

基于动态马尔科夫模型的入侵检测技术研究

尹清波,张汝波,李雪耀,王慧强

(哈尔滨工程大学计算机科学与技术学院,黑龙江哈尔滨 150001)

摘 要: 本文提出了基于动态马尔科夫模型的入侵检测方法. 首先提取特权进程的行为特征,并在此基础上动态构造 Markov 模型. 由动态 Markov 模型产生的状态序列计算状态概率,根据状态序列概率来评价进程行为的异常情况. 利用 Markov 模型的动态构造充分提取特权进程的局部行为特征的相互关系,因此可以在训练数据集有限的条件下使模型更精确、检测能力大大加强. 实验表明该算法准确率高、实时性强、占用系统资源少. 本文所提方法算法简单、预测准确,适合于进行实时检测.

关键词: 入侵检测; 动态马尔科夫模型; 信息安全

中图分类号: TP309 **文献标识码:** A **文章编号:** 0372-2112 (2004) 11-1785-04

Research on Technology of Intrusion Detection Based on Dynamic Markov Model

YIN Qing-bo, ZHANG RU-bo, LI Xue-yao, WANG Hui-qiang

(College Of Computer Science And Technology, Harbin Engineering University, Harbin, Heilongjiang 150001, China)

Abstract: A new method for anomaly intrusion detection is proposed based on dynamic Markov model. At first, behavioral features are extracted from the privileged processes, and then the Markov model is founded dynamically based on the features. The state sequences of dynamic Markov model are analyzed to infer the state probability, which is used to classify the normal or abnormal behavior. Because Markov model is constructed dynamically, it can extract the relationships of local behavioral features of the privileged processes adequately. When the training sets are limited, the method predicts exactly. The experiments show this method is simple, effective and efficient, and can be used in practice to monitor the computer system in real time.

Key words: intrusion detection; dynamic Markov model; information security

1 引言

入侵检测系统是现代计算机网络系统安全系统深层防御的重要组成部分. 它通过对计算机系统中的一些系统信息和系统中用户的一些行为信息进行分析来检测出对系统的入侵.

入侵检测系统可以从多个层次对计算机系统进行异常检测,但其关键是要找到最能代表用户、程序或系统行为的特征,利用提取的特征可以正确分类入侵和正常行为. 1996年,Stephanie Forrest 提出了一种通过监视特权进程的系统调用序列进行实时检测和分析的方法来对入侵行为进行检测. 该方法是将被监视进程的系统调用序列和库中的各个纪录进行匹配,如果匹配的比例比较大,则认为该进程进行的是正常的行为,否则认为是异常行为^[1]. 近几年,基于统计的学习方法得到发展,有基于频率统计、数据挖掘、有限自动机、神经网络、贝叶斯推理、支持向量机、隐 Markov 模型等^[2-7]. 这些方法除了基于隐 Markov 模型比 Forrest 的模型好一些外,其它模型相差不多. 这说明:(1) 系统调用的规律性很强,简单的模型可以

工作的很好;(2) 模型的精确度还有待提高,即现有建模方法还不够精确. 以上的研究工作对基于系统调用的入侵检测方法提出了一些部分解决问题的方法,但都需要有大量的训练数据样本. 在实际应用中,能获得的训练样本数目有限,因此以上方法的检测能力有限并且会产生大量误报.

为了克服以上方法的缺点,本文提出了一种新的入侵检测方法——基于动态马尔科夫模型的入侵检测技术(Dynamic Markov Chain——DMC). 首先从特权进程的行为特征入手,对特权进程产生的系统调用序列提取特征向量来建立正常特征库,并在此基础上动态建立 Markov 模型. 利用 Markov 模型的动态构造来充分提取特权进程的局部行为特征的相互关系,因此可以在训练数据集有限的条件下使模型更精确、检测能力大大加强.

