

# 基于模糊 C 均值聚类的网络入侵检测算法<sup>\*</sup>

杨德刚

(重庆大学计算机科学与工程学院 重庆400030) (重庆师范大学数学与计算机科学学院 重庆400047)

**摘要** 入侵检测已成为网络安全的第二层重要防御线。分析了对新型未知的攻击的入侵检测,提出基于模糊 C 均值聚类的网络入侵检测算法。用 KDD-99 数据集的仿真实验结果表明算法的可行性、有效性和可扩展性,并有效提高了聚类检测的检测率,降低了误检率。

**关键词** 入侵检测,聚类分析,模糊 C 均值聚类

## Research of the Network Intrusion Detection Based on Fuzzy Clustering

YANG De-Gang

(College of Computer Science and Engineering, Chongqing University, Chongqing 400030)

(College of Mathematics and Computer Science, Chongqing Normal University, Chongqing 400047)

**Abstract** The intrusion detection become the second floor defense line of the network security. Analyze to the characteristic of the intrusion detection technique for newly and unknown attack, and brings forward algorithm of Network Intrusion Detection based on Fuzzy C-means Clustering. The result of simulations run on the KDD-99 datasets show the feasible, efficient and extensible for unknown intrusion detection, and increase detection rate of the clustering detection and decrease the false alarms rate.

**Keywords** Intrusion detection, Clustering analysis, Fuzzy c-means clustering algorithm

## 1 引言

入侵检测已成为传统安全技术,如防火墙之后的第二道安全防御线。入侵检测从技术上主要分为误用检测和异常检测两类。前者通过特征匹配检测是否发生入侵,对已知入侵检测准确,速度快,但它不能检测未知入侵;而后者则在建立正常模型基础上,通过实际行为与之比较是否偏离来判断入侵行为的发生。它能检测已知入侵,也能检测未知入侵。因此异常检测成为当前研究的热点领域。

聚类分析用于发现数据实例中的隐性模式和用于检测入侵中有意特征。好的聚类方法是要将一个数据集划分为若干组,并且具备高的组内相似性和低的组间相似性两个特点。因此,聚类分析用于入侵检测是可行的。Portnoy<sup>[3]</sup>最先提出基于聚类分析的入侵检测技术,该技术在数据实例进行正规化处理,然后采用欧氏距离,并用单链法进行聚类,经过标识,通过分类以检测入侵。其检测率可达50%,误检率1%左右。该方法存在自适应和可扩展性不强的特点。本文则根据聚类和入侵实例的特征,提出基于模糊 C 均值聚类的入侵检测算法,仿真实验结果表明该算法能有效提高检测率,降低误检率。

## 2 基于模糊聚类的入侵检测算法

### 2.1 基于聚类分析的入侵检测过程

用聚类分析进行入侵检测其过程如下:

(1)收集实例数据 网络入侵检测系统首先收集网络上传输的数据流,通常用抓包工具来收集网络数据包,其中最常用的如 UNIX 下的 TCPDUMP,它能够监听和接收网络中所

有正在传输的数据包,并把它们记录到文件中。原始的网络数据包并不适合进行聚类分析,需要将原始的网络数据包恢复成 TCP/IP 层的连接记录,其中每个连接记录代表一次 TCP/IP 层的连接,并包含一个网络连接的多个属性,如网络协议、连接起始时间、连接结束时间、服务、源 IP 地址、目的地址、连接终止状态等。

(2)实例数据标准化 由于采集的网络数据的属性值之间的差别可能很大,而且它们可能用不同的单位来度量。为了消除由于度量不同对聚类的影响,应对测量的属性值进行标准化,方法如下:

$$x'_{ij} = \frac{x_{ij} - \bar{x}_j}{S_j}, \begin{cases} i=1, 2, \dots, n \\ j=1, 2, \dots, m \end{cases} \quad (1)$$

$$\text{其中 } \bar{x}_j = \frac{1}{n} \sum_{i=1}^n x_{ij}, S_j = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2}$$

