

文章编号: 2095-2163(2024)03-0174-07

中图分类号: TP393.08

文献标志码: A

基于 GAN 和特征选择技术的入侵检测数据增强

崔子才, 钟伯成, 赵欣阳

(上海工程技术大学 电子电气工程学院, 上海 201620)

摘要: 为了解决传统 GAN 模型的缺陷,更好地扩展网络入侵数据和缓解数据高维性问题,本文提出了 GAN-CS 数据增强模型。对数据进行预处理后,使用改进后的 WGAN-GP 对攻击数据进行增强,生成额外的攻击样本后,使用卡方检验方法选择最能够代表数据集的特征,生成用于分类器训练平衡后的数据集,最后使用多种不同的分类器对数据集进行分类,评估模型效果。本文基于 UNSW-NB15 分别进行了数据增强数据量选择实验、模型可行性实验、模型优越性比较等 3 个维度的实验。结果表明,在多个分类器下,本文提出的模型均表现出比同类模型更好的效果,可以有效提高入侵检测模型的检测性能。**关键词:** 入侵检测; 数据增强; WGAN-GP 算法; UNSW-NB15 数据集

Intrusion detection data augmentation based on GAN and feature selection technique

CUI Zicai, ZHONG Bocheng, ZHAO Xinyang

(School of Electronic and Electrical Engineering, Shanghai University of Engineering Science, Shanghai 201620, China)

Abstract: In order to address the shortcomings of the traditional GAN model, better extend the network intrusion data and alleviate the problem of high dimensionality of data, this paper proposes the GAN-CS data enhancement model. The data is preprocessed and then augmented with attack data using the improved WGAN-GP to generate additional attack samples. Then, the features that best represent the dataset are selected using a Chi-Square test method, and a balanced dataset is generated for classifier training, finally the dataset is classified using a variety of different classifiers to evaluate the model effect. In this paper, experiments based on UNSW-NB15 are conducted to perform three dimensions experiments such as data enhancement data volume selection experiments, model feasibility experiments, and model superiority comparison, respectively, and the results show that the models proposed in this paper all show better results than similar models under multiple classifiers, which can effectively improve the detection performance of intrusion detection models.

Key words: intrusion detection; data enhancement; WGAN-GP algorithm; UNSW-NB15 dataset

0 引言

随着信息技术的不断发展,互联网已在政治、军事、经济、交通等领域发挥着重大作用。与此同时,各种网络攻击行为始终存在于互联网中,不仅可能造成巨大的经济损失,严重时甚至会威胁到国家安全和社会的稳定发展。入侵检测系统(Intrusion Detection System, IDS)是一种积极主动的安全防护技术,通过对网络进行实时监控,能够有效感知网络攻击行为,为安全管理人员提供相应决策。近年来,入侵检测系统已广泛地利用机器学习算法监控恶意活动,如贝叶斯网络、支持向量机、决策树等。此外,

随着深度学习的快速发展,卷积神经网络、循环神经网络等陆续在入侵检测系统中得到广泛应用。

然而,基于深度学习的检测系统高度依赖数据集,数据类间不平衡将严重影响检测的准确率。研究中发现,异常流量数据远小于征程流量数据,并且获取所需供给流量的途径有限。主要包括 3 方面:

(1) 采用真实攻击手段对网络用户进行入侵;

(2) 根据各种攻击代码的行为对其进行建模后,根据模型生成所需的攻击流量;

(3) 从真实存在的攻击事件中获取攻击流量^[1]。

第 1 种方法的网络攻击代码难以获取,第 2 种

基金项目: 国家自然科学基金青年科学基金项目(62102241)。

作者简介: 崔子才(1998-),女,硕士研究生,主要研究方向:网络安全;赵欣阳(1996-),男,硕士研究生,主要研究方向:网络安全。

通讯作者: 钟伯成(1964-),男,博士,教授,主要研究方向:计算机网络安全。Email: bczhong@sues.edu.cn

收稿日期: 2023-03-08

方法模型的准确率会影响攻击流量的可靠性,第3种方法不适合进行大规模的采集工作。数据增强(Data Augmentation)技术通过某些技术手段让有限的的数据产生更多的等价数据,实现数据更加复杂的表征,能够一定程度缓解数据缺乏和类间不平衡等问题。

