

# 一种新的基于模糊聚类 and 免疫原理的入侵监测模型

陶新民, 陈万海, 郭黎利

(1 哈尔滨工程大学信息与通信工程学院, 黑龙江哈尔滨 150001)

**摘 要:** 本文提出了一种新的基于模糊聚类 and 免疫原理相结合的入侵检测模型, 同时文章中对 RPCL 算法进行了改进, 克服了原有 RPCL 算法中不同变量量纲以及变量相互间相关性对算法性能的影响, 同时解决了原有算法对初始分类数敏感的不足. 利用改进后的 RPCL 算法解决了模糊聚类中分类个数不确定的问题, 最后利用遗传算法进行半径的最优调整形成非正常子空间的特征函数. 最后给出试验结果并进行了分析.

**关键词:** 入侵检测; 模糊聚类; 免疫原理; 遗传算法; RPCL 算法; 相关性

**中图分类号:** TN 915 08 **文献标识码:** A **文章编号:** 0372-2112 (2006) 07-1329-04

## A Novel Model of IDS Based on Fuzzy Cluster and Immune Principle

TAO Xin-min, CHEN Wan-hai, GUO Li-li

(1. Communication Technology Institute of HREU, Harbin, Heilongjiang 150001, China)

**Abstract** This paper presents a novel model of intrusion detection based on fuzzy cluster and immune principle. The original RPCL algorithm is modified in order to address the problem of different variability of variables and correlation between variables. The sensitivity to initial number of clusters is also solved in this paper. Especially, this paper uses the extended RPCL algorithm to determine the initial number of clusters in the fuzzy cluster algorithm. The genetic algorithm is used to optimize the radius deviation for the determination of characteristic function of abnormal subspace. Some results are finally reported with some concluding remarks.

**Key words** intrusion detection; fuzzy cluster; immune principle; genetic algorithm; RPCL algorithm; correlation

### 1 引言

随着计算机系统中互联速率的不断增长, 网络安全越来越成为一个重要的挑战. 为了满足这个挑战, 人们设计了入侵监测系统, 其目的是为了安全系统中重要数据的有效性和整合性, 同时也保护计算机网络不受拒绝服务攻击的侵害、未授权信息的公开以及数据的改变和破坏. 入侵检测可以分为两类: 一种基于模式匹配的入侵检测系统, 另一种基于异常发现的入侵检测系统. 前者依赖于匹配以前定义的已知入侵模式的信号特征, 后者依靠与正常行为的偏离程度作为入侵行为的判定条件. 入侵检测所采用的方法也有很多, 其中有基于审计的攻击检测、基于神经网络的攻击检测技术<sup>[1-4]</sup>、基于专家系统的攻击检测技术<sup>[5]</sup>、基于推理的攻击检测技术以及基于遗传和免疫的检测技术等. 本文提出了一个基于模糊聚类分析和免疫原理的入侵检测模型. 人工免疫技术早在 1994 年由 Forrest 提出的, 1999 年由 Dasgupta 和 Hofmeyr 提出将免疫原理应用到计算机安全领域中. Stephanie Forrest 和她在墨西哥

大学的工作组很长时间一直致力于人工免疫系统的研究. 在他们的观点中, 保护计算机免受病毒的侵害本质上就是区分自己和其他非己的通用问题. 这种方法叫做阴性选择算法, 它被用来监控受保护的数据和文件的变化. 本文技术的核心是利用模糊聚类分析和遗传算法来进化能覆盖整个非我模式空间的规则, 即确定非正常样本集合的特征函数. 然后利用人工免疫原理的阴性选择算法进行入侵检测. 模糊聚类<sup>[8]</sup>的思想是建立在 Zadeh 的模糊理论<sup>[8]</sup>的基础上, 同时结合了传统的聚类分析算法. 但是该算法事先要求确定分类个数, 因此需要利用确定最优分类个数的方法. L. Xu 在文章<sup>[3]</sup>中提出了一个确定最优分类个数的方法, 叫做对手惩罚竞争学习算法 (RPCL). 该算法自提出后得到了广泛的应用. 但是我们从文章<sup>[3]</sup>描述的算法以及试验中发现, 当数据集合含有数据量纲以及变量彼此之间具有相关性时, 该算法性能就会明显下降. 同时, 该算法对初始给定的聚类个数十分敏感. 这里, 本文针对以上情况对文章<sup>[3]</sup>提出的 RPCL 算法做了一些改进, 克服了原有 RPCL 算法的不足. 同时, 也解决了模糊聚类分析中预先确

定分类个数的问题.最后本文利用遗传算法对数据进行优化,提高了系统的检测精度.最终确定了非正常样本子空间的特征函数从而实现了检测模型.

