# A Predictive Model of Insider Threat Based on Bayesian Network

Hui Wang[1,2], Yunfeng Wang[1] and Guangcan Yang[1]

[1] Henan Polytechnic University, Jiaozuo, China
[2] Jilin University, Changchun, China

*Abstract*—At present, development of science and technology accelerates the society-informationization, many enterprises follow the trend of era to build internal network for convenient communication, but the increasing network security incidents cause a new understanding about the importance of internal network. The predictive model of insider threat based on Bayesian network is put forward in this paper. In the model, insider behaviors in the process of operation are considered as research objects, resource and intrusion evidence for operation sequence are seen as nodes, and then the network attack graph of Bayesian network is established. The concept of meta-operation, atomic attack and intrusion evidence are put forward in the graph. The node variable, its value and the conditional probability distribution of network attack graph are defined. Based on Bayesian network approximate inference, the improved likelihood weighted algorithm is presented to calculate the parameter and to quantify the insider threat. According to the simulation experiment data analysis, this approach is fast, simple and accurate, and plays an effective role in the process of insider threat prediction and evaluation.

*Index Terms*—Insider Threat, Bayesian Network, Network Attack Graph, Likelihood Weighted Algorithm

## I. INTRODUCTION

With development of science and universality of network, the network security incidents are increasing, these incidents occurred not only in external network of enterprises, but also in internal network. In fact, the danger of insider threat is far greater than external threat. At present, there are many defense and testing tools, and the technical level is higher. Unfortunately, they are mainly used to solve the external threat, the technology for internal network security protection yet to be developed and valued.

Bayesian network is a graphics mode, used to describe the dependent relationship among random variables, and applied to solve uncertainty problems. It has been widely applied on analysis, forecast and defense. Considering the uncertainty, vulnerability and complexity of insider threat, the insider threat prediction model based on Bayesian network is put forward in this paper, graphics mode is used to describe the various states in internal network and to form the network attack graph, and Bayesian network inference algorithm is used to calculate the risk probability.

## II. RELATED RESEARCH

In [1], the authors for the first time used Bayesian network model to evaluate information risk, they adopted the varying probability based on the security state among hosts, but did not measure the security of information system by essence of risk, so their model has a poor expandability.

TVA was used by the authors to model Bayesian network for risk assessment in [2], this modeling method only increased the speed of modeling, unfortunately, how to determine the conditional probability distribution of the model was not mentioned.

In [3], the authors used Bayesian network model to calculate the risk probability of information security. They took the factors of risk probability into planning penetration diagram, with planning penetration diagram as network structure of Bayesian network model, the value of the parameter was obtained by the expert knowledge and the distribution of the maximum entropy. Both the modeling speed and the generating parameter were improved.

Network attack graph was used to describe all the paths that attacker chooses for target, and was took for the safety analysis of system in [4]. There are many shortages by using the method in [4] for calculation, for example, it is complex to calculate the confidence of attack graph nodes, and the support of mathematical theory is scarce.

Combining with previous researches, insider threat is a serious problem in network security, and difficult to resist and management[5, 6]. Network attack graph is an important tool for safety analysis recent years[7], and the technology of generating network attack graph automatically is the hot spot that experts research at home and abroad, so far, it has already achieved good results. Bayesian network is a mature theoretical model to solve uncertainty risk assessment and prediction, and a kind of strict data language with perfect modeling method and inference algorithm [8]. In the aspect of modeling, its language is rigorous, and calculation is accurate. Therefore, the predictive model of insider threat based on Bayesian network is put forward, with network attack graph as model structure of Bayesian network, and the parameters of predictive model of insider threat are calculated by the improved likelihood weighted algorithm. In the process of modeling, the concept of meta-operation, intrusion evidence and atomic attack are put forward, so that the method of modeling is fast, efficient and applicable.

