0.13
2	0.01	-	0.02	0.11	0.1	-0.11	0.01	0.06	0.0	-0.11
3	0.01	-	0.01	0.11	0.00	-0.11	0.01	0.06	0.0	-0.11
4	0.01	-	0.01	0.11	0.00	-0.11	0.01	0.06	0.0	-0.11

From Eq. 1 and Eq. 2, the attitude of a node depends on the merit of its neighbours and vice versa. A fixed-point iteration method is used to obtain the solution. The Merit and Attitude of node  $i$  at iteration  $n$  are denoted by  $A^n(i)$  and  $M^n(i)$  respectively. The proposed algorithm 3 is designed to compute merit and attitude scores for each node in the network. Initially, we start with a Merit and Attitude score of  $-1$  (i.e., the first iteration) in step 1. In step 2a, the merit scores for each node are updated using the attitude scores from the previous iteration. In step 2b, the attitude scores are updated using the newly updated Merit scores in the same iteration. Both Merit and Attitude scores are mutually recursive and are updated until both the scores converges in step 3. The scores of Merit and Attitude from the last iteration are the final scores. In the final step 4, all the nodes whose attitude score is less than zero are added the list  $L$  along with the user's attitude score.

*Example 5:* Table V demonstrates the value of Attitude and Merit that are updated after each iteration by applying algorithm 3 to Figure 6c. The Attitude column of node  $P1$  is blank because there are no outgoing edges from  $P1$ . The last iteration shows the final attitude and merit score of the nodes. It can be seen that, node  $P3$  and  $P5$  are bullies with a confidence score of 0.11 and 0.11, respectively.

**Bully Finding Algorithm****Input:** Bullying Signed Network  $G_s = (V, E, W)$ **Output:** List of bullies and its attitude score  $L = [(u_1, s_1), (u_2, s_2), \dots, (u_{|L|}, s_{|L|})]$ 

- 1) Initialize  $M^0(v) = -1$  and  $A^0(v) = -1, \forall v \in V$ .
- 2) Set iteration index  $i = 1$ 
  - a) For each  $v \in V$  compute merit score  
 $M^i(v) = \frac{1}{2|in(v)|} \sum_{u \in in(v)} (w_{uv})(A^{i-1}(u))$  where  $|in(v)|$  is the number of incoming edges to the node  $v$
  - b) For each  $u \in V$  compute attitude score  
 $A^i(u) = \frac{1}{2|out(u)|} \sum_{v \in out(u)} (w_{uv} + X_{uv})$  where  $|out(u)|$  is the number of outgoing edges from the node  $u$
- 3) If there exist atleast one  $v \in V : M^i(v) \neq M^{i-1}(v)$  or  $A^i(v) \neq A^{i-1}(v)$ 
  - a) Increase the iteration index  $i = i + 1$
  - b) Repeat step 2a & 2b for each iteration
- 4) For each  $v \in V$  add the node and its corresponding attitude score value greater than 0 to the list  $L$
- 5) Return  $L$

Algorithm 3: Bully Finding Algorithm (BFA)

## V. ALGORITHM ANALYSIS

In this section, we show the proof of convergence of the centrality measure and perform complexity analysis of our proposed approach.

## A. Convergence of Centrality Measure

We start the convergence proof by showing the difference between the attitude of a node at any iteration and the infinite iteration is bounded which then leads to convergence by proving the error bound  $\epsilon, \ll 1$ .

After a certain iteration  $t$ , the attitude score of that iteration becomes close to  $A^\infty$ . Since merit of a node can be expressed in terms of attitude of other nodes, this implies that merit values exhibit similar properties.

**Proposition 5.1:** Attitude and Merit (A&M) of a node at any iteration  $n$  and the infinite iteration is bounded by inverse exponential function of  $n$ .

*Proof :* We prove this in Appendix A

## B. Complexity Analysis

**Proposition 5.2:** The overall complexity of our proposed approach in the average case is  $\mathcal{O}(k \times l + \log n)n$ .

*Proof :* We can determine the time complexity of the proposed approach in three phases: constructing conversation graph, constructing bullying signed network and bully finding.