$x'_{ij}$  即为标准化后的实例数据值。通过计算每个特征值与平均值之间的标准偏差,可得到该特征值在正规化空间中的新值。

(3)数据初始化聚类 对收集的连接记录进行标准化,即可进行初始化的聚类过程。能够利用的聚类算法有很多,如可以采用层次的聚类算法和基于模型的聚类方法等。

(4)对初始化聚类集的优化 为了获得更优化的聚类结果,应对初始化的聚类簇进行优化,以尽可能地将相似的实例聚为一类,相异的实例分到不同聚类。

(5)聚类簇的划分 由于正常条件下正常实例占网络行为的大部分,因此可以假设正常行为在数量上远远大于各种攻击行为。当使用聚类算法得到聚类簇后,即利用这一假设进行划分,区分哪些是正常连接数据,哪些是入侵连接数据。

<sup>\*</sup> 基金项目:重庆师范大学校内基金(计算机网络安全)。杨德刚 硕士生,讲师,主要研究方向:数据挖掘、信息安全。



Conf. 1995

- 9 Park J, Sandhu R. Towards Usage Control Models: Beyond Traditional Access Control. SACMAT02, Monterey, California, USA, ACM, 2002
- 10 Park J, Sandhu R. Originator Control in Usage Control. In: 3<sup>rd</sup> International Workshop on Policies for Distributed Systems and Networks (Policy02). June 2002
- 11 Sandhu R, Park J. Usage Control: A Vision for Next Generation Access Control. MMM-ACNS 2003
- 12 Sandhu R, Coyne E J, Feinstein H L, Youman C E. Role-based access control models. IEEE Computer, 1996, 29(2): 38~47
- 13 Sandhu R. Rationale for the RBAC96 family of access control models. In: Proc. of the 1st ACM Workshop on Role-Based Access Control. ACM, 1997
- 14 Barkley J. Comparing Simple Role Based Access Control Models and Access Control Lists. In: Proc. of the Second ACM Workshop on Role-Based Access Control, Nov. 1997. 127~132
- 15 Park J, Sandhu R, Schifalacqua J. Security Architectures for Controlled Digital Information Dissemination. ACSAC, 2000
- 16 Sandhu R. Access Control: The Neglected Frontier. ACISP, 1996
- 17 Chen F, Sandhu R. Constraints for RBAC, ACM RBAC, 1995
- 18 Thomas R, Sandhu R. Conceptual Foundations for a Model of

Task-based Authorizations, CSFW, 1994

- 19 Sandhu R. Design and Implementation of Multilevel Databases. RADC, 1994
- 20 Sandhu R, et al. Role-Based Access Control: A Multi-Dimensional View. ACSAC, 1994
- 21 Thomas R, Sandhu R. Discretionary Access Control in Object-Oriented Databases: Issues and Research Directions. NCSC, 1993
- 22 Jajodia S, Sandhu R. Toward a Multilevel Secure Relational Data Model. SIGMOD, 1991
- 23 Sandhu R. Evaluation by Parts of Trusted Database Management Systems. RADC, 1991
- 24 Sandhu R, Share M. Some Owner-Based Schemes with Dynamic Groups in the Schematic Protection Model. OAKLAND, 1986
- 25 Sandhu R, Bhamidipati V, Munawer Q. The ARBAC97 Model for Role-Based Administration of Roles. ACM Transactions on Information and System Security, 1999, 2(1): 105~135
- 26 Park J, Sandhu R. The UCONABC Usage Control Model. ACM Transactions on Information and System Security (TISSEC), Feb. 2004
- 27 Yovits M C. Database Security. Advances in Computers. Academic Press, 1994, 38: 1~74

(上接第87页)