## III. THE PREDICTIVE MODEL OF INSIDER THREAT BASED ON BAYESIAN NETWORK

From the definition of Bayesian network, the model construction based on Bayesian network has two aspects. In qualitative aspect, it uses a directed acyclic graph to depict interdependent and independent relations among different variable nodes, which is the network structure of model. In quantitative aspect, it uses conditional probability distribution to describe the dependent relation from the child node to his father, which is the parameter of model.

### A. The Definition of Meta-operation

Network attack is made up of different sets of commands or operations, and the set of different commands or operations formed by some kind of methods is called meta-operation [9].

### B. The Definition of Atomic Attack

The plan of SPRINT (Signature Powered Revised Instruction Table) was defined in literature [10, 11], which forecasted insider threat based on the using proposes of users. Before users using system, they should submit their intents according to the form of list < theme, motion, object, period >. The legal set of meta-operation—Mos (material operate set) in internal network is the mapping on the attributes <motion, object> of SPRINT plan which users submitted.

In the process of attacking, the series of operations between two resources (including all kinds of orders and operations) correspond to a subset of Mos—Sub-Mos, and called atomic attack. Atomic attack is the minimum set of meta-operation that the attacker needed from one resource to next resource.

### C. The Definition of Intrusion Evidence

In order to describe the series of operations among attack resources, the concept of intrusion evidence is introduced. Intrusion evidence is the set of the series of operations recorded the attacker from a resource to another resource by log. Intrusion evidence is also a set composed by lots of meta-operations with certain order. Different meta-operation and different order make different intrusion evidence.

The confidence of intrusion evidence is used to measure the size of intrusion evidence. The confidence of intrusion evidence is a probability, which is the set of meta-operation that intrusion evidence covers atomic attack. Here, cover means according to the sequence of meta-operation in set of atomic attack, the meta-operation number in set of intrusion evidence and set of atomic attack is same. Suppose $m_i$ is the set of meta-operation, the set of atomic attack is $\{m_1, m_2, m_3, m_4\}$ and the order is $m_1 \rightarrow m_2 \rightarrow m_3 \rightarrow m_4$; In this way, the cover of the intrusion evidence $\{m_1, m_5, m_6, m_2, m_7\}$ is 2, and the confidence is 50%; the cover of the intrusion evidence $\{m_5, m_6, m_7, m_8\}$ is 0, and the confidence is 0; the cover of the intrusion evidence $\{m_2, m_1, m_5, m_6\}$ is 1, and the confidence is 25%.

For ease of explanation, assuming that the meta-operation number is n in atomic attack, when the cover number of intrusion evidence is less than n, the case is called insufficient evidence; when the cover number of intrusion evidence is equal to n, this case is called sufficient evidence.

Figure 1 is an example to show the relationship among meta-operation, atomic attack and intrusion evidence. $v_1$ and $v_2$ are two resource nodes, and $o_1$ and $o_2$ are two intrusion evidences, the relationship is shown as figure 1: when $V_1$ is occupied, the attacker can occupy $V_2$ via intrusion evidence $o_1$ or $o_2$. Suppose the set of meta-operation M= $\{m_i \mid i=1, 2, 3……, n\}$, and if the attacker wants to occupy resource node v2 from resource node $v_1$, there must be four meta-operations that $m_2 \rightarrow m_4 \rightarrow m_5 \rightarrow m_7$, arrow show the sequence of four sets of meta-operation, then $\{m_2, m_4, m_5, m_7\}$ is atomic attack. When the meta-operation number is less than four or the sequence of these four sets of meta-operation is different, they can not be atomic attack. Intrusion evidence is the combination of different meta-operations, when intrusion evidence $o_1 = \{m_1, m_2, m_3, m_4, m_6\}$, the cover of intrusion evidence is 2, it is called insufficient evidence, and can not achieve attack; When intrusion evidence $o_1 = \{m_1, m_2, m_3, m_4, m_5, m_6, m_7\}$, and the set of intrusion evidence $\supseteq$ the set of atomic attack, so the cover is 4, it is called sufficient evidence, and the resource node $v_2$ can be achieved successfully.
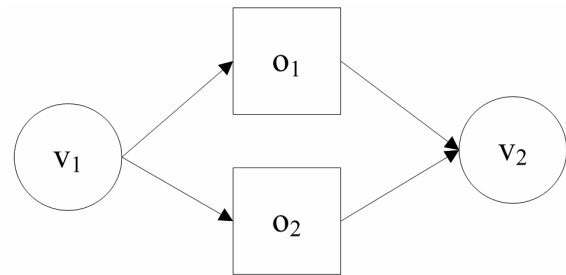


