

基于轻量化机器学习算法的多维数据异常检测研究

贾诗音 李敏 王淼

湖北商贸学院 湖北 武汉 430000

【摘要】：随着物联网、边缘计算的快速发展，工业、网络、金融等领域产生海量多维数据，异常检测成为保障系统稳定运行的关键，但传统算法存在计算复杂度高、资源消耗大等问题，难以适配边缘设备等资源受限场景。本文对基于轻量化机器学习算法的多维数据异常检测方法进行分类分析，包括传统轻量化、轻量化深度学习及混合轻量化三类方法并对比总结；通过设计实验，结合典型数据集对比各类算法性能；结合多领域案例验证方法实用性。研究表明，轻量化算法能在平衡检测精度与资源消耗的同时，满足实时检测需求，为资源受限场景的多维数据异常检测提供有效解决方案。

【关键词】：轻量化机器学习；多维数据；异常检测；边缘计算；模型压缩

DOI:10.12417/2705-0998.26.06.106

在数字化转型加速发展的背景下，工业物联网传感器、网络设备、金融交易系统持续产出海量多维数据，此类数据有着高维性、关联性、动态性等特性，其背后所隐藏的异常信息可能引发设备故障、网络入侵、金融欺诈等风险，异常检测成为数据挖掘与智能运维领域的核心研究项目。传统异常检测算法虽能实现一定的检测效果，但多依赖高性能计算设备，参数量大、计算开销高，无法很好地适配边缘设备、嵌入式系统等资源受限场景的实时检测需求。随着 TinyML、模型压缩等技术的兴起，轻量化机器学习算法凭借低资源消耗、高实时性的优势，成为解决上述难题的主要路径。本文聚焦轻量化机器学习跟多维数据异常检测的融合应用，归纳各类检测方法，用实验对性能进行验证，参照实际场景分析应用价值，为相关领域的研究以及工程化部署给予参考与支持。

1 基于轻量化机器学习算法的多维数据异常检测方法分类与分析

1.1 基于传统轻量化机器学习算法的多维数据异常检测

基于传统机器学习的轻量化异常检测，核心是通过简化模型结构、优化计算流程，在维持基础检测能力的基础上，降低资源使用，适配多维数据场景。该类方法在经典传统算法的基础上，采用特征筛选、样本抽样、参数精简等轻量手段进行改进，无需复杂算力支撑，部署成本不高。典型算法有轻量化 K-Means、剪枝优化的决策树、简化样式的 One-Class SVM 及稀疏孤立森林等。其中，稀疏孤立森林经过裁剪冗余节点、筛选关键特征，很大程度降低内存占用，可直接在低功耗边缘设备上部署；轻量化 K-Means 对距离计算的逻辑加以优化，缩减迭代的数量，加大多维数据聚类的效率，可应用于中小规模多

维数据的点异常检测。这类方法原理易懂，训练速度快，对硬件资源的要求极低，但在处理高维复杂数据时，很难捕获特征间的深层关联，检测精度比较有限，易受数据噪声方面的干扰，更适合数据维度不高、异常模式比较简单的场景^[1]。

1.2 基于轻量化深度学习算法的多维数据异常检测

基于轻量化深度学习的异常检测，利用网络剪枝、量化、知识蒸馏等主要技术完成，对深度学习模型进行压缩优化，兼顾多维数据复杂特征提取能力和资源消耗的管理。这类方法借助深度学习特征学习所具备的优势，解决传统轻量化算法精度不佳的问题，符合高维、复杂的多维数据实际场景。主流做法包含轻量化自编码器(AE)及其相关变体、简化版 GRU/LSTM、微型 CNN 等。例如，轻量化的 GE-GRU-VAE 模型结合图嵌入与注意力机制，精简网络层数与参数数量，在使时间和空间复杂度下降的同时，提升多维时间序列异常检测的精度；Tiny-LSTM 利用权重共享、对隐藏层节点进行裁剪，保有时序特征学习的能力，可部署于微控制器等资源极度受限设备。此类型方法检测的精确程度高、泛化能力较好，但模型优化所面临的难度大，部分算法有可解释性不高的问题，而且训练的过程需要一定算力支撑，适配场景更倾向于高维复杂数据实时检测。

1.3 基于混合轻量化算法的多维数据异常检测

基于混合轻量化算法的异常检测，核心是融合传统轻量化以及轻量化深度学习算法的优势，通过模块协同优化，补足单一算法的不足，达成检测精度、效率与资源消耗的均衡，适合复杂多维数据场景。