类标识属性,其仅供算法结果分析。表1和图1是各数据集中心攻击类型分布情况。

表1 实验数据集结构

	实例数	正常实例数	入侵实例数	攻击类型数据	分类攻击数
数据集1	2100	2000	100	21	4
数据集2	2100	2000	100	11	3
数据集3	2100	2000	100	13	4
数据集4	2100	2000	100	8	2
数据集5	2100	2000	100	13	4

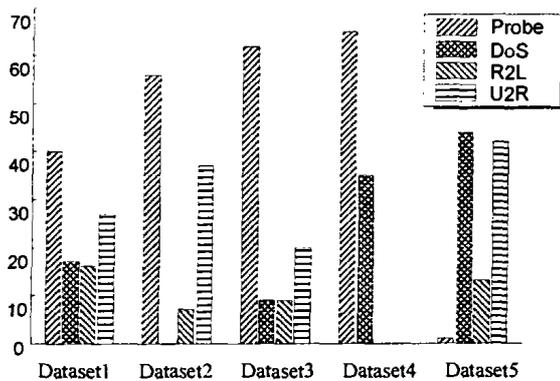


图1 实验数据集中攻击类型分布

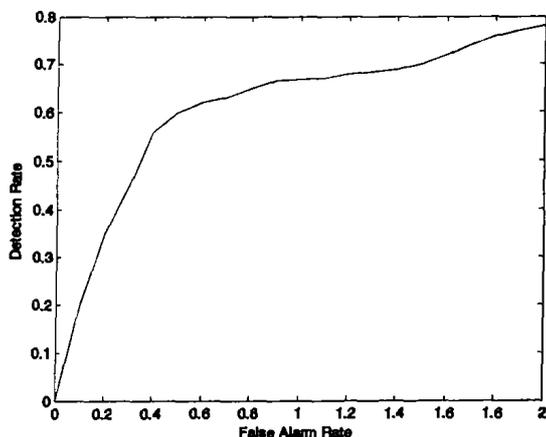


图2 NIDBFCM 检测率和误检测率的相对曲线

针对以上数据集,选用实验平台 Windows2000, Matlab6.5语言编程环境, Intel Celeron 1.70GHz CPU, 256MB 内存。为评价入侵检测算法性能使用两个指标:检测率 DR (Detection Rate) 和误检率 FR (False Alarm Rate)。检测率 DR 是被系统检测到的入侵实例数目与检测数据集中总的入侵实例数目之比;误检率 FR 则为被误判为入侵的正常实例数目与总的正常实例数目之比。这两项性能指标能充分反映算法的入侵检测能力。最好的入侵检测系统应该使 DR 尽可能的大,而 FR 则尽可能的小。图2是上述5组实例数据的平均检测结果,在隶属度改变的情况下,检测率提高的情况下,误检率也会提高。

从图2可知最好性能当 FR=1%时, DR=69%,这充分表明算法对于未知入侵行为检测的可行性和有效性,与 Portnoy 提出的方法相比,在一定程度上提高了检测率,降低了误检率。

**结束语** 本文分析了聚类技术及其在入侵检测中的应用,提出以 FCM 进行入侵检测,分析了检测过程,并用 KDD-99 实验数据集进行仿真实验,其结果表明该算法提高了检测的效率,算法是可行的、有效的。

### 参考文献

- 1 Wenke L. A Data Mining Framework for Building Intrusion Detection Models. In IEEE Symposium on Security and Privacy, 1999
- 2 Lee W, Stolfo S J, Mok K W. Mining in a data-flow environment: experience in intrusion detection. submitted for publication, 1999
- 3 Portnoy L, Eskin E, Stolfo S J. Intrusion detection with unlabeled data using clustering. In: Proc. of ACM CSS Workshop on Data Mining Applied to Security (DMSA-2001). Philadelphia: ACM Press, 2001(11)
- 4 Li Xiangyang. Clustering and Classification Algorithm for Computer Intrusion Detection: [PhD.]. Arizona State University, 2001. 12
- 5 Equihua M. Fuzzy clustering of ecological data J. Ecol, 1990, 78: 561~567
- 6 <http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html>