Figure 1.   The relationship among meta-operation.

### D. The Definition of Network Attack Graph

In order to describe the different attack paths and attack characteristics in attacking process, network attack graph (NAG) is adopted to show the mutual relationship of different attack position in attacking path. The definition of network attack graph is as follow: NAG=（V，V0，G，O，E，P).

1. V is the set of resource nodes, and V=$\{v_i \mid i=1,2……, N\}$, $v_i$ is a single resource node, used to describe the resource that the attacker occupied in attacking process, and its value is True or False. Whether the value of $v_i$ is True shows the attacker occupy the resource successfully or not.

2. $V_0$ is the set of resources that the attacker has occupied in initial state, and also a subset of V. In network attack graph, it is the set of initial nodes.

3. G is the set of target nodes, means the set of target nodes which the attacker wants to achieve finally.

4. O is the set of intrusion evidence nodes, and means the set of series of operations when the attacker has occupied some resources. O = $\{o_i \mid i = 1, 2..., n\}$, $o_i$ is a single intrusion evidence node, and the value of it is the confidence of intrusion evidence, range from [0, 1]. When the confidence of intrusion evidence is less than 100%, it shows the attacker is insufficient evidence, and the next resource node is not occupied; when the confidence of intrusion evidence is 100%, it shows the attacker is sufficient evidence, and the next resource node is occupied.

5. E is the set of directed edges associate with all kinds of nodes. E = (E₁UE₂), in which $E_1 \subseteq V \times O$, $E_2 \subseteq O \times V$. $E_1$ represents the condition only some resources the attacker has occupied can the attack evidence be triggered; $E_2$ shows the attacker can have some resources when attack evidence occurring.

6. P is the conditional probability distribution of each node in network attack graph. $P = (P_1 U P_2)$, $P_1$ is the conditional probability distribution of intrusion evidence nodes, and $P_2$ is the condition probability distribution of resource nodes. Generally, Pre(x) is the father node set of nodes X, and Con(x) is the child node set of nodes X. For the completeness of network attack graph, the relationships of AND and OR among different nodes should be considered, there are different relationships of AND and OR among different nodes Pre(x), and different relationships of AND and OR among different nodes Pre($v_i$). The specific details are described as below:

① The relationship of AND among different nodes Pre($o_i$): it shows only the attacker occupies different resources at the same time, the attacker could implement further attack, complete intrusion evidence, and occupy more resources;

② The relationship of OR among different nodes Pre($o_i$): it shows as long as the attacker occupies any resource, the attacker could implement further attack, complete intrusion evidence, and occupy more resources;

③ The relationship of AND among different nodes Pre($v_i$): it shows if the attacker wants to achieve next resource node, the attacker should complete all the intrusion evidences of each father node.

④ The relationship of OR among different nodes Pre($v_i$): it shows if the attacker completes intrusion evidence of any father node, the attacker can occupy the resource node.

According to definitions above, the network attack graph is completed, and is shown as Figure 2: directed edges represent the relationships between the resources the attacker has owned and the intrusion evidences, and the backbone structure of network attack graph is formed by considering all kinds of situations. $v_1$, $v_2$ and $v_3$ are the initial nodes, $v_7$ is a target node, there are three ways if the attacker wants to occupy $v_7$ (symbol $\wedge$ shows the relationship of AND among different nodes).
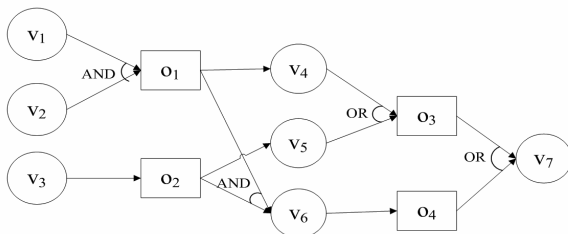


Figure 2. Network Attack Graph.