## 2 模型

### 2.1 改进的 RPCL模型

香港中文大学 L. Xu教授在文献[3]中建议了一种竞争学习的思想,叫做对手惩罚竞争学习算法(RPCL),用来实现动态决定相应聚类个数的方法.它的基本思想是对于每一个输入不仅胜者单元的权值做更新以适应输入,同时它的竞争对手(次胜者)单元也要按一定的比例被降学习处理.这种算法自从推出后,得到了广泛的应用同时获得了良好的效果.但是我们从文章描述的算法以及试验中发现,当数据集合含有数据量纲以及变量彼此之间相关时,该算法性能就会明显下降.同时,该算法对初始给定的聚类个数十分敏感.这里,本文针对以上情况对文献[3]提出的RPCL算法做了一些改进,克服了上述不足.

首先,我们对于分类的数据进行处理,实现去除数据的量纲影响.从文献[3]中发现,原有的RPCL算法利用的是欧式距离来计算彼此间的相似性.这在数据集合中分类是一个圆形时候效果较好,但是当数据的聚类是不同方向椭圆时效果不理想,因此我们需要对数据进行变尺度处理,使原本为椭圆的聚类在新的坐标系下呈现圆形.具体方法如下:

$$\text{样本的均值: } \bar{x} = \frac{1}{N} \sum_{i=1}^N x_i \quad (1)$$

$$\text{样本的方差: } s_x^2 = \frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N-1} \quad (2)$$

样本的标准偏差: 方程(2)的平方根.

对向量中的每一个成员变量作如下标准化处理:

$$Z_i = \frac{x_i - \bar{x}}{s_x} \quad (3)$$

不难看出,变量作如下变化后我们对距离的计算相当于原始空间中的取  $D = \text{diag}(s_1^2, s_2^2, \dots, s_n^2)$  的马氏距离<sup>[7]</sup>.

其次,为了消除向量中变量间相关性的影响,我们需要对数据进行PCA线性变换<sup>[7]</sup>,使得坐标基之间线性无关.经过上述两个方面的处理实现了弥补原RPCL算法中的不足.不难发现,经过两步处理后的数据空间,进行欧式距离的计算相当于在原始空间中马氏距离的计算.

为了克服对初始聚类值敏感性的特点,本文对原始的RPCL进行了如下处理,改进后的RPCL算法如下:

初始化  $p = 0$

给定一个初始分类数值  $K$  ( $K > m$ ),  $m$  为实际的最优分类数值.

(1) 初始化权值向量集合  $\{w_i\}_{i=1}^K$ , 设置  $t = 0$   $t$  为权值向量更新的次数.

(2) 非正常样本数据集中随机选择  $x$ , 从  $i = 1, \dots, K$ :

$$u_i = \begin{cases} 1 & , i = c \\ -1 & , i = r \\ 0 & , \text{others} \end{cases} \quad (4)$$

这里,  $l_c \|x - w_c\|^2 = m_j \ln l_j \|x - w_c\|^2$  和  $l_r \|x - w_c\|^2 = m_j \ln \|x - w_c\|^2$ .

其中,  $l_i = n_i / \sum_{j=1}^K n_j$ , 这里  $c$  代表胜者单元,  $r$  代表次胜者单元,  $n_j$  代表类  $j$  的单元数.  $\|\cdot\|$  代表在变换后空间中的欧式距离.

(3) 更新  $w_{i+1}(t+1) = w_i(t) + \Delta w_i(t)$ , 这里

$$\Delta w_i(t) = \begin{cases} a_c(x - w_i) & , u_i = 1 \\ -a_r(x - w_i) & , u_i = -1 \\ 0 & , \text{other} \end{cases} \quad (5)$$

这里  $0 \leq a_r \leq a_c \leq 1$ , 分别为竞争对手和胜利单元的学习率.

(4) 对于每一个分类, 计算它的含有样本的个数. 如果

小于某一个阈值, 这里设置为  $r \sqrt{\frac{\sum_{i=1}^K (N_i - \bar{N})^2}{N}}$ ,  $\bar{N}$  为平均

每一个类含有的样本数目,  $0 < r < 1$  我们就删除这一分类. 如果删除的分类数目与本次未删除的分类数目的比值小于某一个阈值  $\xi$  记住本次的分类数范围. 同时  $p = p + 1$  如果  $p = 2$  退出循环. 如果  $p = 1$ , 重新设置新的初始值  $K'$ , 再次循环. 最后将两次循环的结果范围进行集合交运算并取其中间值作为最优分类数; 否则,  $t = t + 1$ , 返回(2).

