



# **iJRASET**

International Journal For Research in  
Applied Science and Engineering Technology



---

# **INTERNATIONAL JOURNAL FOR RESEARCH**

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

---

**Volume: 9      Issue: 1      Month of publication: January 2021**

**DOI: <https://doi.org/10.22214/ijraset.2021.32879>**

**[www.ijraset.com](http://www.ijraset.com)**

**Call:  08813907089**

**E-mail ID: [ijraset@gmail.com](mailto:ijraset@gmail.com)**

# A Framework to Predict Social Crimes using Twitter Tweets

Sreya C M

Department of Computer Science and Engineering, APJ Abdul Kalam Technological University

**Abstract:** Now a day's new wave of social media technologies like facebook, blogs, twitter plays an important role in formal and informal communications. In social media like twitter, users share their ideas, thoughts, and news in under 280 characters of text. An increasing amount of data coming from social networks can be used to generate a variety of data patterns for various sorts of investigation like human social behavior, system security, criminology etc. A framework is developed to predict major sorts of social media crimes (Cyber stalking, Cyber bullying, Cyber Hacking, Cyber Harassment, and Cyber Scam) using the data obtained from social media website.

The proposed system contains three modules; data (tweet) pre-processing, classifying model builder and prediction. Data is interpreted using statistical models with Python. To create the prediction model, Multinomial Naïve Bayes (MNB), K-Nearest Neighbors (KNN) or Support Vector Machine model (SVM) classifiers can be employed that classify given data into different classes of crime. Further the accuracy of the system at different levels is measured. Results shows that each of the three algorithm attain the precision, Recall and F-measure above than 0.9.

**Keywords:** Machine learning, Social media crimes, Twitter, Python, Natural Language Processing (NLP)

## I. INTRODUCTION

Twitter may be a unique way of following friends and sending tweets (Twitter messages) unlike the other social media networks because twitter friendship isn't mutual. for instance , you'll follow the celebrities without requiring them to follow you back. Twitter plays a virtual online world for its users. Virtual world interacts sort of a world where location act as an intermediate connection. Commonly used GPS feature in Smartphone and tablet-enabled social media users to connect real-time locations when sending out tweets. for instance , a tweet associated with an event sort of a tornado could be written during a very short time after a user witnessed a tornado was formed. the information might be spread faster than the other electronic media (TV, news or website). Secondly, tweets contain information that would help to judge the particular situation of the event.

In 2018 monthly active Twitter user was 330million that were posting 500 million tweets per day. These massive tweets having diverse dimensions of data, employed by researchers for various sorts of inquiries to predict future trends like future marketing outcomes, forecasting box-office movies revenues, flu spreading diseases, disaster response, crime prediction, forecasting election result etc.

Crimes occur everywhere within the world, for increase rate of crimes enforcement agencies are demanding advanced information systems which will help to scale back the crimes and protect the society. Criminology is that the scientific study of crime to seek out out the causes of crimes by collecting and investigating data. That way natural language Processing may be a good approach for text analysis. Online social networks are widely used by people to share their feelings and opinions leading to bullying and threats. So here trying to predict social media crimes by using twitter data analysis. Multinomial Naïve Bias classifier (MNB), K Nearest Neighbour classifier (KNN) and Support Vector Machine (SVM) models are used for implementing the project. KNN is employed for tweet classification.

## II. LITERATURE SURVEY

In [1] author proposes an event detection approach that utilizes hashtags in tweets, which adopted the feature extraction used in STREAMCUBE and applied a K-means clustering approach to it. The experiments revealed that the K-means approach gives better result than STREAMCUBE in the clustering . A discussion on optimal K values for the K-means approach is also given in this study. In [2] author says that Hackers make extensive utilization of online communities, sharing knowledge, tools, also performing coordination and recruitment activities. This paper presents a set of activities which analyze online communication patterns, including technical discussions and user profiling in order to identify potential facts. [3] Retrieving information from social networks is the earliest step in many data analysis fields such as Natural Language Processing, Sentiment Analysis and Machine Learning.

In [3] author proposes a brand new methodology for collecting historical tweets inside any date range using web scraping techniques bypassing for Twitter API restrictions. [4] Internet age has brought with it a slew of tools and research which enable stalkers, from ex-lovers to end strangers, to follow a person’s life in great detail without their consent. In this article, the author reviews the present literature on the topic and explores the discrepancy between technologies used by stalkers and technologies used against stalkers, then suggests some research avenues which may help correct this imbalance. In [5], a detailed survey on cyberbullying is done. The purpose of this survey is to explore the varied research works performed for detection and prevention on Cyberbullying

### III. PROPOSED SYSTEM AND METHODOLOGY

The proposed system is a framework to predict social media crimes including Cyber stalking, Cyber bullying, Cyber Hacking, Cyber Harassment, and Cyber Scam. The proposed system is composed of three modules. First module is tweet pre-processing that is used as an input to the second module for generating a trained model and final is Prediction. Block diagram of proposed framework is shown in the figure.

