

Robust Algorithm for Identifying Bot Authors on Social Media

Bitao Jin

<https://github.com/Bitao2/CSCI2952Q-Final-Project.git>

Abstract This study investigates the effectiveness of machine learning methods in identifying bot accounts on social media platforms, as well as the robustness of these methods when dealing with various common disguise techniques. The research utilized a publicly available Twitter dataset to compare the performance of three different machine learning models. In addition, the study attempted to compare the effectiveness of bot account detection based on tweet-level information versus user-level information. The experimental results indicate that incorporating user behavior information significantly enhances the accuracy of bot account identification, and this approach also demonstrates very high robustness against different disguise techniques.

Research Background

In recent years, the rapid development of social media has led to a large amount of user-generated content, and social media has become one of the main channels for the public to access information. This phenomenon has also brought about some potential issues. For instance, many social media platforms are inundated with a plethora of automated and semi-automated bot accounts, which can lead to a series of negative impacts, including the dissemination of false or harmful messages and the manipulation of public opinion. Therefore, accurately identifying these bot accounts plays a crucial role in maintaining the public environment of social media platforms. (Aljabri et al. 2023)

In real-world scenarios, researchers have attempted various technical methods to construct detection systems to identify these bot accounts. At the same time, the creators of these accounts continuously adopt various disguise techniques to make their accounts appear as normal users. This ongoing arms race between detection and disguise technologies has driven rapid development in this field and has increased the robustness requirements for detection systems.

This project focuses on the Twitter dataset, examining the performance of various detection technologies on this dataset and their ability to respond to common disguise techniques.

Data Description

The current project utilizes the dataset provided by <https://github.com/rrsr28/Twitter-Bot-Detection/tree/main/Datasets>). The dataset was gathered based on the research conducted by Alarfaj et al. 2023. It comprises a total of 155,758 tweets from 92 Twitter accounts collected over a period from 2011 to 2018. Out of these, there are 47 bot accounts and 45 human accounts, with an average of 1,673 tweets per account.

A comparison of the number and posting frequency of bot and human accounts is illustrated in Figure 1. It is evident that the tweeting frequency is consistent between the two groups, indicating that, solely based on the volume of tweets, these bot accounts do not exhibit any significant irregularities.

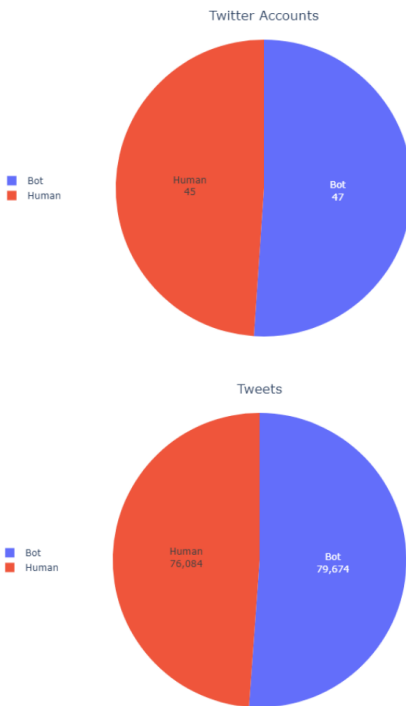


FIGURE 1. *The Accounts Distribution*

The distribution of tweet counts from two groups of users over different time periods is shown in Figure 2. It is evident that prior to 2018, both groups had relatively low posting activity. This is likely due to the data collection period being closer to 2018, making it more challenging to gather complete posting data from earlier years, resulting in certain data gaps.

From the data collected in 2018, it can be observed that the number of tweets sent

by Bot users fluctuates less over time compared to human users. This aligns with the nature of Bot behavior, which tends to post automatically.

