

Association for Information Systems AIS Electronic Library (AISeL)

PACIS 2018 Proceedings

Pacific Asia Conference on Information Systems
(PACIS)

6-26-2018

Cyberbullying Detection on Social Network Services

Pei-Ju Lee

National Chung Cheng University, pjee@mis.ccu.edu.tw

Ya-Han Hu

National Chung Cheng University, yahan.hu@mis.ccu.edu.tw

Kuanchin Chen

Western Michigan University, kc.chen@wmich.edu

J. Michael Tarn

Western Michigan University, mike.tarn@wmich.edu

Lien-En Cheng

National Chung Cheng University, lian555046@gmail.com

Follow this and additional works at: <https://aisel.aisnet.org/pacis2018>

Recommended Citation

Lee, Pei-Ju; Hu, Ya-Han; Chen, Kuanchin; Tarn, J. Michael; and Cheng, Lien-En, "Cyberbullying Detection on Social Network Services" (2018). *PACIS 2018 Proceedings*. 61.

<https://aisel.aisnet.org/pacis2018/61>

This material is brought to you by the Pacific Asia Conference on Information Systems (PACIS) at AIS Electronic Library (AISeL). It has been accepted for inclusion in PACIS 2018 Proceedings by an authorized administrator of AIS Electronic Library (AISeL). For more information, please contact elibrary@aisnet.org.

Cyberbullying Detection on Social Network Services

Research-in-Progress

Pei-Ju Lee*

National Chung Cheng University
No.168, Sec. 1, University Rd.,
Minhsiung, Chiayi, Taiwan
pjlee@mis.ccu.edu.tw

Ya-Han Hu

National Chung Cheng University
No.168, Sec. 1, University Rd.,
Minhsiung, Chiayi, Taiwan
yahan.hu@mis.ccu.edu.tw

Kuanchin Chen

Western Michigan University
1903 W Michigan Ave.,
Kalamazoo, MI, USA
kc.chen@wmich.edu

J. Michael Tarn

Western Michigan University
1903 W Michigan Ave.,
Kalamazoo, MI, USA
mike.tarn@wmich.edu

Lien-En Chen

National Chung Cheng University
No.168, Sec. 1, University Rd., Minhsiung, Chiayi, Taiwan
lian555046@gmail.com

Abstract

Social networks such as Facebook or Twitter promote the communication between people but they also lead to some excessive uses on the Internet such as cyberbullying for malicious users. In addition, the accessibility of the social network also allows cyberbullying to occur at anytime and evoke more harm from other users' dissemination. This study collects cyberbullying cases in Twitter and attempts to establish an auto-detection model of cyberbullying tweets base on the text, readability, sentiment score, and other user information to predict the tweets with harassment and ridicule cyberbullying tweets. The novelty of this study is using the readability analysis that has not been considered in past studies to reflect the author's education level, age, and social status. Three data mining techniques, k-nearest neighbors, support vector machine, and decision tree are used in this study to detect the cyberbullying tweets and select the best performance model for cyberbullying prediction.

Keywords: cyberbullying, readability, Twitter, data mining, classification

Introduction

With the rapid development of the Internet, everyone's communication is no longer limited by the need to be on site. The social network services (SNSs) allow people with the same interests make timely contact or share information on the SNS platforms. Despite these advantages of using SNSs, they also lead to excessive uses such as cyberbullying on the Internet of malicious users (Tokunaga 2010; Van Hee et al. 2015). Belsey (2009) defines cyberbullying as a use of information and communication technology to support malicious, repetitive and hostile behaviors against individuals or groups in order to harm others. Traditional bullying is a form of adolescent violence that attacks a group or an individual who can not be easily counterattacked for a prolonged period of time; often happened in schools and public spaces frequented by young people. The proliferation of cyberbullying is far greater than that of traditional bullying because of the accessibilities of the SNSs

(Grigg 2010). The well-known SNSs such as Facebook or Twitter has become the propagation channels of cyberbullying (Whittaker and Kowalski 2015). The number of active user on the SNSs increasing continually, but at present, the controls of these platforms are few and only monitor the photo posting or user privacy but not the dissemination of cyberbullying.