在已有的研究中,成本函数等算法级别上的解决方案,在IDS数据不平衡问题上的研究较少,部分研究是通过欠采样、过采样等方法解决该问题,但欠采样的方式缩小了整体的样本数量,而过采样的方式又容易引发过拟合问题,并不能较好地处理数据不平衡问题。2014年,Goodfellow等学者^[2]提出了一种新的生成模型生成对抗网络(Generative Adversarial Networks, GAN),通过生成模型和判别模型的相互博弈学习生成高质量样本,能够很好地处理数据类别不平衡问题。

然而,传统的GAN模型对于网络入侵中的离散数据生成效果较差,生成的离散数据不能够以均匀的概率分布。为了解决传统GAN模型的缺陷、更好地扩展网络入侵数据,本文对WGAN-GP模型的网络结构进行改进,结合特征选择算法,提出了GAN-CS模型。将预处理后的数据输入GAN-CS模型,生成的样本数据更加逼真和详细,整体效果更稳定,增强后数据多分类的准确性得到了提高。

1 相关研究

传统的数据增强方法(如:过采样),通过随机过采样和合成过采样来生成新的少数类样本,以均衡数据集中各类别的数量^[3]。为了缓解随机过采样生成的新样本与原样本相似度高问题,Chawla等学者^[4]提出SMOTE算法,通过随机选择少数类样本点作为采样种子点,并使用线性插值的方法生成新的少数类样本。但是,当少数类样本由多个子群组成,并且多数类样本分布在子群中间时,线性插值法会生成与多数类重叠的样本,导致分类性能下降。Barua等学者^[5]提出基于多数类加权的少数类样本过采样技术,用来确定边界的少数类样本。但该方法的性能在很大程度上取决于如何对少数类样本进行分区以及加权,并且所选样本可能存在冗余信息。Bej等学者^[6]提出LoRAS方法,通过基于样本点的凸集生成新的少数类样本。

传统的过采样方法更适于处理低维数据,难以处理高维数据。随着深度学习的不断发展,在不均衡图像数据分类中会经常用到深度生成模型。其中

变分自动编码器(Variational Autoencoder, VAE)^[7],以及生成对抗网络(GAN)被广泛应用。VAE应用于给定的不平衡数据来捕捉特征维度间的关联性,进而获得样本在隐空间上的分布,最后通过解码器获得原始空间上的扩充数据集^[8]。研究是利用最小平方误差,衡量生成样本分布与原始样本分布之间的距离。然而,基于元素点之间的误差无法很好地捕捉数据的真实分布。为了提升VAE的性能,GUO等学者^[9]通过2个高斯分布,分别对多数类和少数类的隐空间变量进行建模,该模型适用于多分类数据。

作为一种数据生成策略,GAN能有效地学习隐空间到原始空间的映射函数。研究者提出了条件GAN(Conditional GAN, cGAN)^[10],用来生成特定类别的样本。基于GAN的数据生成方法中,生成器的输入通常是随机噪声,可能会导致特征高度纠缠并破坏方向相关的特征^[11]。为了缓解该问题,研究人员提出BAGAN^[12],是将AE(Autoencoder)和cGAN集成在一起,将新的隐空间编码作为cGAN的输入。但该方法中,GAN模式崩溃以及梯度消失和爆炸的问题仍无法避免,生成的数据甚至可能导致类别边界变形^[13]。Salem等学者^[14]使用Cycle-GAN将ADFA-LD数据集转换为图像,再使用Cycle-GAN学习正常数据的图像来创建异常数据的图像,将生成的综合异常数据与原始数据一起用于模型的训练。实验结果表明,该方法优于综合采样技术,显示了生成对抗网络在异常生成中的潜力。Yin等学者^[15]提出Bot-GAN模型。该模型根据各种流量的异常行为提取出相关特征,并结合GAN生成假样本,经判别器判别为真样本后,可继续细分为正常或异常样本。实验结果表明,相较于原始检测模型,Bot-GAN模型在测试集上有较高的检测准确率。

然而,入侵检测数据属于非图像类数据,在将其转换成图像处理的过程中可能会带来精度损失,因此本文提出GAN-CS模型解决IDS数据的不平衡问题。

2 关键技术

2.1 生成对抗网络

生成对抗网络(GAN)中包含2个网络。一是生成网络 G ,用于生成假样本;另一个是判别网络 D ,用于判别样本的真假。这2个目标相反的网络不断地进行交替训练,当最后收敛时,如果判别网络再也无法判断出一个样本的来源,也就等价于生成网

