



文章编号:1001-9081(2018)07-1941-05

DOI:10.11772/j.issn.1001-9081.2018010178

## 基于 SMOTE 和深度信念网络的异常检测

沈学利, 覃淑娟\*

(辽宁工程技术大学 电子与信息工程学院, 辽宁 葫芦岛 125105)

(\*通信作者电子邮箱 qinshujuanup@163.com)

**摘要:**针对现有海量非平衡数据集中少数类别样本入侵检测率低的问题,提出了一种基于合成少数类过采样技术(SMOTE)和深度信念网络(DBN)的异常检测(SMOTE-DBN)方法。首先,用SMOTE技术增加了少数类别样本的样本数;然后在预处理后的较平衡数据集上,用非监督的受限玻尔兹曼机(RBM)对预处理后的高维数据进行特征降维;其次,用反向传播(BP)算法微调模型参数,获得预处理后数据的较优低维表示;最后通过softmax分类器对较优低维数据进行分类。KDD1999数据集仿真实验表明,SMOTE优化处理能够提高模型对少数类别样本的检测率,在相同数据集上,SMOTE-DBN方法与DBN方法、支持向量机(SVM)方法相比,检测率分别提高了3.31个百分点和7.34个百分点,误报率分别降低了1.11个百分点和2.67个百分点。

**关键词:**合成少数类过采样技术;深度信念网络;受限玻尔兹曼机;逻辑回归;入侵检测

中图分类号: TP393.08 文献标志码:A

### Anomaly detection based on synthetic minority oversampling technique and deep belief network

SHEN Xueli, QIN Shujuan\*

(School of Electronics and Information Engineering, Liaoning Technical University, Huludao Liaoning 125105, China)

**Abstract:** To solve low detection rate problem of intrusion for a small number of samples in mass unbalanced datasets, an anomaly detection based on Synthetic Minority Oversampling Technique (SMOTE) and Deep Belief Network (DBN), called SMOTE-DBN method, was proposed. Firstly, SMOTE technology was used to increase the number of samples in minority categories. Secondly, on the preprocessed balanced data set, the dimensionality of the preprocessed high-dimensional data was reduced by unsupervised Restricted Boltzmann Machine (RBM). Thirdly, the model parameters were finely tuned by Back Propagation (BP) algorithm to obtain the lower low-dimensional representation of the preprocessed data. Finally, softmax classifier was used to classify the optimal low-dimensional data. The simulation experiment results on KDD1999 show that, compared with DBN method and Support Vector Machine (SVM) method, the detection rate of SMOTE-DBN method is increased by 3.31 and 7.34 percentage points respectively, and the false alarm rate is decreased by 1.11 and 2.67 percentage points respectively.

**Key words:** Synthetic Minority Oversampling Technique (SMOTE); Deep Belief Network (DBN); Restricted Boltzmann Machine (RBM); Logistic Regression (LR); intrusion detection

### 0 引言

随着网络规模的日益扩大和网络攻击的日益增加,入侵检测(Intrusion Detection, ID)依然是人们研究的热点之一。为了提高入侵检测系统(Intrusion Detection System, IDS)对未知网络攻击的识别能力和用户数据的关联性分析能力,许多研究学者将机器学习的方法引入到入侵检测系统中<sup>[1-2]</sup>,如支持向量机(Support Vector Machine, SVM)<sup>[3-5]</sup>在处理小样本数据集时检测率较高,但是由于其时间复杂度(为 $O(n^3)$ )和空间复杂度(为 $O(n^2)$ )的局限性,处理海量数据集时性能较差;神经网络(Neural Network, NN)<sup>[6-7]</sup>在一定程度上具有适应性和可扩展性,但是处理海量数据集时鲁棒性仍有待提高;深度学习(Deep Learning, DL)<sup>[8-10]</sup>能够挖掘高维数据的潜在特征,分类识别能力较强,但是现有的方法没有考虑到少数类别样本的入侵检测问题。而在入侵检测系统中,把提权(User to Root, U2R)攻击识别为正常用户数据,比把拒绝服

务(Denial of Service, DoS)攻击识别为正常用户数据对系统的危害更大,因此,在保证较高检测率和较低误报率的基础上,识别并阻断少数类别的攻击有着重要的现实意义。

