

基于多级代理许可区块链的联邦边缘学习模型

葛丽娜^{1,2,3}, 栗海澳^{1,3}, 王捷^{1,2,3}

(1. 广西民族大学人工智能学院, 广西 南宁 530006; 2. 广西混杂计算与集成电路设计分析重点实验室, 广西 南宁 530006;
3. 广西民族大学网络通信工程重点实验室, 广西 南宁 530006)

摘要: 针对零信任边缘计算环境下联邦学习面临的隐私安全及学习效率低等问题, 提出了一种边缘计算中基于多级代理许可区块链的联邦学习模型, 设计多级代理许可区块链构建联邦边缘学习可信底层环境, 实现分层模型聚合方案缓解模型训练压力, 利用秘密共享和差分隐私设计混合策略增强模型隐私。针对边缘客户端可信度为零或极差的问题, 设计了基于信誉验证的联邦任务节点选择算法, 将正向训练样本及本地模型作为信誉奖励, 完善安全验证方案, 进一步保证模型抵御恶意敌手攻击的有效性。实验结果表明, 在 40% 恶意敌手的攻击下, 相较于现有的先进方案, 所提方案准确率提升了 10%, 以较高的模型准确率实现了较高的隐私安全。

关键词: 联邦学习; 区块链; 数据安全; 隐私保护; 边缘计算

中图分类号: TP18; TP309.2

文献标志码: A

DOI: 10.11959/j.issn.1000-436x.2024072

Federated edge learning model based on multi-level proxy permissioned blockchain

GE Li'na^{1,2,3}, LI Haiao^{1,3}, WANG Jie^{1,2,3}

1. School of Artificial Intelligence, Guangxi Minzu University, Nanning 530006, China

2. Guangxi Key Laboratory of Hybrid Computation and IC Design Analysis, Nanning 530006, China

3. Key Laboratory of Network Communication Engineering, Guangxi Minzu University, Nanning 530006, China

Abstract: Aiming at the problems of privacy security and low learning efficiency faced by federated learning in zero trust edge computing environment, a federated learning model based on multi-level proxy permission blockchain for edge computing was proposed. The multi-level proxy permission blockchain was designed to establish a trusted underlying environment for federated edge learning, and the hierarchical model aggregation scheme was implemented to alleviate the pressure of model training. A hybrid strategy was devised to enhance model privacy using secret sharing and differential privacy. A federated task node selection algorithm based on reputation verification was devised to address the problem of zero or extremely poor credibility of edge clients. Positive training samples and the local model were utilized as reputation rewards to refine the security verification scheme, and further ensure the effectiveness of the model against malicious adversaries. Experimental results show that under the attack of 40% malicious adversaries, compared with the existing advanced schemes, the accuracy of the proposed scheme is improved by 10%, and high privacy security is achieved with high model accuracy.

Keywords: federated learning, blockchain, data security, privacy-preserving, edge computing

收稿日期: 2023-09-08; 修回日期: 2023-12-01

基金项目: 国家自然科学基金资助项目(No.61862007); 广西自然科学基金资助项目(No.2020GXNSFBA297103)

Foundation Items: The National Natural Science Foundation of China (No. 61862007), Guangxi Natural Science Foundation (No.2020GXNSFBA297103)

0 引言

近年来,物联网设备逐渐智能化,存储及计算性能不断增强,被广泛用于可穿戴医疗监控辅助诊断^[1]、智能农业监测^[2]、智能工业控制^[3]等多种传感任务。智能物联网的成功实践离不开大数据基础及机器学习(ML, machine learning)技术的支撑。然而,依赖云端中心集中处理的物联网计算范式逐渐无法承受海量数据的压力,同时,用户数据上传至云端中心可能导致严重的数据安全和隐私泄露问题。边缘计算(EC, edge computing)作为新型分布式计算范式,利用靠近终端用户的边缘服务器来处理 and 存储数据,为物联网应用提供及时且智能的服务。尽管EC辅助物联网提供了特殊和增强的服务质量,但在数据安全性和隐私方面仍存在巨大的风险^[4]。

为了有效应对集中式ML面临的挑战,Mcma-han等^[5]提出联邦学习(FL, federated learning)。FL作为隐私保护分布式ML范式,保证用户数据不出本地,以隐私保护的方式协作学习,是解决数据隐私安全和异构知识融合等问题的主流方案,目前已成功应用于医疗诊断^[6]、金融预测^[7]、推荐系统^[8]等领域。针对EC辅助物联网进行ML模型训练时面临的数据安全和隐私泄露威胁,FL能够提供有效的解决方案。然而,传统FL框架的结构特点及不安全的边缘网络环境导致FL在实际的应用场景中仍然面临隐私安全问题。例如,恶意敌手投放有毒梯度来实施安全攻击,或者通过部分模型梯度信息或全局模型反演重建训练集,导致用户隐私数据泄露等。针对FL中的隐私安全问题,不少专家学者开展了相关研究并提出相应的解决方案^[9]。其中,较主流的隐私保护FL方案主要基于密码学技术,例如安全多方计算^[10]和同态加密^[11],以及基于扰动技术,例如差分隐私^[12]。

安全多方计算在不泄露隐私数据的情况下,允许多个参与方对私有数据进行联合计算,拥有严格的密码学理论证明,计算准确度高。同态加密允许对密文直接进行运算,FL中可以委托第三方对数据进行处理而不泄露隐私。文献[13]设计了一种垂直FL框架,结合用于多方协作建模的环形架构简化通信,防御半诚实攻击。但该方案未充分考虑数据投毒攻击对模型准确率的影响。文献[14]提出一种隐私增强FL框架,使用同态加密有效检测FL中的

投毒行为,增强模型鲁棒性。文献[15]将分布式选择随机梯度下降法与Paillier密码系统相结合,提供了数据保密性。以上方案增强了FL安全性,但面对大规模EC环境中的参与节点可信度为0或极差的情况,未能很好地解决海量数据加密处理的通信资源消耗大、系统效率低以及缺乏参与方身份验证的问题。

基于加密技术的FL产生的加解密运算为整个FL系统带来了巨大的额外开销。对于FL系统而言,设计一个结构简单且轻量级的算法极其重要。差分隐私由于拥有严格的数学理论证明和较小的系统开销常被用于设计高效隐私增强的FL。文献[16]为了保护FL的隐私并保留调整学习性能的能力,提出一种具有时变噪声幅度的差分隐私扰动机制。文献[17]开发了一种混合量子经典模型,保证差分隐私在不显著降低模型准确率的同时保护用户敏感信息。然而以上方案在提升FL模型效用和安全性能的同时,未能很好考虑FL训练效率的问题。文献[18]深入研究了差分隐私预算和本地模型训练的梯度参数之间的关系,设计了2种不同的动态隐私预算分配方式,有效提升了模型的训练效率。考虑在大规模EC环境下的基于差分隐私的FL通常会花费额外的通信资源来提供必要的控制信息,且好奇的聚合服务器和不可信的FL参与方都是潜在的安全隐患,基于差分隐私的FL在模型效用、学习效率及隐私安全方面仍存在诸多威胁。

综上所述,目前的隐私保护FL技术被广泛研究和创新。然而,面对大规模EC环境下的参与节点可信度为0或很差的情况仍存在不足。为解决目前的集中式FL面临的中央服务器单一节点信任、学习效率低、恶意敌手投毒攻击及梯度参数泄露隐私数据的问题,本文提出一种EC环境下的基于多级代理许可区块链的联邦边缘学习模型(MAPBFL, federated edge learning model based on multi-level agent permissioned blockchain)。本文的主要贡献如下。

1) 提出基于多级代理许可区块链的联邦边缘学习模型,设计多级代理许可区块链构建联邦边缘学习可信底层环境,实现分层模型聚合方案以缓解大量模型聚合的压力,引入秘密共享和差分隐私设计混合策略进一步增强分层模型聚合过程的隐私性。

2) 提出基于信誉验证的联邦任务节点选择算法, 引入分层模型验证机制, 设计新型信誉积分奖惩措施, 提高模型对抗恶意敌手攻击的有效性。

3) 对所提模型进行实验评估, 在40%恶意敌手的投毒攻击下, 本文方案相较于对比方案的模型准确率提升了10%, 结果表明本文方案可以在实现较高模型可用性的同时提供更高级别的隐私安全。

1 相关工作

本文针对EC环境下的FL面临的隐私安全和学习效率等挑战展开研究, 总结分析了目前的相关工作进展及本文提出的MAPBFL的研究动机。MAPBFL采用多级代理许可区块链构建可信训练环境, 设计分层模型聚合方案提升模型训练效率, 利用混合隐私策略增强模型隐私性。本文基于信誉验证提出FL客户端选择算法, 提升了MAPBFL识别恶意梯度的有效性。

1.1 联邦边缘学习隐私安全

EC环境下的FL系统可以利用边缘服务器的计算能力和数据进行模型训练而不需要数据离开本地, 保证了本地用户隐私数据的安全。但不可信的网络环境导致本地用户数据仍然面临隐私安全威胁。考虑目前的FL大多采用集中式框架, 对中央服务器的单点依赖问题仍然存在。因此, 研究人员将区块链和FL相结合, 利用区块链的去中心化特性减少单一实体的信任威胁, 利用其身份验证和访问机制构建安全可信的FL环境。本文模型设计了多级代理许可区块链结构用于构建联邦边缘学习可信底层环境, 保证用户数据隐私性的同时实现模型安全聚合。然而, 允许所有设备参与FL过程不是长久可行的方案。边缘客户端在设备性能、数据质量等方面的异构性及可靠性等问题都会影响FL模型训练。因此, 本文提出了一种基于信誉验证的联邦任务节点选择算法, 设计新型奖惩机制, 进一步提升本文方案的可靠性。

