






Iterative Cybercrime Risk and Response of Generative Artificial Intelligence Technology Based on Reinforcement Learning

Ye Huang¹ , Yanwei Xu²  and Juan Luo³ 

^{1,2,3} Cyber Security College, Jiangxi Police Institute, Nanchang, Jiangxi 330000, China, [1HYe16889@163.com](mailto:HYe16889@163.com), 2xywlbqok@163.com, 3luojuan@jxga.edu.cn

Corresponding author: Ye Huang, HYe16889@163.com

Abstract. The rapid development of generative artificial intelligence technology not only improves work efficiency and creates new value but also brings new risks, such as data security, information misdirection, and cybercrime. This paper first analyzes the operational mechanism of generative artificial intelligence technology and then puts forward the corresponding recognition process according to its characteristics. This paper realizes background data mining of generative artificial intelligence technology based on reinforcement learning through a machine learning algorithm and constructs a cybercrime risk model targeting false data and false information through deep learning technology. The experimental results show that the proposed model has better false data injection and false information recognition performance than other models, with a higher recognition rate, better stability, and lower energy consumption, and can continuously realize the cybercrime risk recognition operation. At the same time, the model in this paper can effectively identify information on cybercrime suspicion and abnormal data behavior in different monitoring segments to provide decision-making information for cybercrime risk response.

Keywords: Reinforcement Learning; Production Formula; Artificial Intelligence; Cybercrime Risk; Machine Learning; Deep Learning

DOI: <https://doi.org/10.14733/cadaps.2025.S9.121-134>

1 INTRODUCTION

As a bright new star in the field of artificial intelligence, generative artificial intelligence technology, especially the model based on reinforcement learning represented by ChatGPT, is leading an unprecedented technological revolution. These technologies, through the depth excavation of large-scale data sets and by using the model of a highly complex algorithm achieved from content creation, whether it's accurate and fluent text generation, lifelike image painting, or fascinating video editing, have demonstrated their amazing creativity and application value of [1]. This breakthrough not only greatly enriched the cultural life of mankind but also profoundly changed the face of many industries, such as information communication, education, and

entertainment, and promoted the rapid development of multimodal information processing technology. However, just as any advanced technology is accompanied by the characteristics of double-edged swords, the continuous iteration and popularization of generative artificial intelligence technology has quietly unveiled the tip of the iceberg of its potential criminal risks. Fraudulent activities began to use these technologies to generate realistic false information to mislead the public and defraud money; Data reveal that the risk also will increase because the more complex algorithm means more data processing links, a link any omissions may become the breach of the hacking [2]. In addition, malicious attackers may also tamper with the generated content, spread rumours, undermine social stability, and even threaten national security [3].

The cybercrimes generated in the iterative process of generative artificial intelligence technology are often closely related to people's lives. Generative AI can highly simulate the dialogue mode of real people and make the interlocutor trust, which is easy to use for fraud [4]. There are also a large number of copycat apps with names including "ChatGPT" in the market, which mislead users to click or download through fake links, and then defraud personal information or property [5]. In the face of this kind of cybercrime, people's awareness is not strong, the relevant legal, ethical and social cognitive lag is also a lot of people don't understand the dangers of this kind of cybercrime [6]. At the same time, cybercrimes arising from the iteration of generative artificial intelligence technology, due to its technical characteristics and the concealment of criminal means, make it an extremely challenging task to identify and combat such cybercrimes [7]. Through deep learning and neural network algorithms, it can generate highly realistic and personalized text, pictures, videos and other content. When criminals use these technologies to commit cybercrimes, they are often able to create sophisticated scams that make it difficult for victims to discern authenticity [8]. Due to the covert nature of generative AI technology, criminals can commit cybercrimes without directly revealing their identity. The generative artificial intelligence technology itself is constantly developing, and its algorithms and models are constantly optimized, so that the criminal methods are also rapidly updated. Criminals can quickly grasp new technologies and apply them to criminal activities, making it difficult for law enforcement agencies to keep up with the changing speed of criminal methods [9]. Criminals can use generative AI technology to commit many types of cybercrimes, including fraud, defamation, and infringement of intellectual property rights. The diversification of these criminal methods makes it more difficult to identify and combat them [10]. The content generated by generative AI technology is often digital and easily destroyed or tampered with by criminals. In addition, due to the complexity and anonymity of the network environment, investigation organs often face many difficulties in collecting evidence. Therefore, based on the generating mechanism of artificial intelligence technology, through the collection of machine learning algorithms and deep learning with false information recognition and false data injection as a starting point for emergent artificial intelligence technology iterative cybercrime risk identification, reducing the probability of the cybercrime.

Traditional predictive policing mainly relies on spatial statistics and adaptive regression models, which to some extent improve the efficiency and accuracy of policing work. Reinforcement learning enables AI systems to autonomously learn and make decisions in complex and ever-changing environments through continuous trial and error and optimization strategies, thereby more accurately predicting the likelihood, location, and timing of cybercrime occurrence. Generative AI can not only predict future events but also generate simulated data such as possible cybercrime scenarios and modus operandi. However, generative AI technology based on reinforcement learning has brought a new perspective to cybercrime risk prediction. By combining biometric and multi-factor authentication technologies, stricter identity verification processes are set up for high-risk user accounts to reduce the deception of account cloning. Reinforcement learning models can learn and identify the unique network topology of Sybil attacks, such as abnormally dense node connections, low entropy distributions, etc., effectively blocking the spread of attacks.

