

面向无人机网络攻击检测的两级特征选择方法

邓琬中¹, 文新², 李维皓¹, 张玮石¹, 白猛¹

(1. 华北计算机系统工程研究所, 北京 100083;

2. 中国电子信息产业集团有限公司, 广东 深圳 518057)

摘要: 针对机器学习算力需求高与无人机 (Unmanned Aerial Vehicles, UAVs) 计算资源有限的矛盾, 以及传统基于固定分箱的信息增益特征选择方法存在高判别力特征被低估的缺陷, 提出一种基于两级特征选择方法的无人机入侵检测方案。该方案采用“卡方检验初筛—启发式信息增益搜索精选”方法, 为每个特征自适应确定最优分箱, 从而准确量化其判别能力。同时将特征数量作为超参数, 与 XGBoost 分类器协同优化。基于 UAV-NIDD 数据集的实验表明, 该方法在保持高检测性能的同时, 将模型检测时间减少了约 97.5%。实验验证, 该方案有效平衡了检测精度与计算开销, 可为资源受限的无人机平台提供高效实时的入侵检测能力。

关键词: 无人机; 入侵检测; 特征选择; 信息增益; 自适应离散化

中图分类号: TP393

文献标志码: A

DOI: 10.19358/j.issn.2097-1788.2026.04.005

中文引用格式: 邓琬中, 文新, 李维皓, 等. 面向无人机网络攻击检测的两级特征选择方法[J]. 网络安全与数据治理, 2026, 45(4): 35-44.

英文引用格式: Deng Wanjin, Wen Xin, Li Weihao, et al. A two-stage feature selection method for network attack detection in UAV networks[J]. Cyber Security and Data Governance, 2026, 45(4): 35-44.

A two-stage feature selection method for network attack detection in UAV networks

Deng Wanjin¹, Wen Xin², Li Weihao¹, Zhang Weishi¹, Bai Meng¹

(1. National Computer System Engineering Research Institute of China, Beijing 100083, China;

2. China Electronics Corporation, Shenzhen 518057, China)

Abstract: To address the dual challenges of high computational demands in machine learning and the limited computing resources of unmanned aerial vehicles (UAVs), as well as the drawback of traditional fixed-binning information gain feature selection methods which underestimate highly discriminative features, this paper proposes a UAV intrusion detection scheme based on a two-level feature selection approach. The scheme adopts a "Chi-square test preliminary screening-Heuristic Information Gain Feature Selection (HIS) fine selection" method, which adaptively determines the optimal binning for each feature to accurately quantify its discriminative power. Meanwhile, the number of features is treated as a hyperparameter and jointly optimized with the XGBoost classifier. Experiments on the UAV-NIDD dataset demonstrate that the proposed method maintains high detection performance while reducing model detection time by approximately 97.5%. The results verify that this scheme effectively balances detection accuracy and computational cost, providing an efficient and real-time intrusion detection capability for resource-constrained UAV platforms.

Key words: UAVs; intrusion detection; feature selection; information gain; adaptive discretization

0 引言

随着低空经济的快速发展, 无人机 (Unmanned Aerial Vehicles, UAVs) 在农业、物流、交通和灾害响应等领域的应用日益广泛。这类无人机具备一定的自主性, 能够提升效率并融入日常生活。然而, 无人机因其开放的通信信道、无保护的传感器输入和有限

的计算防护能力, 容易受到安全威胁攻击, 攻击者会采取各种办法来阻碍、扰乱和控制无人机^[1]。Warakulasooriya 等人^[2]提到, 2012 年到 2022 年间的商业无人机受到多种网络攻击的威胁。Wang 等人^[3]总结出无人机群可能受到的攻击有 20 种, 攻击涉及网络、传感器和机器学习模型。这些网络攻击会

导致无人机在任务执行过程中出错, 不仅损害无人机本身的价值, 还可能造成公众安全风险和经济损失。

入侵检测系统 (Intrusion Detection System, IDS) 为缓解 UAVs 的网络安全威胁提供了一种解决方案^[3]。基于机器学习或深度学习算法构建模型, 通过对网络流量进行分类和异常检测, 来监控和识别 UAVs 通信中的异常行为, 精准检测网络攻击行为, 为后续防御措施的触发提供前提条件。决策树^[4]、随机森林^[5]、支持向量机^[6]构建的 IDS 准确率已高达 95%, 而长短时记忆循环神经网络^[7] (Long Short Term Memory and Recurrent Neural Network, LSTM-RNN)、卷积长短时记忆^[8] (Convolutional Long Short-Term Memory, ConvLSTM) 等深度学习算法的引入, 又将准确率推向了 99% 的高峰。然而, 问题也接踵而至: 计算资源有限的 UAVs 无法承载模型对高算力的要求, 且模型本身的高计算复杂度也产生了部署困难、实时性能低的挑战。因此, 特征工程成为解决问题的关键, 通过特征选择减少数据维度、剔除冗余信息, 不仅能减少模型推理的计算开销、降低模型对算力的需求, 还能提升检测准确率, 避免过拟合风险, 从而使 IDS 更适用于无人机网络安全防护。

目前, 已有若干工作围绕 IDS 模型的特征工程进行了研究。在入侵检测领域, 信息增益 (Information Gain, IG) 是广泛采用的特征选择准则, 但传统方法常依赖固定分箱对连续属性进行离散化, 增加了重要特征信息丢失与漏选风险。Shi 等人^[9]提出基于信息损失与粗糙集的动态离散化算法, 验证了固定分箱在 KDDCup99 数据集上的局限性; Hassan^[10]进一步通过实验证明, 固定分箱会显著降低朴素贝叶斯等分类器的检测效果; Zhang 等人^[11]也指出, 固定分箱容易忽略特征间细微差异, 导致部分高判别性特征被漏选。这些研究共同表明, 传统固定分箱方法在 IG 特征选择中存在明显缺陷, 亟需能够提升 IG 特征选择准确性的方案。

为了应对上述挑战, 本文提出了两级特征选择方法, 设计了启发式信息增益搜索 (Heuristic Information Gain Feature Selection, HIS) 算法, 针对传统信息增益算法的核心缺陷, 即统一分箱策略可能导致重要特征因不恰当离散化而被漏选, 提出了自适应优化方案。该方法通过启发式搜索改进 IG 计算, 精准量化特征的判别能力, 结合卡方检验的初步筛选, 实现了统计显著性筛选与信息量优化的协同, 构建了更具系统性和精准性的特征工程流程, 为无人机入侵检测系统中高维数据的高效处理提供了有力支撑。

