# Spammer Detection and Fake User Identification Using Machine Learning

Dr.M.ChinnaRao[1], Mathi EswaraSaiKrishna[2], Meesala SuvarnaBabu[3], Athmuri DheerajBabu[4], Lankapalli Laban[5]

[1]*Professor, Lingaya's Institute of Management and Technology, Madalavarigudem, Vijayawada, Andhra Pradesh*

[2]*UG Scholar, Lingaya's Institute of Management and Technology, Madalavarigudem, Vijayawada, Andhra Pradesh*

[2,3,4,5] *Department of CSE, Lingaya's Institute of Management and Technology, Madalavarigudem, Vijayawada, Andhra Pradesh*

*Abstract*— **Online social Network has rapidly become an online source for acquiring real-time his/her information about users. Twitter is an Online Social Network (OSN) where users can share anything and everything, such as news, opinions, and even their moods. Several arguments can be held over different topics, such as politics, Particular affairs, and important events. When a user tweets something, it is instantly conveyed to her followers, allowing them to outspread the received information at a much broader level. The project proposes the detection of spammers and fake user identification on Twitter using various machine learning algorithms like Random forest, Naive Bayes and extreme machine learning algorithm. Moreover, a taxonomy of the Twitter spam detection approaches is presented that classifies the techniques based on their ability to detect: Fake content, Spam based on URL, Spam in trending topics, Fake users. The presented techniques are also compared based on various features, such as user features, content features, graph features, structure features, and time features. We are hopeful that the presented study will be a useful resource for researchers to find the highlights of recent developments in Twitter spam detection on a single platform**

*Index Terms*— **Arguments, Broader Level, Outspread, Taxonomy**

## I. INTRODUCTION

Social media networking sites like as Twitter, Facebook, MySpace, Instagram and LinkedIn have been gaining huge popularity in the recent era. Twitter is the one of popular and largest networking sites compare to other social media sites. Twitter has been allows social media networking users to post latest news and share messages. The size of the posted messages is no more than 280 characters; such messages are called tweets in Twitter network. Generally online networking sites are being often used by people to express opinions on any product, emotions and beliefs on persons. To post feedback and reviews on purchased products, these networking sites can act as best platforms for users. Now a days 0.13% of messages advertised on Twitter are clicked, whenever users click on these links they accessed into spam data, which are higher than that of email spam [1]. Twitter and other online social networks are mainly used for sharing valuable information, huge user base have made them main target for cybercriminals and socialbots. In social networking sites, we can call spambots as socialbots.

It has gotten very honest to acquire any sort of data from any source over the world by utilizing the Internet. The expanded interest of social destinations grants clients to gather plentiful measure of data and information about clients. Gigantic volumes of information accessible on these destinations likewise draw the consideration of phony clients. Twitter has quickly become an online hotspot for getting ongoing data about clients. Twitter is an Online Social Network (OSN) where clients can share everything without exception, for example, news, sentiments, and even their dispositions. A few contentions can be held over various points, for example, governmental issues, current issues, and significant occasions. At the point when a client tweets something, it is immediately passed on to his/her supporters, permitting them to extend the got data at lot more extensive level. With the development of OSNs, the

need to consider and break down clients' practices in online social stages has escalated. Numerous individuals who don't have a lot of data with respect to the OSNs can undoubtedly be deceived by the fraudsters. There is additionally an interest to battle and spot controls on the individuals who use OSNs just for notices and in this manner spam others' records. As of late, the recognition of spam in long range interpersonal communication destinations pulled in the consideration of scientists. Spam discovery is a troublesome errand in keeping up the security of informal organizations.

People in the world engages in social media weekly once, with half of the people participating every day (48% users). One in six (16%) use social media to get information about an emergency. In the Figure 1 represents how many users are using the social networks are illustrated, facebook as whole is having many users. During an emergency, nearly one third of the people population would use social media to let others know they are safe. Face book is a podium to share news, requests for feedback, queries, and links with an engrossed community that help people a place to share information with each other. Face book contains People-based, groups, or webpage-based accounts and average user spends almost 3 hours per day on Facebook.