2 THE OPERATION MECHANISM OF GENERATIVE ARTIFICIAL INTELLIGENCE TECHNOLOGY BASED ON REINFORCEMENT LEARNING

Generative artificial intelligence is a technology that generates new data by learning its distribution. Based on existing data samples, Sarker [11] researched and generated similar or new data, such as images, text, audio, etc. The core of generative artificial intelligence lies in its ability to capture the inherent patterns and features of data and generate innovative and diverse content based on this. Mainly based on deep neural networks, generative artificial intelligence trains on large-scale datasets, learns and abstracts the basic laws and probability distributions of data, and then generates new data. Its core lies in a profound understanding of data and pattern recognition. Data collection, pre-training, supervised fine-tuning, reinforcement learning, and content generation are the core technical links of reinforcement learning-based generative artificial intelligence technology, which together constitute the complete workflow of generative artificial intelligence. With the increasing maturity of artificial intelligence technology, artificial intelligence is gradually penetrating into all aspects of human life, which represents the arrival of the era of artificial intelligence. However, artificial intelligence is not absolutely safe, and technological advancements have led to greater criminal risks associated with AI-related cybercrimes compared to traditional forms of cybercrime. Meanwhile, due to its deep learning capabilities, Sikder and Harvey [12] are able to make autonomous decisions and can also share human work to carry out activities independently. In situations where artificial intelligence is beyond human control and commits criminal acts, the autonomy and unpredictability of AI's decision-making make it difficult to determine the subject of criminal responsibility in the process of determining criminal liability. The anthropomorphism of artificial intelligence behaviour allows it to replace humans in engaging in difficult production activities, and its powerful data analysis capabilities enable it to help humans store and process massive amounts of data information. In the context of using artificial intelligence as a criminal tool, new criminal methods and scenarios will emerge, resulting in the inability to clearly apply the charges and sentencing ranges in traditional criminal law. In situations where artificial intelligence is used as a criminal target, the data stored and processed in AI will face significant risks. At the same time, this autonomy and unpredictability will increase the management obligations of artificial intelligence designers, manufacturers, and users, and make it difficult to determine the causal relationship between behaviour and results. Therefore, in order to create a safe social environment and ensure the standardized and safe development of artificial intelligence, how to deal with the criminal law risks brought by AI-related cybercrimes has become an urgent problem to be solved in criminal law. Tahmasebi [13] comprehensively and systematically addresses the criminal law risks associated with AI-related cybercrimes. Before proposing specific solutions, it is necessary to clarify that the basic stance of criminal law response is to maintain the modesty of criminal law, abide by the principle of legality of cybercrimes and punishments, adhere to the principle of technical neutrality, and carry out specific responses under the guidance of the basic stance. For the criminal law risks arising from "tool utilization" cybercrimes, it is necessary to fully utilize expanded interpretations to address new criminal behaviours and scenarios within the traditional criminal law framework. For the criminal law risks arising from "criminal object type" cybercrimes, it is necessary to distinguish and punish behaviours in the data storage and data processing stages. Firstly, based on the position of artificial intelligence in cybercrime, cybercrimes involving artificial intelligence should be classified into three categories: "tool utilization type" cybercrimes, "criminal object type" cybercrimes, and "AI out of control type" cybercrimes. Then, based on the analysis of the characteristics of each cybercrime type, the criminal law risks brought by different cybercrime types should be sorted out. If artificial intelligence is allowed to develop freely without regulation, it will inevitably bring greater risks. Reasonably establish the duty of care for artificial intelligence designers, manufacturers, and users, and use normative causal relationship theory to identify the causal relationships involved.

As artificial intelligence increasingly penetrates various fields of society to create more convenient conditions, it has also given rise to various new forms of criminal activities, especially in the field of negligent cybercrimes. Specifically, emerging features such as deep learning,

autonomous manipulation, and open swarm intelligence in artificial intelligence hinder the ability of natural person actors to anticipate the outcomes of subsequent intelligent activities that occur outside of the initial programming process. The fault attribution model centred on foreseeability by Treiblmaier and Rejeb [14] improperly increases the criminal liability risk for developers and other technical personnel or users. There is currently no unified R&D standard and regulatory framework in the field of artificial intelligence, and there are differences in standards between different fields. Specific attention obligations cannot be refined, making it difficult to define objective attention obligations in the field of AI-related cybercrimes. Due to the autonomy of artificial intelligence in enhancing criminal outcomes, neutral intelligence technology implies risks to personal and property safety, making it impossible to reasonably determine the criminal responsibility of negligent offenders based on traditional criminal law, creating practical and theoretical challenges. Based on this, scholars of the new theory of negligence propose the objective and subjective duty of care for negligent cybercrimes in the field of artificial intelligence. The subjective duty of care is related to the ability of the actor to foresee, and based on the individual's subjective ability, it is necessary to answer the question of whether there is a possibility of avoiding the consequences of legal interest infringement caused by artificial intelligence for the actor.

