

基于区块链的隐私保护可信联邦学习模型

朱建明¹⁾张沁楠¹⁾高胜¹⁾丁庆洋²⁾袁丽萍¹⁾

¹⁾(中央财经大学信息学院,北京 100081)

²⁾(北京联合大学管理学院,北京 100020)

摘要 基于联邦学习的智能边缘计算在物联网领域有广泛的应用前景。联邦学习是一种将数据存储在参与节点本地的分布式机器学习框架,可以有效保护智能边缘节点的数据隐私。现有的联邦学习通常将模型训练的中间参数上传至参数服务器实现模型聚合,此过程存在两方面问题:一是中间参数的隐私泄露,现有的隐私保护方案通常采用差分隐私给中间参数增加噪声,但过度加噪会降低聚合模型质量;另一方面,节点的自利性与完全自治化的训练过程可能导致恶意节点上传虚假参数或低质量模型,影响聚合过程与模型质量。基于此,本文将联邦学习中心化的参数服务器构建为去中心化的参数聚合链,利用区块链记录模型训练过程的中间参数作为证据,并激励协作节点进行模型参数验证,惩罚上传虚假参数或低质量模型的参与节点,以约束其自利性。此外,将模型质量作为评估依据,实现中间参数隐私噪声的动态调整以及自适应的模型聚合。原型搭建和仿真实验验证了模型的实用性,证实本模型不仅能增强联邦学习参与节点间的互信,而且能防止中间参数隐私泄露,从而实现隐私保护增强的可信联邦学习模型。

关键词 区块链; 联邦学习; 智能边缘计算; 差分隐私; 共识算法

中图分类号 TP18

Privacy Preserving and Trustworthy Federated Learning Model Based on Blockchain

ZHUJian-Ming¹⁾ZHANGQin-Nan¹⁾GAOSHeng¹⁾DINGQing-Yang²⁾YUANLi-Ping¹⁾

¹⁾(School of information, Central University of Finance and Economics, Beijing 100081)

²⁾(School of management, Beijing Union University, Beijing 100020)

Abstract Intelligent edge computing based on federated learning has a wide application prospect in the field of Internet of things (IoT). However, it is still faced with the dilemma of lacking enough data sources in the current practice of artificial intelligence. In this context, distributed machine learning aggregates edgedevices' raw data into a parameter server for model training, but it easily leads to data privacy leakage and causes excessive storage overhead. In particular, federated learning (FL) is a distributed machine learning framework that stores data locally, which can effectively protect the data privacy of edge intelligent nodes. According to client settings, FL can be classified into twotypes: cross-device FL and cross-silo FL. In cross-deviceFL, a central entity acts as the central parameter server, which is also theowner of the global model. Meanwhile, the participatingnodesas the clients to perform local training. In cross-silo FL, all participating nodesact as the clients to perform local training. In addition, they are also the owners of the global model and can makeuse of the trained global model. In this paper, we focus on cross-device FL, in which intelligence edge devices can provide model training services by sensing the raw data from IoT devices such as intelligence vehicles, smartphones etc. Most of the existing cross-device FL implements model aggregation by uploading the intermediate parameters of model training to the parameter server.

本课题得到国家重点研发计划(No.2017YFB1400700)、国家自然科学基金项目(No.62072487),北京市自然科学基金项目(No.M21036)资助。朱建明, 博士, 教授, 主要研究领域为区块链技术与信息安全.E-mail: zjm@cufe.edu.cn. 张沁楠 (通信作者), 博士研究生, 主要研究领域为区块链与智能边缘计算.E-mail: zhangqnp@163.com. 高胜 (通信作者), 博士, 副教授, 主要研究领域为区块链技术与信息安全.E-mail: sgao@mail.xidian.edu.cn. 丁庆洋, 博士, 主要研究领域为区块链与大数据治理. E-mail: dingqingyang66@163.com. 袁丽萍, 硕士研究生, 主要研究领域为区块链技术与隐私保护.E-mail: yuanliping_cufe@163.com.

There are two problems in this process. On the one hand, there is privacy leakage of intermediate parameters. The existing privacy protection schemes usually use differential privacy to add the noise on intermediate parameters, but excessive noise will reduce the quality of the global model. On the other hand, the training process of node self-interest and full autonomy may lead to malicious nodes uploading false parameters or low-quality models, thus affect the aggregation processes and model quality. In this paper, the centralized parameter server federated learning is constructed as a decentralized parameter aggregation chain, and the intermediate parameters of the model training process recorded on the blockchain as evidence. Moreover, the cooperative nodes are encouraged to verify the model parameters and punishes the participating nodes who upload false parameters or low-quality models so as to restrict their self-interest. In view of above challenges, we take the model quality as the metric to dynamically adjust privacy noise of intermediate parameters and propose a federated adaptive (FedAdp) model aggregation algorithm. The prototype development and experimental simulations show that the proposed FedAdp model aggregation algorithm can achieve higher accuracy of aggregation model when occur poisoning attack. By dynamically adjusting the Laplace random noise, it's realized the tradeoff between privacy protection and the accuracy error of the aggregation model. The experiment of blockchain performance confirmed that our scheme has good practicability. It is proved that the model can not only enhance the mutual trust between the participating nodes of federated learning, but also prevent the privacy disclosure of intermediate parameters, so as to realize the federated learning model with enhanced trust and privacy protection.

Key words blockchain; federated learning; intelligent edge computing; differential privacy; consensus algorithm

1 引言

联邦学习 (Federated Learning)^[1] 是一种协作式机器学习框架, 参与协作的节点利用本地数据训练模型, 通过参数聚合实现多来源数据的预测效果。当前人工智能在实践过程中仍然面临数据来源不足的困境。在医疗领域中, 标注数据需动用 1 万人长达 10 年的时间才能收集到足够多有效的数据^[2]。联邦学习中数据存储在节点本地实现分布式机器学习, 实现了隐私保护的数据协作。随着移动通信技术和智能边缘设备的兴起, 联邦学习在智慧城市^[3]、电子医疗^[4]、无线通讯^[5]、移动边缘网络^[6]等领域有着广泛的应用前景^[7]。目前联邦学习已产生基于同行业数据的横向联邦学习, 以及面向多行业数据的纵向联邦学习与联邦迁移学习^[8], 并与大数据、云计算、区块链、智能边缘计算等前沿技术深度融合, 成为产学研界共同关注的研究热点。

数字经济背景下, 数据合规成为世界趋势。2018 年, 欧盟颁布《通用数据保护条例》(GDPR) 严格规范数据的使用。2017 年 6 月 1 日起我国施行的《中华人民共和国网络安全法》指出不得泄露、篡改用户数据。2019 年 5 月 28 日, 我国国家互联网信息办公室公开《数据安全管理办法(征求意见稿)

稿)》, 并于 2020 年 12 月 1 日发布了《常见类型移动互联网应用程序 (APP) 必要个人信息范围》公开征求意见通知, 可见用户数据的流转和使用必须满足越来越严苛的数据管理条例。此外, 数据要素具有巨大的潜在价值, 但由于行业竞争、利益冲突等因素, 大多仍呈现数据孤岛形式。联邦学习满足数据合规, 并可以解决数据孤岛问题^[9]。

大数据时代个人隐私保护一直备受关注。数据隐私的泄露通常会引起公众不满, 例如 Facebook 数据泄露曾引发大范围抗议活动, 国内求职简历售卖也一度登上了微博热搜。造成数据隐私泄露的主要原因可能是数据在流转过程中丢失, 或者利用数据挖掘技术从海量用户信息中非法获取个人敏感信息^[10]。数据本身的无限复制特性导致数据一旦发生泄露, 数据的流转和使用难以追踪。

联邦学习是当前人工智能背景下实现数据隐私保护的有效办法, 根据初始设置不同 Kairouz 等人^[11]将联邦学习分为跨设备 (Cross-device) 和跨筒仓 (Cross-silo) 两种类型。跨设备联邦学习的全局模型由中心化的参数服务器控制, 跨筒仓联邦学习的每个参与者都可以是全局模型的聚合者和拥有者。本文主要关注跨设备联邦学习, 其基本思想是将模型训练分散在 k 个节点进行, 每轮训练结束后参数服务器收集节点本地模型参数执行模型聚合

算法，并将更新的全局模型参数返回给各节点继续迭代训练直至模型收敛。与分布式机器学习相比，联邦学习具有以下优势：（1）参与节点协作训练，利益共享；（2）数据不出节点本地，保护数据隐私并满足数据合规；（3）模型准确率与数据聚合之后训练的准确率效果相当。

跨设备联邦学习通过参数服务器执行模型聚合算法，但中心化参数服务器可能遭受恶意攻击被截获甚至篡改模型聚合过程的中间参数。此外，中心化参数服务器使参与节点与参数服务器之间存在大量远程数据通信，也导致了数据篡改与隐私泄露风险。为了解决此问题，已有学者提出去除第三方的联邦学习方法^[12]。在该方法中，通过交换用户公钥将模型聚合委托给参与训练的节点一方来执行。然而，现有的联邦学习隐私保护方案仍存在以下两方面问题：

（1）**中间参数隐私泄露**：联邦学习避免了因数据收集而引起的数据泄露问题，但是仍然存在中间参数的隐私泄露^[13]，尤其是梯度隐私泄露^[14]。现有机器学习中的模型攻击^[15]、数据攻击^[16]、推理攻击^{[17][18]}（Inference attack）、后门攻击（Backdoor attacks）^[19]、链接攻击（Linkability attack）^[20]、投毒攻击（Poisoning attack）^[21]等方法都可对中间参数包括梯度数据进行原始数据推断，从而泄露参与节点本地数据的敏感信息。

（2）**节点多方信任问题**：在模型聚合过程中，可能存在半诚实或者恶意的参数服务器与参与节点^[12]。首先，参数服务器可能滥用或泄露数据，通过将模型参数泄露给第三方获得额外收益；此外，参与节点由于自利性考虑，可能会提供虚假参数以提高活跃度与贡献度，存在搭便车行为。

如图 1 所示，假设企业 A、B 因数据匮乏有联合建模需求，企业考虑数据隐私采用联邦学习建模方案。在建模过程中，A、B 无需交换本地数据，只需将加密样本对其，向参数服务器交换中间加密参数。参数服务器执行模型聚合算法，将参数返回给 A、B 进行新一轮模型更新，然后不断迭代此过程直至模型收敛。在此过程中，假设企业 B 遭受恶意攻击而导致数据隐私泄露甚至恶意篡改。由于自利性考虑，企业 B 不愿消耗计算资源参与训练，但是想要从中获利，因此选择搭便车上传未经训练的虚假参数，这两种情况都会影响参数服务器的模型聚合。此外，在执行模型聚合算法的过程中也难以避免程序漏洞，从而影响联邦学习协作训练结果。

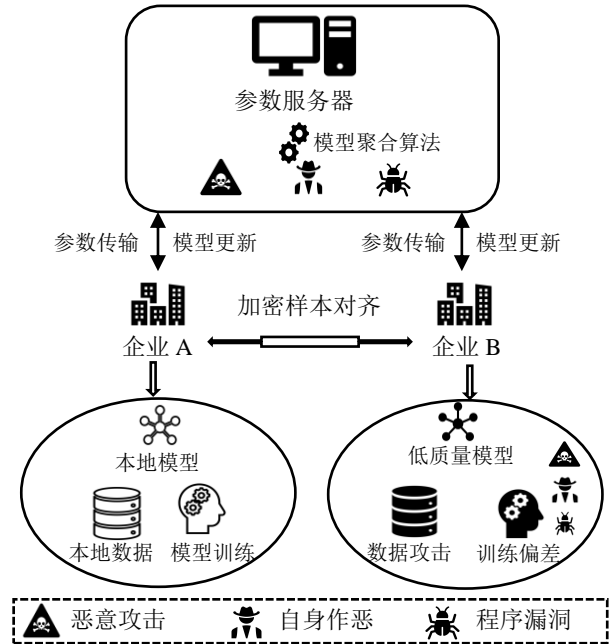


图 1 参与节点提供低质量模型以及不可信的参数服务器

综上所述，为了解决跨设备联邦学习中间参数隐私泄露以及多方互信问题，本文将中心化的参数服务器构建为去中心化的参数聚合链，其中参与训练节点 k 将中间参数上传至区块链，通过共识算法与智能合约进行参数验证和模型聚合，经过共识验证后的聚合模型参数返回给各参与节点进行模型更新。根据模型质量与节点信誉评分调整聚合权重，降低误差较大模型参数所占比重，从而提升聚合模型的准确率。基于这一研究思路，本文基于已有文献^[22]利用差分隐私实现联邦学习中间参数的隐私保护。并设计协作者异步参数审计机制对模型质量评估结果进行共识验证，防止搭便车行为。同时基于模型质量与信誉评分构建自适应的模型聚合算法，并通过智能合约自动触发算法执行。

与现有的基于区块链的联邦学习的工作相比，本文主要贡献归纳为：

（1）**自适应的模型聚合与隐私感知**：本文以交叉熵（Cross Entropy）^[23]为模型质量评估依据，基于模型质量与节点累计信誉值设计自适应的模型聚合算法。其次，基于差分隐私保护联邦学习中间参数隐私，根据模型质量评估结果在中间参数基础上增加不同程度的拉普拉斯（Laplace）^[24]噪声，防止中间参数在传输过程中泄露本地数据隐私。