1. $(v_1 \wedge v_2) \rightarrow o_1 \rightarrow v_4 \rightarrow o_3 \rightarrow v_7$;

2. $v_3 \rightarrow o_2 \rightarrow v_5 \rightarrow o_3 \rightarrow v_7$;

3. $((v_1 \wedge v_2 \rightarrow o_1) \wedge (v_3 \rightarrow o_2)) \rightarrow v_6 \rightarrow o_4 \rightarrow v_7$;

For more explicit relationship among nodes in network attack graph, the line order relation among nodes in network attack graph is introduced. The line order relation

is the process of forming attack path, and described with STEP. Firstly, the partial order set of network attack graph should be found. Here, the concept of in-degree and out-degree are introduced. In directed graph, the number of edges that go toward the node is called the in-degree of the node, and the number of edges go out from the node is called the out-degree of the node. For example, in figure 4, the in-degree of $o_1$ is 2, out-degree is 2, and the in-degree of $v_1$ is 0, out-degree is 1. In the process of judging partial order relation among nodes, the choice of path would be made according to in-degree and out-degree of the node: firstly, finding out the node m with in-degree is 0, then cutting off the directed edges associate with the node m, if the directed edges associate with another node n, the partial order relation < m, n > between node m and node n is recorded. And like the method above, all the partial order relations among nodes in network attack graph can be found out ultimately. Partially ordered set (POS): dotted lines means the directed edges which will be cut off, and the specific process is shown as figure 3：
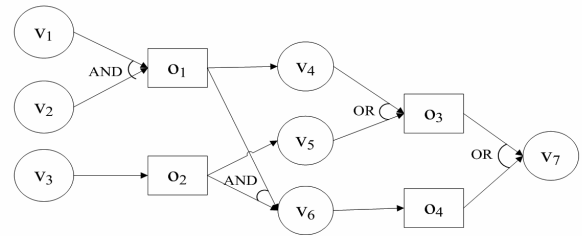


Figure 3. The process of finding POS

① Firstly, $v_1$, the edge $v_1 \rightarrow o_1$ should be cut off. Recording: POS= {< $v_1$, $o_1$>}.

② Secondly, $v_2$, the edge $v_2 \rightarrow o_1$ should be cut off. Recording: POS= {< $v_1$, $o_1$>, < $v_2$, $o_1$>}, $(v_1 \wedge v_2)$.

③ Thirdly, $v_3$, as above. Recording: POS= {< $v_1$, $o_1$>, < $v_2$, $o_1$>, < $v_3$, $o_2$>}, $(v_1 \wedge v_2)$.

④ The edges $o_1 \rightarrow v_4$ and $o_1 \rightarrow v_6$ should be cut off. Recording: POS= {< $v_1$, $o_1$>, < $v_2$, $o_1$>, < $v_3$, $o_2$>, < $o_1$, $v_4$>, < $o_1$, $v_6$>}, $(v_1 \wedge v_2)$, $(o_1 \wedge o_2)$.

⑤ The edges $o_2 \rightarrow v_5$ and $o_2 \rightarrow v_6$ should be cut off. Recording: POS= {< $v_1$, $o_1$>, < $v_2$, $o_1$>, < $v_3$, $o_2$>, < $o_1$, $v_4$>, < $o_1$, $v_6$>, < $o_2$, $v_5$>, < $o_2$, $v_6$>}, $(v_1 \wedge v_2)$, $(o_1 \wedge o_2)$.

