

SPAMMER DETECTION AND FAKE USER IDENTIFICATION ON SOCIAL NETWORKS

¹Mrs.Soujenya Voggu,²B.Sri Nandini,³Kadiyala Devika ⁴Pragathi Vadlamudi

¹Professor, Dept. of CSE-Cyber Security, Geethanjali College of Engineering and Technology Cheeryal (V),
Keesara (M), Medchal(D), Hyderabad, Telangana 501301,

soujenyav.cse@gcet.edu.in

^{2, 3, 4, BTech} Student, Dept. of CSE-Cyber Security, Geethanjali College of Engineering and Technology Cheeryal (V),
Keesara (M), Medchal(D), Hyderabad, Telangana 501301,

srinandhini8181@gmail.com,devikakadiyala62@gmail.com,pragathi.vadlamudi@gmail.com

ABSTRACT:

Social networking sites engage millions of users around the world. The users' interactions with these social sites, such as Twitter and Facebook have a tremendous impact and occasionally undesirable repercussions for daily life. The prominent social networking sites have turned into a target platform for the spammers to disperse a huge amount of irrelevant and deleterious information. Twitter, for example, has become one of the most extravagantly used platforms of all times and therefore allows an unreasonable amount of spam. Fake users send undesired tweets to users to promote services or websites that not only affect legitimate users but also disrupt resource

consumption. Moreover, the possibility of expanding invalid information to users through fake identities has increased that results in the unrolling of harmful content. Recently, the detection of spammers and identification of fake users on Twitter has become a common area of research in contemporary online social Networks (OSNs). In this paper, we perform a review of techniques used for detecting spammers on Twitter. Moreover, a taxonomy of the Twitter spam detection approaches is presented that classifies the techniques based on their ability to detect: (i) fake content, (ii) spam based on URL, (iii) spam in trending topics, and (iv) 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.

Keywords:Spammer Detection, Social Networks

I INTRODUCTION

The project focuses on addressing the pervasive issue of spam and fake user activity on Twitter, one of the most widely utilized social networking platforms. As social media sites like Twitter and Facebook have become integral parts of daily life for millions worldwide, they have also become lucrative targets for spammers disseminating irrelevant and harmful information. This project specifically delves into the challenges posed by spammers on Twitter, examining the undesirable repercussions on legitimate users and the overall disruption of resource consumption. The primary concern revolves around the extravagant use of Twitter, allowing an unreasonable influx of spam. Fake users employ various tactics, such as sending undesired tweets, to promote services or websites. This not only

negatively impacts authentic users but also facilitates the spread of invalid information through fake identities, resulting in the dissemination of harmful content. Given the gravity of the situation, the project undertakes a comprehensive review of existing techniques employed for detecting spammers on Twitter.

II. LITERATURE SURVEY

1.Spam Detection In Twitter Authors: C. Chen, Y. Wang, J. Zhang, Y. Xiang, W. Zhou, And G. Min

Twitter spam has become a critical problem nowadays. Recent works focus on applying machine learning techniques for Twitter spam detection, which make use of the statistical features of tweets. In our labeled tweets data set, however, we observe that the statistical properties of spam tweets vary over time, and thus, the performance of existing machine learning-based classifiers decreases. This issue is referred to as “Twitter Spam Drift”. In order to tackle this problem, we first carry out a deep analysis on the statistical features of one million spam tweets and one million non-spam tweets, and then propose a novel L fun scheme. The proposed scheme can discover “changed” spam tweets from unlabeled

tweets and incorporate them into classifier's training process. A number of experiments are performed to evaluate the proposed scheme. The results show that our proposed L fun scheme can significantly improve the spam detection accuracy in real-world scenarios.

