WILEY | Hindawi

## Research Article

# A Hybrid Alarm Association Method Based on AP Clustering and Causality

**Xiao-ling Tao** [1,2] **Lan Shi** [1] **Feng Zhao** [3] **Shen Lu** [1] **and Yang Peng** [1]

[1]*Guangxi Key Laboratory of Cryptography and Information Security, Guilin University of Electronic Technology, Guilin 541004, China*
[2]*Guangxi Cooperative Innovation Centre of Cloud Computing and Big Data, Guilin University of Electronic Technology, Guilin 541004, China*
[3]*School of Information and Communication, Guilin University of Electronic Technology, Guilin 541004, China*

Correspondence should be addressed to Feng Zhao; zhaofeng@guet.edu.cn

Internet of Things (IoT) brought great convenience to people's daily lives. Meanwhile, the IoT devices are facing severe attacks from hackers and malicious attackers. Hackers and malicious attackers use various methods to invade the Internet of Things system, causing the Internet of Things to face a large number of targeted, concealed, and penetrating potential threats, which makes the privacy problem of the Internet of Things suffers serious challenges. But the existing methods and technologies cannot fully identify the attacker's attack process and protect the privacy of the Internet of Things. Alarm correlation method can construct a complete attack scenario and identify the attacker's intention by alarming the alarm data which provides an effective protection for user privacy. However, the existing alarm correlation methods still have the disadvantages of low correlation accuracy, poor correlation efficiency, and strong dependence on the knowledge base. To address these issues, we propose an alarm correlation method based on Affinity Propagation (AP) clustering algorithm and causal relationship. Our method considers that the alarm data triggered by the same attack process has high similarity characteristics, adopts the AP algorithm to improve the correlation efficiency, and at the same time constructs a complete attack process based on the causal correlation idea. The new alarm correlation method has a high correlation effect and builds a complete attack process to help managers identify attack intentions and prevent attacks.

## 1. Introduction

Smart city and intelligent transportation system improved the people's lifestyle. The Internet of Things (IoT) applications brought great convenience to people's lives [1]. IoT is seen as the third wave and revolution in the development of the global information industry after the advent of computers and the Internet. By the huge market scale and broad industry application prospect, IoT has become the current hot research field. With the continuous change of technology and the advent of 5G networks, the scale and complexity of the Internet of Things continue to increase, and the complex network architecture of heterogeneous integration and interconnection of the Internet of Things is facing increasingly prominent security and efficiency issues [2], and data privacy

has also become one of the most important issues in the Internet of Things [3]. The security issue of the Internet of Things has increasingly become a hot issue that people are concerned about today. According to the white paper on the development of China's network security in 2019 [4], in 2018, the size of China's IoT security market reached 8.82 billion, with a growth rate of 34.7%, which was significantly higher than the industry average growth rate. In recent years, viruses, Trojans, vulnerabilities, spyware, and other attacks and threats against the Internet of things emerge in an endless stream. For example, in 2019, security researchers discovered that popular connected or smart home devices sold by large retailers such as Wal-Mart and Best Buy generally have serious security vulnerabilities and privacy issues. Amazon's Ring also has privacy and security issues, as well

as a series of problems such as the terrifying Mirai botnet continues to maintain a high-speed growth, data upload leakage in intelligent cyber-physical systems [5], and privacy protection of data sharing in the industrial Internet of Things [6]. Various advanced multistep attacks are also appearing more and more frequently. Due to their penetration, pertinence, and concealment, they pose a serious threat to the Internet of Things [7]. In addition to this, the types of attacks on the network are becoming more abundant [8], such as worm attacks, vulnerability attacks, denial of service (DoS) attacks, and phishing attacks. Such a series of intrusions have brought severe challenges to the security and privacy protection of the Internet of Things [9].

In the complex network system, correlation analysis of alarm data is of great importance. It is one of the most effective methods for constructing attack scenarios, allowing managers to intuitively analyse attack trends. The principle of the alarm correlation technology is to dig out the internal connection between the attack events through the correlation analysis and processing of the alarm data [10] and further correlate the alarm information to realize the reconstruction of the attack scenario to help the network manager grasps the entire attack process. Identify the attacker's attack intention, which can effectively prevent network attacks and protect the privacy of the Internet of Things.

There are a series of research in the field of alarm correlation, such as causality-based correlation method [11, 12], data mining method [13, 14], and attribute similarity-based correlation method [15, 16]. However, there are still some problems that need to be solved. First of all, the popularization and diversification of the Internet of Things make the network environment more complicated, and the attacks on the network are also complex and changeable [17]. The existing methods cannot build a more comprehensive attack scene against complex intrusion behaviours. Secondly, because of the high false-positive rate of intrusion detection system, the key attack steps are missing. And the existing methods have strong dependence on the knowledge base, thus affecting the accuracy of correlation results, resulting in low correlation accuracy and poor correlation efficiency. Therefore, how to coordinate the relation between correlation accuracy and correlation efficiency to achieve more ideal effect of alarm correlation is an urgent problem.

*1.1. Contributions.* An efficient alarm correlation method is an effective way to reconstruct attack scenarios for helping network administrators to identify attackers' attack intentions and protect network privacy. Therefore, we propose a hybrid method based on Affinity Propagation (AP) clustering algorithm and causality to correlate alarm data. The main contributions of this paper are summarized as follows:

We improve the similarity calculation method in AP algorithm and use the attribute similarity calculation method to replace the traditional similarity measurement method in AP clustering. According to the different properties of alarm data, we define different similarity calculation functions to calculate their similarity. Combined with the weight of each attribute, we calculate the overall similarity of the alarm. Then, use the AP algorithm to cluster the massive alarm data

and classify the alarm data with higher similarity into the same attack scenario. AP clustering algorithm does not rely on prior knowledge to automatically classify attack scenes, which can greatly improve the correlation efficiency of alarm data.