⑥ The directed edge $v_4 \rightarrow o_3$ should be cut off, and the relationship between $v_4$ and $v_5$ is OR. Recording: POS= {< $v_1$, $o_1$>, < $v_2$, $o_1$>, < $v_3$, $o_2$>, < $o_1$, $v_4$>, < $o_1$, $v_6$>, < $o_2$, $v_5$>, < $o_2$, $v_6$>, < $v_4$, $o_3$>}, $(v_1 \wedge v_2)$, $(o_1 \wedge o_2)$.

⑦ Secondly, $v_5$ and $v_6$, the edge $v_5 \rightarrow o_3$ should be cut off, and the relationship between $v_4$ and $v_5$ is OR. Recording: POS= {< $v_1$, $o_1$>, < $v_2$, $o_1$>, < $v_3$, $o_2$>, < $o_1$, $v_4$>, < $o_1$, $v_6$>, < $o_2$, $v_5$>, < $o_2$, $v_6$>, < $v_4$, $o_3$>, < $v_5$, $o_3$>}, $(v_1 \wedge v_2)$, $(o_1 \wedge o_2)$.

⑧ Thirdly, $v_6$ and $o_3$, the edge $v_6 \rightarrow o_4$ should be cut off. Recording: POS= {< $v_1$, $o_1$>, < $v_2$, $o_1$>, < $v_3$, $o_2$>, < $o_1$, $v_4$>, < $o_1$, $v_6$>, < $o_2$, $v_5$>, < $o_2$, $v_6$>, < $v_4$, $o_3$>, < $v_5$, $o_3$>, < $v_6$, $o_4$>}, $(v_1 \wedge v_2)$, $(o_1 \wedge v_2)$.

⑨ Then, the in-degree of nodes $o_3$ and $o_4$ are 0. Firstly, the edge $o_3 \rightarrow v_7$ should be cut off. Recording: POS= {< $v_1$, $o_1$>, < $v_2$, $o_1$>, < $v_3$, $o_2$>, < $o_1$, $v_4$>, < $o_1$, $v_6$>, < $o_2$, $v_5$>, <

$o_2$, $v_6$>, < $v_4$, $o_3$>, < $v_5$, $o_3$>, < $v_6$, $o_4$>, < $o_3$, $v_7$>}, ($v_1 \wedge v_2$), ($o_1 \wedge o_2$).

⑩ Secondly, $o_4$, the edge $o_4 \rightarrow v_7$ should be cut off. Recording: POS= {< $v_1$, $o_1$>, < $v_2$, $o_1$>, < $v_3$, $o_2$>, < $o_1$, $v_4$>, < $o_1$, $v_6$>, < $o_2$, $v_5$>, < $o_2$, $v_6$>, < $v_4$, $o_3$>, < $v_5$, $o_3$>, < $v_6$, $o_4$>, < $o_3$, $v_7$>, < $o_4$, $v_7$>}, ($v_1 \wedge v_2$), ($o_1 \wedge o_2$).

According to the partial order set (POS), the partial order set among nodes without the relationship is AND can be obtained:

$POS_1$= {< $v_1$, $o_1$>, < $o_1$, $v_4$>, < $v_4$, $o_3$>, < $o_3$, $v_7$>};

$POS_2$= {< $v_2$, $o_1$>, < $o_1$, $v_4$>, < $v_4$, $o_3$>, < $o_3$, $v_7$>};

$POS_3$= {< $v_1$, $o_1$>, < $o_1$, $v_6$>, < $v_6$, $o_4$>, < $o_4$, $v_7$>};

$POS_4$= {< $v_2$, $o_1$>, < $o_1$, $v_6$>, < $v_6$, $o_4$>, < $o_4$, $v_7$>};

$POS_5$= {< $v_3$, $o_2$>, < $o_2$, $v_5$>, < $v_5$, $o_3$>, < $o_3$, $v_7$>};

$POS_6$= {< $v_3$, $o_2$>, < $o_2$, $v_6$>, < $v_6$, $o_4$>, <$o_4$, $v_7$>}.

