

文章编号: 1004 - 7271(2005)02 - 0220 - 05

·研究简报·

基于统计方法的入侵检测模型

A model of statistics-based intrusion detection

任照松, 印润远

(上海水产大学信息学院计算机系, 上海 200090)

REN Zhao-song, YIN Run-yuan

(Department of Computer, College of Information, Shanghai Fisheries University, Shanghai 200090, China)

关键词: 入侵检测; 会话向量; 伯努利向量; 疑义商

Key words: intrusion detection; session vector; bernoulli vector; suspicion quotient

中图分类号: TP 309.5 文献标识码: A

随着 internet 的发展, 信息安全越来越引起人们的注意. 据计算机紧急响应小组 (CERT) 公布的从 1988 年到 1999 年攻击事件增长数据可以看出, 从 1988 年到 1999 年短短十年里, 攻击事件增长将近 20 倍, 与攻击相应的结果就是经济的损失, 据美国商业杂志《信息周刊》公布的一项调查报告称, 黑客攻击和病毒安全问题仅在 2000 年就造成了上万亿美元的经济损失, 在全球范围内每隔数秒就会发生一起网络攻击事件. 由此不难看出: 传统的基于信息加密的被动式信息保护技术已经不能满足人们对信息安全的需求, 主动检测攻击的防御技术变得迫切重要, 这也是近年来入侵检测技术得以迅速发展的原因之一, 本文正是基于这种需要而提出的一种入侵检测模型。

1 入侵检测技术

入侵检测技术是在二十世纪 90 年代后随着网络的迅速发展, 攻击事件和病毒等安全问题的大量出现而迅速发展起来的一种网络安全技术. 按照数据分析手段的不同, 入侵检测可以分为滥用 (misuse intrusion detection) 入侵检测和异常入侵检测 (anomaly intrusion detection)^[1]. 在异常入侵检测中, 最广泛使用的较为成熟技术是统计分析, IDES 系统实现了最早的基于主机的统计模型. 另一种主要的异常检测技术是神经网络技术^[2]. 此外, 如基于贝叶斯网络的异常检测方法, 基于模式预测的异常检测方法, 基于数据挖掘^[3]的异常检测方法以及基于计算机免疫学的检测方法^[4]也相继出现. 对于滥用入侵检测也有多种检测方法, 如专家系统 (expert system), 特征分析 (signature analysis), 状态转移分析 (state-transition analysis) 等. 在国内, 入侵检测技术目前正处于发展时期, 很多入侵检测技术相继出现, 但基本处于理论研究阶段, 在市场很少有成型的实用产品出现, 和国外的水平差距仍然很大。

2 基于统计方法的入侵检测模型

本文提出的这个模型可以用图 1 表示. 计算机系统 (如 Unix) 里都有审计记录^[5], 它记录了所有用

收稿日期: 2004-08-27

作者简介: 任照松 (1980 -), 男, 安徽舒城人, 硕士研究生, 专业方向为网络技术应用与网络安全。

通讯作者: 印润远 (1953 -), 男, 上海市人, 副教授, 从事计算机网络安全、网络技术应用方面的研究. Tel: 021 - 65711417, E-mail: ryyin

@shfu.edu.cn.