（2）**区块链模型参数审计与共识**：利用区块链去中心化和强信任属性，对模型参数进行协作者参数审计，有效识别恶意节点的投毒攻击与搭便车行为。本文提出节点贡献度证明（Proofofcontribution，

PoC) 共识算法, 降低高贡献度节点共识挖矿难度, 减少计算资源浪费, 提升节点参与公平性。

(3) **支持轻量级智能边缘计算:** 目前已有的基于区块链的联邦学习工作^{[25][26]}大多将联邦学习节点构建为点对点 (Peer-to-Peer) 的区块链结构, 本文将区块链和联邦学习分层构建, 适用于轻量级智能边缘计算场景中的联邦学习。其中, 区块链作为模型聚合引擎, 通过智能合约实现模型参数的自动化可配置聚合操作, 降低了中心化参数服务的恶意攻击与计算偏差风险。轻量级智能边缘计算节点仅需完成本地模型训练和参数更新, 无需承担区块链节点面临的数据冗余和共识过程中的通信开销。

(4) **仿真实验与安全性分析:** 使用 MNIST 和 CIFAR-10 数据集, 基于拉普拉斯差分隐私测试在固定隐私噪声与动态隐私噪声情况下, MLP 与 CNN 模型训练的准确率对比情况。通过模拟投毒攻击, 进一步展示自适应模型聚合算法的抵御效果。利用 PythonFlask 搭建区块链平台, 为联邦学习提供去中心化的模型聚合服务。大量实验表明, 当参与节点使用本模型进行联邦学习模型训练, 计算时间开销与存储开销较小, 说明本模型有较好的实用性。

2 相关工作

2.1 联邦学习的隐私保护

联邦学习最早被 Google 提出用于用户输入法预测, 智能手机设备在本地进行模型训练后将参数上传至云服务器进行模型更新^[27]。虽然在联邦学习过程中数据保留在节点本地, 但仍有可能泄露用户隐私。Shmatiko^[28]等人提出过度学习的概念, 发现模型训练过程可能隐式学习到不属于学习目标的数据属性。Melis^[17]等人通过推理攻击发现中间梯度可能暴露训练数据的敏感信息。Hitaj^[16]等人指出通过生成对抗网络 (Generative Adversarial Network, GAN) 可以学习到差分隐私保护后参数中的重要信息。Orekondy^[20]等人发现通过链接攻击可以得到中间梯度包含的重要数据特征。Dillenberger^[29]等人讨论了联邦学习的安全问题, 地理上分散的节点向集中的服务器发送更新参数时不仅容易受到网络攻击, 而且存在着隐私泄露风险。

现有的联邦学习参数隐私保护方案主要包括差分隐私 (Differential Privacy, DP), 同态加密 (Homomorphic Encryption, HE) 和安全多方计算

(Secure Multi-party Computation, SMC) 等, 但同态加密和安全多方计算由于密钥管理和多方通信面临性能瓶颈问题, 而差分隐私性能影响较小。Shokri^[30]等人提出差分隐私可以将模型参数扰动后传输, 但是其代价是模型准确率的降低, 并且对于参与节点较少的模型影响更加明显。Choudhury^[31]等人提出了在联邦学习的上下文中提供动态差分隐私的方法, 该方法旨在实现节点效用最大化与模型性能提升, 同时支持 GDPR 和 HIPAA 法案要求的隐私级别。Aono^[32]等人提出采用加性同态加密保护深度学习中训练数据隐私性, 但是具有参与节点诚实的强假设。Bagdasaryan^[33]等人研究发现不诚实的参与方可能降低联邦学习聚合模型质量。Truex^[34]等人将差分隐私与安全多方计算结合, 能够在保障隐私的情况下, 随着参与方的增加而减少噪声注入的增长, 同时保证训练结果的可信性。HybridAlpha^[35]是一种采用基于功能加密的安全多方计算联邦学习协议, 实验证明 HybridAlpha 可以减少约 68% 的训练时间, 并减少 92% 的数据传输量, 同时实现与现有解决方案相同的模型准确率和隐私保护级别。

2.2 基于区块链的联邦学习可信增强

区块链^{[36][37]}通过密码技术和点对点通讯形成去中心化的分布式账本, 利用智能合约实现可编程的自动化可信交易, 在实现去中心化数据共享的同时确保数据的不可篡改。在结合联邦学习与区块链研究中, 一个主要方向是引入区块链来构建联邦学习分布式可信计算框架, 从而增强节点间的互信与聚合模型的可信性。引入区块链的主要原因如下: 1) 联邦学习中节点本身和相互通信的过程中都可能遭到恶意攻击; 2) 参与联邦学习的节点可能存在自身作恶或者因缺乏合理的激励机制而产生搭便车行为, 区块链的引入则可以有效解决此问题。

目前, 已有一些学者将区块链和联邦学习结合, 以解决联邦学习应用过程中的数据安全性问题。如果联邦学习的中心化参数服务器存在单点失效和恶意攻击, 其他参与节点将无法获知全局模型。FLChain^[38]提出采用区块链实现不可篡改的局部模型参数更新方案, 设计了一个基于区块链点对点网络的联邦学习架构, 将全局模型以默克尔帕特里夏树 (Merkle Patricia Tree, MPT) 结构存入区块链。Awan^[39]等人提出了一个基于区块链的隐私保护框架, 消除了参与者的半诚实假设, 并采用加密技

术来保护数据隐私,但仍存在参数服务器单点失效风险。BAFFLE^[40]利用去中心化的智能合约来协调联邦学习中的模型聚合和参数更新任务,通过将全局参数空间分解为不同的块,通过遵循评分和投标策略提高计算性能。FL-Block^[41]是一种基于区块链的联邦学习方案,通过采用分布式哈希表来保证区块生成效率,利用工作量证明共识机制维护全局模型一致性。Kim^[42]等人将区块链引入联邦学习解决单点失效问题,并分析了最佳区块生成率。但忽略了中间参数隐私保护。Zhao^[43]等人提出将基于声誉的联邦学习系统应用于移动边缘计算,并引入差分隐私来保护客户的敏感信息。许多大型互联网公司已经开展了区块链和联邦学习的研发工作并且代码开源,其中包括 IBM¹, 微软²等。整体来看,现有研究大多忽略了噪声误差对全局模型的影响,以及合作式模型训练中的参与积极性与公平性。

2.3 基于区块链的联邦学习节点贡献度激励

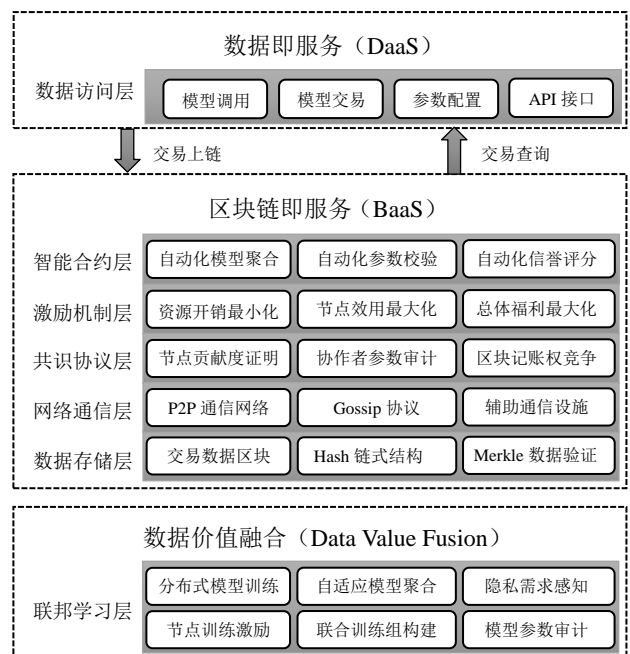
在现有的联邦学习中,假设参与节点拥有足够的本地数据用于模型训练并且愿意参与协作训练过程,但这与实际情况往往不符。比如,在电子医疗领域,医药公司或研究机构很难收集到用户的医疗数据,由于数据的特殊敏感性,大多医院不愿共享病人的私人数据。此外,数据拥有方也担心无法达到预期收益而不愿参与联邦学习。

现有研究已将区块链应用于联邦学习的参与节点激励中,用于提升节点的参与积极性与公平性。Deepchain^[44]是一个基于区块链激励的深度学习框架,利用区块链和密码学实现了隐私保护的分布式深度学习。通过 Deepcoin 平台资产奖励参与方对模型训练的贡献,以提高节点参与的活跃性。此外,Deepchain 通过智能合约实现自动化超时检查和参数验证以实现激励的公平性。Kang^[45]等人提出了一种基于区块链的信誉值评估管理方案,利用多权重主观逻辑模型计算参与节点信誉值并通过契约理论实现联邦学习中可靠参与者的选择与激励。DP-AFL^[46]是一种用于车联网的异步联邦学习算法,利用分布式异步更新方案避免中心化模型聚合的安全威胁。Kim^[25]等人认为在缺乏有效激励措施的情况下设备可能存在欺骗行为,因此要求区块链网络节点在共识过程中验证参数结果,对不诚实行为给予一定的惩罚机制。Lu^[47]等人提出在物联网中

将联邦学习节点构建为区块链网络,通过差分隐私实现了隐私保护的数据共享方案,并提出一种基于模型质量证明(Proof of Quality, PoQ)的区块链共识算法用于降低节点计算资源的开销,但该方案仍难以避免投毒攻击以及节点的搭便车行为。

在已有研究的基础上,我们将区块链作为模型聚合计算引擎,支持轻量级边缘节点的参与。根据模型训练损失与历史累计信誉评分自适应调整本地模型聚合权重,提升聚合结果的可信性。通过基于拉普拉斯机制的差分隐私实现联邦学习的中间梯度参数的动态隐私保护。设计基于节点贡献度的共识算法,降低节点计算资源开销,提升参与积极性和对聚合模型的贡献度。

基于区块链的可信联邦学习模型的研究架构如下图2所示。联邦学习层通过分布式本地模型训练、自适应模型聚合、隐私需求感知等实现隐私保护的数据价值融合;区块链服务层通过链式区块结构共享模型参数,保证数据的不可篡改和可追溯特性,并通过激励机制和智能合约实现参与节点的公平性与可信的自动化模型聚合;基于区块链的隐私保护可信联邦学习实现数据驱动的个性化模型调用与交易服务,用户通过参数配置和 API 接口即可访问基于区块链和联邦学习提供的安全可信数据模型服务。在数字经济时代,数据即服务(DaaS)可以促进数据拥有方的紧密合作,发挥海量数据背后的潜在价值,降低各类用户构建数据模型的投入成本,为大数据交易市场提供隐私保护的可靠模型构建及交易范式。



1 <https://github.com/IBM/federated-learning-lib>

2 <https://github.com/microsoft/0xDeCA10B>

图2 基于区块链的隐私保护可信联邦学习模型研究架构图

3 预备知识

本文将区块链作为可信模型参数共享平台，记录本地训练的中间梯度参数和训练损失，并通过区块链节点执行智能合约实现联邦学习模型聚合，降低中心化参数服务器的恶意攻击与不可信风险。为了防止在联邦学习数据传输过程中参与节点身份信息暴露以及数据隐私泄露，我们首先通过数据声明创建联合训练小组，数据声明部分利用伪公钥地址加密和数字签名技术防止攻击者非法获取节点的本地数据以及用户抵赖行为。

3.1 数据声明

传统联邦学习通过参与节点本地数据合作训练获得统一的全局模型。其中，参与节点 P 代表模型的数据拥有方，对本地数据具有完全的自治权，中心化的参数服务器负责聚合训练节点的中间参数形成全局模型。任务的目标函数可定义为^[48]：

$$\min_m F(m), F(m) := \sum_{k=1}^N w_k F_k(m) \quad \#(1)$$

其中， N 代表参与联邦学习训练节点总量， F_k 为第 k 个训练节点的本地目标函数， w_k 为对应节点的影响权重，满足 $w_k \geq 0$ 且 $\sum_{k=1}^N w_k = 1$ 。 $F_k(m)$ 通常被定义为基于本地数据训练模型 m 的经验风险。由于数据量和隐私性考虑，参与节点本地数据不能直接存储在区块链上。当一个数据提供者参与到联邦学习过程中，会生成唯一的分布式身份标识（DID）存储在链上，同时保存数据配置文件参数，包括数据的关键词、类型、大小、Hash等。配置文件参数以合约账户的键值对形式存储在区块链上，其中键保存参与节点的身份标识信息，值保存配置文件参数，通过MPT树^[49]进行数据的存储与验证。

表1 文中符号描述

符号	描述
pk_p^{psu}	参与节点 P 的伪公钥地址
skp	参与节点 P 的密钥
g_i	联合训练小组 i 的素数生成器
q	随机选择的大素数
Z_q^*	$\{1, 2, \dots, q-1\}$
G_1	循环倍增的素数序列组