After dividing the attack scenarios, we sort the alarm data in the same attack scenario in the order of attack time and then associate the alarm data of the same attack event in the same attack scenario according to the principle of causality between the attack sequences. Finally, build a complete attack process.

*1.2. Organization.* The remainder of this paper is organized as follows. The related work is introduced in Section 2. We review the conception of AP clustering algorithm and attribute similarity calculation in Section 3. Our scheme is given in Section 4. Section 5 analyses the effectiveness of our proposed method on the honey pot dataset. We conclude this paper in Section 6.

## 2. Related Work

Most of the attacks launched by attackers on the network are composed of multiple interrelated attack actions, these attacks involve multiple intrusions [18]. Alarm correlation is a technology that extracts effective attribute information from a large amount of original alarm data, connects the alarms induced in the same attack step based on certain rules, and reconstructs the attack process, which can effectively identify attack intentions and reduce repeated alarms. In recent years, the correlation analysis of alarm data has always been a research hotspot in the field of network security [19]. So far, many researchers have done a lot of research on alarm data from the perspectives of causality, data mining, and attribute similarity.

The correlation method based on causality is the most common correlation method. It does not require the support of an expert knowledge base and performs related analysis on the alarm data according to the premise and possible consequences of the attack type [20]. Literature [21] proposed an alert correlation framework (RTECA), the type of framework extracts causality based on Bayesian networks in offline mode and constructs an attack graph. Aiming at the problem that the existing association methods fail to identify many distributed attacks, a real-time alarm correlation method based on attack planning graph (APG) is proposed in reference [22]. This method establishes an attack graph model according to attack types and causality. In order to obtain effective network intrusion alarm information and reveal the intention of attackers, the literature [23] proposed a method to construct attack scenarios based on single-value causality graphs, which constructs attack scenes based on causal graph and can correctly reflect the real hacker intrusion process. The above methods have a strong dependence on prior knowledge. Once an intermediate link in the attack step is missing, a complete attack scene cannot be constructed.

Association method based on data mining is a research hotspot in recent years, it does not need expert knowledge and prior knowledge. It can automatically mine data through

statistical methods to find attack scenarios. Both literature [24] and literature [25] extract complex attack scenarios by mining frequent attack sequences, which can effectively mine attack scenarios and discover more valuable attack patterns. Aiming at the problem that causal knowledge is difficult to obtain automatically in alarm correlation analysis, in reference [26], the transition probability matrix between different attack types is automatically mined based on Markov property, thereby constructs causal knowledge of each attack scenario. In literature [27], a plot mining algorithm is used to discover possible combinations of alarms, and then, a supervised decision tree (DT) learning method is used to detect multistep attack scenarios. The literature [28] proposes an alarm correlation framework based on Markov chain, the framework combines statistics and mining techniques to correlate alerts. Although the abovementioned correlation method does not require a large amount of knowledge base and prior knowledge, there are still defects in the statistical analysis process of large amount of calculation and low accuracy.

The alarm correlation method based on attribute similarity is to judge whether there is correlation between alarms by comparing the alarm similarity and the set threshold. The literature [29] uses the similarity of alarms to determine the causal relationship between alarms and reconstructs the attack scene through the evidence in alarms. This method can quickly and incrementally reconstruct known and unknown attack schemes without expert intervention. Literature [30] determines the causal relationship between attack events by calculating the similarity between attacks, thereby constructing attack paths. Mining association rules from the perspective of alarm timing, literature [31] proposes an alarm correlation method based on block similarity that converts the alarm data sequence into a time node sequence and improves the maximum correlation coefficient method to enhance the correlation accuracy. The alarm association method based on attribute similarity has the advantages of simple algorithm and strong real-time performance, but there is no standard for attribute similarity. The final association result is greatly affected by the parameters such as similarity weight coefficient, and the association result cannot show the relationship between attacks very well.

## 3. Preliminaries

In this section, we review the conception of AP clustering algorithm and attribute similarity calculation.

*3.1. AP Clustering Algorithm.* The AP clustering algorithm is a graph-based clustering algorithm, it was first proposed by Frey and Dueck [32] in Science Journal in 2007. The algorithm regards all samples in the dataset as possible cluster centres and transmits information through iterations between data points. In the process of iteration, the iteration information for each point continues to be updated until m specific cluster centres are produced to achieve the corresponding classification. The basic idea is as follows:

The AP algorithm takes the similarity matrix $S$ formed by similarity among data $N$ points as input for clustering analy-

sis. It uses $s(i, j)$ to represent similarity between node $i$ and node $j$ and introduces the concept of reference $P$ to represent reference degree of data points as clustering centre. The reference degree of point $i$ is expressed as $P(i)$ or $s(i, i)$, and the larger the value, the more likely point $i$ is to be the cluster centre. Because the AP algorithm considers that each data point is likely to be the cluster centre, so all $P$ take the same value, and the final number of clusters is greatly affected by the value of the reference degree. Generally, the median or minimum value of the input similarity value is used as the value of $P$. At the same time, by setting the damping factor ($\lambda$), it avoids data shock during the clustering process and achieves a better convergence effect. Its value range is [0,1].

The AP algorithm also introduces the two concepts of responsibility and availability and realizes the transfer and update of data points by iteratively updating the responsibility matrix and the availability matrix and then obtains the final cluster centre point. The formula used in AP clustering algorithm is given below.