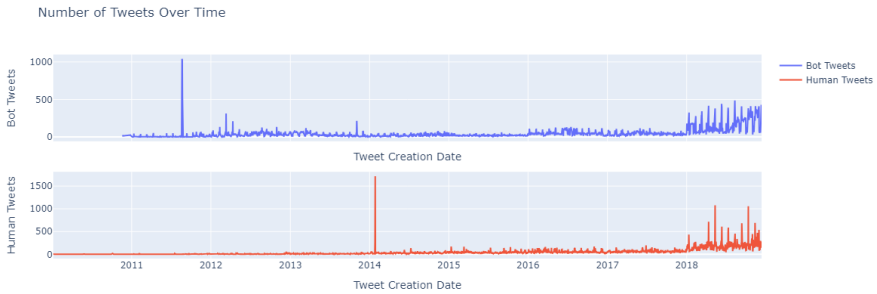


FIGURE 2. *The Tweets Over Time*

The semantic sentiment analysis of tweets from two groups of users is illustrated in Figure 3. From the perspective of tweet sentiment, the most significant difference between the two user groups is that Bot users tend to post a much higher proportion of neutral tweets compared to human users. This indicates a substantial difference in the content shared by these two groups. Therefore, a suitable strategy for identifying Bot users on Twitter is to focus on analyzing the content of their tweets.



FIGURE 3. *The Tweets Sentiment*

Research Methodology

Based on the analysis of the characteristics of the experimental data, this project aims to study the effectiveness of text classification models in identifying Bot users within the Twitter environment. In this task, there are two classification categories: determining whether a tweet is generated by a Bot (label 1) or a human (label 0).

After constructing the text classification model, this experiment will explore two types of Bot identification processes. The first process focuses on identifying Bots at the tweet level, while the second process examines user-level Bot identification. Below, we will sequentially outline the specific methods and procedures employed in

this experiment.

Introduction to Text Classification Tasks

The text classification task is one of the common techniques in natural language processing, widely applied in various scenarios involving text processing, such as text topic analysis, sentiment analysis of product reviews, and user intent analysis in intelligent customer service systems.

A standard text classification task typically includes the following steps:

1. **Text Data Preprocessing:** Cleaning the text data to be classified, such as removing invalid characters, filtering out meaningless text, and eliminating privacy information.
2. **Feature Extraction:** Converting the text into a vector format that machine learning models can process. Common feature extraction methods include Bag of Words, TF-IDF, etc.
3. **Model Training:** Choosing a suitable machine learning model for text classification and setting appropriate parameters for training, such as Naive Bayes, Support Vector Machines, Logistic Regression, Random Forest, etc.
4. **Model Evaluation:** Using appropriate evaluation metrics and a testing dataset to assess the classification performance of the model. Common evaluation metrics for classification tasks include accuracy, precision, recall, and F1-score, etc.

Introduction to TF-IDF Features

In this task, TF-IDF is used as the method for feature extraction. TF-IDF features are widely applied in information retrieval and text mining tasks. The calculation formula is as follows:

$$\text{TF-IDF} = \text{TF} \times \text{IDF} \quad (1)$$

In this formula, TF represents the frequency of a word occurring in a document; IDF indicates the importance of the word across the entire dataset, which is calculated as the total number of documents in the dataset divided by the number of documents containing that word. Generally speaking, words that appear in all documents are usually considered less important.

Introduction to Classification Models

In this experiment, three classification models are employed to determine whether a tweet was sent by a bot user. The models used are Naive Bayes (Bayes 1968), Logistic Regression (LaValley 2008), and Random Forest (Breiman 2001). Below is an overview of each model and its underlying principles.

Naive Bayes Model

Naive Bayes is a classification method based on Bayes' theorem. It assumes that the influence of each feature on the result is independent of the others. While this assumption does not always hold in real-world scenarios, it simplifies the calculations, making them efficient. Naive Bayes is especially suited for text classification tasks, such as spam detection, as it can quickly train and predict on large datasets.