**Constructing conversation graphs phase.** In the constructing conversation phase, the runtime complexity is the time taken to construct  $m$  conversations from  $n$  tweets and then to generate graphs from the constructed conversations.

Initial sorting of tweets uses merge sort which takes a computational time of  $\mathcal{O}(n \log n)$ . The conversation is constructed by doing a binary search on the  $DID$  and  $STD$  of the conversation tweet and the current tweet respectively, leading to  $m$  conversations with a computational time of  $\mathcal{O}(n \log n)$ . The cost for generating graph from the conversations is  $\mathcal{O}(m)$ . Therefore the average computational cost to construct conversation graphs is

$$\mathcal{O}(n \log n + n \log n + m) = \mathcal{O}(n \log n + m)$$

**Constructing bullying signed network phase.** In the constructing bullying signed network phase, we traverse though each conversation graph where the bullying score is calculated for each node with respect to the edges with same order. For each conversation graph  $m$ , the maximum number of nodes in the worst case is  $k$ . Therefore the total computational cost is  $\mathcal{O}(n.k + m.k)$

**Bully finding phase.** In the bully finding phase, the runtime is the time taken to detect the bullying users using attitude and merit centrality. For each  $l$  number of iteration, A&M centrality touches each edge atmost twice. Therefore the average case in detecting bullies in each iteration is  $\mathcal{O}(2n.k)$  and for the given  $l$  iteration it is  $\mathcal{O}(n.k.l)$

Therefore, the overall complexity of our proposed approach in the average case is:

$$\mathcal{O}((k.l + \log n)n + k.m) = \mathcal{O}(k.l + \log n)n \text{ since } m, k \ll n.$$

## VI. EXPERIMENTAL EVALUATION

In this section, we evaluate the performance of the proposed algorithms. First, we present the data used in our evaluation, second we discuss the implementation details, and the way we process it to build ground truth. Finally, we present the experimental results which include determining the coefficients  $\alpha, \beta$  and  $\gamma$ , utility and scalability.

## A. Dataset

In this paper, we rely on Twitter's Streaming API, which provides free access to 1% of all tweets. The API returns each tweet in a JSON format, with the content of the tweet, metadata (e.g., creation time, source ID, destination ID, reply/retweet, etc.) as well as information about the poster (e.g., username, followers, friends). To prevent our own bias, we first randomly chose 5000 interconnected users and collected all the tweets in JSON format totaling 5.6M within a six



month time-frame between May and October 2017. We then extracted features like username, text, replyname, mentions and network based features like source ID, destination ID from the Twitter JSON. There were about 2% of the tweets which were in languages other than English. When examining the users, about 90% of their geographical location were in USA, 6% of the users' location were in United Kingdom and remaining 4% were from Ecuador, Japan and China.

### B. Implementation and Setup

We implemented our algorithm in Java, and our experiments were conducted on a machine equipped with an Intel(R) Core(TM) i7-8550U CPU @ 2.00GHz processor and 16.0 GB RAM, running Windows 10 64-bit operating system.

We employed Amazon *Mechanical Turk* (MTurk) workers to respond to an online survey that we developed. We provided 2700 surveys with each survey consisting of 10 conversations. Each survey was assigned to three workers to classify the bullying behavior of the users in the conversations according to four predefined labels (strongly positive, likely positive, likely negative and strongly negative) to avoid biased interpretation of bullies. Overall, the workers rated 27000 conversations containing 1700 users, which were extracted from the set of raw Twitter data by using algorithm 1. The MTurk UI enables requesters to create and publish surveys (HITS) in a batch when processing many HITS of the same type thus saving time. For our study, we created a csv file that contained 2700 HITS. MTurk automatically created a separate HIT for each set of conversation in the csv file as shown in Figure 7. The results to rate each users involved in the set of conversations were obtained from the workers. A significant share of the participants for the survey came from USA, Canada, Europe and India. There was not a marked variation in the rating provided by the workers. There were about 7978 strongly negative, 47426 likely negative, 56704 likely positive and 23762 strongly positive user interactions. Some of these users appear in few conversations therefore we collect these ratings based on the users and number of workers and compute using a metric to identify 569 users as bullies. Finally, the results are normalized to form the ground truth. We analyze and compute the ground truth in a metric which results in bullying and non-bullying users. Having the computed results as ground truth, we evaluated the performance measure by experimenting the proposed algorithms results w.r.t. the number of users increasing linearly from 500 to 1700.