本文的主要工作如下:

(1) 提出了新的无人机网络攻击检测方案 TSFS, 解决了机器学习模型的高检测率对算力的高要求和低算力环境之间的矛盾。本文通过特征选择与机器学习模型相结合, 在保证检测性能的前提下大幅降低计算开销和模型复杂度, 使 TSFS 能够在资源受限的边缘设备上实时运行, 解决了高性能检测与低功耗部署之间的核心冲突, 使其能够部署在算力受限的无人机平台上。

(2) 为了解决所提出的问题, 设计了两级特征选择方法, 对每个候选特征自适应搜索出局部最优分箱方案, 准确量化其真实判别能力。同时将所选特征的数量作为超参数纳入模型训练的调优过程, 通过贝叶斯优化实现特征选择与 XGBoost 分类器的协同优化, 实现了以模型性能为导向的特征评估和筛选。

(3) 实验结果表明, 相较传统信息增益方法, 本文的两级特征选择方法能够选出判别能力更强的特征子集, 显著提升分类模型的准确率、召回率和 F1 分数, 并且在准确率超过 95% 的同时, 模型训练时间和推理时间大幅减少, 能够满足实时检测需求。与其他机器学习算法相比, 本文方法在计算效率方面具有明显优势, 更适合部署在资源受限的无人机平台上, 验证了所提方案在实际应用场景中的可行性和有效性。

1 相关工作

近年来, 信息增益及其改进方法在入侵检测特征选择领域得到了广泛应用与深入研究。主流研究主要围绕 IG 的有效性和传统 IG 缺陷的改进展开。

在有效性方面, 一部分研究关注于 IG 对检测率和检测精度的提升。Stiawan 等人^[12]在 CICIDS-2017 数据集上进行特征分析研究, 系统验证了 IG 在异常检测中的有效性。Dicholkar 等人^[13]则进一步将 IG 用于 DoS 攻击检测, 显著提升了检测精度。另外, Sarvari 等人^[14]将 IG 与粒子群优化结合, 优化了 SVM 模型参数, 显著提升了异常入侵检测精度。基于 UNSW-NB15 数据集上的堆叠模型与 IG 特征选择研究^[15]则进一步提升了多类攻击识别率。另一部分研究还额外关注了 IG 对特征降维、鲁棒性等实用性能的提升。Nimbalkar 等人^[16]针对 IoT 环境提出 IG 与增益比的联合特征选择, 在 IoT-BoT 和 KDDCup99 数据集上有效降低特征维度并提升检测率。Fang 等人^[17]将 IG 与机器学习结合, 验证了其在实际网络入侵场景的实用性。Niranjan 等人^[18]提出了 EBJRV 集成分类模型, 以 IG 为核心特征选择手段, 通过多模型投票策略提高入侵检测的鲁

棒性。

在传统 IG 缺陷改进方面, Dong 等人^[19]将信息增益比用于缓解多值特征偏好问题, 提升特征子集的判别力。而 Abdullah 等人^[20]运用 IGRF-RFE 混合递归特征消除方法一定程度上弥补了传统 IG 的不足, 进一步提升了多类攻击检测性能。Anagha 等人^[21]对比了 IG 与 SHAP (SHapley Additive exPlanations, 夏普利加性解释)、相关性等方法, 指出 IG 在部分场景下的不足, 并提出改进思路。这些研究发现了传统 IG 的一些不足之处并提出了改进方法, 却并未注意到传统 IG 在固定分箱策略下的缺陷。

信息增益及其改进型特征选择方法在提升入侵检测效率与精度方面发挥了重要作用, 但针对固定分箱造成重要特征漏选风险这一关键挑战, 尚未出现应对之策。

2 问题描述

2.1 多分类问题建模

无人机网络攻击检测的核心任务是从网络流量数据中识别恶意攻击行为, 区分正常流量与多种攻击类型, 例如 DDoS 攻击、劫持攻击、中间人攻击等。本文将前述实际问题建模为监督学习框架下的多分类问题。给定训练数据集 $D = \{(x_i, y_i)_{i=1}^N\}$, 其中 $x_i = (X_{1i}, X_{2i}, \dots, X_{ni}) \in \mathbf{R}^n$ 表示第 i 个样本的特征向量, $y_i \in \{0, 1, 2, \dots, C-1\}$ 表示其类别标签, 且 $y_i = 0$ 代表正常流量, $y_i \in \{1, 2, \dots, C-1\}$ 代表不同类型的攻击, n 为原始特征维度, N 为样本总数, C 为类别个数。检测任务是学习一个分类函数 $f: \mathbf{R}^n \rightarrow$

$\{0, 1, \dots, C-1\}$, 使其在测试集上的分类误差最小。

该多分类问题的完整求解流程如图 1 所示, 包括特征工程、模型训练和模型部署三个阶段。在特征工程阶段, 需要从原始流量数据包中提取统计特征, 并通过特征选择方法降低维度。在模型训练阶段, 基于选出的特征子集训练分类模型并优化超参数。在模型部署阶段, 对实时流量进行特征提取和分类推理。本文的研究重心聚焦于特征工程阶段中的特征选择环节, 旨在从高维原始特征空间中筛选出具有强判别能力的低维特征子集, 为下游分类模型提供高质量输入, 从而在保证检测精度的前提下显著降低计算开销, 满足实时检测的需求。

2.2 传统方法的局限性

传统 IG 方法通过量化特征对类别不确定性的降低程度来评估其判别能力, 具备理论扎实、计算高效的优势。然而, 传统 IG 在处理连续特征时存在关键缺陷。如图 2 所示, 计算离散化后的 IG 值, 需对连续特征进行分箱处理, 现有方法普遍采用统一的分箱策略。这种“一刀切”方案忽略了特征分布的差异性——不同特征可能呈正态分布、偏态分布或存在离群值, 统一分箱策略无法适配所有特征。甚至, 若分箱点设置不当, 同类样本被人为分散到不同箱中, 将导致条件熵增大, IG 值被系统性低估。这些缺陷致使高判别能力的特征可能因分箱不当在排序中靠后而被错误淘汰, 导致特征子集代表性不足, 制约了模型性能上限。

