

# Detecting Cyberbullying activities Over Social Media

John Hani, Mohamed Nashaat, Mostafa Ahmed, Zeyad Emad

September 18, 2018

## Abstract

Cyberbullying, defined by Cambridge dictionary, as the use of technology to harm or offend other people by sending them obnoxious messages. Without doubt, cyberbullying is one of the most critical topics nowadays, because it has a huge impact on our society. Cyberbullying causes big problems to the society and must be prevented as fast as possible. Lately, there have been some serious work done to help in the detection and prevention of the cyberbullying. Although, its not easy to detect it on social media. There are many challenges to detect cyberbullying automatically like: how to make the AI understand the meaning of a post or How to deal with anonymous users . Our aim in this project is to build an application that Detects cyberbullying using sentiment analysis and some deep learning techniques.

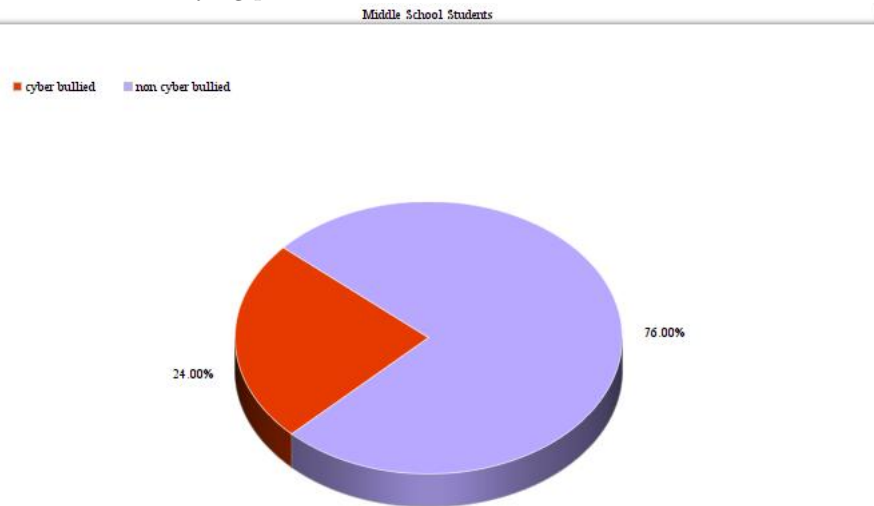
## 1 Introduction

### 1.1 Background

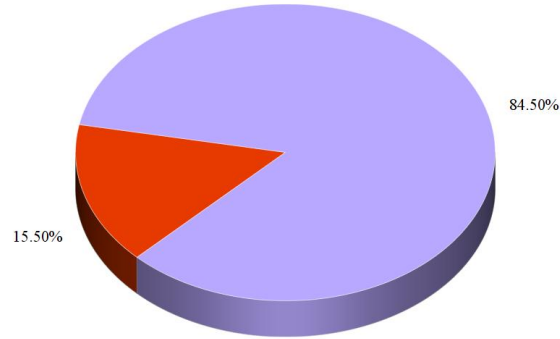
Social media has been increased strongly in the Middle East in the last decade. According to social media in the Middle East 2017 the number of the active social media users is 93 million people a day. As we know the social media is good place for communication, sharing information and maintain the old relationships. In the other hand it have many bad impacts on the society specially the teenagers. One of the biggest bad impacts is cyberbullying. According to the Center for Disease Control, students who are bullied are more likely to experience low self-esteem and isolation, perform poorly in school, have few friends in school, have a negative view of school, experience physical symptoms (such as headaches, stomachaches, or problems sleeping), and to experience mental health issues (such as depression, suicidal thoughts, and anxiety).so Our aim is to detect the cyberbullying in the messages , report this users and to inform parents about the this case using natural language processing and neural networks , classifiers and sentiment analysis.