### 2.2 Fuzzy聚类模型

由改进后的RPCL算法获得最优分类数和中心点集合后, 我们利用模糊聚类分析算法进行聚类计算. 其算法如下:

给定非正常样本的数据集:  $z_k = [z_{1k}, z_{2k}, \dots, z_{nk}]^T \in$

$R^n$ ,  $k = 1, \dots, N$ . 寻找: 分类矩阵:  $U = \begin{bmatrix} u_{11} & \dots & u_{1N} \\ \vdots & \vdots & \vdots \\ u_{c1} & \dots & u_{cN} \end{bmatrix}$  和聚

类中心:  $V = \{v_1, v_2, \dots, v_c\}$ ,  $v_i \in R^n$  目标函数:

$$J(Z; V; U; A) = \sum_{i=1}^c \sum_{j=1}^N u_{ij}^m d_{A_i}^2(z_j, v_i) \quad (6)$$

约束于:  $0 \leq u_{ij} \leq 1$   $i = 1, \dots, c$ ,  $j = 1, \dots, N$  成员函数度  $0 < \sum_{j=1}^N u_{ij} < 1$ ,  $i = 1, \dots, c$  没有聚类为空  $\sum_{i=1}^c u_{ij} = 1$ ,  $j = 1, \dots, N$  整个成员函数和循环运算:

(1) 计算聚类的原型:

$$v_i = \frac{\sum_{k=1}^N u_{ik}^m z_k}{\sum_{k=1}^N u_{ik}^m} \quad (7)$$

(2) 计算距离:

$$d_{ik} = (z_k - v_i)^T (z_k - v_i) \quad (8)$$

(3)更新分类矩阵

$$u_{ik} = \frac{1}{\sum_{j=1}^c (d_{ik} / d_{jk})^{1/(m-1)}} \tag{9}$$

(4)直到  $\| \Delta U \| < \varepsilon$

这里,  $m$  代表模糊度, 一般取值为  $m = 2$

最后, 通过计算每一类中的距离中心点最大的样本差值作为半径组成一个半径集合  $R = \{r_1, r_2, \dots, r_c\}$ ,  $r \in R$  以及确定的中心点集合  $V = \{v_1, v_2, \dots, v_c\}$ ,  $v_i \in R^n$  构成一个描述非正常行为属性的空间。

2.3 遗传算法确定模型

本文利用遗传算法确定半径偏差  $Rd = \{rd_1, rd_2, \dots, rd_c\}$  的选择。这里, 我们对半径集合  $R$  中每一个半径的优化范围分别设置成  $(0, r_i/2)$ 。对于每一个优化变量进行 8 进制编码。目标函数设置成:

$$fitness_k = volume(R) - C \cdot num\_error(R) \tag{10}$$

其中,  $volume(R)$  代表空间的体积:

$$volume(R) = \prod_{i=1}^c (r_i + rd_i) \tag{11}$$

$num\_error(R)$  = 出错样本数 / 总的测试样本数。C 代表惩罚因子。

然后设置每一个群体的个数, 和进化世代的数目, 经过选择、交叉、变异, 最终确定最优半径偏差集合  $Rd_{optimal}$ 。

2.4 探测器模型

正常子空间 (确定):  $S$  为特征向量集合,  $Self \subseteq S$  代表系统的正常状态。它的补空间  $Non\_Self$  被定义为  $Non\_Self = S - Self$  我们利用  $Self$  集合的特征函数表示

$$x_{self}(x): [x_{min}, x_{max}]^n \rightarrow \{0, 1\}$$
$$x_{self}(x) = \begin{cases} true, & x \in Self \\ false, & x \in Non\_Self \end{cases} \tag{12}$$

本文利用了免疫原理中的阴性选择算法, 即利用上述建议的模型得到的确定的中心点集合  $V = \{v_1, v_2, \dots, v_c\}$ ,  $v_i \in R^n$  和  $R = \{r_1, r_2, \dots, r_c\}$ ,  $r \in R$  以及半径偏差  $Rd = \{rd_1, rd_2, \dots, rd_c\}$  来确定非正常子空间特征函数。如果行为特征满足非正常样本空间的特征函数, 我们确认为一个入侵行为。非己空间的特征函数规则如下:

$R^1$ : If  $Cond_1$  then  $non\_self$

$\vdots$

$R^c$ : If  $Cond_c$  then  $non\_self$

这里,

(1)  $Cond_i = \|x - v_i\| \leq r_i + rd_i$

(2)  $x$  线性变换后空间中对应的特征向量,  $\|\cdot\|$  代表变换后的空间的欧式距离。每一个规则定义了变换后空间中的一个超球面, 这样这些规则的集合最大限度地利用超球面覆盖整个非正常样本空间。这样非正常样本集合的特征函数利用  $R = \{R^1, \dots, R^c\}$  定义如下:

$$x_{non-self}(x) = \begin{cases} 1, & \text{如果存在 } R^j \in R, \text{ 使得 } x \in R^j \\ 0, & \text{其它} \end{cases}$$

这里,  $x \in R^j$  表示满足规则  $R^j$  中的条件部分。

系统框图 1

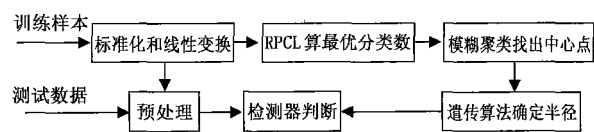


图 1 IDS 系统流程图

3 试验

我们利用了 1999 DARPA 入侵检测评估数据, 这里对数据进行了一些处理, 选择了正常的数据集 800 个、2200 个非正常数据集, 共 3000 个数据集来训练和测试我们的特征空间选取算法。我们针对该数据进行两次聚类实验, 在第一次实验中, 我们设置初始的分类数为 400, 最终满足条件为 [83 70]; 在第二次实验中, 我们设置初始的分类数为 500, 最终满足条件为 [78 66]; 这两个实验结果的交叉均值为 74。从表 1 中可以看出, 传统的 RPCL 的方法对初始值的敏感性, 同时通过采用 David-Bound-Index 最优分类的比较可知, 改进的 RPCL 方法克服了传统 RPCL 方法的不足并提高了分类的性能。

表 1 改进的 RPCL 分类试验结果比较

| 方法       | 初始分类数 | 分类范围    | David Index 值 |
|----------|-------|---------|---------------|
| 改进的 RPCL | 400   | [83 70] | 1.03          |
| 改进的 RPCL | 500   | [78 66] |               |
| 传统的 RPCL | 300   | 50      | 3.06          |
| 传统的 RPCL | 400   | 170     | 8.04          |
| 传统的 RPCL | 500   | 120     | 6.07          |

本文利用 DARPA 99 入侵检测评估数据中含 3000 个非正常样本的训练数据集, 利用改进的 RPCL 算法进行计算。结果如下:

表 2 改进后的 RPCL 试验结果

| 试验次数 | 分类次数 |     |     |    |    | 最优分类范围  |
|------|------|-----|-----|----|----|---------|
| 1    | 400  | 159 | 131 | 83 | 70 | [83 70] |
| 2    | 500  | 252 | 110 | 78 | 66 | [78 66] |

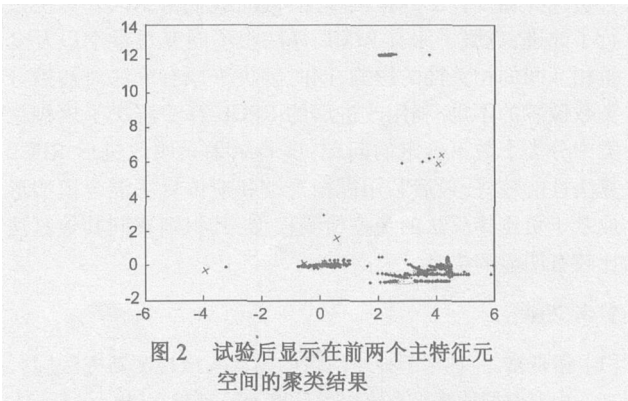


图 2 试验后显示在前两个主特征元空间的聚类结果

取两者的交集并求出中间值 74 作为最优分类解。同时利用模糊聚类算法计算出中心点和半径, 试验中发现当将经过改进的 RPCL 算法算出的中心点作为模糊聚类算

法的初始值时,很快收敛. 试验中的循环次数为 18 次,分类结果如图 2 所示.

然后我们利用本文中提出的遗传算法决定半径的偏差,群体的个数为 50 进化世代的数目 400 为了运算方便,这里我们将目标函数取负值,即求得最小值. 试验结果如图 3 所示.

