



BOTWISE : HARNESSING BIGDATA TO DETECT TWITTER BOTS

¹MOHAMMAD KOWSAR, ²MRS. SYED ZAHADA

¹ PG Scholar in the department of MCA at QIS College of Engineering & Technology (AUTONOMOUS), Vengamukkapalem, Ongole- 523272, Prakasam Dt., AP., India.

² Assistant Professor in the department of MCA at QIS College of Engineering & Technology (AUTONOMOUS), Vengamukkapalem, Ongole- 523272, Prakasam Dt., AP., India.

ABSTRACT:

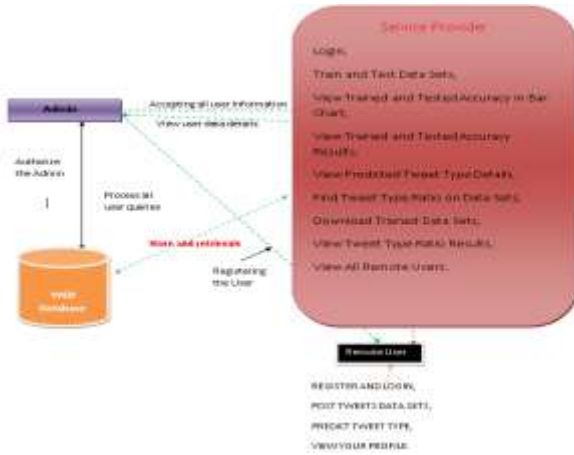
In nowadays world lots of people like a businessman's, Media, politicians, etc., uses Twitter daily & have become an important part of life. Twitter is one of the favourite social networking sites which let the individuals to express their sentiment on various topics like politics, sports, stock market, entertainment etc. It is one of the quickest means of transmission information. It extremely affects people's viewpoint. An increasing number of people on twitter but hide their identity for malignant purpose. It is dangerous for other users hence the necessity for identifying the twitter bots. So it is essential that tweets are sent by authentic users and not by twitter bots. A twitter bot transmits spam subject matters. Therefore detecting of bots helps to determine spam messages. The characteristics of twitter accounts are utilized as Features in machine learning algorithms to label users as genuine or fake. In this paper, we used three machine learning algorithms to detect the account is fake or real, which are Decision Tree, Random Forest, and Multinomial Naive Bayes The classification performance of the algorithms is compared with their accuracy. The accuracy given by the Decision tree algorithm is 93%, the Random Forest algorithm is 90% and the Multinomial Naive Bayes is 89%. Hence it is seen that the Decision tree gives more accuracy as compared to Random Forest and Multinomial Naive Bayes

INDEX: Social media, Twitter, big data analytics, shallow learning, deep learning, tweet-based bot detection.

INTRODUCTION:

Twitter is one of the most popular micro-blogging social media platforms that has millions of users. Due to its popularity, Twitter has been targeted by different attacks such as spreading rumors, phishing links, and malware. Tweet-based botnets represent a serious threat to users as they can launch large-scale attacks and manipulation campaigns. To deal with these threats, big data analytics techniques, particularly shallow and deep learning techniques have been leveraged in order to accurately distinguish between human accounts and tweet-based bot accounts. In this paper, we discuss existing techniques, and provide a taxonomy that classifies the state-of-the-art of tweet-based bot detection techniques. We also describe the shallow and deep learning techniques for tweet-based bot detection, along with their performance results. Finally, we present and discuss the challenges and open issues in the area of tweet-based bot detection.

SYSTEM ARCHITECTURE



METHODOLOGY

Data Preprocessing:

Once collected, the Twitter data undergoes preprocessing to ensure its quality and consistency. This involves tasks such as removing duplicates, handling missing values, and cleaning noisy or irrelevant data. Additionally, text data may be processed using techniques like tokenization, stemming, and stop-word removal to prepare it for analysis.