Both the relationship of $v_1$ and $v_2$ and the relationship of $o_1$ and $o_2$ are AND in the graph. Here, the relationship of AND represents when both of the two nodes are ready, the partial order relation with the next node could be formed. So the improved partial order set forms the line sequence set, and the line sequence set of the network attack graph in this paper is shown as follows:

$STEP_1$= {< {$v_1$, $v_2$}, $o_1$>, < $o_1$, $v_4$>, < $v_4$, $o_3$>, < $o_3$, $v_7$>};

$STEP_2$= {< {< {$v_1$, $v_2$}, $o_1$>, <$v_3$, $o_2$>}, $v_6$>, <$v_6$, $o_4$>, < $o_4$, $v_7$>};

$STEP_3$= {< $v_3$, $o_2$>, < $o_2$, $v_5$>, < $v_5$, $o_3$>, < $o_3$, $v_7$>};

Compare $STEP_1$, $STEP_2$ and $STEP_3$ with three paths the network attack graph concluded above, it can be known that the line order relation and the actual attack path is consistent. At the same time, it proves that this method is correct and effective. The thought will be expressed as the core of the algorithm 2.

*E. The Calculation of Bayesian Network Model*

The specific operation is: sampling sequentially according to the partial order relation of nodes in network attack graph, in the process of sampling, if the sampled variable is an evidence variable which has already known, then sampling by the distribution of $p(X)$, and setting the sampling result as the observation of evidence variable; If the sampled variable is not an evidence variable, then sampling by the distribution of $p(X | \pi(X))$, and the sampling result is the value of the sample. Likelihood weighted algorithm does not waste any sample, and makes every sampling result of evidence variable is consistent with the value of evidence variable of the posterior probability distribution, so that every sample has a role to play effectively.

The formula of calculating posterior probability of likelihood weighted algorithm is:

P(Q=q/E=e)≈

$$\frac{\text{The sum of probabilities of } E = e \text{ in sampling results in the condition of } Q = q}{\text{The sum of probabilities of } E = e \text{ in sampling results}} \quad (1)$$

Due to the AND and OR relations between the nodes in Bayesian network attack graph, the topology order would be affected among nodes when calculating the likelihood weighted, so the improved algorithm 1 and 2 of likelihood weighted algorithm is put forward in this paper.

Algorithm 1: Likelihood Weighted (NAG, m, E, e, Q, q)

Based on Bayesian network attack graph and known evidence variable, the posterior probability distribution of query variable can be found out.

Input: Bayesian network attack graph (NAG), the sample m, the evidence variable E, the value of evidence variable e, the query variable Q, the value of query variable q.

Output: the approximation of $p(Q = q | E = e)$

1：A←STEP（NAG，D）；

2：$\omega_e$←0；$\omega_{q,e}$←0;

3：for (i=1 to m)

4：$D_i$←0；

5：for (each variable X of A)

6：if (X∈E)

7：x ← the observed value of X；

8：else

9：x ← the sampling result of P (X| $\pi$ (X))；

10：end if

11：end for

12：$D_i$ ← $D_i \cup$ {X=x}；

13：$\omega_i \leftarrow \prod_{X \in E} P(X | \pi(X)) | D_i$；

14：$\omega_e \leftarrow \omega_e + \omega_i$；

15：if (D is consistent with D=q)

16：$\omega_{q,e} \leftarrow \omega_{q,e} + \omega_i$；

17：end if

18：end for

19：return $\omega_{q,e} / \omega_e$

Algorithm 2: Line order relation STEP (NAG, SETP$_i$)

According to the Bayesian network attack graph (NAG), the node set of ordered attack path D can be found out.

Input: network attack graph (NAG), in-degree, out-degree, partial order relation (POS), line order relation STEP, arbitrary node X and Y, set M with the relationship of AND;

Output: the set of nodes with line order relation of Bayesian network;

1：POS ← $\varnothing$；POSi ← $\varnothing$；SETPi ← $\varnothing$；M← $\varnothing$；

2：for（each node variable X with $\lambda_1$ =0）；

3：All the nodes variable Y which has the directed edge with node X in NAG should be found out；

4：If (the out-degree of Y is greater than 1, and the two edges toward to Y is the relationship of AND.