2. Automatically Identifying Fake News in Popular Twitter Threads Author: C. Buntain And J. Gol Beck

Information quality in social media is an increasingly important issue, but web-scale data hinders experts' ability to assess and correct much of the inaccurate content, or "fake news," present in these platforms. This paper develops a method for automating fake news detection on Twitter by learning to predict accuracy assessments in two credibility-focused Twitter datasets CREDBANK, a crowd sourced dataset of accuracy assessments for events in Twitter, and PHEME, a dataset of potential rumors in Twitter and journalistic assessments of their accuracies. We apply this method to Twitter content sourced from Buzz Feed's fake news dataset and show models trained against crowd sourced workers outperform models based on journalists assessment and models trained on a pooled dataset of both crowd sourced workers and journalists. All three datasets, aligned into a uniform format, are

also publicly available. A feature analysis then identifies features that are most predictive for crowd sourced and journalistic accuracy assessments, results of which are consistent with prior work. We close with a discussion contrasting accuracy and credibility and why models of non-experts outperform models of journalists for fake news detection in Twitter.

3. A Performance Evaluation Of Machine Learning-Based Streaming Spam Tweets Detection Author: C. Chen, J. Zhang, Y. Xie, Y. Xiang,W. Zhou, M. M. Hassan, A. Alelaiwi

The popularity of Twitter attracts more and more spammers. Spammers send unwanted tweets totwitter users to promote websites or services, which are harmful to normal users. In order to stop spammers, researchers have proposed a number of mechanisms. The focus of recent works is on the application of machine learning techniques into Twitter spam detection. However, tweets are retrieved in a streaming way, and Twitter provides the Streaming API for developers and researchers to access public tweets in real time. There lacks a performance evaluation of existing machine learning-based streaming spam detection methods. In this paper, we bridged the gap by carrying out a performance evaluation, which was

from three different aspects of data, feature, and model. A big groundtruth of over 600 million public tweets was created by using a commercial URL-based security tool. For real-time spam detection, we further extracted 12 lightweight features for tweet representation. Spam detection was then transformed to a binary classification problem in the feature space and can be solved by conventional machine learning algorithms. We evaluated the impact of different factors to the spam detection performance, which included spam to no spam ratio, feature discretization, training data size, data sampling, time-related data, and machine learning algorithms. The results show the streaming spam tweet detection is still a big challenge and a robust detection technique should take into account the three aspects of data, feature, and model

III SYSTEM ANALYSIS

EXISTING SYSTEM

- Sahami et al. proposed textual and non-textual and domain-specific features and learned naiveBayes classifier to segregate spam emails from legitimate ones. Schafer proposed metadata- based approaches to detect botnets based on compromised email accounts to diffuse mail spams. Spam campaigns on Facebook were analyzed by

Gao et al. using a similarity graph based on semantic similarity between posts and URLs that point to the same destination. • Furthermore, they extracted clusters from a similarity graph, wherein each cluster represents a specific spam campaign. Upon analysis, they determined that most spam sources were hijacked accounts, which exploited the trust of users to redirect legitimate users to phishing sites.

- Yang et al. and Ahmed and Abolish used content- and interaction based attributes for learning classifiers to segregate spammers from benign users on different OSNs.

- Yang et al. and Ahmed and Abolish analyzed the contribution of each feature to spammer detection, whereas Yang et al. conducted an in-depth empirical analysis of the evasive tactics practiced by spammers to bypass detection systems. They also tested the robustness of newly devised features.

- Zhu et al. used a matrix factorization technique to find the latent features from the sparse activity matrix and adopted social regularization to learn the spam discriminating power of the classifier on the Renner network, one of the most popular OSNs in China. Another spammer detection

approach in social media was proposed by Tan ET a.

Disadvantages

- There are no Hybrid techniques to classify different spams behaviors.
- There is no spambot detection techniques

PROPOSED SYSTEM

- In the proposed system, the system proposes a Fake User Identification approach for detecting social spam bots in Twitter, which utilizes an amalgamation of metadata-, content-, interaction-, and community based features. In the analysis of characterizing features of existing approaches, most network-based features are not defined using user followers and underlying community structures, thereby disregarding the fact that the reputation of user in a network is inherited from the followers (rather than from the ones user is following) and community members. Therefore, the system emphasizes the use of followers and community structures to define the network-based features of a user.
- The system classifies set of features into several broad categories, namely, (i) fake content, (ii) spam based on URL, (iii) spam

in trending topics, and (iv) fake users, wherein the network category is further classified into interaction and community based features. Metadata features are extracted from available additional information regarding the tweets of a user, whereas content-based features aim to observe the message posting behavior of a user and the quality of the text that the user uses in posts. Network based features are extracted from user interaction network.