这类方法通常采用“特征提取+异常识别”的双模块结构，传统轻量化算法负责快速筛选明显异常、降低计算费用，轻量化深度学习算法负责精准识别模糊异常、提高检测精度，二者协作实现高效的检测。常见典型方案有小波特征提取搭配轻量化 SVM、图嵌入与 GRU-VAE 混合模型、稀疏孤立森林和 Tiny-LSTM 组合模型等^[2]。其中，稀疏孤立森林快速把正常样本排除，降低后续的计算量，Tiny-LSTM 对筛选后的可疑样本做精准细致检测，极大提高检测的效率和精确性。这类方法综合性能是最好的，能适应复杂度有差异的多维数据，但模型结构相对复杂，需平衡各模块资源消耗以及协同效率，部署难度稍高于单一轻量化算法，更适用于高要求的复杂检测场景。

1.4 各类方法的对比总结

综合对比上述三类轻量化异常检测方法，三者在检测精度、资源消耗、适用场景上各有侧重，核心差异集中在性能和资源的平衡能力。传统轻量化算法的优势是部署便捷、资源消耗极低，参数量与推理延迟都处在最低的水平，训练成本低，可用于数据维度低、异常模式简单的中小规模场景，但检测精度不高，难以应对复杂多维数据。轻量化深度学习算法检测精度高，泛化能力强，可有效捕获多维数据的深层相关特性，适配高维复杂的数据场景，但模型优化所面临的难度大，部分算法的可解释性差，资源消耗高于传统轻量化的算法。混合轻量化算法兼顾前两类方法的优势，检测精度与轻量化深度学习算法的精度接近，资源消耗在两者中间的范围，可适配复杂多维的数据局面，但模型结构复杂，部署以及优化的成本偏高。在实际使用中，要把场景的资源约束、数据复杂程度以及检测精度需求结合起来，选出合适的检测方式。

2 实验对比与分析

2.1 实验环境与数据集

实验设置“边缘设备+参考服务器”双环境，兼顾资源有限场景和性能对比需求，保证实验结果的实用性和全面性。边缘设备选用 NVIDIA Jetson Nano 与 ARM Cortex-M7，模拟工业、边缘计算等真实部署场景，前者负责中等资源约束下模型的运转，后者适配低功耗、低内存的极端资源场景；参考服务器为 Intel Core i7-12700H、16GB 内存，用于对比轻量化算法和传统非轻量化算法性能上的不同点。软件环境所采用的是 Python 语言，基于 TensorFlow Lite、PyTorch Mobile 框架达成轻量化模型的部署工作，凭借 Scikit-learn、TensorBoard 实现数据处理和性能的监控^[3]。数据集选用三类公开标准数据集，分别和工业、网络安全、金融三大核心应用层面对应：SWaT 和 WADI 工业数据集（含多维传感器数据，标注设备异常）、KDD Cup 99 网络流量数据集（含攻击类型标注）、信用卡欺诈数据集（多维交易特征），所有的数据集都已预处理，去掉无效值以及噪声，保证实验数据的有效性和代表性。

2.2 评估指标

实验运用两类核心评估指标，综合衡量算法的检测性能和轻量化属性，防止单一指标造成的片面性。第一类为检测性能指标，包括精确率、召回率、F1 值以及检测的准确率，其中精确率对异常检测的准确性进行衡量，召回率度量异常识别的完整性，F1 值把二者性能平衡起来，检测准确率传达整体识别的效果，四类指标一起评估算法异常识别能力^[4]。第二类是轻量化指标，包括参数量、浮点运算次数（FLOPs）、推理延迟和内存占用，参数量和 FLOPs 呈现模型的复杂度，推理延迟对实时检测能力进行衡量，内存占用呈现出模型对硬件资源的需求，直接与边缘设备部署适配性相关。此外，补充稳定性方面的指标，凭借添加数据噪声、模拟数据分布的变化，估量算法抗干扰和泛化能力，保障实验结果能对应实际应用中的复杂数据场景，为算法工程化部署给出可靠的参考内容。

2.3 实验对比与分析

实验选取三类方法中的典型算法进行对比，包括传统的轻量化稀疏孤立森林、轻量化深度学习的 GE-GRU-VAE、混合轻量化的稀疏孤立森林与 Tiny-LSTM，同时将传统非轻量化算法引入用作对照。实验所得到的结果显示，传统轻量化算法所消耗的资源是最少的，稀疏孤立森林参数量小于 10MB，推理延迟是每样本小于 1ms，但 F1 值普遍小于 0.8，抗噪声能力较弱；轻量化深度学习算法有着最优的精度，GE-GRU-VAE 模型 F1 值在 0.95 以上，参数量仅是传统深度学习模型参数量的 1.7%，适配边缘设备，但其推理延迟比传统轻量化算法略高；混合轻量化算法的综合性能最优，F1 值接近 0.96，资源消耗处于两类方法资源消耗之间，具有较强的抗干扰能力。与传统非轻量化算法相比，三类轻量化算法参数量、推理延迟均有一个数量级以上的降低，且检测精度损失控制在 5% 以内，充分验证了轻量化算法在资源受限场景中的可行性与优越性，给不同场景下算法的选取提供明确依据，为不同资源约束条件提供最优算法选型方案。

