

doi: 10.7690/bgzdh.2024.08.006

基于集成学习的物联网攻击检测方法

窦佳恩¹, 张瑛瑛², 陈 玮²

(1.重庆邮电大学通信与信息工程学院, 重庆 400065; 2.重庆邮电大学理学院, 重庆 400065)

摘要: 针对物联网入侵检测方法在检测精度和资源消耗方面存在的问题, 提出基于集成学习轻量级梯度提升机(light gradient boosting machine, LightGBM)和随机森林(random forest, RF)模型的物联网攻击检测方法。采用网格搜索(grid search, GS)和启发式算法对模型的超参数进行调优超参数优化, 使用合成少数过采样技术(synthetic minority over-sampling technique, SMOTE)进行数据增强, 解决物联网数据集标签不平衡问题。实验结果表明: 通过 SMOTE 技术后, 基于网格搜索下的轻量级梯度提升机(grid search-light gradient boosting machine, GS-LightGBM)模型准确率达 99.91%, 且在不平衡数据集上表现优异; 在资源消耗方面, 基于遗传算法下的随机森林(genetic algorithm-random forest, GA-RF)模型在准确率 99.88%的情况下, 平均推理时间达到微秒级别, 推理时占用内存不到 1kB, 在大部分资源受限的物联网设备上能实现高效运行。

关键词: 物联网; 网络安全; 集成学习; 超参数优化; SMOTE

中图分类号: 文献标志码: A

IoT Attack Detection Method Based on Ensemble Learning

Dou Jiaen¹, Zhang Yingying², Chen Wei²

(1. School of Communication and Information Engineering, Chongqing University of Posts and Telecommunications, Chongqing 400065, China;

2. School of Science, Chongqing University of Posts and Telecommunications, Chongqing 400065, China)

Abstract: In addressing the issues of detection accuracy and resource consumption in IoT intrusion detection methods, an IoT attack detection method based on ensemble learning lightweight gradient boosting machine (LightGBM) and random forest (RF) models is proposed. Grid search (GS) and heuristic algorithms are used to tune the hyperparameters of the models, and the synthetic minority over-sampling technique (SMOTE) is employed for data augmentation to address the imbalance in IoT dataset labels. Experimental results show that after applying the SMOTE technique, the grid search-light gradient boosting machine (GS-LightGBM) model achieves an accuracy of 99.91% and performs excellently on imbalanced datasets. In terms of resource consumption, the genetic algorithm-random forest (GA-RF) model achieves an accuracy of 99.88%, with an average inference time at the microsecond level and memory usage of less than 1kB during inference, enabling efficient operation on most resource-constrained IoT devices.

Keywords: IoT; network security; ensemble learning; hyperparameter optimization; SMOTE

0 引言

随着物联网(IoT)技术的快速发展,物联网设备在多个领域被广泛应用。然而,由于资源受限和部署环境复杂,物联网设备常成为网络攻击的目标^[1],因此,物联网入侵检测成为关键研究方向。

近年来,集成学习和深度学习方法在网络入侵检测系统(intrusion detection system, IDS)中取得了显著成果。Kumar 等^[2]提出了一种分布式入侵检测系统,通过在分布式雾节点上训练随机森林(RF)和极限梯度提升(extreme gradient boosting, XGBoost)模型进行检测, XGBoost 在二分类攻击检测中表现优异,而 RF 在多攻击检测中表现更好。文献^[3]提出了基于轻量级梯度提升机(LightGBM)的物联网

入侵检测模型。结果表明: LightGBM 在继承梯度提升树优势的基础上,更加轻量级,对拒绝服务(DoS)和分布式拒绝服务(DDoS)攻击有较强检测能力。资源消耗方面,文献^[4]中的 FlowGuard 方案虽然能准确检测 DDoS 攻击流,但由于复杂的分类流程,推理时间较长。除此之外,卷积神经网络(convolutional neural network, CNN)^[5]和多层感知器(multilayer perceptron, MLP)^[6]等深度学习算法也广泛应用于该领域。这些算法中,超参数往往控制着模型的训练过程和复杂度,影响最终性能。现有研究中,超参数设置通常通过实验调整,缺乏系统性且计算资源消耗大。此外,物联网数据集通常存在数据不平衡问题,导致模型在少量数据标签类