针对上述问题,本文提出了一种基于合成少数类过采样技术(Synthetic Minority Oversampling Technique, SMOTE)和深度信念网络的异常检测(anomaly detection based on SMOTE and Deep Belief Network, SMOTE-DBN)方法,在保证其他类别样本检测率的前提下,能够提高少数类别样本的检测率,同时降低误报率。

### 1 异常检测框架

基于SMOTE和深信网的异常检测框架包含三部分内容,如图1所示。

1) 数据预处理。如图1(a)部分所示,通过合成少数类过采样技术(SMOTE)降低数据集的不平衡度,再将数据集中的符号型特征数据数值化,并对数据型特征数据进行归一化

收稿日期:2018-01-21;修回日期:2018-03-27;录用日期:2018-03-28。 基金项目:国家自然科学基金资助项目(61602227)。

作者简介:沈学利(1969—),男,江苏连云港人,教授,硕士,主要研究方向:信息安全、网络安全;覃淑娟(1993—),女,新疆塔城人,硕士研究生,主要研究方向:网络安全。



处理,详见3.1节内容。

2) 数据特征降维。如图1(b)部分所示,将预处理后的数据集用深度信念网络(Deep Belief Network, DBN)进行特征抽取,先用受限玻尔兹曼机(Restricted Boltzmann Machine, RBM)对数据集自底向上进行预训练,获得模型的初始参数,再用BP(Back Propagation)网络微调模型参数,获得较优模型参数,更好地将原始高维数据映射至低维数据,详见2.3、2.4节内容。

3) 逻辑回归(Logistic Regression, LR)分类器。如图1(c)部分所示,通过softmax逻辑回归分类器,对较优低维数据集进行5种用户数据状态的识别,详见2.5节内容。

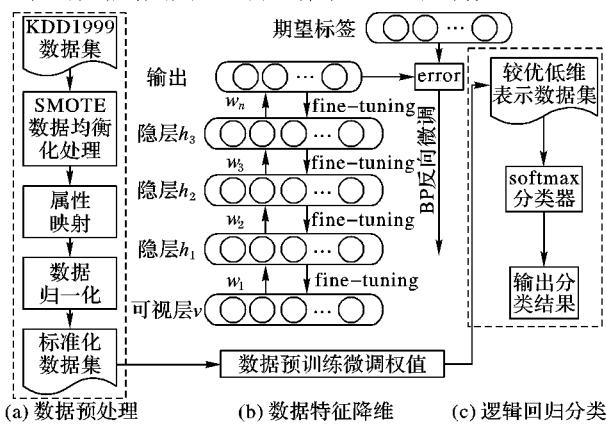


图1 基于SMOTE-DBN模型的异常检测框架

Fig. 1 Anomaly detection framework based on SMOTE-DBN model

## 2 相关算法

### 2.1 SMOTE

SMOTE算法是一种典型的过取样方法<sup>[11-12]</sup>,它用少数类样本控制人工样本的生成与分布,实现均衡数据集的目的。核心思想是在某少数类别样本中随机地选取一个样本点,并在其最近邻的k个样本之间,插入n个人工合成的少数类别样本,从而增加少数类别样本的数量,均衡化数据集。

由于入侵检测基准数据集中的数据分布很不均匀,现有检测方法对少数类别样本的检测率很低<sup>[13-14]</sup>,因此采用SMOTE方法来消除非均衡样本集对检测精度的影响。

此外,由于SMOTE选取样本的随机性,可能会选取在样本集边缘的样本点进行近邻插值,引起模糊样本边界的问题。为了避免新合成的样本点具有极少的少数类样本特征,致使数据集的原始分布改变,要尽可能地选取不在样本边缘的样本点,K-means算法<sup>[15]</sup>能有效解决这个问题。用K-means计算出样本点的簇心m,选取簇心的k个近邻进行插值操作,得出新样本x<sub>new</sub>:

$$x_{\text{new}} = m + \text{rand}(0,1) * (x - m) \quad (1)$$

其中x为簇心m的近邻样本,rand(0,1)表示0~1的随机数。

插入新样本的步骤如算法1所示。

算法1 插入新样本。

输入:原始少数类训练样本集,需要合成的样本数n,循环变量k。

输出:少数类训练样本集。

for t = 1, 2, ..., k

用K-means算法记录少数类样本的簇心m

for i = 1, 2, ..., n

随机选取簇心m的近邻样本点x,用式(1)在x与m之间插入新的样本点

end

### 2.2 DBN模型

DBN<sup>[16]</sup>是由若干层非监督的RBM网络和单层BP神经网络构成的深层神经网络。训练模型的主要步骤如下:

1) 用对比分歧(Contrastive Divergence, CD)算法<sup>[17]</sup>单独无监督地训练每一层RBM网络,确保特征向量映射到不同特征空间时,尽可能多地保留特征信息。

2) BP网络接收RBM的低维输出特征向量作为输入特征向量,有监督地训练实体关系分类器。由于每一层RBM网络只能确保自身层内的权值对该层特征向量映射达到最优,并不是对整个DBN的特征向量映射达到最优,所以反向传播网络将错误信息自顶向下传播至每一层RBM,微调整整个DBN。RBM网络训练模型的过程可以看作对一个深层BP网络权值参数的初始化,使DBN克服了BP网络因随机初始化权值参数而容易陷入局部最优和训练时间长的缺点。

通过上述步骤,能够构建出具有多隐藏层的非线性网络结构,挖掘海量数据集的潜在特征,从而学习出高维数据的较优低维表示,得到更易分类的入侵检测数据特征。

### 2.3 预训练

RBM是DBN的核心模块之一<sup>[18]</sup>,由可见层单元( $v$ )和隐藏层单元( $h$ )构成。可见层与隐藏层的层内无连接,层级之间全连接。如图2所示,可见层单元为 $v = (v_1, v_2, \dots, v_m)$ 描述输入数据的特征,隐藏层单元为 $h = (h_1, h_2, \dots, h_n)$ ,通过学习输入数据的特征自动生成。

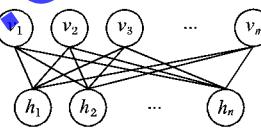


图2 RBM结构

Fig. 2 Structure of RBM

已知 $v$ 的情况下,隐藏层节点的条件概率分布满足:

$$\begin{cases} P(h_j = 1 | v) = e^{c_j + W_{jv}v_i} / (1 + e^{c_j + W_{jv}v_i}) = \text{sig}\left(\sum_i W_{ji}v_i + c_j\right) \\ P(h_j = 0 | v) = 1 - P(h_j = 1 | v) \end{cases} \quad (2)$$

同理,在已知 $h$ 的情况下,可见层节点的条件概率分布满足:

$$\begin{cases} P(v_i = 1 | h) = \text{sig}\left(\sum_j w_{ij}h_j + b_i\right) \\ P(v_i = 0 | h) = 1 - P(v_i = 1 | h) \end{cases} \quad (3)$$

关于RBM建立的能量函数为:

$$E(v, h | \theta) = - \sum_{i=1}^m \sum_{j=1}^n W_{ji}v_ih_j - \sum_{i=1}^m b_i v_i - \sum_{j=1}^n c_j h_j \quad (4)$$

其中: $\theta = \{W, b, c\}$ 为RBM的模型参数, $W$ 为可见层到隐藏层间的权值连接矩阵, $b$ 和 $c$ 分别表示可见层和隐藏层上的乘性偏置。

基于能量函数,可以建立 $v, h$ 的联合分布函数:

$$\begin{cases} P(v, h) = e^{-E(v, h | \theta)} / Z \\ Z = \sum_{v, h} e^{-E(v, h | \theta)} \end{cases} \quad (5)$$

为了求得联合概率分布的最大值,更新模型参数,本文采用CD算法获取样本。首先,初始可见单元状态被设置为一个训练样本,并由初始可见单元 $v$ 学习得到第一层隐藏层单元 $h^1$ ,获得后验概率 $P(h^1 | v)$ 。再由隐藏层单元 $h^1$ 确定每个可见单元取值为1的概率,重构获得新的可见层单元 $v'$ 。接着采用梯度下降法求解参数,训练样本的梯度为:



$$\begin{cases} \frac{\partial \ln P(\mathbf{v}, \mathbf{h})}{\partial \boldsymbol{\theta}} = \sum_{\mathbf{h}} P(\mathbf{h} \mid \mathbf{v}, \boldsymbol{\theta}) D - \sum_{\mathbf{v}, \mathbf{h}} P(\mathbf{h} \mid \mathbf{v}, \boldsymbol{\theta}) D \\ D = \frac{\partial E(\mathbf{v}, \mathbf{h} \mid \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \end{cases} \quad (6)$$