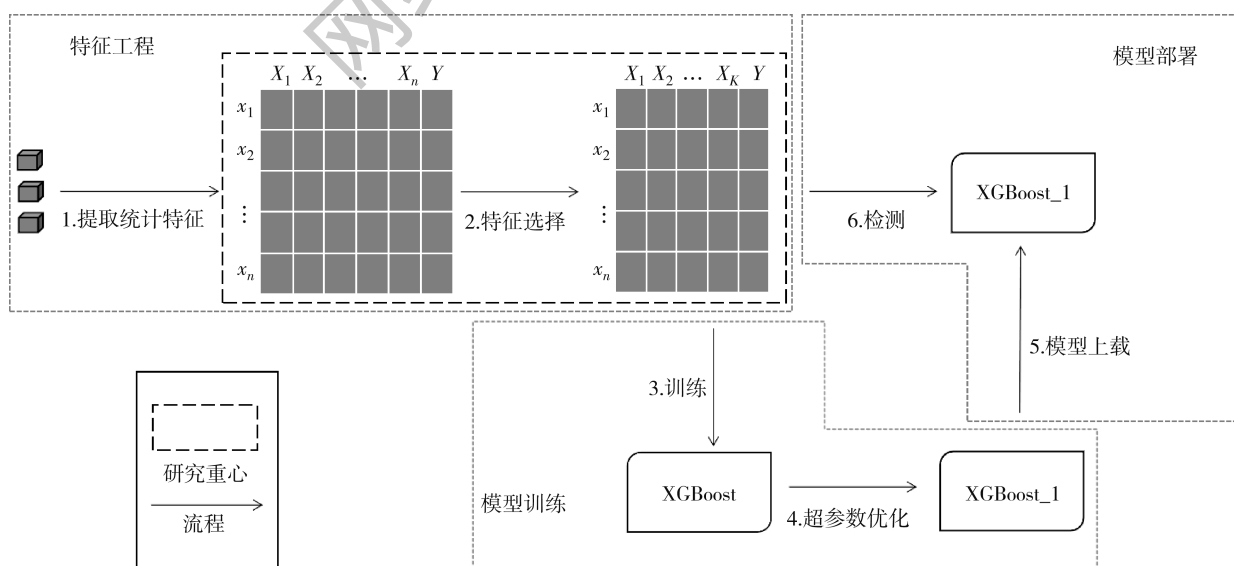


图 1 多分类问题求解流程

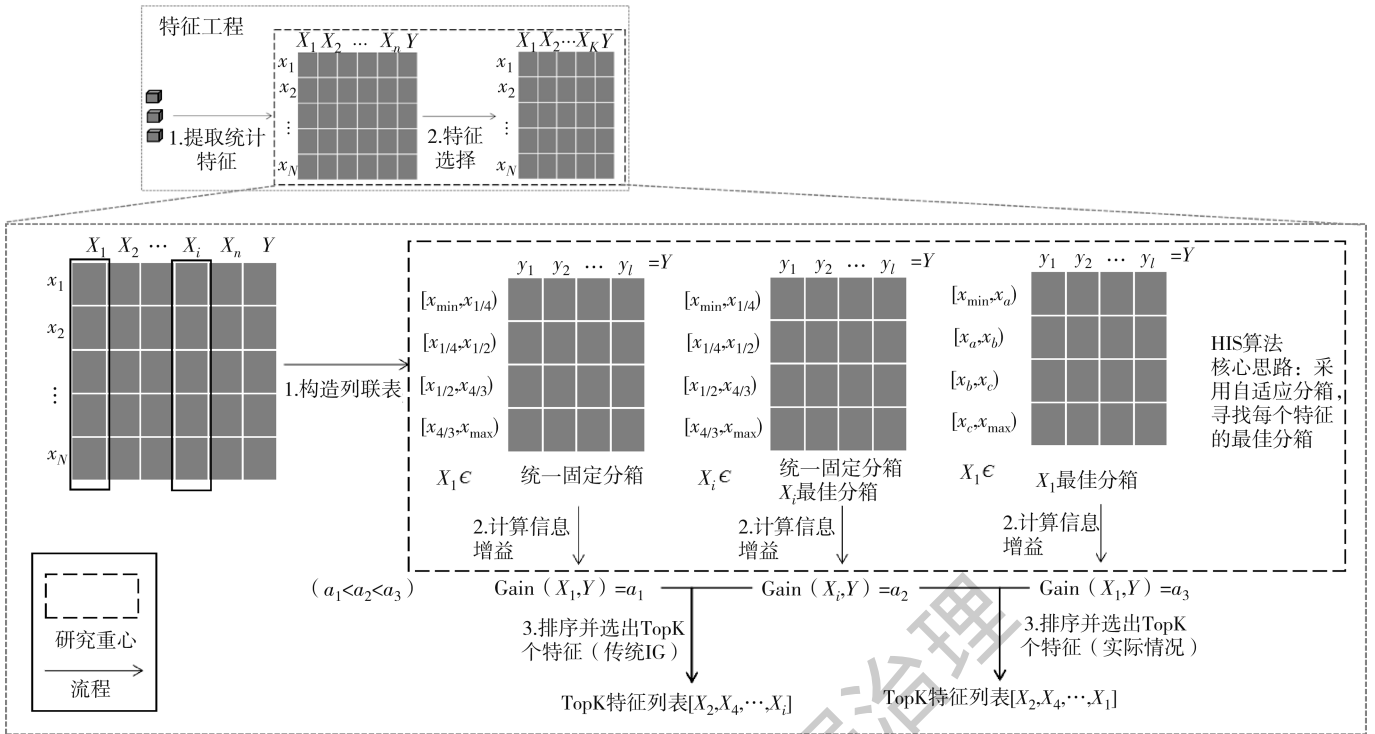


图2 传统 IG 的关键缺陷和 HIS 改进

此外,从协同优化角度看,传统方法将特征选择与模型训练视为独立串行环节,并且特征数量等关键参数依赖经验法则或反复试错,缺乏联合优化机制。这种分离式范式忽略了特征选择对模型性能的深层影响及训练过程中的反馈信息,难以实现全局最优。在无人机网络攻击检测场景中,仅依靠 IG 值决策,无法权衡分类性能与计算效率的多目标需求,将导致特征子集理论指标优异但实际部署时无法满足系统综合性要求。