Advantages

- A novel study that uses community-based features with other feature categories, including metadata, content, and interaction, for detecting automated spammers.
- Used Hybrid technique to classify spammers such as random forest, decision tree, and Bayesian network.

IV IMPLEMENTATION

Architecture:

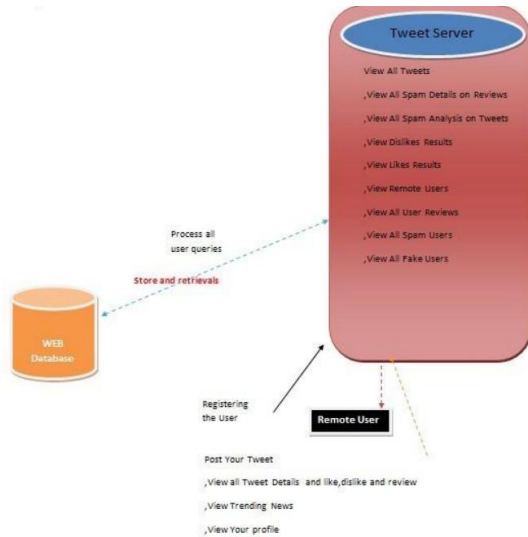


Fig-1. Architectures of the system model

Data Collection Module:

Collects user data such as account creation details, activity logs, IP addresses, device information, etc.

Feature Extraction Module:

Extracts relevant features from the collected data, such as frequency of posting, patterns of interaction, language usage, etc.

Machine Learning Models:

Utilizes machine learning algorithms for classification tasks, trained on labeled data to identify spam or fake users based on extracted features.

Natural Language Processing (NLP) Module:

Analyzes textual data for sentiment, grammar, and coherence to identify suspicious patterns in language usage. Behavioral Analysis Module: Analyzes user behavior patterns such as clickthrough rates, session duration, time of activity, etc., to detect anomalies indicative of spam or fake accounts.

Social Network Analysis Module:

Design:

Implement a Python class named SocialNetworkAnalyzer for analyzing user connections and network structures.

Methods:

Build network (user connections): Construct a graph representing user connections. detect_suspicious_clusters (graph): Identify suspicious clusters or patterns within the network. Dependencies: NetworkX or igraph for graph manipulation and analysis.

Decision Making Module:

Design:

Develop a Python class named Decision Maker to integrate results from different modules and make final classification decisions.

Methods:

Combine results (results):

Combine results from various analysis modules to make classification decisions. Dependencies: None specific, as it depends on the outputs of other modules.

System Administration Module:

Design:

Create a Python script or class responsible for system administration tasks such as user interface, logging, and reporting.

Functionality:

Provide a user interface for system administrators to interact with the system, log activities, and generate

Reports. Dependencies: Flask or Django for web-based user interfaces, logging module 18 for activity logging, matplotlib or seaborn for visualization

MODULES

1. Graphical User Interface (GUI):

This module encompasses the user interface designed using Django, providing an interactive platform for users to engage with the system. It includes features for input, display, and user interaction related to spam detection and classification.

2. Algorithm Development:

This module involves the implementation of algorithms in Python to extract metadata, content, interaction, and community-based features for spam detection. The algorithms utilize a hybrid approach, incorporating random forest, decision tree, and Bayesian network classifiers for efficient detection.

3. Database Management:

This module is responsible for managing the storage and retrieval of data related to users, tweets, and features. MySQL is employed as the database management system, ensuring efficient data organic.

Algorithms

1. Machine Learning Algorithms

- **Supervised Learning:** Classification algorithms such as Support Vector Machines (SVM), Decision Trees, Random Forest, Naive Bayes, and Logistic Regression are used to classify users or content as spam or legitimate.
- **Unsupervised Learning:** Clustering algorithms like K-Means, DBSCAN, or

anomaly detection techniques can identify patterns of behavior that are characteristic of spammers or fake users.

- **Deep Learning:** Neural networks, especially deep learning architectures like Convolutional Neural Networks (CNNs) or Recurrent Neural Networks (RNNs), can be used to learn complex patterns in user behavior or content for spam detection.