户使用计算机资源的情况,不管是合法用户还是非法用户,只要设置合理,它都能详细地记录用户的活动,而这些审计记录就是我们入侵检测模型的实施依据。

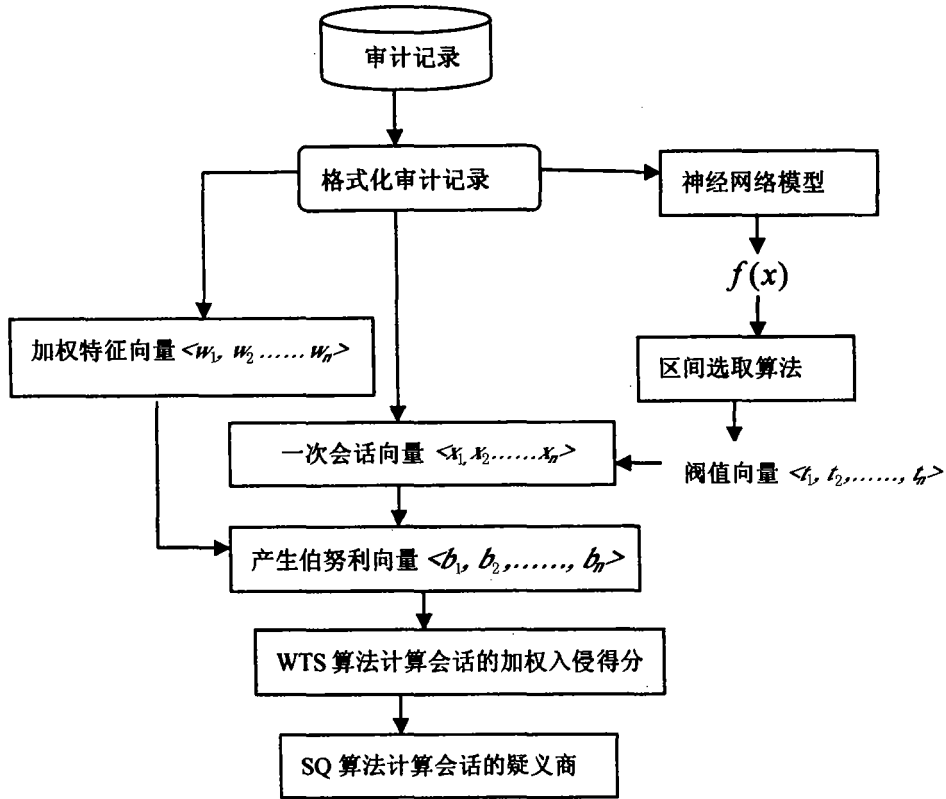


图 1 基于统计方法的入侵检测模型

Fig. 1 A model of statistics-based intrusion detection

2.1 审计记录的格式化

各个系统产生的审计记录格式各不相同,为了本文中提出的模型的处理需要,要求对系统的审计记录进行格式转化,转化为需要的二维关系表格的形式,列表代表审计时考察的各种属性,如用户 ID,活动类型,CPU 使用率,网络连接时间,会话时间,读文件属性,安全 I/O 失效次数等,行表示一条审计记录,它记录了用户的一次会话信息,一次会话信息记录了用户登录到退出这段时间内的所有活动信息。

2.2 用神经网络模型来训练数据

这个模型是用来检测系统中出现的异常现象的,因此必须知道什么情况才是系统的正常状态,而描绘一次会话活动的信息是会话向量,这样,对于会话向量的每个属性值应该有一个相应的正常取值区间,对任一个会话向量 $X = \langle x_1, x_2, \dots, x_n \rangle$,其属性 x_i 属于正常区间 $t_i = [t_{i, \min}, t_{i, \max}]$,这意味着在一次正常的会话中,属性 x_i 属于该属性的正常区间 $[t_{i, \min}, t_{i, \max}]$ 。由每个属性的正常区间 t_i 组成阈值向量 $T = \langle t_1, t_2, \dots, t_n \rangle$,这个向量是要从系统的正常会话的审计记录里提取出来,在本文的模型里使用神经网络技术来训练这些审计记录,获得每个属性的正常区间,从而得到阈值向量。上面也提到,阈值向量里每个属性代表一个区间,因此在这里使用一个二元组来描述这个区间,同时出于实际情况的考虑, x_i 属性值应遵从高斯分布,也就是正常会话活动向量的 90% 以上属性值都落在相应的区间内,也即属性的正常区间,而非正常活动向量的属性值则落在对应的区间外,这个区间我们称之为入侵可能区间。如图 1 所示,在神经网络模型^[6]里,用格式化审计记录后所得的二维表的一列作为一次输入,输出的为一个满足这一列数据分布要求的函数 $f(x)$ 。再由这个函数 $f(x)$ 经由下面区间选取算法获得属性

的正常区间。

2.3 区间选取算法

如果神经网络训练的数据足够多的话,上步输出的函数 $f(x)$ 应该满足高斯分布(正态分布)。假设 A, B 为 $f(x)$ 分布中当 $f(x) \rightarrow 0$ 时, x 的两个极限取值, C 点为 $f(x)$ 的最大值, 可以用 $f'(x) = 0$ 来确定。下面给出一个算法来确定 t_i , 假设会话向量 $X = \langle x_1, x_2, \dots, x_n \rangle$ 的第 i 个属性的高斯分布函数为 $f(x)$, A_i, B_i, C_i 分别代表 A, B, C 三点。

下面的算法将给出 t_i 值。

$$m = \int_{A_i}^{B_i} f(x) dx$$

For $i = 1$ to count /* 使用 count 次循环是让最后得到的 t_i 值在中值附近 */
随机产生一个 $t_{i, \min} \in (A_i, C_i)$