络可以生成符合真实数据分布的样本。

GAN的对抗博弈可以通过判别函数 $D(X): R^n \rightarrow [0, 1]$ 和生成函数 $G: R^d \rightarrow R^n$ 之间目标函数的极大极小值来进行数学化表示。生成器 G 将随机样本 $z \in R^d$ 分布 γ , 转化为生成样本 $G(z)$ 。判别器 D 试图将其与来自分布 μ 的训练样本区分开来, 而 G 试图使生成的样本在分布上与训练样本相似。GAN 解决的极小极大值的描述如下所示:

$$\min_G \max_D V(D, G) = \min_G \max_D (E_{x \sim \mu} [\log D(x)] + E_{z \sim \gamma} [\log(1 - D(G(z)))]) \quad (1)$$

其中, E 表示关于下标中指定分布的期望值。

对于给定的生成器 G , $\max_D V(D, G)$ 优化判别器 D , 以区分生成的样本 $G(z)$ 。其原理是尝试将高值分配给来自分布 μ 的真实样本, 并将低值分配给生成的样本 $G(z)$ 。相反, 对于给定的判别器 D , $\min_G V(D, G)$ 优化 G , 使得生成的样本 $G(z)$ 将试图误导判别器 D 以分配高值。GAN 算法描述如下。

算法1 生成对抗网络的训练过程

输入 训练集 D , 对抗训练迭代次数 T , 每次判别网络的训练迭代次数 K , 小批量样本数量 M
输出 生成对抗网络 $G(z, \theta)$

1. 随机初始化 θ 和 φ
2. For $t \leftarrow 1$ to T do
3. // 训练判别器网络 $D(x, \varphi)$
4. For $k \leftarrow 1$ to K do
5. // 采集小批量训练成本
6. 从训练集合 D 中采集 M 个样本 $\{x^{(m)}\}, 1 \leq m \leq M$
7. 从分布 $N(0, 1)$ 中采集 M 个样本 $\{z^{(m)}\}$, 这里 $1 \leq m \leq M$
8. 使用随机梯度上升更新 φ , 梯度为:
9. $\frac{\partial}{\partial \varphi} [\frac{1}{M} \sum_{m=1}^M (\log D(x^{(m)}, \varphi) + \log(1 - D(G(z^{(m)}), \theta), \varphi))]$
10. End for
11. // 训练判别器网络 $D(z, \theta)$
12. 从分布 $N(0, 1)$ 中采集 M 个样本 $\{z^{(m)}\}, 1 \leq m \leq M$
13. 使用随机梯度上升更新 θ , 梯度为:
14. $\frac{\partial}{\partial \theta} [\frac{1}{M} \sum_{m=1}^M (D(G(z^{(m)}), \theta), \varphi)]$
15. End for

2.2 WGAN-GP

在实际训练中, GAN 经常遇到的问题: 一是模

式崩溃, 生成器生成非常窄的分布, 仅覆盖数据分布中的单一模式, 即生成器只能生成非常相似的样本; 二是没有指标可以表征收敛情况。总之, 判别器越好, 生成器梯度消失越严重。

在判别器最优的前提下, 把原始 GAN 定义的生成器 $loss$ 等价变换为最小化真实分布与生成分布之间的 JS 散度, 最小化生成器的 $loss$ 即近似于最小化真实分布与生成分布之间的 JS 散度。若希望 2 个分布之间的 JS 散度越小, 通过优化 JS 散度就能将生成分布转化为真实分布, 最终实现以假乱真; 若 2 个分布完全没有重叠的部分, 或者相互重叠的部分可忽略, 则两者之间的 JS 散度就一直是 $\log 2$ 。

因此, 原始 GAN 问题的根源可以归结为 2 点: 等价优化的距离衡量 (JS 散度) 不合理, 以及生成器随机初始化后的生成分布很难与真实分布有不可忽略的重叠。

基于以上问题, 研究者提出了 Wasserstein 距离, 即 Earth - Move 距离衡量 2 个分布之间的距离^[16]。其优越性在于, 即使 2 个分布没有任何重叠, 也可以反映两者之间的距离, 公式如下。

$$W(P_r, P_g) = \frac{1}{K} \sup_{\|f\|_{L \leq K}} E_{x \sim P_r} [f(x)] - E_{x \sim P_g} [f(x)] \quad (2)$$