获得模型参数的更新规则:

$$\boldsymbol{\theta}^{n+1} = \boldsymbol{\theta}^n + \varepsilon \frac{\partial \ln P(\mathbf{v}, \mathbf{h})}{\partial \boldsymbol{\theta}} \quad (7)$$

预训练过程如算法 2 所示。

算法 2 预训练过程。

输入: 可见层特征变量  $\mathbf{v} = (v_1, v_2, \dots, v_m)$ , 初始权重  $\mathbf{W}$ , 乘性偏置  $\mathbf{b}, \mathbf{c}$ , 学习率  $\varepsilon$ , 迭代次数  $k$ 。

输出: RBM 的模型参数  $\boldsymbol{\theta} = \{\mathbf{W}, \mathbf{b}, \mathbf{c}\}$ 。

for  $t = 1, 2, \dots, k$

for  $j = 1, 2, \dots, n$

由式(2)计算每个隐藏层的特征向量的值  $h_j^{(t)} = P(h_j \mid \mathbf{v}^{(t)})$   
for  $i = 1, 2, \dots, m$

由式(3)计算每个可见层的特征向量的值  $v_i^{(t+1)} = P(v_i \mid h^{(t)})$

for  $i = 1, 2, \dots, m$

for  $j = 1, 2, \dots, n$

由式(7)更新模型参数  $\boldsymbol{\theta} = \{\mathbf{W}, \mathbf{b}, \mathbf{c}\}$

end

## 2.4 BP 微调权重

BP 神经网络是带监督的前馈神经网络<sup>[19]</sup>, 有监督的训练经过预训练的 DBN 模型, 利用输出误差自顶向下地估计每一层 RBM 的传播误差, 获得更优的模型参数。BP 微调权重过程如算法 3 所示。

算法 3 BP 微调权重过程。

输入: 可见层特征变量  $\mathbf{v} = (v_1, v_2, \dots, v_m)$ , 预训练得到的模型参数  $\boldsymbol{\theta} = \{\mathbf{W}, \mathbf{b}, \mathbf{c}\}$ , 迭代次数  $k$ , 学习率  $\varepsilon$ 。

输出: 微调后的模型参数  $\boldsymbol{\theta} = \{\mathbf{W}, \mathbf{b}, \mathbf{c}\}$ 。

for  $t = 1, 2, \dots, k$

对所有  $v_i$  的输出单元  $o_i$ , 计算其误差梯度  $\sigma_i$  ( $e_i$  为期望输出):  
 $\sigma_i = o_i(1 - o_i)(e_i - o_i)$  (8)

对所有隐藏层单元  $h_j$ , 计算其误差梯度  $\sigma_j$ , 并更新模型参数  $\boldsymbol{\theta}$ :

$\sigma_j = o_j(1 - o_j) \sum_i \theta_{ij} \sigma_i$  (9)

$\theta_{ji}^{t+1} = \theta_{ji}^t + \varepsilon \theta_{ji}$  (10)

end

## 2.5 softmax 分类器

测试数据集中有五种用户数据状态<sup>[20~21]</sup>, 分别为正常状态(Normal)、拒绝服务(Denial of Service, DoS)攻击、远程未授权(Remote to Local, R2L)攻击、提权(User to Root, U2R)攻击、端口扫描(Probing), 依序标记为 1~5, 如表 1 所示。

由表 1 可知, 数据集中有多类用户数据状态, 而 softmax 分类器能够适应多分类问题, 且相较于 SVM 等分类器结构简单, 因此, 构建一个 softmax 分类器, 对训练后获得的较优低维表示的数据进行分类。