The emerging features of deep learning in artificial intelligence lead to data black boxes within the scope of autonomous manipulation processes. Yue and Shyu [15] determined that violating the subjective duty of care through individual ability exceeded the scope of ordinary people's ability. The autonomous activities of artificial intelligence lead to complex causal relationships while also blurring the empirical direction of causal relationships. Different attribution paths are constructed at different levels and categories, which can be specifically divided into legislative thinking paths and interpretive improvement paths. Interpreters, within the basic logic of current criminal law regulation of cybercrimes, use reasonable theoretical explanations to eliminate or alleviate the dilemma of attribution, which not only meets the practical risks faced by artificial intelligence but also conforms to the principle of criminal law modesty, and supports the viewpoint of this article. Faced with the new forms of infringement of legal interests brought about by the emerging field of artificial intelligence, the criminal law community has extensively explored the criminal responsibility of artificial intelligence with the aim of addressing the dilemma of fault attribution related to artificial intelligence. This breaks down the data barriers in traditional administrative structures and lays a solid foundation for establishing an efficient and collaborative cybercrime management system.

3 ITERATIVE CYBERCRIME RISK IDENTIFICATION MODEL BASED ON ARTIFICIAL INTELLIGENCE TECHNOLOGY BASED ON MACHINE LEARNING AND DEEP LEARNING

3.1 Background Feature Mining of Generative Artificial Intelligence Technology Based on Reinforcement Learning

In the reinforcement learning phase, a carefully constructed reward function of the quality evaluation model answers "the referee"; it can, according to the humans, give positive or negative feedback with a clear reward signal model. After receiving these signals, the model will intelligently adjust its internal strategies and parameters to generate answers more closely to human expectations the next time, thus improving the quality and accuracy of the answers. This kind of based on human feedback mechanism of reinforcement learning not only greatly enhances the adaptability and flexibility of the GPT model makes it able to be quickly adapted to different areas and the scene needs and also makes the model more accurately capture and reflect humanity's preferences and needs. After careful polishing of deep reinforcement learning, the language model highly simulates the human thinking mode and can respond to the user's input or instruction nimbly and generate the expected language output. This process fully demonstrates the outstanding capabilities of generative AI in the field of mass data processing and analysis. What should not be ignored, however, is that in the process of model building and training, research and development of personal preferences and potential bias will inevitably infiltrate among them,

quietly to generate a type of artificial intelligence to set the limitation of internal boundary. These limitations may be reflected in the tendency to homogenize generated content, the unconscious transmission of bias, and a lack of innovative thinking and creativity. Figure 1 shows the schematic diagram of the operational mechanism of generative artificial intelligence technology for reinforcement learning.

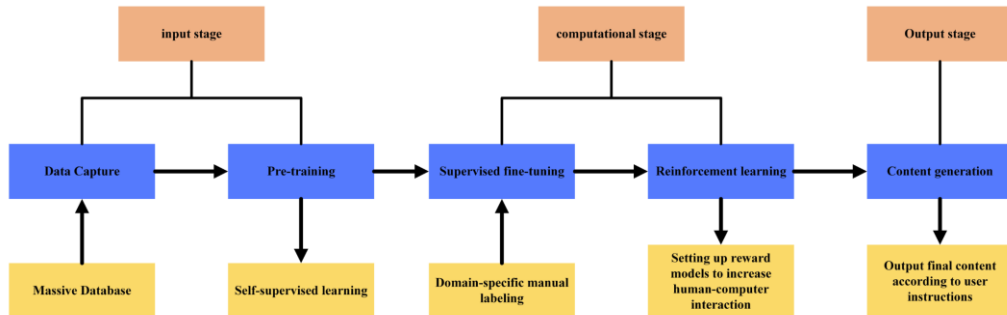


Figure 1: Schematic diagram of the operational mechanism of generative artificial intelligence technology based on reinforcement learning.

Generated based on the operation mechanism of artificial intelligence technology, the identification model of a train of thought of this article is based on the objective basis, by machine learning and neural network method to optimize the depth of the fusion effect, which includes the two core processes. The first process uses machine learning algorithms to dig deep into the background features of generative AI content creation and conduct detailed analysis to accurately determine whether the research preconditions are met. The process then uses the high sensitivity of the neural network model to capture and identify signs of negligent criminal activity to further verify that the second prerequisite is met. On this basis, the predictive model makes forward-looking assumptions about potential hazards to scientifically estimate their possible negative effects. The second process focuses on the intelligent identification of key users and selects the user groups requiring special attention from the massive user data through the analysis ability of machine learning and neural networks. The user is then set to an automatic prompt object, through personalized information push or warning, to improve overall safety and efficiency.