其中, P_r 为最小化真实分布, P_g 为生成分布。

WGAN 算法描述如下。

算法2 WGAN 的训练过程

输入 学习率 α , 裁剪参数 c , 对抗训练迭代次数 T , 每次判别网络的训练迭代次数 K 。初始临界参数 w_0 , 初始生成器参数 θ_0

1. while θ 未收敛 do
2. for $t = 0, 1, \dots, K$ do
3. $\{x^{(i)}\}_{i=1}^T \sim P_r$ a batch from real data
4. $\{z^{(i)}\}_{i=1}^T \sim p(z)$ a batch of priors
5. $g_w \leftarrow \nabla_w [\frac{1}{T} \sum_{i=1}^T f_w(x^{(i)}) - \frac{1}{T} \sum_{i=1}^T f_w(g_\theta(z^{(i)}))]$
6. $w \leftarrow w + \alpha \cdot RMSProp(w, g_w)$
7. $w \leftarrow clip(w, -c, c)$
8. End for
9. $\{z^{(i)}\}_{i=1}^T \sim p(z)$ a batch of prior samples
10. $g_\theta \leftarrow -\nabla_\theta [\frac{1}{T} \sum_{i=1}^T f_w(g_\theta(z^{(i)}))]$
11. $\theta \leftarrow \theta - \alpha \cdot RMSProp(\theta, g_\theta)$
12. $\frac{\partial}{\partial \theta} [\frac{1}{M} \sum_{m=1}^M (D(G(z^{(m)}), \theta), \varphi)]$

13. End while

然而, WGAN 中权重裁剪的实现方式存在 2 个重要问题:

(1) 判别器的 loss 希望尽可能拉大真假样本的分数差, 实验发现基本上最终权重集中在两端, 这样参数的多样性会减少, 使判别器得到的神经网络学习一个简单的映射函数, 造成巨大的浪费。

(2) 容易导致梯度消失或者梯度爆炸, 若把裁剪阈值设得较小, 每经过一个网络, 梯度就会变小, 多级之后会成为指数衰减; 反之则会导致指数爆炸。

为此, WGAN-GP 引入了梯度惩罚项 (Gradient Penalty)^[17]。当且仅当一个可微函数的梯度范数 (Gradient Norm) 在任意处都不超过 1 时, 该函数满

足 1-Lipschitz 条件。损失函数公式如下:

$$L = \frac{E}{x \sim P_g} [D(\tilde{x})] - \frac{E}{x \sim P_r} [D(x)] + \lambda \frac{E}{x \sim P_g} [(\| \nabla_x \|_2 - 1)^2] \quad (3)$$

3 GAN-CS 模型

本文提出 GAN-CS 模型对不平衡数据进行增强, 并对模型的性能进行评估。模型框架如图 1 所示。由图 1 可知, 模型包括数据预处理、数据增强、特征选择及性能评估四部分。