1.2 基于多级代理许可区块链的联邦边缘学习

FL和区块链的结合可以将智能EC网络转变为去中心化、隐私安全增强的系统。文献[19]设计了基于区块链的隐私保护拜占庭鲁棒FL方案, 使用区块链促进模型透明的流程和法规的实施, 采用余弦相似度和同态加密保证模型安全聚合。然而该方案在处理海量用户数据时会带来不小的计算开销。

同时, 该方案对于资源异构的客户端设备缺乏一定的资源调整优化策略和参与协作训练的激励机制。文献[20]构建了一种激励感知区块链辅助协作机制, 增强了边缘节点在安全保障下参与协作的意愿, 并联合优化卸载和缓存决策及计算和通信资源分配方案最小化边缘节点完成FL任务的总成本。但是该方案并未考虑FL本地参与方的可信验证。可信度未知的边缘客户端的随时参与和退出对FL系统的隐私安全性提出了更高的要求。文献[21]提出一种基于区块链的可验证模型用于保护FL, 该模型结合可信执行平台来保护客户端本地模型训练, 并采用多重签名驱动的全局模型验证来保证模型的可验证性。然而, 该方案并未考虑潜在的诚实但好奇的FL任务节点对训练数据隐私性的威胁。

为了有效提升零信任EC环境下的FL系统的隐私安全性和学习效率, 本文方案利用多级代理许可区块链构建联邦边缘学习可信底层环境, 解决了中央服务器单一节点的信任威胁。然后, 利用分层模型聚合方案来验证低质量的模型更新, 提升FL训练效率的同时增强模型对抗恶意梯度的鲁棒性。最后, 设计基于秘密共享和差分隐私的混合策略保证共享模型的隐私性。

1.3 基于信誉验证的联邦任务节点选择

公平的客户端选择算法可以使用身份验证、信誉评估等技术阻止恶意客户端参与训练, 减少恶意攻击对模型的危害^[22]。因此, Xu等^[23]提出自适应客户端选择策略减轻恶意敌手对模型性能造成的影响。然而该方案并未考虑海量异构数据对FL稳健性的影响。为此, Pene等^[24]开发了一种异构客户端选择的激励设计方案, 利用合作博弈论和基于客户异质性水平的动态聚类方法来克服客户端选择和数据异质性对FL稳健性的影响。同时, 针对如何设计FL客户端选择算法来提升FL通信效率的问题, Yang等^[25]提出了一种基于全局后验和局部倾斜分布之间的核化斯坦因差异的客户端选择方案, 有效减少了通信开销。然而, 在大规模、零信任的EC环境下提升FL系统的数据安全性和隐私性涉及多方面因素。同时, 由于目前的网络攻击逐渐多元化、复杂化, 如何在兼顾系统开销和效率的同时开发更高隐私安全性能的FL客户端选择算法仍然存在挑战。

为了进一步增强所提方案的安全性, 本文通过

新型奖励机制设计了基于信誉验证的联邦任务节点选择算法，通过分层模型验证机制设计和新型信誉积分奖惩措施来提升本文方案对抗恶意敌手破坏模型可用性攻击的有效性。

2 系统模型和问题描述

为解决零信任 EC 环境下的集中式 FL 框架训练海量数据时所面临的设备异构、有限算力、单点故障、多角度的数据安全攻击和较差隐私性能等问题，本文结合区块链和 EC 设计了隐私安全增强的联邦边缘学习模型，分析了系统在安全假设下所面临的威胁模型，并给出本文方案的设计目标。

2.1 模型总体架构

如图 1 所示，系统模型自下而上由边缘设备接入层、边缘服务层、代理区块链层、核心区块链层组成，包括智能终端设备、高性能终端设备、边缘服务器、代理区块链和核心区块链等实体。其中，系统模型各层设计及相关实体功能介绍如下。

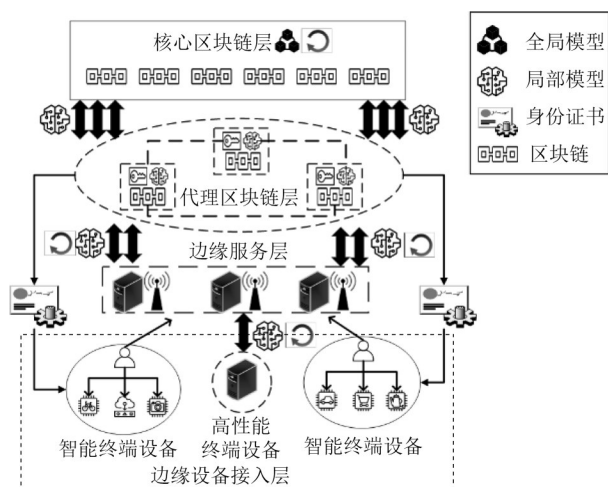


图 1 系统模型

1) 边缘设备接入层。该层主要功能是收集边缘终端数据，由各类物联网智能终端设备组成，包括智能手机、手环、摄像头、全球定位系统 (GPS, global positioning system) 及高性能智能终端等，同时负责区域边缘服务器的选举。

2) 边缘服务层。该层主要功能是协助边缘终端设备训练本地模型，由边缘服务器和部分高性能边缘终端设备组成，同时负责加密和传输数据处理结果。

3) 代理区块链层。该层主要功能是实现局部

模型安全聚合，由各个信任域的代理区块链组成。各个代理区块链负责所属区域的设备注册和认证，控制访问权限，同时维护公共参数、验证信息等关键数据。

4) 核心区块链层。该层主要功能是实现全局模型安全聚合，负责全域加密数据解密、正确性证明检验、共识机制，管理全域关键参数和系统密钥，同时负责与代理区块链协同构建不同层级实体之间的安全数据共享及全域数据整合等。

2.2 安全假设

假设 1 本文假设参与 FL 训练的边缘客户端可信度未知，即边缘客户端在向区块链申请注册之前没有经过任何安全审查。

假设 2 本文假设参与 FL 训练的边缘服务器协助边缘客户端完成本地模型训练是诚实可靠的。

假设 3 本文假设 FL 模型质量验证方是半诚实的，虽然遵循 FL 训练规则，但其好奇其余协作方的敏感数据。

2.3 威胁模型

在零信任的 EC 环境下，系统模型可能面临的数据安全攻击和隐私泄露威胁如下。

中央服务器单一节点信任威胁。中央服务器处理大规模数据可能导致系统效率下降，一旦发生故障，整个 FL 过程就会失败。

数据安全攻击。恶意敌手可能投放精心设计的中毒样本或上传虚假模型破坏 FL 全局模型的可用性和完整性。

隐私泄露威胁。恶意攻击者可能通过部分模型梯度信息或全局模型反演重建训练集，导致用户隐私数据泄露。

2.4 设计目标

本文方案的设计目标是在抵御恶意敌手攻击的同时，以较高隐私性能保护 FL 过程安全。接下来将讨论 MAPBFL 应对 2.3 节所述威胁模型时能够实现的安全目标。

1) 去中心化信任。多级代理区块链取代了集中式 FL 的中央服务器，有效缓解了大规模本地模型的处理压力和单点故障威胁。

2) 模型安全性。代理区块链执行信誉验证算法，确保训练中的恶意敌手被准确识别，同时使用验证模型和全局模型分别对本地模型和局部模型进行正确性检验。

3) 数据隐私性。边缘服务器采用秘密共享对本地模型进行加密,代理区块链向聚合后的局部模型中添加差分隐私噪声保证聚合过程的隐私性。

3 联邦边缘学习模型

3.1 多级代理许可区块链

为了实现安全高效的FL模型训练过程,本文模型引入代理层将单一类型的许可区块链设计为多级结构,主要分为代理区块链层和核心区块链层,并以多级代理许可区块链作为FL中央服务器。具体设计如下。

代理区块链层主要负责设备的注册认证、局部模型聚合及全局模型下发等任务。代理区块链层包括边缘客户端设备、边缘服务器、认证服务器以及边缘验证器等。所有边缘客户端设备和边缘服务器分别向所属区域的认证服务器申请注册认证。注册成功的设备由认证服务器进行管理,并指定活动权限,例如本地模型训练。边缘验证器则负责所属区域的本地模型安全聚合。代理区块链层采用分布式公证人机制与核心区块链层进行数据交互,传输与FL训练相关的学习参数等数据。

核心区块链层主要负责本地模型验证、记录共享数据信息和全局模型聚合等任务。核心区块链包括任务发布者、在线验证设备和云端服务器等。任务发布者将FL模型计算任务发布到区块链网络,边缘客户端设备等数据拥有者可以监听区块链网络发布的任务事件进行任务注册,或者选择继续监听下一个计算任务。在线验证器主要负责验证代理区块链层上传的本地模型参数的正确性,并由云端服务器对正确的本地模型参数进行聚合优化。