2 DMC 方法的基本原理

直觉上,一定的系统调用排列应对应一定的程序功能,即程序行为的局部规律性应很强. 特权进程通常完成特定的、有

收稿日期:2003-09-15;修回日期:2004-03-23

基金项目:国防预先研究项目(No. 413150702);哈尔滨工程大学基础研究基金(No. HEUF04084)

限的行为, 所以其行为在时间和空间上比其他用户程序更稳定.

Markov 模型是一种简化的、高效的随机模型. Markov 模型假定在已知系统现在所处的状态的情况下, 系统将来的演变与过去无关. 特权进程的行为与其他用户进程相比是特定的、有限的、更稳定的, 其运行过程中功能的转换与衔接也是有限的、比较稳定的. 进程的正常操作和异常操作, 其系统调用序列的概率分布是不同的, 因此通过分析系统调用序列的统计性质就可以利用异常检测方法来判断该进程是否为入侵. 计算机系统进程发出的系统调用序列的变化可以被近似地看作符合 Markov 模型. 计算机所发出的当前系统调用序列只与前一时刻发出的系统调用序列直接相关, 而和前一时刻以前发出的系统调用序列不相关. 因此可以用 Markov 模型来描述进程的整体行为.

可以将进程发出的系统调用序列的不同组合看作模型的不同状态, 进程运行过程中功能的转换与衔接 (对应不同组合的系统调用序列) 可以作为状态间的转移. 因此, 进程的行为可以看作是 Markov 模型的一系列状态之间的转换. 则可以利用一系列状态的转换概率值来判别系统调用序列是否为异常.

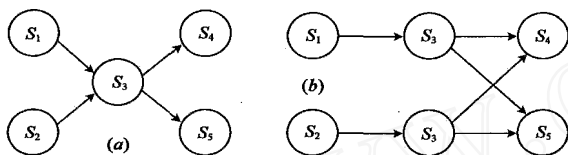


图 1 (a) 克隆前的状态转换; (b) 克隆后的状态转换

若在 Markov 模型中有状态 S_1, S_2, S_3, S_4, S_5 且相互关系如图 1 (a) 所示. 由图示知, 当状态 S_1, S_2 转移到 S_3 时, 以前的状态及转换信息就丢失了. 在状态 S_3 无法确定推知前一状态是 S_1 或 S_2 . 而由 S_3 转移到 S_4 还是 S_5 却可能与以前的状态 S_1 或 S_2 有关. 了解这一相关性是否存在的一种简单的办法是复制状态 S_3 , 产生一个新状态 S_3 , 称状态复制为克隆. 状态克隆后建立一个新的 Markov 模型, 如图 1 (b). 状态克隆后如果 S_3 的前一个状态与其下一个状态没有相关性, 则损失的信息很少. 克隆后的模型更复杂并且状态转移概率的估计对统计量更敏感. 如果 S_3 前后状态转换的相关性存在的话, 则 Markov 模型状态转移概率的估计值的准确性将有明显提高. 利用克隆方法动态构造的马尔科夫模型, 称为动态马尔科夫模型. 称像 S_3 一样具有多输入、多输出的状态为母状态.

3 DMC 算法

用 Σ 表示所有系统调用的集合, 用 O_{tr} 表示所有训练序列的集合, 用 O_{ic} 表示所有检测序列的集合. 因为特权进程的行为与其他用户进程相比是特定的、有限的、稳定的, 其运行过程中功能的转换与衔接也是有限的、稳定的. 因此, 其系统调用序列构成的集合是所有系统调用构成的序列的集合的真子集, 记为 Σ^* , 则 $O_{tr} \subseteq \Sigma^*$. 对于任意序列 Σ^* , 用 $| \Sigma^* |$ 代表序列的长度, 且 i 表示 Σ^* 的前 i 个调用子序列, $[i]$ 表示中第 i 个系统调用.