2.3 本文的优化目标与问题分解

本文要解决的核心问题可形式化描述为:给定原始特征集 $F = \{x_1, x_2, \dots, x_n\}$ 和带标签的训练样本 D ,目标是选出优化后的特征子集 $F^* \subset F$,满足 $|F^*| = K \ll n$,使得 F^* 基于训练的分类模型 $M(F^*)$ 在分类性能和计算效率上达到最优平衡。

$$F^* = \underset{F' \subset F, |F'|=K}{\operatorname{argmax}} [\alpha \cdot \operatorname{Acc}(M(F')) + \beta \cdot F_1(M(F')) - \gamma \cdot T(M(F'))] \quad (1)$$

其中 $\operatorname{Acc}(\cdot)$ 表示准确率, $F_1(\cdot)$ 表示 F1 分数, $T(\cdot)$ 表示推理时间的总开销, α, β, γ 为权重系数。该优化问题需要满足以下约束条件:(1) 特征子集 F^* 具有强判别能力,能够充分表征不同类别的差异;(2) 特征间冗余度低,避免信息重复;(3) 特征数量 K 在合理范围内,兼顾模型性能与实时性

要求。

为求解上述优化问题,本文将其分解为两个递进的子问题。首先,针对原始高维特征集,采用卡方检验进行统计显著性初筛,快速剔除与类别标签独立的无关特征,从 n 个原始特征中选出 m 个候选特征构成集合 $F'(m < n)$ 。其次,针对传统信息增益方法因固定分箱策略导致高贡献特征被低估和遗漏的问题,设计 HIS 算法,为每个候选特征自适应寻找局部最优分箱策略。对于特征, HIS 算法在所有可能的分箱策略集合中搜索最优方案,通过计算所有候选特征在各自最优分箱下的最大 IG 值并按降序排列,选出前 K 个 IG 值最大的特征构成最优特征子集 F^* ,并将 K 作为超参数纳入分类模型的优化范围,通过网络搜索与交叉验证同时优化特征选择参数和 XGBoost 模型的训练参数,实现特征选择与模型训练的协同优化,让下游模型的实际性能反向指导特征选择过程。

据此,本文构建了“卡方检验初筛—HIS 算法精选”的两级特征选择方法,将特征工程与模型训练深度融合,为无人机网络攻击检测提供了一套系统化的优化方案。该方法的核心创新在于 HIS 算法的自适应分箱机制,能够充分挖掘每个特征的真实判别能力,避免传统方法因分箱不当导致的特征遗漏问题。

3 基于两级特征选择方法的 UAVs 网络攻击检测方案 TSFS

3.1 卡方检验

卡方检验是对已经归一化到 $[0, 1]$ 区间的数据进行特征与标签的独立性分析, 其步骤为:

(1) 构造列联表, 按 $[0, 0.25)$, $[0.25, 0.5)$, $[0.5, 0.75)$, $[0.75, 1]$ 对归一化后的连续值进行分箱将其离散化为有限类别, 统计每个特征类别与标签类别组合的样本个数; 以特征类别为行, 标签类别为列, 填充统计出的频数。

(2) 确定显著水平 α 、自由度 df 和临界值 b , 按照业界统一标准, α 通常取值 0.05, df 通过 $(\text{列联表行数} - 1) \times (\text{列联表列数} - 1)$ 计算得到, b 依据 α 和 df 的取值在卡方分布表中查找得到。

(3) 计算卡方统计量, 卡方检验的统计量 X^2 为:

$$X^2 = \sum \sum \frac{(O_{ij} - E_{ij})^2}{E_{ij}} \quad (2)$$

其中 O_{ij} 表示在第 i 行、第 j 列的实际频数, E_{ij} 表示在第 i 行、第 j 列的期望频数, 期望频数为:

$$E_{ij} = \frac{(R_i \cdot C_j)}{N} \quad (3)$$

其中 R_i 表示第 i 行的总数, C_j 表示第 j 列的总数, N 表示总的样本数量。

(4) 比较统计量 X^2 与临界值 b 的大小, 若 $X^2 < b$ 则特征与标签相互独立, 去除该特征; 反之则特征与标签不独立, 保留该特征。

实验验证, 经过卡方检验初筛后, 数据集 UAV-NIDD^[22] 的特征数量从编码后的 50 个减少至 48 个。

3.2 HIS 算法

HIS 是使用启发式搜索改进后的 IG 特征选择算法, 保证了计算出的每个特征的 IG 值都是该特征在不同分箱策略下局部最大的 IG 值, 一定程度上降低了贡献大的特征因不恰当的分箱而被漏选的可能性。与传统的 IG 特征选择算法相比, HIS 能够选出更具代表性的最优特征子集。IG 是特征选择中常用的一种方法, 利用熵和条件熵的概念, 通过计算特征对样本按标签分类的影响来选择最佳特征。熵表示数据集的不确定性, 其定义为:

$$H(I) = - \sum (p \cdot \log_2 p) \quad (4)$$

其中 p 表示数据集中每个标签的概率, 通常由每个标签中的样本数量和样本总数计算得出。对于每个特征 X , 计算该特征给定条件下的数据集的条件熵。条件熵表示在特征 X 的给定条件下的数据集的不确定性, 其

定义为:

$$H(I|X) = \sum H(I|X=x_i) \left(\frac{|I_{X=x_i}|}{|I|} \right) \quad (5)$$

其中 $H(I|X=x_i)$ 表示特征 X 取值为 x_i 条件下的数据集条件熵, $|I|$ 表示数据集的样本总数, $|I_{X=x_i}|$ 表示特征 X 下属于第 i 个属性的样本数。特征 X 在数据集上的信息增益 $\text{Gain}(I, X)$ 表示特征 X 对样本分类的影响, 其定义为:

$$\text{Gain}(I, X) = H(I) - H(I|X) \quad (6)$$

传统的 IG 特征选择算法采用统一的分箱策略对特征连续值进行离散化, 算法的伪代码如算法 1 所示。

算法 1 传统 IG 特征选择算法

输入: 数据集 D , 特征集 F

输出: 降序排序的特征子集 F_k

- 1) 初始化 $F_k, S[\]$
- 2) 计算整体熵 $H(Y) = - \sum_{y \in Y} p(y) \cdot \log_2 p(y)$
- 3) for X_i in F :
- 4) 计算条件熵 $H(Y|X_i)$
- 5) 计算信息增益 $\text{Gain}(Y, X_i)$
- 6) $S \leftarrow (X_i, \text{Gain}(Y, X_i))$
- 7) end for
- 8) 按信息增益降序排序 S
- 9) $F_k \leftarrow \text{topK}(S)$