## 1.2 Market Motivation

The market is really in need for application to detect cyberbullying According to pacers national bullying prevention center



■ cyber bullied    ■ non cyber bullied



So the community needs this application and the market should help the community to have solutions to their problems. Also one of our market motivation is that the UNICEF is now making a great campaign to eliminate bullying in Egypt.

### 1.3 Academic motivation

Our work is motivated by the previous work by the field. Rui Zheo, ET AL developed enhanced bag of words[7].

### 1.4 Problem Definitions

Increase the accuracy of the detection of cyberbullying in messages and reporting this users.

## 2 Project Description

Our goal is to build an application that's able to detect and prevent cyberbullying using sentiment and contextual analysis and deep learning techniques. First, we want to make the system able to understand the posts in social media to detect emotions in the posts and decide whether there is cyberbullying or not. Moreover, the system should be able to differentiate between sarcastic posts and

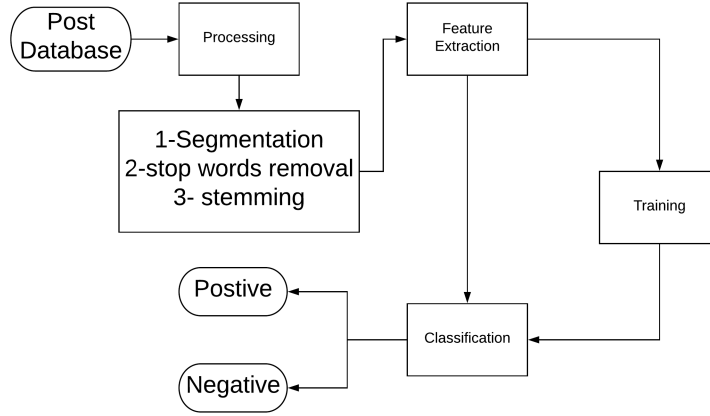
cyberbully posts by understanding the context. Finally, when a cyberbullying is detected a warning appears to the parents of bullied child or to administrator of the social media platforms to take the appropriate action. The system doesn't have any extra costs.

## 2.1 Scope

The system will cover in its scope:

1. Sentiment analysis of text combined with contextual analysis.
2. The system will work on word level analysis and also phrase level analysis.
3. The system will use deep learning classifier.

## 2.2 Project Overview



## 3 Similar System Information

1. Sentiment Informed Cyberbullying Detection in Social Media [3]:
  - (a) In this paper the researchers were motivated by psychological and sociological findings, wanted to investigate the relationship between sentiment information and Cyberbullying behaviors.
  - (b) The main problem is to use sentiment analysis to detect cyberbullying and deal with 2 problems: short, noisy and unstructured content information and the obfuscation of the obnoxious words by the users.
  - (c) Researchers proposed a principle learning framework called (SICD) and they study whether sentiment information is particularly correlated with cyberbullying behaviors and how to deal with short and unstructured content.

- (d) Researchers conducted extensive experiments on two real-world datasets. The experimental results show the effectiveness of the proposed model as well as the impact of sentiment information.
- (e) This Paper is going to help us in the sentiment analysis section as they done many experiments that investigate the effectiveness of the sentiment analysis on cyberbullying.

## 2. Automatic Detection of Cyberbullying on Social Networks based on Bullying Features[7]:

- (a) They made this program because the increasing of social media which increase the cyber bullying that give bad impacts on children and teenagers such as depression and suicidal thoughts.
- (b) The main problem that BOW is every word is independent from the other that fail to see the sentence as all.
- (c) They made a framework that detect the cyber bullying, based on word embeddings they made a list of insulting words then they assign to them weights. After this they concatenate latent semantic feature with bag of words then they classified them with SVM.