The core mission of generative AI background feature mining is to accurately measure the proportion of AIGC technical elements in the information environment to achieve early warning and scientific assessment of potential destructive effects of activities. The process begins with the extensive collection of information sources, the extraction of mathematical indicators closely related to the AIGC context, and the comprehensive integration of technologies to ensure the full coverage and high accuracy of the collected data. Then, using advanced feature selection mechanisms, the data is dissected to precisely identify and extract the subset of features that are critical to the analysis. At the same time, this process also involves the elimination of insignificant or repetitive feature information, aiming to simplify the analysis complexity, speed up the processing speed, and improve the overall analysis performance. The selected eigenvalues are then introduced into a carefully constructed computational model that incorporates complex algorithmic logic to deeply mine the hidden information and key patterns in the AIGC background features to produce valuable mining results. To ensure the reliability and universality of the results, the whole mining process has experienced many tests and iterations and constantly optimized the calculation process and parameter configuration until a set of accurate and widely applicable universal standards has been formed. The establishment of this standard not only provides a solid foundation for the damage impact assessment of current AIGC activities but also points out the direction for future analysis and development in related fields. Table 1 shows the background feature mining standards for generative artificial intelligence.

| <i>Feature category</i> | <i>Feature paradigm</i> | <i>The value range of the feature paradigm</i> | <i>Feature paradigm</i> | <i>The value range of the feature paradigm</i> | <i>Feature paradigm</i> | <i>The value range of the feature paradigm</i> |
|-------------------------|-------------------------|--|--------------------------|---|----------------------------|--|
| Content ratio | Lower proportion | The content ratio is 0-33.33% | Medium proportion | The content ratio is 33.33%-66.66% | Relatively high proportion | Content 66.66% - 99.99% |
| Number of tools | Thin quantity | The number of tools ranges from 0 to 2 | Considerable quantity | The number of tools is between 3 and 5 | Abundant in quantity | The number of tools is more than 6 |
| Technology type | Single type | Technical model of pure language class | Various types | Language and image technology model | Type Complex | Language images and audio and video technology |
| Formation rate | Slower speed | The average speed is below 5 s/article | Medium speed | Average speed in article 6 to 8 s/a | faster | The average speed is higher than 8s/ bar |
| Modal scale | Scale limitation | The generated content modes are below 2 types | multi-scale | Generated content mode between 2-3 classes | Grand scale | The generated content modes are higher than 3 types |
| Affective granularity | Coarse grain size | Content emotion type is lower than 2 categories | Particle size in general | Content emotional type between 2-3 categories | Finer granularity | The content emotion type is higher than the 3 categories |
| Subject intensity | Weak strength | The content contains less than two topics | Ordinary-strength | Content The number of topics ranges from 2 to 3 | Intense intensity | The number of content topics is greater than 3 |
| Object element | oligo element | The content mentions less than two objects | Element confounding | The content refers to objects between 2-3 | Numerous elements | The content mentions that object 3 higher |
| Escape sequence | Escape deletion | Content semantic conversion is less than 2 times | Escape scarce | Content semantic conversion between 2-3 times | Frequent escape | Content semantic conversion is higher than 3 times |

Table 1: Generative artificial intelligence background feature mining standards.

In this paper, the filter method in machine learning is selected for generative artificial intelligence background feature mining. Pearson correlation coefficient expression is shown as (1):

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \tag{1}$$

Type, represented as a sample quantity n , Characteristics of X The observed value is expressed as x_i . The observed value of the target variable Y is expressed as y_i . The corresponding mean is expressed as \bar{x} and \bar{y} .

The statistical expression of the chi-square test is shown in (2):

$$\chi^2 = \sum_{i=1}^k \sum_{j=1}^l \frac{(O_{ij} - E_{ij})^2}{E_{ij}} \tag{2}$$

Where the number of categories of feature X is k , the number of classes of the target variable is l , the observed frequency is denoted O_{ij} , note for the expected frequency E_{ij} .

The mutual information expression is shown in (3):

$$I(X;Y) = \sum_{x \in X} \sum_{y \in Y} p(x,y) \log\left(\frac{p(x,y)}{p(x)p(y)}\right) \tag{3}$$

X and Y The joint probability distribution of $p(x,y)$, the edge probability distributions of the two are respectively expressed as $p(x), p(y)$.

3.2 False Information and Data Recognition Module Based on Deep Learning

Aiming at the characteristics of generative artificial intelligence, this paper adopts a triplet convolutional twin Transformer network and CatBoost model to realize false information and data recognition, which combines the local feature extraction capability of the convolutional neural network, the global context modelling capability of Transformer and the similarity comparison capability of twin network. This network structure can capture both local details and global context information when processing image, text, or sequence data and compare the similarities between different inputs through the twin structure. As shown in Figure 2.

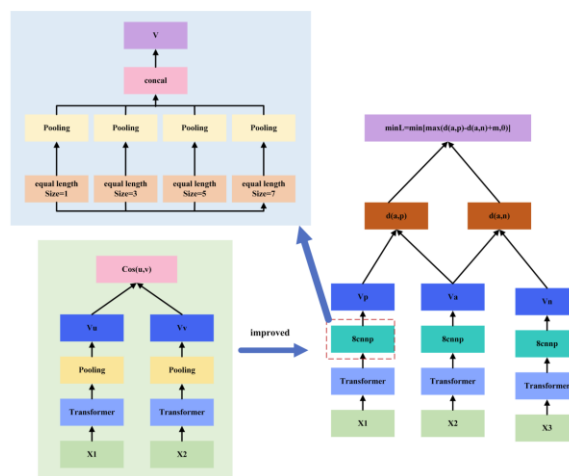


Figure 2: Network structure diagram of the triplet convolutional twin transformer.