本文设计的 HIS 算法利用启发式搜索进行优化, 通过寻找对于每个特征而言局部最佳的分箱策略, 计算并保存在该分箱策略下的局部最大信息增益值, 降低重要特征被漏选的概率。启发式搜索的初始分箱策略与卡方检验的统一分箱策略相同。此外, 将 $[0, 1]$ 区间的三个分割点设置为初始分割点, 再设置一个常量 0.05, 后续随机选择一个分割点, 该分割点加、减常量后形成新分割点, 新分割点和两个旧分割点组成两个邻近分箱策略。启发式搜索的终止条件为当前分箱策略下的信息增益值大于两个邻近分箱策略下的信息增益值。HIS 算法的伪代码如算法 2 所示。

算法 2 HIS 算法

输入: 数据集 D , 特征集 F

输出: 降序排序的特征子集 F^*

函数 TS

输入: 三个分割点 a_1, a_2, a_3

输出: a_1, a_2, a_3 组成的分箱策略下的条件熵

- 1) 以升序排列 a_1, a_2, a_3
- 2) 按 a_1, a_2, a_3 进行分箱

```

3) 计算条件熵  $H(Y | X_i) = \sum H(Y | X_i = x) \left( \frac{|Y_{X_i=x_i}|}{|Y|} \right)$ 
4) 返回  $H(Y | X_i)$ 
结束函数
1) 初始化  $F^*$ ,  $S$  []
2) 计算整体熵  $H(Y) = - \sum_{y \in Y} p(y) \cdot \log_2 p(y)$ 
3) for  $X_i$  in  $F$ :
4)   初始化分箱策略  $a_1 = 0.25, a_2 = 0.5, a_3 = 0.75,$ 
 $b = 0.05$ 
5)   初始化条件熵  $m = \text{TS}(a_1, a_2, a_3)$ 
6)   while true do
7)     随机选择一个分割点  $a = \text{random}(a_1, a_2, a_3)$ 
8)     计算邻近条件熵  $H_1 = \text{TS}(\dots, a + b, \dots), H_2 =$ 
 $\text{TS}(\dots, a - b, \dots)$ 
9)     if  $\min(H_1, H_2, m) == m$  then break
10)    else if  $\min(H_1, H_2, m) == H_1$  then
11)       $m = H_1$ 
12)       $a = a + b$ 
13)    else
14)       $m = H_2$ 
15)       $a = a - b$ 
16)    end if
17)  end while
18) 计算信息增益值  $\text{Gain}(Y, X_i) = H(Y) - m$ 
19)  $S \leftarrow (X_i, \text{Gain}(Y, X_i))$ 
20) end for
21)  $F^* \leftarrow \text{topK}(S)$ 

```

注: $\text{TS}(\dots, a + b, \dots)$ 表示随机选出的分割点加 0.05 后计算出的条件熵, $\text{TS}(\dots, a - b, \dots)$ 表示随机选出的分割点减 0.05 后计算出的条件熵

综上, HIS 算法是一种针对传统 IG 特征选择方法的改进算法, 其核心作用在于通过为每个特征自适应寻找局部最优分箱策略, 计算并保留该特征在不同离散化方案下的最大信息增益值, 从而有效降低重要特征因不恰当分箱而被漏选的风险。该算法突破了传统特征选择中固定预处理范式的局限, 为特征工程提供了新的优化维度, 使选出的特征子集更具代表性和判别能力, 并能够直接提升下游机器学习模型的分率准确率和泛化性能, 对连续特征占比高的数据集尤为有效。

3.3 模型训练和测试

本研究选用 XGBoost 作为核心检测模型。XGBoost 是一种高效的梯度提升决策树 (Gradient Boosting Decision Tree, GBDT), 通过集成多棵浅层决策树构建的模型, 每轮迭代用新树拟合当前模型的残差以逐步降低

损失。相较于经典 GBDT, XGBoost 在三个方面实现了显著增强。一是优化策略方面, 在目标函数中引入惩罚叶子数量与权重的正则化项, 并利用一阶梯度与二阶梯度 Hessian 信息评估分裂收益, 提升拟合的稳定性和准确性。二是工程加速方面, 采用直方图近似算法和特征级并行加速最佳切分点搜索, 原生支持稀疏输入和缺失值的自适应处理, 通过列存储、块结构和外存计算优化内存使用, 支持多线程和 GPU 加速。三是过拟合控制方面, 集成学习率收缩、行列子采样、L1/L2 正则化、分裂增益阈值和早停等机制。与随机森林的并行“多数表决”机制相比, XGBoost 采用串行“残差纠错”策略显著降低偏差; 与经典 GBDT 相比, 二阶梯度信息和工程优化使其训练速度更快、稳定性更强; 与 LightGBM 相比, XGBoost 在小数据集或噪声环境下更稳定且生态更完善; 与 CatBoost 相比, XGBoost 在数值特征或预编码特征上表现更优。综合准确率、运行效率和工程适配性等优势, 本研究采用 XGBoost 分别基于网络数据训练集构建模型, 并经过超参数优化后得到 XGBoost_1 作为检测模型。

模型训练在离线环境中进行, XGBoost_1 使用包含网络攻击的公开数据集训练。首先对原始数据集进行归一化处理以统一量纲, 然后应用两级特征选择方法保留特征子集, 再进行标准化将数据转换为均值为 0、方差为 1 的分布。标准化后的训练数据集用于模型训练, 训练过程采用贝叶斯优化方法寻找最佳超参数组合。需要特别指出的是, 特征子集的特征数量 K 也作为超参数纳入优化范围, 以匹配最优特征子集。超参数优化的评价指标设定为模型准确性、F1 分数和模型检测时间。本文完成 XGBoost_1 的训练和超参数优化, 并将其用于验证两级特征选择方法的有效性。模型测试流程与部署后的检测流程基本一致, 测试数据经归一化处理, 保留训练阶段确定的特征子集并进行标准化, 标准化后的测试数据输入模型进行分类, 正常数据直接标记输出, 攻击数据附加对应攻击类型标签后输出, 以此评估模型性能。