H_1	抗冲突的 hash 函数 1
H_2	抗冲突的 hash 函数 2

开始训练之前，联邦学习任务发布者需要进行联合训练小组的建立。首先每个训练节点提前声明所拥有的数据信息，聚合相关关系较强的本地模型。为了保证参与节点的数据隐私和匿名性，本文借鉴 Deepchain^[44]采用伪公钥地址进行数据加密，伪公钥地址定义如下：

$$pk_p^{psu} \in \{g_1^{skp}, g_2^{skp}, \dots, g_n^{skp}\} \quad \#(2)$$

参与节点 P 选择密钥因子 $skp \in Z_q^*$ ，并且生成 n 个公钥地址 $g_i^{skp} \in G_1$ ， $i \in [1, n]$ 。参与节点 P_i 声明数据资产 $data_{P_i}$ ，并使用伪公钥地址 $pk_{P_i}^{psu}$ 进行加密，加密过程可表示为：

$$pk_{P_i}^{psu}(data_{P_i}) = \left\{ \left(g^{H_1(data_{P_i})}, \sigma_{j_{P_i}} \right), data_{key} \right\} \quad (3)$$

$$s.t. \sigma_{j_{P_i}} = \left(H_2(j) \cdot g^{H_1(data_{j_{P_i}})} \right)^{H_1(data_{P_i})}$$

其中 $data_{P_i}$ 经过Hash运算后可以防止信息传输过程中的数据篡改，将 $data_{P_i}$ 分别存储在 l 个区块中， $data_{j_{P_i}}$ 代表第 j 个区块中的内容， $\sigma_{j_{P_i}}$ 由 $data_{j_{P_i}}$ 经hash运算得到， $data_{key}$ 代表数据声明包含的参数信息，并将其作为联合训练组构成依据。

联邦学习目前的难点之一是因数据非独立同分布（Non-IID）^[50]导致的数据样本无法对齐问题。如果参与节点训练目标差异过大也无法进行协作训练。因此，联邦学习开始之前需要自动化聚集相关性较强的数据提供者。我们考虑将所有参与节点根据加密数据声明中的 $data_{key}$ 通过文本分类算法Word2Vec^[51]分为不同的联合训练组，同组的节点聚合目标相似。在训练过程中，每个节点维护一份检索日志 $log(i)$ ，当接收到数据模型访问请求时，节点首先检索本地日志，如匹配到所需模型结果直接返回，无需重复建模。

在区块链执行模型聚合算法过程中，各节点在本地训练之后将中间梯度参数上传至区块链网络中的共识节点进行模型聚合计算，各节点竞争记账权将达成共识的聚合模型结果记录上链。参与训练节点得到本地模型之后广播给联合训练组中的协作节点进行模型质量验证，协作节点利用自身本地数据验证模型，并将验证结果增加数据签名后进行广播。验证之后的模型参数通过自适应模型聚合算法形成本轮聚合模型 m_G ，由区块链网络中的主节

点返回给联合训练组中的参与节点，开始进行下一轮的模型更新迭代。

3.2 数据信息熵

联邦学习训练过程中需要多来源数据进行联合建模，数据的信息量决定整体聚合模型的准确率。根据经典的信息论原理，数据信息熵^[52]可用于衡量数据集中包含的信息量，因此可将其作为数据提供者的贡献度评判依据，从而提升联邦学习联合建模的准入门槛，从源头上防止低质量模型的引入，同时可作为模型收益公平分配的依据。我们首先对数据信息熵的相关概念进行界定。

定义 1. (数据元组) 对于给定的数据集 D ，元组 t 被定义为 D 中的一条记录 r 的非空子集，即 $t \subseteq r$ 且 $t \neq \emptyset$ 。

定义 2. (元组集合) 元组集合 Tup 是一系列元组 $\{t_{i_1}, t_{i_2}, \dots, t_{i_k}\}$ 的集合。因此，元组集合是数据集 D 的非空子集，即 $Tup \in D$ 且 $Tup \neq \emptyset$ 。元组集合可以是数据集的子集也可以是数据集本身。

定义 3. (数据信息熵) 数据信息熵是数据信息量测量的基础方法，可以测量单个元组集合的信息量，元组集合是信息熵度量的最小单位。对于一个有 n 个元组的元组集合，其信息熵 H_{ind} 的定义为：

$$H_{ind}(Tup) = - \sum_{t_i \in Tup} p(t_i) \log_b p(t_i) \quad \#(4)$$

其中， $p(t_i)$ 是数据元组 t_i 的概率密度函数， b 是对数函数的基，当 $b = 2$ 时，信息量的度量单位是比特 (Bit)，即信息论中常用的数值单位。

除此之外，对于连续型的数据集在计算信息熵时，由于样本数量有限很难准确计算其概率密度函数，因此可以采用窗口函数的概率密度来估算连续型数据集的信息熵。对于连续型数据集 $x = \{x_1, x_2, \dots, x_n\}$ ，如果真实概率密度函数是 $p(x)$ ，其估计密度函数 $\hat{p}(t_i)$ 定义为：

$$\hat{p}(t_i) = \frac{1}{n} \int_{i=1}^n \phi(x - x_i, h) \quad \#(5)$$

其中 ϕ 是窗口函数， h 是窗口高度，在本文中采用高斯函数作为窗口函数：

$$\phi(z, h) = \frac{1}{(2\pi)^{\frac{d}{2}} h^d \|\sigma\|^{\frac{1}{2}}} \exp\left(-\frac{z^T \sigma^{-1} z}{2h^2}\right) \quad \#(6)$$

其中， z 是数据集的 d 维随机向量， σ 是协方差矩阵。

3.3 激励机制

联邦学习涉及多个相互独立的参与方，建立一个合理的激励机制有利于公平地分配利益，有效降低节点搭便车风险。在信息不对称的情况下，激励机制的设计需要满足个人理性 (Individual Rationality, IR) 与激励相容 (Incentive Compatibility, IC)，以确保参与节点得到充分激励^[53]。联邦学习中每个参与节点训练目标是实现节点效用最大化，每个节点希望从协作式的联邦学习中获得比单一节点训练更优的模型质量。此外，平台通过对外提供数据模型交易获得收益，参与节点 P_k 可以分得一定的平台奖励 $R(P_k)$ 。节点参与训练的主要成本是其模型训练过程中存储资源与通信资源开销。因此，参与节点 P_k 的效用函数 $U(P_k)$ 定义如下：

$$U(P_k) = R(P_k) + \sum_{\tau=t}^T (Q_{\tau}(M_G) - Q_{\tau}(M_{P_k}) - C_{\tau}(P_k)) \quad \#(7)$$

其中， τ 是模型迭代次数， $Q_{\tau}(M_G)$ 是经过联邦学习模型聚合之后的全局模型质量， $Q_{\tau}(M_{P_k})$ 是单节点训练的模型质量， $C_{\tau}(P_k)$ 是模型在 τ 次迭代过程中节点 P_k 的资源开销。由于个人理性考虑，参与者个人效用满足非负性，即 $U(P_k) \geq 0$ 。

由此，联邦学习中参与节点激励机制的目标函数可以定义为：

$$U^*(P_k) = \arg \max_{m_k \in \{m_k(t): t < T\}} U(P_k) \quad \#(8)$$

$$s. t. \forall P_k \in P, k \in (1, 2, \dots, N)$$

激励机制需要保证参与节点的一致性和活跃性^[44]。一致性指在模型训练过程中所有参与节点的共同贡献决定联合训练模型的质量，因此需要所有参与节点尽可能提供高质量的模型参数。活跃性通过激励策略吸引参与节点积极参与合作训练过程，使各节点趋于效用最大化目标。激励机制可以增加参与节点的活跃度，并且能根据节点贡献获得效用最大化，满足激励相容。

从整体福利角度考虑，激励机制的目标是达到帕累托最优的资源配置，每个参与节点的福利状态取决于节点在联邦学习中的参与贡献度。假设参与节点 P_k 在模型训练过程中的贡献度为 $C(P_k)$ ，模型收益之后所分配的奖励为 $R(P_k)$ ，奖励形式以平台积分形式发放，用于支付使用数据模型的费用。因此，所有参与节点的总体福利 $welfare(P)$ 满足：

$$\begin{aligned} \text{welfare}(P) &= \sum_{k=0}^N R(P_k) - C(P_k) \\ \text{s. t. } C(P_k) &= \theta \cdot C(P_k) \#(9) \end{aligned}$$

其中, 奖励 $R(P_k)$ 根据节点贡献度 $C(P_k)$ 分配, 贡献度越高的节点所获得奖励越多。节点贡献度 $C(P_k)$ 主要由三部分构成: 数据贡献度 $C_D(P_k)$ 、在线资源开销 $C_R(P_k)$ 以及模型质量贡献度 $C_Q(P_k)$ 。其中数据贡献度主要通过数据信息熵 $H_{ind}(P_k)$ 与数据量在总体数据中的占比 C_k 衡量, α 为权重参数。模型质量贡献度通过节点 P_k 提供的本地模型质量 $Q(M_{P_k})$ 判断, 因此节点 P_k 在一次模型训练中的贡献度可以定义为:

$$\begin{aligned} C(P_k) &= C_D(P_k) + C_R(P_k) + C_Q(P_k) \\ &= \alpha C_k + H_{ind}(d_k) + C_R(P_k) + Q(M_{P_k}) \quad (10) \end{aligned}$$

节点在线资源开销主要通过 CPU 资源消耗进行估算。假设本地模型训练过程中每次迭代输入数据样本量大小相同, CPU 周期频率为 f_k , 则节点 P_k 一次本地模型迭代的 CPU 资源消耗可以定义为^[54]:

$$C_R(P_k) = \zeta c_k s_k f_k^2 \#(11)$$

其中, ζ 是节点计算芯片组的有效电容参数^[55], c_k 是执行一个数据样本所需的 CPU 周期, s_k 是本地训练所需的数据样本大小。

节点诚实度 θ 可通过节点在历史任务中的表现判断, 结合主观逻辑模型^[56]通过节点参与历史任务中诚实行为、恶意行为以及恶意行为被修正的概率来预测节点在当前本地模型训练中的诚实度表现。假设参与节点 P_k 诚实表现的概率为 $Pr_c(P_k)$, 节点出现恶意行为被修正的概率为 $Pr_v(P_k)$, 参与节点不诚实的概率为 $Pr_{vc}(P_k) = Pr_v(P_k) \cdot (1 - Pr_c(P_k))$, 如果发现节点的不诚实行为由平台在用户押金中扣除罚金 f_p , 由此节点 P_k 在训练过程中表现的数值关系可描述为:

$$\begin{aligned} Pr_c(P_k) &= R(P_k)(1 - Pr_{vc}(P_k)) - f_p \cdot Pr_{vc}(P_k) \\ &\quad - C(P_k)(Pr_c(P_k) + Pr_v(P_k)) \quad (12) \end{aligned}$$

参与节点的诚实表现与训练过程中奖励与罚金相关, 其中诚实度 θ 、奖励 $R(P_k)$ 、罚金 f_p 在满足 $f_p/R(P_k) > (1 - Pr_{vc}(P_k))/Pr_{vc}(P_k)$ 条件下的数值关系可以通过定理 1 确定。

定理 1. 如果满足 $f_p/R(P_k) > (1 - Pr_{vc}(P_k))/Pr_{vc}(P_k)$, 其中 $Pr_{vc}(P_k) = Pr_v(P_k) \cdot (1 - \theta)$, 则参与者表现诚实的概率至多为 θ , 即诚实度为 θ 。

证明. 由 $Pr_c'(P_k) < \theta$, 不失一般性, 当 $\theta = 0$, 需证明 $Pr_c'(P_k) < 0$ 。

$$\begin{aligned} Pr_c'(P_k) &= R(P_k)(1 - Pr_{vc}'(P_k)) - f_p \cdot Pr_{vc}'(P_k) \\ &\quad - C(P_k) \cdot (Pr_c'(P_k) + Pr_v'(P_k)) \\ &< 0 \end{aligned}$$

由 $f_p/R(P_k) > (1 - Pr_{vc}'(P_k))/Pr_{vc}'(P_k)$

$$\begin{aligned} \text{则 } Pr_c'(P_k)/R(P_k) &= (1 - Pr_{vc}'(P_k)) - f_p/R(P_k) \cdot \\ &Pr_{vc}'(P_k) - C(P_k) \cdot (Pr_c'(P_k) + Pr_v'(P_k)) < \\ &-C(P_k) \cdot (Pr_c'(P_k) + Pr_v'(P_k)) \end{aligned}$$

由 $C(P_k)$ 、 $Pr_c'(P_k)$ 与 $Pr_v'(P_k)$ 的非负性, 可得:
 $Pr_c'(P_k) < 0$

证毕。

3.4 差分隐私

差分隐私 (Differential privacy, DP)^[57]可以在最大化查询准确率的情况下保护数据隐私, 一般通过增加随机化噪声^[58]避免攻击者获取原始数据。模型训练过程中的梯度数据在传输过程中可能泄露数据隐私, 差分隐私适用于深度学习的参数隐私保护^[59]。本地差分隐私 (Local differential privacy, LDP)^[60]可以实现本地训练过程的隐私保护。不同于传统的中心化差分隐私, LDP 关注数据收集过程的隐私保护, 免去了可信第三方假设, 并且可以抵御具备先验知识的攻击者, 因此 LDP 更适用于联邦学习中间参数隐私保护。差分隐私相关介绍如下:

定义 4. (差分隐私) 假设有两个相邻数据集 D 和 D' , 有且仅有一条数据不同, 即 $|D \Delta D'| \leq 1$, 输出服从某一分布的随机化算法 A 满足 ϵ 差分隐私, 则 D 和 D' 作为 A 的输入得到的输出分布 O 难以区分, 即:

$$\Pr\{A(D) = O\} \leq e^\epsilon \cdot \Pr\{A(D') = O\} \#(13)$$

其中, $\epsilon > 0$ 是差分隐私预算, 隐私预算越小, 则增加噪声越大, 隐私保护级别越高, 数据可用性越低。差分隐私满足两个算法组合特性, 分别是序列组合性^[61]和并行组合性^[62]。

定义 5. (序列组合性) 对于数据集 D 和 n 个随机化算法 $\{A_i\}, i \in [1, n]$, 如果 $A_i(D)$ 满足 ϵ_i 差分隐私, 则 $\{A_i\}$ 在 D 上的顺序序列组合满足 $\sum_i^n \epsilon_i$ 差分隐私。表明多个算法同时作用在同一数据集上时, 总体隐私预算是各隐私预算之和。

定义 6. (并行组合性) 将数据集 D 分成 n 个互不相交的集合 $\{D_i\}, i \in [1, n]$, 每个集合分别作用一个随机算法 $\{A_i\}, i \in [1, n]$, 如果 A_i 满足 ϵ_i 差分隐私, 则 $\{A_i\}$ 在 D 上的并行序列组合满足 $\max_i(\epsilon_i)$ 差分隐私。表明多个算法作用于一个数据集的不相交子集

上时，总体隐私预算是各隐私预算中的最大值。

定义 7. (松弛差分隐私) 为了增强差分隐私的实用性，Dwork^[63]在差分隐私定义中增加松弛项 δ ，采用较小的隐私预算得到更高的隐私保护级别：

$$\Pr\{A(D) = O\} \leq e^\epsilon \cdot \Pr\{A(D') = O\} + \delta \quad \#(14)$$

定义 8. (全局敏感度) 对于任意的查询函数 $f(D): D \rightarrow R^d$ ， $f(D)$ 的全局敏感度 Δf_{GS} 定义如下：

$$\Delta f_{GS} = \max_{D, D'} \|f(D) - f(D')\|_1 \quad \#(15)$$

$$s.t. \quad |D \Delta D'| \leq 1$$

全局敏感度是查询函数 f 在数据集 D 与相邻数据集 D' 查询时 $f(D)$ 与 $f(D')$ 之间最大曼哈顿距离，全局敏感度只和查询函数 f 相关，反应的是查询函数 f 在一对相邻数据集 D 和 D' 查询时最大的变化范围。 $f(D)$ 的全局敏感度越大，则在相同隐私预算 ϵ 情况下，参数增加的噪声也越大。

定义 9. (局部敏感度) 对于任意的查询函数 $f(D): D \rightarrow R^d$ ， $f(D)$ 的局部敏感度 Δf_{LS} 定义如下：

$$\Delta f_{LS} = \max_D \|f(D) - f(D')\|_1 \quad \#(16)$$

$$s.t. \quad |D \Delta D'| \leq 1$$

局部敏感度是查询函数 f 只在数据集 D' 查询时 $f(D)$ 与 $f(D')$ 之间最大曼哈顿距离，局部敏感度与数据集 D 相关，反应的是查询函数 f 在 D' 查询时最大的变化范围。

定义 10. (拉普拉斯机制) 对于任意函数 $f(D): D \rightarrow R^d$ ，隐私算法 M 满足 ϵ 差分隐私，则

$$M(D) = f(D) + \langle \text{Lap}\left(\mu, b = \frac{\Delta f}{\epsilon}\right) \rangle \quad \#(17)$$

其中， $\text{Lap}(\cdot)$ 是拉普拉斯噪声， μ ， b 分别是 Laplace 分布的位置参数和尺度参数。

4 可信联邦学习模型聚合与隐私保护

4.1 本地模型训练

传统分布式的机器学习需要将所有参与节点的本地数据收集在中心化服务器上执行模型训练。联邦学习可以将数据保留在参与节点本地进行分布式机器学习，从源头上避免了本地数据的隐私泄露。假设将联邦学习的 N 个参与节点定义为 $P = \{P_1, P_2, \dots, P_N\}$ ，对于单个节点 P_k ，其本地数据

标签为 $d_k = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ ，其中 x_i 是输入参数， y_i 是期望输出。联邦学习所有参与节点标签数据集定义为 $D = \cup d_i$ 。假设参与节点 P_k 本地训练的模型参数为 $m_k = \{m_1, m_2, \dots, m_m\}$ ，联邦学习模型训练的目标是得到全局训练模型 $M_G = h_m(x^i)$ ，使得与所有参与节点标签数据 D 的损失函数 $L(M_G)$ 最小。单节点 P_k 对于数据标签 d_k 损失函数定义为：

$$L(h_{m_k}(x^i)) = \frac{1}{|d_k|} \sum_{j \in d_k} f_j(h_{m_k}(x^i), y_j) \quad \#(18)$$

其中 $f_j(h_{m_k}(x^i), y_j)$ 是数据标签 (x_j, y_j) 基于模型 $h_{m_k}(x^i)$ 的损失函数。参与节点 P_k 的在 T 轮迭代过程中训练目标是在隐私保护条件下最优化本地模型 $h_{m_k}^*(x^i)$ 使得其损失函数最小，即：

$$h_{m_k}^*(x^i) = \arg \min_{m_k \in \{m_k(t); t < T\}} L(h_{m_k}(x^i)) \quad \#(19)$$

$$s.t. \quad \Pr(m_k \in R_d) \leq e^\epsilon \Pr(m'_k \in R_d)$$

$$\forall P_k \in P, k \in (1, 2, \dots, N)$$

其中 $m_k(t)$ 是在 t 轮联合训练中的参数集， T 是参数更新迭代次数的最大值。 $\Pr(m_k \in R_d) \leq e^\epsilon \Pr(m'_k \in R_d)$ 是更新参数 m_k 的差分隐私条件。

节点的本地训练采用随机梯度下降算法 (Stochastic Gradient Descent, SGD) 用于最小化损失函数。SGD 可以在目标函数相反的梯度方向更新参数以实现最小化损失函数，其中梯度的计算公式如下：

$$\nabla L(h_{m_k}(x^i)) = \frac{\partial L(h_{m_k}(x^i))}{\partial m_k} \quad \#(20)$$

对于节点 P_k 在 t 轮迭代过程中，模型参数更新可以定义为：

$$m_k(t) = m_k(t-1) + \alpha_t \cdot \nabla L(h_{m_k}(x^i)) \quad \#(21)$$

其中， α_t 代表向相反梯度移动的步长，即学习率 (learning rate)，通过向损失函数相反梯度方向移动从而不断逼近模型的最优结果。

4.2 PoC 共识算法

在联邦学习模型聚合过程中需要区块链节点的多方共识验证，然而现有的共识算法并不完全适用于此场景。工作量证明 (Proof of Work, PoW)^[36] 不仅会消耗大量节点算力，也不利于轻量级边缘节

点的参与。用户权益证明 (Proof of Stake, PoS)^[64] 中不在线节点也可以积累币龄, 可能导致参与节点的搭便车行为。本文基于联邦学习多节点协作训练, 提出一种适用于联邦学习模型聚合的共识算法——节点贡献度证明 (Proof of Contribution, PoC), 不仅降低了节点的计算开销, 而且提升了公平性。

PoC 共识算法结合节点在线时间、本地模型质量与数据贡献度三个方面来竞争区块链账本的记账权。利用节点在联邦学习协作中的贡献度达成共识, 不仅可以更好地利用区块链节点的计算和通信资源, 也可以激励高贡献度节点的参与, 保证合作训练的公平性。当接收到区块链记账请求 Req 时, 根据区块链网络中节点的贡献度从交易区块链中选择一个主节点。在一段时间内, 选择贡献度最大的节点作为主节点 P_{leader} 接收交易记账请求, 并负责将更新参数返回。所有区块链节点设置初始贡献度 OC , 通过选择性参与联邦学习更新贡献度。

在线时间作为节点初始贡献度 OC 评估的主要原因如下: 区块链系统的稳健性受到在线节点数量的影响, 所有在线节点共同参与区块链系统维护和链上数据的存储。本模型中节点的在线时间贡献度 (Online Contribution, OC) 计算公式如下:

$$OC = \beta \log(T_l - T_a - T_{off}) \quad \#(22)$$

其中 OC 是节点的在线贡献值; β 是在线时间系数, 控制在线时间贡献度的比例, T_l 是现有区块链的最后一个区块的时间戳; T_a 是每个节点第一次加入区块链网络的时间戳; T_{off} 是节点的下线时间间隔。

参与联邦学习模型质量评估的主要依据是本地训练过程中的交叉熵 (Cross Entropy)^[23], 用于衡量模型输出向量与真实结果间的偏差。假设联邦学习所有参与节点有一个用于验证模型质量的标签数据集 $d_u = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$, 则对于标签数据 (x_i, y_i) 的交叉熵 $H(f(x_i), y_i)$ 可通过下式计算:

$$H(f(x_i), y_i) = - \sum_{x_i} y_i \log f(x_i) \quad \#(23)$$

其中, y_i 代表标签数据期望输出, $f(x_i)$ 代表模型的预测结果。交叉熵 $H(f(x_i), y_i)$ 的值越小, 则代表模型预测的概率分布越接近真实结果, 即模型训练质量越高。

联邦学习中参与者 P_k 的数据贡献度主要由数据信息熵 $H_{ind}(d_k)$ 与数据量贡献度决定。数据量贡献度 C_k 的评估根据参与节点 P_k 与其他参与节点的数据量比值确定, 即:

$$C_k = \frac{|d_k|}{\sum_{i=1}^N |d_i|} \quad \#(24)$$

因此, 参与节点 P_k 本地训练模型的数据贡献度 C_D 可以定义为:

$$C_D = \alpha \frac{|d_k|}{\sum_{i=1}^N |d_i|} - \sum_{t_i \in T_{up}} p(t_i) \log_b p(t_i) \quad \#(25)$$

此外, 根据 Shannon^[52] 数据信息熵的非负性质, 共识过程中节点的贡献度证明 $proof_C$ 可通过下式计算 (其中 ρ 为调节系数):

$$proof_C = OC + \rho \frac{C_D}{H(f(x_i), y_i)} \quad \#(26)$$

共识过程开始之后, 联合训练组中贡献度最高的节点被选为主节点接受交易请求。节点进行本地模型训练, 将本地训练损失 $H(f(x_i), y_i)$ 广播给联合训练组中的其他节点, 并将交易 Tx 打包成区块 $B_k = \{ \langle H_{k-1}, H_M, t_i \rangle, Tx[m_i, H(f(x_i), y_i)], proof_C \}$, 其中 $\langle \rangle$ 中的内容是区块头, H_{k-1} 代表上一区块的哈希值, H_M 代表区块体 Merkle 根的哈希值, t_i 代表时间戳, Tx 代表交易内容, m_i 代表 P_k 本地模型训练参数, $H(f(x_i), y_i)$ 代表交叉熵, m_i 和 $H(f(x_i), y_i)$ 以交易形式记录在链上, $proof_C$ 代表节点贡献度证明, 根据节点贡献调整挖矿难度 D_m , 通过挖矿竞争记账权, 节点 P_k 的挖矿难度 D_m^k 调整方法为:

$$D_m^k = \frac{1}{proof_{C^k}} \quad \#(27)$$

节点 P_k 广播 B_k 给其他参与节点进行共识验证, 验证内容包括区块头、时间戳、数字签名, 并利用标签数据集 d_u 对模型 m_i 进行参数验证。

算法 1. PoC 共识算法.