求 $t_{i, \max}$ 使得在 $n = \int_{t_{i, \min}}^{t_{i, \max}} f(x) dx$ 下, 满足条件 $n/m \in (0.9, 1)$

$$u_i(i) = (t_{i, \min}, t_{i, \max})$$

Next i

$$t_i = \sum_{i=1}^{\text{count}} u_i(i) / \text{count}$$

上述过程重复就可以得到阈值向量 $\langle t_1, t_2, \dots, t_n \rangle$ 。

2.4 加权特征向量

加权特征向量(weighted feature vector)^[7] $W = \langle w_1, w_2, \dots, w_n \rangle$ 表示对于某一种类型的入侵 I , 属性 x_i 对入侵 I 的贡献有多大, 这个贡献力用 w_i 衡量, 也即对于一次确定为入侵 I 的会话, 其属性 x_i 处于入侵可能区间的概率为 w_i , 如果有 $w_i > w_j (i \neq j)$, 并且第 i, j 个属性的值都分别落在阈值 t_i, t_j 的范围外, 那么第 i 个属性的值在确定一个特别类型的入侵中更有用。加权特征向量是从审计记录中使用统计分析的方法产生的, 首先必须收集一段时间内的格式化的审计记录, 把所有已知入侵会话分类, 对于每一种入侵类型, 分别用下面的算法计算其入侵向量, 计算入侵类型 I 的加权特征向量算法:

初始加权特征向量 $W^1 = \langle w_1, w_2, \dots, w_n \rangle$ 使得其属性 $w_i = 0$;

For $i = 1$ to n

each R in 审计记录里入侵类型 I 的所有会话向量组(其个数记为 N)

If ($R - > x_i \notin t_i$) then $w_i = w_i + 1$

Next

$$w_i = w_i / N$$

Next i

2.5 计算伯努利向量

伯努利向量(Bernoulli Vector) $B = \langle b_1, b_2, \dots, b_n \rangle$ 就是一个简单的二进制向量, 它描述了: 对于一次会话活动, 其相应的会话向量中那些属性的值落在阈值之外, 也就是说, 若第 i 属性 x 点值落在范围 t_i 之外的话, $b_i = 1$, 否则 $b_i = 0$ 从而这个过程可以描述为:

for $i = 1$ to n

$$b_i = \begin{cases} 0 & \text{若 } t_{i, \min} < x_i < t_{i, \max} \\ 1 & \text{否则} \end{cases}$$

2.6 加权入侵得分的计算

加权入侵得分 (Weighted Intrusion Score) 是指: 对于一个某种入侵类型 I , 某次会话属于入侵类型 I 的得分是多少。有了计算加权入侵得分, 就能够根据所有会话的加权入侵得分的分布情况, 可以给出一个给定的会话分配一个疑义商, 下面将给出它的计算方法。

它可由下式来计算:

$$Wis = \sum_{i=1}^n b_i * W_i$$

2.7 计算疑义商

疑义商 (Suspicion Quotient)^[7] 是描述一个会话和其他会话相比, 它相似于某种入侵的程度。如果有两个会话, 会话 1 的疑义商大于会话 2 的疑义商, 那么我们就认为会话 1 比会话 2 对该种入侵有更大的相似性。

一个会话的疑义商可以通过会话集中有百分之几的会话的疑义商小于等于当前会话的疑义商来确定, 这个百分比就是当前会话的疑义商。换句话说, 如果所有的会话按照它们加权入侵得分来排序的话, 在排好序的序列中, 某一特定会话的所在的位置就是该会话的疑义商。如若有一个会话的疑义商是 0.97, 则意味着只有占 0.03 的会话比该会话更可疑。这样我们就可以通过疑义商的大小来判断一个会话是不是一个特定的入侵。下面我们将给出疑义商的计算算法。