3.2 基于多级代理许可区块链的联邦边缘学习模型训练

MAPBFL的训练流程如图2所示。

1) 本地模型训练。首先,代理区块链完成边缘设备认证注册,并颁发许可证书。其次,代理区块链根据设备信誉积分选取本次FL任务参与方。最后,边缘终端设备 k 委托所属区域的边缘服务器使用本地数据集通过随机梯度下降训练本地模型,并通过多次训练迭代保证较优的模型性能。其中,利用随机梯度下降算法更新本地模型梯度为

$$L_t^k = L_{t-1}^k - \eta \nabla F_k(L_{t-1}^k, D_s) \quad (1)$$

其中, L_t^k 为 k 在第 t 轮训练时的本地模型; D_s 为 k 的

S_k 批次的训练数据量; $F_k(L_{t-1}^k, D_s)$ 为局部目标函数,通常定义为 D_s 的经验风险; $\nabla F_k(L_{t-1}^k, D_s)$ 为 $F_k(\cdot)$ 在 $t-1$ 轮训练时的梯度; η 为学习率。

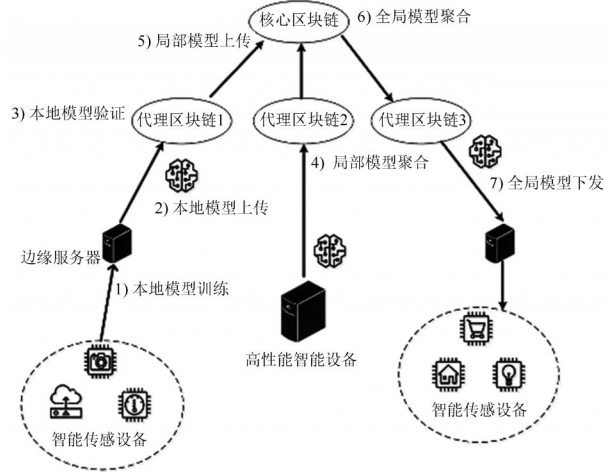


图2 MAPBFL的训练流程

2) 本地模型上传。边缘服务器使用加法秘密共享和认证加密保证 L_t^k 在代理区块链聚合过程中的隐私性,通过加法秘密共享对 L_t^k 进行加密得到 \tilde{L}_t^k ,并以事务 $(\tilde{L}_t^k, \text{Com}(\tilde{L}_t^k))$ 的形式上传至代理区块链。其中, $\text{Com}(\tilde{L}_t^k)$ 表示边缘服务器对 \tilde{L}_t^k 的承诺。边缘服务器使用私钥 key_k 对 \tilde{L}_t^k 进行签名,用于验证事务的真实性和完整性,确保事务不被恶意篡改。

边缘服务器使用加法秘密共享对 L_t^k 进行加密,表示为

$$\begin{aligned} & \{x_i \text{random}[0, p]\}_{i=1,2,\dots,u-1} \\ & \tilde{L}_{t,i}^k = L_t^k + \sum_{i=1}^{u-1} x_i \text{mod } p \\ & \tilde{L}_t^k = \{\tilde{L}_{t,1}^k, \tilde{L}_{t,2}^k, \dots, \tilde{L}_{t,u}^k\} \end{aligned} \quad (2)$$

其中, p 为大素数; x_i 为系统选取的随机数,用于秘密份额计算; $\tilde{L}_{t,i}^k$ 为 L_t^k 经过秘密分割后的第 i 个秘密份额; u 为分割后的秘密片段总数,等于代理区块链中的验证方数量。边缘服务器作为秘密分发者通过安全信道将 $(x_i, \tilde{L}_{t,i}^k)$ 分别发送给代理区块链中的验证方。

边缘服务器使用本地私钥 key_k 对 \tilde{L}_t^k 进行签名,表示为

$$(\tilde{L}_t^k, E_k) = \text{hash}_{\text{key}_k}(\tilde{L}_t^k) \quad (3)$$

其中,使用 $\text{hash}_{\text{key}_k}(\tilde{L}_t^k)$ 对 \tilde{L}_t^k 进行签名,用于通信

过程中验证 \widetilde{L}_t^k 的完整性和正确性。

3) 本地模型验证。代理区块链对接收的 \widetilde{L}_t^k 进行局部聚合之前要经过边缘验证器 $v \in V$ 的合法验证, V 表示边缘验证器集合。首先, v 收到边缘服务器上传的事务 $(\widetilde{L}_t^k, \text{Com}(\widetilde{L}_t^k))$, 检查事务的正确性证明并进行解密, 其中, 至少需要 $\frac{2}{3}$ 的验证方提供其私有秘密片段协同解密。其次, v 将解密后的事务广播至在线验证器进行评估, 超过 $\frac{2}{3}$ 的在线验证器投票为真, 相关学习设备将授予指定的奖励 R_t^k ; 否则将被识别为有毒模型, 并扣除相关学习设备的现有信誉积分, 一旦设备信誉积分低于预定阈值, 将被禁止参与下一轮FL训练。

代理区块链对超过 $\frac{2}{3}$ 的参与方提供的秘密份额进行秘密重构来恢复 L_t^k , 表示为

$$L_t^k = \sum_{i=1}^{\frac{2u}{3}} \widetilde{L}_{t,i}^k \prod_{j \neq i, j=1}^{\frac{2u}{3}} \frac{x_j}{x_j - x_i} \bmod p \quad (4)$$

4) 局部模型聚合。代理区块链将通过验证的 L_t^k 进行局部聚合, 并利用差分隐私向聚合后的局部模型添加高斯噪声来保护模型梯度隐私。

在局部模型更新中添加高斯噪声扰动可表示为

$$\begin{aligned} \widetilde{L}_t^w &= L_t^w + N(\sigma^2) \\ \sigma^2 &= \frac{2\Delta f^2 \log\left(\frac{1.25}{\delta}\right)}{\varepsilon^2} \end{aligned} \quad (5)$$

其中, L_t^w 表示第 t 轮训练中代理区块链执行局部聚合后得到的模型; $N(\sigma^2)$ 表示服从均值为0且方差为 σ^2 的高斯分布抽样; ε 表示隐私预算, 隐私预算越小, 噪声越大, 隐私保护效果越好, 但模型精度越差; δ 是松弛项, 表示违反严格差分隐私的最大容忍概率; \widetilde{L}_t^w 表示加入高斯噪声之后的局部模型更新; Δf 表示灵敏度值, 其定义如下

$$\Delta f = \max_{D, D'} \|f(D) - f(D')\|_2 \quad (6)$$

其中, 灵敏度值和查询函数 $f(\cdot)$ 有关, 表示 $f(D)$ 和 $f(D')$ 之间最大的曼哈顿距离。

5) 局部模型上传。代理区块链将 \widetilde{L}_t^w 上传至核心区区块链, 并生成有效的数字签名。

6) 全局模型聚合。核心区区块链检验 \widetilde{L}_t^w 的正确性, 并使用联邦平均算法进行全局聚合, 表示为

$$G_t \leftarrow G_{t-1} + \frac{1}{n} \sum_{w=1}^n \widetilde{L}_t^w \quad (7)$$

其中, n 为代理区块链上传的局部模型总数, G_{t-1} 为第 $t-1$ 轮训练中的全局模型。

7) 全局模型下发。核心区区块链判断聚合后的全局模型是否达到预设的收敛精度或最大迭代次数, 如果达到, 则终止本次FL任务; 否则开始新一轮的FL训练。

3.3 基于信誉验证的联邦任务节点选择算法

为了解决MAPBFL中的分层模型验证及设备信誉积分奖惩问题, 提升模型识别恶意敌手的准确率。本文基于权益证明机制(PoS, proof of stake)算法思想设计了基于信誉验证的联邦任务节点选择算法, 根据不同设备的学习任务进行评估, 决定进行奖励或惩罚。FL任务节点累计奖励越多, 则其信誉权重越大, 下一轮被选中参与FL任务的可能性越大。具体算法流程如算法1所示。

算法1 基于信誉计算的可靠高效联邦任务节点选择算法

输入 初始网络状态和任务信息, 初始全局模型 G_0 和局部模型 L_0 , 相关设备集合 E , 初始联邦任务节点选择序列 $\text{Ep} \subset E$, 边缘终端设备 $k \in K$ 且 $K \subset \text{Ep}$, 边缘验证器 $v \in V$ 且 $V \subset \text{Ep}$ 和在线验证器 $m \in M$ 且 $M \subset \text{Ep}$, 事务 $(L_t^k, \text{Com}(L_t^k))$ 和 $(L_t^v, \text{Com}(L_t^v))$, 数据集 D_s 及单位奖励 r

输出 联邦任务节点选择序列 Ep

- 1) 随机初始化联邦节点选择序列 Ep ;
- 2) for $i = 1$ to Ep do
- 3) for $j = 1$ to K do
- 4) if (超过 $\frac{2}{3}$ 的验证方验证合格)
- 5) $r_t^k = \text{epoch}_t^k \text{sample}_k r$;
- 6) $R_t^k = R_{t-1}^k + r_t^k$;
- 7) else
- 8) $r_t^k = -\text{epoch}_t^k \text{sample}_k r$;
- 9) $R_t^k = R_{t-1}^k + r_t^k$;
- 10) end if
- 11) 将 k 的状态动作和信誉总积分 R_t^k 存储在代理区块链中;
- 12) end for j
- 13) for $m = 1$ to M do
- 14) $r_t^{m_Verify} = |(L_t^k, \text{Com}(L_t^v))| r$;
- 15) $r_t^{m_Vote(L_t^k)} = |m_Vote(L_t^k)| r$;
- 16) $r_t^{m_Vote} = (|r_t^{m_Verify}| + |r_t^{m_Vote(L_t^k)}|)$;

