不平衡数据环境下基于 GRU-CNN 模型的网络安全检测

熊天运¹ 韦 富¹ XIONG Tianyun WEI Fu

摘要

为解决现有方法在处理不平衡数据时性能欠佳、网络安全检测延迟较长以及误检率较高等问题,文章提出一种在不平衡数据环境下,基于 GRU-CNN 模型的网络安全检测方法。应用混合采样算法处理网络运行不平衡数据,采用并行结构改进 GRU-CNN 模型,通过 GRU 捕捉网络运行数据中的时序特征,利用 CNN 提取网络运行数据中的空间特征,融合处理时序特征与空间特征,计算其与已有网络异常特征集合之间相关系数,与制定阈值进行比较,从而实现网络安全的有效检测。实验结果显示,该方法应用后不同类别编号数据集合中的数据量极为相近,网络安全检测延迟最小值为 3 min,网络安全误检率最小值为 3.2%。

关键词

网络安全检测; 改进 GRU-CNN 模型; 数据特征提取与融合; 不平衡数据环境; 损失函数

doi: 10.3969/j.issn.1672-9528.2025.02.018

0 引言

不平衡数据环境下的网络安全检测应用模型往往难以同 时满足实时性和准确性的要求,需要在应用模型性能、实时 性及其准确性之间进行权衡和折中[1-2]。针对网络安全检测方 法进行深入的研究, 对保护个人隐私、维护国家安全和社会 稳定、保障经济活动正常运行等方面具有重要意义。文献[3] 针对网络安全事件标注数据匮乏的问题,通过提示问答的方 式获取事件表示知识,生成伪标注数据以增强模型训练;文 献[4]将深度学习与入侵检测相结合,利用深层胶囊网络提 取高维特征,混合注意力机制关注重要特征,双重路由算法 捕捉向量特征并进行聚类; 文献 [5] 针对网络安全态势感知 算法检测准确率低、误差大的问题,采用大数据分析方法分 解入侵信息特征。现有的方法虽然降低误差并提高预测精确 度,但对大数据处理要求较高,若数据量庞大或特征复杂, 将增加计算难度和时间成本。为了满足网络安全运行需求, 本文提出了一种不平衡数据环境下基于 GRU-CNN 模型的网 络安全检测方法研究。

1 不平衡数据处理

在本文方法探究过程中,不平衡数据处理是首要任务, 混合采样算法便是解决这一问题的有效手段之一,具体过程 如下:

(1)设置网络运行不平衡数据集合为 $\hat{X} = \{X_1, X_2, \cdots, X_n, \cdots, X_n\}$,N表示数据类别的总数量。以集合数据量为依据,制定

1. 广西理工职业技术学院 广西崇左 532200

多数类数据集合与少数类数据集合判定规则,表达式为:

$$\begin{cases}
Q_i \ge \hat{Q} & X_i \to 3 & \text{3} \\
Q_i < \hat{Q} & X_i \to 3 & \text{3} \\
\end{pmatrix}$$
(1)

式中: Q_i 表示平衡数据集合 X_i 数据量; \hat{Q} 表示多数类数据集合与少数类数据集合判定阈值,其需要根据实际网络运行不平衡数据情况进行具体的设置。

依据式 (1) 确定多数类数据集合为 $\{X_1, X_2, \dots, X_L\}$,少数类数据集合为 $\{X_1, X_2, \dots, X_S\}$ 。其中,L、S 表示多数类数据集合与少数类数据集合的数量,并且 $L+S=N^{[6]}$ 。

(2)采用 Tomek Links 算法对多数类数据集合进行 欠采样处理,以数据集合 $X_L=\{X_{L1},X_{L2},\cdots,X_{Lj},\cdots,X_{LN}\}$ (N_L 表示数据总量)为例,计算任意两个数据之间的距 离 $d(X_{Lj},X_{Lk})$,若其小于数据 X_{Lj} 与剩余全部数据之间的距离 $d(X_{Lj},X_{Lp})$, $p=1,2,\cdots,N_L$ 且 $p\neq j,k$,则认定 (X_{Lj},X_{Lk}) 是一个 Tomek Link^[7]。依据上述定义遍历数据集合 X_L ,确定其 Tomek Link 集合,记为 T,从而明确待移除的数据,表达式为:

$$X_L' = \left\{ X_{Lj} \left| \left(X_{Lj}, X_{Lk} \right) \in T \right. \right\} \tag{2}$$