为验证模型生成的数据质量, 分别使用处理后的数据集和混合数据集训练相同的分类器, 通过比较 2 种数据集下的多分类结果, 评估本文所提模型的性能。

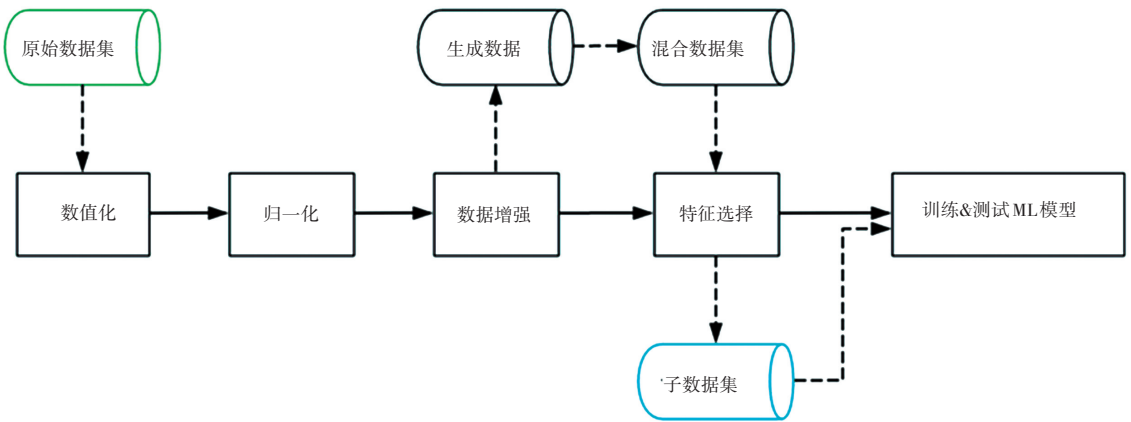


图 1 模型框架

Fig. 1 Model framework diagram

3.1 数据预处理

入侵检测数据通常包含非数字特征, 例如协议和状态等。为了能更好地被计算机识别和处理, 需先通过 one-hot 编码方法将入侵检测数据集中存在的离散型数据进行数值化, 然后将所有数字特征进行归一化, 以保证消除数据的可读性和消除异常值。这些非数字特征需要转换为数字特征以适合本文模型。非数字特征映射到 0 到 S - 1 之间的整数值, 其中 S 是符号数。

不同维度的数据特征尺度不一致会影响入侵检测的结果。需要对数据进行归一化处理, 以消除指标之间的维度影响。除了攻击类型标签, 将所有特征缩放到 [0, 1]。min-max 归一化用于线性缩放数据值, 如下所示:

$$x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}} \quad (4)$$

其中, x 为归一化前的值; x' 为归一化后的值; x_{max} 为样本数据的最大值; x_{min} 为样本数据的最小值。

将预处理后的数据集通过改进后的 WGAN-GP 模型进行扩充, 与预处理数据集混合, 形成混合数据集。

3.2 特征选择

特征选择是一种数据降维方法, 常用于处理高维、复杂的数据。从所有的特征中, 选择出有意义、对模型有帮助的特征, 以避免必须将所有特征都导入模型去训练的情况, 从而提升模型的训练速度和效率, 并提升准确率。

卡方检验是以 χ^2 分布为基础的一种常用假设检验方法。方法的无效假设 H_0 是: 观察频数与期望频数没有差别^[18]。该检验的基本思想是: 首先假设 H_0 成立, 基于此前提计算出 χ^2 值, 则表示观察值与

理论值之间的偏离程度。根据 χ^2 分布及自由度可以确定在 H_0 假设成立的情况下获得当前统计量及更极端情况的概率 P 。如果当前统计量大于 P 值,说明观察值与理论值偏离程度太大,应当拒绝无效假设,表示比较资料之间有显著差异;否则就不能拒绝无效假设,尚不能认为样本所代表的实际情况和理论假设有差别。其公式如下:

$$\chi^2 = \sum \frac{(A - E)^2}{E} \quad (5)$$

其中, A 为实际值, E 为理论值。

4 实验结果与分析

4.1 UNSW-NB15 数据集

UNSW-NB15 数据集^[19]由澳大利亚网络安全中心的网络靶场实验室创建。该数据集包含各种新颖的攻击,因此已广泛用于入侵检测。其中包含9种类型的攻击来模拟真实网络环境,即 Fuzzers、Analysis、Backdoor、DoS、Exploits、Generic、Reconnaissance、Shellcode 和 Worms。UNSW-NB15 数据集包含一个训练集和一个测试集。训练集有82 332条记录,测试集有175 341条记录。UNSW-NB15 数据集构成如图2所示。由图2可知,不同种类流量数据分布严重不均,且数据间存在数量的差异。如 Worms 在数据集中的分布仅占不足0.1%,而 Normal 的数量超过整体数据集中的1/3。

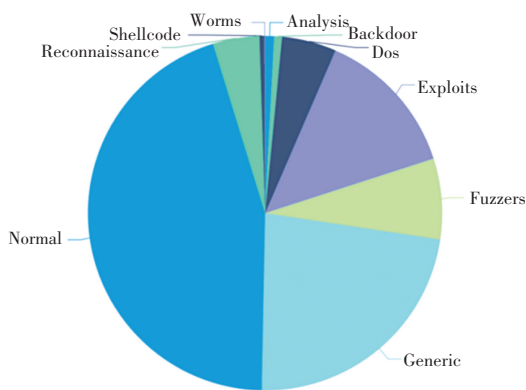


图2 UNSW-NB15 数据集构成

Fig. 2 UNSW-NB15 dataset

4.2 评价指标

数据在二分类问题中可被分为正样本和负样本,并将数据按照真实类别和预测类别划分为4种类型:当样本预测为正且实际为正时的真阳性

(TP),当样本预测为负且实际为负时的真阴性(TN),当样本预测为正但实际为负时的假阳性(FP),当样本预测为负但实际上为正时的假阴性(FN)。

为了评价本文所提模型的性能,采用识别准确率($Accuracy$)、精确率($Precision$)、召回率($Recall$)、 F 值($F - Measure$)作为评价指标。