When using Naive Bayes for classification, we compare the posterior probabilities of different categories. The calculation formula is as follows:

$$P(C_k | X) \propto P(C_k) \cdot \prod_{i=1}^n P(x_i | C_k) \quad (2)$$

By identifying the category C_k that maximizes the above expression, the classification task is accomplished.

Logistic Regression Model

Logistic regression is a linear model used for binary classification. Despite its name containing the word "regression," it is primarily used for classification purposes. It maps a linear combination of input feature values to a probability using an S-shaped (sigmoid) function, allowing for effective classification. Logistic regression is straightforward and effectively handles the relationship between two classes, commonly applied in prediction tasks across finance, healthcare, and other fields.

Random Forest Model

Random Forest is an ensemble model composed of multiple decision trees. It constructs many decision trees, each trained on a subset of the data, and then lets these trees "vote" to decide the final classification. This approach enhances prediction accuracy and helps control overfitting. Random Forest is a powerful model suitable for both classification and regression tasks, and due to its flexibility and robustness, it is widely used in various data analysis applications.

Bot Identification Process Based on Tweets

In this experiment, we explore the first approach for identifying Bot accounts, which is based on the characteristics of individual tweets. Each tweet is treated as a standalone document. If we determine that a tweet was sent by a Bot, it leads us to conclude that the account responsible for sending that tweet is indeed a Bot account. The detailed processing workflow is illustrated in Figure (4).

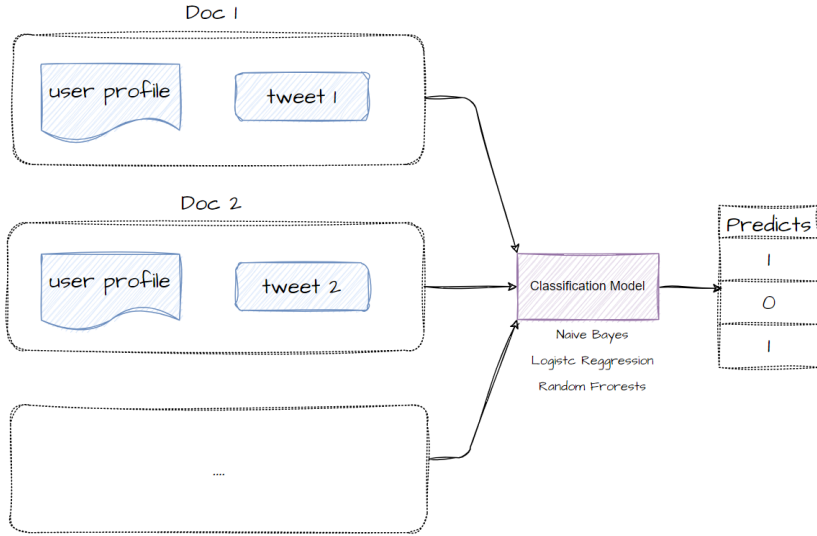


FIGURE 4. *Tweet Based Bot Detection Process*

User-Dimension Based Bot Detection Process

The bot detection process based on tweet dimensions has several drawbacks, including:

- 1) Not every tweet from bot users exhibits characteristics of bots, which may lead to misjudgment.
- 2) Evaluating individual tweets fails to consider the user's behavior from a holistic perspective, hindering a comprehensive assessment based on these behaviors.

To enhance the robustness of bot detection, we propose a user-dimension based bot detection approach, the detailed process of which is illustrated in Figure 5.