A core part of the Transformer model is the self-attention mechanism, which allows the model to process each element in the input sequence by comparing it to other elements in order to process each element correctly in different contexts, as shown in (4):

$$attention(Q, K, V) = \text{soft max}(QK^T / \sqrt{d_k})V \tag{4}$$

The word vector input in the self-attention mechanism will have a query vector, a key vector, and a value vector, and their corresponding matrices are respectively represented in the formula Q, K, V, d_k On behalf of the key dimension vector.

Hypothesis $(K, U) = [(k_1, u_1), (k_2, u_2), \dots, (k_i, u_i)]$, it said the key value of input vector format, query vector described as q when the attention function as shown in formula (5):

$$atta((K, U), q) = \sum_{i=1}^N \frac{\exp(s(k_i, q))}{\sum_j \exp(s(k_j, q))} u_i \tag{5}$$

Formula, $s(k_i, q)$ represents the attention-scoring function.

Let the input sample description be $X = [x_1, x_2, \dots, x_i]$ common attention, the scoring function contains four additive models, respectively, the dot product model, scaling the dot product model, and bilinear model, the formula as shown in (6) - (9):

$$s(x, q) = b^T \tanh(Wx + Dq) \tag{6}$$

$$s(x, q) = x^T q \tag{7}$$

$$s(x, q) = \frac{x^T q}{\sqrt{L}} \tag{8}$$

$$s(x, q) = x^T Wq \tag{9}$$

Type, the dimension of the input vector is represented as L The learnable parameters are expressed as W, D, b .

Triplet improvement focuses on loss function optimization, as shown in (10) for triplet loss function calculation formula:

$$\min L = \min[\max(d(a, p) - d(a, n) + \text{margin}, 0)] \tag{10}$$

Where, the loss value is expressed as L , the target sample is represented as a , samples with the same label are represented as p , the spatial distance between the two is expressed as $d(a, p)$, Samples that are different from the target sample label are represented as n , the spatial distance between the two is expressed as $d(a, n)$, the boundary value is denoted by margin . The process of triplet convolution is shown in Figure 3.

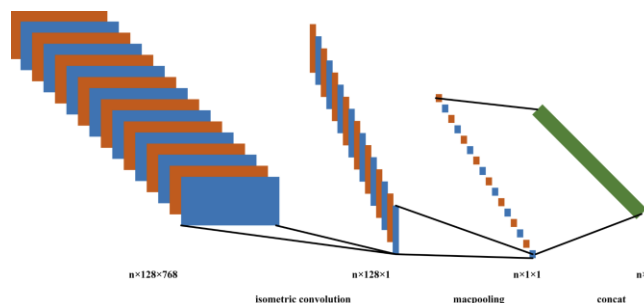


Figure 3: Triplet convolution convolution process.

CatBoost is an integrated learning model based on a gradient lifting decision tree, which is particularly good at processing categorical features and has the advantages of automatic feature scaling, strong robustness, and fast training speed. CatBoost simplifies data preprocessing by directly processing categorical features through an ordered catalogue classification algorithm without complex feature transformation. In gradient lift, model $F(x)$ By iteratively adding a weak learner $h_t(x)$ to build, as shown in (11):

$$F_t(x) = F_{t-1}(x) + \alpha_t h_t(x) \quad (11)$$

Where, the learning rate is expressed as α , in a sequence of t if learning it is expressed as $h_t(x)$.

Based on the calculation formula of target variable statistics such as shown in (12):

$$TS_lable = \frac{count(fea = t \& y = 1)}{count(fea = t)} \quad (12)$$

Where, the coded value is denoted as TS_lable , characterized by t the sample size is expressed as $count(fea = t)$.

The improved calculation formula is shown in (13):

$$TS_lable_k = \frac{count_k(fea = t \& y = 1) + \alpha p}{count_k(fea = t) + \alpha} \quad (13)$$

Where, the sequence number is k the sample encoding value TS_lable_k , under this ordering condition, the number of sequences is k The matching number of previous samples $count_k$. A priori experience is expressed as p and its weight coefficient is expressed as a .

4 EXPERIMENTAL RESULTS AND ANALYSIS

The generative artificial intelligence technology based on reinforcement learning will increase a lot of potential cybercrime risks due to the improvement of technology in the iterative process and provide strong technical support for criminals. Therefore, this type of cybercrime is diverse and wide in scope. In order to verify the identification and early warning ability of the iterative cybercrime risk identification model based on artificial intelligence technology based on machine learning and deep learning, this paper simulated the two most common high-risk cyber crime scenes in life for simulation experiments.

