

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

Mr. K. Pazhanivel¹, Ajai Kumar. B², Mageshwaran. M³, Dhivakar. K⁴

Assistant Professor, Department of Science and Computer Engineering¹

Students, Department of Science and Computer Engineering^{2,3,4}

Anjalai Ammal Mahalingam Engineering College, Kovilvenni, Tiruvarur, Tamil Nadu, India

Abstract: *The volume of audio visual 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. The disruption of false news sources and bot accounts for disseminating fake news is a possibility in this scenario. Applied research has been supported by promotional information as well as sensitive stuff over the network. Artificial Intelligence will be used to automatically assess the trustworthiness of social media accounts (AI). In this research, we describe a multilingual strategy to using Deep Learning to solve the bot identification problem on Twitter. End-users can utilise machine learning (ML) methodologies to assess the trustworthiness of a Twitter account. To achieve so, a number of tests were carried out using cutting-edge Multilingual Language Models. Construct an encoding of the user account's text-based features, which is then concatenated with the rest of the metadata to build a potential input vector on top of a Bot-DenseNet Dense Network. As a result, this article evaluates the language constraint from prior experiments where the encoding of the language was limited. Only the metadata information or the metadata information along with some other information was examined by the user account. properties of fundamental semantic text The Bot-DenseNet also generates a low-dimensional representation of the data. Within the Information Retrieval (IR) framework, a user account can be utilised for any application.*

Keywords: Bot Detector, Deep Learning, Feature Representation, Language Models, Misinformation Detection, Social Media Mining, Transfer Learning, Transformers

I. 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., 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.

However, the rapid rise of social media platforms has fueled a desire to sway people's opinions on specific issues by disseminating propaganda or skewed facts. Bots, which have been widely reported in various studies [31], [32], [40], are automatic systems capable of generating and spreading multimedia information throughout the network without the supervision of a human being. 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 for consumers on such platforms by decreasing privacy threats while also enhancing platform trustworthiness.

As a result, this article intends to advance the state-of-the-art in this field by presenting a novel method for automatically (i) encoding an input user account as a low-dimensional feature vector regardless of its language, and (ii) encoding an input user account as a low-dimensional feature vector.

(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).

II. LITERATURE REVIEW

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 behaviour 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 categorise bots and non-bots accounts. In the year 2016, BotOrNot was proposed in [12] as a tool for automatically detecting bots on 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. Furthermore, [38] presented a technique for detecting Twitter bots using a huge quantity of metadata from the account to conduct the categorization. Several scientific studies, such as [27], [44], [45], have lately included additional annotated samples to support this research, including certain strategies for reaching a higher level of accuracy by selectively picking a selection of training samples that better generalise the problem. [22] uses a language-independent approach to uncover potential traits that can be used to distinguish between human and bot accounts. After that, the model is trained and verified using over 8000 data distributed in an imbalanced manner, and its accuracy approaches 98 %

Furthermore, writers in [32] suggested a 2D Convolutional Neural Network. Detecting bots from human accounts using a network model based on user-generated material, including gender (male, female).female account) that is bilingual (Spanish and English). The authors of [40] investigate a similar purpose. As major N-Grams, both Word and Character N-Grams are used. To accomplish the classification, you'll need certain features. Recently, a new approach to the problem was proposed. new altmetrics data to be proposed by writers in [5], Investigate social networks, which are analysed and put to use. construct a Graph Convolutional Network (GCN) that achieves. This challenge required above 70% accuracy. Authors, on the other hand, [34] proposed a novel one-class classifier to improve Twitter bot identification without requiring any prior knowledge. concerning 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.

III. TRANSFORMERS MODEL FOR USER ACCOUNT ENCODING

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 behaviour. 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 usage, (ii) level of reach, and (iii) profile data.

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 utilised for any IR purpose in a more efficient manner.

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.

The many tools and stages required to meet the objectives of this initial phase of the system are depicted in an exemplary block diagram. 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 through a multilingual pre-trained LM model to provide 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.

IV. PROPOSED SYSTEM

We describe a multilingual method to addressing 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.