The update formula of responsibility matrix $R$ is as follows:

$$r_{t+1}(i, k) = \begin{cases} s(i, k) - \max_{j \neq k}\{a_t(i, j) + r_t(i, j)\}, i \neq k, \\ s(i, k) - \max_{j \neq k}\{s(i, j)\}, i = k. \end{cases}$$

(1)

The update formula of availability matrix $A$ is as follows:

$$a_{t+1}(i, k) = \begin{cases} \min\left\{0, r_{t+1}(k, k) + \sum_{j \neq i,k} \max\{r_{t+1}(j, k), 0\}\right\}, i \neq k, \\ \sum_{j \neq k}\max\{r_{t+1}(j, k), 0\}, i = k. \end{cases}$$

(2)

At the same time, in order to avoid the problem of data oscillation in the process of matrix update, AP algorithm attenuates the above two formulas by setting damping coefficient, and the update formula is as follows:

$$R_{t+1}(i, k) = \lambda * r_t(i, k) + (1 - \lambda) * r_{t+1}(i, k).$$

(3)

$r_{t+1}(i, k)$ represents the responsibility of point $i$ and point $k$ after the $t + 1$th update, and $R_{t+1}(i, k)$ represents the degree of responsibility after attenuation.

$$A_{t+1}(i, k) = \lambda * a_t(i, k) + (1 - \lambda) * a_{t+1}(i, k).$$

(4)

$a_{t+1}(i, k)$ represents the availability degree after the $t + 1$ th times update, and $A_{t+1}(i, k)$ represents the availability after attenuation.

The flow of AP algorithm is as follows:

*Step 1.* Set the initialized responsibility and availability matrix as 0 matrix and set parameters damping factor and maximum iteration times MaxIterNum.

*Step 2.* Input the data to calculate the similarity matrix *S* and then calculate the median value of the similarity matrix and assign it to the parameter preference.

*Step 3.* Use formula (1) to update the responsibility matrix.

*Step 4.* Use formula (2) to update the availability matrix.

*Step 5.* Attenuate formula (1) and formula (2) according to the attenuation coefficient.

*Step 6.* Check whether the clustering result meets the termination condition, if it is satisfied, the algorithm ends and the result is output; otherwise, it returns to step 3 for the next iteration.

*Step 7.* When the algorithm is finished, output the final cluster centre and the dataset of the classified categories.

*3.2. Attribute Similarity Calculation.* Within a certain time-threshold, the alarm data belonging to the same attack scenario must have certain relations in IP address, port, and alarm type. Therefore, when clustering, we use the attribute selection method in literature [33] for reference and conduct correlation analysis from the four attributes of attack type, IP, port, and time.

*3.2.1. Similarity of Attack Types.* For the alarm type attribute, if the two alarm data alert$_i$ and alert$_j$ have the same alarm type, set their similarity to 1; otherwise, it is 0. The calculation formula is as follows:

$$\text{sim}_{\text{type}} = \begin{cases} 0, \text{alter}_i.\text{type} \neq \text{alter}_j.\text{type}, \\ 1, \text{alter}_i.\text{type} = \text{alter}_j.\text{type}. \end{cases} \quad (5)$$

*3.2.2. IP Address Similarity.* The IP address in the alarm log is expressed in decimal form. Therefore, it is necessary to convert the IP address into a binary form first and then calculate the similarity by comparing the same consecutive prefix digits [34]. The calculation formula is as follows:

$$\text{sim}_{\text{ip}} = \frac{r}{32}, \quad (6)$$

where *r* represents the two alarm data alert$_i$ and alert$_j$ from high to low, the IP address is the same number of consecutive digits.

*3.2.3. Port Similarity.* The port number is a Boolean attribute. If two alarm ports are identical, the similarity is considered as 1; otherwise, it is 0. The calculation formula is as follows:

$$\text{sim}_{\text{port}} = \begin{cases} 0, \text{alter}_i.\text{port} \neq \text{alter}_j.\text{port}, \\ 1, \text{alter}_i.\text{port} = \text{alter}_j.\text{port}. \end{cases} \quad (7)$$

*3.2.4. Time Similarity.* For the calculation of time similarity, we first compare the date attributes and then use the sigmoid function to calculate the time similarity for alarms with the same date attributes. Otherwise, the similarity is 0. The calculation formula is as follows:

$$\text{sim}_{\text{timestamp}} = \begin{cases} 0, \text{alter}_i.\text{date} \neq \text{alter}_j.\text{date}, \\ \dfrac{1}{1 + e^t}, \text{alter}_i.\text{date} = \text{alter}_j.\text{date}. \end{cases} \quad (8)$$

That $t = |t_i - t_j|/60$.

After calculating the similarity of four attributes of attack type, IP, port, and time, the overall similarity between alarm data is obtained by taking the weighted average. The formula for calculating the overall similarity between two alarms is as follows:

$$\text{sim}\left(\text{alter}_i, \text{alter}_j\right) = \sum_{l=1}^{6} \text{sim}_l * \omega_l, \quad (9)$$

where sim$_l$ indicates the similarity of each attribute, the $\omega_l$ represents the weight corresponding to each attribute, and the subscript of 6 indicates that the formula is weighted by six attributes, which are attack type, timestamp, source IP address, source port number, destination IP address, and destination port number. The weight of each attribute is determined by principal component analysis based on the idea of reference [35].

## 4. Our Scheme

In this section, we propose the hybrid alarm correlation method based on AP clustering algorithm and causality, it mainly includes three phases: (1) alarm data preprocessing, (2) attack scene division based on AP clustering, and (3) constructing attack process graph based on causality. Our method is based on the idea that the alarm data with high similarity after preprocessing are aggregated into the same cluster by using AP clustering algorithm, so as to realize the division of attack scenarios. Then, the alarm data in the same attack scenario are further correlated and analysed by using causal correlation method to restore the attack process. Our method can restore attack process without setting attack knowledge base, and it well shows the logical relationship among alarm information, eliminates redundant data, improves the correlation accuracy, and realizes multistep attack restoration. The overall flow of the algorithm is shown in Figure 1.