输入: 共识节点 P_c , 验证节点 i , 超时时间 T

输出: 新区块 $\{Block_j\}$

FOR 参与节点 $P_i \in P_c$ & $timestamp \leq TDO$

IF $i = 0$ THEN

P_k 本地训练并广播 m_0 与 $H(f(x_i), y_i)$

计算初始贡献度证明 $proof_C = OC$

ELSE

P_k 本地训练并广播 m_i 与 $H(f(x_i), y_i)$

P_k 计算贡献度证明 $Proof_C$ 并广播 B_k

共识节点 P_j 验证区块 B_k

ENDIF

ENDFOR

挖矿难度赋值 $D_m = 1/Proof_C$, 竞争完成记账

RETURN{Block_j}

4.3 协作者异步参数审计

为了防止参与节点自身作恶上传虚假的训练损失，我们提出了协作者异步参数审计机制。在每轮迭代过程中，训练协作者对接收到的模型参数利用自己本地数据计算模型误差。经过异步参数审计之后，模型质量评估不单依靠节点自身计算的训练误差，还需考虑协作者对模型参数的审计结果，从而获取到更加真实公平的本地模型质量评估结果。

在联合训练过程中的协作者间的模型质量审计可以看作区块链上的一笔交易，交易Tx包括参与节点ID、模型参数 m_i ，训练损失 $L(m_i)$ 、时间戳，即， $Tx = \{ID, m_i, L(m_i), t\}$ 。为了保证交易的安全性，采用验证节点的公钥 PK_j 对交易进行加密，私钥 SK_j 对交易进行数字签名，组成加密的交易信息 $PK_j\{E(SK_j(Tx), PK_j)\}$ ，并发送给联合训练组中的其他节点进行验证，被激励的节点解密交易数据后利用自己的本地数据验证模型参数 m_i ，并将模型验证结果 $L_{P_j}(m_i)$ 以交易形式广播给联合训练组中的其他节点。所有节点在一定时间内计算接受到的验证结果并计算 P_k 本次联合训练的模型质量，经过共识验证后的质量验证结果 $L^u(P_k)$ 可通过下式计算：

$$L^u(P_k) = L_{P_k}(m_i) + \frac{1}{n-1} \sum_{j \neq k}^{n-1} L_{P_j}(m_i) \quad (28)$$

其中 $L_{P_k}(m_i)$ 是参与节点 P_k 本地训练模型的损失， $L_{P_j}(m_i)$ 是除参与者 P_k 之外的其他节点用自己本地数据对模型 m_i 计算得到的损失。

基于PoC共识算法，区块链主节点 P_k 广播 B_k 给其他参与节点 P_j 进行共识验证，并利用本地数据对模型参数 m_i 计算损失进行协作者参数审计，避免参与节点上传虚假模型损失，实现了可信的联邦学习。此外，区块链的链式存储结构可以实现交易追溯审计，从而约束节点的不诚实行为。协作节点参与模型验证的原因主要包括两方面：（1）协作者上传低质量的模型会影响到整体聚合模型质量（2）参与验证节点会受到平台积分激励，并增加竞争记账权的贡献度。

4.4 中间参数差分隐私保护

为了实现本地训练过程中的差分隐私，本文基

于LDP采用拉普拉斯机制为模型更新过程的中间参数增加噪声，可以形式化为：

$$\begin{aligned} \tilde{m}_k(t) = & m_k(t-1) + \alpha_t \cdot (\nabla L(h_{m_k}(x^i)) \\ & + 1/(1 + e^{-H_k(f(x_i), y_i)}) \cdot \\ & \langle \text{Lap}(s/\epsilon) \rangle) \end{aligned} \quad (29)$$

其中 $\langle \text{Lap}(s/\epsilon) \rangle$ 是拉普拉斯噪声， s 是局部敏感度， ϵ 是隐私预算， $H_k(f(x_i), y_i)$ 是模型质量评估的交叉熵，值越小则模型质量越高，需要减小噪声对模型参数造成的扰动。 $1/(1 + e^{-H_k(f(x_i), y_i)})$ 是动态隐私噪声调节系数，通过交叉熵控制噪声大小，经过归一化函数计算之后，控制噪声调节系数 $\gamma = 1/(1 + e^{-H_k(f(x_i), y_i)}) \in [0, 1]$ 。节点在每轮训练结束之后更新加噪之后的训练参数。隐私预算约束限定为 ϵ ，经过 T 轮迭代之后，总的预算约束 $\epsilon = \sum_{t=1}^T \epsilon_t$ ，在第 t 轮迭代过程中隐私预算 $\epsilon_t = \frac{\epsilon}{T}$ ， $1 \leq t \leq T$ ，如果累计的预算约束超过 ϵ ，则训练终止。

针对联邦学习模型聚合过程中的中间梯度增加基于拉普拉斯机制的差分隐私保护，降低了因推理攻击而导致的本地数据隐私泄露风险。每次模型迭代训练满足了 ϵ -差分隐私，根据本地差分隐私的定义，满足了数据收集过程的差分隐私保护。对于整体模型聚合过程中的隐私预算成本，可将隐私预算 ϵ 拆分至每轮迭代过程中，基于差分隐私的序列可组合性^[65]，隐私预算满足 $\epsilon_1 + \epsilon_2 + \dots + \epsilon_T \leq \epsilon$ ，因此聚合模型满足 ϵ -差分隐私。

分析本地差分隐私对中间梯度参数产生的误差影响，LDP通过噪声扰动实现隐私保护，噪声大小主要通过隐私预算和敏感度控制，隐私预算与噪声大小成反比，而敏感度与噪声大小成正比。因此在控制噪声误差时可以考虑动态调整隐私预算和噪声敏感度控制对参数造成的误差影响。目前已有的差分隐私噪声误差优化主要通过最小化查询结果的绝对误差^[66]。分析本地差分隐私产生的渐近误差边界，当数据项个数为 d ，LDP数据扰动带来的渐近误差边界与候选值 k ，隐私预算 ϵ ，参与者个数 n 有关，对于本地差分隐私中的代表性方法RAPPOR^[67]，其噪声扰动带来的渐进误差边界为 $O\left(\frac{dk}{\epsilon\sqrt{n}}\right)$ ，属于算法本身的固有误差^[68]。

4.5 自适应模型聚合算法

区块链节点收集到参与训练节点的本地模型之后进行模型聚合, 现有的模型聚合方法主要是联邦平均 (Federated Averaging, FedAvg) [69], 但联邦平均没有考虑到低质量模型参与聚合带来的全局模型质量降低问题。本文基于 FedAvg 提出了自适应模型聚合算法 (Federated Adaptive Averaging, FedAdp), 根据模型质量评估结果以及节点在协作过程中的信誉值评分调整模型聚合权重, 增加高质量模型参数在聚合模型中的贡献度占比, 从而提升聚合模型的准确率。

模型质量评估主要依靠节点 P_k 在第 i 轮本地训练的模型损失确定, 我们采用交叉熵 $H(f(x_i), y_i)$ 来评估本地模型损失。自适应模型聚合算法的质量评估权重 $Q_i(P_k)$ 可以定义为:

$$Q_i(P_k) = 1 - \frac{H_k(f(x_i), y_i)}{\sum_{j=1}^N H_j(f(x_i), y_i)} \quad \#(30)$$

为了避免参与节点提供虚假的模型质量评估结果, 本文采用协作者异步参数审计机制用于参数验证, 节点 P_k 广播模型参数 m_i 与交叉熵 $H(f(x_i), y_i)$, 协作者节点利用本地标签数据集 d_u 验证模型参数, 从而约束节点的自利性。

每次模型聚合各节点的质量评估权重都被记录在区块链上, 节点 P_k 在聚合过程中的信誉表现可以用一个累计的信誉分数 $S(P_k)$ 来表示:

$$S(P_k) = \sum_{i=1}^{\tau} \frac{1}{1 + e^{-\log(Q_i(P_k))}} / \tau \quad \#(31)$$

其中, τ 是节点 P_k 参与模型聚合的次数, $Q_i(P_k)$ 是节点 P_k 在参与第 i 次模型聚合中区块链上记录的质量评估权重, 我们将 $S(P_k)$ 作为节点 P_k 在联邦学习中的历史累计信誉值, 作为联邦学习模型聚合与参与节点选择的依据。

至此, 自适应模型聚合算法中全局模型更新计算过程可以表示为:

$$m_G(t) = m_G(t-1) + \frac{1}{N} \sum_{k=1}^N S(P_k) \cdot Q_i(P_k) \cdot \tilde{m}_k(t) \quad \#(32)$$

其中, N 是参与模型聚合节点个数, $Q_i(P_k)$ 是节点 P_k 在第 i 次模型聚合中本地模型的质量评估权重, $S(P_k)$ 是节点 P_k 的历史信誉值评分。

聚合模型 $m_G(t)$ 达到收敛的条件是当 $m_G(t)$ 的交叉熵小于预先设定的值 H_T 或迭代次数达到最大迭代次数阈值 T , 即:

$$M_G \leftarrow m_G(t) \quad \#(33)$$

$$s.t. H(m_i(t)) < H_T \mid t = T$$

算法 2. 自适应模型聚合智能合约。

输入: 模型聚合请求 Req , 参与节点 P , 迭代次数 i

输出: 全局训练模型 M_G

WHILE $H(m_G(t)) \geq H_T \& t < TDO$

FOR 参与节点 $p_i \in PDO$

P_k 本地训练模型并加噪 $\tilde{m}_k(t)$

P_k 计算模型损失 $L_{P_k}(\tilde{m}_k(t))$

广播模型参数 $\tilde{m}_k(t)$ 与损失 $L_{P_k}(\tilde{m}_k(t))$

协作节点 P_j 参数审计并广播 $L^u(P_k)$

协作节点 P_j 进行参数更新 $m_G(t)$

ENDFOR

协作节点 P_j 执行 PoC 共识并记账

ENDWHILE

RETURN M_G

5 实验分析及讨论

本文实验代码基于 Python (V3.6.10)、PyTorch (V0.4.1) 实现联邦学习协作式模型训练, 通过线程池模拟多节点分布式并行训练。实验测试数据集选用机器学习中被广泛用于图像识别任务的真实图片数据集 MNIST^[70] 和 CIFAR-10^[71]。MNIST 包含 60000 个训练样本和 10000 个测试样本, CIFAR-10 包含 50000 个训练样本和 10000 个测试样本。实验过程中通过代码自动化加载 MNIST 和 CIFAR-10 数据集提供训练数据, 预先根据参与节点数等分训练数据集, 模拟横向联邦学习过程。本实验采用随机梯度下降算法 (Stochastic Gradient Descent, SGD) 迭代优化本地模型, 神经网络模型选用多层感知机 (Multi-Layer Perceptron, MLP)^[72] 和卷积神经网络 (Convolutional Neural Networks, CNN)^[73] 进行全局模型的迭代训练。模型聚合算法采用本文所提出的自适应模型聚合算法实现全局模型的参数更新。实验环境为 Intel (R) Core (TM) i7-9700 CPU 3.00GHz 16GB RAM, 操作系统为 Windows10。实验过程中的参数设置见表 2。

表 2 实验参数设置

参数名称	默认取值
学习率 α_t	0.01
SGD 动量	0.5
Laplace 位置偏移 μ	0

数据批大小	64
本地模型 epochs	5
激活函数	ReLU

5.1 隐私保护可信联邦学习模型性能分析

本文联邦学习隐私保护方案通过给参与节点本地模型的中间梯度参数增加动态调整的拉普拉斯噪声实现，实验中调用 Numpy (V1.15.4) 随机噪声函数 (random.laplace) 产生拉普拉斯噪声，其位置偏移参数 $\mu = 0$ 。参与节点本地模型训练的学习率 $\alpha_t = 0.01$ ，动量参数 $Momentum = 0.5$ ，默认数据批大小 $batch\ size = 64$ 。实验过程中我们将数据集等分为 100 份，每个节点拥有 600 条 MNIST 数据集和 500 条 CIFAR-10 数据集，参与训练节点分别选择 10、30、60、100，依次对比 MLP 模型与 CNN 模型在不同数据集下的训练准确率与训练损失情况。

我们首先采用三层全连接的多层感知机模型 (MLP) 进行联邦学习模型训练，MLP 隐层包含 64 个神经元，激活函数为 ReLU。我们对比了 MLP 模型在 MNIST 数据集与 CIFAR-10 数据集的训练损失和模型准确率，并改变参与者数量观察训练损失与模型准确率变化情况，对比结果如图 3 所示。从 MLP 模型训练损失(MNIST 数据集)结果图中，我们发现随着全局模型聚合次数的增加，训练损失呈现

下降趋势，参与节点越多损失下降越多，当参与节点超过 30 个时，训练损失收敛于 0.5 以下；从 MLP 模型训练准确率(MNIST 数据集)结果图中，我们发现随着全局模型聚合次数的增加，训练准确率呈现上升趋势，10 个参与节点损失震荡明显，当参与节点超过 30 个时，模型准确率在 85% 以上，参与节点越多，准确率收敛越快也越平稳，MNIST 数据集收敛速度比 CIFAR-10 数据集更快，MNIST 数据集全局聚合 20 次以上模型准确率基本趋于收敛。MLP 模型在 CIFAR-10 数据集下的训练损失和模型准确率结果显示，模型损失与准确率趋势结果和 MNIST 数据集一致，但是 CIFAR-10 数据集下的模型损失稳定在 1.6 以下，100 个参与节点训练的模型准确率比 MNIST 数据集降低了 42.55%。

此外，我们采用卷积神经网络 (CNN) 测试训练效果，CNN 模型卷积核大小为 5×5 ，步长为 1，激活函数为 ReLU。从 CNN 模型结果图中显示 CNN 模型整体训练质量高于 MLP 模型，在 MNIST 数据集下，CNN 模型准确率比 MLP 高 2.69%，而在 CIFAR-10 数据集下，CNN 模型准确率比 MLP 高 7.43%。CNN 模型在 MNIST 数据集训练准确率同样明显高于 CIFAR-10 数据集，100 个参与节点训练的 CNN 模型，MNIST 数据集训练准确率比 CIFAR-10 数据集高 37.81%。