(1)准确率($Accuracy$)。研究推得的定义公式为:

$$A = \frac{TP + TN}{TP + FP + TN + FN} \quad (6)$$

(2)精确率($Precision$)。精确率又称查准率,是常用的评价指标,就是计算所有被预测为正的样本中实际为正样本的概率,公式如下:

$$P = \frac{TP}{TP + FP} \quad (7)$$

(3)召回率($Recall$)。召回率又称查全率,就是计算实际为正的样本中被预测为正样本的概率,公式如下:

$$R = \frac{TP}{TP + FN} \quad (8)$$

其中, FN (False Negative)表示将正样本预测为负样本的数量。

(4) $F1$ 值($F1 - value$)。为避免精确率和召回率相矛盾的情况,需要将两者进行综合考虑,最常见的方法为 F 值,就是计算精确率和召回率的加权调和平均,公式如下:

$$F = \frac{(\alpha^2 + 1) \times P \times R}{\alpha^2 \times (P + R)} \quad (9)$$

其中,当参数 $\alpha = 1$ 时,即为 $F1$ 值。

4.3 实验结果与分析

本文所有实验结果均由10次实验后取平均值得到。研究内容分述如下。

(1)实验一:数据增强数据量选择实验。为更好地体现本文提出模型的数据增强的效果,在使用同一分类器(实验选择决策树)的前提下,分别将不同攻击类别的数据由原来的数据量增加10 000、增加至10 000、20 000和37 000(Normal的数据量为37 000)(见表1),并对混合后数据集中的准确率和原数据集的准确率做了比较(见表2)。由此可见,当攻击样本数据增加至20 000时,分类器的准确率最高。

表 1 数据增强前后数据量对比

Table 1 Comparison of data volume before and after data augmentation

攻击类别	增强前数量	增强10 000	增强至10 000	增强至20 000	增强至37 000
Analysis	677	10 677	10 000	20 000	37 000
Backdoor	583	10 583	10 000	20 000	37 000
Dos	4 089	14 089	10 000	20 000	37 000
Exploits	11 132	11 132	11 132	20 000	37 000
Fuzzers	6 062	16 062	10 000	20 000	37 000
Generic	18 871	18 871	18 871	20 000	37 000
Normal	37 000	37 000	37 000	37 000	37 000
Reconnaissance	3 496	13 496	10 000	20 000	37 000
Shellcode	378	10 378	10 000	20 000	37 000
Worms	44	10 044	10 000	20 000	37 000

表 2 数据增强前后准确率对比 (使用决策树分类器)

Table 2 Comparison of accuracy before and after data augmentation (using decision tree classifier)

数据量	准确率
增强前	0.498 3
增强 10 000	0.696 5
增强至 10 000	0.465 7
增强至 20 000	0.702 5
增强至 37 000	0.699 8

数据集和混合数据集进行多分类,比较结果如图 3 所示。由图 3 可知,每个分类器的准确率都有不同程度的提高,表明本文所提模型可以提高分类器的整体性能,数据增强模型是有效的,且效果明显。所有分类器的召回率也都有不同程度的提高,表明本文生成的样本提高了攻击样本的多样性,从而增强了分类器知识学习的泛化性。由于不同的检测器具有不同的学习能力,因此影响程度存在一定差异。最后,分类器的 $F1$ 值表明分类器的整体性能得到了有效的提升。

