

A Deep Learning Approach for Robust Detection of Bots in Twitter Using Transformers

ELANCHELIYAN.R¹, ARAVINDHAN.R², RAHUL.G³, BRINDA.B.M⁴

^{1,2,3}Undergraduate student, Department of Computer Science and Engineering, Paavai College of Engineering

⁴Assistant Professor, Department of Computer Science and Engineering, Paavai College of Engineering

Abstract: - The volume of audiovisual content produced on social networks has increased tremendously in recent decades, and this information is quickly spread and consumed by a large number of people. In this situation, the disruption of fake news providers and bot accounts for spreading propaganda and sensitive content over the network has sparked practical research into using Artificial Intelligence to automatically assess the reliability of social media accounts (AI). In this research, we provide a multilingual method to solving the bot identification task in Twitter using deep learning (DL) approaches to assist end-users in determining the legitimacy of a particular Twitter account. To do so, a series of experiments were carried out using state-of-the-art Multilingual Language Models to generate an encoding of the user account's text-based features, which were then concatenated with the rest of the metadata to create a potential input vector on top of a Dense Network called Bot-DenseNet. As a result, this article evaluates the language limitation from earlier research that just evaluated the metadata information or the metadata information together with certain basic semantic text aspects when encoding the user account. Furthermore, the Bot-DenseNet generates a low-dimensional representation of the user account that may be used in any Information Retrieval (IR) application.

Key words: Deep Learning, Transformers

1. INTRODUCTION

Due to the benefits of posting, disseminating, and exchanging enormous volumes of multimedia material throughout the network, social media platforms such as Twitter and Facebook have developed a large level of popularity and influence among millions of users in recent years. As a result, as noted in [22], these platforms enable users to construct a digital community, which has enabled users to not only discover and embrace new relationships, but also to maintain and strengthen old ones. On the other hand, due to both their large influence on people's lifestyles and their evolution as a potential communication tool, these platforms have exponentially increased their appeal for marketing and commercial purposes by analyzing user behavior and opinion in various topics or events such as political elections. As a result, numerous research studies in the social media sector have been developed for various goals, such as sentiment analysis [35], traffic control [48], and consumer behavior mining [4].

Furthermore, with the disruptive expansion of Artificial Intelligence (AI) algorithms, detecting bots or untrustworthy sources has become a critical topic that has to be explored. It

sparked numerous studies and publications with the goal of developing robust autonomous systems to improve the quality of experience of consumers on such platforms by minimizing privacy threats while also boosting platform trustworthiness.

As a result, the goal of this paper is to advance the state-of-the-art in this field by proposing a novel method for automatically

(i) Encoding an input user account as a low-dimensional feature vector regardless of its language, and

(ii) Using a Deep Neural Network (DNN) named Bot-DenseNet, identifying the input encoding vector as a suspicious bot account with a particular likelihood.

(iii) Creating a low-dimensional embedding that represents the user account's original input encoding vector and can be used for any other Information Retrieval purpose (IR).

Our research focuses on identifying Bot accounts on Twitter by looking at three key factors: the account's activity level, popularity, and profile information.

The global collection of characteristics may be divided into two modalities: metadata and text-based descriptors, the latter of which is encoded using new Language Model Embedding's (LME) to overcome the language constraint that has plagued earlier research.

The rest of the paper is structured as follows: Previous studies and investigations into the Bot identification framework are discussed in Section II. Section III highlights the essential components of the proposed multilingual approach, including Section III-A, which describes how the input encoding vector is generated, and Section III-B, which summarizes the suggested Bot-DenseNet model's architecture. Then, in Section IV, you'll learn about the various tests that were carried out, as well as the results and breakthroughs that were made. Section V summarizes the overall conclusions as well as future study.

II. RELATED WORK

Artificial intelligence (AI) techniques such as Deep Learning (DL) and Machine Learning (ML) methods have recently gained popularity and interest in many applied research and industry services related to social media analysis, where sentiment analysis and text classification have been the central focus of these investigations, particularly for search engines and recommender systems.

As writers point out in [8,] the essence of sentiment analysis is extracting an aspect term from an input sentence to determine its polarity as positive, neutral, or negative, and it is usually solved as a multi-class classification problem. Furthermore, sentiment analysis has been widely employed in multiple studies for both reviews and user opinions analysis in online commercial platforms [16], [47], as well as user behavior mining in social media platforms like Twitter [6, [29], [37], [42].