如式(11)所示, 对测试数据集进行五种用户数据状态的识别:

$$h_{\boldsymbol{\theta}'}(\mathbf{x}'^{(i)}) = \begin{bmatrix} P(y^{(i)} = 1 \mid \mathbf{x}'^{(i)}; \boldsymbol{\theta}') \\ P(y^{(i)} = 2 \mid \mathbf{x}'^{(i)}; \boldsymbol{\theta}') \\ \vdots \\ P(y^{(i)} = 5 \mid \mathbf{x}'^{(i)}; \boldsymbol{\theta}') \end{bmatrix} = \frac{1}{\sum_{j=1}^5 e^{\boldsymbol{\theta}_j^T \mathbf{x}'^{(i)}}} \begin{bmatrix} e^{\boldsymbol{\theta}_1^T \mathbf{x}'^{(i)}} \\ e^{\boldsymbol{\theta}_2^T \mathbf{x}'^{(i)}} \\ \vdots \\ e^{\boldsymbol{\theta}_5^T \mathbf{x}'^{(i)}} \end{bmatrix} \quad (11)$$

其中:  $\boldsymbol{\theta}' = \{\mathbf{W}', \mathbf{b}'\}$  为模型参数,  $\mathbf{W}'$  表示权值矩阵,  $\mathbf{b}'$  表示加性偏置。

表 1 测试数据集分布

Tab. 1 Distribution of test dataset

攻击类别	攻击子类型	标签
Normal	'Normal'	1
DoS	Back, Land, Neptun, Pod, Smurf, Teardrop, Mailbomb, Processtable, Udpstorm, Apache2, Worm	2
R2L	Guess_passwd, Ftp_write, Imap, Phf, Multihop, Warezmaster, Warezclient, Xlock, Xsnoop, Snmpguess, Snmpgetattack, Httpstunnel	3
U2R	Satan, IPsweep, Nmap, Portsweep, Mscan, Saint	4
Probing	Buffer_overflow, Loadmodul, Rootkit, Perl, Sqattack, Xterm, Ps	5

将要分类的较优低维数据  $\mathbf{x}'$  输入到一套超平面中, 每个超平面代表一个类, 以输入的数据到第  $j$  类超平面的距离表示该数据属于第  $j$  类的概率, 概率最大的类即为数据的所属类别:

$$P(y = j \mid \mathbf{x}'^{(i)}, \boldsymbol{\theta}') = \text{softmax}_j(\mathbf{W}' \mathbf{x}'^{(i)} + \mathbf{b}')$$
 (12)

## 3 实验验证

### 3.1 数据预处理

本文采用 KDD 1999 数据集<sup>[22]</sup>作为测试数据集。该数据集中的每一项数据共有 41 项特征属性和 1 项标签属性, 特征属性包括传输控制协议(Transmission Control Protocol, TCP)基本连接特征(No. 1~No. 9)、TCP 连接的内容特征(No. 10~No. 22)、基于时间的网络流量特征(No. 23~No. 31)以及基于主机的网络流量统计特征(No. 32~No. 41), 特征属性的类型分别为连续型(Continuous, C)和离散型(Symbolic, S)<sup>[23]</sup>, 如表 2 所示。实验所选取的数据集如表 3 所示。

表 2 数据集特征

Tab. 2 Features of dataset

类型	特征
C	duration(1), src_bytes(5), dst_bytes(6), wrong_fragment(8), urgent(9), host(10), num_failed_logins(11), num_compromised(13), root_shell(14), su_attempted(15), num_root(16), num_file_creations(17), num_shells(18), num_access_files(19), num_outbound_cmds(20), count(23), srv_count(24), serror_rate(25), srv_serror_rate(26), rerror_rate(27), srv_rerror_rate(28), same_srv_rate(30), srv_diff_host_rate(31), dst_host_count(32), dst_host_srv_count(33), dst_host_same_srv_rate(34), dst_host_diff_srv_rate(35), dst_host_same_src_port_rate(36), dst_host_srv_diff_host_rate(37), dst_host_serror_rate(38), dst_host_srv_serror_rate(39), dst_host_rerror_rate(40), dst_host_srv_rerror_rate(41)
S	protocol_type(2), service(3), flag(4), land(7), logged_in(12), is_host_login(21), is_guest_login(22)

数据预处理分 3 个步骤。

1) 降低数据集的不平衡度。

由表 3 可知, KDD 1999 数据集的数据状态分布很不均衡, 训练集中样本 U2R 的数量远小于 DoS 和 Normal 的样本数, 因此, 本文采用 SMOTE 技术, 将 U2R 的样本数增大至原来的 10 倍, 以均衡样本数。