式中: X,表示待移除数据。

则经过欠采样处理后的多数类数据集合表示为:

$$Y_L = X_L - X'_{L, \text{ delete}} \tag{3}$$

式中: Y, 表示处理后的多数类数据集合。

(3)少数类数据集合应用 SMOTE 算法进行过采样处理,以数据集合 X_s ={ X_{S1} , X_{S2} , ···, X_{Sr} , ···, X_{SN} }(N_S 表示集合数据总量)为例,以数据 X_{Sr} 为起始点,找出其 K 个最近邻数据 \hat{X}_{Sr} ,计算起始点与近邻之间的差值 \hat{X}_{Sr} ~ X_{Sr} ,以此为基础,

生成新的少数类数据[8],表达式为:

$$X_S' = X_{Si} + \zeta \times \left(\hat{X}_{Si} - X_{Si}\right) \tag{4}$$

式中: X'_s 表示新的少数类数据: ζ 表示在 $0\sim1$ 范围内的随机常数。

则经过采样处理后的少数类数据集合表示为:

$$Y_S = X_S - X'_{S, \text{ new}} \tag{5}$$

式中:Y。表示处理后的少数类数据集合。

(4) 在完成欠采样和过采样后,检查处理后的不平衡数据集是否达到平衡状态,即多数类数据集合和少数类数据集合内部数据数量是否接近或相等,表达式为:

$$N(Y_{t}) \approx N(Y_{s}) \tag{6}$$

式中: $N(Y_L)$ 表示 Y_L 内部的数据数量; $N(Y_S)$ 表示 Y_S 内部的数据数量。

若满足式 (6) ,则输出不平衡数据处理结果 $\hat{Y} = \{Y_1, Y_2, \dots, Y_n, \dots, Y_N\}$;若不满足式 (6) ,则重复进行第 (2) (3) 步,直至满足式 (6) 为止。

上述过程应用混合采样算法(Tomek Links 算法与SMOTE 算法)完成了不平衡数据的处理,为网络安全检测目标的实现奠定坚实的数据基础。

2 改进 GRU-CNN 模型设计

2.1 改进 GRU-CNN 模型架构

GRU-CNN模型是网络安全检测应用的主要模型之一,而传统GRU-CNN模型采用的是串行结构,CNN和GRU按顺序串联在一起。首先,CNN从输入数据中提取特征,然后将这些特征输入到GRU中进行序列处理。此结构适用于数据具有明显时序特征且需要逐步提取和处理的场景。然而,串行结构可能面临信息损失的问题,因为CNN提取的特征在传递给GRU时可能会丢失一些细节信息,从而降低网络安全检测的精准性。为了改善传统GRU-CNN模型存在的问题,本文方法对GRU-CNN模型进行改进,采用并行结构对CNN和GRU进行连接^[9]。并行结构允许CNN和GRU同时处理输入数据,分别提取网络运行数据的空间特征和时序特征,然后将空间特征和时序特征进行融合以进行最终的网络安全检测。

2.2 数据特征提取与融合

在改进 GRU-CNN 模型中,GRU 负责捕捉网络运行数据中的时序特征,CNN 负责提取网络运行数据中的空间特征,对两者进行融合处理,获取更加全面、精准的网络运行数据特征,为网络安全检测结果的获取提供便利 [10]。

改进 GRU-CNN 模型中 CNN 应用流程为:

Step 1: 输入处理好的不平衡数据 $\hat{Y} = \{Y_1, Y_2, \dots, Y_n, \dots, Y_N\}$,对其进行标准化处理 [11],表达式为:

$$\widetilde{Y}_{i} = \frac{Y_{i} - \min(Y_{i})}{\max(Y_{i}) - \min(Y_{i})}$$
(7)