- 17) $R_t^{m_Vote} = R_{t-1}^{m_Vote} + r_t^{m_Vote}$;
- 18) 存储 m 当前的状态动作和信誉总分 $R_t^{m_Vote}$;
- 19) end for m
- 20) for $r = 1$ to V do
- 21) $r_t^{v_verify} = |(L_t^k, \text{Com}(L_t^k))|r$;
- 22) $R_t^{v_verify} = R_{t-1}^{v_verify} + r_t^{v_verify}$;
- 23) 存储 v 当前的状态动作和信誉总分 $R_t^{v_verify}$;
- 24) end for r
- 25) 更新当前网络和各设备状态信息;
- 26) 更新 FL 网络所有设备节点状态, 代理区块链为新注册的边缘节点设备发放初始信誉总分;
- 27) 代理区块链根据设备信誉积分重新为各个设备节点分配任务角色;
- 28) end for i
- 29) return 联邦任务节点选择序列 Ep

下面将详细介绍 FL 各任务节点信誉积分的计算。

1) 边缘终端设备提供本地数据奖励。首轮通信中随机挑选 v 对边缘服务器使用 k 的本地数据训练得到的 L_t^k 进行验证。

首先, v 将 L_t^k 广播至所有 m 进行投票, 其中, 赞成票记为 $P_L_t^k$, 反对票记为 $N_L_t^k$ 。然后, v 收集所有 m 发回的投票结果, 并给出最终的投票意见, 若赞成票数超过总票数的 $\frac{2}{3}$, 则为该轮中提供高质量本地数据的 k 发放相应的信誉积分奖励 r_t^k , r_t^k 主要与样本数据 se_k 和训练批次 epo_t^k 有关。若小于总票数的 $\frac{2}{3}$, 则扣除 k 的信誉积分。

k 提供本地训练集的信誉积分如式(8)所示, 信誉总分如式(9)所示。

$$r_t^k = \begin{cases} epo_t^k se_k r, P_L_t^k > \frac{2}{3} (P_L_t^k + N_L_t^k) \\ -epo_t^k se_k r, P_L_t^k \leq \frac{2}{3} (P_L_t^k + N_L_t^k) \end{cases} \quad (8)$$

$$R_t^k = \begin{cases} R_{t-1}^k + r_t^k, P_L_t^k > \frac{2}{3} (P_L_t^k + N_L_t^k) \\ R_{t-1}^k + r_t^k, P_L_t^k \leq \frac{2}{3} (P_L_t^k + N_L_t^k) \end{cases} \quad (9)$$

$$\frac{1}{n} \sum_{i=1}^n R_{i-1}^k, 1 \leq i \leq n, \text{且 } i \in \mathbb{Z}_+$$

其中, r 表示参与任务的基础奖励, 即单位奖励,

$epo_t^k se_k r$ 表示对本轮训练任务提供优质本地数据的 k 的信誉奖励, $-epo_t^k se_k r$ 表示对本轮训练任务提供低质量本地数据的 k 的信誉惩罚。对于新注册的任务参与者, 取当前 FL 所有任务参与者的信誉积分平均值 $\frac{1}{n} \sum_{i=1}^n R_{i-1}^k$ 作为其初始信誉积分。一轮 FL 任务结束之后, 如果 k 当前的信誉积分小于最小信誉阈值, 则无法参与下一轮 FL 训练。

2) 边缘验证器验证本地模型奖励。 v 验证事务 $(L_t^k, \text{Com}(L_t^k))$ 的数字签名, 并获得相应的验证奖励 $r_t^{v_verify}$, 具体奖励为

$$r_t^{v_verify} = |(L_t^k, \text{Com}(L_t^k))|r \quad (10)$$

v 当前的信誉总分 $R_t^{v_verify}$ 为上一轮训练的信誉积分 $R_{t-1}^{v_verify}$ 加上本轮训练所获得的信誉积分 $r_t^{v_verify}$, 计算式为

$$R_t^{v_verify} = R_{t-1}^{v_verify} + r_t^{v_verify} \quad (11)$$

3) 在线验证器投票奖励。 $m \in M$ 对 v 广播的验证事务 $(L_t^k, \text{Com}(L_t^k))$ 进行正确性验证, M 为在线验证器集合, 其中 $\text{Com}(L_t^k)$ 表示 v 对接收到的梯度正确性所做出的承诺。 m 对从事务 $(L_t^k, \text{Com}(L_t^k))$ 的签名验证中提取的 L_t^k 进行投票, 将投票结果进行签名并发回给 v , 即发送回复事件 $m_Vote(L_t^k)$ 。 v 检验 $m_Vote(L_t^k)$ 的数字签名, 并给出投票结果。

若 m 的投票意见与 v 给出的投票结果相同, 则给予 m 相应的信誉积分奖励。未通过签名验证的 $m_Vote(L_t^k)$ 将被视为无效投票。具体来说, 在第 t 轮训练中, m 验证模型所获得总奖励如式(12)所示, 当前信誉总分如式(13)所示。

$$r_t^{m_Vote} = (|r_t^{m_Verify}| + |r_t^{m_Vote(L_t^k)}|)$$

$$r_t^{m_Verify} = |(L_t^k, \text{Com}(L_t^k))|r$$

$$r_t^{m_Vote(L_t^k)} = |m_Vote(L_t^k)|r \quad (12)$$

$$R_t^{m_Vote} = R_{t-1}^{m_Vote} + r_t^{m_Vote} \quad (13)$$

其中, $|r_t^{m_Verify}|$ 表示第 t 轮训练中 m 正确验证 $(L_t^k, \text{Com}(L_t^k))$ 的奖励, $|r_t^{m_Vote(L_t^k)}|$ 表示第 t 轮训练中 m 对 L_t^k 进行正确投票的奖励。对于未通过签名验证的事务, v 不会将其封装的模型参数发送给其他 m 进行投票, 阻止了恶意敌手上传虚假模型破坏全局模型精度。

4 理论分析

本节从理论上证明了本文所提模型的隐私性和安全性。

4.1 隐私性分析

为了保证分层聚合过程中模型梯度的隐私性，MAPBFL 引入秘密共享和差分隐私设计了混合安全策略，本节将在半诚实敌手威胁模型中进一步证明混合安全策略的隐私性。

定义 1 秘密恢复阈值。若少于 $\frac{2}{3}$ 的秘密片段拥有方协同合作则不能得到关于秘密 S 的任何信息，任意超过 $\frac{2}{3}$ 的秘密片段拥有方协同合作可以根据式(4)恢复 S 。

定义 2 组合性^[26]。若 $M = \{M_1, M_2, \dots, M_n\}$ 是由一组满足 ϵ_i 的差分隐私算法 $M_i (1 \leq i \leq n)$ 组成，对于数据集 D ，如果 $M_i(D)$ 满足 ϵ_i 差分隐私，则所有算法构成的集合 M 在 D 上满足 $\sum_i \epsilon_i$ 差分隐私。

定理 1 在不超过 $\frac{2}{3}$ 的验证方合谋的情况下，边缘服务器采用秘密共享加密本地模型可以抵御半诚实敌手窃取诚实客户端真实梯度的攻击。

证明 假设在第 t 轮训练中，验证方是半诚实敌手。边缘服务器将加密后的本地模型参数分成 n 份，分别发送给 n 个 v_i ，由于 v_i 只知道第 i 份秘密值，根据定义 1，少于 $\frac{2}{3}$ 的验证方无法单独恢复完整秘密值。当其与剩余 $\left\lfloor \frac{2}{3}(n-1) \right\rfloor$ 的 v_i 合谋时，可获得合谋后的秘密份额集合 $\{D_i^{m,i}\}_{i=1}^{\lfloor \frac{2}{3}n \rfloor}$ ，根据定义 1 设定的最低参数恢复阈值，当且仅当超过 $\frac{2}{3}$ 的 v_i 相互合谋才能恢复秘密值，则上半诚实敌手数量配置无法获取边缘服务器的真实本地模型梯度。

证毕。

定理 2 在半诚实敌手威胁模型下，即使负责全局模型聚合的计算节点是半诚实的，代理区块链层的本地差分隐私方案设计仍然可以保证全局模型训练过程的隐私性。

证明 假设在第 t 轮训练中，负责全局模型聚合的计算节点为半诚实敌手。代理区块链对局部模型参数采用高斯机制加噪，如式(5)所示，每轮模型训练均满足 ϵ -差分隐私，其中，总的隐私预算约束满足 $\epsilon = \sum_i \epsilon_i$ ，则第 t 轮训练中的隐私预算 $\epsilon_t = \frac{\epsilon}{t}$ 。

根据本地差分隐私方案，每轮模型迭代保证所

有局部模型加入高斯噪声，满足全局模型聚合过程的隐私性。根据定义 2 可知，模型训练的总隐私预算 ϵ 是每轮训练中的隐私预算之和，其中， $\epsilon = \epsilon_1 + \epsilon_2 + \dots + \epsilon_n$ ， n 为模型训练总轮数，因此全局模型聚合过程满足 ϵ -差分隐私，负责全局模型聚合的节点只能获得经过差分隐私加噪后的局部模型，无法从局部模型中获取原始梯度信息，从而保证了全局模型聚合过程的隐私性。