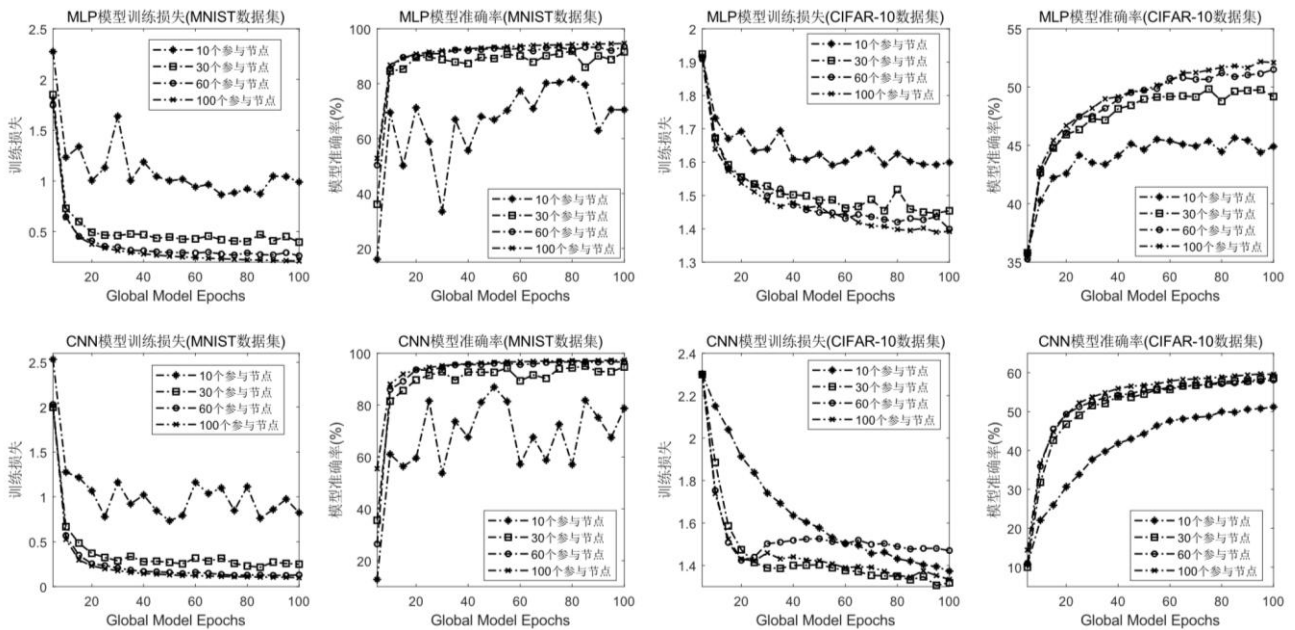


图 3 MLP 和 CNN 模型不同参与节点数量在 MNIST 数据集与 CIFAR-10 数据集训练损失与准确率对比

接下来我们进行差分隐私噪声误差分析。我们采用 MLP 模型，设置训练的节点数为 30 个，对比不同隐私预算与不同加噪方案对模型准确率的影响。我们首先采用传统的拉普拉斯噪声，设置隐私预算 ϵ 分别为 0.1、0.3、1、5。表 3 记录了中间参数加噪的模型训练准确率与不加噪时模型准确率

的平均误差，可以发现随着隐私预算的增加，MLP 模型准确率

表 3 MLP 模型差分隐私保护准确率平均误差

隐私预算	固定噪声平均误差 (%)	动态噪声平均误差 (%)
0.1	2.7031	2.6089
0.3	2.2842	2.1816
1	1.0578	0.4326
5	0.4675	0.0284

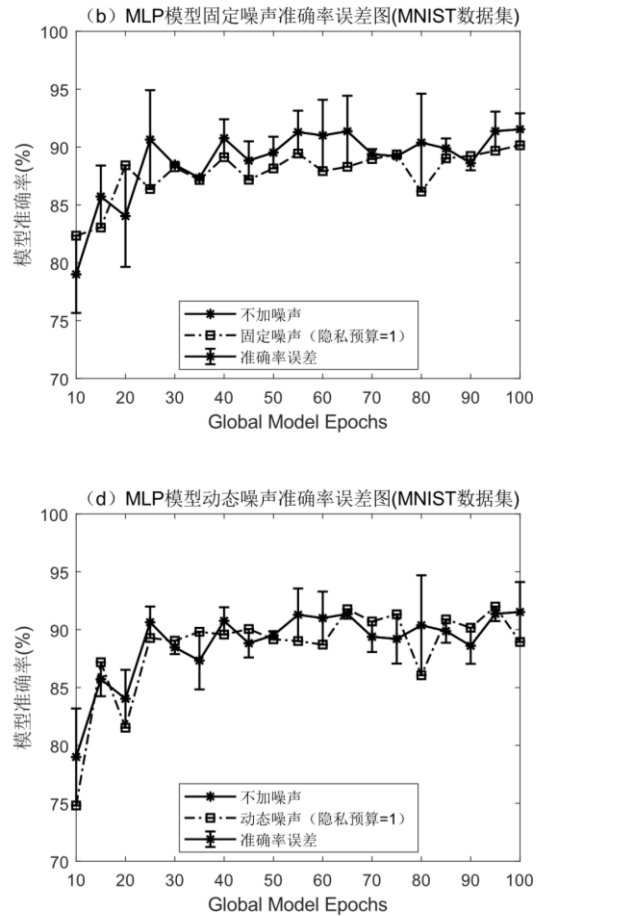
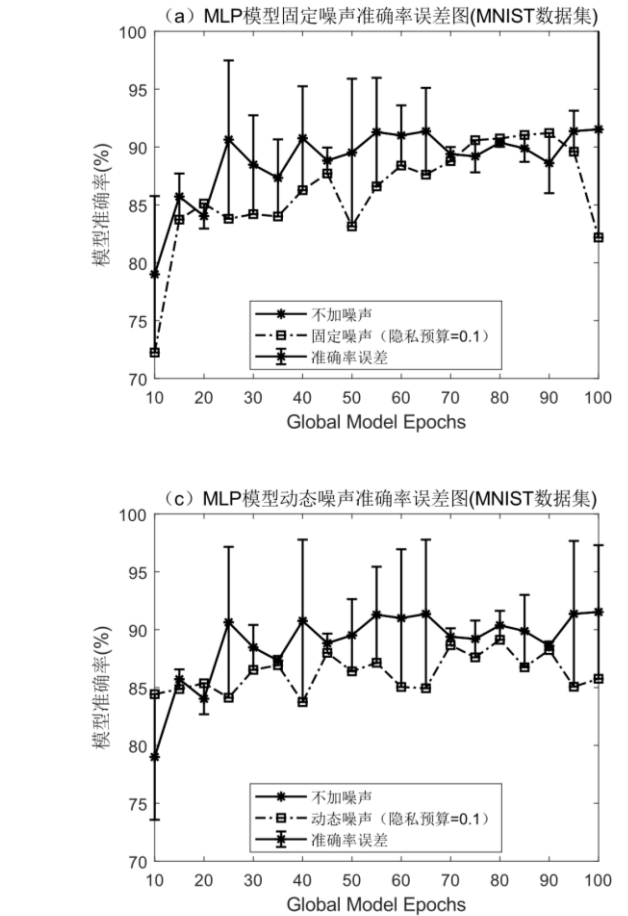


图 4 不同隐私保护方案在设置不同的初始隐私预算情况下对模型准确率的误差影响

此外，我们采用联邦学习中常见的投毒攻击 (Poisoning attack) 模拟恶意节点攻击行为。参考已有文献^[45]的做法，修改参与节点本地数据集标签 (labels) 模拟投毒攻击行为，攻击强度 $\rho \in [0,1]$ 是数据标签被修改的比例。在实验中我们采用 CNN

模型，数据集选用 MNIST 数据集训练模型，并将数据集等分为 100 份，一共包含 60 个参与节点，恶意节点数分别为 10 个和 20 个，攻击强度分别取 0.1、0.5、1。CNN 模型模拟投毒攻击的准确率平均误差如下表 4 所示，从表中我们可以看到随着攻击

强度增加，模型准确率平均误差增加。当恶意节点数增加时，模型准确率平均误差增加明显。

表 4 CNN 模型投毒攻击准确率平均误差

攻击强度	10 个恶意节点 (%)	20 个恶意节点 (%)
0.1	2.2716	8.2458
0.5	3.2468	8.6605
1	3.5547	10.0268

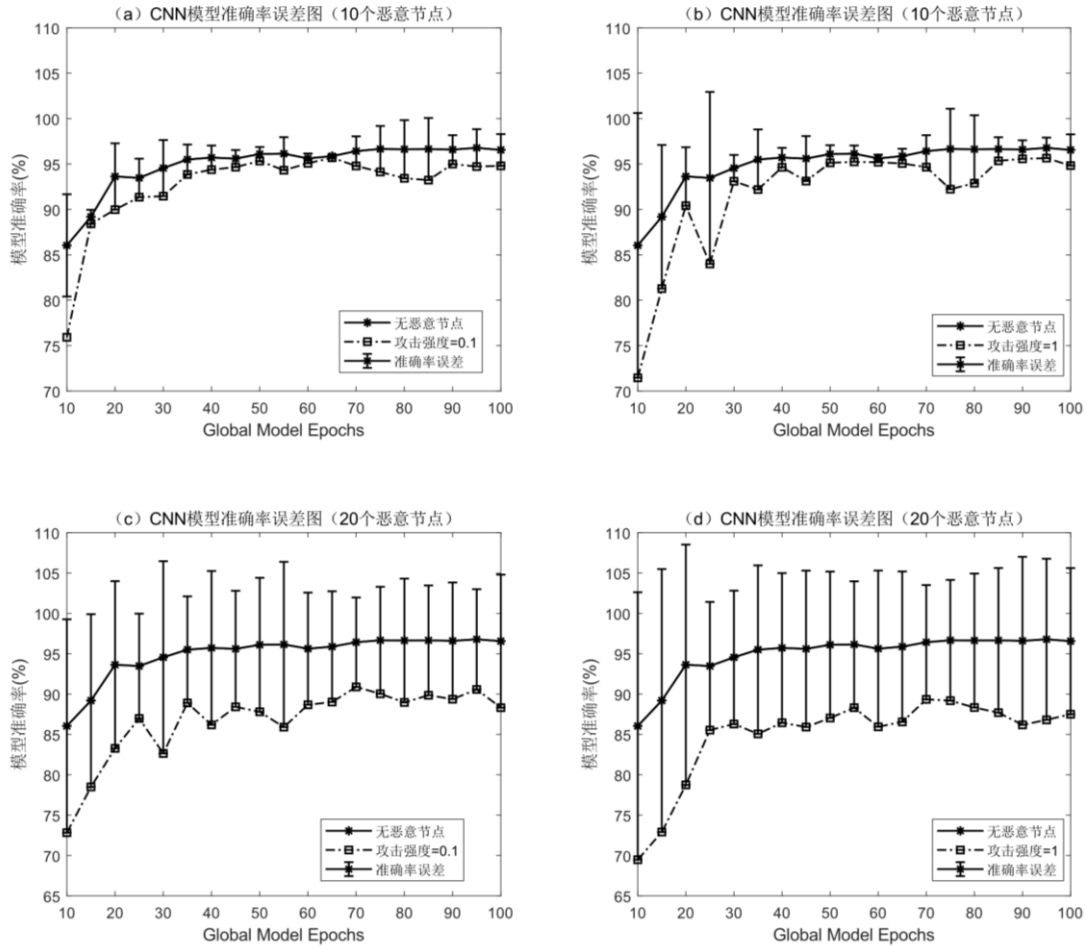


图 5 投毒攻击在不同攻击强度与不同恶意节点数量情况下对模型准确率的误差影响

我们进一步展示在抵御投毒攻击时，本文提出的自适应聚合算法与 FedAvg 聚合算法的模型准确率对比结果，如图 6 所示。我们分别对比不同恶意节点数量下，自适应聚合算法与 FedAvg 聚合算法的准确率变化情况。实验设置共有 60 个节点参与训练，攻击强度 $\rho = 0.1$ ，当有 10 个恶意节点参与时，自适应模型聚合算法准确率相较于 FedAvg 聚

我们对比了投毒攻击在不同攻击强度与不同恶意节点数量情况下对模型准确率的误差影响，如图 5 所示，从图中我们可以发现随着攻击强度与恶意节点数的增加，模型准确率误差增加明显。并且恶意节点数增加对模型准确率的误差影响大于数据标签修改时攻击强度的增加。

合算法在 100 次模型聚合中准确率平均提升了 2.778%。随着恶意节点增加至 20 个，准确率提升效果更加明显，准确率平均提升 4.851%。因此，本文提出的自适应模型聚合算法，由于在模型聚合过程中及时动态调整了高质量模型的聚合权重，因此相较于 FedAvg 聚合算法能更好的抵御恶意节点的投毒攻击。