Furthermore, in the last decade, the continuous growth of social media platforms such as Twitter and Facebook, as well as the widespread dissemination of non-trusted information on them, has sparked applied research to automatically identify these non-trusted sources, which in many cases correspond to non-human or Bot accounts. [9] Proposed one of the first studies in this subject, which used a random forest technique to identify bots and non-bots accounts using a manually annotated dataset with roughly 2000 samples using a random forest approach to categories bots and non-bots accounts. BotOrNot [12] was introduced in 2016 as a tool to automatically detect bots in Twitter based on similarities between social bot features. This model has sparked further

research in the sector, with this service even being used to automatically annotate data from Twitter.

The authors of [11] annotated over 8000 accounts and developed a classifier that obtained a high degree of accuracy for such a large number of examples. [38] Also provided a technique for detecting Twitter bots using a huge quantity of metadata from the account to conduct the classification.

Several recent scientific studies, such as [27], [44], [45], have included additional annotated samples to support this research, including some strategies for improving accuracy by selectively picking a selection of training samples that better generalize the problem. In [22], potential features to distinguish between human and bot accounts are identified using a language-agnostic technique.

The model is then trained and verified using over 8000 imbalanced samples, and its accuracy exceeds 98 percent.

Furthermore, authors in [32] proposed a 2D Convolutional Network model based on user-generated materials for recognizing bots from real accounts, including gender (male, female account) and language (Spanish, English). Authors in [40] investigate a similar goal, using both Word and Character N-Grams as major features to perform categorization.

Authors in [5] suggested a different approach to the problem, in which novel altimetry's data used to examine social networks are analyzed and used to train a Graph Convolutional Network (GCN) that achieves above 70% accuracy in this task. Authors in [34], on the other hand, presented a revolutionary one-class classifier to improve Twitter bot detection without requiring any prior information about them. By the time their trials were completed, most of the aforementioned methodologies were limited due to a lack of significant amounts of annotated data for this specific purpose. This difficulty is also mentioned in [45], thus this research took into account all currently accessible public datasets in order to construct a system that uses the most up-to-date, newest, and relevant state-of-the-art annotated data from Twitter. Furthermore, while many systems leverage both metadata and text-based features from user accounts, the text-based features are either extracted at a lexical level or only cover a few languages, such as Spanish or English.

Unlike previous research, our proposed model uses novel multilingual Language Models (LM) to encode all text-based features of an input user account, including transformer models like the so-called BERT [14] or Contextual string

embedding's proposed in [2.] An input vector for the user account is obtained by concatenating both the metadata set of features and the output vector provided by these LMs. Finally, based on the aforementioned input vector, this study presents a Dense-based DL model to produce both the account's final decision and a low-dimensional embedding of the user.

III. A MULTILINGUAL APPROACH FOR USER ACCOUNT ENCODING VIA TRANSFORMERS

As previously stated, our technology is a multilingual approach capable of better identifying suspicious Twitter accounts based on a set of indicators that are independent of the account's language. More specifically, the methodology for developing the entire system can be divided into two stages: (i) a preprocessing stage in which a multilingual input vector for the user account is generated, and (ii) a final decision system in which patterns in the input vector generated during the first stage are used to determine whether the account has normal or abnormal behavior.

Furthermore, the former procedure is in charge of getting a large number of annotated Twitter accounts in a binary format, with the positive class indicating that the account is a Bot and the negative class indicating that it is a human account. Following that, numerous features from each Twitter account were collected in order to improve some relevant factors, including (i) level of activity, (ii) level of popularity, and (iii) profile information.

Finally, this first stage culminates in the creation of an input vector including both textual and metadata information for each Twitter account by merging all of the features. Section III-A explains the entire process in order to provide all of the implementation specifics.

The latter method, described in Section III-B, is in charge of using Deep Neural Networks to automatically find patterns in the input encoding vector in order to accurately discriminate between bots and human Twitter accounts (DNNs). Furthermore, because of its low-dimensional character, this technique automatically gets a low-dimensional feature representation of the input vector, which may be utilized for any IR purpose in a more efficient manner.