*4.1. Alarm Preprocessing.* Different intrusion detection systems generate different formats of alarm data according to the abnormal conditions of the network environment. These data cannot be directly used for correlation analysis, so attribute filtering and normalization processing of alarm logs in different formats are the bases of subsequent work. We use Intrusion Detection Information Exchange Format (IDMEF) [36] to extract eight attributes from the original log and standardize the format of the original alarm data and define the alarm data as a seven-tuple. The meaning of each attribute is shown in Table 1.
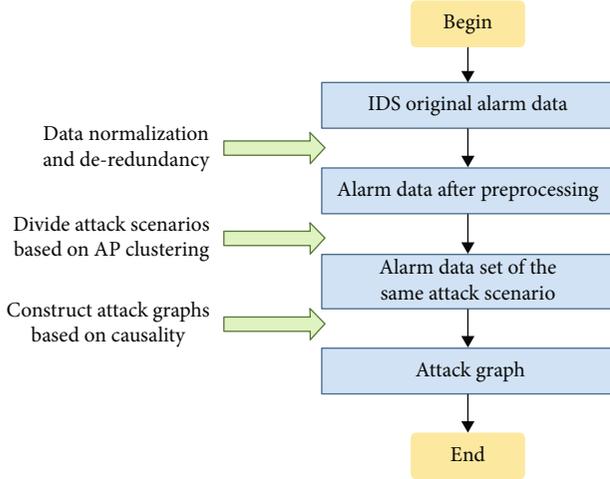
Figure 1: Flow chart of alarm correlation method.

Table 1: Alarm data attribute.

| Attribute | Meaning |
| --- | --- |
| Signature | Characteristic string |
| Type | Alarm category |
| Date | Alarm date |
| Timestamp | Alarm timestamp |
| Src_ip | Source IP |
| Src_port | Source port |
| Des_ip | Destination IP |
| Des_port | Destination port |

Due to the large number of repeated and redundant alarms in the alarm data, it is difficult for us to obtain valuable alarm information from the massive alarm data [37]. In order to solve this problem, we deduplicate and merge the original alarm data. We set a time threshold, and under the condition of not exceeding this time threshold, in addition to the signature attribute, we merge the alarm data with high similarity of other attributes to remove duplicates and add alert_id attribute to the deduplicated data to prepare for the follow-up work.

*4.2. Attack Scene Division Based on AP Clustering.* In the process of alarm correlation analysis, we use AP clustering algorithm to divide attack scenarios. The AP algorithm divides the massive and disordered alarm logs into a collection of attack scenarios with small intraclass spacing and large interclass spacing without prior knowledge, and it does not need to set the clustering number in advance, nor does it need to randomly select the initial clustering centre. It overcomes the defect of the traditional clustering algorithm that is sensitive to the initial conditions. What is more, compared with the K-means clustering algorithm, its fitting degree is much higher than that of the K-means algorithm, and the squared error of the results is also smaller.

The standard AP clustering algorithm uses Euclidean distance as the similarity calculation criterion, but the type of

the alarm log generated by the intrusion detection system is string type, and there is a certain connection between the attributes of the alarm data, and the relative importance of each attribute field is not the same. It is difficult to calculate its similarity by using the distance calculation formula, and it will also destroy the connection between the alarms. Therefore, we improve the similarity matrix calculation method of AP clustering algorithm. According to the attribute similarity calculation method given in Section 3, we calculate the attribute similarity and then use the AP clustering algorithm to aggregate the alarms with higher attribute similarity. In order to make it easier to understand our attack scenario division method, we give an algorithm flow to illustrate the construction of our method, as shown in Algorithm 1.

As described in Algorithm 1, Alert is an alarm dataset, and Scene is an attack scene set based on AP clustering partition. Firstly, the similarity matrix is obtained based on the above attribute similarity calculation method, and initialize the attraction matrix and the attribution matrix. Then, according to the requirements of the AP method, the similarity is taken as a negative value. Finally, the alarm data with high similarity is clustered into the same attack scenario according to the AP method.

*4.3. Constructing Attack Process Graph Based on Causality.* Every attack has its premise and corresponding consequences. That is, the previous attack is the precondition of the next attack, and the next attack is the consequence of the previous attack. For example, in a multistep attack, before launching an attack on the target, the intruder first scans and detects the target, finds the vulnerabilities in the target, and then starts the attack based on the vulnerabilities. Each of these attack steps can be regarded as a prerequisite for the next attack step, and the next attack step can be regarded as the consequence of the previous attack. Therefore, a complete attack sequence can be obtained by connecting the premise and result of alarm according to causality, which is based on causality. The method of dividing attack scenes was introduced, and the attack scenes were divided. Based on the previous work, this part will analyse the alarm data of the same attack scene by causal association method [38] and then construct an attack graph. The flow chart is shown in Figure 2.

In the multistep attack, the attack with correlation occurred in a short period of time and the alarm data with causality existed in the order of time, and the attack premise and the attack result have corresponding relations in IP and port attributes, that is, the destination IP address and destination port number of the attack premise must be the same as the source IP address and source port number of the attack result. The specific implementation process of association is as follows:

*Step 1.* Read the alarm dataset after the cluster processing sequentially and conduct correlation analysis of the data in each attack scene in turn.