3.1 进程行为的特征提取与特征库的建立

设 k 为滑动窗口大小, 用滑动窗口在序列 Σ^* 上滑动, 步长为 1 个系统调用, 则产生 $| \Sigma^* | - k + 1$ 个短序列.

将短序列 $A_i = (a_1, a_2, \dots, a_k)$ 作为特征向量. 若 O_{tr} , 将特征向量 (a_1, a_2, \dots, a_k) 作为纪录建库, 即若 A_i 为一个特征向量, 则正常特征库为 $N_A = \{A_1, A_2, \dots, A_n\}$, 用 $| N_A |$ 表示特征库中的特征向量的数量.

3.2 进程行为的动态马尔科夫模型的训练

步骤 1 马尔科夫模型状态及转移关系的确定

将特征库中的每个特征向量 $A_i = (a_1, a_2, \dots, a_k)$ 作为 Markov 模型的状态, 即状态向量 S_i . 用计数器记录每个状态出现次数 N_i 、状态转换次序及次数 N_{ij} .

步骤 2 构造动态马尔科夫模型

对步骤 1 产生的马尔科夫模型进行动态优化. 通过对一步状态转换次序及次数 N_{ij} 进行分析, 找到每一个具有多输入、多输出的母状态 S . 对每个母状态的输入状态个数进行计数, 记为 $| S |$. 对母状态进行克隆, 产生 $| S | - 1$ 个新状态. 建立克隆状态映射表, 记录母状态 S 克隆前后的状态映射关系, 如表 1 所示.

表 1 克隆状态映射表

前一状态 (输入状态)	当前状态	克隆后的状态
S	S	S

步骤 3 对步骤 2 产生的马尔科夫模型重新进行训练

在训练过程中需建立两个指针 P_{first} 和 P_{now} , 分别记录状态转换过程中的前一状态和当前状态. 将 P_{first} 和 P_{now} 与克隆状态映射表中的前两项进行匹配, 若能成功将 P_{now} 的值修改为 S .

Markov 模型状态转移概率矩阵 P 和初始化分布 π 可由对训练序列的状态及转换的观察统计求得. 本文中用状态转移的频率来近似概率计算. 重新用计数器记录每个状态出现次数 N_i 、状态转换次序及次数 N_{ij} . 在训练序列中不可能包含正常序列的所有排列情况, 即训练不可能是充分的. 这样有必要把状态分为两部分: 必要状态、补充状态. 若 O_{tr} , 则相应求得的 $(S_i = A_i) \quad N_A$, 称 S_i 为必要状态; 若 $(S_j \notin O_{tr} \quad (\Sigma^*))$, 则有 $(\forall A_i \in N_A) \quad (S_j = A_i)$, 称 S_j 为补充状态. 为了正常库及状态规模扩展方便, 设状态 S_0 表示补充状态. 则 Markov 模型的必要状态一步状态转移概率:

$$p_{ij} = \frac{N_{ij}}{N_i}, \quad i = 0, j = 0 \text{ 且 } \sum_j p_{ij} = 1$$

补充状态一步状态转移概率:

$$p_{i0} = p_{0j} = \quad (1)$$

式中 N_{ij} 表示状态 i 向状态 j 转移的次数; N_i 表示由状态 i 转移到任意一个状态的次数. 因为补充状态 S_0 在训练过程中没有出现过, 因此有理由相信其出现的概率非常小, 可使 $p_{i0} = \min(p_{ij})/10, \quad 0 < i, j \leq | N_A |$.

经过大量实验观察发现, 只要滑动窗口 k 设置得当, 初始状态就为某一固定状态, 因此可设 $\pi = [1, 0, \dots, 0]$.