**Table 8: Precision (%), Recall (%), and F1 Scores (%) for compared methods. Bold face indicates best performance.**

Measures	BoW	sBoW	LSA	LDA	EBow
Precision	75.6	75.7	75.9	74.0	<b>76.8</b>
Recall	77.8	78.3	78.2	76.5	<b>79.4</b>
F1 Score	76.6	76.9	77.0	74.9	<b>78.0</b>



- (d)
  - (e) It is important to us because they concatenate bag of words with latent semantic feature.
- ## 3. Cybercrime detection in online communications: The experimental Case of cyberbullying detection in the Twitter network [1]:
- (a) The bad effects of social media like cyberbullying that make the cyberbullied person suffering from many things such as suicidal thoughts and depression.
  - (b) They dont have word embeddings or sentiment analysis they rely their work on classification.
  - (c) Their model takes network, tweet content, activity and user features from tweets then they train random forest with SMOTE classifiers to classify cyberbullying and non-cyberbullying.
  - (d) Results: under the receiver operating characteristic (ROC) curve (AUC) of 0.943 fmeasure of 0.936 using random forest with SMOTE.

- (e) It is important to us because they use hybrid classifiers which one of them is random forest and we plan to use these methods.
4. Unsupervised Cyber Bullying Detection in Social Networks [4]:
- (a) While cyber bullying is a well-studied problem from a social point of view, only recently it has attracted the attention of computer scientists, especially towards automatic detection tasks. For this reason, only relatively few articles on the subject and very few datasets are available.
- (b) We proposed to adopt an unsupervised approach to detect cyber bully traces over social networks.

**TABLE I**  
**RESULTS OBTAINED ON FORMSPRING.ME DATASET**

<b>Precision</b>	<b>Accuracy</b>	<b>Recall</b>	<b>F1</b>	<b>Method</b>
0.72	0.73	0.69	0.71	GHSOM
0.60	-	0.40	-	C4.5
-	-	0.67	-	SVM

(c)

**TABLE II**  
**AVERAGE RESULTS OBTAINED ON YOUTUBE DATASET.**

<b>Precision</b>	<b>Accuracy</b>	<b>Recall</b>	<b>F1</b>	<b>Method</b>
0.60	0.69	0.94	0.74	GHSOM

**TABLE III**  
**AVERAGE RESULTS OBTAINED ON TWITTER DATASET.**

<b>Precision</b>	<b>Accuracy</b>	<b>Recall</b>	<b>F1</b>	<b>Method</b>
0.81	0.72	0.26	0.4	GHSOM
-	0.67	-	-	Naive Bayes

- (d) We now know multiple sources that we can setup as our data sets ( YouTube, twitter, FormSpring)
5. Experts and Machines against Bullies: A Hybrid Approach to Detect Cyberbullies [2]:

- (a) Most of the technical studies have focused on the detection of cyber-bullying through identifying harassing comments rather than preventing the incidents by detecting the bullies.
  - (b) Proposed methods: we introduce the three types of models used for calculating and assigning the bulliness score to the social network users: a multi-criteria evaluation system, a set of machine learning models and two hybrid models that combine the two.
  - (c) Machine Learning Approaches: We used three well-known machine learning methods, which use pre-labelled training data for automatic learning: a Naive Bayes classifier, a classifier based on decision trees and Support Vector Machines (SVM) with a linear kernel
  - (d) Results: The discrimination capacity of the MCES was 0.72.
6. Cyberbullying System Detection and Analysis [6]:
- (a) Cyber-bullying has recently been reported as one that causes tremendous damage to society and economy.
  - (b) The system relies on the detection of three basic natural language components corresponding to Insults, Swears and Second Person.
  - (c) Proposed Methods: the whole is greater than the sum of its parts. A combination of modestly accurate features coming from heterogeneous data modalities can outperform methods that employ a single modality.