2) 字符型特征数值化。

用属性映射法将字符型特征数据数值化, 分别为 TCP、用



户数据报协议(User Datagram Protocol, UDP)、网际控制报文协议(Internet Control Message Protocol, ICMP),如表4所示。

### 3) 数据型特征归一化。

将数值化后的数据取对数,再根据式(13)归一化到[0,1]区间内:

$$y = (y - \min) / (\max - \min) \quad (13)$$

其中: $y$ 为属性值, $\min$ 为对应特征属性的最小值, $\max$ 为对应特征属性的最大值。

表3 实验数据集

Tab. 3 Experimental dataset

攻击类别	样本数		攻击类别	样本数	
	训练集	测试集		训练集	测试集
Normal	4 000	3 000	U2R	50	200
DoS	2 000	2 000	Probing	2 000	2 000
R2L	1 000	1 000			

表4 字符型特征数值化

Tab. 4 Attribute mapping of character type

协议类型	二进制	协议类型	二进制	协议类型	二进制
TCP	(1,0,0)	ICMP	(0,0,1)	UDP	(0,1,0)

## 3.2 实验评价标准

评价标准定义如下。

TP(True Positive):样本正确判断为正类的样本数。

TN(True Negative):样本正确判断为负类的样本数。

FP(False Positive):样本错误判断为负类的实际正类样本数。

FN(False Negative):样本错误判断为正类的实际负类样本数。

则检测率(Detection Rate, DR)、误报率(False Alarm, FA)、精确率(Accuracy, AC)分别如下:

$$DR = TN / (TN + FN) \quad (14)$$

$$FA = FP / (TP + FP) \quad (15)$$

$$AC = (TP + TN) / (TP + FP + TN + FN) \quad (16)$$

## 3.3 实验分析

实验环境:Windows 7(64位)操作系统,Intel Core i5-5200U CPU @ 2.2 GHz,4 GB RBM,Python3.5。

实验内容:

- 1) 设置实验参数。
- 2) 在相同分类方法的基础上验证SMOTE技术对异常入侵检测的影响。
- 3) 在相同数据集上分析不同分类技术对异常入侵检测的影响。

### 3.3.1 实验参数设置

实验过程中,用DBN对选取的数据集进行训练,由于DBN的参数设置会影响到模型的训练结果,根据文献[24-25]对模型的部分参数进行了设置,训练参数如表5所示,同时通过固定参数,验证了微调的迭代次数对检测率结果的影响,如图3所示。

由图3可知,当迭代次数大于100时,精确率曲线逐渐趋于平缓。为了避免过拟合,后续实验中选取微调的迭代次数为100。

### 3.3.2 SMOTE 算法的有效性验证

为了验证SMOTE算法的有效性,将经过SMOTE技术处

理前后的数据集在DBN算法上进行验证。实验结果表明,经过SMOTE预处理的数据集相较于未经过SMOTE的数据集在精确率方面提高了2.01个百分点,检测率结果如图4所示,DoS的检测率有所降低,但是对少数类样本U2R的检测率有明显提高,其他类别样本的检测率与未经过SMOTE处理的数据集检测率相当。

表5 实验参数

Tab. 5 Experimental parameters

实验参数	值	实验参数	值
DBN输入层节点数	122	预处理阶段迭代次数	20
第一层隐藏层节点数	92	预处理阶段学习率	0.001
第二层隐藏层节点数	70	微调阶段学习率	0.001
第三层隐藏层节点数	50	softmax类别数k	5
DBN输出层节点数	10		

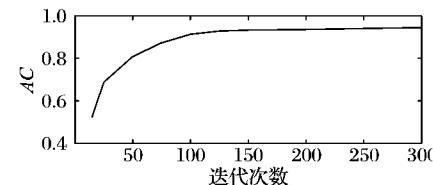


图3 精确率随微调迭代次数的变化

Fig. 3 Accuracy curve changing with number of fine-tuning iterations

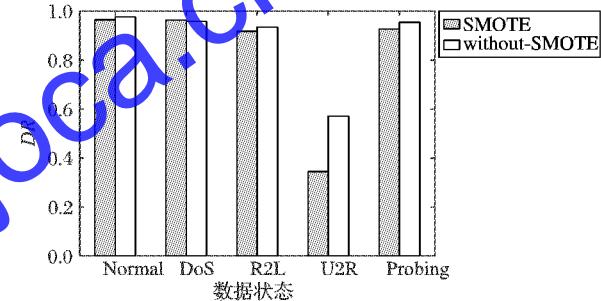


图4 SMOTE 处理前后检测率对比

Fig. 4 Comparison of detection rate before and after SMOTE processing

### 3.3.3 对比实验

将SMOTE-DBN方法与DBN和SVM方法在相同数据集上进行对比实验,如表6所示,SMOTE-DBN方法的检测率相对略高于DBN和SVM方法,且在误报率方面相对较低。