### C. Determining optimal values for coefficients $\alpha, \beta$ and $\gamma$

Recall that  $I_{uv} = \beta \times SA + \gamma \times CS$  in algorithm 1 and  $S_{uv} = I_{uv} + \alpha(I_{uv} - S_{vu})$  in algorithm 2. To determine the coefficient  $\beta$  and  $\gamma$  for bullying indicator  $I$  and  $\alpha$  for the bullying score  $S$ , we generate input tweets of varying length and performed experiment for different values of  $\alpha, \beta$  and  $\gamma$ . That is, with 5.7 million tweets dataset we did experiment for the tweets ranging from 1M, 2M, 3M, 4M and 5M for different  $\alpha, \beta$  and  $\gamma$  values. After experimenting with different values, we found that the coefficient values of  $\beta \geq 0.6$ ,  $\gamma \leq 0.4$  and  $\alpha \leq 0.6$  to provide the greatest accuracy. The accuracy was

measured with  $\beta \geq 0.6$  and  $\gamma \leq 0.4$  for every  $\alpha \leq 0.6$  with respect to the ground truth, using the F1 Measure [41].

Figure 8 depicts the optimal values for the coefficients  $\alpha, \beta$  and  $\gamma$  with respect to the  $\beta$  and  $\gamma$  values, which are set from 60 to 90 and 40 to 10 respectively. We use three different  $\alpha$  values for every bullying indicator coefficients  $\beta$  and  $\gamma$ , which varies from 0.4 to 0.6. In our approach, we observe that the F1 measure increases linearly when the coefficients  $\beta$  increases and  $\gamma$  decreases. We also observe that when we increase the  $\alpha$  value, the F1 measure increases in all the cases indicating that the sentiment analysis has more impact on the bullying indicator than the cosine similarity. This is because, sentiment analysis analyzes not only the text but also emojis, emoticons and, determining the cosine similarity alone hurts the performance. Hence, we take advantage of both Sentiment Analysis and Cosine. Similarly, the response to a tweet has a direct effect on the bullying score.

### D. Utility

We briefly introduce our evaluation metrics that will be used to determine the accuracy of our approach.

- *AccuracyCM* [42]

The accuracy measure is the ratio of the number of bully users detected to the total number of bullies. It does not perform well with imbalanced data sets.

$$AccuracyCM = \frac{\# \text{ of detected bullies}}{\text{total number of bullies}}$$

- *Precision and Recall* [43]

Precision and Recall are evaluation metrics used in binary classification tasks. Precision is the measure of exactness and recall is the measure of completeness. They are defined as follows:

$$Precision = \frac{\# \text{ of true bullies detected}}{\text{total number of detected users}}$$

$$Recall = \frac{\# \text{ of true bullies detected}}{\text{total number of true bullies}}$$

In simple terms, high precision means that an algorithm returned substantially more bully users, while high recall means that an algorithm returned most of the bullies.

- *F1 Measure* [41]

F1 Measure is the Harmonic Mean between precision and recall. The range for F1 is [0, 1]. It measures how many bullies are identified correctly and how robust it is. Mathematically, it can be expressed as :

$$F1 = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

F1 Measure attempts to find a balance between precision and recall. The greater the F1 Measure, the better is the performance of our approach.

To determine the accuracy of our proposed centrality measure, Attitude and Merit, we compare all the evaluation metrics discussed above with respect to the number of users increasing linearly from 500 to 1700 users. Figure 9 illustrates the utility values of the metrics (accuracyCM, precision, recall and F1

Conversation 1:

P1 : This might be the most pathetic thing I have ever read.  
P2 to P1 : @P1 Cry me a river woman.  
P1 to P2 : @P2 This has the makings of a country song: My dog died my man left me my car's in the shop Trump blocked me on P2185.  
P2 to P1 : @P1 Now I'm drinking whiskey and tears.  
P1 to P2 : @P2 ??????