MULTILINGUAL USER ENCODING VIA TRANSFORMERS

As stated in Section III, a first step is required to aggregate various key elements from Twitter accounts in order to create

a robust multilingual encoding representation that may be used as potential inputs for classification purposes across DNNs.

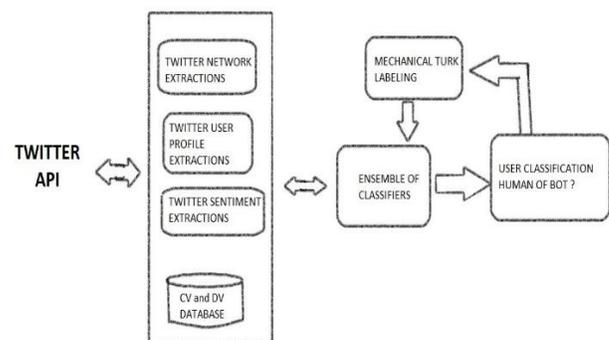
Figure 1 depicts an illustrated block diagram of this process, illustrating the various tools and stages required to meet the system's first phase's objectives.

More specifically, the Twitter API is used in the first stage to retrieve all of the data from the set of users U . Following that, each account is represented as a vector with two modalities in mind: text-based and metadata features. The former set is fed into a multilingual pre-trained LM model to generate a feature vector representation of the text information, which includes the account's description, username, and language. The final stage is a concatenation of both modalities into a single feature vector x , which encapsulates the information of an input user account.

1) DATASET GENERATION

Table 1 lists a number of public datasets that address the bot identification challenge from a binary classification standpoint. Furthermore, some of these datasets were previously used to train and assess [12]'s Botometer (formerly BotOrNot) service.

However, as described by the authors in [7], [28], bot account production is constantly changing over time, and some of the given accounts have already been suspended by Twitter. As a result, preparation is required to improve the usefulness of this enormous collection of statistics by eliminating IDs from accounts that have already been deactivated by Twitter. This is especially important because numerous earlier Bot detectors haven't been updated to reflect the new habits and features that bot accounts may have, making them.



SYSTEM ARCHITECTURE

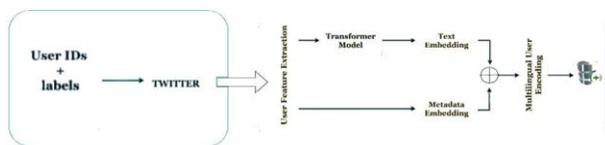
However, due to Twitter's policy restrictions, the datasets only contain the identification of the Twitter account and no other significant information. As a result, an extra data crawling procedure using the Twitter API is used to collect more information about the available accounts for the research. The following information is gathered, as stated in Section III:

- Popularity metrics include the total number of friends and followers, as well as the total number of likes and comments.
- Activity variables include: account age, average tweets per day, tweets & favorites counts, and account creation date.
- Screen name, description, language, location, verified indication, and default profile indicator are all features of the profile information.

The final complete dataset is made up of 37438 Twitter accounts after crawling and preprocessing the data. The remaining 25013 accounts were tagged as human accounts. Bots number 12425.

2) GENERATION OF INPUT USER ENCODING

The development of a user encoding vector based on the aforementioned collection of features to serve as input to the proposed deep learning model is a critical aspect of this first stage. The proposed solutions in prior relevant works [7], [12], [13], [20] had two fundamental constraints: 1) they used metadata-oriented approaches to extract text-based features at a semantic level in terms of Natural Language Processing, or 2) they used more complex NLP procedures based on N-grams or DL solutions, but only supported a restricted number of languages while completing the analysis.



Our method, on the other hand, addresses the aforementioned constraints by combining important metadata characteristics with sophisticated models capable of converting text-based features into vectors regardless of the input text's language.

3) TEXT-BASED FEATURES ENCODING

A specific function $g(u)$ is required for the production of the text-based vector from an input user account u_i . Various state-of-the-art sentence-level encoders from various NLP frameworks were researched and explored. The Flair

framework [2] was used to combine state-of-the-art Word Embedding's (WE) with Transformers [39], [43] for extracting strong document embedding's from text-based characteristics. In this work, the following main families of embedding's were used:

I Contextual string embedding's [3], which are taught without having any explicit concept of words, and so model words as sequences of characters. Furthermore, the context of words is provided by the surrounding text.