证毕。

4.2 正确性分析

本节主要对多级代理区块链能否验证本地模型和局部模型的正确性进行分析。

定理 3 双层模型验证机制使用验证方的验证模型检验本地模型的正确性，使用上一轮的全局模型检验局部模型的正确性。

证明 首先分析代理区块链层使用验证方的验证模型检验边缘服务器本地模型的正确性。理想情况下，代理区块链为了验证边缘服务器本地模型 L_t^k 的正确性，使用测试集对 L_t^k 的精度 $P(L_t^k)$ 和全局模型 G_{t-1}^k 的精度 $P(G_t^k)$ 进行评估和比较，假设 L_t^k 遭受恶意敌手的投毒攻击，则相较于 $P(G_t^k)$ ， $P(L_t^k)$ 严重下降。假设边缘服务器是诚实的且未遭受恶意敌手攻击，则相较于 $P(G_t^k)$ ，本轮训练中的本地模型 L_t^k 可能取得近似 $P(G_t^k)$ 的模型精度。

然而，验证方对于边缘服务器上传的 $\{L_t^k, G_{t-1}^k\}$ 是不能完全信任的，且无法获取边缘服务器本地数据集，为了保证代理区块链准确验证本地模型的正确性，MAPBFL 要求验证方在每次验证开始之前先执行一轮合法局部训练，将获得的验证模型 L_t^v 用于本轮边缘服务器上传的 L_t^k 的正确性验证。验证方收到 L_t^k 后，计算 L_t^k 和 L_t^v 的精度差值，并与预先设定的精度阈值 S 相比较，若 $|P(L_t^k) - P(L_t^v)| \leq S$ ，则验证方判定该本地模型是正确的，否则判定该本地模型为恶意模型。

为了验证代理区块链上传的局部模型的正确性，MAPBFL 使用上一轮全局模型作为验证模型，计算代理区块链上传的局部模型精度和全局模型精度的差值，并与预先设定的精度阈值 S 进行比较，若精度差值小于 S ，核心区块链则判定相关局部模型是诚实的，并允许加入全局模型聚合，否则判定为恶意模型，确保全局模型聚合的正确性。

证毕。

5 实验仿真及性能分析

本节首先介绍了MAPBFL的仿真实验环境,并从模型准确率和时间开销两方面对MAPBFL的性能进行评估,最后分析了基于信誉验证的联邦任务节点选择算法的效果。

5.1 实验设置

1) 实验环境

本文所有实验的服务器环境配置如下: Intel (R) Core(TM) i7-7700 CPU, 3.60 GHz, 32 GB 内存, Window10 操作系统。实验使用 Pycharm 搭建 Python 开发环境,使用 Python 编程语言实现区块链功能模块,使用 Python3.7 中的 Pytorch2.0.1 搭建 FL 框架,采用卷积神经网络作为训练模型,其中每个模块中的卷积核大小、神经元个数等超参数需要根据实验效果来确定。

2) 实验数据集

本文分别基于 MNIST 和 CIFAR-10 数据集对文中所提的 MAPBFL 进行性能评估。为了构造真实的 FL 分布式训练场景,所有实验均设置了 20 个边缘设备用于 MAPBFL,并将整个 MNIST 和 CIFAR-10 训练集随机划分为大小相等且无样本重叠的 20 个子集,每个子集都随机分发给边缘设备训练本地模型。

3) 实验参数

实验设置 FL 全局迭代次数固定为 100 轮,每轮迭代包括 1 次全局聚合和 5 次本地训练,本地训练学习率为 0.005,本地数据批次大小为 10,单位奖励 r 的默认值为 1,最小信誉积分阈值取值范围为 0.01~0.1。

4) 对比方案

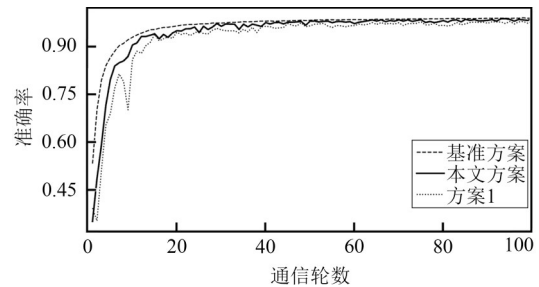
为了显示本文所提方案的先进性,本文采用传统联邦平均算法 (FedAvg, federated gradient average) [5] 作为基准方案,并与方案 1 [27]、方案 2 [19]、方案 3 [28] 在模型准确率、平均时间开销两方面进行对比。其中,方案 1 使用区块链审计客户端更新以抵抗恶意设备的攻击,并采用差分隐私技术保护特征数据隐私。方案 2 使用余弦相似度来判断恶意客户端上传的恶意梯度,增强模型抵抗中毒攻击的有效性。方案 3 采用分组策略汇总梯度,并对不同分组梯度进行拜占庭鲁棒聚合,实现 FL 鲁棒性的增强。

5.2 MAPBFL 模型性能分析

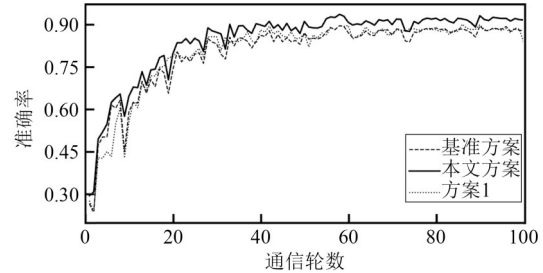
1) 准确率分析

为了评估本文方案的准确率,实验首先在没有

恶意设备参与模型训练的情况下,在 MNIST 和 CIFAR-10 数据集上对基准方案、本文方案、方案 1 的模型准确率进行对比,实验中的准确率定义为正确分类的样本数量占样本总数的百分比。实验在每次训练中均设置了 20 个边缘设备执行 100 轮通信,每轮通信结束后产生的全局模型准确率均被记录,实验结果取平均值。其中,参与每轮训练的 20 个边缘设备被固定划分为 15 个边缘客户端设备、3 个验证器以及 2 个矿工设备,所有设备均可进行本地模型训练。在 MNIST 和 CIFAR-10 数据集上不同方案的模型准确率对比如图 3 所示。



(a) MNIST 数据集上不同方案的模型准确率对比



(b) CIFAR-10 数据集上不同方案的模型准确率对比

图3 不同方案的模型准确率对比(没有恶意设备参与模型训练)

如图 3(a) 所示,在 MNIST 数据集上,初始几轮通信中,基准方案的模型准确率最高,本文方案次之,方案 1 较低,这表明受初始参数随机性的影响,本文方案在前几轮通信中的模型准确率偏低,但模型训练效果的提升是最快的。在第 18 轮通信时,本文方案的模型准确率达到 92.3%,与基准方案的模型准确率基本持平,模型训练效果的提升速度与基准方案基本一致,而方案 1 在第 23 轮通信之后,模型准确率才逐渐与本文方案和基准方案相近。在第 35 轮通信之后,本文方案与基准方案所得模型逐渐收敛,模型准确率达到 97.2%,而方案 1 的模型准确率出现波动,模型收敛速度较慢。

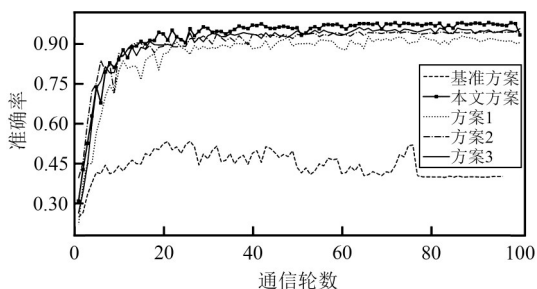
从图 3 可以看出,本文方案在使用秘密共享和差

分隐私技术提升模型隐私安全性能的同时，可以保持与基准方案相近的高准确率，同时实现了比方案1更快的模型收敛速率和更高的模型准确率。如图3(b)所示，本文方案在CIFAR-10数据集上获得了最高的模型准确率90.8%。所以本文方案可以在保证模型较高准确率的情况下，增强模型隐私安全性能。

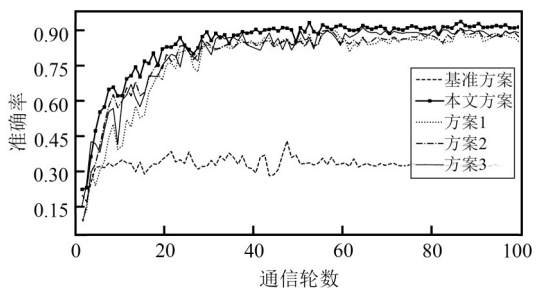
2) 抵御恶意敌手攻击能力分析

为了验证本文方案抵御恶意敌手攻击模型可用性的优势，本节实验分别设置了不同数量的恶意设备参与使用基准方案、本文方案、方案1、方案2、方案3进行的模型训练，其中，恶意设备使用的攻击方法默认是高斯攻击。

在 20% 的恶意设备参与训练的情况下，在 MNIST 和 CIFAR-10 数据集上不同方案的模型准确率对比如图 4 所示。如图 4(a) 所示，基准方案的准确率仅有 40%，因为基准方案没有额外的隐私安全机制，20% 的恶意设备共享虚假的模型参数对全局模型精度造成巨大影响，基准方案的准确率大幅下降。方案 1、方案 2 和方案 3 的模型准确率出现波动，分别维持在 86%、90% 和 91%，然而本文方案的模型准确率最终稳定在 94.3%。如图 4(b) 所示，在 CIFAR-10 数据集上，相较于其余对比方案，本文方案仍然保持较优的模型准确率，收敛后的模型准确率维持在 90%，可以有效抵抗 20% 的恶意设备破坏模型训练的影响。