5：M←M∪ {$X_1 \wedge X_2$};

6：end if；

7：POS ← POS∪ {<X, Y>}

8：$\lambda_2$ (the out-degree of X) ← 0；

9：end for；

10：In the set of POS, all the partial order set POS$_i$ regardless of the relationship of AND should be found out;

11：In the set of POS$_i$, all the line sequence set SETP$_i$ with the relationship of AND should be found out;

12：return SETP$_i$；

IV. THE EXPERIMENTAL DATA AND ANALYSIS

The experimental environment is an internal network composed by a small LAN with Windows system. For figure 3, the variables are sampled by random number generator. For example: the random variable X which obeys the Bernoulli distribution, and its value is 1 or 0. Assuming that P(X=0) =p, P(X=1) =1-p. The process of sampling P(X) is as follows: The real number x which belongs to [0, 1] generated by random number generator,

If the value of x belongs to [0, p], the result of sample is X=0; otherwise, X=1.

The sample is analyzed according to the different sampling time. The number of samples in each group is 30000, and Ni is the sampling instant. The probability of $p(v_7 | v_1, v_2, v_3)$ is the statistics of sampling result, and p= $p(v_7 = true | v_1 = true, v_2 = true, v_3 = true)$ is the searching probability. There are four sampling processes in this experiment, $p_i$ is the combination of probability in each sampling, and the sampling results are shown as table 1.

Table 1. Sampling Results of Likelihood Weighted Algorithm

|    | N1     | N2     | N3     | N4     | N5     | N6     | N7     | N8     | N9     | N10    |
|----|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|
| P1 | 0.2789 | 0.3331 | 0.4965 | 0.5425 | 0.6061 | 0.6778 | 0.7475 | 0.7982 | 0.8536 | 0.9527 |
| P2 | 0.1927 | 0.265  | 0.3024 | 0.4521 | 0.5347 | 0.6024 | 0.7189 | 0.8034 | 0.8156 | 0.83   |
| P3 | 0.1249 | 0.3024 | 0.3654 | 0.4536 | 0.5367 | 0.6024 | 0.6789 | 0.7246 | 0.79   | 0.82   |
| P4 | 0.25   | 0.3547 | 0.4821 | 0.5278 | 0.6196 | 0.7056 | 0.7598 | 0.8236 | 0.8956 | 0.92   |

By contrasting the sampling results, we can see that the attack probability of attackers is increased gradually with the different sampling time, and it is consistent with the actual attack. With the advancement of attacking progress, the resources that attackers own are increasingly, and the targets are getting closer, so the possibility of threat is increasing for system.

The attack processes based on the predictive model of insider threat are shown as figure 4 and figure 5. Assuming the alarm value set by manager is 0.85, when the attack probability reaches 0.85, the attack behaviors will be prevented. As can be seen from figure 4 and figure 5, this model could detect attack behaviors of insiders in real time, and quantify the attack probability by posterior probability. Thus, it can provide references to managers, and help them make right decision.

## V. CONCLUSION

Insider threat is one of the most serious challenges that enterprises faced with at present. Insiders have far greater privileges than external attackers, so they can easily escape from the alarm equipments of internal network, bypass various firewall, intrusion detection and access control, and finally implement attack behaviors successfully without being found. The research of insider threat is still at the primary stage at home and abroad.