[i] Describes the JW300 Dataset, which was used to train the employed model. Multi-forward and multi-backward embedding's are used in this investigation. Their outputs have a 2048-dimensional dimension.

(ii) BERT embedding's, which were conceived and developed by [14] and are based on a bidirectional transformer architecture [39], [43]. The so-called Bert-base-multilingual-cased was used in this investigation.

(iii) RoBERTa, an adaptive variation of the BERT embedding whose purpose is to increase performance in longer sequences or when large amounts of data are present, as suggested by [41]. We used the so-called roberta-large-mnli pre-trained model in this situation.

4) METADATA-BASED FEATURES ENCODING

On the other hand, throughout function $h(u, I)$, all of the corresponding metadata properties from an input user account are properly preprocessed and encoded so that neural networks can read them. Furthermore, as noted in Section III-A, they are concatenated with the aforementioned text-based features.

This set of features, in particular, contains all information linked to the popularity and activity of the user account.

BOT-DenseNet

After obtaining the collection of input user vectors designated by U , a second method is necessary to determine whether an account is a bot or a human.

To achieve this purpose, we present Bot-DenseNet, a Deep Fully-connected based neural network capable of identifying robust decision boundaries artificial on hidden patterns in input vectors in order to better distinguish bot accounts on Twitter.

IV. EXPERIMENTAL RESULTS

One of the main goals of this paper is to analyse different input feature vectors based on Transformers, as well as other unique ways to assess the performance of the same DNN model, using an ablation study.

Furthermore, the DL architecture was trained using supervised learning and a binary classification, with the Positive class referring to Bots and the Negative class referring to Human accounts.

Due to the dataset's unbalanced constraint, two primary actions have been taken: (i) the F1-score metric as the main measurement of the system's performance because it balances both precision and recall metrics in a single value and provides more realistic information about the model's capability to detect both Positive and Negative classes than the classical accuracy metric [26], (ii) the F1-score metric as the main measurement of the system's performance because it balances both precision and recall metrics in a single value and provides more realistic information about the model's capability to detect both.

A. TRAINING & VALIDATING BOT-DENSENET

The goal of these experiments is to discover the best text embedding to layer on top of the Bot-DenseNet, together with the remaining metadata feature vector, in order to determine the best decision boundaries for downstream tasks like the one described in this paper.

Table 3 summarizes the results acquired throughout the training and validation steps for all potential input feature vectors.

Because the dataset is unbalanced, the F1-score metric plays a critical role in the system's evaluation in order to objectively measure the performance of identifying bots in a social network like Twitter, where only a small percentage of total accounts belong to the bot category, as described in previous studies [7], [27], and [34]. The F1-score was computed using both the recall and the precision at this specific epoch because the model is trained utilizing an Early-Stopping callback to terminate the process in the epoch when the loss function in the validation set is no longer lowering.

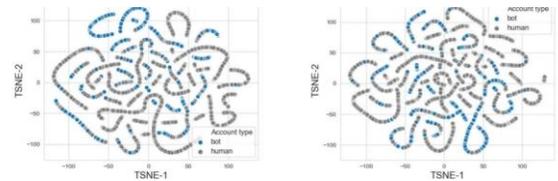
B. SELF-SUPERVISED USER EMBEDDING

Based on a diverse set of input features, our suggested model is capable of detecting fraudulent Twitter accounts.

However, as is widely recognized in the Deep Learning framework, intermediate layers are usually an adequate embedded representation of the inputs that, due to their low-dimensional character, can be used in other downstream tasks such as text categorization or similarity analysis more efficiently. As a result, after training our proposed model, it creates a self-supervised representation of an input Twitter account, because the intermediate hidden layer automatically learns such a representation.

C. DECISION MAKING CRITERIA FOR BOT-DENSENET

Following the completion of the aforementioned studies, a final judgement on the model's best configuration must be made. The following criteria were evaluated in order to do so: (i) The model's performance in terms of F1-score in both training and validation to provide objective criteria when making the final decision;



(ii) The model's simplicity in terms of both trainable parameters and the length of the input feature vector as noted in Table 4; (iii) The simplicity of the final decision boundaries to distinguish between Bot and Human accounts as assessed by observing the low-dimensional embedding's distribution in Table 4; (iv) The simplicity of the final decision boundaries to distinguish between Bot This is an important issue to consider in order to present a strong model that can be generalized in future applications.