Previous studies of cyberbullying detection indicated that vulgarity words or word sentiment analysis are the keys to detection (Ptaszynski et al. 2010; Xu et al. 2012). Because of the malicious, repetitive and hostile behaviors of the attackers against the individuals or groups, it is also one of the obvious features of cyberbullying to observe the frequency of use of attackers and victims on the SNSs (Al-Garadi et al. 2016; Balakrishnan 2015; Chatzakou et al. 2017; Park et al. 2014). Crimes of cyberbullying also decrease as the age or educational levels of users increase (Tokunaga 2010; Al-Garadi et al. 2016). Since many SNSs do not reveal individual user's demographic information or verify most of this self-disclosed information, cyberbullying detection methods that rely only on demographic information is unlikely to be reliable. However, demographic traces (e.g. socioeconomic status, education level, age) may be derived or even inferred from the readability of what was posted (Fang et al. 2016; Due et al. 2009). Therefore, this study looks into users features and texture features such as textual readability to decipher author's demographic traces in order to establish an automatic cyberbullying posts prediction model. As of this writing, none of the existing studies explored the importance of textual readability for its relationship with cyberbullying.

Related Work

Cyberbullying may occur at any time and users of the SNSs may continue to watch it or help its dissemination by turning themselves into bystanders or perpetrators and spread it more widely (Grigg 2010). Kontostathis et al. (2013) indicated that using traditional data preprocessing methods including conversing all words into lowercase, deleting numbers or special symbols ignores the emotions of excessive online communication when people using capitalized letters or emoji (smiley faces, angry, etc.) to emphasize their feelings. Many past cyberbullying studies consistently pointed out that vulgar word and pronouns have a positive impact on cyberbullying identification. Most cyberbullying messages are constructed using aggressive words attacking other people and using pronouns such as "you" and tend to frustrate the bully-victims (Chavan and Shylaja 2015; Chen et al. 2012; Dinakar et al. 2011; Nahar et al. 2013; Nandhini and Sheeba 2015; Yin et al. 2009).

Kowalski and Limber (2013)'s study investigated 931 students in grades six through twelve and analyzed their experiences of cyberbullying and traditional bullying; the results show that the probability of being cyberbullying than traditional bullying has been gradually increasing and focusing on middle-school students since they spend more time on the Internet. Balakrishnan (2015) surveyed users of 17 to 30 years old for the experiences of cyberbullying; their study revealed that young users spend considerable time using Internet and users who spend more than 2-5 hours per day were more persecuted by cyberbullying than users spend less than one hour per day. Past related cyberbullying studies indicated that personal information of users such as age or education level is a critical feature of cyberbullying attackers but the information may not be recorded on every SNSs. Readability analysis is the use of basic statistical methods, such as using text length, word length, the number of words in a sentence, the number of syllables in a word to infer the author's level of education or knowledge. Readability analysis measures the author's education level by measuring complex or difficult of words in the article, for example, the Gunning Fog Index (Gunning, 1969) that record words of 3 or more syllables and average words per sentence. These studies infer that the author's reading skills and writing skills are related; and the author's educational level or even social status can be inferred by analyzing the author's textual readability (Fang et al. 2016; Weren et al. 2013).

Sentiment analysis is often used in natural language processing and text mining to calculate emotions levels and attitudes of words from users. VanHee et al. (2015) extracted the linguistic features of cyberbullying from Ask.fm using the Bag-of-word and sentiment polarity features (i.e. positive or negative sentiment, vulgar words) to detect cyberbullying types such as threats or insult. Many studies consider sentiment analysis as an important tool of cyberbullying detection to explore whether the

attackers insult others with obtrusive, passive, or negative emotions (Chatzakou et al. 2017; Chen et al. 2012; Dinakar et al. 2011; Nahar et al. 2013; Soundar and Ponesakki 2016; Yin et al. 2009).