2. Natural Language Processing (NLP)

Techniques:

- Analyzing the content of messages or user profiles using techniques like sentiment analysis, topic modeling, or language models to detect spammy or fake content.
- Named Entity Recognition (NER) can identify suspicious entities or patterns in text that might indicate spam or fake users.

3. Graph-based Algorithms:

- Analyzing the network structure of interactions between users or entities to detect patterns indicative of spamming behavior. Graph algorithms like PageRank, community detection, or centrality measures can be useful for this purpose.

4. Behavioral Analysis:

- Monitoring user behavior such as posting frequency, time of activity, interaction patterns, etc., to identify anomalies or patterns consistent with spamming or fake accounts.
- User engagement metrics like the ratio of followers to following, the rate of interaction with other users, or sudden spikes in activity can signal spammy behavior.

5. Rule-based Systems:

- Utilizing predefined rules or heuristics based on known characteristics of spam or fake users. These rules can be simple thresholds (e.g., maximum number of posts per hour) or more sophisticated logic based on the analysis of historical data.

6. Human-in-the-loop Systems:

- Incorporating human judgment or feedback into the detection process, where suspicious cases flagged by automated algorithms are reviewed by human moderators for final decision making. This approach can help in handling complex

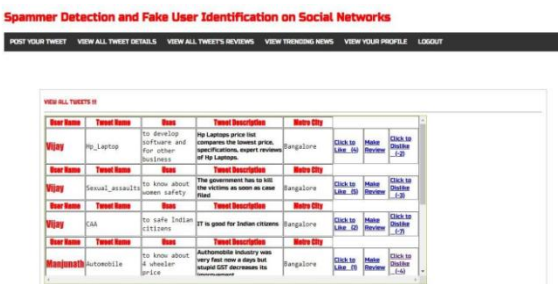
cases and adapting to evolving spamming techniques.

V RESULT AND DISCUSSION

User Login

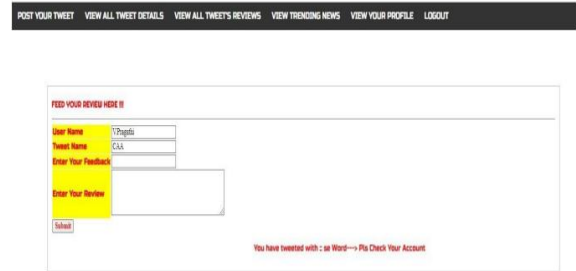


View all Tweets

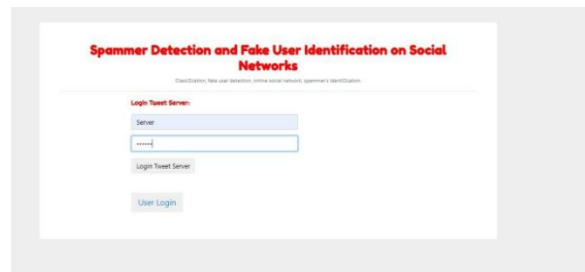


Review any Tweet:

Spammer Detection and Fake User Identification on Social Networks

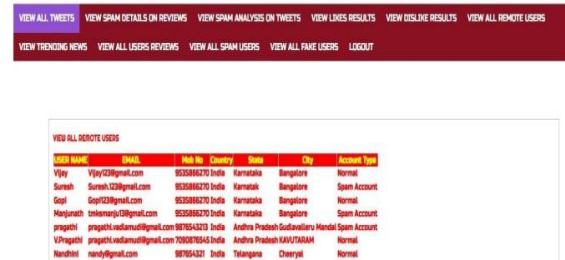


Server Login



View all Remote users

Spammer Detection and Fake User Identification on Social Networks



Display posts based on spam type

SELECT SPAM TYPE:

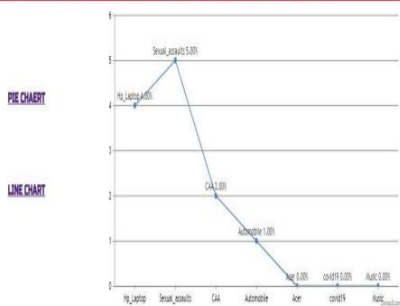
VIEW SPAM ANALYSIS ON CLIENT POSTS II

USER NAME	TWEET NAME	TWEET DESC	USERS	SPAM TYPE	LOCATION NAME
Vijay	Sexual_abuse	The government has to kill the victims as soon as case filed to know about women safety		Offensive	Bangalore
pragathi	con499	to kill the peoples	to kill the peoples	Offensive	Gustavalary Mandi