Select the behavior (sentiment) expressed by each person

P1 : ☐ Strongly Negative ☐ Likely Negative ☐ LikelyPositive ☐ Strongly Positive

P2 : ☐ Strongly Negative ☐ Likely Negative ☐ LikelyPositive ☐ Strongly Positive

Figure 7: Sample user interface of Amazon Mechanical Turk survey (positive - appropriate behavior, negative - inappropriate behavior).

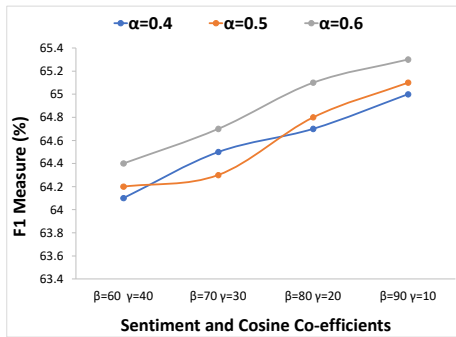


Figure 8: Determining optimal values for coefficients  $\alpha, \beta$  and  $\gamma$ .

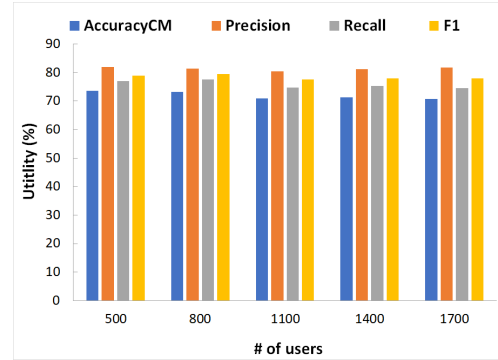


Figure 9: Utility with respect to the number of users

Measure) with respect to the number of users generated from the algorithm 2 as the input.

For the number of users ranging from 500 to 1700, we observed that the AccuracyCM metric ranged between 70.8% to 73.6% and can be biased in the case of unbalanced datasets, however it produces better results when false positives (is an error in bullies detection in which a detection result improperly indicates that a user is bully, when in reality the user is not a bully) and false negatives (is an error in which a test detection improperly indicates that a user is not bully, when in reality the user is a bully) are almost even. In this case of uneven distribution of data, we measure the accuracy with F1 Measure, which ranges from 77.5% to 79.4% while the precision and recall center around 81% and 76% respectively. Therefore from the Figure 9 it can be seen that, the precision outperform other metrics i.e., higher the precision means our algorithm identifies more bullies precisely among the total number of users. The percentages mentioned above for all the metrics remained almost consistent even with the increase in the number of users.

Next, we compare the performance of our proposed centrality measure Attitude and Merit with the research work done by Mishra and Bhattacharya [34] - Bias and Deserve which is explain in the section IV-C. We compare the F1-score in term of accuracy achieved with respect to the number of users generated from the algorithm 2 as the input. Figure 10 elucidates the comparison of the centrality measures w.r.t. the number of users increasing linearly from 500 to 1700 users.

In our approach, we observed that  $A\&M$  has an accuracy of about 80%. Also, our centrality measures outperform  $BAD$

in all the cases i.e., number of users. As the number of users increased from 500 to 1700, the accuracy of Bias and Deserve decreased from 65% to 60%, whereas the proposed centrality measures Attitude and Merit stays consistent. There can be multiple reasons behind it. First of all, the bias score of a node with highly positive bias decreases when it has an outgoing edge with positive weight whereas in  $A\&M$ , the Attitude score increases when a positive node has an outgoing edge with positive weight. Next, when calculating the deserve for a node, the bias value is taken in range of  $[0, 1]$  whereas in  $A\&M$ , merit is calculated with the attitude value in the range of  $[-1, 1]$ . Furthermore,  $BAD$  does not perform well when a node has fewer outgoing and incoming edges. Nevertheless it is still outperformed by the  $A\&M$  centrality.