In the detection of cyberbullying articles, this study uses sentiment analysis to measure the sentiment scores of text on tweets and divides them into negative, positive, or neutral sentiment since negative sentiment may involve bully activities with negative behaviors and words; uses subjectivity analysis tools to detect overly subjective and irrational tweets as well as if there are words in tweets that contain subjective as well as negative sentiment; and uses the readability scores to represent the users' demographic aspects. This study aims to establish a cyberbullying detector for Twitter, organize the important features in the past research, consider readability analysis and sentiment analysis, and adopt data mining techniques trying to increase the accuracy of the cyberbullying predict model.

Methodology

This study uses Python 3 to crawl tweets on Twitter and builds the dataset needed for the cyberbullying prediction model. The data preprocessing and feature extraction are performed for this study. The suspected cyberbullying tweets or non-harmful tweets are manually labeled by experts. The Weka 3.8 (Hall et al. 2009) was used to construct and evaluate three prediction models, namely k-Nearest Neighbors (KNN), Support Vector Machine (SVM) and C4.5 Decision Tree (DT). The performance of these models was compared and the best model was selected to report further insights. Figure 1 is the research framework.

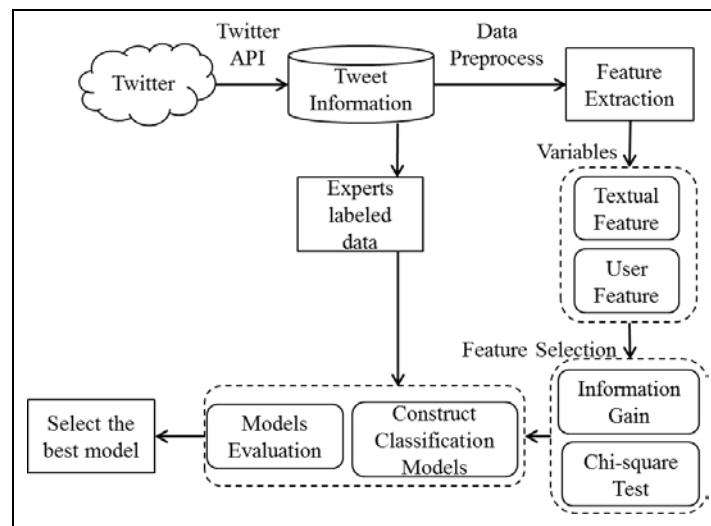


Figure 1 – Research Framework

User and tweet data collection and preprocessing

This research collects information on Twitter, which is currently one of the major SNSs in the world (Anger and Kittl 2011; Mangaonkar et al. 2015). Users can send and read 280-character tweets message and the tweets may contain words, images, emojis, tag other users (@username), or hashtag (i.e. texts to highlight words or sentences). This study attempts to use the Twitter API to collect tweets in English and the senders' user profiles and returns each data in JSON format. Tweets downloaded from twitter will be compared against hate speech from a popular online repository called Hatebase.org - an application that provides a web interface for interactive searches and open-API. This research uses highly controversial words as seeds to collect suspect cyberbullying tweets.

In the preprocessing, this study follows the steps below: (1) convert all tagged users (@USERNAME) in all tweets to USER_NAME; (2) remove URLs, punctuations, and numbers in tweet texts; (3) use Python 3 to connect Google search engine for spelling check service to correct users' spelling errors or using repeated letters such as correcting "NOOOOOOOO !!" to "NO" (4) perform word Segmentation using the Stanford CoreNLP tool that developed by the Stanford Natural Language Processing Team

based on the Recursive Neural Network (RNN) to separate each word, perform stemming and Parts of Speech Tagging (POS) (Manning et al. 2014), so as to restore the original formats of words due to different grammatical patterns in order to increase the accuracy of using text variables.

Variables

In this study, the independent variables were divided into two research areas for data processing. (1) textual features: variables generated by texts such as text readability, sentiment scores, and vulgar words, and (2) user features: the time to last tweet, the total number of tweets of the user, and other records of the user-behavior-related variables. The classification of the above two types of variables will be explored in the following sections.