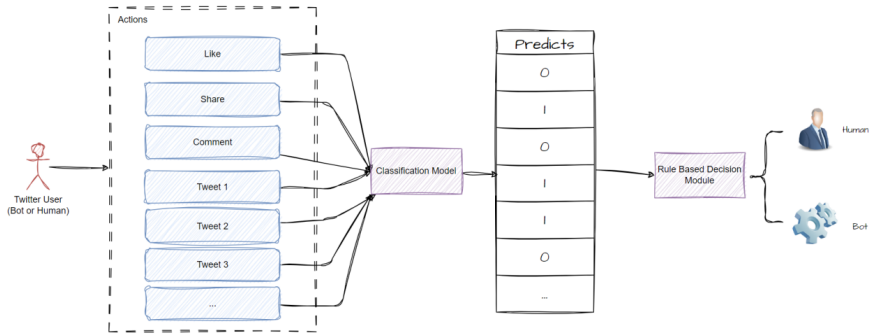


FIGURE 5. *User Based Bot Detection Process*

This process encompasses the following features:

1. It is essential to consider the overall behavior of a user comprehensively. 2. A Twitter user may exhibit a range of behaviors, including:

- Posting tweets
- Engaging in interactions such as liking, retweeting, and saving posts.

The process consists of two core submodules: the first is a classification model, which can directly reuse the models employed in the tweet-dimension based bot detection process. The second is the Rule-Based Decision Module, which integrates the detection results of individual behaviors to yield a final judgment. In this study, we adopt the ratio of individual behaviors identified as bot-like as the determining criterion. If the proportion of a user's actions identified as bots surpasses a certain threshold, that user will be classified as a bot user, as illustrated in Figure 6.

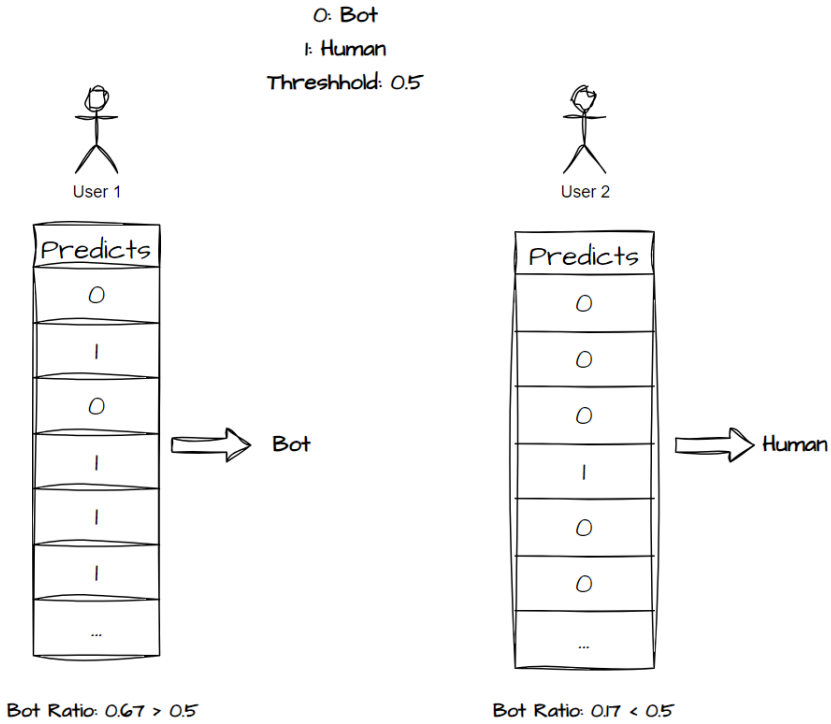


FIGURE 6. Rule-Based Decision Module

Experiment on Disguise Strategies

In real-world scenarios, creators of Bot accounts often employ certain disguise strategies to evade detection by monitoring systems. (Ferrara 2018) Common disguise strategies include:

1. **Mimicking Normal Posting Behavior:** Bot accounts typically simulate the posting style and frequency of ordinary users. For example, they may regularly publish relatively normal content and maintain reasonable time intervals between posts in order to avoid being identified as exhibiting abnormal activity. This content may mimic trending topics or popular trends, further helping them blend into the community.
2. **Acquiring Inactive Real Accounts Through Purchase or Hacking:** Some creators resort to purchasing forgotten or inactive real accounts from the black market, which they then utilize for Bot operations. This method effectively reduces the risk of detection since such accounts usually possess a certain history

and follower base, making them less likely to be immediately identified.