在下面的算法中, 我们要使用一个数据结构—数据立方体, 一种三维的数据聚集模型, 在这里我们使用 $Intrusion_Type(k)$, $attribute(i)$, $Score(j)$ 表示数据立方体的三维。对于一个任意的会话向量, $p_{k,i,j}$ 表示对于第 k 种入侵类型来说, 使用前 i 个属性的权值时, 一个会话具有加权入侵得分 (以下简称为 WIS) j 的可能性, max_Score 是所有可能的入侵得分中的最大值, $p_{k,r_i}(0)$ 定义为对于第 k 种入侵类型来说, 会话向量的第 i 个属性在其阈值范围内的可能性 (此时有 $b_i * w_i = 0$), 类似地, $p_{k,r_i}(1)$ 定义为对于第 k 种入侵类型来说, 会话向量的第 i 个属性在其阈值范围外的可能性 (此时有 $b_i * w_i = w_i$), 也就是 $p_{k,r_i}(1) = 1 - p_{k,r_i}(0)$, 参数 $\{p_{k,r_i}(0), i = 0, 1, \dots, n\}$ 是可由用户定义的, 一般设置为 0.9 左右^[8]。算法描述如下:

```

Pk0,0 = 1.0
for k = 1 to N
  for j = -1.0 to max_Score /* max_Score 是所有人侵得分的最大值 */
    pk,0,j = 0.0
  for k = 1 to N
    for i = 0 to n
      for j = 0 to max_Score
        Pk,i,j = Pk,i,j * Pk,r_i(0) + Pk,i-1,j-w_i * Pk,r_i(1)
      next j
    next i
  next k

```

这个算法结束后, 它就构造了一个三维的数据立方体, 最后我们就可以用这个数据立方体来计算一个会话对应每种入侵类型的疑义商, 表示如下:

```

For k = 1 to N
  Suspicious Quotient SQ[k] =  $\sum_{j=0}^{max\_score} P_{k,n,j}$ 

```

它就是判断入侵的依据, 其值越大意味着相应的会话表示入侵的可能性越大。我们只要根据领域专家的经验, 设置相应的阈值就可以自动地实现入侵报警。

3 讨论

下面就基于统计方法的入侵检测模型的实现问题进行简单的讨论。

就实际应用来看:这个模型在工作前需要一些初始化数据,这个初始化数据就是从格式化后的审计记录中提取得到的每种入侵类型的加权特征向量,这个工作可以在模型实时检测前进行,一次处理多次使用,以后只有新的入侵类型被发现后才需要更新这个审计记录集,且这个更新也就是加入该种入侵类型的加权特征向量。在模型启动工作后,需要不断地收集系统中处于活动状态的会话的会话向量,把它们作为模型的输入,计算其对应每种入侵类型的加权入侵得分和相应的疑义商,然后每一个会话向量取其具有最大值的疑义商,根据领域专家设定的相应类型入侵疑义商的阈值,判断会话是否构成一次入侵。如果超过这个阈值,就发出对应类型的入侵报警。否则,一个会话的一次检测结束。

模型的计算性能方面:在模型启动后,模型运行的计算性能是由两个因素决定的,即会话向量加权入侵得分的计算和对应疑义商的产生。假设系统中已知入侵类型的个数为 N ,加权特征向量取两位有效数字,则对一个会话来说,计算其加权入侵得分的时间复杂性可表示为: $O(n \times N)$,疑义商则为 $O(10^2 \times N + 10^2 \times n \times N)$,如果系统中平均处于活动状态的会话个数为 t ,则这个模型工作时的计算性能(或复杂性)可以用 $O(t \times 10^2 \times n \times N)$ 粗略衡量,也即模型工作中的计算性能是由四个因素共同决定的:系统中处于活动状态的个数,加权特征向量的有效位数,会话向量的维数以及已知入侵类型数目。

参考文献:

- [1] 唐正军.入侵检测技术导论[M].北京:机械工业出版社,2004.4.
- [2] 马锐,刘玉树,杜彦辉.基于神经网络专家系统的入侵检测方法[J].计算机工程与应用.2004,40(2):151-153.
- [3] 胡敏,潘学增,平玲娣.基于数据挖掘的实时入侵检测技术的研究[J].计算机应用研究.2004,21(1):105-108.
- [4] Fabio Gonzalez, Dipankar Dasgupta. Neuro-Immune and Self-Organizing Map Approaches to Anomaly Detection: A Comparison[A]. The Proceedings of the 1st International Conference on Artificial Immune System[C], 2002,(19):101-103.
- [5] Aurobindo Sundaram. An Introduction to Intrusion Detection[R], 1996, 40(1):17-21.
- [6] 戴葵(译).神经网络设计[M].北京:机械工业出版社,2002.9.
- [7] Biswanath Mukherjee, L Todd Heberlein, Karl N Levitt. Network Intrusion Detection[J].IEEE NetWork, 1994,5(2):42-47.
- [8] Smaha S E. HayStack: An Intrusion Detection System[A], Proc. IEEE 4th Aerospace Computer Security Conference[C]. Orlando, FL, 1988, 17(1):70-73.