## II. RELATED WORK

Shivangi Ghee Wala et.al, proposed OSNs also made a variety of attempts to secure sensitive details from multiple privacy risks. However, designers believe that there is now a lack of such a conceptual framework to which data protection instruments must be developed, in spite of a relevance of these proposals. A threat concept should be at a heart of this model. We therefore suggest a risk measure for OSNs throughout this project. Their intention is to connect danger level to social network users so that others may consider how dangerous it would have been to have connections with them while sharing private details. They quantify risk thresholds on a basis of similarities and profit indicators, taking a consumer danger attitudes into consideration. In particular, we adopt an active risk estimation teaching approaches in which user risk behavior is learned from a few necessary user interactions. This same process of risk assessment mentioned in this article has also been designed and tested on actual data.

Rohit Kumar Kaliyar et.al, proposed a method titled as Fake News Detection Using a Deep Neural Network. While there has been detailed analysis and review of an effects of online forums with person-to-person conversation forms, integrating electronic communication platforms in co-located classes has been studied far less carefully. This study examined a perspectives and expectations of middle school students regarding two separate modes of conversation in co-located classrooms: face to face (F2F) and synchronous, computer-mediated contact (CMC). Which study is accessible in French? They therefore differentiate between students who are considered to be involved in F2F classroom conversations and those who are generally silent. These results emphasize a benefit of CMC against F2F conversations in co-locations and demonstrate that F2F and computer-mediated communication are experienced differently by different students ("active" and "silent"). Sybil cyber threats and computer network breaches cut significant security implications.

Aditi Gupta et.al, proposed a strategy named as Towards Distinguishing Counterfeit Client Records in Facebook. Individuals are profoundly subject to OSNs which have pulled in light of a legitimate concern for digital crooks for completing various pernicious exercises. A whole industry of bootleg market administrations has developed which offers counterfeit records-based administrations available to be purchased. They, in this manner, in our work, center around distinguishing counterfeit records on an extremely well known (and hard for information assortment) online interpersonal organization, Facebook. Key commitments of our work are as per a following. A main commitment has been assortment of information identified with genuine and counterfeit records on Facebook. Because of severe protection settings and consistently advancing Programming interface of Facebook with every variant including more limitations, gathering client accounts information turned into a significant test. Their subsequent commitment is a utilization of client channel data on Facebook to comprehend client profile action and distinguishing a broad arrangement of 17 highlights which assume a key part in segregating counterfeit clients on Facebook with

genuine clients. Third commitment is a utilization these highlights and distinguishing a key AI based classifiers who perform well in recognition task out of an aggregate of 12 classifiers utilized. Fourth commitment is a recognizing which kind of exercises (like, remark, labeling, sharing, and so on).

Dr. Priyanka Harjule et.al, proposed a technique named as Unwavering quality of News. This present project's motivation is to explore an ideas, approaches and calculations for distinguishing counterfeit news stories and their makers from online web-based media stages and evaluating their exhibition. This paper presents two models for location of phony news. First by text order where distinctive classifier models were applied and it was discovered that RNN (LSTM) gave a best exactness of 93 %. Second by swarm examination where Boundary tuning technique gave a best exactness of 80 %. In current occasions, in view of a headway of online media stages, counterfeit news identifying with various purposes has been expanding step by step. Counterfeit News on a web is characterized as a manufactured article with a goal to delude, normally for benefitting. Counterfeit news and tricks have been there since before a coming of a Web. Tricks have existed for quite a while, since a "Incomparable moon scam" distributed in 1835. Alongside a expansion in a utilization of online media stages like Facebook, Twitter and so on word gets out quickly among a large number of clients inside an exceptionally limited capacity to focus time.

FAIZA MASOOD et.al, proposed a technique named as Spammer Discovery and Phony Client ID on Informal organizations. In this paper, they play out a survey of strategies utilized for distinguishing spammers on Twitter. We are confident that an introduced investigation will be a helpful asset for analysts to discover a feature of late advancements in Twitter spam recognition on a solitary stage. Informal communication destinations connect with a great many clients around a globe. A client' associations with these social destinations, for example, Twitter and Facebook have a gigantic effect and periodically unfortunate repercussions for everyday life. An unmistakable interpersonal interaction locale have transformed into an objective stage for a spammers to scatter an enormous measure of unessential and pernicious data. Twitter, for instance, has gotten one of a most excessively utilized foundation everything being equal and along these lines permits an absurd measure of spam.