3. **Social Engineering:** By masquerading as ordinary users and interacting with others, Bots can gradually establish trust. This includes actions such as liking, commenting, and sharing, thereby pretending to be an active community member, which makes it more challenging to uncover their Bot identity.

In this task, we aim to incorporate some disguise techniques into the test dataset to verify the robustness of the detection system when faced with these disguise techniques. Specifically, this experiment employs the following procedure to simulate the disguise strategies of Bot accounts:

1. In the posts of Bot accounts, normal user posting content is mixed in at uniform time intervals according to a certain proportion. This involves concatenating a portion of historical tweet data from normal accounts with new tweet data from Bot accounts to construct new fake accounts, thereby simulating those accounts that were improperly acquired for the construction of Bot disguises.
2. A time period is used as a window for the determination of Bot accounts. If an account exhibits a higher proportion of Bot behavior within a certain time period, it is classified as a Bot account. The specific process is illustrated in Figure 7.

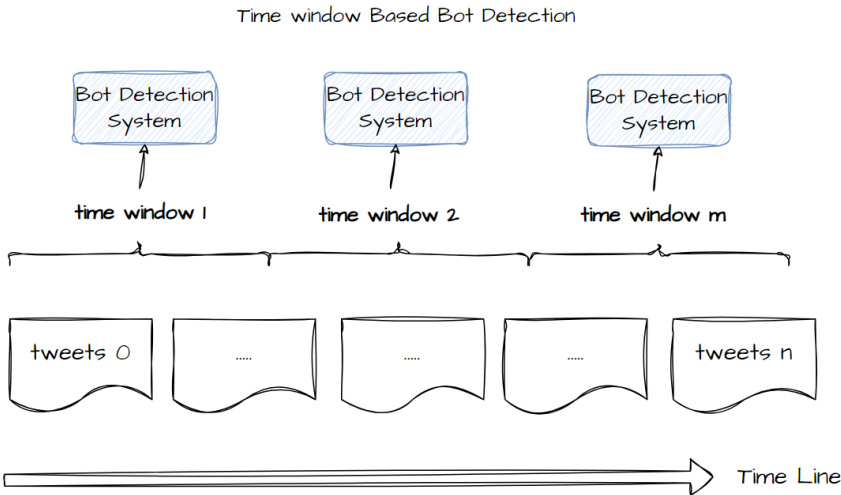


FIGURE 7. Time-Window Based Process

Experimental Comparison

Evaluation Metrics

This experiment uses accuracy as the evaluation metric, validating the algorithmic strategies by splitting the dataset into an 80:20 ratio.

Comparison of Different Classification Models

The performance comparison of different machine learning models on the dataset in this scenario is illustrated in Figure 8. From the experimental results, it can be seen that the logistic regression model exhibits the best performance in terms of accuracy, achieving 66.8%. This indicates that logistic regression is capable of effectively capturing the features within the data, making it suitable for the classification task in this dataset. In contrast, the accuracy of the random forest model is 61.9%, slightly lower than that of the logistic regression model. Although random forests typically demonstrate greater robustness in handling complex data, their performance did not exceed that of logistic regression on this specific dataset. This may be related to feature selection or the distribution of features within the dataset. As for the naive Bayes model, it achieved an accuracy of 60.0%, the lowest among all models. This suggests that the naive Bayes model may be constrained by its fundamental assumptions on this dataset, particularly when the independence assumption between features is not satisfied.