4.1 Simulation Results of False Data Injection Attack Identification

Fake data injection attacks are network security problems, and the iteration of artificial intelligence technology based on machine learning and deep learning enables attackers to design more covert ways to inject fake data. An attacker can quickly generate a large number of false data and carry on the optimization to bypass the existing security protection measures. This allows fake data injection attacks to cause significant damage to targeted systems in a short period of time. In this paper, two other false data injection attack identification methods are selected for comparative experiments. As shown in Figure 4, the three methods calculate the cost test results. The results in the figure show that there is a positive correlation between the calculation cost of the three recognition methods and the number of samples. The RN - LSTM algorithm computational overhead is the highest value, and this paper models the computational overhead value as obviously lower than the other two models. It shows that this model can be completed in a shorter period of time to identify false data injection attacks.

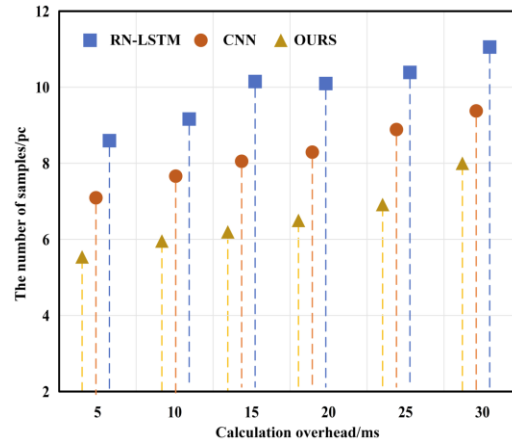


Figure 4: Three kinds of methods to calculate overhead test results.

As shown in Figure 5 for three different recognition models, as a result of the energy consumption contrast experiments after many tests, each test of fake packet size will increase on the basis of the previous number. The results show that with the increase in the size of false packets, the energy effect of the three recognition models is improved to different degrees. Among them, the RN-LSTM algorithm consumes the most energy due to its large structure. Therefore, it has the highest energy consumption and the fastest improvement speed when the scale of false data continues to increase. The energy consumption of the model in this paper is the lowest and increases slowly with the increase in the number of false packets. This shows that in the application of false data recognition, the proposed model can last longer.

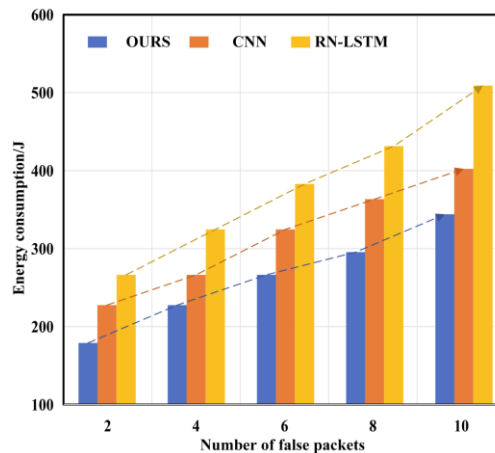


Figure 5: Comparison of energy consumption of three different recognition models.

As shown in Figure 6 for the three methods of the false data injection attacks of recognition rate comparing the results, the experiment is also a model identification of the inspection after many experiments, the number of samples for each test will be based on the previous quantity increases continuously. Results show that the three methods of recognition rate and the relationship between the sample size have a negative correlation. In the case of a small number of samples, the recognition rate of the other two models is higher than 95%, and there is little difference

between the recognition rate and that of the model in this paper. However, with the increase in the number of samples, the recognition rate of the two models has a significant decline, and the decrease is large, and the gap between the recognition rate and the model in this paper is increasing. In this paper, the model recognition rate is above 91%, and with the increase in sample size, its recognition rate lower amplitude is far less than other models. This shows that the model can effectively maintain a high recognition rate of false data injection and ensure the reliability and effectiveness of the corresponding data analysis.

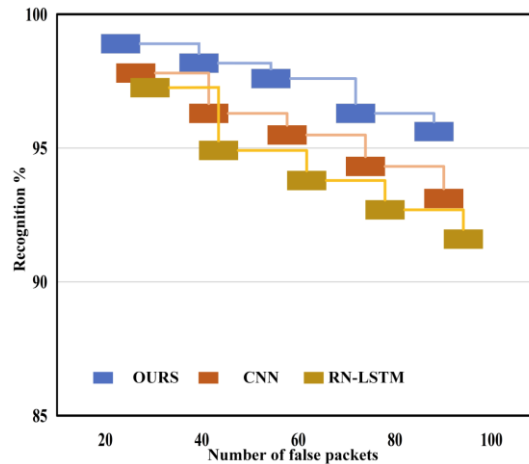


Figure 6: Comparison of recognition rates of fake data injection attacks by three methods.

4.2 Simulation Results of False Information Recognition

The iteration of generative artificial intelligence technology based on reinforcement learning makes it easier for criminals to use AI technology to synthesize false information such as faces, voices, pictures, or videos, and the generated false information is often highly realistic and confusing, making it difficult for ordinary users to distinguish its authenticity. This low threshold characteristic greatly reduces the cost of false information cybercrime so that more criminals can participate in such criminal activities. Therefore, through simulation experiments, this paper tested the identification performance of artificial intelligence based on machine learning and deep learning iterative cybercrime risk recognition model on false information. In this paper, five models are selected for comparative experiments, and the results are shown in Figure 7. As shown in the figure, it can be seen from the results of all indicators that the indicators of the model in this paper are significantly different from other models, and the rankings of F1 and AUC in the figure are basically consistent. This shows that the integrated model in this paper is superior to other models of false information and has better recognition performance.