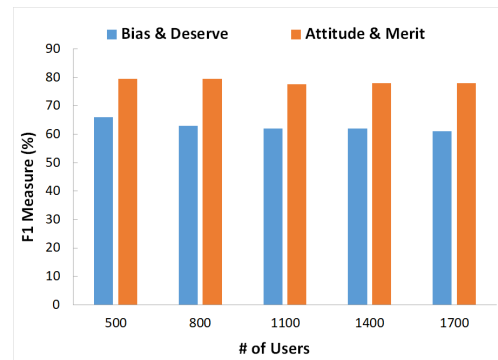


Figure 10: Comparative evaluation of the proposed centrality measure Attitude and Merit with Bias and Deserve

We also compare the accuracy of our BullyNet algorithm with Chatzakou *et al.* [24], Zhao *et al.* [18], and Singh *et*

al. [19]. As shown in Table VI, the *F1 score* of BullyNet outperforms all the other methods. However, *Precision* and *Recall* of BullyNet are outperformed by [19] and [18], respectively.

Table VI: Performance comparisons of different methods.

	Precision	Recall	F1 Score
Chatzakou <i>et al.</i> [24]	75	53	79
Zhao <i>et al.</i> [18]	76.8	79.4	78.0
Singh <i>et al.</i> [19]	82	53	64
<b>BullyNet</b>	<b>81.3</b>	<b>77.6</b>	<b>79.4</b>

### E. Scalability

We measure the scalability of BullyNet with respect to the number of tweets and observe the run-times of our three algorithms: conversation graphs generation, bullying signed network generation and bully finding with optimal values for coefficients  $\alpha, \beta$  and  $\gamma$  set at 0.6, 0.9 and 0.1 respectively.

We observed that running a dataset with 1M records takes up to 8 min for the BullyNet algorithm and the runtime increases linearly as the record size increases linearly from 1M to 5M. Figure 11 depicts the runtime for the records size from 1M to 5M for each dataset. We also observed that the most dominant algorithm of our experiment is conversation graphs generation which took the majority of run time i.e approximately 70% of total execution time of the three algorithms. This is due to the fact that the conversation graphs have to calculate sentiment analysis and cosine similarity for each tweet and then calculate the corresponding bullying indicator  $I$  as the edge weight for each conversation graph.

We observed that there is a linear increase in total run-time with an increase in number of tweets. However, we also observed that the bullying signed network generation algorithm (algorithm 2) runtime didn't grow linearly with the increase in records, rather it tends to remain constant. This is because, there are  $k$  number of nodes in  $m$  conversation graphs. Therefore, to calculate the bullying score for each graphs it takes  $\mathcal{O}(k)$  and does not affect the run-time with the growth in number of tweets. We can observe that similar to the first algorithm, the runtime of the third algorithm also increases linearly with record size. The variation is attributed to the increase in the number of users in each tweet resulting in corresponding increase in computation time for centrality measures.

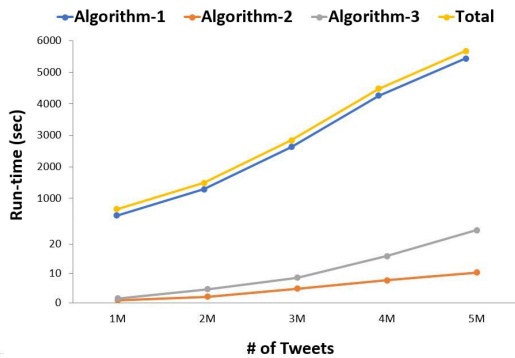


Figure 11: Scalability with respect to the number of tweets

## VII. CONCLUSION AND FUTURE WORK

Although the digital revolution and the rise of social media enabled great advances in communication platforms and social interactions, a wider proliferation of harmful behavior known as bullying has also emerged. This paper presents a novel framework of BullyNet to identify bully users from the Twitter social network. We performed extensive research on mining signed networks for better understanding of the relationships between users in social media, to build a signed network (SN) based on bullying tendencies. We observed that by constructing conversations based on the context as well as content, we could effectively identify the emotions and the behavior behind bullying. In our experimental study, the evaluation of our proposed centrality measures to detect bullies from signed network, we achieved around 80% accuracy with 81% precision in identifying bullies for various cases.