式中: Y表示标准化处理后的网络运行数据; $min(\cdot)$ 与 $max(\cdot)$ 表示最小值与最大值。

Step 2: 在输入数据Y_i上应用多个卷积核进行卷积操作,通过顺序滑动提取网络运行数据的空间特征。需要注意的是,每个卷积核关注输入数据的不同空间特征,生成相应的空间特征图。卷积操作表示为:

$$\chi = \sum_{m} \sum_{i=1}^{N} \check{Y}_{i} \cdot \alpha (m, n)$$
(8)

式中: χ 表示卷积层输出结果: $\alpha(m,n)$ 表示卷积核函数。其中,m 与 n 代表卷积核的尺寸大小。

Step 3: 对卷积层的输出 χ 进行下采样,以减少空间特征数据维度,同时保留重要特征。池化操作主要是提取空间特征的最大值 χ' ,表达式为:

$$\chi' = \max(\chi) \tag{9}$$

改进 GRU-CNN 模型中 GRU 应用流程为:

Step 1: 将处理好的不平衡数据 $\hat{y} = \{Y_1, Y_2, \dots, Y_n, \dots, Y_N\}$ 按时间序列展开,形成一系列特征向量,作为 GRU 的输入,用于捕捉网络运行数据中的时序依赖关系。

Step 2: GRU 利用更新门和重置门对数据流动进行控制。其中,更新门操作公式为:

$$z_{t} = \sigma \left(W_{z} \cdot \left\lceil h_{t-1}, \hat{Y}_{t} \right\rceil \right) \tag{10}$$

式中: z, 表示更新门输出结果; \hat{Y} ,表示当前时刻输入网络运行数据; W_Z 表示更新门权重矩阵; $\sigma(\cdot)$ 表示 sigmoid 函数; h_{t-1} 表示前一时刻的隐藏状态。

重置门操作公式为:

$$r_{t} = \sigma \left(W_{r} \cdot \left[h_{t-1}, \hat{Y}_{t} \right] \right) \tag{11}$$

式中: r_ι 表示重置门输出结果; W_r 表示重置门权重矩阵。

候选隐藏状态的计算公式为:

$$\tilde{h}_{t} = \tanh\left(W \cdot \left\lceil r_{t} \odot h_{t-1}, \hat{Y}_{t} \right\rceil\right) \tag{12}$$

式中: \tilde{h} 表示候选隐藏状态; W表示候选权重矩阵; \odot 代表元素乘法。

则最终隐藏状态的计算公式为:

$$h_{t} = (1 - z_{t}) \odot h_{t-1} + z_{t} \odot \tilde{h_{t}}$$

$$\tag{13}$$

Step 3: GRU 通过迭代更新隐藏状态,逐步整合时间序列中的信息,从而获得网络运行数据的时序特征,即为 $H = \{h, t = 1, 2, \dots, T\}$,T表示网络安全检测周期[12]。

融合网络运行数据空间特征与时序特征,表达式为:

$$\delta = \omega_1 \chi' + \omega_2 H \tag{14}$$

式中: δ 表示网络运行数据融合特征; ω_1 与 ω_2 表示特征融合

系数。

2.3 网络安全检测结果获取

以 2.2 节融合后网络运行数据特征 δ 为依据 $^{[13-16]}$,计算其与已有网络异常特征集合 λ 的相关系数,表达式为:

$$\rho(\delta,\lambda) = \frac{\delta \cap \lambda}{\delta \cup \lambda} \tag{15}$$

当相关系数 $\rho(\delta,\lambda)$ 大于或等于阈值(需根据网络实际运行情况来制定)时 $[^{17-18]}$,判定当前时刻网络运行状态为异常;当相关系数 $\rho(\delta,\lambda)$ 小于阈值时,判定当前时刻网络运行状态为正常 $[^{19-20]}$ 。

综上所述,在 GRU-CNN 模型的改进及其应用下,完成了网络运行状态(异常或者正常)的判定,实现了研究目标。

3 实验与结果分析

3.1 改进 GRU-CNN 模型训练

本文方法采用改进 GRU-CNN 模型实现网络安全检测,其涉及参数较多,需对其进行提前训练,以此保障改进 GRU-CNN 模型性得到最佳发挥。准备训练数据集与测试数据集,如表1所示。