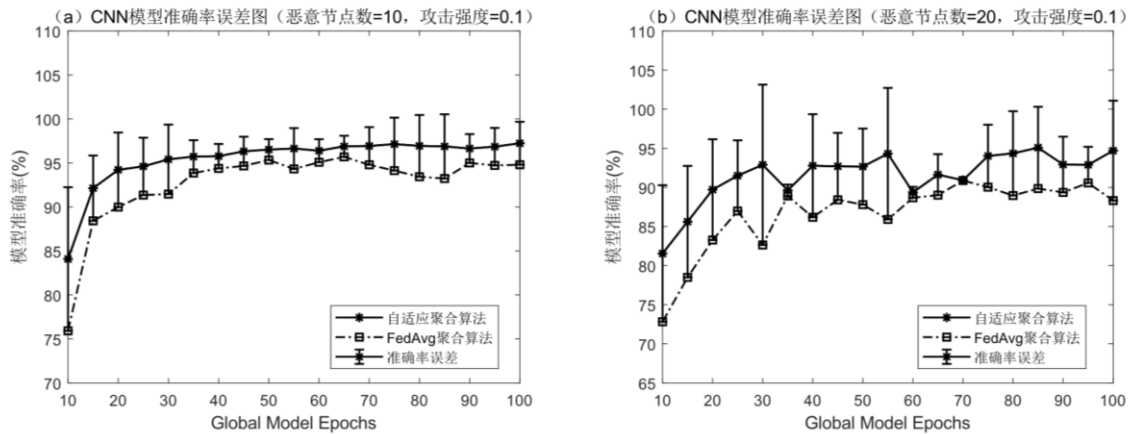


图 6 不同恶意节点数量情况下自适应聚合算法与 FedAvg 聚合算法的准确率对比结果

通过以上实验评估,我们发现联邦学习随着参与节点数量的增加,全局模型准确率有明显提升。采用本地差分隐私对中间梯度参数动态加噪后可以实现隐私保护的同时,降低对模型准确率的影响。通过模拟联邦学习投毒攻击,我们发现恶意节点的增加带来的模型准确率误差影响要大于修改数据标签时攻击强度增加带来的误差影响,并且本文提出的自适应模型聚合算法提升了联邦学习抵御投毒攻击的能力。尽管联邦学习增加了多节点中间参数通信时间,但多来源数据扩大了分布式机器学习的数据规模,并且免去传统分布式机器学习传输参与节点本地数据的通信和存储开销。

5.2 区块链模型聚合服务性能分析

这部分实验主要评估区块链模型聚合服务所需的平均存储开销和时延。本实验基于 Python Flask (V0.12.2) 框架实现区块链服务,对外提供创建交易 (new transactions)、链状态查询(chain)、节点注册(node register)、节点冲突解决 (resolve) 等 Web 服务。在搭建的区块链网络中,支持多节点或多端口的节点加入区块链网络。在区块链网络中,假设每个区块存储 10 次模型聚合交易,如果当前区块链长度|Blockchain|=100,即存在 1000 次历史聚合模型参数交易。

性能测试采用 Siege (V4.0.4) 高性能压力测试工具并结合 Postman 请求访问工具对区块链模型聚合服务进行存储开销和计算时延的测试。首先将模

型聚合服务以智能合约形式部署在区块链上,并对外提供 POST 访问功能。通过开放多端口模拟多节点参与的区块链网络,采用本文所提出的基于贡献度证明 (PoC) 共识机制竞争记账权,矿工挖矿难度根据参与节点的贡献度动态调整,贡献度高的节点挖矿难度低,最终以最长链为基准解决记账冲突。

我们首先分析每个区块存储的交易数量对区块链更新所需要的通信开销以及生成新区块所需的平均计算时延和存储开销的影响。在本实验的区块链网络中,区块中模型聚合交易数量从 10 增加至 100,模型聚合所需的平均计算时延和存储开销变化情况如图 7 (a) 所示,随着聚合交易量的增加,通信开销明显增加,各节点需要消耗更多的存储开销记录聚合模型参数,因此平均存储开销基本呈现线性增长,但是平均计算时延呈现波动增长趋势。此外,我们分析了当前区块链长度对模型聚合更新时所需的通信开销与平均计算时延。分别统计区块链长度从 10 增加至 100 的过程中存储开销与计算时延的变化情况,如图 7 (b) 所示。区块链存储开销呈线性增长趋势,平均计算时延波动增长受到挖矿难度随机性的影响。

随之,我们分析了挖矿难度对模型聚合性能的影响。通过参与节点的贡献度动态调整挖矿难度系数从 1 增加至 8 的过程中存储开销与计算时延的变化情况,如图 7 (c) 所示。区块链存储开销呈线性增长趋势,平均计算时延随挖矿难度成正向增长,

本课题得到国家重点研发计划(No.2017YFB1400700)、国家自然科学基金项目(No.62072487),北京市自然科学基金项目(No.M21036)资助。朱建明,博士,教授,主要研究领域为区块链技术与信息安全.E-mail: zjm@cufe.edu.cn. 张沁楠 (通信作者), 博士研究生, 主要研究领域为区块链与智能边缘计算.E-mail: zhangqnp@163.com.高胜 (通信作者), 博士, 副教授, 主要研究领域为区块链技术与信息安全.E-mail: sgao@mail.xidian.edu.cn.丁庆洋, 博士, 主要研究领域为区块链与大数据治理. E-mail: dingqingyang66@163.com.袁丽萍, 硕士研究生, 主要研究领域为区块链技术与隐私保护.E-mail: yuanliping_cufe@163.com.

当挖矿难度增加至 5 之后, 平均时延陡然增长, 当挖矿难度低于 5, 平均时延小于 200ms, 对区块链模型聚合服务性能影响较小。

我们发现通过对共识算法进行改进之后, 采用区块链为联邦学习提供模型聚合服务, 其性能可以实现与独立的联邦学习性能相当, 一定条件下区块

链的引入不会造成整体系统的性能瓶颈, 这与 Preuveneers D^[74]等人的研究成果结果一致。因此, 基于区块链分布式模型聚合计算引擎的联邦学习, 其性能在可接受范围内, 并不会因为引入区块链而对模型训练过程造成显著的性能影响, 方案具有一定的应用拓展性和实用性。

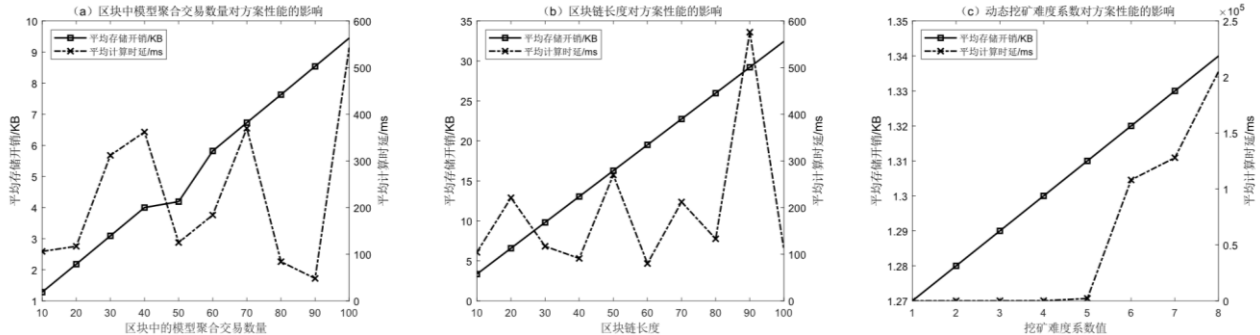


图 7 区块链模型聚合服务性能分析

5.3 方案性能对比与安全性分析

这部分实验我们对比了三种方案下的联邦学习模型准确率效果, 本文方案采用动态差分隐私实现中间梯度隐私保护, 通过区块链记录模型聚合参数及质量评估结果。以本地模型质量与历史信誉评分为依据实现自适应的模型聚合, 优化了传统的联邦平均模型聚合算法。最后通过区块链平台执行智能合约实现自动化模型聚合计算。通过与方案 1^[75]和方案 2^[43]的对比, 说明本方案的实用性。方案 1 重点关注到物联网边缘设备中联邦学习的差分隐私保护, 该方案采用边缘设备和云服务器进行联邦学习训练深度神经网络, 并通过拉普拉斯随机噪声对边缘设备传输到云服务器的数据进行扰动, 实现差分隐私保护。方案 2 是区块链与联邦学习结合的移动边缘计算案例, 采用基于信誉值的联邦学习参与者选择, 并通过在批处理规范层增加差分隐私噪声实现联邦学习隐私保护。三种方案准确率受数据批大小影响对比结果如下图 8 所示。从图中可以看出, 三种方案准确率差别不大, 随着批大小的增加模型准确率都有明显降低。这是因为批大小影响了节点本地模型训练次数, 批大小越小则训练次数越多, 从而可以获得更好的本地模型质量。此外, 通过减少批大小, 可以加快全局模型的收敛速度。

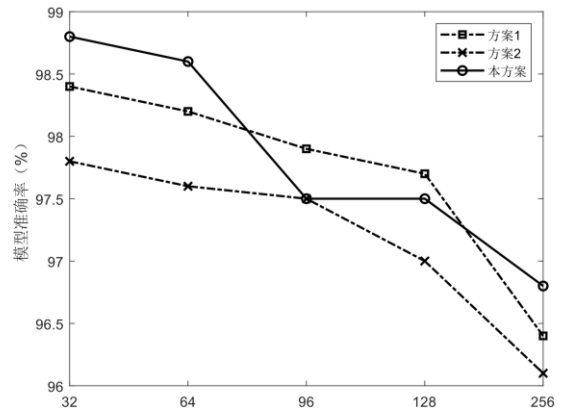


图 8 方案对比与批大小影响分析

综上所述, 区块链作为联邦学中的模型聚合计算引擎, 可以支持轻量级智能边缘节点参与到基于区块链的联邦学习模型训练中, 差分隐私保护可以实现中间参数的隐私保护, 该方案所需的计算时延和存储开销极为有限, 证明了本方案能以区块链为模型聚合引擎实现联邦学习中隐私保护, 并且能够高效执行分布式模型聚合算法, 具有较好的实用性。

分析本方案的安全性, 在联邦学习的数据声明部分采用伪公钥地址加密对数据声明信息进行匿名化加密, 保证数据的隐私和安全。设计激励机制约束节点的自利性, 通过协作者参数审计降低节点搭便车行为的风险。本文提出的 PoC 共识算法可以利用节点的在线时间、模型训练质量以及数据贡献度来争夺记账权, 不仅避免大量边缘节点计算资源的浪费, 还可以激励节点在联邦学习中有更高的贡献度。区块链作为联邦学习模型计算引擎, 可以有多方面的优势: 1) 去中心化架构避免单点失效风

险, 增强了系统的稳健性; 2) 多节点共识验证降低了跨设备联邦学习计算错误与投毒攻击的风险; 3) 区块链保存加噪后的中间参数与质量评估结果, 有利于节点行为的追溯审计与监管; 4) 聚合模型上链实现高效可信的模型数据共享。

6 结语

随着区块链技术的兴起, 为社会创建了可信的数据通信与存储基础设施, 基于联邦学习的边缘智能计算在物联网领域有广泛的应用前景。联邦学习将数据存储在与节点本地进行分布式机器学习, 有效保护数据隐私, 但是中心化的参数服务器仍存在恶意攻击与中间参数隐私泄露风险。因此, 本文利用区块链记录模型训练过程中的中间参数, 并激励协作节点进行参数审计验证。根据本地模型质量评估结果与节点历史信誉评分调整聚合模型的权重系数, 实现自适应的模型聚合。本文针对联邦学习场景提出了基于节点贡献度证明的共识机制, 根据模型贡献度调整挖矿难度系数, 降低边缘节点的计算资源开销。为了保护中间参数隐私, 本文根据模型质量加入不同程度的拉普拉斯噪声, 实现中间参数的隐私感知。最后通过原型搭建和实验仿真验证了本模型的实用性。

在接下来的工作中, 我们将从以下两方面开展研究: 1) 进一步量化评估联邦学习数据提供者的贡献度, 在即将到来的数据模型交易时代, 将交易收益合理回馈给数据提供方需要更深入的研究。2) 区块链现有的性能难以适应当前互联网时代的计算规模, 进一步提升区块链性能并探索更广阔的应用是区块链未来发展的重要方向。此外, 在保护用户隐私和安全的前提下, 联邦学习的训练效率及准确率仍是应用落地的一个瓶颈, 高效稳健的方案仍需更加深入的研究探索。

致谢*感谢对本文提出建议的所有评审专家。*

参考文献

[1] Konečný J, McMahan H B, Yu F X, et al. Federated learning: Strategies for improving communication efficiency. arXiv preprint arXiv:1610.05492, 2016.

[2] 微众银行 AI 项目组. 联邦学习白皮书 V2.0.2020. Artificial intelligence project team of Webank. Federated learning white

paper v2.0 . 2020 (in Chinese)

[3] Jiang J C, Kantarci B, Oktug S, et al. Federated Learning in Smart City Sensing: Challenges and Opportunities. *Sensors*, 2020, 20(21): 6230.

[4] Xu J, Wang F. Federated Learning for Healthcare Informatics. arXiv preprint arXiv:1911.06270, 2019.

[5] Niknam S, Dhillon H S, Reed J H. Federated learning for wireless communications: Motivation, opportunities, and challenges. *IEEE Communications Magazine*, 2020, 58(6): 46-51.

[6] Lim W Y B, Luong N C, Hoang D T, et al. Federated learning in mobile edge networks: A comprehensive survey. *IEEE Communications Surveys & Tutorials*, 2020, 22(3): 2031-2063.

[7] Li T, Sahu A K, Talwalkar A, et al. Federated learning: Challenges, methods, and future directions. *IEEE Signal Processing Magazine*, 2020, 37(3): 50-60.

[8] Yang Q, Liu Y, Chen T, et al. Federated machine learning: Concept and applications. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 2019, 10(2): 1-19.

[9] Liang Y, Guo Y, Gong Y, et al. An isolated data island benchmark suite for federated learning. arXiv preprint arXiv:2008.07257, 2020.