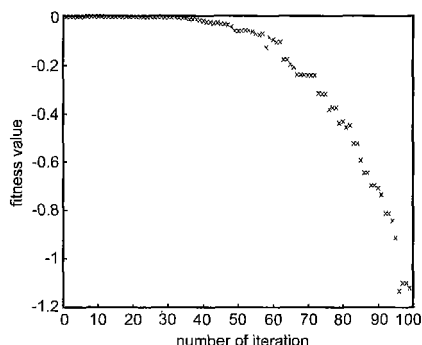


图 3 遗传算法迭代结果

利用上面获得的中心点集合  $V = \{v_1, v_2, \dots, v_c\}$ ,  $v_i \in R^n$  和  $R = \{r_1, r_2, \dots, r_c\}$ ,  $r \in R$  以及半径偏差  $Rd = \{rd_1, rd_2, \dots, rd_c\}$  的试验数据来确定非正常子空间的特征函数,形成检测模型.

以下是对检测模型的试验结果:

表 3 检测试验结果

| 数据集 识别率 (%)                                     | MLP   | 建议方法      |
|---|-------|-----------|
| 1500 个正常样本                                      | 0.87  | 0.998     |
| 含 1331 个 netpunch, satan, portswEEP, ipswEEP 样本 | 0.823 | 0.9082707 |

从上面的试验结果我们发现,识别率很高.同时,试验中发现对非正常样本的识别率低于正常样本的识别率,这里可以通过加大非正常学习样本的数量来进一步提高系统的辨识能力.

## 4 结论

本文提出了一种新的基于模糊聚类 and 遗传算法相结合负选择免疫算法的检测模型,同时我们对 RPCL 算法进行了改进,克服了原有 RPCL 算法中不同变量量纲以及变量相互间的相关性的影响,同时解决了原有算法对初始分类数敏感的不足.利用改进后的 RPCL 算法解决了模糊聚类中分类个数不确定的问题,试验结果表明改进后 RPCL 算法性能较好.最后利用模糊聚类和遗传算法混合模型形成基于负选择算法的免疫检测模型,其识别率同其他算法比较有明显的提高.

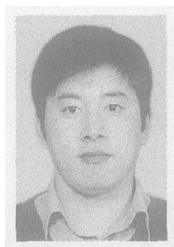
## 参考文献:

- [1] 伍良富. 一种基于神经网络的黑客入侵检测新方法 [J]. 小型微型计算机系统. 2003 08(04): 132-138  
Wu Liangfu A new method for internet intrusion detection

based on neural networks [J]. MiniMicro Systems 2003, 08(04): 132-138 (in Chinese)

- [2] 肖道举, 毛辉, 陈晓苏. BP 神经网络在入侵检测中的应用. 华中科技大学学报 (自然科学版) [J]. 2003, 05(6): 201-204  
Xiao Daoju Mao Hui Chen Xiaosu The application of BP-neural network in the intrusion detection [J]. Journal of Huazhong University of Science and Technology (Nature Science Edition), 2003 05(6): 201-204 (in Chinese)
- [3] L Xu Krzyzak E Oja Rival penalized competitive learning for clustering analysis RBF net and curve detection [J]. IEEE Transactions on Neural Networks 1993, 4(4): 636-649
- [4] 王景新, 戴葵, 宋辉. 基于神经网络的入侵检测系统. 计算机工程与科学, 2003 06(09): 723-726  
Wang Jingxin Dai Kui Song Hui Intrusion detection based on neural networks [J]. Computer Engineering and Science 2003 06(09): 723-726 (in Chinese)
- [5] 王丽娜, 董晓梅, 于戈, 王东. 基于进化神经网络的入侵检测方法 [J]. 东北大学学报 (自然科学版), 2002, 02(06): 540-543  
Wang Lina Dong Xiaomei Yu Ge Wang Dong Method of evolutionary neural network-based intrusion detection [J]. Journal of Northeastern University (Natural Science), 2002 02(06): 540-543 (in Chinese)
- [6] L Xu Rival penalized competitive learning finite mixture and multisets clustering [A]. Intentional Joint Conference on Neural Networks [C]. Alaska Anchorage Press 1998 2525-2530
- [7] M W Mak C K Li X Li Maximum likelihood estimation of elliptical basis function parameters with applications to speaker verification [A]. International Conference on Signal Processing [C]. Beijing Tsinghua University Press 1998 1287-1290
- [8] Zadeh L A. Fuzzy Sets [M]. MIT: Information and Control 1965, 8 338-353

## 作者简介:



陶新民 男, 1973 年生, 博士, 副教授, 主要研究方向为网络与信息安全.  
E-mail: taoxinmin@163.com

陈万海 男, 1963 年生, 教授, 主要研究方向为遥测遥感、扩频通信、深空通信和信息安全.