Feature	Acc.	Prec.	Reca.	F1-me	F2-me
<i>tf-Idf</i>	97,3%	31,2%	68,4%	42,85%	55,23%
<i>LIWC</i>	76,4%	28,4%	57,1%	32,56%	41,97%
<i>Depen</i>	67,5%	27,3%	60,6%	37,64%	48,72%
<i>tf-Idf+LIWC</i>	97,8%	42,4%	75,1%	54,20%	65,01%
<i>LIWC + Depen</i>	82,1%	38,4%	69,5%	49,47%	59,81%
<i>tf-Idf+Depen</i>	97,9%	58,9%	78,4%	67,26%	73,53%
<i>All features</i>	<b>99,4%</b>	<b>69,0%</b>	<b>84,9%</b>	<b>76,13%</b>	<b>81,15%</b>

- (d)
  - (e) This work opens up new direction for future research through using advanced parser, dimension reduction and taking into account users profile in order to strengthen the detection capabilities.
7. Common Sense Reasoning for Detection, Prevention, and Mitigation of Cyberbullying [5]:
- (a) Cyberbullying or harassment on social networks is as much a threat to the viability of online social networks for youth today as spam once was to email in the early days of the internet.

- (b) Proposed models: To detect explicit bullying language pertaining to (1) sexuality, (2) race and culture and (3) intelligence. Binary classifiers outperform their multiclass counterparts: JRip and Support Vector Machines were the best performing in terms of accuracy and kappa values.

	Naive Bayes			Rule-based JRip			Tree-based J48			SVM (poly-2 kernel)		
	Acc.	F1	kappa	Acc.	F1	kappa	Acc.	F1	Kappa	Acc.	F1	kappa
Sexuality	66%	0.67	0.657	<b>80%</b>	0.76	0.598	63%	0.57	0.573	66%	0.77	<b>0.79</b>
Race and Culture	66%	0.52	0.789	<b>68%</b>	0.55	0.789	63%	0.48	0.657	66%	0.63	<b>0.71</b>
Intelligence	72%	0.46	0.467	<b>70%</b>	0.51	0.512	70%	0.51	0.568	72%	0.58	<b>0.72</b>
Mixture	63%	0.57	0.445	63%	0.60	0.507	61%	0.58	0.456	66%	0.63	0.653

- (c)
- (d) Future work: They are currently embarking on the use of a family of latent variable models to model, understand and predict self-harm in adolescents, a phenomenon that is not very well understood in the field of abnormal psychology.

### 3.1 Similar System Description

Cyberbullying is gender specific which means that each gender has their own set of preferences when it comes to terminologies, and each gender specific language features into account improves the discrimination capacity of classifiers. Females tend to use relational style of aggression such cutting a person out of a group, while males use more threatening expressions and profane words. These words, including their abbreviations and acronyms, are treated as a single text and compared in relation to the whole post. Their methodology however also includes the usage of personal pronouns and second person pronoun and last but not least the IFIDF.



### 3.2 Comparison with Proposed Project

	Previous System	Our System
Accuracy	ranging between 34% and 66%	At least as low as the previous system
Methodology	<ol style="list-style-type: none"> <li>1. SVM</li> <li>2. Bag of words</li> <li>3. TD-IDF</li> <li>4. Profane words</li> <li>5. Second person pronouns</li> <li>6. Other personal pronouns</li> <li>7. The weight of the words in each sentence.</li> </ol>	<ol style="list-style-type: none"> <li>1. Sentiment and contextual features analysis</li> <li>2. Bag of words</li> <li>3. Syntactic features</li> <li>4. Semantic features</li> <li>5. Sentiment features</li> <li>6. Social features</li> <li>7. Linguistic Inquiry and Word Count</li> <li>8. TF-IDF</li> <li>9. Unusual capitalization</li> <li>10. Dependency features</li> <li>11. Lexicons and stemming</li> <li>12. Machine learning <ol style="list-style-type: none"> <li>a. SVM</li> <li>b. Naive Bayes</li> <li>c. Decision Tree</li> </ol> </li> <li>13. Hybrid classifiers</li> <li>14. Deep learning</li> </ol>
Application	No Application	A graphical user interface will be used for furthermore illustration
Dataset	Small scale of dataset	Large scale of dataset