*Step 2.* According to the idea of causality, sort the alarm data in each attack scenario in chronological order.

```
Input: Alarm dataset Alert = {a_1, a_2,···,a_n}.
Output: Attack scenario set Scene = {scene_1, scene_2,···,scene_n}.
1    Calculate the similarity matrix →Similarity = [sim_11, sim_12,···sim_nn].
2    Calculate the responsibility matrix→R = [r_11, r_12,···r_nn].
3    Calculate the availability matrix→A = [a_11, a_12 ··· a_nn].
4    Update R matrix and A matrix iteratively
5       if Convergence(cluster)
6          output cluster
7       else
8          return 4
9       end if
10   return Scene = {scene_1, scene_2,···,scene_n}.
```

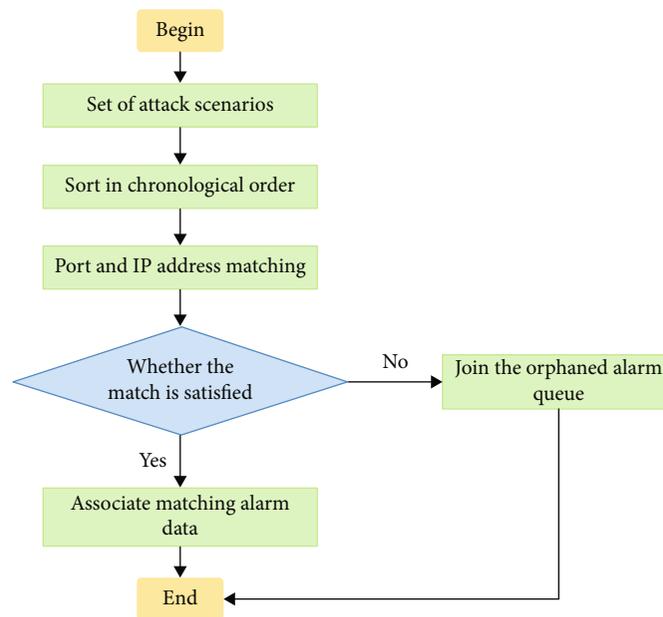ALGORITHM 1: Attack scenario division based on AP clustering.



FIGURE 2: Flow chart of alarm association based on causal relationship.

*Step 3.* Match the first piece of data $alert_1$ in the set with other data one by one according to the step size of 1. Within a certain time threshold, if the corresponding relationship between IP and port attributes is satisfied, that is, if the target IP address of the first data is the same as the source IP address of an alarm data, and the source port number and destination port number are the same, then, the two pieces of alert data are associated. If no qualified data is found after searching, it will be listed $alert_1$ as an isolated alarm.

*Step 4.* Sequentially execute the operations in step 3 on the remaining alarm data until all alarms are analysed.

*Step 5.* After completing the above steps to obtain the associated alarm, use Graph-viz to draw the attack graph.

## 5. Performance Evaluation

*5.1. Experimental Environment.* We use the Python 3.6 programming with PyCharm Community 2017.3 version in Windows 10. Use the Scikit-learn library to simply and efficiently process alarm data files. After the alarm association, we use the Graph-viz drawing tool to visually display the attack scene in the form of an attack graph.

*5.2. Experimental Dataset.* We use the honeypot dataset to verify the effectiveness of our proposed method in alarm correlation and the ability to construct attack scenarios. The honeypot dataset is obtained by simulating real network system and then using network decoy technology to lure intruders to launch attacks and capture the attack data [39]. Honeypot is essentially a kind of intelligence gathering system to trap attackers. All the actions of accessing honeypot system are the attacks of intruders. Through the correlation analysis of the honeypot data, all the activity information of the intruder in the system can be restored, which is convenient for the security management personnel to analyse the attacker's data and take corresponding measures to improve the protection capability of the real network system. These types of information include operating systems, brute force

network attacks, host vulnerabilities, and port scanning. In one aspect, the types of attacks included are shown in Table 2.

*5.3. Experimental Results and Analysis.* In this section, we divide the verification experiment into two parts: construction of attack graphs and correlation efficiency analysis.

*5.3.1. Build Attack Graphs.* After the attack scenarios are divided by the AP clustering algorithm, the alarm data in each attack scenario is used to find out the correlation between the alarms according to the causal relationship. If the following relationships are satisfied, it indicates that there is a connection between the two alarms. If the following relationship is satisfied, it is an isolated alarm.

(1) Within a certain period of time, the occurrence time of $alert_i$ precedes the occurrence time of $alert_j$

(2) $alert_i$'s destination IP is the same as $alert_j$'s source IP address

(3) $alert_i$'s destination port number is the same as $alert_j$'s source port number

After obtaining the associated alarm data based on the idea of causal association, we use the drawing software Graph-viz to construct an attack graph on these alarm data. Below we have selected several representative attack graphs for analysis.

As shown in Figure 3, we restored a distributed attack. The attack figure describes the process of a target host being attacked by multiple hackers. Multiple intruders perform SYN scan or FIN scan on the target host within the same time period, obtain active port information through the returned message, and then capture the host and obtain advanced permissions. Finally, using the host as a springboard, launch different types of distributed attacks on different hosts in the network.

As shown in Figure 4, it depicts an attack source launching the same type of attack on multiple target hosts at the same time. In these target hosts, a pair of pairwise combination is used to launch a centralized attack on the same host, and then a single step attack is implemented.

As shown in Figure 5, we restore a distributed port attack process. The attack source launches distributed attacks on the same port of different target hosts, controls these puppet machines through remote login, and uses the current host as the host to search for the target port to initiate local or remote attacks. Finally, use buffer overflow attacks to destroy the target host.