3.4 超参数优化

贝叶斯优化是一种基于贝叶斯统计的全局优化方法, 用于高效搜索复杂黑箱函数的最优参数。它通过构建一个代理模型来近似目标函数的分布, 利用先前评估结果更新后验概率分布。然后, 通过采集函数选择下一个采样点, 以平衡探索不确定区域和利用已知高性能区域。这种迭代过程逐步逼近全局最优解, 特别适用于机器学习模型的超参数调优。相比于网格搜

索,它能够动态调整搜索策略,智能地避开无效区域,效率更高,并且能够更好地探索全局最优,降低了错过最优解的可能性。所以,本文选择贝叶斯优化对 XGBoost 模型进行超参数调优。

本文将贝叶斯优化的目标函数设计为一个复合分数,综合考虑模型性能、推理时间和模型大小,以平衡准确性和效率。该目标函数的数学表达式为:

$$\text{score} = \frac{(\text{accuracy} + f1_{\text{macro}})}{2} - \frac{\text{val}_{\text{time}}}{100} - \frac{\text{model}_{\text{size}}}{1 \times 10^6} \quad (7)$$

其中 score 为目标函数的值, accuracy 为模型的准确率, $f1_{\text{macro}}$ 为模型的 F1 分数,后两者基于验证集计算得到。 val_{time} 为模型检测时间,通过时间计时量化; $\text{model}_{\text{size}}$ 为模型的大小,通过 pickle 序列化字节数量化,这两者均用于对高计算和高存储开销的配置施加惩罚。

参考 XGBoost 的官方文档和文献 [23],本文选取了 XGBoost 模型的 5 个参数进行优化,并额外添加输入模型的特征个数作为第 6 个优化的参数,以便探索最佳的特征个数,同时还参考典型实践设置了优化边界。具体的参数和优化边界如表 1 所示。

表 1 XGBoost 模型调优参数和优化边界表

参数名	含义	优化边界
max_depth	树的最大深度	[3, 10]
Learning_rate	学习率	[0.01, 0.3]
n_estimators	树的数量	[50, 200]
subsample	子采样率	[0.5, 1.0]
min_child_weight	叶子节点最小权重	[1, 10]
n_features	输入模型的特征个数	[1, 49]

经过 5 次初始随机采样和 20 次迭代后,获取到的模型最佳参数组合如表 2 所示。

表 2 XGBoost 模型最佳参数表

参数名	最佳值
max_depth	3
Learning_rate	0.3
n_estimators	50
subsample	1.0
min_child_weight	10
n_features	43

4 实验验证

本次实验旨在利用 UAV-NIDD 数据集和 XGBoost_1

模型验证两级特征选择方法的有效性。实验在一个配备 4 核处理器、4 GB RAM 的虚拟环境中进行,运行 Ubuntu-22.04 操作系统(基于 Linux 内核),使用 Python 3.9 和 Scikit-learn、XGBoost、Numpy、Pandas、Scipy 库。训练数据集和测试数据集由 UAV-NIDD 数据集以 8:2 的比例划分获得,使用的 XGBoost_1 模型已按照贝叶斯优化获取到的最佳参数组合进行设置,且已利用训练数据集完成训练。

实验设置有三个对照组,分别为无特征选择组(简称 NOT 组)、单 HIS 特征选择组(简称单 HIS 组)、传统 IG 特征选择组(简称 TI 组)。NOT 组不进行特征选择,单 HIS 组只运行 HIS 算法进行特征选择, TI 组只运行传统 IG 特征选择算法进行特征选择,与采用两级特征选择方法的实验组(简称两级组)进行对比。实验选取 HIS 算法的运行时间以及后续模型的准确率、F1 分数和运行时间等指标进行量化分析。其中模型的运行时间为 XGBoost_1 模型使用测试数据集进行测试的时间。

4.1 功能试验

从模型准确性的角度分析,如图 3(a)所示, NOT 组、单 HIS 组、TI 组和两级组的模型准确率分别为 96.87%、98.64%、97.21%、98.92%,再结合图 3(b),模型 F1 分数分别为 96.72%、98.59%、97.05%、98.88%,两级组和单 HIS 组在模型准确率和 F1 分数上均优于其他两组。这证明无论是 HIS 还是两级特征选择方法对模型准确性的负面影响都小于不做特征选择和传统 IG,两者均没有漏选重要特征,且两级特征选择方法挑选出的特征更具代表性。

4.2 性能试验

从时间上分析,如图 4(a)所示,两级组的算法运行时间为 10.161 s,小于单 HIS 组的算法运行时间,进一步分析得出两级组的算法运行时间比单 HIS 组降低了约 22.5%,证明在运行 HIS 算法前使用卡方检验进行初筛能够有效降低算法的运行时间。传统 IG 算法的时间复杂度为 $O(n)$,而两级特征选择算法时间复杂度为 $O(n) + O(hn)$,其中 h 为平均每个特征的搜索次数, n 为特征数量。显然,传统 IG 比两级特征选择算法快,但后者不是针对前者运行时间的改进,所以比较两者的运行时间没有意义。

从模型检测时间上看,如图 4(b)所示,单 HIS 组的模型检测时间为 1 254.459 ms,两级组的模型检测时间为 51.51 ms,后者小于 NOT 组和 TI 组,前者小于 TI 组但大于 NOT 组。这是因为前者未剔除无关的冗余