Line Chart for Spam Analysis

Spammer Detection and Fake User Identification on Social Networks

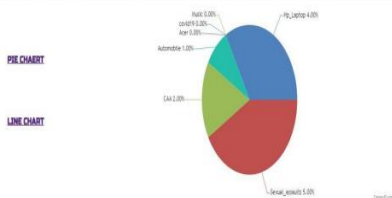
VIEW ALL TWEETS VIEW SPAM DETAILS ON REVIEWS VIEW SPAM ANALYSIS ON TWEETS VIEW LINES RESULTS VIEW ONLINE RESULTS VIEW ALL REMOTE USERS
VIEW TRENDING NEWS VIEW ALL USERS REVIEWS VIEW ALL SPAM USERS VIEW ALL FAKE USERS LOGOUT



Pie Chart for Spam Analysis

Spammer Detection and Fake User Identification on Social Networks

VIEW ALL TWEETS VIEW SPAM DETAILS ON REVIEWS VIEW SPAM ANALYSIS ON TWEETS VIEW LINES RESULTS VIEW ONLINE RESULTS VIEW ALL REMOTE USERS
VIEW TRENDING NEWS VIEW ALL USERS REVIEWS VIEW ALL SPAM USERS VIEW ALL FAKE USERS LOGOUT



View all use reviews

Spammer Detection and Fake User Identification on Social Networks

VIEW ALL TWEETS VIEW SPAM DETAILS ON REVIEWS VIEW SPAM ANALYSIS ON TWEETS VIEW LINES RESULTS VIEW ONLINE RESULTS VIEW ALL REMOTE USERS
VIEW TRENDING NEWS VIEW ALL USERS REVIEWS VIEW ALL SPAM USERS VIEW ALL FAKE USERS LOGOUT

VIEW ALL USED REVIEWS II

User Name	Tweet Name	Review	Spam/Review Analysis	Review Date and Time	Feedback
Suresh	Pg_Latire	hp fuck gives only less hour battery back up	Sexual	2019-12-24 18:01:28.865273	Can improve battery back up
Suresh	CIA	This scheme is not good and what a fuck scheme announced.	Sexual	2019-12-25 16:33:44.742187	Not satisfied
Opal	CIA	This is excellent scheme which secures Indian Citizens	Positive	2019-12-25 16:35:29.588976	Can change few things
Opal	CIA	What a fuck scheme is this	Sexual	2019-12-25 17:27:21.674884	Not satisfied

View all spam users

VIEW ALL SPAM USERS II

USER NAME	EMAIL	Mobile No	Country	State	City	Account Type	SPAM Reason
Suresh	Suresh.C2@gmail.com	933586270	India	Karnatak	Bangalore	Spam Account	Tweeted with Offensive Word
Manjunath	mkmanjya1@gmail.com	933586270	India	Karnataka	Bangalore	Spam Account	Tweeted with Volgar Word
pragathi	pragathivadamsu@gmail.com	987654321	India	Andhra Pradesh	Gustavalary Mandi	Spam Account	Tweeted with Offensive Word

VI CONCLUSION

In conclusion, the evolution of spammer detection and fake user identification represents a vital pursuit in safeguarding online communities and platforms. Through a combination of advanced technologies such as machine learning, natural language processing, and behavioral analysis, significant strides have been made in identifying and mitigating spamming activities and fake accounts. However, the

battle against malicious actors is ongoing, necessitating continuous innovation and adaptation of detection methods to keep pace with evolving tactics. Furthermore, collaboration between automated systems, human moderation, and user participation remains crucial for maintaining the integrity and trustworthiness of online environments. As these detection mechanisms continue to evolve, the goal remains steadfast: to create safer, more authentic digital spaces where users can engage with confidence and trust in the authenticity of their interactions. In conclusion, the evolution of spammer detection and fake user identification represents a vital pursuit in safeguarding online communities and platforms. Through a combination of advanced technologies such as machine learning, natural language processing, and behavioral analysis, significant strides have been made in identifying and mitigating spamming activities and fake accounts. However, the battle against malicious actors is ongoing, necessitating continuous innovation and adaptation of detection methods to keep pace with evolving tactics. Furthermore, collaboration between automated systems, human moderation, and user participation remains crucial for maintaining the integrity and trustworthiness of online environments.