Figure 8 shows the comparison results of the structural performance of different models. Compared with other models, the structural performance indexes of the model in this paper are the highest, which indicates that the combined structural performance of the model in this paper has good stability, can effectively realize the identification of false information, and provides good data support for practical applications.

In order to test the application performance of the model further, this paper conducts corresponding cybercrime risk monitoring for four different monitoring segments through simulation experiments, and the results are shown in Figure 9. The diagram, according to the results of four monitoring periods for suspected cybercrimes, is more than 50%, is a high risk, and monitoring section D is most at risk.

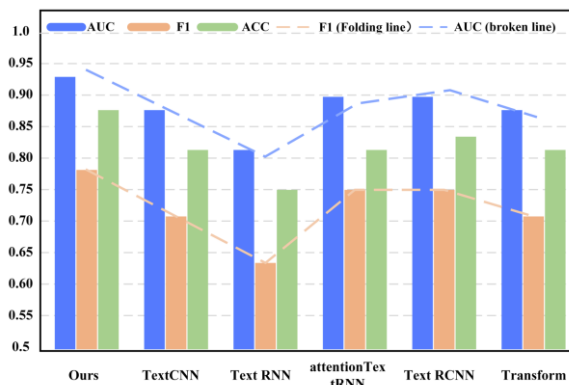


Figure 7: Performance comparison results of different models.

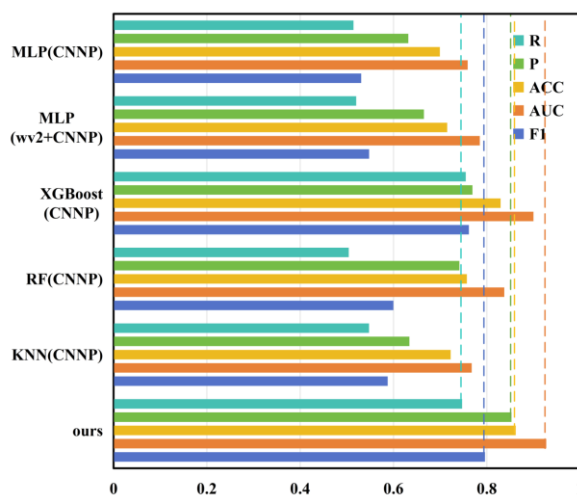


Figure 8: Comparison results of structural performance of different models.

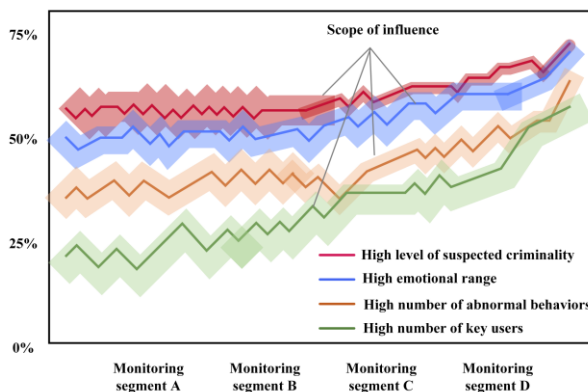


Figure 9: Abnormal behaviours of different monitoring segments and identification results of key users.

In terms of emotional amplitude, the emotional amplitude of the four monitoring segments showed a state of great change. In terms of the number of abnormal behaviours and the number of end users, the number of monitoring segment D is the largest, which has reached more than 45%, while the number of monitoring segment A is relatively small, remaining at about 25%. This shows that the model in this paper can effectively identify data changes related to false information, and analyze abnormal behaviors and key users. At the same time, it can also present the situation of different detection terminals and take corresponding measures according to the actual situation to strengthen monitoring and reduce the risk of cybercrime.

5 CONCLUSIONS

Although the generative artificial intelligence technology based on reinforcement learning has brought great convenience and value to human society, its potential criminal risk cannot be ignored. In this paper, through the systematic analysis of the operational mechanism of generative artificial intelligence technology, through the mining of its background data, the collection of machine learning and deep learning technologies to build a cybercrime risk identification model against false data and false information. The experimental results show that this model is significantly superior to other models in false data injection and false information recognition performance, and shows good stability, with the highest recognition rate. In the simulation experiment, the suspected cybercrime probability and other information of different monitoring segments can be identified, and the abnormal behavior and key user behavior among them can be identified so as to provide strategic information for the corresponding cybercrime risk. In terms of coping strategies, although the model in this paper can provide certain technical support and improve the identification of cybercrime risk, it has limited effect on the entire iterative cybercrime risk identification of artificial intelligence technology and the coping strategies at the level of law and supervision are also needed as a guarantee. Governments and international organizations should strengthen the formulation of policies and regulations on generative AI and clarify the application boundaries and operational responsibilities of such technologies. Establish a regulatory agency or department, that is responsible for the supervision of the application of emergent artificial intelligence and law enforcement work. Strengthen daily inspections and regular inspections of generative AI applications to detect and deal with potential criminal risks in a timely manner. We will increase penalties for violations of laws and regulations to form an effective deterrent effect. In addition, you also need to improve public awareness, strengthen public awareness of the emergent artificial intelligence technology and understanding, and improve the ability to distinguish false information. The potential risks and safe use knowledge of generative AI technology should be widely publicized through media, the Internet, and other channels to guide the public to use relevant technology correctly.