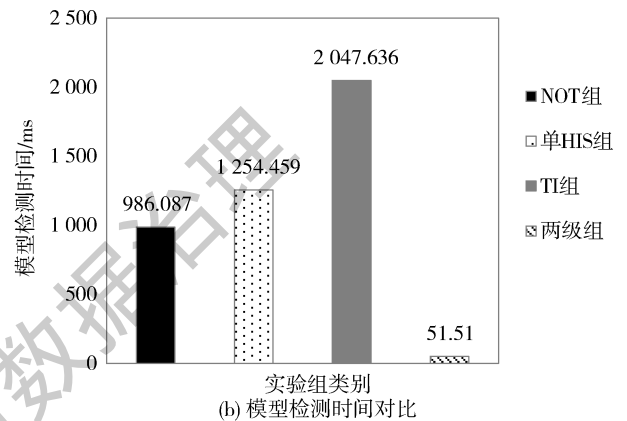
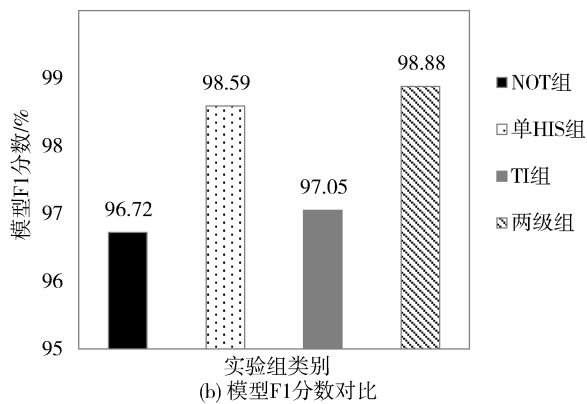
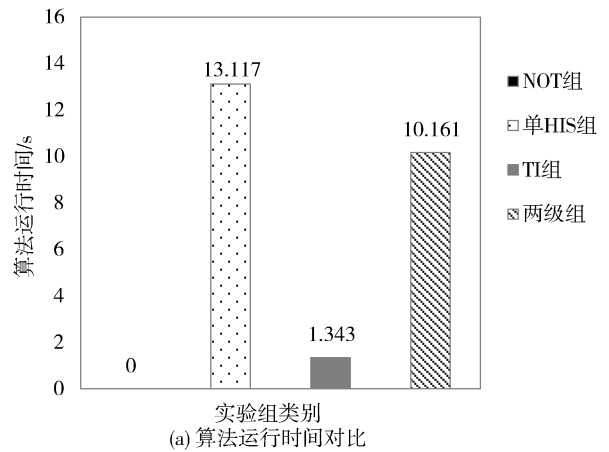
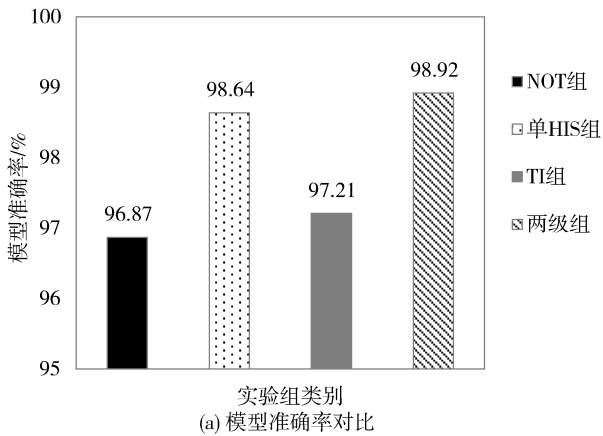


图3 不同实验组在准确性指标上的对比

图4 不同实验组在时间指标上的对比

特征, 导致时间有所增加。进一步分析得出单 HIS 组相比于 TI 组模型检测时间减少了约 38.7%, 两级组相比于 NOT 组和 TI 组分别减少了约 94.8% 和 97.5%, 证明 HIS 算法与传统 IG 相比一定程度上缩短了模型的检测时间, 而两级特征选择方法能够大幅缩短模型的检测时间, 比前两者性能更好。

由上述实验结果可得出, 两级特征选择方法能够挑选出最具代表性的特征, 减少对模型预测精度的影响, 且能够降低模型的检测时间。这验证了两级特征选择方法的有效性, 有利于未来将模型部署到资源受限的无人机上。仅使用 HIS 算法比传统 IG 更有效, 但由于其未剔除无关特征, 检测时间比无特征选择增加 27.2%。

同时, 本文还将提出的 TSFS 方案与业内其他方案

进行了对比, 如表 3 所示。对比发现, TSFS 方案在准确率、F1 分数和模型检测时间上均优于其他方案。

5 结论

本文针对无人机入侵检测系统在资源受限环境下面临的高性能检测与低算力部署环境之间的矛盾, 提出了一种基于两级特征选择方法与 XGBoost 模型相结合的解决方案。针对传统信息增益方法因固定分箱策略导致高判别能力特征被系统性低估和漏选的核心缺陷, 本文设计了“卡方检验初筛—HIS 算法精选”的两级特征选择方法, 通过为每个候选特征自适应搜索局部最优分箱方案, 准确量化其真实判别能力, 并将特征数量作为超参数纳入贝叶斯优化过程, 实现了特征选择与模型训练的协同优化。

表3 部分相关工作的方案和 TSFS 对比

方法	核心特征选择/模型设计	准确率/%	F1 分数/%	模型检测时间/ms
堆叠机器学习 ^[15]	堆叠算法	96.24	96.29	—
杜鹃搜索算法 ^[14]	Cuckoo Search Algorithm	98.81	—	—
孤立森林 ^[22]	Isolation Forest	71.98	72.32	2 703.8
前馈神经网络 ^[22]	Feedforward neural network	87.32	86.92	25 063
TSFS	两级特征选择 + XGBoost	98.92	98.88	51.51

实验结果表明,相较于传统方法,本文提出的两级特征选择方法在保持准确率达 98.92%、F1 分数达 98.88% 的高检测精度的同时,将模型运行时间降低约 97.5%,且算法运行时间相比单一使用 HIS 降低约 22.5%,验证了该方法在分类性能与计算效率之间实现了良好平衡,为无人机入侵检测系统在边缘设备上的实时部署提供了可行方案。

未来研究可从以下方向展开:(1)当前工作将两级特征选择算法置于线下,无法适应演化后的攻击模式,未来可以结合在线学习或持续学习等技术探索在线自适应特征选择机制,使系统能够根据实时流量特征动态调整特征子集,以应对攻击模式演化;(2)当前的启发式搜索工作仅动态调整分割点,并未对箱子个数进行动态调整,未来可将箱子个数也用于生成邻近分箱策略,为寻找最佳分箱提供更大的空间;(3)当前工作仅在数据层面进行压缩,未来可进一步研究模型压缩与量化策略,探索在更严苛的资源约束下的部署可能性,为无人机网络安全防护提供更全面的技术支撑。