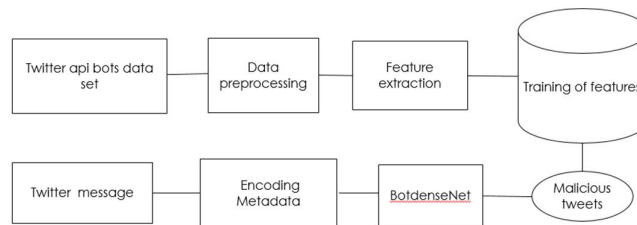


Figure 1: System Architecture of Proposed model

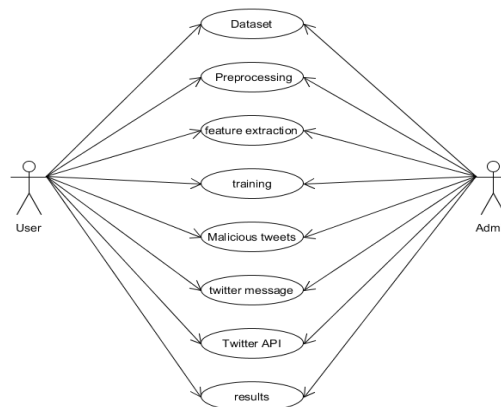


Figure 2: Use case diagram of proposed model

V. MODULE DESCRIPTION

Proposed System Modules

- Dataset Collection
- Pre-processing
- Feature Extraction
- Training

5.1 Dataset Collection

There are several public datasets that address the bot identification problem from a binary classification perspective. Furthermore, some of these datasets were previously used to train and evaluate the [12] Botometer (formerly BotOrNot) service. However, as noted by the authors in [7], [28], the production of bot accounts fluctuates over time, and some of the given accounts have already been suspended by Twitter. 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 & favourites 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 consists of 37438 Twitter accounts, 25013 of which were identified as human after crawling and preprocessing the data. accounts and the remaining 12425 are bots.

5.2 Data Pre-processing

The development of a user encoding vector based on the aforementioned collection is a critical aspect of this first stage. a set of features that will be used as input into the proposed deep learning model. 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.

A specific user account is represented as $u_i = [u_{t I} \ u_{z I}]$ when given an input set of Users $U = u_1, u_2, \dots, u_m$. $I = 1, \dots, m$, where $u_{t I}$ denotes the text-based feature vector and $u_{z I}$ is the metadata-based feature vector. Our proposed method uses a mapping function $f(u)$ to create a new set of Users $U = u_1, u_2, \dots, u_m$, where $u_I = f(u_i) = g(u_{t I}) \parallel h(u_{z I})$.

We use the symbol \parallel to express a concatenation operation between these two vectors in this situation. The system only considers information originating from the same target object, which is why a concatenation layer is used at the end of this process:

the account of the user. Other widely used Collaborative Recommender System alternatives, such as computing the outer product [19], were not considered for this approach because, in those scenarios, the information from the embeddings comes from two different sources: Users and Items, and the goal of the outer product is to catch similarities and discrepancies between these two sets.

5.3 Feature Extraction

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 Embeddings (WE) with Transformers [39], [43] for extracting strong document embeddings from text-based characteristics. In this work, the following main families of embeddings

were used. I Contextual string embeddings [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.

[1] describes the JW300 Dataset, which was used to train the employed model. Multi-forward and multi-backward embeddings are used in this investigation. Their outputs have a 2048-dimensional dimension.

[14] conceived and developed BERT (Bidirectional Encoder Representations from Transformers) embeddings, which are based on a bidirectional encoder. [39], [43] are examples of transformer architecture. 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

Furthermore, several experiments involving three different solutions were carried out. The first method is based on employing one or more stacked embeddings, similar to the method provided by [25], to encode all text-based information from the user account at a sentence-level representation. Following that, a document-level representation is calculated using a Pooling model, which takes the average of all the stacked sentence-level embeddings.