(a) MNIST 数据集上不同方案的模型准确率对比

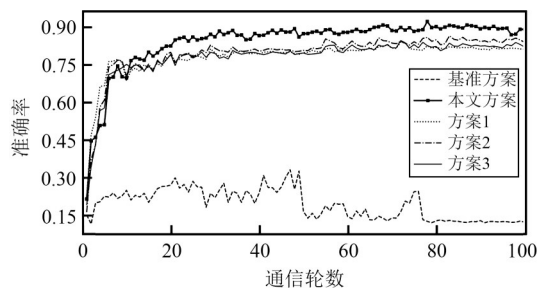


(b) CIFAR-10 数据集上不同方案的模型准确率对比

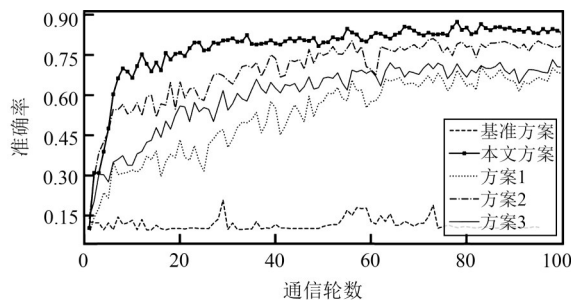
图 4 不同方案的模型准确率对比(20%的恶意设备参与训练)

如图 4 所示，本文方案、方案 1、方案 2 和方案 3 在不同通信轮数的准确率会出现一定幅度的波动，这是因为识别恶意设备和恶意梯度需要时间。其中，本文方案采用客户端信誉积分验证来识别恶意客户端，然而更新全局设备的信誉积分需要一定时间，所以当恶意设备出现时，模型准确率会出现暂时的骤然下降，当全局设备信誉值更新完成后，恶意设备被阻止参与后续的模型训练，此时模型准确率会逐渐提升。同时，从模型准确率波动幅度来看，本文方案的模型准确率受恶意设备影响小于对比方案，本文方案训练效果的提升速率相较于对比方案更加稳定，这是因为本文方案相较于对比方案更多地考虑了验证方成为恶意设备的可能性。本文方案不仅考虑模型提供方的信誉奖励和惩罚，还增加了验证方的信誉评估，通过结合样本数量和模型参数来设计新型的信誉奖惩机制，帮助模型快速识别恶意验证方及恶意客户端。

在 40% 的恶意设备参与训练的情况下，在 MNIST 和 CIFAR-10 数据集上不同方案的模型准确率对比如图 5 所示。



(a) MNIST 数据集上不同方案的模型准确率对比



(b) CIFAR-10 数据集上不同方案的模型准确率对比

图 5 不同方案的模型准确率对比(40%的恶意设备参与训练)

如图 5(a) 所示，在 MNIST 数据集上，随着系统中恶意设备的数量增加至 40%，基准方案的模型准确率逐渐下降，在第 78 轮通信之后，模型准确率低

于10%，失去可用性。此时，本文方案和对比方案训练模型仍能保持较优的模型准确率，本文方案从第40轮通信开始，模型准确率保持在88.2%，且模型准确率逐渐趋于稳定，而对比方案的模型准确率受恶意设备数量增加的影响较大。方案1和方案3的模型准确率稳定后保持在75%左右，模型仍具有一定的可用性，但是模型准确率相较于本文方案低了约13.2%。方案2在第60轮通信之后，模型准确率稳定在78%左右，相较于本文方案低了约10%。如图5(b)所示，在CIFAR-10数据集上，本文方案在抵御40%恶意敌手攻击的同时仍能保持81.3%的模型准确率，比对比方案中效果最优的方案2高出3.2%。

实验结果表明，在恶意设备逐渐增加的情况下，本文方案性能高于基准方案和其余对比方案，能够获得较高准确率的全局模型。

3) 时间开销分析

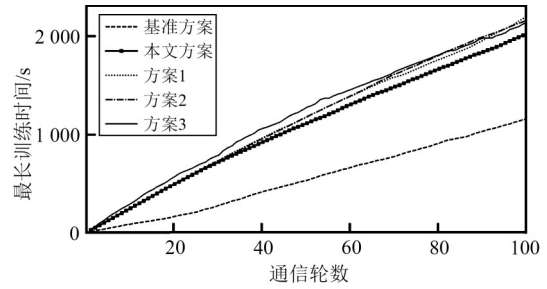
为了分析本文方案实际的运行效率，本节实验测试了基准方案、本文方案、方案1、方案2和方案3在MNIST和CIFAR-10数据集上执行100轮通信所耗费的最长训练时间 train_time ，每轮通信的最长时间 $\text{time}_t^{\text{slowest}_m}$ 选取本轮训练速度最慢的设备完成本轮训练所花费的时间，则不同方案的最长训练时间为

$$\text{train_time} = \sum_{t=1}^{100} \text{time}_t^{\text{slowest}_m} \quad (14)$$

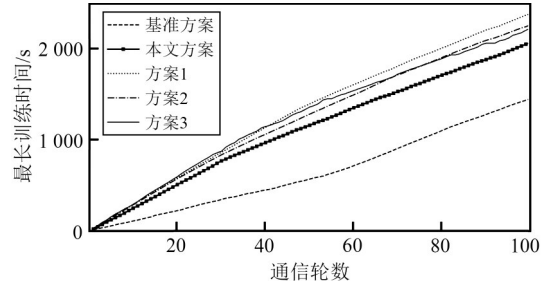
不同方案在MNIST和CIFAR-10数据集上完成100轮通信的最长训练时间如图6所示。如图6(a)所示，对于MNIST数据集，基准方案完成100轮通信所用时间最少，需要约1200s，本文方案需要约1800s，是基准方案的1.5倍，方案1、方案2和方案3的用时接近，需要约2200s，是基准方案的1.83倍。在MNIST数据集上，相较于对比方案，本文方案可以取得较优的模型训练速率。如图6(b)所示，在CIFAR-10数据集上，本文方案执行100轮通信所耗费的最长训练时间为2000s，与对比方案2和方案3所用的时间成本相近，但仍优于方案1、方案2和方案3。

从实验结果可知，本文方案使用区块链构建可信底层环境，引入秘密共享和差分隐私安全策略，增加验证机制和信誉计算会为系统带来额外的时间成本，但同时为零信任的边缘计算环境提供了更高的隐私性和安全性。相较于基准方案在没有任何安全策略下的时间成本，本文方案时间成本仅高于基

准方案1.5倍，考虑本文方案隐私性和安全性的优势，增加的时间成本是可以接受的。



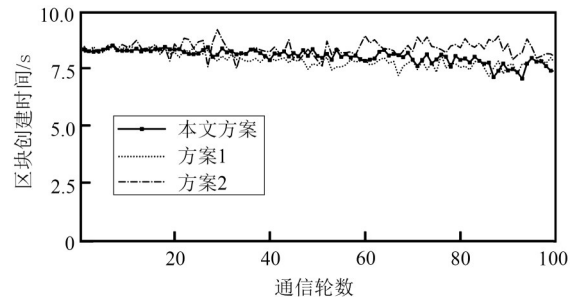
(a) MNIST数据集上不同方案的最长训练时间对比



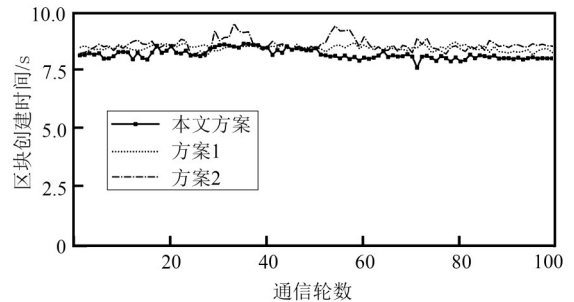
(b) CIFAR-10数据集上不同方案的最长训练时间对比

图6 不同方案的最长训练时间对比

在MNIST和CIFAR-10数据集上，本文方案、方案1和方案2每轮通信中区块创建时间对比如图7所示。



(a) MNIST数据集上不同方案每轮通信中区块创建时间对比



(b) CIFAR-10数据集上不同方案每轮通信中区块创建时间对比

图7 不同方案每轮通信中区块创建时间对比

如图 7(a)所示, 在 MNIST 数据集上, 本文方案在每轮通信中创建区块的时间平均为 8 s, 与方案 1 所用时间相近, 而方案 2 在每轮通信中创建区块的时间平均为 8.3 s。与对比方案 1 和方案 2 相比, 本文方案可以保证比较高效的交易确认, 提升了系统效率。如图 7(b)所示, 在 CIFAR-10 数据集上, 本文方案在每轮通信中创建区块的时间平均为 8.4 s, 方案 1 和方案 2 的用时均略高于本文方案, 因此, 在保证更高隐私性和安全性的同时, 本文方案是可取的。

由上述实验结果可知, 本文方案相较于基准方案和对比方案在保证较高模型隐私安全性能的同时缓解了隐私安全机制设计给模型训练效率带来的压力, 保证了较优的模型训练效率。

5.3 基于信誉验证的联邦任务节点选择算法效果分析

为了验证本文设计的基于信誉验证的联邦任务节点选择算法的效果, 本节实验分别评估了模型训练中本文算法识别恶意设备的效果, 以及恶意设备在模型训练过程中的信誉增长情况。