Textual Features

The textual features in this study are divided into the following seven categories:

1. The basic structure: the text structure of each tweet such as the number of character, syllabus, word, and sentence; average number of syllabus per word, average number of word per sentence; and the special case of Uppercase in English which usually with emotion and emphasis tone (Chatzakou et al. 2017; Dadvar et al. 2013; Huang et al. 2014) are recorded. In addition to documenting the use of Hashtags, the number of URLs and the number of photos in each tweet are also recorded. Therefore, this study also considers words in Hashtag and calculates its ratio in all the text.
2. Vulgar words: This study uses the same vulgar-language dictionaries used in previous cyberbullying studies that include 350 curses and negative emotional words (Al-Garadi et al. 2016; Kontostathis et al. 2013; Zhao et al. 2016), and measures the number of occurrences of vulgar words in each tweet.
3. Person pronoun: In many cyberbullying articles, when offenders intentionally attack or harass others, they tend to use personal pronouns (Yin et al. 2009; Al-Garadi et al. 2016; Chatzakou et al. 2017; Chavan and Shylaja 2015; Chen et al. 2012; Dadvar et al. 2013; Dinakar et al. 2011; Nahar et al. 2013; Nandhini and Sheeba 2015; Yin et al. 2009). This study records the number of second person pronouns and third person pronouns appearing in each tweet such as "You", "yourself", or "USER_NAME".
4. Sentiment scores: The Stanford CoreNLP tool is used to conduct a sentiment analysis for each tweet. The tool uses the Stanford Sentiment Treebank semantic library (Socher et al. 2013), which manually labels sentimental polarity and contains 215,154 phrases as well as 11,855 sentences. Their study through the proposed training model and Recursive Neural Tensor Network (RNTN) to predict sentimental polarity of each word. The sentiment score of each tweet is recorded in this study.
5. Negative subjectivity: Negative subjectivity analysis uses OpinionFinder (Riloff and Wiebe 2003; Wiebe and Riloff 2005; Wilson et al. 2005) text mining tools to provide sentimental judgments and subjective detections. The OpinionFinder's dictionary, MPQA Corpus, records subjective, POS, and sentimental polarity. This study aims at the occurrence number of strong subjective word.
6. Text readability: Author's reading and writing skills are related; in addition, by analyzing the author's textual readability the author's educational level or social status may be inferred (Fang et al. 2016; Weren et al. 2013). As most SNSs do not record demographic information or even verify the validity of user self-disclosed demographics, readability of what users wrote is a surrogate measure to reflect one's demographic background. In order to avoid bias in using single readability metric, this study calculates three different readability metrics including the Flesch-Kincaid Grade Level (Kincaid et al. 1981), Gunning Fog Index (Gunning 1969), Simple Measure Of Gobbledygook (McLaughlin 1969) for each tweet.

The index score indicating the education level of the author of American grade level.

- (a) Flesch-Kincaid Grade Level (Kincaid et al. 1981)

$$FKGL = 0.39 \left(\frac{\text{total words}}{\text{total sentences}} \right) + 11.8 \left(\frac{\text{total syllables}}{\text{total words}} \right) - 15.59$$

The *total words* stand for the total number of words, *total syllables* stand for the total number of syllables, and *total sentences* stand for the total number of sentences in the tweet.

(b) Gunning Fog Index (Gunning 1969)

$$FOG = 0.4 * \left[\left(\frac{\text{Words}}{\text{Sentences}} \right) + 100 \left(\frac{\text{Complex words}}{\text{Words}} \right) \right]$$

The *words* represent the number of words, *sentences* represent the number of sentences, and *complex words* represent the number of words with three or above syllables in the tweet.

(c) Simple Measure Of Gobbledygook (McLaughlin 1969)

$$SMOG = 1.0430 \sqrt{\text{total polysyllables} \times \frac{30}{\text{total sentences}}} + 3.1291$$

The *total polysyllables* indicate the number of words with three or above syllables and the *total sentences* indicate the number of sentences in the tweet.