(2) 实验二:模型有效性实验。在实验一的基础上,使用决策树、随机森林、KNN 和神经网络对原

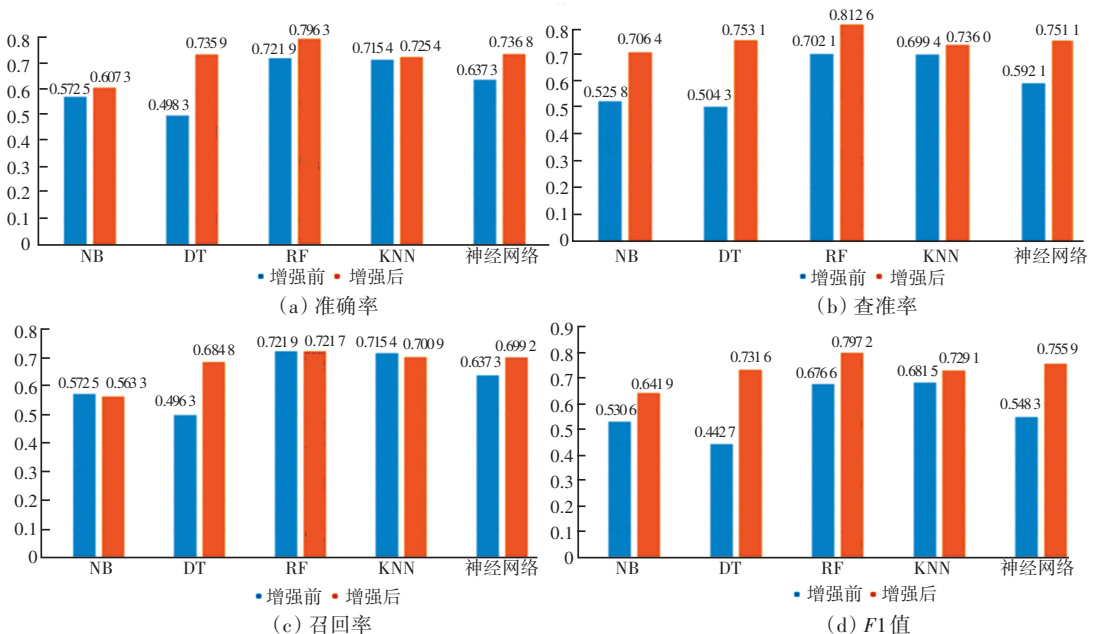


图 3 增强前后对比

Fig. 3 Comparison before and after enhancement

(3) 实验三:模型优越性比较实验。为了证明本文提出的数据增强模型 GAN-CS 优于同类模型,

本实验将 WGAN-GP 和 GAN-FS 作为对比对象。实验结果见表 3。

表3 本文模型与其他模型在 UNSW-NB15 数据增强效果的比较

Table 3 Comparison of the enhancement effect of this model with other models in UNSW-NB15 data

分类器	NB		DT		RF		KNN		神经网络	
	Accuracy	F1	Accuracy	F1	Accuracy	F1	Accuracy	F1	Accuracy	F1
WGAN-GP ^[17]	0.563 3	0.523 6	0.701 5	0.684 8	0.721 7	0.683 5	0.700 9	0.674 8	0.679 2	0.586 6
GAN-FS ^[20]	0.597 0	0.538 1	0.686 0	0.659 7	0.720 6	0.682 8	0.700 9	0.674 8	0.683 2	0.589 7
GAN-CS	0.607 3	0.541 9	0.735 9	0.701 6	0.796 3	0.737 2	0.725 4	0.689 1	0.736 8	0.615 9

从表3中数据可以看出,基于本文提出的模型的分器性能高于其他方案,且GAN-CS在DT和RF上表现更好。DT模型的思想是使用信息熵作为度量来构建熵下降最快的树,WGAN-GP算法基于原始分布生成样本,增加了样本的多样性,并通过特征选择去除了不必要的特征,因此,经过过采样后,决策树可以更好地对样本进行分类。与DT相比,RF是综合学习算法,学习能力更优异,可以提高分类性能。试验结果还表明RF的性能普遍高于DT。

5 结束语

为缓解数据缺乏和类间不平衡等问题,本文提出了GAN-CS数据增强模型。将预处理过的数据集通过改进后的WGAN-GP模型中进行增强,生成后的数据和原数据集混合后得到的混合数据集经特征选择后用于训练入侵检测多分类器。经过3种不同维度的实验可以得出,本文方法提高了入侵检测模型的性能,并且优于其他同类方案。虽然好的数据增强模型能够提升入侵检测的性能,但分类器又影响了入侵检测模型性能的提升,因此设计一个好的分类器将成为下一步的研究方向。

参考文献

[1] 陈家浩,王轶骏,吕诚.一种基于Python符号执行的自动化网络攻击流量获取方法[J].计算机应用与软件,2019,36(2):294-307.