Counterfeit clients send undesired tweets to clients to advance administrations or sites that influence authentic clients as well as disturb asset utilization. In this paper, they played out a survey of strategies utilized for identifying spammers on Twitter. Furthermore, they likewise introduced scientific categorization of Twitter spam location draws near and sorted them as phony substance identification, URL based spam discovery, spam recognition in moving points, and phony client recognition procedures. They additionally looked at an introduced procedure dependent on a few highlights, for example, client highlights, content highlights, chart highlights, structure highlights, and time highlights. In addition, a method was additionally looked at regarding their predefined objectives and datasets utilized. It is foreseen that an introduced survey will assist specialists with finding a data on best-in-class Twitter spam recognition strategies in a merged structure.

Regardless of an improvement of proficient and successful methodologies for a spam discovery and phony client recognizable proof on Twitter there are as yet certain open regions that require impressive consideration by a specialist. An issue is quickly featured as under.

## III. PROPOSED METHODOLOGY

We performed out a survey of procedures utilized for recognizing spammers on Twitter. What's more, we additionally introduced scientific classification of Twitter spam location draws near and sorted them as phony substance recognition, URL based spam discovery, spam identification in moving subjects, and phony client discovery procedures. We additionally looked at the introduced procedures dependent on a few highlights, for example, client highlights, content highlights, diagram highlights, structure highlights, and time highlights. In addition, the methods were likewise analyzed as far as their predetermined objectives and datasets utilized.

A data was taken from Tweepy and partitioned it into preparing set and test set, using ML algorithm and preparing information we've prepared our classifier model. So as to utilize ML model to recognize counterfeit twitter accounts, we required a marked

assortment of clients, pre named phony or certifiable. We get constant information from tweepy API which comprise of 2798 preparing set and 578 test set. A dataset is separated into 70% (training set) and 30% (test set) on which information exploratory examination has been done just as it is additionally investigated to highlight extraction and highlight designing, both preparing and test set are spared in CSV format.

We did Exploratory Data Analysis. In which we found out all a NULL values in a dataset. This is a heat map of all NULL values. It is evident from a heat map that most users don''t have their location, description or URL mentioned in their profile. Above are two plots for Bots friend''s vs. Followers'' and „Non-Bots Friends vs. followers'', they depict a general trend on a characteristics of bots i.e. they possess more followers count compared to friends. On a other hand, Non-Bots have generally an equal number of both friends and followers with some slight variance.

Key Points of Proposed Methodology

- Online social networks, such as Facebook, Twitter, and Weibo have played an important role in people''s common life. Most existing social network platforms, however, face the challenges of dealing with undesirable users and their malicious spam activities that disseminate content, malware, viruses, etc. to the legitimate users of the service.

- The spreading of spam degrades user experience and also negatively impacts server-side functions such as data mining, user behavior analysis, and resource recommendation. In this paper, an extreme learning machine (ELM)-based supervised machine is proposed for effective spammer detection.

- A set of features is then extracted from message content and user behavior and applies them to the ELM-based spammer classification algorithm.

Advantages of Proposed Methodology

These IDs are extraordinary 64-piece unsigned whole numbers, which depend on schedule, rather than being successive. A full ID is made out of a

timestamp, a specialist number, and a succession number. Thus it is used to classify the fake users.

- The proposed methodology also focuses on the friends_count, it ought to be in appropriate proportion with a follower_count. As a friends_count is subject to users count; these two properties are emphatically corresponded.