参考文献

- [1] NING Z, HU H, WANG X, et al. Joint user association, interference cancellation and power control for multi-IRS assisted UAV communications[J]. *IEEE Transactions on Wireless Communications*, 2024, 23 (10): 13408 - 13423.
- [2] WARNAKULASOORIYA K, SEGEV A. Attacks, detection, and prevention on commercial drones: a review[C]//*Proceedings of 2024 International Conference on Image Processing and Robotics (ICIPRoB)*, 2024: 1 - 6.
- [3] WANG X, ZHAO Z, YI L, et al. A survey on security of UAV swarm networks; attacks and countermeasures[J]. *ACM Computing Surveys*, 2025, 57 (3): 1 - 37.
- [4] SHRESTHA R, OMIDKAR A, ROUDI S A, et al. Machine-learning-enabled intrusion detection system for cellular connected UAV networks[J]. *Multidisciplinary Digital Publishing Institute (MDPI), Electronics*, 2021, 10 (13): 1549.
- [5] IHEKORONYE V U, AJAKWE S O, KIM D S, et al. Hierarchical intrusion detection system for secured military drone network: a perspicacious approach[C]//*Proceedings of 2022 IEEE Military Communications Conference (MILCOM)*, 2022: 336 - 341.
- [6] TUFEKCI B, QUACH V, TUNC C, et al. DUDE-IDS: a framework for efficiently detecting network-related drone cyberattacks[C]//*Proceedings of 2024 11th International Conference on Internet of Things: Systems, Management and Security (IOTSMS)*, 2024: 240 - 247.
- [7] RAMADAN R A, EMARA A H, AL-SAREM M, et al. Internet of drones intrusion detection using deep learning[J]. *Multidisciplinary Digital Publishing Institute (MDPI), Electronics*, 2021, 10 (21): 2633.
- [8] ALZHRANI A. Novel approach for intrusion detection attacks on small drones using convLSTM model[J]. *IEEE Access*, 2024, 12: 149238 - 149253.
- [9] SHI Z, XIA Y, WU F, et al. The discretization algorithm for rough data and its application to intrusion detection[J]. *Journal of Networks*, 2014, 9 (6): 1380 - 1387.
- [10] HASSAN D. Supervised versus unsupervised discretization for improving network intrusion detection[J]. *International Journal of Computer Science and Information Security (IJCSIS)*, 2016, 14 (10).
- [11] ZHANG Y, REN X, ZHANG J. Intrusion detection method based on information gain and ReliefF feature selection[C]//*Proceedings of 2019 International Joint Conference on Neural Networks (IJCNN)*. *IEEE*, 2019: 1 - 5.
- [12] KURNIABUDI, STIAWAN D, DARMAWIJOYO, et al. CIC-IDS - 2017 dataset feature analysis with information gain for anomaly detection [J]. *IEEE Access*, 2020, 8: 132911 - 132921.
- [13] DICHOLKAR S V, NIRMAL J H. DoS attack detection using feature selection with Information Gain and ML classification[C]//*Proceedings of 2024 Second International Conference on Advances in Information Technology (ICAIT)*. *IEEE*, 2024, 1: 1 - 6.
- [14] SARVARI S, SANI N F M, HANAPI Z M, et al. An efficient anomaly intrusion detection method with feature selection and evolutionary neural network [J]. *IEEE Access*, 2020, 8: 70651 - 70663.
- [15] KABIR M H, RAJIB M S, RAHMAN A S M T, et al. Network intrusion detection using UNSW-NB15 dataset; stacking machine learning based approach [C]//*Proceedings of 2022 International Conference on Advancement in Electrical and Electronic Engineering (ICAEET)*. *IEEE*, 2022: 1 - 6.
- [16] NIMBALKAR P, KSHIRSAGAR D. Feature selection for intrusion detection system in Internet-of-Things (IoT) [J]. *ICT Express*, 2021, 7 (2): 177 - 181.
- [17] FANG W, TAN X, WILBUR D. Application of intrusion detection technology in network safety based on machine learning[J]. *Safety Science*, 2020, 124: 104604.
- [18] NIRANJAN A, PRAKASH A, VEENA N, et al. EBJRV: an ensemble of bagging, J48 and random committee by voting for efficient classification of intrusions [C]//*Proceedings of 2017 IEEE International WIE Conference on Electrical and Computer Engineering (WIECON-ECE)*. *IEEE*, 2017: 51 - 54.
- [19] DONG R H, YAN H H, ZHANG Q Y. An intrusion detection model for wireless sensor network based on information gain ra-

- tio and bagging algorithm[J]. International Journal of Network Security, 2020, 22 (2): 218 – 230.
- [20] ABDULLAH M, ALSHANNAQ A, BALAMASH A, et al. Enhanced intrusion detection system using feature selection method and ensemble learning algorithms[J]. International Journal of Computer Science and Information Security (IJCSIS), 2018, 16 (2): 48 – 55.
- [21] ANAGHA A S, THOMAS C, NARAYANASWAMY B. Optimized intrusion predictions through feature selection methods[J]. Computers & Security, 2025, 157: 104541.
- [22] HADI H J, CAO Y, KHAN M K, et al. UAV-NIDD: a dynamic dataset for cybersecurity and intrusion detection in UAV networks[J]. IEEE Transactions on Network Science and Engineering, 2025, 12 (4): 2739 – 2757.
- [23] SNOEK J, LAROCHELLE H, ADAMS R P. Practical bayesian optimization of machine learning algorithms[C]//Proceedings of the 26th International Conference on Neural Information Processing Systems (NIPS'12), 2012: 2951 – 2959.
- (收稿日期: 2026 – 01 – 22)

作者简介:

邓婉巾 (1999 –), 通信作者, 女, 硕士研究生, 主要研究方向: 无人机网络攻击检测、传感器攻击检测。E-mail: 3204657580@qq.com。

文新 (1986 –), 男, 硕士, 正高级工程师, 主要研究方向: 网络空间对抗、人工智能。

李维皓 (1990 –), 女, 博士, 高级工程师, 主要研究方向: 网络安全、隐私保护。

网络安全与数据治理

版权声明

凡《网络安全与数据治理》录用的文章，如作者没有关于汇编权、翻译权、印刷权及电子版的复制权、信息网络传播权与发行权等版权的特殊声明，即视作该文章署名作者同意将该文章的汇编权、翻译权、印刷权及电子版的复制权、信息网络传播权与发行权授予本刊，本刊有权授权本刊合作数据库、合作媒体等合作伙伴使用。同时，本刊支付的稿酬已包含上述使用的费用，特此声明。

《网络安全与数据治理》编辑部

www.pcachina.com