首先, 设置 30% 的恶意设备参与训练, 其中包括恶意客户端设备和恶意验证方设备。验证方设备接收客户端设备上传的本地模型, 并在 MNIST 和 CIFAR-10 数据集上评估该本地模型效果, 评估过程如下。验证方设备设定验证阈值为 0.05, 并计算局部验证模型和客户端设备上传的本地模型的精度差值 $accuracy_i^{m,v}$, 如式 (15) 所示, 通过比较 $accuracy_i^{m,v}$ 和验证阈值的大小来评估客户端设备上传的本地模型质量。

$$accuracy_i^{m,v} = L_i^v - L_i^m \quad (15)$$

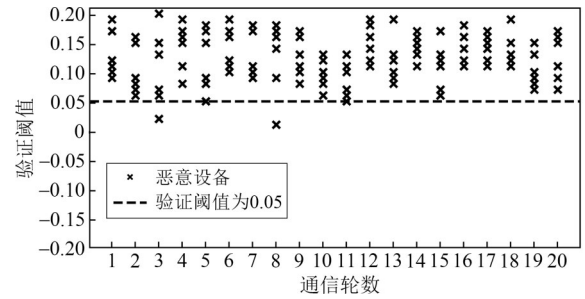
其中, L_i^v 为第 t 轮训练中的验证方设备的局部验证模型, L_i^m 为客户端设备 m 在第 t 轮训练中上传的本地模型。

其次, 本节实验设置了 3 名矿工使用上一轮聚合的全局模型 G_{t-1} 来评估边缘验证器上传的局部模型, 并计算模型精度差值 $accuracy_i^{v,s}$, 如式 (16) 所示, 最终 3 名矿工根据验证阈值给出投票结果。

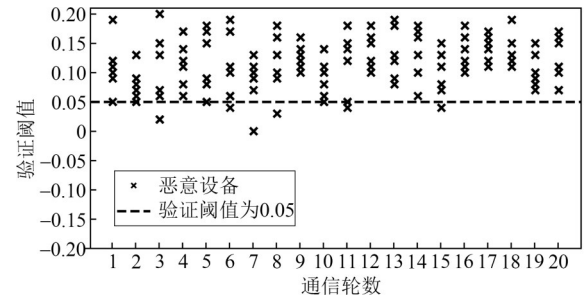
$$accuracy_i^{v,s} = G_{t-1} - L_i^v \quad (16)$$

在 MNIST 和 CIFAR-10 数据集上进行 6 次模型训练的前 20 轮通信中系统识别的恶意设备的效果如图 8 所示。经过多轮实验证明, 当验证阈值设定为 0.05 时, 绝大多数恶意设备可以被有效识别, 此

时算法识别恶意设备的性能达到最优, 即算法设定精度差值大于 0.05 时, 上传的局部模型均被视为低质量的恶意模型, 提供该模型的相关设备也将被扣除相应信誉积分, 并标记为恶意设备。



(a) MNIST 数据集上系统识别恶意设备的效果



(b) CIFAR-10 数据集上系统识别恶意设备的效果

图 8 不同数据集上系统识别恶意设备的效果

图 9 和图 10 分别显示了在 MNIST 和 CIFAR-10 数据集上进行模型训练时的每轮通信的设备信誉积分累计情况。本节实验分别在不添加验证机制和添加本文算法所提的信誉验证机制的情况下, 验证了模型训练过程中各设备信誉积分增长的情况。在 MNIST 数据集上, 如图 9(a)所示, 在未设置验证机制的情况下, 恶意设备的信誉积分随着模型的迭代训练在不断增加。如图 9(b)所示, 加入本文算法设计的验证机制后, 在初始通信轮数中恶意设备的信誉积分便停止增长。如图 9(c)所示, 随着模型训练的加入验证机制后的第 2 次训练中, 恶意设备数量减少, 并在早期阶段停止信誉增长。在第 4~6 次训练中, 如图 9(d)~图 9(f)所示, 恶意设备已不再被选中参与模型训练, 只有诚实设备的信誉积分会随着训练的进行而不断增长。在 CIFAR-10 数据集上, 如图 10(b)~图 10(d)所示, 本文算法在前 3 次训练的初始通信轮数中很好地阻止了恶意设备信誉积分的增长, 说明本文算法在 CIFAR-10 数据集上仍然可以取得较好的恶意设备识别效果。

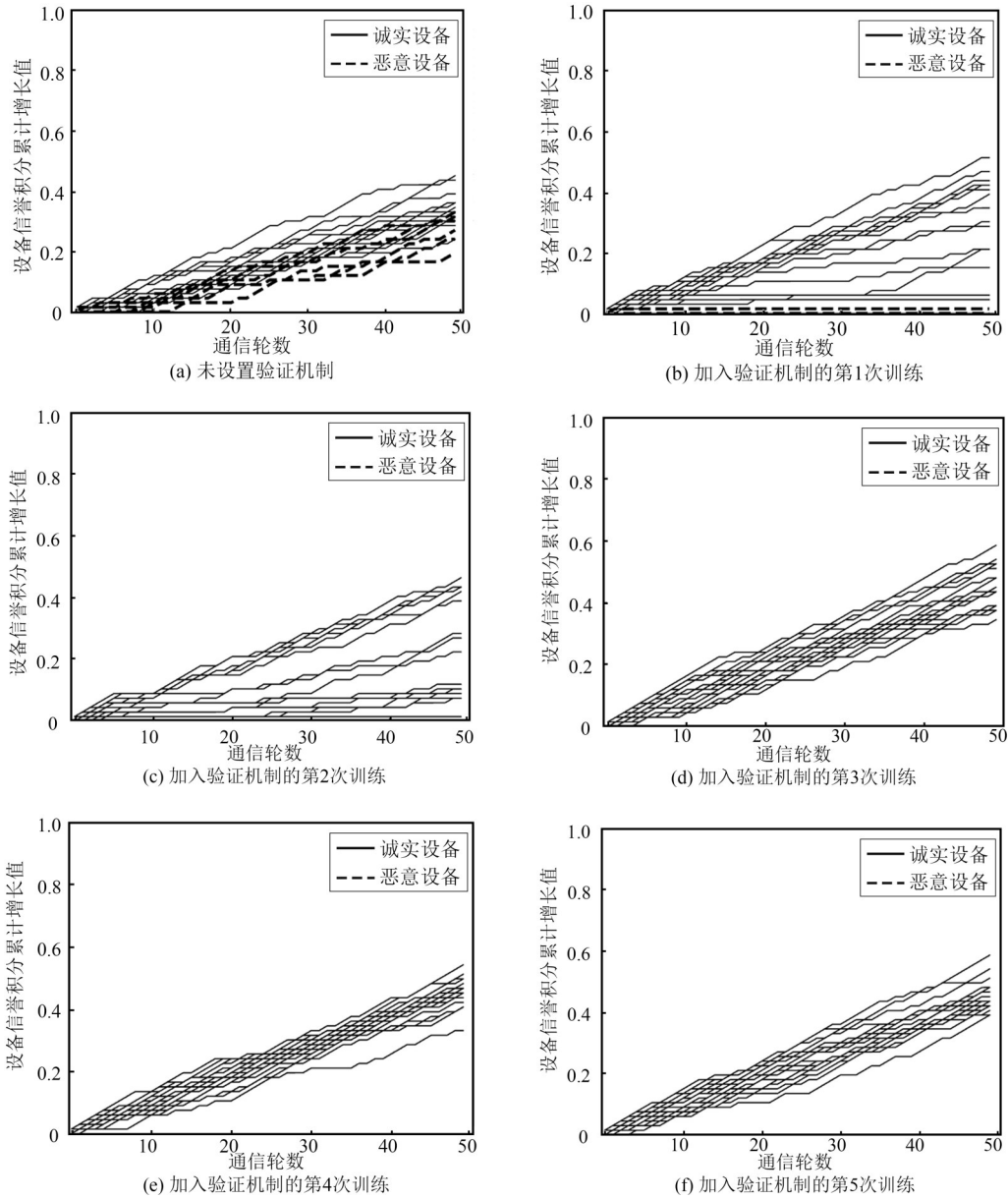


图9 MNIST数据集上进行模型训练时的每轮通信的设备信誉积分累计情况

由上述实验结果可知, 本文设计的基于信誉验证的联邦任务节点选择算法可以准确发现训练过程中的恶意设备, 并有效阻止其参与下一轮的学习过程, 使用此算法可以明显提升模型鲁棒性和收敛速率, 保证模型安全。

6 结束语

针对目前零信任边缘计算下的联邦学习仍面临数据安全攻击和隐私泄露的问题, 本文提出了基于多级代理许可区块链的联邦边缘学习模型, 为联邦边缘学习构建了可信底层环境, 保证边缘用户的可信接入。然后设计了分层模型聚合, 有效缓解了大

规模边缘数据的处理压力, 并通过秘密共享和差分隐私混合策略的设计进一步增强了分层模型聚合的安全性和隐私性。进一步地, 本文设计了一种基于信誉验证的联邦任务节点选择算法, 有效提高了模型鲁棒性。根据实验结果可知, 相较于现有的先进方案, 本文方案能够在40%的恶意敌手攻击下, 将模型准确率提升10%, 实现了较高隐私安全性能的模型训练。未来研究工作将进一步完善本文信誉评分机制在模型鲁棒性能提升工作中的应用, 研究更加高效的混合安全策略, 适用于更高效和更高安全性要求的边缘计算场景。