## 4 Project Management and Deliverable

### 4.1 Tasks and time plan

Task Name	Start Time	Finish
Idea Discussion	1/8/2018	1/8/2018
Idea Research	1/8/2018	13/9/2018
Proposal Writing	13/9/2018	16/9/2018
Implementing Prototype	16/9/2018	17/9/2018
Delivering Rehearsal	18/9/2018	18/9/2018
Delivering Proposal	18/9/2018	26/9/2018
Doing Survey	10/10/2018	20/10/2018
Implementing Second Prototype	20/10/2018	25/10/2018
Writing SRS	25/10/2018	30/10/2018
Implementing	30/10/2018	25/11/2018
Preparing For External Examiner	25/11/2018	3/12/2018
Implementing	3/12/2018	18/1/2019
Writing SDD	18/1/2019	1/2/2019
Implementing	1/2/2019	1/4/2019
Preparing For Implementation Evaluation	1/4/2019	25/4/2019
Writing 8 Pages Paper	25/4/2019	28/4/2019
Finalizing Implementation	28/4/2019	7/5/2019
Writing Final Thesis	10/5/2019	25/5/2019
Presenting Final Thesis	25/6/2019	25/6/2019

## 4.2 Budget and Resources Costs

No Costs is Needed

## 4.3 Supportive Documents

- 4.3.1 "A Multilingual System for Cyberbullying Detection: Arabic Content Detection using Machine Learning" By Batoul Haidar, Maroun Chamoun, Ahmed Serhrouchni
- 4.3.2 "Cyberbullying Detection using Time Series Modeling" By Nektaria Potha, Manolis Maragoudakis
- 4.3.3 "Cybercrime detection in online communications: The experimental case of cyberbullying detection in the Twitter network" By Mohammed Ali Al-garadi , Kasturi Dewi Varathan , Sri Devi Ravana
- 4.3.4 "Machine Learning Approach for Detection of Cyber-Aggressive Comments by Peers on Social Media Network" By Vikas S Chavan , Shylaja S S
- 4.3.5 "Machine learning and semantic analysis of in-game chat for cyber bullying" By Shane Murnion , William J Buchanan , Adrian Smales , Gordon Russell
- 4.3.6 "Unsupervised Cyber Bullying Detection in Social Networks" By Michele Di Capua , Emanuel Di Nardo, Alfredo Petrosino
- 4.3.7 "Cyberbullying System Detection and Analysis" By Yee Jang Foong , Mourad Oussalah

## 5 References

### References

- [1] M. A. Al-garadi, K. D. Varathan, and S. D. Ravana, "Cybercrime detection in online communications: The experimental case of cyberbullying detection in the twitter network," *Computers in Human Behavior*, vol. 63, pp. 433–443, 2016.
- [2] M. Dadvar, D. Trieschnigg, and F. de Jong, "Experts and machines against bullies: A hybrid approach to detect cyberbullies," in *Canadian Conference on Artificial Intelligence*. Springer, 2014, pp. 275–281.
- [3] H. Dani, J. Li, and H. Liu, "Sentiment informed cyberbullying detection in social media," in *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer, 2017, pp. 52–67.

- [4] M. Di Capua, E. Di Nardo, and A. Petrosino, “Unsupervised cyber bullying detection in social networks,” in *Pattern Recognition (ICPR), 2016 23rd International Conference on*. IEEE, 2016, pp. 432–437.
- [5] K. Dinakar, B. Jones, C. Havasi, H. Lieberman, and R. Picard, “Common sense reasoning for detection, prevention, and mitigation of cyberbullying,” *ACM Transactions on Interactive Intelligent Systems (TiiS)*, vol. 2, no. 3, p. 18, 2012.
- [6] Y. J. Foong and M. Oussalah, “Cyberbullying system detection and analysis,” in *2017 European Intelligence and Security Informatics Conference (EISIC)*. IEEE, 2017, pp. 40–46.
- [7] R. Zhao, A. Zhou, and K. Mao, “Automatic detection of cyberbullying on social networks based on bullying features,” in *Proceedings of the 17th international conference on distributed computing and networking*. ACM, 2016, p. 43.