There are still several open questions deserving further investigation. First, our approach focuses on extracting emotions and behavior from texts and emojis in tweets. However, it would be interesting to investigate images and videos, given that many users use them to bully others. Second, it does not distinguish between bully and aggressive users. Devising new algorithms or techniques to distinguish bullies from aggressors would prove critical in better identification of cyberbullies. Another topic of interest would be to study the relationship between conversation graph dynamics and geographic location and how these dynamics are affected by the geographic dispersion of the users? Are the proximity increase the bullying behaviour?

## APPENDIX A

### CONVERGENCE OF CENTRALITY MEASURE

*Proposition 5.1 :* Attitude and Merit (A&M) of a node at any iteration  $n$  and the infinite iteration is bounded by inverse exponential function of  $n$ .

*Proof :* By using mathematical induction we prove the convergence of attitude. Given its definition, the attitude score  $A^\infty(i)$  and  $A^{t+1}(i)$  can be written as,

$$A^\infty(i) = \left| \frac{1}{2|out(i)|} \sum_{j \in out(i)} \{w_{ij} \pm \frac{1}{2|in(j)|} \sum_{k \in in(j)} (w_{kj} \times A^\infty(k))\} \right|$$

$$A^{n+1}(i) = \left| \frac{1}{2|out(i)|} \sum_{j \in out(i)} \{w_{ij} \pm \frac{1}{2|in(j)|} \sum_{k \in in(j)} (w_{kj} \times A^n(k))\} \right|$$

*Base case:* For  $n=1$ , we have

$$= \left| \frac{1}{2|out(i)|} \sum_{j \in out(i)} \{w_{ij} \pm \frac{1}{2|in(j)|} \sum_{k \in in(j)} w_{kj} (A^\infty(k) - A^0(k))\} \right|$$

$$\leq \frac{1}{2|out(i)|} \sum_{j \in out(i)} \{ |w_{ij}| \pm \frac{1}{2|in(j)|} \sum_{k \in in(j)} |w_{kj}| |A^\infty(k) - A^0(k)| \}$$

$$[\because |x.y| \leq |x||y| \text{ \& } |w_{ij}| \text{ and } |w_{kj}| \leq 1]$$

$$\leq \frac{1}{2|out(i)|} \sum_{j \in out(i)} \left\{ \frac{1}{2|in(j)|} \sum_{k \in in(j)} |(A^\infty(k)) - A^0(k)| \right\}$$

$$\leq \frac{1}{2|out(i)|} \sum_{j \in out(i)} \left\{ \frac{1}{2|in(j)|} \sum_{k \in in(j)} 2 \right\}$$

Since  $A(k) \in [-1, +1]$ , we have  $|A^\infty(k) - A^0(k)| \leq 2$

$$\leq \frac{1}{2|out(i)|} \sum_{j \in out(i)} \left\{ \frac{1}{2|in(j)|} 2|in(j)| \right\} = \frac{1}{2}$$

**Induction step :** We assume the bound to be true for  $A^n(i)$  so, by the hypothesis  $|A^\infty(i) - A^n(i)| \leq \frac{1}{2^{n+2}}$ . In the  $(n+1)^{th}$  iteration,

$$|A^\infty(i) - A^n(i)|$$

$$= \left| \frac{1}{2|out(i)|} \sum_{j \in out(i)} \left\{ w_{ij} \pm \frac{1}{2|in(j)|} \sum_{k \in in(j)} w_{kj} (A^\infty(k)) - A^n(k) \right\} \right|$$

$$\leq \frac{1}{2|out(i)|} \sum_{j \in out(i)} \left\{ \frac{1}{2|in(j)|} \sum_{k \in in(j)} |(A^\infty(k)) - A^n(k)| \right\}$$

$$\leq \frac{1}{2|out(i)|} \sum_{j \in out(i)} \left\{ \frac{1}{2|in(j)|} \sum_{k \in in(j)} \frac{1}{2^n} \right\} = \frac{1}{2^{n+2}}$$

Therefore the error is bounded by an inverse exponential function. Thus, we conclude that a convergence has been achieved in determining the measures 'attitude' and 'merit'.