3.3 检测方法

检测算法是建立在连续 l 个状态序列的转换概率上的,

即异常行为要实现入侵,其在行为的某一局部要实现特定的、与正常序列可区别的功能,依据概率值可以区别出正常、异常情况.与建立 Markov 模型的过程相同,得到要检测序列的连续 l 个状态,求其连续个状态序列的转换概率.

$$P_l(S_{i-l+1}, \dots, S_i) = \prod_{t=i-l+1}^i P_{S_{t-1}S_t} \quad (2)$$

存在两种异常情况:

(1) $S_{i-1} \neq S_i$ 且 $(\forall A_i \in N_A) (S_i \in A_i)$, 即状态 S_i 不属于必要状态,则 $S_i = S_0$.

(2) $S_{i-1} = S_i$ 且 $P(S_{i-1}, S_i) = 0$, 则 $P(S_{i-1}, S_i) = 0$.

为了便于在线检测,可用如下递推算算法:

$$P_t(S_{t-l+1}, \dots, S_t) = \begin{cases} P_{t-1} * P_{S_{t-1}S_t} \\ P_{S_{t-l}S_{t-1}} \end{cases}, t > l \quad (3)$$

4 试验及其结果分析

本试验采用美国新墨西哥大学计算机科学系提供的特权进程在正常及入侵过程中的系统调用序列.试验中,把 LPRCP 数据分为两部分:一类为训练数据,共 600 个正常进程的系统调用序列数据,用于建立正常进程的模型;另一类为测试数据,共有 600 个正常进程,1000 个异常进程,这些测试数据用来测试模型与算法的有效性.

为了评价 DMC 模型的检测性能,对相关的统计量作如下定义,并且为了与其它方法相比较将误检率分为两种情况进行定义:

检测率 $TPR(\text{True Positive Rate}) = \frac{\text{检测出的异常序列数}}{\text{异常序列总数}}$;

误检率 1 - $FPR(\text{False Positive Rate}) = \frac{\text{正常序列被误报为异常序列数}}{\text{正常序列总数}}$;

误检率 2 - $SFPR(\text{Subsequence False Positive Rate}) = \frac{\text{正常短序列被误报为异常短序列数}}{\text{正常短序列总数}}$.

在训练数据固定的情况下,正常行为的特征库的规模(特征向量的个数)随短序列长度 k (滑动窗口大小)的变化而变化的情况如图 2.

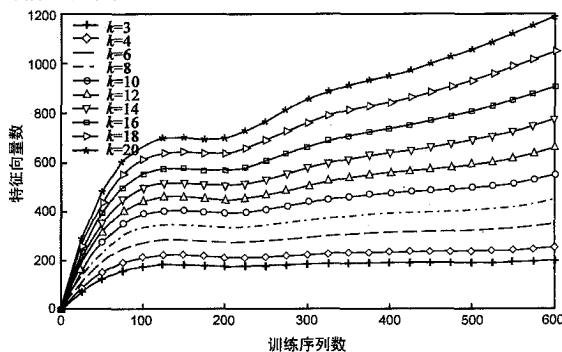


图 2 特征库的规模随短序列长度 k 的变化情况

由图 2 可知,刚开始特征库的增长很快,但随后逐渐减慢.当短序列长度 $2 < k < 12$ 时,且训练序列数大于 500 时,特征库规模几乎不再增长.由此可知,正常行为是有限的;且在训练数据有限的情况下,合理利用滑动窗口大小和训练集

关系,可以建立起一个较合理的模型来描述进程的行为.

在这一模型下, $k=3$ 时的状态数、误检率 FPR 与训练数据集的关系如图 3 所示.由图 3 可知,随着训练数据集的增大,最初状态数迅速增大、误检率迅速降低;当训练数据集增大到一定程度时,状态数增长缓慢、误检率不再有明显的降低.