User Features

The frequencies of SNSs usages of cyberbullying attackers and victims are also the major characteristics of cyberbullying since the cyberbullying attackers usually conduct repetitive and hostile behaviors against individuals or groups (Belsey 2009). The SNS usage frequency can be determined from the frequency of user tweets or replies to others' tweets (Al-Garadi et al. 2016; Balakrishnan 2015; Chatzakou et al. 2017; Park et al. 2014). This study uses the Twitter API to extract user information including the total number of tweets, whether the user use default image in their profile, whether the user account has been authenticated, the number of words in self-introduction and sentiment scores in self-introduction are calculated using the Stanford CoreNLP tool (Manning et al. 2014).

Data Mining Techniques

This study uses three kinds of classifier which are k-Nearest Neighbors (KNN), Support Vector Machine (SVM), and C4.5 Decision Trees (DT) and uses the automatically adjusted default values for analyzation. The KNN classification algorithm mainly compares the objects to be predicted with their K nearest neighbors to determine their own categories. The SVM algorithm is often used in image processing, predictive classification, and supervised linear binary classification. The principle is to find the hyperplane in the vector space formed by the training data and divide the data into two groups. The DT presents the attribute relations in the form of a simple tree structure and deduces the rules between attributes and categories. Quinlan (1986) proposed the C4.5 algorithm, which use the gain ratio to calculate the gain of each attribute and select the highest attribute as the division attribute of the root node.

Experimental Design and Evaluation Methods

The data mining software Weka 3.8 is used to construct the models of KNN, SVM, and DT of cyberbullying prediction. The information gain and chi-square test are used for feature selection from textual features and user features and set up a new feature set for the models. The 10-fold cross-validation is used for evaluations of the three models.

The prediction accuracy of these generated cyberbullying models by different algorithms will be evaluated using the confusion Matrix: the *true positive (TP)* represents the actual cyberbullying tweets which also be predicted as cyberbullying; the *false positive (FP)* represents the actual non-cyberbullying tweets which be predicted as cyberbullying; the *true negative (TN)* represents the actual non-cyberbullying tweets which be predicted as non-cyberbullying; and the *false negative (FN)*

represents the actual cyberbullying tweets which be predicted as non-cyberbullying. The prediction accuracy parameters are defined as following: $precision = TP / (TP+FP)$, the rate at which the prediction model properly classifies cyberbullying tweets; $recall = TP / (TP+FN)$, the rate of the model properly catches the actual cyberbullying tweets; and $F-measure = (2*TP) / (2*TP+FP+FN)$, the rate of classification model accuracy.

Conclusion

This study is expected to collect user information (e.g. total number of tweets, self-introduction, recent 20 tweets) and text information (e.g. subjectivity, sentiment analysis, person pronouns) from Twitter to establish a cyberbullying detection model. This study adopts text readability as a novel textual feature that has not been used in past cyberbullying related literature and attempts to increase the detection accuracy of the suspected cyberbullying tweet. We use the data mining techniques including k-nearest neighbors, support vector machine, and decision tree for detections and select the best model of cyberbullying prediction. This study contributes to predicting the bully tweets through the models timely which avoid harmful posts appearing on the social network service and help the manager of the website to remove the vicious posts as soon as possible in order to prevent the proliferation of cyberbullying.

Acknowledgments

This research was supported in part by the Ministry of Science and Technology of the Republic of China (grant number MOST 106-2410-H-194-021).