6 ACKNOWLEDGEMENT

Science and technology research project of Jiangxi Provincial Department of Education, Research on cybercrime risk and response of iteration of generative artificial intelligence.

Ye Huang, <https://orcid.org/0009-0009-9749-5743>

Yanwei Xu, <https://orcid.org/0009-0004-6688-3907>

Juan Luo, <https://orcid.org/0009-0005-6813-2223>

REFERENCES

- [1] Baak, M.; Koopman, R.; Snoek, H.; Klous, S.: A new correlation coefficient between categorical, ordinal and interval variables with Pearson characteristics, *Computational Statistics & Data Analysis*, 152(1), 2020, 107043. <https://doi.org/10.1016/j.csda.2020.107043>

- [2] Baltuttis, D.; Teubner, T.; Adam, M.-T.: A typology of cybersecurity behavior among knowledge workers, *Computers & Security*, 140(1), 2024, 103741. <https://doi.org/10.1016/j.cose.2024.103741>
- [3] Berk, R.-A.: Artificial intelligence, predictive policing, and risk assessment for law enforcement, *Annual Review of Criminology*, 4(1), 2021, 209-237. <https://doi.org/10.1146/annurev-criminol-051520-012342>
- [4] Chen, W.-D.; Murtazashvili, I.: Blockchains for emergency and crisis management, *Public Administration Review*, 83(5), 2023, 1409-1414. <https://doi.org/10.1111/puar.13647>
- [5] Dai, D.; Boroomand, S.: A review of artificial intelligence to enhance the security of big data systems: state-of-art, methodologies, applications, and challenges, *Archives of Computational Methods in Engineering*, 29(2), 2022, 1291-1309. <https://doi.org/10.1007/s11831-021-09628-0>
- [6] Hajek, P.; Barushka, A.; Munk, M.: Fake consumer review detection using deep neural networks integrating word embeddings and emotion mining, *Neural Computing and Applications*, 32(23), 2020, 17259-17274. <https://doi.org/10.1007/s00521-020-04757-2>
- [7] Hoffman, C.-J.; Howell, C.-J.; Perkins, R.-C.; Maimon, D.; Antonaccio, O.: Predicting new hackers' criminal careers: A group-based trajectory approach, *Computers & Security*, 137(1), 2024, 103649. <https://doi.org/10.1016/j.cose.2023.103649>
- [8] Jethava, G.; Rao, U.-P.: Exploring security and trust mechanisms in online social networks: an extensive review, *Computers & Security*, 140(1), 2024, 103790. <https://doi.org/10.1016/j.cose.2024.103790>
- [9] Nama, M.; Nath, A.; Bechra, N.; Bhatia, J.; Tanwar, S.; Chaturvedi, M.; Sadoun, B.: Machine learning-based traffic scheduling techniques for intelligent transportation system: Opportunities and challenges, *International Journal of Communication Systems*, 34(9), 2021, e4814. <https://doi.org/10.1002/dac.4814>
- [10] Peker, I.; Ar, I.-M.; Erol, I.; Searcy, C.: Leveraging blockchain in response to a pandemic through disaster risk management: an IF-MCDM framework, *Operations Management Research*, 16(2), 2023, 2023642-667. <https://doi.org/10.1007/s12063-022-00340-1>
- [11] Sarker, I.-H.: Multi-aspects AI-based modeling and adversarial learning for cybersecurity intelligence and robustness: A comprehensive overview, *Security and Privacy*, 6(5), 2023, e295. <https://doi.org/10.1002/spy2.295>
- [12] Sikder, A.-S.; Harvey, K.: Techno-resilience: unraveling the impact of cutting-edge information technology in crisis management and emergency response for enhanced disaster preparedness and response efficiency.: IT in crisis management and emergency response, *International Journal of Imminent Science & Technology*, 1(1), 2021, 138-169. <https://doi.org/10.1016/j.ijinfomgt.2019.102049>
- [13] Tahmasebi, M.: Beyond defense: proactive approaches to disaster recovery and threat intelligence in modern enterprises, *Journal of Information Security*, 15(2), 2024, 106-133. <https://doi.org/10.4236/jis.2024.152008>
- [14] Treiblmaier, H.; Rejeb, A.: Exploring blockchain for disaster prevention and relief: A comprehensive framework based on industry case studies, *Journal of Business Logistics*, 44(4), 2023, 550-582. <https://doi.org/10.1111/jbl.12345>
- [15] Yue, Y.; Shyu, J.-Z.: A paradigm shift in crisis management: The nexus of AGI-driven intelligence fusion networks and blockchain trustworthiness, *Journal of Contingencies and Crisis Management*, 32(1), 2024, e12541. <https://doi.org/10.1111/1468-5973.12541>