V. CONCLUSION

It has been detailed a robust approach for detecting Bots in Twitter accounts. Transfer learning approaches have been used in this study to extract compact multilingual representations of text-based attributes linked with user accounts using sophisticated state-of-the-art NLP models such as Transformers. Several constraints linked to processing text-based features to improve the input feature vector from several languages were addressed as a result of this work.”

VI. FUTURE WORKS

Furthermore, a final classifier termed Bot-Dense Net was trained and validated using a huge collection of samples gathered via the Twitter API by combining text encodings with extra metadata on top of a dense-based neural network. To acquire a single vector indicating the text-based properties of the user account, numerous tests were undertaken utilising various combinations of Word Embeddings, document embedding (Pooling and LSTMs), and Transformers. Following that, a detailed comparison of the proposed classifier's performance when using these approaches of Language Models as part of the input has been presented in order to determine which input vector provides the best result in terms of performance simplicity and feasibility in the generation of decision boundaries.

REFERENCES

- [1]. "JW300: A wide-coverage parallel corpus for low resource languages," in Proc. 57th Annu. Meeting Assoc. Comput. Linguistics. Association for Computational Linguistics, Florence, Italy, July 2019, pp. 3204–3210.
- [2]. A. Akbik, T. Bergmann, D. Blythe, K. Rasul, S. Schweter, and A. Vollgraf, "FLAIR: An easy-to-use framework for state-of-the-art NLP," in Proc. North Amer. Chapter Assoc. Comput. Linguistics, Demonstrations, 2019, pp. 54–59.
- [3]. "Contextual string embeddings for sequence tagging," A. Akbik, D. Blythe, and R. Vollgraf, in Proc. 27th Int. Conf. Comput. Linguistics, 2018, pp. 1638–1649.
- [4]. Twitter sentiment analysis with a deep neural network: An upgraded technique leveraging user behavioural information," Cogn. Syst. Res., vol. 54, pp. 50–61, May 2019.
- [5]. "Bot prediction on social networks of Twitter in altmetrics using deep graph convolutional networks," Soft Comput., vol. 24, pp. 11109–11120, Jan. 2020.
- [6]. "Bot prediction on social networks of Twitter in altmetrics using deep graph convolutional networks," Soft Comput., vol. 24, pp. 11109–11120, Jan. 2020.
- [7]. "Character level embedding with deep convo lyminal neural network for text normalisation of unstructured data for Twitter sentiment analysis," M. Arora and V. Kansal, Social Netw. Anal. Mining, vol. 9, no. 1, p. 12, Dec. 2019.
- [8]. "Identification of credulous users on Twitter," in Proc. 34th ACM/SIGAPP Symp. Appl. Comput., A. Balestrucci, R. De Nicola, O. Inverso, and C. Trubiani.
- [9]. "Various techniques to aspect-based senti ment analysis," A. Bhoi and S. Joshi, arXiv:1805.01984, 2018. [Online]. The abstract is available at <http://arxiv.org/abs/1805.01984>.
- [10]. "Detecting automation of Twitter accounts: Are you a human, bot, or cyborg?" Z. Chu, S. Gianvecchio, H. Wang, and S. Jajodia, "Detecting automation of Twitter accounts: Are you a human, bot, or cyborg?" Nov./Dec. 2012, IEEE Trans. Depend. Sec. Comput., vol. 9, no. 6, pp. 811–824.
- [11]. "Recurrent batch normalisation," T. Cooijmans, N. Ballas, C. Laurent, Gülçehre, and A. Courville, arXiv:1603.09025, 2016. [Online]. The abstract is available at <http://arxiv.org/abs/1603.09025>.
- [12]. S. Cresci, R. Di Pietro, M. Petrocchi, A. Spognardi, and M. Tesconi, "The paradigm-shift of social spambots: Evidence, theories, and tools for the arms race," in Proc. 26th International Conference on World Wide Web Companion, 2017, pp. 963–972.
- [13]. "BotOrNot: A system to analyse social bots," C. A. Davis, O. Varol, E. Ferrara, A. Flammini, and F. Menczer, in Proc. 25th Int. Conf. Companion World Wide, 2016, pp. 273–274.
- [14]. A. Davoudi, A. Z. Klein, A. Sarker, and G. Gonzalez-Hernandez, AMIA Summits Transl. Sci. Proc., vol. 2020, p. 136, May 2020.
- [15]. "BERT: Pre-training of deep bidirectional transformers for language interpretation," [14] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, arXiv:1810.04805, 2018. [Online]. The abstract is available at <http://arxiv.org/abs/1810.04805>.
- [16]. J. Diesner, E. Ferrari, and G. Xu, in Proc. IEEE/ACM Int. Conf. Adv. Social Network Analysis and Mining, Sydney, NSW, Australia: ACM, August 2017. [Online], doi: 10.1145/3110025, <https://dblp.org/rec/bib/conf/asunam/2017>.
- [17]. C. D. Santos and M. Gatti, "Deep convolutional neural networks for sentiment analysis of short texts," in Proceedings of the 25th International Conference on Computer Linguistics (COLING), 2014, pp. 69–78.
- [18]. L. Floridi and M. Chiriatti, "GPT-3: Its Nature, Scope, Limits, and Consequences," in "GPT-3: Its Nature, Scope, Limits, and Consequences," in "GPT-3: Its Nature, Scope, Limits, and Nov. 2020, Minds Mach., vol. 30, pp. 681–694.
- [19]. "Maximum mean discrepancy is aware of adversarial attacks," R. Gao, F. Liu, J. Zhang, B. Han, T. Liu, G. Niu, and M. Sugiyama, arXiv:2010.11415, 2020. [Online]. The abstract is available at <http://arxiv.org/abs/2010.11415>.
- [20]. "Outer product based neural collaborative filtering," X. He, X. Du, X. Wang, F. Tian, J. Tang, and T.-S. Chua, arXiv:1808.03912, 2018. [Available online at: <http://arxiv.org/abs/1808.03912>]
- [21]. J. Im, E. Chandrasekharan, J. Sargent, P. Lighthammer, T. Denby, A. Bhargava, L. Hemphill, D. Jurgens, and E. Gilbert, "Still out there: Modeling and recognising Russian troll accounts on Twitter," in Proc. 12th ACM Conf. Web Sci., July 2020, pp. 1–10.
- [22]. "Batch normalisation: Accelerating deep network training by decreasing internal covariate shift," S. Ioffe and C. Szegedy,