The second method involves training a recurrent network with a Long Short-Term Memory (LSTM) across all of the word embeddings required to construct the sentence-level encoding. Finally, the document-level embedding is produced directly using an intermediary layer from a pre-trained Transformer model. Additionally, the multilingual BERT transformer pre-trained model BERT-base-multilingual-cased as well as the RoBERTa pre-trained model roberta-largemnli were used. The former provides a 768-dimensional embedded vector, while the latter produces a 1024-dimensional embedded representation. The official Hugging Face repository has all of the details on the aforementioned pre-trained models.

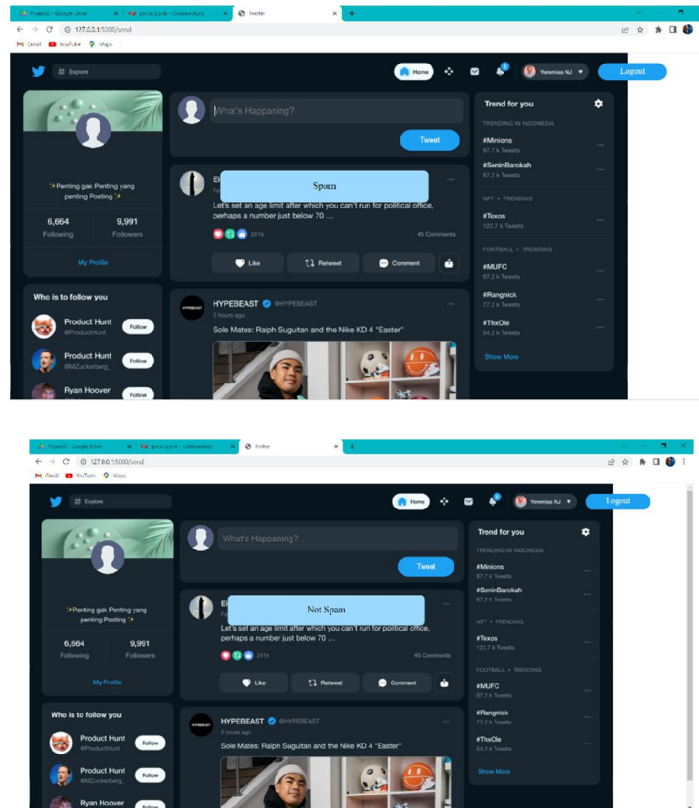
5.4 Training

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.

Summarises 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 success of recognising bots in a social network like Twitter, where only a small percentage of total accounts correspond to the bot category as defined in prior research [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 utilising an Early-Stopping callback to terminate the process in the epoch when the loss function in the validation set is no longer lowering. In addition, each row specifies the pre-trained LM used (Flair, BERT, RoBERTa, and so on), as well as the method used to construct the final input user encoding (through Pooling, bidirectional LSTM, or directly the embedding derived from an intermediate layer of the Transformer model). Table 3 shows that when the text embeddings are combined directly, the outcomes are as follows:

The proposed model Bot-DenseNet achieves higher performance in terms of F1-score in both the training and validation sets, thanks to metadata characteristics collected from intermediary levels of the Transformers. When employing Pooling or LSTMs to generate the final text embeddings, on the other hand, the F1-score metric in the training phase is greater than in the validation phase, indicating that the model is overfitting. As a result, the system's performance when analysing unseen observations in future forecasts may suffer as a result of this problem.

VI. OUTPUT BY THE PROPOSED SYSTEM



VII. RESULT

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. metric as the main measurement of the system's Due to the dataset's unbalanced constraint, two primary actions have been taken: I the F1-score metric as the major measurement of the system's performance because it balances both classes, as revised in [26]; (ii) the F1-score performance because it balances both classes during the training and evaluation of the system.

The precision and recall metrics are combined into a single value, giving more accurate information about the model's ability to recognise both Positive and Negative classes than the traditional accuracy statistic.