*5.3.2. Correlation Efficiency Analysis.* The correlation ratio and false alarm rate are reasonable indicators to verify the validity of alarm correlation. The false alarm rate refers to the ratio of false alarms that are not generated by real attacks to the total number of alarms. The correlation ratio refers to the ratio of the number of alarms generated by real attacks and the number of correctly associated alarms to

TABLE 2: Types of honeypot data attacks.

| Attack type | Quantity |
| --- | --- |
| Portmap-request-mountd | 111 |
| Web-cgi | 10 |
| Ping zeros | 51 |
| SYN FIN scan | 47 |
| DNS-version-query | 116 |
| DNS-zone-transfer | 3989 |
| Large-icmp | 286 |
| Ping Microsoft Windows | 14 |
| RPC-rpcinfo-query | 24 |
| Spp_portscan | 838 |
| SourcePortTraffic-53-tcp | 26 |
| Ping Nmap 2.36BETA | 459 |
| Socks-probe | 2627 |
| Telnet-login-incorrect | 397 |
| PING-ICMP time exceeded | 12 |
| IDS118-MISC-traceroute ICMP | 2360 |
| PING-ICMP destination unreachable | 709 |
| IDS212–MISC | 1487 |
| NAMED Iquery probe | 146 |
| RPC-portmap-request-status | 67 |
| MISC-Source Port Traffic 53 TCP | 60 |
| SMTP-expn-root | 786 |
| Portmap-request-mountd | 111 |

the total number of alarms. The calculation formulas are as follows:

$$FAR = \frac{NIA}{TNA} \times 100\%, \quad (10)$$

where FAR represents the false alarm rate, NIA represents the number of false alarms in isolated alarms, that is, the number of false alarms that did not participate in the correlation, and TNA represents the total number of alarms.

$$CR = \frac{NPA}{TNA} \times 100\%, \quad (11)$$

where CR represents the correlation ratio, NPA represents the number of alarms correctly participating in the association, and TNA represents the total number of alarms.

This paper selects two different alarm correlation analysis methods proposed in literature [40] and literature [41] to compare with the method proposed in this paper. Among them, literature [40] and literature [41] use a single alarm correlation method. It can be seen from Table 3 that our multitype mixed alarm correlation method is the method with the best correlation effect, and the correlation ratio reaches 96.7%, which is higher than the single correlation method. In addition, associating alarm data in attack scenarios based on AP clustering can find out more internal logical connections between alarms and reduce isolated alarms. The false
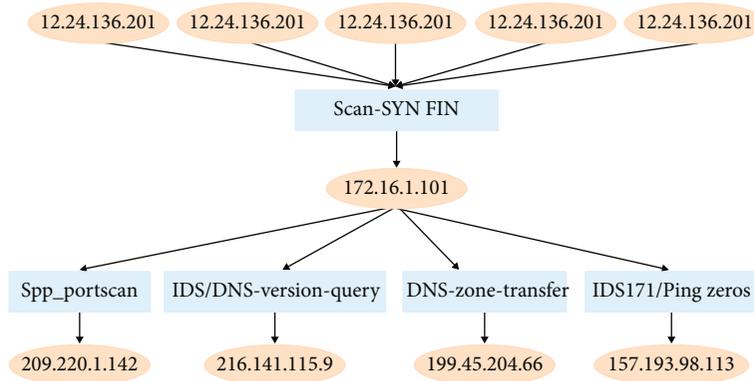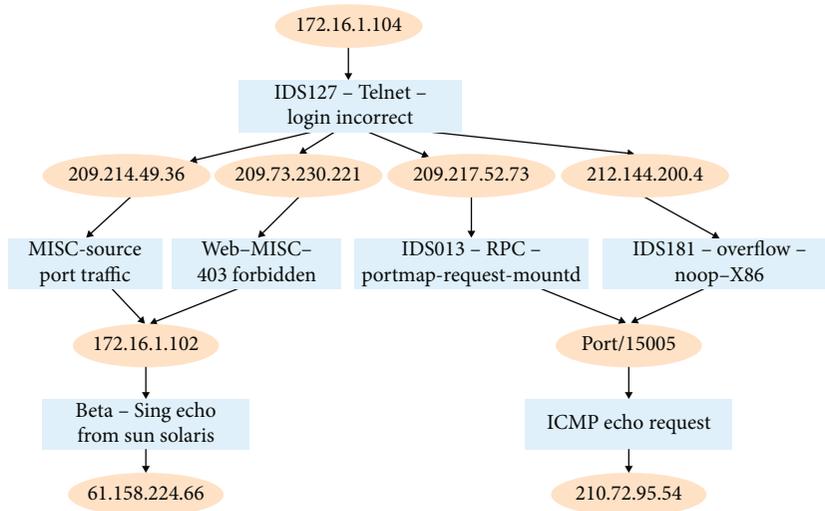
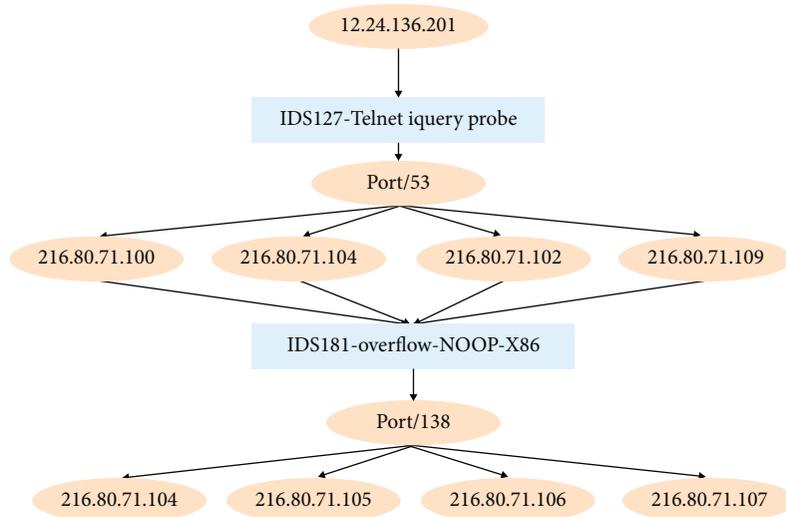FIGURE 3: Attack figure 1.