- [45]. "Arming the public with artificial intelligence to defeat social bots," K. Yang, O. Varol, C. A. Davis, E. Ferrara, A. Flammini, and F. Menczer," Jan. 2019, *Hum. Behav. Emerg. Technol.*, vol. 1, no. 1, pp. 48–61.
- [46]. "Scalable and generalizable social bot detection through data selection," in *Proc. AAAI Conf. Artif. Intell.*, vol. 34, 2020, pp. 1096–1103. [45] K.-C. Yang, O. Varol, P.-M. Hui, and F. Menczer, "Scalable and generalizable social bot detection through data selection," in *Proc. AAAI Conf. Artif. Intell.*
- [47]. "Understanding deep learning needs rethinking generalisation," say C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals. 1611.03530, arXiv:1611.03530. [Online]. The abstract is available at <http://arxiv.org/abs/1611.03530>.
- [48]. "Multi-layer attention based CNN for target-dependent sentiment categorization," S. Zhang, X. Xu, Y. Pang, and J. Han," pp. 2089–2103, *Neural Process. Lett.*, vol. 51, no. 3, June 2020.
- [49]. "Context-based prediction for road traffic state using trajectory pattern mining and recurrent convolutional neural networks," J. Zhu, C. Huang, M. Yang, and G. P. Cheong Fung," Jan. 2019, *Inf. Sci.*, vol. 473, pp. 190–201.

BIOGRAPHIES

ELANCHELIYAN R is an undergraduate student, department of computer science and engineering in paavai college of engineering.

ARAVINDHAN R is an undergraduate student, department of computer science and engineering in paavai college of engineering.

RAHUL G is an undergraduate student, department of computer science and engineering in paavai college of engineering.

Mrs. BRINDA B.M is an Assistant professor, department of computer science and engineering in paavai college of engineering.