VIII. CONCLUSION

A robust approach for detecting Bots in Twitter accounts is discussed in this paper. Transfer learning approaches have been used in this study to extract compact multilingual representations of text-based attributes linked with user accounts using powerful state-of-the-art NLP models such as Transformers. Several restrictions linked to processing text-based features to improve the input feature vector from multiple languages were minimised as a result of this. Furthermore, a final classifier termed Bot-DenseNet was trained and verified 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, multiple tests were undertaken utilising various combinations of Word Embeddings, document embeddings (Pooling and LSTMs), and Transformers. Following that, a detailed comparison of the proposed classifier's performance when using these approaches of Language Models as input has been presented in order to determine which input vector provides the best result in terms of performance simplicity in the generation of decision boundaries and feasibility.

REFERENCES

- [1]. Ž. Agić and I. Vulić, “JW300: A wide-coverage parallel corpus for lowresource languages,” in Proc. 57th Annu. Meeting Assoc. Comput. Linguistics. Florence, Italy: Association for Computational Linguistics, Jul. 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. Conf. North Amer. Chapter Assoc. Comput. Linguistics, Demonstrations, 2019, pp. 54–59.
- [3]. A. Akbik, D. Blythe, and R. Vollgraf, “Contextual string embeddings for sequence labeling,” in Proc. 27th Int. Conf. Comput. Linguistics, 2018, pp. 1638–1649.
- [4]. A. S. M. Alharbi and E. de Doncker, “Twitter sentiment analysis with a deep neural network: An enhanced approach using user behavioral information,” *Cogn. Syst. Res.*, vol. 54, pp. 50–61, May 2019.
- [5]. N. R. Aljohani, A. Fayoumi, and S.-U. Hassan, “Bot prediction on social networks of Twitter in altmetrics using deep graph convolutional networks,” *Soft Comput.*, vol. 24, pp. 11109–11120, Jan. 2020.
- [6]. M. Arora and V. Kansal, “Character level embedding with deep convolutional neural network for text normalization of unstructured data for Twitter sentiment analysis,” *Social Netw. Anal. Mining*, vol. 9, no. 1, p. 12, Dec. 2019.
- [7]. A. Balestrucci, R. De Nicola, O. Inverso, and C. Trubiani, “Identification of credulous users on Twitter,” in Proc. 34th ACM/SIGAPP Symp. Appl. Comput., Apr. 2019, pp. 2096–2103.
- [8]. A. Bhoi and S. Joshi, “Various approaches to aspect-based sentiment analysis,” 2018, arXiv:1805.01984. [Online]. Available: <http://arxiv.org/abs/1805.01984>
- [9]. Z. Chu, S. Gianvecchio, H. Wang, and S. Jajodia, “Detecting automation of Twitter accounts: Are you a human, bot, or cyborg?” *IEEE Trans. Depend. Sec. Comput.*, vol. 9, no. 6, pp. 811–824, Nov./Dec. 2012.
- [10]. T. Cooijmans, N. Ballas, C. Laurent, Ç. Gülçehre, and A. Courville, “Recurrent batch normalization,” 2016, arXiv:1603.09025. [Online]. Available: <http://arxiv.org/abs/1603.09025>
- [11]. 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 Int. Conf. World Wide Web Companion, 2017, pp. 963–972.
- [12]. C. A. Davis, O. Varol, E. Ferrara, A. Flammini, and F. Menczer, “BotOrNot: A system to evaluate social bots,” in Proc. 25th Int. Conf. Companion World Wide, 2016, pp. 273–274.
- [13]. A. Davoudi, A. Z. Klein, A. Sarker, and G. Gonzalez-Hernandez, “Towards automatic bot detection in twitter for health-related tasks,” *AMIA Summits Transl. Sci. Proc.*, vol. 2020, p. 136, May 2020.
- [14]. J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “BERT: Pre-training of deep bidirectional transformers for language understanding,” 2018, arXiv:1810.04805. [Online]. Available: <http://arxiv.org/abs/1810.04805>
- [15]. J. Diesner, E. Ferrari, and G. Xu, in Proc. IEEE/ACM Int. Conf. Adv. Social Netw. Anal. Mining. Sydney, NSW, Australia: ACM, Aug. 2017. [Online]. Available: <https://dblp.org/rec/bib/conf/asunam/2017>, doi: 10.1145/3110025.
- [16]. C. D. Santos and M. Gatti, “Deep convolutional neural networks for sentiment analysis of short texts,” in Proc. 25th Int. Conf. Comput. Linguistics (COLING), 2014, pp. 69–78.
- [17]. L. Floridi and M. Chiriatti, “GPT-3: Its nature, scope, limits, and consequences,” *Minds Mach.*, vol. 30, pp. 681–694, Nov. 2020.
- [18]. R. Gao, F. Liu, J. Zhang, B. Han, T. Liu, G. Niu, and M. Sugiyama, “Maximum mean discrepancy is aware of adversarial attacks,” 2020, arXiv:2010.11415. [Online]. Available: <http://arxiv.org/abs/2010.11415>
- [19]. X. He, X. Du, X. Wang, F. Tian, J. Tang, and T.-S. Chua, “Outer productbased neural collaborative filtering,” 2018, arXiv:1808.03912. [Online]. Available: <http://arxiv.org/abs/1808.03912>
- [20]. J. Im, E. Chandrasekharan, J. Sargent, P. Lighthammer, T. Denby, A. Bhargava, L. Hemphill, D. Jurgens, and E. Gilbert, “Still out there: Modeling and identifying Russian troll accounts on Twitter,” in Proc. 12th ACM Conf. Web Sci., Jul. 2020, pp. 1–10.
- [21]. S. Ioffe and C. Szegedy, “Batch normalization: Accelerating deep network training by reducing internal covariate shift,” 2015, arXiv:1502.03167. [Online]. Available: <http://arxiv.org/abs/1502.03167>