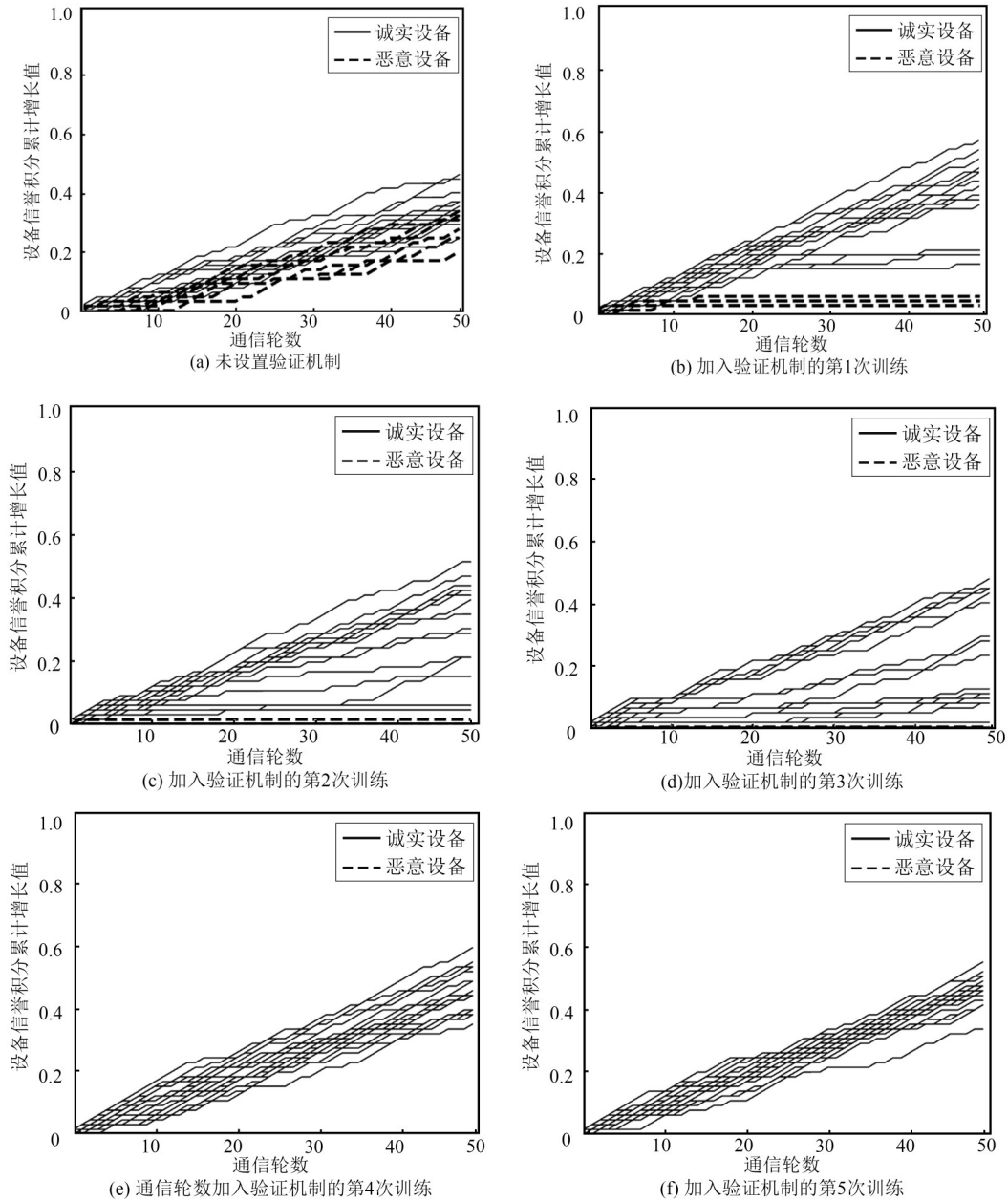


图 10 CIFAR-10数据集上进行模型训练时的每轮通信的设备信誉积分累计情况

参考文献:

[1] AHMAD J M, ZHANG W J, KHAN F, et al. Lightweight and smart data fusion approaches for wearable devices of the Internet of Medical Things[J]. Information Fusion, 2024, 103: 102076.

[2] REHMAN A, SABA T, KASHIF M, et al. A revisit of Internet of things technologies for monitoring and control strategies in smart agriculture[J]. Agronomy, 2022, 12(1): 127-130.

[3] TRUONG H T, TA B P, LE Q A, et al. Light-weight federated learning-based anomaly detection for time-series data in industrial control systems[J]. Computers in Industry, 2022, 140: 103692.

[4] SHANG S, LI X, LU R X, et al. A privacy-preserving multidimensional

range query scheme for edge-supported industrial IoT[J]. IEEE Internet of Things Journal, 2022, 9(16): 15285-15296.

[5] MCMAHAN B, MOORE E, RAMAGE D, et al. Communication-efficient learning of deep networks from decentralized data[C]//Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (AISTATS). New York: PMLR, 2017: 1273-1282.

[6] SU Y F, HUANG C W, ZHU W W, et al. Multi-party diabetes mellitus risk prediction based on secure federated learning[J]. Biomedical Signal Processing and Control, 2023, 85: 104881.

[7] BYRD D, POLYCHRONIADOU A. Differentially private secure multi-party computation for federated learning in financial applications[C]//Proceedings of the First ACM International Conference on AI in Fi-

- nance. New York: ACM Press, 2020: 1-9.
- [8] VYAS J, DAS D, CHAUDHURY S. Federated learning based driver recommendation for next generation transportation system[J]. Expert Systems with Applications, 2023, 225: 119951.
- [9] GE L N, LI H A, WANG X, et al. A review of secure federated learning: privacy leakage threats, protection technologies, challenges and future directions[J]. Neurocomputing, 2023, 561: 126897.
- [10] GOLDREICH O. Secure multi-party computation[J]. Manuscript (Preliminary Version), 1998, 78: 110.
- [11] RIVEST R L, ADLEMAN L, DERTOUZOS M L. On data banks and privacy homomorphisms[J]. Foundations of Secure Computation, 1978, 4(11): 169-180.
- [12] DWORK C, MCSHERRY F, NISSIM K, et al. Calibrating noise to sensitivity in private data analysis[C]//Proceedings of the 2006 Proceedings of the Third Conference on Theory of Cryptography. Berlin: Springer, 2006: 265-284.
- [13] LI J L, YAN T J, REN P C. VFL-R: a novel framework for multi-party in vertical federated learning[J]. Applied Intelligence, 2023, 53(10): 12399-12415.
- [14] SCHNEIDER T, SURESH A, YALAME H. Comments on "Privacy-enhanced federated learning against poisoning adversaries" [J]. IEEE Transactions on Information Forensics and Security, 2023, 18: 1407-1409.
- [15] MA X, ZHOU Y Q, WANG L H, et al. Privacy-preserving Byzantine-robust federated learning[J]. Computer Standards & Interfaces, 2022, 80: 103561.
- [16] YUAN X, NI W, DING M, et al. Amplitude-varying perturbation for balancing privacy and utility in federated learning[J]. IEEE Transactions on Information Forensics and Security, 2023, 18: 1884-1897.
- [17] WATKINS W M, CHEN S Y C, YOO S. Quantum machine learning with differential privacy[J]. Scientific Reports, 2023, 13(1): 2453.
- [18] LIU C, TIAN Y L, TANG J C, et al. A novel local differential privacy federated learning under multi-privacy regimes[J]. Expert Systems with Applications, 2023, 227: 120266.
- [19] MIAO Y, LIU Z, LI H, et al. Privacy-preserving Byzantine-robust federated learning via blockchain systems[J]. IEEE Transactions on Information Forensics and Security, 2022, 17: 2848-2861.
- [20] WANG Q, CHEN S G, WU M. Incentive-aware blockchain-assisted intelligent edge caching and computation offloading for IoT[J]. Engineering, 2023, 31: 127-138.
- [21] KALAPAAKING A P, KHALIL I, ATIQUZZAMAN M. Blockchain-enabled and multisignature-powered verifiable model for securing federated learning systems[J]. IEEE Internet of Things Journal, 2023, 10(24): 21410-21420.
- [22] MAYHOUB S, SHAMI T. A review of client selection methods in federated learning[J]. Archives of Computational Methods in Engineering, 2023, 31: 1129-1152.
- [23] XU X L, NIU S S, WANG Z, et al. Client selection based weighted federated few-shot learning[J]. Applied Soft Computing, 2022, 128: 109488.
- [24] PENE P, LIAO W, YU W. Incentive design for heterogeneous client selection: a robust federated learning approach[J]. IEEE Internet of Things Journal, 2023, 11(4): 5939-5950.
- [25] YANG J R, LIU Y, KASSAB R. Client selection for federated Bayesian learning[J]. IEEE Journal on Selected Areas in Communications, 2023, 41(4): 915-928.
- [26] MCSHERRY F D. Privacy integrated queries: an extensible platform for privacy-preserving data analysis[C]//Proceedings of the 2009 ACM SIGMOD International Conference on Management of data. New York: ACM Press, 2009: 19-30.
- [27] ZHAO Y, ZHAO J, JIANG L, et al. Privacy-preserving blockchain-based federated learning for IoT devices[J]. IEEE Internet of Things Journal, 2020, 8(3): 1817-1829.
- [28] LI Y L, YUAN D, SANI A S, et al. Enhancing federated learning robustness in adversarial environment through clustering Non-IID features[J]. Computers & Security, 2023, 132: 103319.

[作者简介]



葛丽娜 (1969-), 女, 广西环江人, 博士, 广西民族大学教授、硕士生导师, 主要研究方向为信息安全、物联网、区块链和智能计算等。



栗海澳 (1999-), 男, 河南信阳人, 广西民族大学硕士生, 主要研究方向为网络与信息安全、联邦学习和区块链。



王捷 (1995-), 女, 广西柳州人, 广西民族大学讲师, 主要研究方向为网络与信息安全 and 区块链。