FIGURE 4: Attack figure 2.



FIGURE 5: Attack figure 3.

TABLE 3: Comparison of correlation ratio and false alarm rate.

| Correlation analysis method | FAR | CR |
| --- | --- | --- |
| The method presented in this paper | 2.1% | 96.7% |
| Method in literature [40] | 10.7% | 83.6% |
| Method in literature [41] | 4.5% | 93.2% |

alarm rate is only 2.1%, which is much smaller than other comparison algorithms. It shows that our method can effectively find out the correlation between alarm data and restore the complete attack scenario.

## 6. Conclusions

In this paper, we propose an alarm correlation method based on AP clustering algorithm and causality. Our method fully considers the logical relationship of each alarm information in the relevant attributes which analyses the characteristics of multistep attack alarm information and combines the shortcomings of existing alarm correlation methods to propose an attack scenario division method based on AP clustering. The experiment result showed that our method can achieve a correlation efficiency with 96.7% and can fully restore the attack process and construct a complete attack graph.

## Data Availability

The data used to support the findings of this study are available from the corresponding author upon request.

## Conflicts of Interest

The authors declare that there are no conflicts of interest regarding the publication of this study.

## Acknowledgments

## References

[1] Y. B. Wang, "Opportunities and challenges facing the security development of the Internet of Things," *Information Security and Communication Confidentiality*, vol. 39, no. 6, pp. 7–12, 2017.

[2] X. Cheng, J. L. Zhang, and B. Chen, "Cyber situation comprehension for IoT systems based on APT alerts and logs correlation," *Sensors*, vol. 19, no. 18, p. 4045, 2019.

[3] Z. Cai and Z. He, "Trading private range counting over big IoT data," in *2019 IEEE 39th International Conference on Distributed Computing Systems (ICDCS)*, pp. 144–153, Dallas, TX, USA, 2019.

[4] CCID Consulting, "2019 China cybersecurity development white paper," *China Computer News*, p. 6, 2019.

[5] Z. P. Cai and X. Zheng, "A private and efficient mechanism for data uploading in smart cyber-physical systems," *IEEE Transactions on Network Science and Engineering*, vol. 7, no. 2, pp. 766–775, 2020.

[6] X. Zheng and Z. P. Cai, "Privacy-preserved data sharing towards multiple parties in industrial IoTs," *IEEE Journal on Selected Areas in Communications*, vol. 38, no. 5, pp. 968–979, 2020.

[7] X. Cheng, J. L. Zhang, Y. F. Tu, and B. Chen, "Cyber situation perception for Internet of Things systems based onzero-dayattack activities recognition within advanced persistent threat," *Concurrency and Computation: Practice and Experience*, no. e6001, 2020.

[8] X. L. Tao, Y. Peng, F. Zhao, P. Zhao, and Y. Wang, "A parallel algorithm for network traffic anomaly detection based on isolation forest," *International Journal of Distributed Sensor Networks*, vol. 14, no. 11, 2018.

[9] Y. C. Yang, L. F. Wu, G. S. Yin, L. J. Li, and H. Zhao, "A survey on security and privacy issues in Internet-of-Things," *IEEE Internet of Things Journal*, vol. 4, no. 5, pp. 1250–1258, 2017.

[10] F. Valeur, G. Vigna, C. Kruegel, and R. Kemmerer, "Comprehensive approach to intrusion detection alert correlation," *IEEE Transactions on Dependable and Secure Computing*, vol. 1, no. 3, pp. 146–169, 2004.

[11] X. Qin and W. Lee, "Statistical causality analysis of INFOSEC alert data," in *International Symposium on Research in Attacks, Intrusions, and Defenses (RAID)*, pp. 73–93, Pittsburgh, USA, 2003.

[12] J. Zhang, X. P. Li, H. J. Wang, J. Q. Li, and B. Yu, "A real - time alarm correlation method based on attack plan diagram," *Computer Application*, vol. 36, no. 6, pp. 1538–1543, 2016.

[13] Z. Li, J. Lei, L. Wang, and D. Li, "A data mining approach to generating network attack graph for intrusion prediction," in *2007 4th International Conference on Fuzzy Systems and Knowledge Discovery (FSKD)*, pp. 307–311, Haikou, China, 2007.

[14] A. Ramaki, M. Amini, and R. Ebrahimi Atani, "RTECA: real time episode correlation algorithm for multi-step attack scenarios detection," *Computers & Security*, vol. 49, pp. 206–219, 2015.

[15] H. S. Gao and Y. M. Li, "An association analysis method of ASON alarm based on hierarchical attribute similarity clustering," *Science and Technology and Engineering*, vol. 15, no. 6, pp. 210–214+225, 2015.

[16] D. P. Hostiadi, M. D. Susila, and R. R. Huizen, "A new alert correlation model based on similarity approach," in *2019 1st International Conference on Cybernetics and Intelligent System (ICORIS)*, pp. 133–137, Denpasar, Indonesia, 2019.

[17] X. Tao, Y. Peng, F. Zhao, S. F. Wang, and Z. Liu, "An improved parallel network traffic anomaly detection method based on bagging and GRU," in *2020 15th International Conference on Wireless Algorithms, Systems, and Applications (WASA)*, pp. 420–431, Qingdao, China, 2020.