- The proposed implementation onthe Twitter list is actually a rundown of individuals on Twitter that are by one way or another associated. They may have a place with a specific classification (for example news associations), or they may be associated through their substance (for example identified with planting), or they may even be associated through an occasion that they're all joining in (for example Oscars).
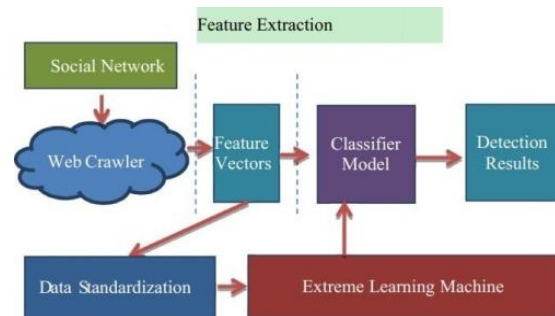


Figure 2: Extreme Learning Machine

- The proposed method has capable of Profile analysis; Profile that who has not given enough data on a profile are in all likelihood has a place with bot or spammers.

- Spammers and individuals who pester others regularly use accounts without a profile picture used for characterization and relapse errands quite well.
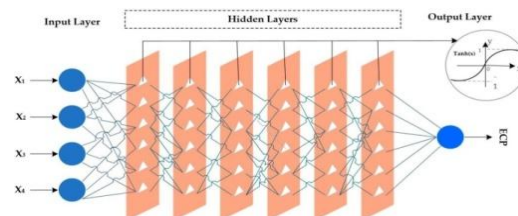


Figure 1: Spammer detection model

Figure 1 illustrates the basic concept of the proposed spammer detection model. In this

solution, training data are converted into a series of feature vectors that consists of a set of formulated attribute values. These vectors construct the input value of a supervised machine learning algorithm. After training, a classification model is applied to distinguish whether the specific user belongs to either a normal user or spammer. Because spammers and non- spammers have different social behaviors, it is capable to distinguish abnormal behaviors from legitimate ones. In this paper, we used a model based on 18 features, which were the following: the number of followees, the number of followers, the number of messages, the number of friends following each other, the number of favorites, the number of created days, fraction of followees per followers, fraction of original messages, number of messages per day, the average number of reposts, the average number of comments, average number of likes, the average number of URLs, the average number of pictures, the average number of hashtags, the average number of user mentions, fraction of messages containing URLs, and fraction of messages containing pictures.

Extreme Learning Machine (ELM) is a proficient networking system which is motivated by the biological neural networks. ELM consists a huge data which are consistent in various networks to consent the understanding among the networks are represented as neurons are operated similarly. Each neuron is associated with one another by association interface linked with weight that contains data such as input signal (input image) used to resolve a specific issue due to the obtained weight typically stimulate or restrain the signal is impart. Every neuron comprises of inner state named as an activation signal. Output signals are generated by integrating the input signals and activation rule and forwarded to remaining units. Generally, the neural networks are employed to diminish the noisy data and it has potential to classify patterns on which they have not been trained. ELM comprises of 3 layers includes input layer, a hidden layer, and output layer shown in figure 2. The output of a neuron is either directly or indirectly fed back to its input through linked neurons used in complex pattern recognition tasks.

Input layer – The input layer comprises of a number of neurons related to the number of inputs to the neuronal network. It has passive nodes which not participate in the actual signal modification, but only transmits the signal to another layer.

Hidden layer – The hidden layer consists of a random number of layers with an arbitrary number of neurons which participate in the signal modification.

Output layer - The number of neurons in the output layer relates to the number of the output values of the neural network and the nodes are active ones.

Transfer Function - A neuron more often gets numerous inputs with own relative weight, which needs on the processing element's summation function. Weights are adaptive coefficients that decide the intensity of the input signal as enlisted by the artificial neuron and it is a measure of an input's connection strength. The inputs and weights are considered as vectors are represented as (u1, u2, u3...un) and the weights are (w1,w2,w3...wn). The input and weighting coefficients are integrated into various ways before forwarding to the transfer function. The summation function also selects the minimum, maximum, majority, product or several normalizing algorithms. The neural inputs are combined to determine the network hypothesis. In some cases, the summation functions have an additional activation function applied to the output before it is passed to the transfer function to permit the output to vary with respect to time.

The consequence of the summation work is changed to an output by a transfer function. The summation can be contrasted with some edge to choose the neural output. When the summation is superior than the threshold value, at that point it creates a signal and it is fewer than the threshold, no signal (or some inhibitory signal) is created in both response are noteworthy. The threshold or transfer function is non-linear..