表1 训练数据集与测试数据集表

组别	训练数据集		测试数据集	
	类别编号	数据量/MB	类别编号	数据量/MB
1	1	2485	1	689
	2	562	2	1021
	3	451	3	1578
	4	1859	4	2103
	5	1523	5	567
	6	2574	6	784
	7	698	7	1245
	8	1475	8	1340
2	1	985	1	1124
	2	1526	2	1578
	3	567	3	965
	4	452	4	2310
	5	1245	5	1974
3	1	2415	1	985
	2	1245	2	1120
	3	634	3	1356
	4	895	4	1021
	5	1020	5	784
	6	754	6	2169
	7	1211	7	1079

为了综合考虑空间特征和时序特征的影响,通过设计一个加权混合损失函数,将 CNN 部分和 GRU 部分的损失按照一定比例组合起来,以此来确定改进 GRU-CNN 模型训练终止条件,表达式为:

$$\min L_{total} = \mu L_{CNN} + (1 - \mu) L_{GRU} \tag{16}$$

式中: L_{total} 表示改进 GRU-CNN 模型的加权混合损失函数; μ 表示权重系数,其主要根据任务需求进行调整; L_{CNN} 与 L_{GRU} 代表 CNN 和 GRU 部分对应的损失函数。

当改进 GRU-CNN 模型输出加权混合损失函数达到最小值时,表明其性能达到最佳状态,确定改进 GRU-CNN 模型内部参数的最佳取值,为后续的实验进行提供一定的便利。

3.2 网络安全检测性能分析

3.2.1 网络安全检测延迟结果分析

使用病毒攻击实验网络,将其攻击开始时间记为 0 min,应用本文方法、方法 1、方法 2 与方法 3 对网络安全进行检测,如图 2 所示。

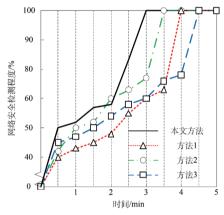


图 2 网络安全检测延迟结果示意图

本文方法网络安全检测延迟为 3 min, 方法 1 网络安全检测延迟为 4 min, 方法 2 网络安全检测延迟为 3.5 min, 方法 3 网络安全检测延迟为 4.5 min。通过比较发现,本文方法网络安全检测延迟时间更短,说明本文方法网络安全检测效率更高。

3.2.2 网络安全误检率结果分析

本文方法、方法 1、方法 2 与方法 3 应用后, 网络安全 误检率结果如图 3 所示。

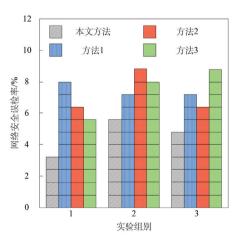


图 3 网络安全误检率结果示意图

在不同实验组别背景下,本文方法应用后网络安全误检 率均低于方法 1、方法 2 与方法 3, 其最小值达到了 3.2%。 因为本文方法应用改进 GRU-CNN 模型, 其融合门控循环单 元(GRU)与卷积神经网络(CNN)的优势。GRU作为一 种轻量级的循环神经网络变体,通过更新门和重置门有效捕 捉序列数据中的长期依赖关系,这对于识别网络中的复杂攻 击模式至关重要。

4 结语

网络安全检测是保护关键信息基础设施免受恶意攻击的 首要防线。随着数字化转型的加速,政府机构、金融机构、 医疗机构及众多关键行业越来越依赖于信息技术系统来支撑 日常运营和服务提供。上述系统存储和处理的往往是高度敏 感和机密的数据,一旦遭受黑客入侵、数据泄露或恶意软件 的感染,不仅会导致服务中断、财产损失,还可能对国家安全、 社会稳定和个人隐私构成严重威胁。因此,通过网络安全检 测及时发现并应对潜在威胁, 是确保这些系统稳定运行和数 据安全的基石。已有方法由于应用模型的自身缺陷, 无法获 得精准的网络安全检测结果, 故提出不平衡数据环境下基于 GRU-CNN 模型的网络安全检测方法研究。

参考文献:

- [1] 李群, 王超, 任天宇. 基于渗透测试的多层次网络安全入 侵检测方法 [J]. 沈阳工业大学学报, 2022, 44(4):372-377.
- [2] 钱爱娟, 樊昕, 董笑菊, 等. 基于社区发现的网络异常检测 方法 [J]. 计算机学报, 2022, 45(4):825-837.
- [3] 汤萌萌,郭渊博,张晗,等.基于提示问答数据增强的小 样本网络安全事件检测方法 [J]. 通信学报, 2024, 45(8):62-74.
- [4] 尹晟霖,张兴兰,左利宇.双重路由深层胶囊网络的入侵 检测系统 [J]. 计算机研究与发展, 2022, 59(2):418-429.
- [5] 张婷婷, 王智强. 基于反向传播算法的网络安全态势感知 仿真 [J]. 计算机仿真, 2024, 41(3):436-440.
- [6] 刘宁,朱波,阴艳超,等.一种混合 CGAN 与 SMOTEENN 的不平衡数据处理方法 [J]. 控制与决策, 2023, 38(9):2614-2621.
- [7] 叶志威,张晓龙,林晓丽.一种面向药物-靶点相互作用 预测的不平衡数据处理方法 [J]. 武汉科技大学学报, 2022, 45(1):68-74.
- [8] 高冰, 顾兆军, 周景贤, 等. 面向 ICS 不平衡数据的重叠 区混合采样方法 [J]. 计算机工程与应用, 2023, 59(19):305-
- [9] 姚芳,汤俊豪,陈盛华,等.基于ISSA-CNN-GRU模型

- 的电动汽车充电负荷预测方法[J]. 电力系统保护与控制, 2023, 51(16):158-167.
- [10] 杨飞,郝晓莉,杨建,等.基于多车型 CNN-GRU 性能 预测模型的轨道状态评价 [J]. 西南交通大学学报, 2023, 58(2):322-331.
- [11] 王力, 李志新, 张亦弛. 基于红外的 SSA-CNN-GRU 电路 板芯片故障诊断 [J]. 激光与红外, 2023, 53(4):556-565.
- [12] 胡瑾, 雷文晔, 卢有琦, 等. 基于 1D CNN-GRU 的日光温 室温度预测模型研究 [J]. 农业机械学报, 2023, 54(8):339-346
- [13] 钱倍奇, 陈谦, 张政伟, 等. 基于异构数据特征级融合的 多任务暂态稳定评估 [J]. 电力系统自动化, 2023, 47(9): 118-128.
- [14] 吴欣玥、廖家仪、张晓荣. 基于多源数据融合的成都市职 住空间特征及影响因素研究 [J]. 规划师, 2023, 39(1):120-127.
- [15] 郝峰,方冰,祁炜雯,等.基于数据融合并行特征提取 的调峰电源设备状态评估方法[J]. 水电能源科学, 2023, 41(5):203-206.
- [16] 王静, 丁卫平, 尹涛, 等. 基于多模态模糊特征融合 的脑龄协同预测算法[J]. 模式识别与人工智能, 2024, 37(7):613-625.
- [17] 黄闽南, 范佳铭, 王一山, 等. 基于时域曲线相关性鉴别 的分布式光纤扰动定位算法研究[J]. 仪表技术与传感器, 2024(2):93-97.
- [18] 姜吉光, 侯小龙, 苏成志, 等. 基于 EMD-PSO 优化阈值 的 COD 光谱去噪研究 [J]. 激光杂志, 2023, 44(7):51-56.
- [19] 樊超阳, 李朝锋, 杨苏辉, 等. CEEMDAN 联合小波阈值 算法在水下激光雷达中抑制散射杂波的应用 [J]. 物理学 报,2023,72(22):100-107.
- [20] 陈晋市, 张淼淼, 王普长, 等. 基于主成分分析阈值选取 的挖掘机主泵载荷谱外推[J]. 吉林大学学报(工学版), 2023, 53(2):355-363.

【作者简介】

熊天运(1999-), 男, 广西田林人, 本科, 研究方向: 计算机方向。

韦富(2000-) 男, 广西来宾人, 本科, 研究方向: 计 算机方向。

(收稿日期: 2024-10-23)