表6 SMOTE-DBN与DBN、SVM实验结果对比 %

Tab. 6 Comparison of experimental results among SMOTE-DBN, DBN and SVM %

实验方法	DR	FA	实验方法	DR	FA
SMOTE-DBN	93.67	2.36	SVM	86.33	5.03
DBN	90.36	3.47			

## 4 结语

本文提出了一种基于SMOTE和深度信念网络的异常检测方法,提高了入侵检测的数据分析能力。通过SMOTE处理技术,均衡化非均衡数据集,在一定程度上解决了分类器倾向于将用户数据归类到多数类类别样本的问题。同时结合softmax算法改进了DBN算法,并与DBN和SVM方法进行对比实验。实验结果表明,SMOTE-DBN算法的性能相对较优,对高维数据有很强的特征提取能力和信息识别能力,可应用于网络分布复杂的环境下;但DBN中的结构参数为人工设



置,不一定是最优的网络结构,因此如何选取合理的网络参数是下一步解决的问题。

#### 参考文献 (References)

- [1] TSAI C F, HSU Y F, LIN C Y, et al. Intrusion detection by machine learning: a review [J]. *Expert Systems with Applications*, 2009, 36(10): 11994–12000.
- [2] DAS S, NENE M J. A survey on types of machine learning techniques in intrusion prevention systems [C]// Proceedings of the 2017 International Conference on Wireless Communications, Signal Processing and Networking. Piscataway, NJ: IEEE, 2017: 2296–2299.
- [3] CHAND N, MISHRA P, KRISHNA C R, et al. A comparative analysis of SVM and its stacking with other classification algorithm for intrusion detection [C]// Proceedings of the 2016 International Conference on Advances in Computing, Communication & Automation. Piscataway, NJ: IEEE, 2016: 1–6.
- [4] ABUROMMAN A A, REAZ M B I. Ensemble of binary SVM classifiers based on PCA and LDA feature extraction for intrusion detection [C]// Proceedings of the 2017 Advanced Information Management, Communicates, Electronic and Automation Control Conference. Piscataway, NJ: IEEE, 2017: 636–640.
- [5] TENG S, WU N, ZHU H, et al. SVM-DT-based adaptive and collaborative intrusion detection [J]. *IEEE/CAA Journal of Automatica Sinica*, 2017, 5(1): 108–118.
- [6] DENG C, QIAO H. Network security intrusion detection system based on incremental improved convolutional neural network model [C]// Proceedings of the 2017 International Conference on Communication and Electronics Systems. Piscataway, NJ: IEEE, 2017: 1–5.
- [7] 杨雅辉, 黄海珍, 沈晴霓, 等. 基于增量式 GHSOM 神经网络模型的入侵检测研究 [J]. *计算机学报*, 2014, 37(5): 1216–1224. (YANG Y H, HUANG H Z, SHEN Q N, et al. Research on intrusion detection based on incremental GHSOM [J]. *Chinese Journal of Computers*, 2014, 37(5): 1216–1224.)
- [8] 杨昆朋. 基于深度学习的入侵检测[D]. 北京: 北京交通大学, 2015: 31–47. (YANG K P. *Intrusion detection based on deep Learning* [D]. Beijing: Beijing Jiaotong University, 2015: 31–47.)
- [9] GAO N, GAO L, HE Y, et al. Intrusion detection model based on deep belief nets [J]. *Journal of Southeast University (English Edition)*, 2015, 31(3): 339–346.
- [10] ALOM M Z, BONTUPALLI V R, TAH A T M. Intrusion detection using deep belief networks [C]// Proceedings of the 2016 Aerospace and Electronics Conference. Piscataway, NJ: IEEE, 2016: 339–344.
- [11] 霍玉丹, 谷琼, 蔡之华, 等. 基于遗传算法改进的少数类样本合成过采样技术的非平衡数据集分类算法 [J]. *计算机应用*, 2015, 35(1): 121–124. (HUO Y D, GU Q, CAI Z H, et al. Classification method for imbalance dataset based on genetic algorithm improved synthetic minority over-sampling technique [J]. *Journal of Computer Applications*, 2015, 35(1): 121–124.)
- [12] DEMIDOVA L, KLYUEVA I. SVM classification: optimization with the SMOTE algorithm for the class imbalance problem [C]// Proceedings of the 2017 Embedded Computing. Piscataway, NJ: IEEE, 2017: 1–4.