□

## REFERENCES

- [1] J. Tang, C. Aggarwal, and H. Liu, "Recommendations in signed social networks," in *Proceedings of the International Conference on WWW*, 2016, pp. 31–40.
- [2] D. Liben-Nowell and J. Kleinberg, "The link-prediction problem for social networks," *Proceedings of the ASIS&T*, vol. 58, no. 7, pp. 1019–1031, 2007.
- [3] U. Brandes and D. Wagner, "Analysis and visualization of social networks," in *Graph drawing software*, 2004, pp. 321–340.
- [4] X. Hu, J. Tang, H. Gao, and H. Liu, "Social spammer detection with sentiment information," in *Proceedings of IEEE ICDM*, pp. 180–189, 2014.
- [5] E. E. Buckels, P. D. Trapnell, and D. L. Paulhus, *Trolls just want to have fun*, 2014, pp. 67:97–102.
- [6] S. Kumar, F. Spezzano, and V. Subrahmanian, "Accurately detecting trolls in slashdot zoo via decluttering," in *Proceedings of IEEE/ACM ASONAM*, 2014, pp. 188–195.
- [7] J. W. Patchin and S. Hinduja, "2016 cyberbullying data," 2017.
- [8] C. R. Center, "https://cyberbullying.org/bullying-laws."
- [9] D. Cartwright and F. Harary, "Structural balance: a generalization of heider's theory," *Psychological review*, vol. 63, no. 5, p. 277, 1956.
- [10] J. Leskovec, D. Huttenlocher, and J. Kleinberg, "Signed networks in social media," in *Proceedings of the SIGCHI CHI*, 2010, pp. 1361–1370.
- [11] R. Plutchik, "A general psychoevolutionary theory of emotion," in *Theories of emotion*, 1980, pp. 3–33.
- [12] W. Medhat, A. Hassan, and H. Korashy, "Sentiment analysis algorithms and applications: A survey," *Proceedings of the Ain Shams engineering journal*, vol. 5, no. 4, pp. 1093–1113, 2014.
- [13] L. Tang and H. Liu, "Community detection and mining in social media," *Synthesis lectures on data mining and knowledge discovery*, vol. 2, no. 1, pp. 1–137, 2010.
- [14] S. Bhagat, G. Cormode, and S. Muthukrishnan, "Node classification in social networks," in *Social network data analytics*, 2011, pp. 115–148.
- [15] J. Tang, Y. Chang, C. Aggarwal, and H. Liu, "A survey of signed network mining in social media," in *Proceedings of the ACM Comput. Surv.*, no. 3, pp. 42:1–42:37, 2016.
- [16] J. Kunegis, J. Preusse, and F. Schwagereit, "What is the added value of negative links in online social networks?" in *Proceedings of the International Conference on WWW*, 2013, pp. 727–736.
- [17] Z. Wu, C. C. Aggarwal, and J. Sun, "The troll-trust model for ranking in signed networks," in *Proceedings of the ACM International Conference on WSDM*, 2016, pp. 447–456.
- [18] R. Zhao, A. Zhou, and K. Mao, "Automatic detection of cyberbullying on social networks based on bullying features," in *Proceedings of the ICDN*, 2016.
- [19] V. K. Singh, Q. Huang, and P. K. Atrey, "Cyberbullying detection using probabilistic socio-textual information fusion," in *Proceedings of the IEEE/ACM ASONAM*, pp. 884–887, 2016.
- [20] H. Hosseinmardi, S. A. Mattson, R. I. Rafiq, R. Han, Q. Lv, and S. Mishra, "Detection of cyberbullying incidents on the instagram social network," in *Proceedings of the CoRR*, 2015.
- [21] J.-M. Xu, X. Zhu, and A. Bellmore, "Fast learning for sentiment analysis on bullying," in *Proceedings of the First International WISDOM*, 2012, pp. 10:1–10:6.
- [22] A. C. Squicciarini, S. M. Rajtmajer, Y. Liu, and C. H. Griffin, "Identification and characterization of cyberbullying dynamics in an online social network," in *Proceedings of the IEEE/ACM ASONAM*, 2015, pp. 280–285.
- [23] P. Galan-Garcia, J. De La Puerta, C. Gómez, I. Santos, and P. Bringas, "Supervised machine learning for the detection of troll profiles in twitter social network: Application to a real case of cyberbullying," vol. 24, pp. 42–53, 2014.