Ref.	Dataset	Tweets	Training Set	Testing Set	Pre-processing	Features	Classifier	Architecture Approach	Accuracy	FP	TP	Precision
[24]	Cisco-2017	11.4m	8.36m	N/A	GoFE-SMOTE	Tweet & account Metadata	LSTM	Decision tree / Supervised	96%	N/A	N/A	96%
[25]	Cisco-2017	11.4m	4.92m	2.80m	Globe	Tweet Text/Tweet & Account Metadata	R/LSTM	Decision tree / Unsupervised	95%	6%	84%	65%
[26]	Own	48,334	42,392	1,800	Word embedding	Tweet Metadata	Autoencoder LSTM	Generative / Unsupervised	87%	N/A	N/A	85%
[27]	Dataset	5,600	5,000	500	DeepLick	Tweet Text & metadata	CNN + LSTM	Decision tree / Supervised	N/A	N/A	N/A	88.4%
[28]	CLEF-2019	N/A	1,877	1,240	Word embedding	Tweet Text & metadata	CNN	Decision tree / Unsupervised	85%	N/A	N/A	87%
[29]	Twitter	20	500	25,677	Spam detection	IDs, screen name, location	Bayesian classification	Supervised	N/A	N/A	N/A	89%
[30]	TwitterFaire Project, Social Hierarchy, User Popularity Rank	9,907,680, 5,612,306, 150,378	N/A	N/A	Social bot detection	Over-based, user profile-based and social graph-based attributes	Deep Q-Learning (DQL)	Unsupervised	95%	N/A	N/A	88%
[31]	ASW EC2	AM	N/A	N/A	N/A	Tweet statistics + Category vector + Sentiment + LDA	Graph neural network	Unsupervised	89%	N/A	N/A	N/A
[32]	Sex and Gender Profiling 2015	42,000	280,000	144,000	Human tweet distribution	R-LSTM	Deepbox	Unsupervised	76.6%	N/A	N/A	N/A
[33]	CLEF 2019	3,018	2,875	1,240	Twitter bot	N/A	Convolutional neural network	Unsupervised	N/A	N/A	N/A	83.2%
[34]	ISIS dataset	9M	N/A	N/A	N/A	N/A	Deep neural network	Decision tree / Unsupervised and Semi-supervised	92%	N/A	N/A	90%

Table 1: Summary of Deep learning-based detection method

Classification Accuracy:

Classification accuracy is defined the numbers of the correct predictions are divided by the total number of inputs samples or total number of predictions made. It is mathematically given as:

$$CA = \frac{(NCP)}{(TNI)}$$

Allowing to this experiment the performance of artificial neural network is capable with highest accuracy is 98.9%. Decision tree is much closed result to ANN, having 92.7% accuracy.

Confusion Matrix:

Confusion matrix is doing well performance for the binary classification methods; we are also



In this module, the admin can view the list of users who all registered. In this, the admin can view the user’s details such as, user name, email, address and admin authorize the users.

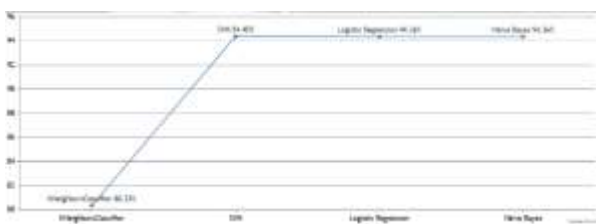
Remote User:

In this module, there are n numbers of users are present. User should register before doing any operations. Once user registers, their details will be stored to the database. After registration successful, he has to login by using authorized user name and password. Once Login is successful user will do some operations like REGISTER AND LOGIN, PREDICT TWEET TYPE, VIEW YOUR PROFILE.

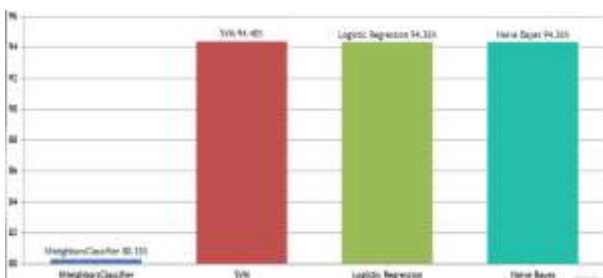
RESULTS ANALYSIS

Tweet Type	Ratio
Bot	48.0
Quality	52.0

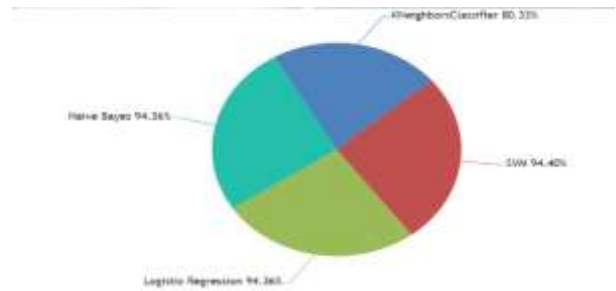
Prediction Ratio Details



Line Chart Prediction Results



Bar Chart Prediction Results



Pie Chart Prediction Results

CONCLUSION

Twitter is one of the most popular social media platforms that allows connecting people and helps organizations reaching out to customers. Tweet-based botnet can compromise Twitter and create malicious accounts to launch large-scale attacks and manipulation campaigns. In this review, we have focused on big data analytics, especially shallow and deep learning to fight against tweet-based botnets, and to accurately distinguish between human accounts and tweet-based bot accounts. We have discussed related surveys, and have also provided a taxonomy that classifies the state-of-the-art tweet-based bot detection techniques up to 2020. In addition, the shallow and deep learning techniques are described for tweet-based bot detection, along with their performance results. Finally, we presented and discussed the open issues and future research challenges.