## IV. SIMULATIONAL RESULTS

Table 1: Comparison of existing and proposed approaches in terms of obtained accuracy and precision values.

| Method | Accuracy | Precision |
|---|---|---|
| Random forest | 60 | 53.4 |

| Naïve bayes | 60.66 | 55 |
| SVM | 86.66 | 43.33 |
| Extreme machine learning | 87.5 | 82.5 |

Table 2: Obtained values of recall and F-measure using existing random forest, naïve Bayes, SVM and proposed EML approaches.

| Method | Accuracy | Precision |
| --- | --- | --- |
| Random forest | 60 | 53.4 |
| Naïve bayes | 60.66 | 55 |
| SVM | 86.66 | 43.33 |
| Extreme machine learning | 87.5 | 82.5 |

## V. CONCLUSION

In this paper, we performed a review of techniques used for detecting spammers on Twitter. In addition, we also presented taxonomy of Twitter spam detection approaches and categorized them as fake content detection, URL based spam detection, spam detection in trending topics, and fake user detection techniques. We also compared the presented techniques based on several features, such as user features, content features, graph features, structure features, and time features. Moreover, the techniques were also compared in terms of their specified goals and datasets used. It is anticipated that the presented review will help researchers find the information on state-of-the-art Twitter spam detection techniques in a consolidated form. Despite the development of efficient and effective approaches for the spam detection and fake user identification on Twitter, there are still certain open areas that require considerable attention by the researchers. The issues are briefly highlighted as under: False news identification on social media networks is an issue that needs to be explored because of the serious repercussions of such news at individual as well as collective level. Another associated topic that is worth investigating is the identification of rumor sources on social media. Although a few studies based on statistical methods have already been conducted to detect the sources of rumors, more sophisticated approaches, e.g., social network based approaches, can be applied because of their proven effectiveness

## REFERENCES

[1] Social Networks Analysis and Mining (ASONAM) 2018 Aug 28 (pp. 1191-1198). IEEE.

[2] Pakaya FN, Ibrohim MO, Budi I. Malicious Gheewala S, Patel R. ML based Twitter Spam account detection: a review. In2018 Second International Conference on Computing Methodologies and Communication (ICCMC) 2018 Feb 15 (pp. 79-84). IEEE.

[3] Kaliyar RK. Fake news detection using a deep neural network. In2018 4th International Conference on Computing Communication and Automation (ICCCA) 2018 Dec 14 (pp. 1-7). IEEE.

[4] Erşahin B, Aktaş Ö, Kılınç D, Akyol C. Twitter fake account detection. In2017 International Conference on Computer Science and Engineering (UBMK) 2017 Oct 5 (pp. 388-392). IEEE.

[5] Gupta A, Kaushal R. Towards detecting fake user accounts in Facebook. In2017 ISEA Asia Security and Privacy (ISEASP) 2017 (pp. 1-6). IEEE.

[6] Alom Z, Carminati B, Ferrari E. Detecting spam accounts on Twitter. In2018 IEEE/ACM International Conference on Advances in Account Detection on Twitter Based on Tweet Account Features using Machine Learning. In2019 Fourth International Conference on Informatics and Computing (ICIC) 2019 Oct 16 (pp. 1-5). IEEE.

[7] Jardaneh G, Abdelhaq H, Buzz M, Johnson D. Classifying Arabic tweets based on credibility using content and user features. In2019 IEEE Jordan International Joint Conference on Electrical Engineering and Information Technology (JEEIT) 2019 Apr 9 (pp. 596-601). IEEE.

[8] Harjule P, Sharma A, Chouhan S, Joshi S. Reliability of News. In2020 3rd International Conference on Emerging Technologies in Computer Engineering: ML and Internet of Things (ICETCE) 2020 Feb 7 (pp. 165-170). IEEE.

[9] Simon NT, Elias S. Detection of fake followers using feature ratio in self-organizing maps. In2017 IEEE SmartWorld, Ubiquitous Intelligence & Computing, Advanced & Trusted Computed, Scalable Computing &

Communications, Cloud & Big Data Computing, Internet of People and Smart City Innovation(SmartWorld/SCALCOM/UIC/ATC/ CBDCom/IOP/SCI) 2017 Aug 4 (pp. 1-5). IEEE.