Bayesian network is mainly to solve the uncertainty problems, and the factors of insider threat are complex, vulnerable and uncertain. So, the network structure of insider threat is described with Bayesian network model in this paper, and the likelihood weighted algorithm is used to calculate the posterior probability of network attack graph. Thus, the network structure of insider is analyzed qualitatively and quantitatively. Lastly, proved by experimental results, the model put forward in this paper can predict insider threat accurately and effectively.
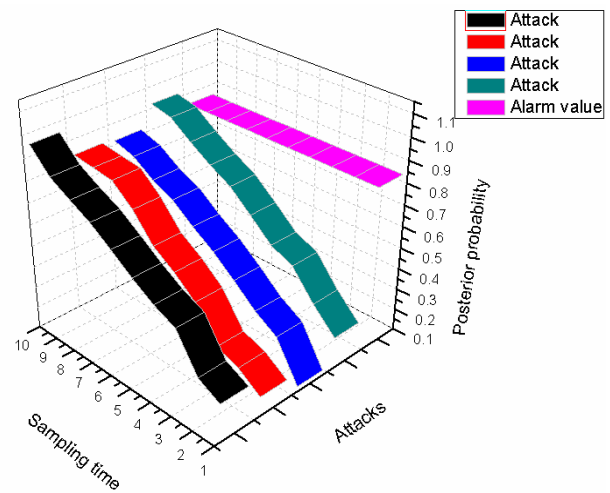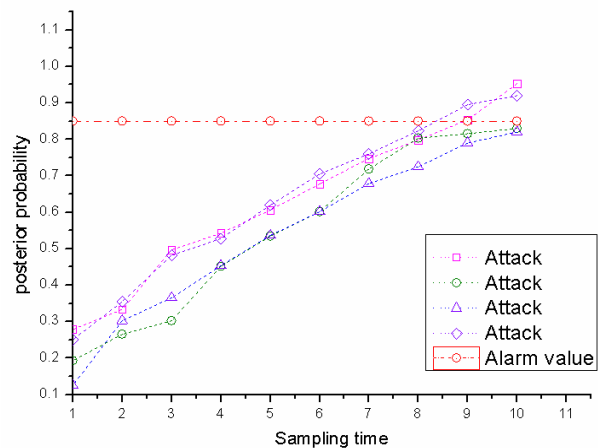
## ACKNOWLEDGMENT

Figure 4. Attack Processes



Figure 5. Attack Processes

## REFERENCES

[1] Y. liu and H. Man, Network vulnerability assessment using Bayesian networks.SPIE'05.Orlando FL USA：ACM Press, 2002, pp. 217-224.

[2] Marcel F., etc., Measuring network security by using attack graphs of Bayesian network-based. IEEE International Computer Software and Applications Conference. Washington USA, 2011, pp. 698-703.

[3] Wang Z.Z., etc., The calculation of risk probability in information security based on Bayesian network model . Journal of Electronic, 2011 (S1).

[4] Zhang ShaoJun, etc, The calculation of confidence based on Bayesian inference in attack graph [J]. Journal of Software, September 2010.

[5] Brancik and Kenneth, The optimization of situational awareness for insider threat detection. CODASPY'11, 2011, pp. 231-235.

[6] Nelli and Suraj, Role-based differentiation for insider detection algorithms. Proceedings of the ACM Conference on Computer and Communications Security, 2010, pp.55-62.

[7] Wang GuoYu, etc, The modeling method based on network attack graph. Journal of National University of Defense Technology (NUDT), April 2009.

[8] Yang Jian, etc, A method of threat assessment based on Bayesian network . Journal of PLA University of Science and Technology, Nov 2010.

[9] Wang Hui and Liu ShuFen, An extensible model for the prediction of insider threat . Chinese Journal of Computers, August 2006.

AUTHORS

**Hui Wang** is with the College of Computer Science and Technology, Henan Polytechnic University, Jiaozuo, China (e-mail: wanghui_jsj@hpu.edu.cn).

**Yunfeng Wang** is with the College of Computer Science and Technology, Henan Polytechnic University, Jiaozuo, China (e-mail: wanghwyh@yahoo.com).

**Guangcan Yang** is with the College of Computer Science and Technology, Henan Polytechnic University, Jiaozuo, China (e-mail: jsjyjs@hpu.edu.cn).