- [24] D. Chatzakou, N. Kourtellis, J. Blackburn, E. De Cristofaro, G. Stringhini, and A. Vakali, "Mean birds: Detecting aggression and bullying on twitter," in *Proceedings of the ACM on WebSci*, 2017, pp. 13–22.
- [25] L. Cheng, J. Li, Y. N. Silva, D. L. Hall, and H. Liu, "Xbully: Cyberbullying detection within a multi-modal context," in *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining*, 2019, pp. 339–347.
- [26] H.-T. Kao, S. Yan, D. Huang, N. Bartley, H. Hosseinmardi, and E. Ferrara, "Understanding cyberbullying on instagram and ask. fm via social role detection," in *Companion Proceedings of The 2019 World Wide Web Conference*, 2019, pp. 183–188.
- [27] L. Cheng, R. Guo, Y. Silva, D. Hall, and H. Liu, "Hierarchical attention networks for cyberbullying detection on the instagram social network," in *Proceedings of the 2019 SIAM International Conference on Data Mining*. SIAM, 2019, pp. 235–243.
- [28] L. Cheng, J. Li, Y. N. Silva, D. L. Hall, and H. Liu, "Pi-bully: Personalized cyberbullying detection with peer influence," in *IJCAI*, 2019, pp. 5829–5835.
- [29] J. Tang, C. Aggarwal, and H. Liu, "Node classification in signed social networks," in *Proceedings of the SIAM ICDM*, 2016, pp. 54–62.
- [30] S. Brin and L. Page, "The anatomy of a large-scale hypertextual web search engine," *Computer Networks*, vol. 30, no. 1–7, pp. 107–117, 1998.
- [31] M. Shahriari and M. Jalili, "Ranking nodes in signed social networks," *Social network analysis and mining*, vol. 4, no. 1, p. 172, 2014.
- [32] C. d. Kerchove and P. V. Dooren, "The pagetrust algorithm: How to rank web pages when negative links are allowed?" in *Proceedings of the SIAM International Conference on Data Mining*, 2008, pp. 346–352.
- [33] P. Bonacich and P. Lloyd, "Calculating status with negative relations," *Social networks*, vol. 26, no. 4, pp. 331–338, 2004.
- [34] A. Mishra and A. Bhattacharya, "Finding the bias and prestige of nodes in networks based on trust scores," in *Proceedings of the international conference on WWW*, 2011.
- [35] M. Dadvar, D. Trieschnigg, R. Ordelman, and F. de Jong, "Improving cyberbullying detection with user context," in *Proceedings of the European Conference on IR*, 2013, pp. 693–696.
- [36] K. Dinakar, R. Reichart, and H. Lieberman, "Modeling the detection of textual cyberbullying," in *Proceedings of the international AAAI WSM*, 2011.
- [37] Y. Chen, Y. Zhou, S. Zhu, and H. Xu, "Detecting offensive language in social media to protect adolescent online safety," in *Proceedings of the PASSAT and SCSM*, 2012, pp. 71–80.
- [38] L. F. Ortega, J. A. Troyano, F. L. Cruz, C. G. Vallejo, and F. Enríquez, "Propagation of trust and distrust for the detection of trolls in a social network," *Computer Networks*, vol. 56, no. 12, pp. 2884–2895, 2012.
- [39] C. J. Hutto and E. Gilbert, "Vader: A parsimonious rule-based model for sentiment analysis of social media text," in *Proceedings of the international AAAI WSM*, 2014.
- [40] J. Han, J. Pei, and M. Kamber, *Data mining: concepts and techniques*, 2011.
- [41] Y. Sasaki *et al.*, "The truth of the f-measure," *Teach Tutor mater*, pp. 1–5, 2007.
- [42] C. E. Metz, "Basic principles of roc analysis," in *Seminars in nuclear medicine*, 1978, pp. 283–298.
- [43] J. W. Perry, K. Allen, and M. M. Berry, "Machine literature searching x. machine language; factors underlying its design and development," *American Documentation (pre-1986)*, p. 242, 1955.