3 应用场景与案例分析

3.1 工业物联网领域

工业物联网领域中，多维数据异常检测是做到设备故障预警、保障生产有序的核心手段，该领域一般采用边缘部署模式，对算法资源消耗和实时性要求十分高，轻量化机器学习算法可良好适配。工业场景中，传感器实时获取设备的温度、振动、电流、压力等多维度的相关数据，异常数据一般预示设备故障，传统算法设置部署于云端，有着响应延迟高、网络传输压力大的问题。某智能制造企业采用轻量化 GE-GRU-VAE 模型，部署到工业边缘网关，对生产流水线设备多维时序数据进行实时异常检测。该模型开展量化、剪枝优化工作后，参数量仅 0.12M，内存使用不到 8MB，样本推理的延迟低于 5ms，较传统云端检

测方案而言,响应速度提升80%以上。实际开展应用的过程中,该模型顺利捕捉到设备早期故障信号,故障漏报率降低了35%,有效减少生产停机损失,验证了轻量化算法在工业物联网场景的实用性与经济性,为智能制造企业降本增效提供了可靠技术路径。

3.2 网络安全领域

网络安全领域的主要需求是实时辨别网络流量中的异常行为(如端口扫描、DDoS攻击、异常访问等),网络流量有着高维、实时、海量的特征,且多会部署在边缘防火墙、网关等资源受约束的设备,轻量化算法是最优的方案。传统网络异常检测算法的复杂度过高,无法很好适配边缘设备实时检测需求,时常出现漏报、误报且响应滞后。某网络安全企业采用轻量化混合模型,部署边缘防火墙,实时对包含网络带宽、连接数、数据包类型的多维流量数据分析。该模型凭借特征筛选和模型压缩,内存占用下降幅度是70%,推理延迟控制在10ms这个范围以内,不依赖于云端的算力。正式部署之后,该模型能快速识别不同的网络异常行为,检测准确率达94%,有效缓解云端数据传输所带来的压力,提高网络安全防护的实时性,满足中小型企业、边缘节点的网络安全防护要求^[5]。

3.3 金融风控领域

在金融风控领域中,多维数据异常检测主要是识别信用卡盗刷、虚假交易等欺诈行为,交易数据含有金额、地点、频率、

设备信息等多维特征,要求算法同时兼顾实时性、高精度和低资源消耗,轻量化算法可满足该场景的需求。传统金融风控算法一般部署在云端,有着交易延迟大、算力成本高等问题,无法很好适配实时交易监控场景。某商业银行采用轻量化决策树剪枝算法,对信用卡交易多维数据做实时的异常监测,模型部署于金融边缘节点当中,经过参数精简以及特征优化后,参数量下降了60%,推理延迟显著降低,可做到交易实时监督^[6]。实际应用中,该模型在让检测准确率达到92%的前提下,能够识别虚假交易以及盗刷等异常行为,使误报率处于3%以内,大幅降低云端的算力成本,给金融机构提供高效、经济的风控解决方案,保障用户资金不受损失,同时降低交易风险,提升用户体验,助力银行数字化转型与智能风控体系建设。

4 结语

综上所述,本文系统地对基于轻量化机器学习算法的多维数据异常检测三类核心方法进行了梳理,经实验对比验证了各算法性能之间的差异,结合工业物联网、网络安全、金融风控三个领域相关案例,证实了轻量化算法在资源受限场景中的实用性和优越性。研究发现,轻量化算法可切实平衡检测精度和资源消耗,解决传统算法适配边缘设备存在的痛点,但目前仍有精度和轻量化平衡不足、工程化部署适配性差等方面的问题。未来可把研究重点放在高精度轻量化模型设计、跨设备适配优化等方向,进一步促使该技术在更多领域实现工程化推广落地,为多维数据异常检测给出更高效、更经济的解决方案。

参考文献:

- [1] 彭朝琴,熊思成,李奇聪,等.基于注意力机制的融合预测改进自编码的EMA多维数据异常检测方法[J].北京航空航天大学学报,1-15.
- [2] 霍纬纲,吴艺凝.基于三重生成对抗的多维时间序列异常检测[J].计算机工程与设计,2025,46(05):1304-1310.
- [3] 蒋融融,顾国民,刘洋,等.基于工业系统多维传感数据的异常检测与诊断[J].浙江工业大学学报,2024,52(06):621-632+665.
- [4] 李少波,王岩,杨磊,等.面向多维数据与随机噪声的无人机飞行数据异常检测方法[J].中国惯性技术学报,2024,32(07):733-742.
- [5] 李泽宇,乔钢柱,张苗苗.基于对偶对抗学习的多维时间序列异常检测[J].中北大学学报(自然科学版),2024,45(02):205-212.
- [6] 宋世军,樊敏.基于谱聚类的多维数据集异常数据检测方法[J].吉林大学学报(工学版),2023,53(10):2917-2922.