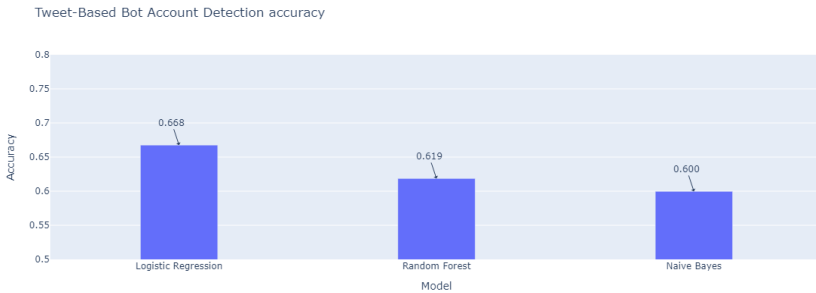


FIGURE 8. *Tweet-Based Process Experiments*

User-Dimension Based Recognition Effectiveness

The performance of three models at different thresholds in the user-dimension based identification process is illustrated in Figure 9. From the experimental results, it can be observed that when the threshold is set between 0.7 and 0.8, the bot detection based on the user dimension achieves the best performance. Among the models evaluated, Logistic Regression remains the most effective. As showing in Figure 10 Utilizing this model leads to an overall improvement of approximately 20 percentage points in recognition performance compared to the tweet-based strategy, with an accuracy metric reaching 86.8%.



FIGURE 9. *User-Based Process Experiments*

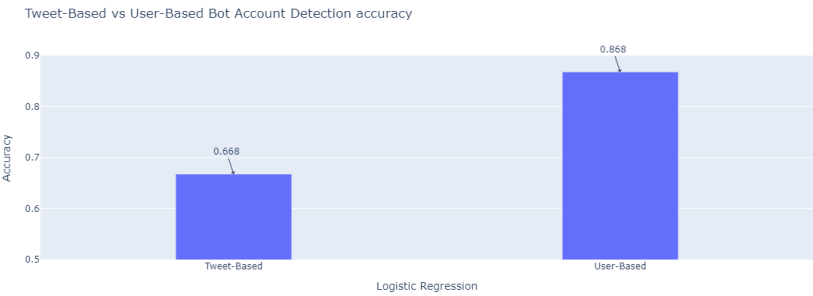


FIGURE 10. *"Tweet-Base Process VS User-Based Process" Experiments*

Experiment on the Recognition Effectiveness of Disguise Technologies

Comparative Effectiveness Against Disguise Technology 1

The experimental results for Disguise Technology 1 are illustrated in Figure ?? . It can be observed that when tweets sent by Bot accounts randomly include 25% of tweets from human accounts, both strategies exhibit a certain degree of decline in effectiveness. However, as long as an appropriate threshold is selected, the user-based recognition strategy can still maintain stability at around 82.5%. The decline magnitude for the two strategies is -6.13% and -4.95% , respectively, indicating that the decline in effectiveness of the user-based recognition strategy is smaller, demonstrating higher robustness.

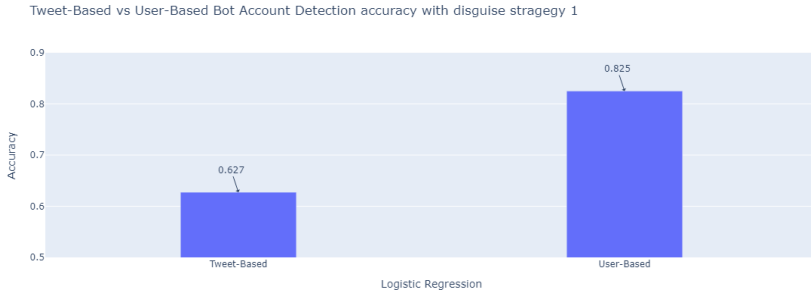


FIGURE 11. *"Tweet-Base Process VS User-Based Process" Disguise Technology 1*

Comparative Effectiveness Against Disguise Technology 2

The experimental results for Disguise Technology 2 are illustrated in Figure 12. It can be observed that in cases where accounts that were once normal human accounts have been subsequently transformed into Bot accounts, both tweet-based and user-based methods experience a significant drop in effectiveness. However, by incorporating a time-window recognition strategy based on the user dimension, this method exhibits excellent robustness against this type of disguise technology. As showing in Figure 13 When the decision threshold is set between 0.7 and 0.8, the overall recognition performance remains almost consistent with conditions where disguises were not applied.