收稿日期: 2024-04-23; 修回日期: 2024-05-25

基金项目: 重庆市自然科学基金项目(CSTB2023NSCQMSX0435)

第一作者: 窦佳恩(2003—),男,陕西人。

别上的预测能力降低,影响整体性能。综上,提高超参数设置的有效性,解决标签不平衡问题,提高边缘设备部署的可行性,对于提升物联网攻击检测系统的可靠性具有重要意义。

为解决上述问题,笔者基于集成学习模型 LightGBM 和 RF 进行研究,主要贡献如下:

1) 为提高超参数的有效性,使用网格搜索(grid search, GS),粒子群优化(particle swarm optimization, PSO)、贝叶斯优化(Bayesian optimization, BO)和遗传算法(genetic algorithm, GA)对 LightGBM 和 RF 模型进行超参数调优,并与网格搜索下的 XGBoost, MLP 等模型进行对比,说明了集成学习模型在此类问题上的有效性以及超参数调优的重要性。

2) 针对标签类别不平衡问题,采用合成少数过采样技术(SMOTE)进行数据增强。经验证,数据增强后的 GS-LightGBM 模型在整体检测精度和少数类别标签上的准确率均表现优异。

3) 结合物联网技术特点,分析模型的实际部署。结果表明,GA-RF 模型在保持高准确率的同时,在模型大小、推理内存占用和时延方面表现良好,具有高实用价值。

1 基于集成学习的物联网攻击检测模型

1.1 数据介绍与预处理

本研究利用 RT-IoT2022 数据集进行实验,该数据集由多种实时物联网设备生成,包括 ThingSpeak-LED 和 MQTT-Temp,涵盖多种通信协议和网络服务,模拟真实物联网环境^[7]。数据集包含 83 个特征和 12 种标签,共记录了 123 117 条数据。在使用前,对数据进行预处理,包括编码非数值型列,将其转换为数值表示,并将所有特征标准化,使其均值为 0、方差为 1,确保特征在相同尺度上进行比较。

1.2 相关理论

1.2.1 超参数调优

超参数调优是在模型训练前设置超参数,以控制学习过程和模型复杂度。常用方法包括网格搜索(GS)和启发式算法。GS 是在定义的参数空间内进行穷举搜索,评估每个参数组合的性能并选择最优组合,虽简单易用但扩大搜索范围后计算成本高。启发式算法如贝叶斯优化(BO)、粒子群优化(PSO)、差分进化(DE)和遗传算法(GA)通过模拟生物进化或群体行为进行高效参数搜索。

1.2.2 SMOTE 技术

SMOTE 技术通过生成合成样本来增加少数类样本,从而平衡数据集。具体步骤如下:1) 使用最近邻算法找到每个少数类样本的 k 个近邻样本。2) 从 k 个近邻样本中随机选择一个,生成新的合成样本。3) 根据需要生成的合成样本数量,重复以上过程。

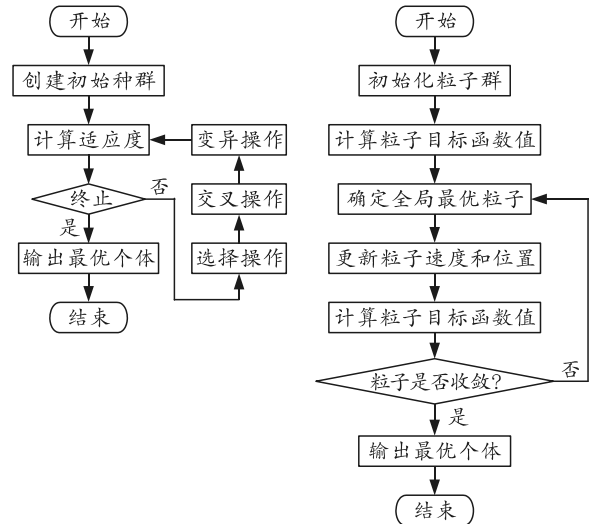


图 1 遗传算法和粒子群算法流程

1.3 模型构建

笔者基于集成学习 LightGBM 和 RF 建立最优模型,并验证部署可行性,具体的实施步骤如下:

1) 数据预处理:分离特征列与标签,对字符串类型特征进行编码,数据归一化。

2) 对 LightGBM 和 RF 模型进行超参数调优。首先,使用网格搜索进行参数调优,并与其他模型比较。然后,使用启发式算法对 RF 和 LightGBM 进一步调优。模型训练中,将数据集按 7:3 分为训练集和测试集,并采用交叉验证。以准确率为目标函数寻找最优超参数组合。

3) 使用 SMOTE 方法对训练集进行数据增强,分析其在不同算法上对不平衡数据集的优化效果。

4) 部署分析:结合物联网设备特征,分析模型的资源占用,验证部署可行性。

表 2 混淆矩阵

分类	预测正类	预测负类别
实际正类	TP(真正例)	TF(假反例)
实际负类	FP(假正例)	TN(真反例)

1.4 评估标准

该研究评估各模型分类能力的指标包括:准确率(Accuracy, A),每种标签的精确率(Precision,

P)、召回率 (Recall, R) 和 F1 分数 (F1-score, F1) 的宏平均 (macro-average) 及加权平均 (weighted average)。这些指标均通过混淆矩阵计算得出。

Accuracy 是正确分类样本数占总样本数的比例。Precision 是被预测为某类样本中实际属于该类的比例。Recall 是实际为某类样本中被正确预测的

比例。F1 是精确率与召回率的调和平均值。计算公式如下：

$$A=(TP+TN)/(TP+FP+TN+FN) \quad (1)$$

$$P=TP/(TP+FP) \quad (2)$$

$$R=TP/(TP+FN) \quad (3)$$

$$F1=(2 \times P \times R)/(P+R) \quad (4)$$

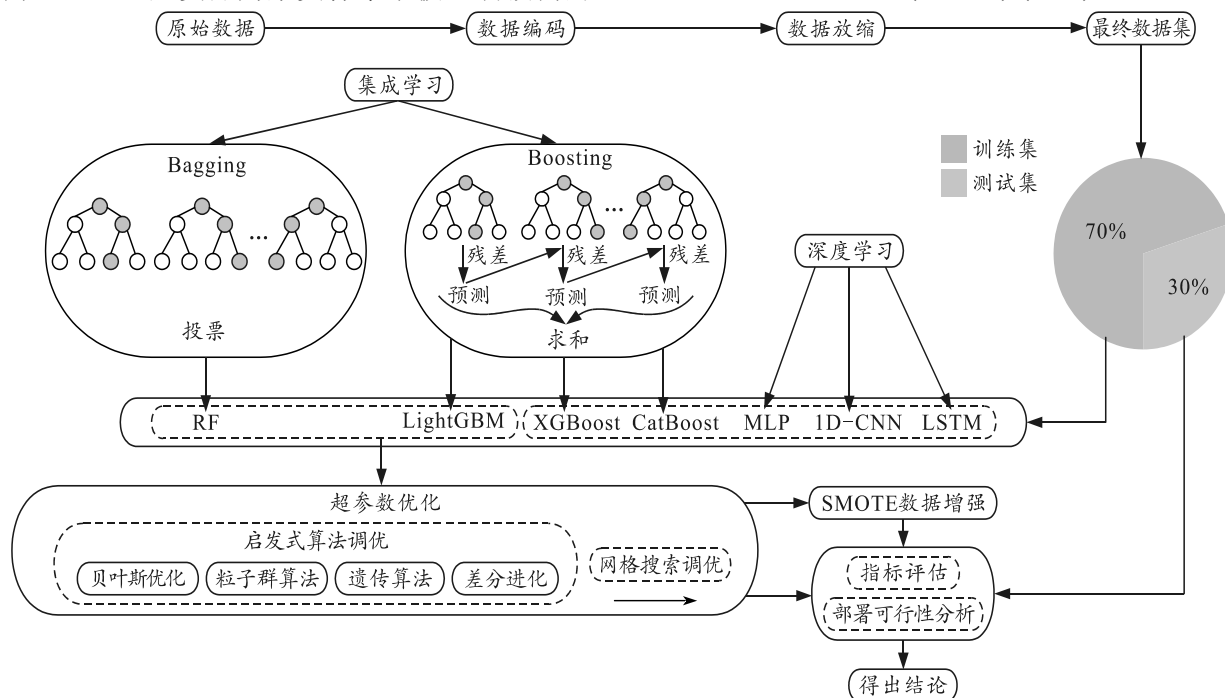


图 2 模型建立流程

2 实验结果及分析

2.1 网格搜索优化超参数

对 RF 和 LightGBM 模型进行网格搜索超参数优化，并与其他模型对比 (表 3)。RF、XGBoost、CatBoost 和 LightGBM 在准确率和加权平均指标上表现优异，说明其对整体数据的分类能力较强。而 MLP、1D-CNN 和 LSTM 在处理不平衡数据时表现稍逊，特别是在少数类数据的识别上，宏平均召回率较低。由于该数据集特征独立且无时间依赖性，LSTM 和 1D-CNN 的复杂结构未能提高性能。MLP 虽然能处理非时间序列数据，但表现仍不如集成学习模型。集成学习模型通过集成多个基学习器，减

少过拟合，提高性能，尤其是 XGBoost 和 LightGBM 在准确率上表现最佳。LightGBM 由于使用基于直方图的决策树算法，提高了训练效率，处理大规模数据时训练速度更快，内存消耗更低。

2.2 基于启发式算法的集成学习超参数优化

对 RF 和 LightGBM 模型使用启发式算法进行超参数调优 (表 5 和表 6)。结果表明，经过调优后，模型准确率相较于网格搜索均有所提升。其中，GA 调优的 RF 和 LightGBM 模型准确率分别上升了 0.03% 和 0.02%。模型的加权平均指标均有所提升，说明整体分类能力增强。然而，RF 的宏平均 F1 指标下降，表明其对少量数据类别的识别能力降低。