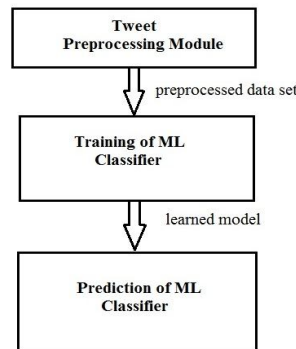


Fig. 1: Block diagram of Proposed Framework

#### A. Tweet Preprocessing

The data gathered from Twitter is unstructured and noisy Following preprocessing steps are performed to remove noise form data.

- 1) *Tweets Tokenization*: The first step in data preprocessing is tweets segmentation or tweets tokenization. Word tokenizing is a basic unit for text analysis
- 2) *Removal of Punctual Marks*: The second phase in tweets processing is the removal of punctual marks such as period, semicolon, comma, question mark, ellipsis, exclamation point, quotation marks, parentheses and apostrophe from the dataset.
- 3) *Stop Word Elimination*: The most frequently used words in the text are called stop words. The words which come frequently in the text has very low worth.
- 4) *Lower-case conversion*: Text analysis treats both upper and lower case letters equally. The size of feature words is increased if we deal with both upper case and lowercase letter in our training corpus.
- 5) *Stemming*: Another important step in data preprocessing is stemming. Stemming reduces the word to its base form by removing affixes that save time and space.

#### B. Feature Selection

The process of selecting an appropriate feature form a huge collection of processed data. There are several objectives of feature selection for instance, it reduces the training time, easier to interpret, improve the accuracy and enhance the performance of model if effective attributes are selected from the corpus. The TF-IDF approach is used to select feature words from Social crime dataset.

#### C. Algorithms

Three supervised machine learning classifiers used in this project are Multinomial Naïve Bayes, Support Vector machine and K Nearest Neighbors:

- 1) *Multinomial Naïve Bayes*: Naïve Bayes classifier is a simple, faster, efficient and easy to implement classifier.
- 2) *K- Nearest Neighbour*: KNN stand for K-Nearest Neighbors and is a well known supervised machine learning classifier used for tweets classification.
- 3) *Support Vector Machine*: Support Vector Machine (SVM) is widely used for short text classification based on the principle of structured data.



Pre-processed dataset is divided into two sets as 70% training and 30% testing. Effectiveness and accuracy of each machine learning algorithm are measured by using standard evaluation criteria such as precision, recall, and F-measure. TF-IDF approach is used to select feature words from Social crime dataset. TF-IDF stand for Term frequency- The inverse Document Frequency which means that if a term comes frequently in a particular document it has less importance and less weighting in that document. Scikit is a python library used to extract features from the text in the numerical form. Vectorization is that the process which transforms a set of text documents into numerical feature vectors. The proposed system gives a better accuracy with the help of supervised machine learning classifiers such as K Nearest Neighbor (KNN) , Multinomial Naïve Bayes (MNB) classifier and Support Vector Machine model. The interpretation of data can be done with python libraries like pandas, numpy and scikit-learn.

#### IV. CONCLUSIONS

The aim of this project is to predict major social media crimes through twitter data analysis. The proposed system gives a better accuracy with the help of supervised machine learning classifiers such as KNN ,MNB and SVM. Comparative analysis has been performed between these three supervised machine learning classifiers. This system can identify individuals that potentially will be involved in an act of crime and also can predict someone's sentiments through twitter tweets. More crime classes can be added to make the system efficient and robust.

#### REFERENCES

- [1] Yang, Shih-Feng, and Julia Taylor Rayz. "An event detection approach based on Twitter hashtags." arXiv preprint arXiv:1804.11243 (2018).
- [2] Babko-Malaya, Olga, et al. "Detection of hacking behaviors and communication patterns on social media." 2017 IEEE international conference on big data (Big Data). IEEE, 2017.
- [3] Hernandez-Suarez, Aldo, Gabriel Sanchez-Perez, Karina Toscano-Medina, Victor Martinez-Hernandez, Victor Sanchez, and Héctor Perez-Meana. "A web scraping methodology for bypassing twitter API restrictions." arXiv preprint arXiv:1803.09875 (2018).
- [4] Eterovic-Soric, Brett, et al. "Stalking the stalkers—detecting and deterring stalking behaviours using technology: A review." Computers & security 70 (2017): 278-289.
- [5] Krithika, V., and V. Priya. "A Detailed Survey On Cyberbullying in Social Networks." 2020 International Conference on Emerging Trends in Information Technology and Engineering (ic-ETITE). IEEE, 2020.



10.22214/IJRASET



45.98



IMPACT FACTOR:  
7.129



IMPACT FACTOR:  
7.429



# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24\*7 Support on Whatsapp)