[2] GOODFELLOW I J, POUGET-ABADIE J, MIRZA M, et al. Generative adversarial nets[C]//Proceedings of the 27th International Conference on Neural Information Processing Systems. Cambridge: MIT Press, 2014: 2672-2680.

[3] 王馨月.生成式数据增强的不平衡数据分类方法研究[D].北京:北京交通大学,2021.

[4] CHAWLA N V, BOWYER K W, HALL L O, et al. SMOTE: synthetic minority over-sampling technique [J]. Journal of Artificial Intelligence Research, 2002, 16: 321-357.

[5] BARUA S, ISIAM M M, YAO Xin, et al. MWMOTE—majority weighted minority oversampling technique for imbalanced data set learning[J]. IEEE Transactions on Knowledge and Data Engineering, 2012, 26: 405-425.

[6] BEJ S, NAREK D, MARKUS W, et al. LoRAS: An oversampling approach for imbalanced datasets[J]. Machine Learning, 2021, 110: 279-301.

[7] KINGMA D P, WELING M. Auto-encoding variational bayes[J]. arXiv preprint arXiv:1312.6114, 2013.

[8] WAN Zhiqiang, ZHANG Yazhou, HE Haibo. Variational autoencoder based synthetic data generation for imbalanced learning[C]//2017 IEEE Symposium Series on Computational Intelligence (SSCI). Honolulu, USA:IEEE, 2017: 1-7.

[9] GUO Ting, ZHU Xingquan, WANG Yang, et al. Discriminative sample generation for deep imbalanced learning [C]// Twenty-Eighth International Joint Conference on Artificial Intelligence. Macao, China: International Joint Conferences on Artificial Intelligence Organization, 2019: 2406-2412.

[10] GAUTHIER J. Conditional generative adversarial nets for convolutional facegeneration[J]. Class Project for Stanford CS231N: Convolutional Neural Networks for Visual Recognition, 2014(5): 2.

[11] NAZRUL H, BHATTACHARYYA D K, KALITA J K. Botnet in DDoS attacks: Trends and challenges[J]. IEEE Communications Surveys & Tutorials, 2015, 17(4): 2242-2270.

[12] MARIANI G, SCHEIDEGGER F, ISTRATE R, et al. BAGAN: Data augmentation with balancing GAN[J]. arXiv preprint arXiv: 1803.09655, 2018.

[13] SANTURKAR S, SCHMIDT L, MADRY A. A classification-based study of covariate shift in GAN distributions [C]// International Conference on Machine Learning. Stockholm, Sweden :PMLR, 2018:4480-4489.

[14] SALEM M, TAHERI S, YUAN J S. Anomaly generation using generative adversarial networks in host-based intrusion detection [C]// 2018 9th IEEE Annual Ubiquitous Computing, Electronics & Mobile Communication Conference (UEMCON). New York, USA:IEEE, 2018: 683-687.

[15] YIN Chuanlong, ZHU Yuelei, LIU Shengli. An enhancing framework for botnet detection using generative adversarial networks[C]// 2018 International Conference on Artificial Intelligence and Big Data (ICAIBD). Chengdu, China:IEEE, 2018:228-234.

[16] ARJOVSKY M, CHINTALA S, BOTTOU L. Wasserstein generative adversarial networks [C]//International Conference on Machine Learning. Sydney :PMLR, 2017: 214-223.

[17] GULRAJANI I, AHMED F, ARJOVSKY M, et al. Improved training of Wasserstein GANs [C]// Proceedings of the 31st International Conference on Neural Information Processing Systems. Long Beach: Curran Associates Inc., 2017: 5769-5779.

[18] 陈湛,梁雪春.基于基尼指标和卡方检验的特征选择方法[J].计算机工程与设计,2019,40(8):2342-2345,2360.

[19] ZOGHI Z, SERPEN G. UNSW-NB15 computer security dataset: Analysis through visualization [J]. arXiv preprint arXiv: 2101.05067, 2021.

[20] LIU Xiaodong, LI T, Zhang Runzi, et al. A GAN and feature selection-based oversampling technique for intrusion detection[J]. Security and Communication Networks, 2021(1):1-15.