表 3 网格搜索下模型性能指标

模型	A	Macro-average			Weighted Average		
		P	R	F1	P	R	F1
RF	0.998 5	0.996 3	0.956 9	0.974 3	0.998 5	0.998 5	0.998 5
LightGBM	0.998 7	0.954 3	0.960 3	0.954 8	0.998 8	0.998 8	0.998 8
CatBoost	0.998 5	0.975 2	0.957 2	0.965 5	0.998 6	0.998 6	0.998 6
XGBoost	0.998 7	0.983 5	0.961 5	0.970 7	0.998 7	0.998 7	0.998 7
MLP	0.993 7	0.946 9	0.904 5	0.918 2	0.994 1	0.993 8	0.993 7
1D-CNN	0.993 3	0.959 3	0.893 4	0.907 2	0.993 8	0.993 4	0.993 3
LSTM	0.994 9	0.968 9	0.905 2	0.932 8	0.995 0	0.995 0	0.994 9

表 4 启发式算法的超参数寻优范围

模型	超参数	范围
RF	决策树的数量	[50, 200]
	树的最大深度	[10, 100]
	构建每个节点时,从所有特征中考虑的特征数量比例	[0.1, 1.0]
LightGBM	每个树的叶子节点数量	[30, 100]
	学习率	[0.01, 0.2]
	基学习器的数量	[50, 200]

表 5 RF 模型启发式算法调优性能指标

RF 优化 算法	A	Macro-average			Weighted Average		
		P	R	F1	P	R	F1
BO	0.998 7	0.976 9	0.958 7	0.965 7	0.998 8	0.998 8	0.998 8
DE	0.998 7	0.977 0	0.959 7	0.966 2	0.998 8	0.998 8	0.998 8
GA	0.998 8	0.967 8	0.961 0	0.963 4	0.998 9	0.998 9	0.998 9
PSO	0.998 7	0.967 7	0.959 7	0.962 7	0.998 7	0.998 7	0.998 7

表 6 LightGBM 模型启发式算法调优性能指标

LightGBM 优化算法	A	Macro-average			Weighted Average		
		P	R	F1	P	R	F1
BO	0.998 8	0.966 5	0.960 7	0.962 6	0.998 9	0.998 9	0.998 9
GA	0.998 9	0.967 6	0.960 7	0.963 2	0.998 9	0.998 9	0.998 9
PSO	0.998 8	0.958 5	0.960 4	0.954 8	0.998 9	0.998 8	0.998 8