- [13] ALRAWASHDEH K, PURDY C. Toward an online anomaly intrusion detection system based on deep learning [C]// Proceedings of the 2017 International Conference on Machine Learning and Applications. Piscataway, NJ: IEEE, 2017: 195–200.
- [14] POTLURI S, DIEDRICH C. Accelerated deep neural networks for enhanced intrusion detection system [C]// Proceedings of the 2016 International Conference on Emerging Technologies and Factory Automation. Piscataway, NJ: IEEE, 2016: 1–8.
- [15] 金建国. 聚类方法综述 [J]. *计算机科学*, 2014, 41(b11): 288–293. (JIN J G. Review of clustering method [J]. *Computer Science*, 2014, 41(b11): 288–293.)
- [16] HINTON G E, OSINDERO S, TEH Y W. A fast learning algorithm for deep belief nets [J]. *Neural Computation*, 2006, 18(7): 1527–1554.
- [17] HINTON G. Training products of experts by minimizing contrastive divergence [J]. *Neural Computation*, 2002, 14(8): 1771–1800.
- [18] FISCHER A, IGEL C. An introduction to restricted Boltzmann machines [C]// CIARP 2012: Proceedings of the 17th Iberoamerican Congress on Progress in Pattern Recognition, Image Analysis, Computer Vision, and Applications. Berlin: Springer, 2012: 14–36.
- [19] ZHANG H, LI B. Application of an improved multi-layer BP neural network algorithm in intrusion detection [C]// Proceedings of the 2016 Sixth International Conference on Instrumentation & Measurement, Computer, Communication and Control. Piscataway, NJ: IEEE, 2016: 619–622.
- [20] 孔令智. 基于网络异常的入侵检测算法研究 [D]. 北京: 北京交通大学, 2017: 38–39. (KONG L Z. Research on intrusion detection algorithm based on network anomaly [D]. Beijing: Beijing Jiaotong University, 2017: 38–39.)
- [21] 陈虹, 万广雪, 肖振久. 基于优化数据处理的深度信念网络模型的入侵检测方法 [J]. *计算机应用*, 2017, 37(6): 1636–1643. (CHEN H, WAN G X, XIAO Z J. Intrusion detection method of deep belief network model based on optimization of data processing [J]. *Journal of Computer Applications*, 2017, 37(6): 1636–1643.)
- [22] STOLFO S J, FAN W, LEE W, et al. Cost-based modeling for fraud and intrusion detection: results from the JAM project [C]// DISCEX '00: Proceedings of the 2000 DARPA Information Survivability Conference and Exposition. Piscataway, NJ: IEEE, 2000: 130–144.
- [23] YIN C L, ZHU Y F, FEI J L, et al. A deep learning approach for intrusion detection using recurrent neural networks [J]. *IEEE Access*, 2017, 5: 21954–21961.
- [24] GAO N, GAO L, GAO Q, et al. An intrusion detection model based on deep belief networks [C]// Proceedings of the 2nd International Conference on Advanced Cloud and Big Data. Washington, DC: IEEE Computer Society, 2014: 247–252.
- [25] 高妮, 贺毅岳, 高岭. 海量数据环境下用于入侵检测的深度学习方法 [J]. *计算机应用研究*, 2018, 35(4): 1197–1200. (GAO N, HE Y Y, GAO L. Deep learning method for intrusion detection in massive data [J]. *Applications Research of Computers*, 2018, 35(4): 1197–1200.)

This work is partially supported by the National Natural Science Foundation of China (61602227).

**SHEN Xueli**, born in 1969, M. S., professor. His research interests include information security, cyber security.

**QIN Shujuan**, born in 1993, M. S. candidate. Her research interests include cyber security.