References

- Al-garadi, M. A., Varathan, K. D., and Ravana, S. D. 2016. "Cybercrime Detection in Online Communications: The Experimental Case of Cyberbullying Detection in the Twitter Network," *Computers in Human Behavior* (63), pp. 433-443.
- Anger, I., and Kittl, C. 2011. "Measuring Influence on Twitter," *Proceedings of the 11th International Conference on Knowledge Management and Knowledge Technologies: ACM*, p. 31.
- Balakrishnan, V. 2015. "Cyberbullying among Young Adults in Malaysia: The Roles of Gender, Age and Internet Frequency," *Computers in Human Behavior* (46), pp. 149-157.
- Belsey, B. 2009. "Cyberbullying." from <http://www.cyberbullying.ca/>
- Chatzakou, D., Kourtellis, N., Blackburn, J., De Cristofaro, E., Stringhini, G., and Vakali, A. 2017. "Mean Birds: Detecting Aggression and Bullying on Twitter," *Proceedings of the 2017 ACM on Web Science Conference: ACM*, pp. 13-22.
- Chavan, V. S., and Shylaja, S. 2015. "Machine Learning Approach for Detection of Cyber-Aggressive Comments by Peers on Social Media Network," *Advances in computing, communications and informatics (ICACCI), 2015 International Conference on: IEEE*, pp. 2354-2358.
- Chen, Y., Zhang, L., Michelony, A., and Zhang, Y. 2012. "4is of Social Bully Filtering: Identity, Inference, Influence, and Intervention," *Proceedings of the 21st ACM international conference on Information and knowledge management: ACM*, pp. 2677-2679.
- Dadvar, M., Trieschnigg, D., Ordelman, R., and de Jong, F. 2013. "Improving Cyberbullying Detection with User Context," *European Conference on Information Retrieval: Springer*, pp. 693-696.
- Dinakar, K., Reichart, R., and Lieberman, H. 2011. "Modeling the Detection of Textual Cyberbullying," *The Social Mobile Web* (11:02), pp. 11-17.
- Due, P., Merlo, J., Harel-Fisch, Y., Damsgaard, M. T., Holstein, B. E., Hetland, J., Currie, C., Gabhainn, S. N., de Matos, M. G., and Lynch, J. 2009. "Socioeconomic Inequality in Exposure to Bullying During Adolescence: A Comparative, Cross-Sectional, Multilevel Study in 35 Countries," *Am J Public Health* (99:5), pp. 907-914.
- Fang, B., Ye, Q., Kucukusta, D., and Law, R. 2016. "Analysis of the Perceived Value of Online Tourism Reviews: Influence of Readability and Reviewer Characteristics," *Tourism Management* (52), pp. 498-506.