FUTURE ENHANCEMENT

Feature Extraction: Various features can be extracted from Twitter data, such as tweet frequency, account age, posting patterns, follower count, friends count, profile information, language used, sentiment analysis of tweets, and more. Python libraries like



NLTK (Natural Language Toolkit) and TextBlob can be used for text analysis and sentiment analysis.

Bot Detection Rules: Developers can also create rules-based approaches to detect bots. This involves defining specific criteria or thresholds based on patterns observed in bot behavior. Python provides tools for implementing these rules efficiently.

Visualization: Python libraries like Matplotlib and Seaborn can be used to visualize the data and analysis results, making it easier to interpret and communicate findings.

REFERENCES

- 1) Jorge R, Javier M, Raúl M, Octavio L and Armando L, 2020, A one-class classification approach for bot detection on twitter, *Computers and Security*, 91, pp. 1-14.
- 2) Rodríguez-Ruiz J, Mata-Sánchez J, Monroy R, Loyola-González O, and López-Cuevas A, 2020, A one-class classification approach for bot detection on twitter, *Computers and Security*, 91, pp. 1-14.
- 3) Rahman M, Likhon A, Rahman A and Choudhury M, 2019, Detection of fake identities on twitter using supervised machine learning, PhD dissertation, Brac University.
- 4) Beskow D and Carley K, 2018, Bot conversations are different: leveraging network metrics for bot detection in twitter, *Proc. Int. Conf. On Advances in Social Networks Analysis and Mining (Barcelona, Spain)*, pp. 825-832.
- 5) Knauth J, 2019, Language-agnostic twitter-bot detection, *Proc. Int. Conf. on Recent Advances in Natural Language Processing (Varna, Bulgaria)*, pp. 550-558.
- 6) Daouadi K, Rebaï R, and Amous I, 2020, Real-time bot detection from twitter using the twitterbot+ framework, *Journal of Universal Computer Science*, 26, pp. 496-507 [4].
- 7) Dabiri S, and Heaslip K, 2019, Developing a twitter-based traffic event detection model using deep learning architectures, *Expert Systems with Applications*, 118, pp. 425-439.
- 8) Gabryel M, Damaševičius R, and Przybyszewski K, 2018, Application of the bag-of-words algorithm in classification the quality of sales leads, *Proc. Int. Conf. on Artificial Intelligence and Soft Computing*, Springer (Cham), pp. 615- 622.
- 9) Sahoo S, and Gupta B, 2019, Hybrid approach for detection of malicious profiles in twitter, *Computers and Electrical Engineering*, 76, pp. 65-81.
- 10) Beğenilmiş E, and Uskudarli S, 2018, Organized behavior classification of tweet sets using supervised learning methods, *Proc. Int. Conf. on Web Intelligence, Mining and Semantics (Novi Sad Serbia)*, pp. 1-9.
- 11) Ali D, Abadi M, and Dadfarnia M, 2018, Socialbothunter: botnet detection in twitter-like social networking services using semi-supervised collective classification,



Proc. Int. Conf. on Dependable, Autonomic and Secure Computing; Pervasive Intelligence and Computing, and Big Data Intelligence and Computing and Cyber Science and Technology Congress (Athens, Greece), pp. 496-503.

- 12) Torres J, Comesaña C, and García-Nieto P, 2019, Machine learning techniques applied to cybersecurity, International Journal of Machine Learning and Cybernetics, 10, pp.2823-2836 [10].
- 13) Satija T, and Kar N, 2019, Detecting malicious twitter bots using machine learning, Proc. Int. Conf. on Computational Intelligence, Security and Internet of Things, Springer (Singapore), pp. 182-194.
- 14) Kenta M, Takashi I, and Masayuki G, 2011, A proposal of extended cosine measure for distance metric learning in text classification, Proc. Int. Conf. on Systems, Man, and Cybernetics, Anchorage (AK, USA), pp. 1741-1746.
- 15) Kaggle, accessed on May 23rd 2020, <https://www.kaggle.com/charvijain27/detecting-twitter-bot-data>.

AUTHOR PROFILE:



Mrs. Syed Zahada, currently working as an Assistant Professor in the Department of MCA, QIS College of Engineering and Technology, Ongole, Andhra Pradesh. She did her MCA from Azad college

of computers, Hyderabad, Affiliated to Osmania University. Her area of interest is Machine Learning, Artificial intelligence, Cloud Computing and Programming language.



Ms. Mohammad Kowsar, currently pursuing Master of Computer Applications at QIS College of engineering and Technology (Autonomous), Ongole, Andhra Pradesh. She Completed B.Sc. in Statistics from Sri Harshini Degree College, Ongole, Andhra Pradesh. Her areas of interest are Machine learning & Cloud computing.