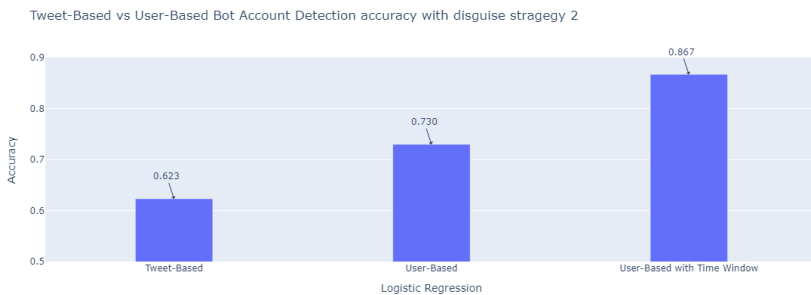


FIGURE 12. *Tweet-Base Process VS User-Based Process VS User-Based Time-Window Process*

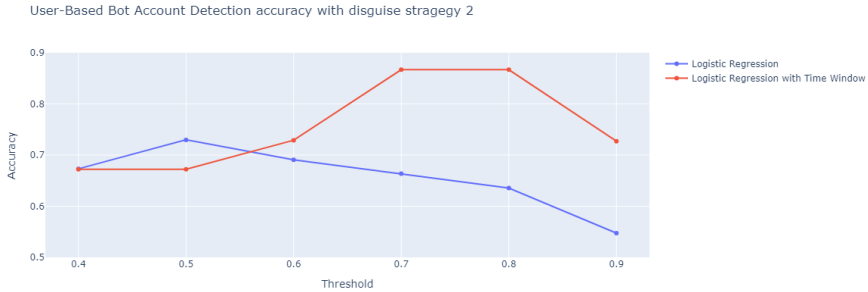


FIGURE 13. User-Based Process Experiments On Disguise Technology 2

Summary and Future Research Directions

This project primarily investigates the application of machine learning algorithms in the identification of Bot accounts, as well as methods to enhance the robustness of identification systems.

Firstly, the project compares the performance of different machine learning models on Twitter data, while also experimenting with the effectiveness of Twitter Bot account identification based on user-level dimensions. It contrasts the performance of user-level processes with that of tweet-level processes, thereby demonstrating the significance of comprehensively considering user behavior information in the task of Bot account detection.

In addition, in response to common Bot account masquerading techniques, this project also experiments with the robustness of these strategies under different masquerading methods. By incorporating a time-window-based Bot account detection approach, the user-level detection system exhibits relatively stable performance across various masquerading techniques.

In future research, in-depth analysis of user behavior could be pursued. Moreover, integrating various link information from social networks and attempting to further optimize Bot account detection performance and robustness based on graph results could be beneficial.

References

- Alarfaj, Fawaz Khaled, Hassaan Ahmad, Hikmat Ullah Khan, Abdullah Mohammed Alomair, Naif Almusallam, and Muzamil Ahmed. 2023. Twitter bot detection using diverse content features and applying machine learning algorithms. *Sustainability* 15 (8):6662.
- Aljabri, Malak, Rachid Zagrouba, Afrah Shaahid, Fatima Alnasser, Asalah Saleh, and Dorieh M Alomari. 2023. Machine learning-based social media bot detection: a comprehensive literature review. *Social Network Analysis and Mining* 13 (1):20.

- Bayes, Thomas. 1968. Naive bayes classifier. *Article Sources and Contributors*, 1–9.
- Breiman, Leo. 2001. Random forests. *Machine learning* 45:5–32.
- Ferrara, Emilio. 2018. Measuring social spam and the effect of bots on information diffusion in social media. *Complex spreading phenomena in social systems: Influence and contagion in real-world social networks*, 229–255.
- LaValley, Michael P. 2008. Logistic regression. *Circulation* 117 (18):2395–2399.