- Grigg, D. W. 2010. "Cyber-Aggression: Definition and Concept of Cyberbullying," *Australian Journal of Guidance and Counselling* (20:02), pp. 143-156.
- Gunning, R. 1969. "The Fog Index after Twenty Years," *Journal of Business Communication* (6:2), pp. 3-13.
- Hall, M., Frank, E., Holmes, G., Pfahringer, B., Reutemann, P., and Witten, I. H. 2009. "The Weka Data Mining Software: An Update," *ACM SIGKDD explorations newsletter* (11:1), pp. 10-18.
- Huang, Q., Singh, V. K., and Atrey, P. K. 2014. "Cyber Bullying Detection Using Social and Textual Analysis," in: *Proceedings of the 3rd International Workshop on Socially-Aware Multimedia - SAM '14*, pp. 3-6.
- Kincaid, J. P., Aagard, J. A., Hara, J. W. O., and Cottrell, L. K. 1981. "Computer Readability Editing System," *IEEE Transactions on Professional Communication* (PC-24:1), pp. 38-42.
- Kontostathis, A., Reynolds, K., Garron, A., and Edwards, L. 2013. "Detecting Cyberbullying: Query Terms and Techniques," in: *Proceedings of the 5th Annual ACM Web Science Conference*. Paris, France: ACM, pp. 195-204.
- Kowalski, R. M., and Limber, S. P. 2013. "Psychological, Physical, and Academic Correlates of Cyberbullying and Traditional Bullying," *J Adolesc Health* (53:1 Suppl), pp. S13-20.
- Mangaonkar, A., Hayrapetian, A., and Raje, R. 2015. "Collaborative Detection of Cyberbullying Behavior in Twitter Data," *2015 IEEE International Conference on Electro/Information Technology (EIT)*, pp. 611-616.
- Manning, C., Surdeanu, M., Bauer, J., Finkel, J., Bethard, S., and McClosky, D. 2014. "The Stanford CoreNlp Natural Language Processing Toolkit," *Proceedings of 52nd annual meeting of the association for computational linguistics: system demonstrations*, pp. 55-60.
- Mc Laughlin, G. H. 1969. "Smog Grading-a New Readability Formula," *Journal of reading* (12:8), pp. 639-646.
- Nahar, V., Li, X., and Pang, C. 2013. "An Effective Approach for Cyberbullying Detection," *Communications in Information Science and Management Engineering* (3:5), p. 238.
- Nandhini, B. S., and Sheeba, J. I. 2015. "Online Social Network Bullying Detection Using Intelligence Techniques," *Procedia Computer Science* (45), pp. 485-492.
- Park, S., Na, E. Y., and Kim, E. M. 2014. "The Relationship between Online Activities, Netiquette and Cyberbullying," *Children and Youth Services Review* (42), pp. 74-81.
- Ptaszynski, M., Masui, F., Nitta, T., Hatakeyama, S., Kimura, Y., Rzepka, R., and Araki, K. 2016. "Sustainable Cyberbullying Detection with Category-Maximized Relevance of Harmful Phrases and Double-Filtered Automatic Optimization," *International Journal of Child-Computer Interaction* (8), pp. 15-30.
- Quinlan, J. R. 1986. "Induction of Decision Trees," *Machine learning* (1:1), pp. 81-106.
- Riloff, E., and Wiebe, J. 2003. "Learning Extraction Patterns for Subjective Expressions," *Proceedings of the 2003 conference on Empirical methods in natural language processing: Association for Computational Linguistics*, pp. 105-112.
- Socher, R., Perelygin, A., Wu, J., Chuang, J., Manning, C. D., Ng, A., and Potts, C. 2013. "Recursive Deep Models for Semantic Compositionality over a Sentiment Treebank," *Proceedings of the 2013 conference on empirical methods in natural language processing*, pp. 1631-1642.
- Soundar, K. R., and Ponesakki, P. 2016. "Cyberbullying Detection Based on Text Representation," *International Journal of Engineering Science* (6:10), pp. 2776-2785.
- Tokunaga, R. S. 2010. "Following You Home from School: A Critical Review and Synthesis of Research on Cyberbullying Victimization," *Computers in Human Behavior* (26:3), pp. 277-287.
- Van Hee, C., Lefever, E., Verhoeven, B., Mennes, J., Desmet, B., De Pauw, G., Daelemans, W., and Hoste, V. 2015. "Automatic Detection and Prevention of Cyberbullying," *International Conference on Human and Social Analytics (HUSO 2015): IARIA*, pp. 13-18.
- Weren, E. R., Moreira, V. P., and Oliveira, J. 2013. "Using Simple Content Features for the Author Profiling Task," *Notebook for PAN at Cross-Language Evaluation Forum. Valencia, Spain*.
- Whittaker, E., and Kowalski, R. M. 2015. "Cyberbullying Via Social Media," *Journal of School Violence* (14:1), pp. 11-29.
- Xu, J.-M., Zhu, X., and Bellmore, A. 2012. "Fast Learning for Sentiment Analysis on Bullying," *Proceedings of the First International Workshop on Issues of Sentiment Discovery and Opinion Mining: ACM*, p. 10.

- Yin, D., Xue, Z., Hong, L., Davison, B. D., Kontostathis, A., and Edwards, L. 2009. "Detection of Harassment on Web 2.0," *Proceedings of the Content Analysis in the WEB 2.0 Workshop*, pp. 1-7.
- Zhao, R., Zhou, A., and Mao, K. 2016. "Automatic Detection of Cyberbullying on Social Networks Based on Bullying Features," in: *Proceedings of the 17th International Conference on Distributed Computing and Networking - ICDCN '16*. pp. 1-6.