- [22]. J. Knauth, "Language-agnostic Twitter-bot detection," in Proc. Int. Conf. Recent Adv. Natural Lang. Process. (RANLP), 2019, pp. 550–558.
- [23]. Z. Lin, S. Mu, F. Huang, K. A. Mateen, M. Wang, W. Gao, and J. Jia, "A unified matrix-based convolutional neural network for finegrained image classification of wheat leaf diseases," IEEE Access, vol. 7, pp. 11570–11590, 2019.
- [24]. F. Liu, W. Xu, J. Lu, G. Zhang, A. Gretton, and D. J. Sutherland, "Learning deep Kernels for non-parametric two-sample tests," 2020, arXiv:2002.09116. [Online]. Available: <http://arxiv.org/abs/2002.09116>
- [25]. Y. Liu, P. Dmitriev, Y. Huang, A. Brooks, and L. Dong, "An evaluation of transfer learning for classifying sales engagement emails at large scale," 2019, arXiv:1905.01971. [Online]. Available: <http://arxiv.org/abs/1905.01971>
- [26]. P. Lynn, "The advantage and disadvantage of implicitly stratified sampling," Methods, Data, Analyses, vol. 13, no. 2, p. 14, 2019.
- [27]. M. Mazza, S. Cresci, M. Avvenuti, W. Quattrociocchi, and M. Tesconi, "RTbust: Exploiting temporal patterns for botnet detection on Twitter," in Proc. 10th ACM Conf. Web Sci., 2019, pp. 183–192.
- [28]. A. Minnich, N. Chavoshi, D. Koutra, and A. Mueen, "BotWalk: Efficient adaptive exploration of Twitter bot networks," in Proc. IEEE/ACM Int. Conf. Adv. Social Netw. Anal. Mining, Jul. 2017, pp. 467–474.
- [29]. U. Naseem, I. Razzak, K. Musial, and M. Imran, "Transformer based deep intelligent contextual embedding for Twitter sentiment analysis," Future Gener. Comput. Syst., vol. 113, pp. 58–69, Dec. 2020.
- [30]. M. Orliński and N. Jankowski, "Fast t-SNE algorithm with forest of balanced LSH trees and hybrid computation of repulsive forces," Knowl.- Based Syst., vol. 206, Oct. 2020, Art. no. 106318.
- [31]. J. Pizarro, "Using N-grams to detect bots on Twitter," in Proc. CLEF, Working Notes, 2019, pp. 1–10.
- [32]. M. Polignano, M. G. de Pinto, P. Lops, and G. Semeraro, "Identification of bot accounts in Twitter using 2D CNNs on user-generated contents," in Proc. CLEF, Working Notes, 2019, pp. 1–11.
- [33]. C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, "Exploring the limits of transfer learning with a unified text-to-text transformer," J. Mach. Learn. Res., vol. 21, no. 140, pp. 1–67, 2020.
- [34]. J. Rodríguez-Ruiz, J. I. Mata-Sánchez, R. Monroy, O. Loyola-González, and A. López-Cuevas, "A one-class classification approach for bot detection on Twitter," Comput. Secur., vol. 91, Apr. 2020, Art. no. 101715.
- [35]. K. Shuang, H. Guo, Z. Zhang, J. Loo, and S. Su, "A word-building method based on neural network for text classification," J. Exp. Theor. Artif. Intell., vol. 31, no. 3, pp. 455–474, May 2019.