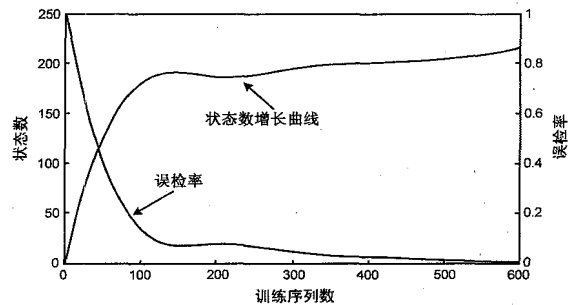


图 3 $k=3$ 时的状态数与误检率的关系

表 2 给出了 DMC 的检测结果与 S. Mukkamala^[5]、Forrest^[1,6]等人得出的研究结果的比较.

表 2 DMC 与其它几种入侵检测方法的性能比较

	k	最高 TPR	最低 SFPR
stide	6	96.9 %	0.0008
r-stide	6	96.9 %	0.0075
RIPPER	6	95.3 %	0.0016
HMM	6	96.9 %	0.0003
SVM	6	99.87 %	0.0003
DMC	3	100 %	0.00016
DMC	6	100 %	0.00016614

5 结论

本文首先从特权进程的行为特征入手,提取特征向量来建立正常特征库,并在此基础上动态建立 Markov 模型.利用 Markov 模型的动态构造能够充分提取特权进程的局部行为特征的相互关系,因此可以在训练数据集有限的条件下使模型更精确、检测能力大大加强.最后将系统调用序列转变为状态序列,求取定长状态序列概率来判断是否异常.

实验结果表明,本方法对入侵进行检测非常有效.但仍有很多问题需要进一步研究,如状态序列长度的确定等.这些问题还有待于进一步的研究.

参考文献:

[1] S Forrest, S A Hofmeyr, A Somayaji, T A Longstaff. A sense of self for unix processes[A]. In 1996 IEEE Symposium on Security and Privacy [C]. Los Alamitos, CA: IEEE Computer Society Press, 1996. 120 - 128.

[2] T Lane, C E Brodley. Temporal sequence learning and data reduction for anomaly detection[J]. ACM Transactions on Information and System Security, 1999, 2(3): 295 - 331.

- [3] Lee W,Stolfo S J. Data mining approaches for intrusion detection[A]. Proceedings of the 7th USENIX Security Symposium[C]. San Antonio. Texas :the USENIX Association ,1998. 26 - 29.
- [4] K Ilgun ,R Kemmerer ,P Porras. State transition analysis :A rule-based intrusion detection approach[J]. IEEE Transactions on Software Engineering ,1995 ,21(3) :181 - 199.
- [5] S Mukkamala ,GJanowski ,A H Sung. Intrusion detection using neural networks and support vector machines[A]. Proceedings of IEEE International Joint Conference on Neural Networks[C]. Hawaii ,2002. 1702 - 1707.
- [6] Warrender C ,Forrest S ,Pearlmutter B. Detecting intrusion using system calls :Alternative data models[A]. IEEE Symposium on Security and Privacy[C]. Oakland ,USA ;1999. 133 - 145.
- [7] 谭小彬,王卫平,奚宏生,殷保群. 计算机系统入侵检测的隐马尔可夫模型[J]. 计算机研究与发展 ,2003 ,40(2) :245 - 250.

作者简介:



尹清波 男,1975 年出生于黑龙江省齐齐哈尔市,2004 年获哈尔滨工程大学计算机科学与技术学院计算机应用技术专业工学硕士学位,现为博士研究生,主要研究方向为信息系统安全理论与技术、机器学习、模式识别。



张汝波 男,1963 年出生于吉林省吉林市,哈尔滨工程大学计算机科学与技术学院教授、博士生导师,主要研究方向为信息系统安全理论与技术、机器学习、模式识别与人工智能。

李雪耀 男,出生于浙江省奉化市,哈尔滨工程大学计算机科学与技术学院教授,主要研究方向模式识别与人工智能、语音信号处理。

王慧强 男,出生于河南省沈秋,哈尔滨工程大学计算机科学与技术学院教授,博士生导师,主要研究方向为信息系统安全理论与技术。

www.cnki.net