[18] J. Navarro, A. Deruyver, and P. Parrend, "A systematic survey on multi-step attack detection," *Computers & Security*, vol. 76, pp. 214–249, 2018.

[19] X. Fu and L. Xie, "Research on security alarm association technology," *Computer Science*, vol. 37, no. 5, pp. 9–14+29, 2010.

[20] L. Cheng, Y. Wang, and X. K. Ma, "GSLAC: A general scalable and low-overhead alert correlation method," in *2016 IEEE Trustcom/BigDataSE/ISPA*, pp. 316–323, Tianjin, China, 2016.

[21] A. Ahmadian Ramaki and A. Rasoolzadegan, "Causal knowledge analysis for detecting and modeling multi-step attacks,"

*Security and Communication Networks*, vol. 9, no. 18, pp. 6042–6065, 2016.

[22] S. Haas and M. Fischer, "GAC: graph-based alert correlation for the detection of distributed multi-step attacks," in *SAC 2018: Symposium on Applied Computing*, pp. 979–988, Pau, France, 2018.

[23] C. Y. Zhang and X. Wu, "Intrusion scenario dynamic correlation algorithm based on single value causality diagram," *Advanced Materials Research*, vol. 926-930, pp. 3063–3067, 2014.

[24] K. Y. Li, Y. Li, J. Y. Liu, R. Zhang, and X. Duan, "Attack pattern mining algorithm based on security log," in *2017 IEEE International Conference on Intelligence and Security Informatics (ISI)*, pp. 205–205, Beijing, China, 2017.

[25] F. Faraji Daneshgar and M. Abbaspour, "Extracting fuzzy attack patterns using an online fuzzy adaptive alert correlation framework," *Security and Communication Networks*, vol. 9, no. 14, pp. 2245–2260, 2016.

[26] X. W. Feng, D. X. Wang, M. H. Huang, and J. Li, "A method for mining causal knowledge based on Markov properties," *Computer Research and Development*, vol. 51, no. 11, pp. 2493–2504, 2014.

[27] M. Soleimani and A. Ghorbani, "Multi-layer episode filtering for the multi-step attack detection," *Computer Communications*, vol. 35, no. 11, pp. 1368–1379, 2012.

[28] Y. Zhang, S. Zhao, and J. Zhang, "RTMA: real time mining algorithm for multi-step attack scenarios reconstruction," in *2019 IEEE 21st International Conference on High Performance Computing and Communications; IEEE 17th International Conference on Smart City; IEEE 5th International Conference on Data Science and Systems (HPCC/SmartCity/DSS)*, pp. 2103–2110, Zhangjiajie, China, 2019.

[29] M. Barzegar and M. Shajari, "Attack scenario reconstruction using intrusion semantics," *Expert Systems with Applications*, vol. 108, pp. 119–133, 2018.

[30] J.-w. Tian, X. Li, Z. Tian, and W.-h. Qi, "Network attack path reconstruction based on similarity computation," in *2017 13th International Conference on Natural Computation, Fuzzy Systems and Knowledge Discovery (ICNC-FSKD)*, pp. 2457–2461, Guilin, China, 2017.

[31] B. Yang, J. J. Li, C. Qi, H. G. Li, and Y. He, "Novel correlation analysis of alarms based on block matching similarities," *Industrial & Engineering Chemistry Research*, vol. 58, no. 22, pp. 9465–9472, 2019.

[32] B. Frey and D. Dueck, "Clustering by passing messages between data points," *Science*, vol. 315, no. 5814, pp. 972–976, 2007.

[33] N. H. Yang, H. Q. Yu, Z. L. Qian, and H. Sun, "Modeling and quantitatively predicting software security based on stochastic Petri nets," *Mathematical and Computer Modelling*, vol. 55, no. 1-2, pp. 102–112, 2012.

[34] S. H. Ahmadinejad, S. Jalili, and M. Abadi, "A hybrid model for correlating alerts of known and unknown attack scenarios and updating attack graphs," *Computer Networks*, vol. 55, no. 9, pp. 2221–2240, 2011.

[35] L. Q. Zhou and W. X. Wei, "Intrusion detection method based on principal component analysis and Simhash," *Computer and Digital Engineering*, vol. 43, no. 7, pp. 1291–1294, 2015.

[36] A. Baláž, N. Ádám, E. Pietriková, and B. Madoš, "ModSecurity IDMEF module," in *2018 IEEE 16th World Symposium on Applied Machine Intelligence and Informatics (SAMI)*, pp. 43–48, Kosice and Herlany, Slovakia, 2018.

[37] X. L. Tao, Y. M. Gong, and F. Zhao, "An OSSEC alarm data aggregation method based on classification," *Computer Engineering and Design*, vol. 41, no. 4, pp. 908–914, 2020.

[38] P. Ning, C. Yun, and D. Reeves, "Analyzing intensive intrusion alerts via correlation," in *International Symposium on Research in Attacks, Intrusions, and Defenses (RAID)*, pp. 74–94, Zurich, Switzerland, 2002.

[39] D. Q. Yang, W. M. Liu, and Z. Yu, "Research on active defence application based on honeypot," *Journal of Network and Information Security*, vol. 4, no. 1, p. 57, 2018.

[40] S. Wang, G. M. Tang, J. H. Wang, Y. F. Sun, and G. Kou, "Construction method of attack scenario based on causal knowledge network," *Computer Research and Development*, vol. 55, no. 12, pp. 2620–2636, 2018.

[41] C. T. Kawakani, S. B. Junior, and R. S. Miani, "Intrusion alert correlation to support security management," in *Brazilian Symposium on Information Systems (SBSI)*, pp. 313–320, Florianópolis, Brazil, 2016.