- [36]. N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A simple way to prevent neural networks from overfitting," J. Mach. Learn. Res., vol. 15, no. 1, pp. 1929–1958, 2014.
- [37]. D. Stojanovski, G. Strezoski, G. Madjarov, and I. Dimitrovski, "Twitter sentiment analysis using deep convolutional neural network," in Proc. Int. Conf. Hybrid Artif. Intell. Syst. Springer, 2015, pp. 726–737. [Online]. Available: <https://scholar.googleusercontent.com/scholar.bib?q=info:HnIU7VyTzLUJ:scholar.google.com/&output=citation&scisdr=CgXc4k0kELTt-pJoVIM:AAGBfm0AAAAAYGxtTIO0qf0SoEojztYZqYNU1uzAmqAp&sci sig=AAGBfm0AAAAAYGxtTlfX vsJzQ3eCFjVQwVDi0pipTQma&scisf=4&ct=citation&cd=-1&hl=es>
- [38]. O. Varol, E. Ferrara, C. A. Davis, F. Menczer, and A. Flammini, "Online human-bot interactions: Detection, estimation, and characterization," 2017, arXiv:1703.03107. [Online]. Available: <http://arxiv.org/abs/1703.03107>
- [39]. A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in Proc. Adv. Neural Inf. Process. Syst., 2017, pp. 5998–6008.
- [40]. I. Vogel and P. Jiang, "Bot and gender identification in Twitter using word and character N-grams," in Proc. CLEF, Working Notes, 2019, pp. 1–9.
- [41]. B. Wang and C.-C. J. Kuo, "SBERT-WK: A sentence embedding method by dissecting bert-based word models," 2020, arXiv:2002.06652. [Online]. Available: <http://arxiv.org/abs/2002.06652>
- [42]. L. Wang, J. Niu, and S. Yu, "SentiDiff: Combining textual information and sentiment diffusion patterns for Twitter sentiment analysis," IEEE Trans. Knowl. Data Eng., vol. 32, no. 10, pp. 2026–2039, Oct. 2020.

- [43]. T. Wolf et al., “HuggingFace’s transformers: State-of-the-art natural language processing,” 2019, arXiv:1910.03771. [Online]. Available: <http://arxiv.org/abs/1910.03771>
- [44]. K. Yang, O. Varol, C. A. Davis, E. Ferrara, A. Flammini, and F. Menczer, “Arming the public with artificial intelligence to counter social bots,” *Hum. Behav. Emerg. Technol.*, vol. 1, no. 1, pp. 48–61, Jan. 2019.
- [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.*, vol. 34, 2020, pp. 1096–1103.
- [46]. C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals, “Understanding deep learning requires rethinking generalization,” 2016, arXiv:1611.03530. [Online]. Available: <http://arxiv.org/abs/1611.03530>
- [47]. S. Zhang, X. Xu, Y. Pang, and J. Han, “Multi-layer attention based CNN for target-dependent sentiment classification,” *Neural Process. Lett.*, vol. 51, no. 3, pp. 2089–2103, Jun. 2020.
- [48]. J. Zhu, C. Huang, M. Yang, and G. P. Cheong Fung, “Context-based prediction for road traffic state using trajectory pattern mining and recurrent convolutional neural networks,” *Inf. Sci.*, vol. 473, pp. 190–201, Jan. 2019