As these detection mechanisms continue to evolve, the goal remains steadfast: to create safer, more authentic digital spaces where users can engage with confidence and trust in the authenticity of their interactions.

FUTURE ENHANCEMENT

The Future enhancement of spammer detection and fake user identification systems is imperative to maintain the integrity and trustworthiness of online platforms. By leveraging advanced techniques such as machine learning, natural language processing, and behavioral analysis, platforms can effectively detect and mitigate spamming behavior and identify fake accounts. Incorporating relocation analysis, semantic understanding, and real-time monitoring further fortifies these systems against sophisticated spamming tactics. Moreover, collaborative efforts between automated algorithms and human moderators, coupled with user education initiatives, foster a robust ecosystem where users can engage safely and authentically.

As technology evolves and spammers adapt, ongoing innovation and adaptation remain essential to stay ahead of emerging threats and preserve the quality and reliability of online interactions. Despite the development of efficient and successful ways for spam

detection and fake user identification on Twitter, there are still certain gaps in the study that need to be addressed. The following are a few of the issues: Because of the substantial ramifications of false news on an individual and communal level, false news identification on social media networks is a subject that needs to be investigated. The identification of rumor origins on social media is another related topic worth researching. Although a few studies using statistical methods to discover the origin of rumors have already been undertaken, more complex approaches, such as social network-based approaches, can be used due to their demonstrated efficiency

VII REFERENCES

1. C. Chen, S. Wen, J. Zhang, Y. Xiang, J. Oliver, A. Alelaiwi, and M. M. Hassan, "Investigating the deceptive information in Twitter spam," *Future Gener. Comput. Syst.*, vol. 72, pp. 319–326, Jul. 2017.
2. M. Babcock, R. A. V. Cox, and S. Kumar, "Diffusion of pro- and antifalse information tweets: The black panther movie case," *Comput. Math. Org. Theory*, vol. 25, no. 1, pp. 72–84, Mar. 2019.
3. M. Washha, A. Qaroush, M. Mezghani, and F. Sedes, "A topic-based hidden Markov model for real-time spam tweets filtering," *Procedia Comput. Sci.*, vol. 112, pp. 833–843, Jan. 2017.
4. Spammers on Twitter based on content and social interaction," in *Proc. Int. Conf. Netw. Inf. Syst. Comput.*, pp. 413–417, Jan. 2015.
5. Prasadu Peddi (2015) "A review of the academic achievement of students utilizing large-scale data analysis", ISSN: 2057-5688, Vol 7, Issue 1, pp: 28-35.
6. Prasadu Peddi (2015) "A machine learning method intended to predict a student's academic achievement", ISSN: 2366-1313, Vol 1, issue 2, pp:23-37.

AUTHORS

Mrs. Soujanya Voggu, Professor Dept. of CSE-Cyber Security, Geethanjali College of Engineering and Technology Cheeryal (V), Keesara (M), Medchal(D), Hyderabad, Telangana 501301.

Email: soujenyav.cse@gcet.edu.in

Miss. B.Sri Nandini, Dept. of CSE-Cyber Security, Geethanjali College of Engineering and Technology Cheeryal (V), Keesara (M), Medchal(D), Hyderabad, Telangana 501301.

Email: srinandhini8181@gmail.com

Miss.Kadiyala Devika, Dept. of CSE-Cyber Security, Geethanjali College of Engineering and

Technology Cheeryal (V), Keesara (M),
Medchal(D), Hyderabad, Telangana 501301.

Email: devikakadiyala62@gmail.com

Miss. Pragathi Vadlamudi, Dept. of CSE-
Cyber Security, Geethanjali College of Engineering
and Technology Cheeryal (V), Keesara (M),
Medchal(D), Hyderabad, Telangana 501301.

Email: pragathi.vadlamudi@gmail.com