2.3 SMOTE 数据增强

研究发现, GA-RF 和 GA-LightGBM 模型整体性能较好,但由于数据类别不平衡,其在某些标签类别上仍有提升空间。通过 SMOTE 增强训练集数据后,模型性能如表 7 所示。GS-LightGBM 模型表现最佳,各项指标均有所提高。具体来看,其中少数标签类别“NMAP_FIN_SCAN”和“Metasploit_Brute_Force_SSH”的 F1 指标分别上涨 2.9%和 19.56%。这表明 SMOTE 在处理不平衡数据时效果显著。然而,在 GA-RF 和 GA-LightGBM 模型中,训练集上的准确率、宏平均和加权平均指标均达到了 0.999 9,但测试集性能下降,表明训练集数据增强后,启发式算法调优容易导致过拟合。

表 7 SMOTE 数据增强后模型性能指标

模型	A	Macro-average			Weighted Average		
		P	R	F1	P	R	F1
GA-RF	0.998 8	0.960 8	0.960 8	0.959 1	0.998 8	0.998 8	0.998 8
GA-LightGBM	0.998 8	0.945 5	0.961 9	0.951 6	0.998 8	0.998 8	0.998 8
GS-RF	0.998 8	0.985 1	0.959 8	0.970 7	0.998 8	0.998 8	0.998 7
GS-LightGBM	0.999 1	0.973 3	0.987 9	0.980 2	0.999 2	0.999 2	0.999 2

2.4 部署可行性分析

在部署物联网攻击检测模型时,终端设备负责推理和数据收集,可降低响应时间,满足实时性要求。模型训练因需要大量计算资源和存储空间,通常在中央服务器或云端进行。随着物联网设备性能提升和 5G mMTC 的发展,常见处理器如 ARM Cortex-A53/A72 主频在 1-2.5 GHz 之间,用于高性

能设备如网关和边缘计算设备^[8]。

本实验使用 Intel(R) Core(TM) i5-10210U CPU @ 1.60GHz 2.11 GHz 处理器,单线程模拟部署模型推理。结果如表 8 所示,优化后的 GA-RF 模型大小仅 3.45 MB,满足大多数物联网设备的资源限制,在 99.88%准确率下,平均推理时间达到微秒级,内存占用不到 1kB,部署效率高。尽管 GA-LightGBM 和 GS-LightGBM(SMOTE)模型准确率略高于 RF,但其推理时占用资源远大于 RF 模型。

表 8 模型部署模拟指标

算法	准确率	模型大小 (MB)	推理过程中 内存使用(MB)	每个样本平均 推理时间(μs)
GA-LightGBM	0.998 9	38.80	0.76	171
GS-LightGBM (SMOTE)	0.999 1	6.26	2.71	4 213
GA-RF	0.998 8	3.45	0.00	4
GS-RF(SMOTE)	0.998 8	10.24	0.00	2

3 结束语

笔者研究了基于 LightGBM 和 RF 模型的物联网攻击检测方法,通过网络搜索和启发式算法进行超参数优化,并使用 SMOTE 技术解决标签不平衡问题。实验结果显示,经过优化后的模型在准确率、精确率和 F1 分数等指标上表现出色,尤其是 GA-RF 和 GA-LightGBM 模型整体性能最优,证明了启发式算法在超参数调优中的有效性。使用 SMOTE 增强后,GS-LightGBM 模型在不平衡数据上的性能显著提升,准确率达到 99.91%。在终端部署分析中,GA-RF 模型在推理时间和内存占用方面表现优异,适合实际应用。在相同的数据集下,相较于量化自动编码器(QAE-f16)^[7]和粒子群优化-深度学习(PSO-DL)^[9]模型,SMOTE 增强后的 GS-LightGBM 模型准确率分别提升了 2.66%和 5.91%。

现在仍有很多物联网设备(如低端传感器节点)内存只有几 kB,未来的研究应注重模型轻量化和资源优化,以提升实用性和部署效率。结合先进的超参数优化算法和数据增强方法,机器学习模型将在物联网安全领域发挥重要作用。

参考文献:

[1] 张玉清,周威,彭安妮. 物联网安全综述[J]. 计算机研究与发展, 2017, 54(10): 2130-2143.

[2] Kumar R, Kumar P, Tripathi R, et al. A distributed intrusion detection system to detect DDoS attacks in blockchain-enabled IoT network[J]. Journal of Parallel and Distributed Computing, 2022, 164: 55-68.