[10] Norberg P A, Horne D R, Horne D A. The privacy paradox: Personal information disclosure intentions versus behaviors. *Journal of consumer affairs*, 2007, 41(1): 100-126.

[11] P. Kairouz et al., "Advances and open problems in federated learning," *Foundations and Trends in Machine Learning*, 2021,4, (1):1-121.

[12] Liu Y, Kang Y, Xing C, et al. A secure federated transfer learning framework. *IEEE Intelligent Systems*, 2020, 35(4): 70-82.

[13] Kairouz P, McMahan H B, Avent B, et al. Advances and open problems in federated learning. arXiv preprint arXiv:1912.04977, 2019.

[14] Melis L, Song C, De Cristofaro E, et al. Exploiting unintended feature leakage in collaborative learning//*Proceedings of the 2019 IEEE Symposium on Security and Privacy (SP)*. IEEE, San Francisco, USA, 2019: 691-706.

[15] Yan X, Cui B, Xu Y, et al. A method of information protection for collaborative deep learning under gan model attack. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 2019.(Early Access)

[16] Hitaj B, Ateniese G, Perez-Cruz F. Deep models under the GAN: information leakage from collaborative deep learning//*Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*. Dallas, USA,2017: 603-618.

[17] Melis L, Song C, De Cristofaro E, et al. Inference attacks against collaborative learning. arXiv preprint arXiv:1805.04049, 2018, 13.

[18] Pyrgelis A, Troncoso C, De Cristofaro E. Knock knock, who's there? Membership inference on aggregate location data. arXiv preprint arXiv:1708.06145, 2017.

[19] Bagdasaryan E, Veit A, Hua Y, et al. How to backdoor federated learning//*Proceedings of the International Conference on Artificial Intelligence and Statistics*. Palermo, Italy, 2020: 2938-2948.

- [20] Orekondy T, Oh S J, Schiele B, et al. Understanding and controlling user linkability in decentralized learning. arXiv preprint arXiv:1805.05838, 2018.
- [21] Kang J, Xiong Z, Niyato D, et al. Reliable federated learning for mobile networks. *IEEE Wireless Communications*, 2020, 27(2): 72-80.
- [22] Geyer R C, Klein T, Nabi M. Differentially private federated learning: A client level perspective. arXiv preprint arXiv:1712.07557, 2017.
- [23] De Boer P T, Kroese D P, Mannor S, et al. A tutorial on the cross-entropy method. *Annals of operations research*, 2005, 134(1): 19-67.
- [24] Sarathy R, Muralidhar K. Evaluating Laplace noise addition to satisfy differential privacy for numeric data. *Transactions on Data Privacy*, 2011, 4(1): 1-17.
- [25] Kim H, Park J, Bennis M, et al. On-device federated learning via blockchain and its latency analysis. arXiv preprint arXiv:1808.03949, 2018.
- [26] Majeed U, Hong C S. FLchain: Federated learning via MEC-enabled blockchain network//Proceedings of the 2019 20th Asia-Pacific Network Operations and Management Symposium (APNOMS). Matsue, Japan, 2019: 1-4.
- [27] Hard A, Rao K, Mathews R, et al. Federated learning for mobile keyboard prediction. arXiv preprint arXiv:1811.03604, 2018.
- [28] Song C, Shmatikov V. Overlearning reveals sensitive attributes. arXiv preprint arXiv:1905.11742, 2019.
- [29] Dillenberger D N, Novotny P, Zhang Q, et al. Blockchain analytics and artificial intelligence. *IBM Journal of Research and Development*, 2019, 63(2/3): 5: 1-5: 14.
- [30] Shokri R, Shmatikov V. Privacy-preserving deep learning//Proceedings of the 22nd ACM SIGSAC conference on computer and communications security. Denver, USA, 2015: 1310-1321.
- [31] Choudhury O, Gkoulalas-Divanis A, Salonidis T, et al. A Syntactic Approach for Privacy-Preserving Federated Learning//ECAI 2020. IOS Press, 2020: 1762-1769.
- [32] Aono Y, Hayashi T, Wang L, et al. Privacy-preserving deep learning via additively homomorphic encryption. *IEEE Transactions on Information Forensics and Security*, 2017, 13(5): 1333-1345.
- [33] Bagdasaryan E, Poursaeed O, Shmatikov V. Differential privacy has disparate impact on model accuracy. *Advances in Neural Information Processing Systems*, 2019, 32: 15479-15488.
- [34] Truex S, Baracaldo N, Anwar A, et al. A hybrid approach to privacy-preserving federated learning//Proceedings of the 12th ACM Workshop on Artificial Intelligence and Security. London, UK, 2019: 1-11.
- [35] Xu R, Baracaldo N, Zhou Y, et al. Hybridalpha: An efficient approach for privacy-preserving federated learning//Proceedings of the 12th ACM Workshop on Artificial Intelligence and Security. London, UK, 2019: 13-23.
- [36] Nakamoto S, Bitcoin A. A peer-to-peer electronic cash system[J]. Bitcoin.—URL: <https://bitcoin.org/bitcoin.pdf>, 2008, 4.
- [37] Shao Q F, Jin C Q, Zhao Z, et al. Blockchain: Architecture and research progress. *Chinese Journal of Computers*, 2017, 40 (1) :1-21 (in Chinese)
(邵奇峰, 金澈清, 张召, 等. 区块链技术: 架构及进展. 计算机学报, 2018, 41(5): 969-988.)
- [38] Bao X, Su C, Xiong Y, et al. Flchain: A blockchain for auditable federated learning with trust and incentive//Proceedings of the 2019 5th International Conference on Big Data Computing and Communications (BIGCOM). QingDao, China, 2019: 151-159.
- [39] Awan S, Li F, Luo B, et al. Poster: A reliable and accountable privacy-preserving federated learning framework using the blockchain//Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security. London, UK, 2019: 2561-2563.
- [40] Ramanan P, Nakayama K, Sharma R. BAFFLE: Blockchain based aggregator free federated learning. arXiv preprint arXiv:1909.07452, 2019.
- [41] Qu Y, Gao L, Luan T H, et al. Decentralized privacy using blockchain-enabled federated learning in fog computing. *IEEE Internet of Things Journal*, 2020, 7(6): 5171-5183.
- [42] Kim H, Park J, Bennis M, et al. Blockchain-enabled federated learning. *IEEE Communications Letters*, 2019, 24(6): 1279-1283.
- [43] Zhao Y, Zhao J, Jiang L, et al. Mobile edge computing, blockchain and reputation-based crowdsourcing iot federated learning: A secure, decentralized and privacy-preserving system. arXiv preprint arXiv:1906.10893, 2019.
- [44] Weng J, Weng J, Zhang J, et al. Deepchain: Auditable and privacy-preserving deep learning with blockchain-based incentive. *IEEE Transactions on Dependable and Secure Computing*, 2019.(Early Access)
- [45] Kang J, Xiong Z, Niyato D, et al. Incentive mechanism for reliable federated learning: A joint optimization approach to combining reputation and contract theory. *IEEE Internet of Things Journal*, 2019, 6(6): 10700-10714.
- [46] Lu Y, Huang X, Dai Y, et al. Differentially private asynchronous federated learning for mobile edge computing in urban informatics. *IEEE Transactions on Industrial Informatics*, 2019, 16(3): 2134-2143.
- [47] Lu Y, Huang X, Dai Y, et al. Blockchain and federated learning for privacy-preserved data sharing in industrial IoT. *IEEE Transactions on Industrial Informatics*, 2019, 16(6): 4177-4186.
- [48] <https://blog.ml.cmu.edu/2019/11/12/federated-learning-challenges-methods-and-future-directions/>
- [49] Mathiyalahan S, Manivannan S, Nagasundaram M, et al. Data integrity verification using MPT (Merkle Patricia Tree) in cloud computing. *International Journal of Engineering Technology*, 2018, 7(2.24): 500-503.

- [50] Wang H, Kaplan Z, Niu D, et al. Optimizing federated learning on non-iid data with reinforcement learning//Proceedings of the IEEE INFOCOM 2020-IEEE Conference on Computer Communications. Toronto, Canada,2020: 1698-1707.
- [51] Mikolov T, Chen K, Corrado G, et al. Efficient estimation of word representations in vector space. arXiv preprint arXiv:1301.3781, 2013.
- [52] Shannon C E. A mathematical theory of communication. ACM SIGMOBILE mobile computing and communications review, 2001, 5(1): 3-55.
- [53] Zhang Y, Liu L, Gu Y, et al. Offloading in software defined network at edge with information asymmetry: A contract theoretical approach. Journal of Signal Processing Systems, 2016, 83(2): 241-253.
- [54] Tran N H, Bao W, Zomaya A, et al. Federated learning over wireless networks: Optimization model design and analysis//Proceedings of the IEEE INFOCOM 2019-IEEE Conference on Computer Communications. Paris, France, 2019: 1387-1395.
- [55] Kang J, Xiong Z, Niyato D, et al. Incentive design for efficient federated learning in mobile networks: A contract theory approach//Proceedings of the 2019 IEEE VTS Asia Pacific Wireless Communications Symposium (APWCS). Singapore, 2019: 1-5.
- [56] Kang J, Xiong Z, Niyato D, et al. Toward secure blockchain-enabled Internet of Vehicles: Optimizing consensus management using reputation and contract theory. IEEE Transactions on Vehicular Technology, 2019, 68(3): 2906-2920.
- [57] Dwork C. Differential privacy: A survey of results//Proceedings of the International conference on theory and applications of models of computation. Berlin, Germany:Springer, 2008: 1-19.
- [58] Dwork C, Kenthapadi K, McSherry F, et al. Our data, ourselves: Privacy via distributed noise generation//Proceedings of the Annual International Conference on the Theory and Applications of Cryptographic Techniques. Berlin,Germany:Springer, 2006: 486-503.
- [59] Abadi M, Chu A, Goodfellow I, et al. Deep learning with differential privacy//Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security. Vienna, Austria, 2016: 308-318.
- [60] Qin Z, Yu T, Yang Y, et al. Generating synthetic decentralized social graphs with local differential privacy//Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security. Dallas, USA,2017: 425-438.
- [61] McSherry F, Talwar K. Mechanism design via differential privacy //Proceedings of the 48th Annual IEEE Symposium on Foundations of Computer Science (FOCS'07). Rhode Island, USA, 2007: 94-103.
- [62] McSherry F D. Privacy integrated queries: an extensible platform for privacy-preserving data analysis//Proceedings of the 2009 ACM SIGMOD International Conference on Management of data.Rhode Island, USA,2009: 19-30.
- [63] Dwork, Cynthia, and Aaron Roth. The algorithmic foundations of differential privacy. Foundations and Trends® in Theoretical Computer Science .2014(2014):211-407.
- [64] Bentov I, Lee C, Mizrahi A, et al. Proof of activity: Extending bitcoin's proof of work via proof of stake [extended abstract] y. ACM SIGMETRICS Performance Evaluation Review, 2014, 42(3): 34-37.
- [65] Qin Z, Yang Y, Yu T, et al. Heavy hitter estimation over set-valued data with local differential privacy//Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security. Vienna, Austria,2016: 192-203.
- [66] Xiao X, Bender G, Hay M, et al. iReduct: Differential privacy with reduced relative errors//Proceedings of the 2011 ACM SIGMOD International Conference on Management of data. Athens, Greece,2011: 229-240.
- [67] Erlingsson Ú, Pihur V, Korolova A. Rappor: Randomized aggregatable privacy-preserving ordinal response//Proceedings of the 2014 ACM SIGSAC conference on computer and communications security. Scottsdale Arizona, USA,2014: 1054-1067.
- [68] Ye QQ, Meng XF, Zhu MJ, Huo Z. Survey on local differential privacy. Journal of Software, 2018,29(7):1981–2005 (in Chinese).
叶青青,孟小峰,朱敏杰,霍峥.本地化差分隐私研究综述.软件学报,2018,29(7):1981–2005.
- [69] Li X, Huang K, Yang W, et al. On the convergence of fedavg on non-iid data. arXiv preprint arXiv:1907.02189, 2019.
- [70] LeCun Y. The MNIST database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- [71] Krizhevsky A, Hinton G. Convolutional deep belief networks on cifar-10. Unpublished manuscript, 2010, 40(7): 1-9.
- [72] Pal S K, Mitra S. Multilayer perceptron, fuzzy sets, classification. IEEE Transactions on Neural Networks, 1992,3(5):683-697.
- [73] Krizhevsky A, Sutskever I, Hinton G E. Imagenet classification with deep convolutional neural networks. Communications of the ACM, 2017, 60(6): 84-90.
- [74] Preuveneers D, Rimmer V, Tsingenopoulos I, et al. Chained anomaly detection models for federated learning: an intrusion detection case study. Applied Sciences, 2018, 8(12): 26-63.
- [75] Jiang L, Lou X, Tan R, et al. Differentially Private Collaborative Learning for the IoT Edge//Proceedings of the International Conference on Embedded Wireless Systems and Networks. Beijing, China, 2019: 341-346.



Zhu Jianming, Ph.D., professor. His research interests include blockchain and information security.

Zhang Qinnan, Ph.D. candidate. Her research interests include blockchain and intelligent edge computing.

Gao Sheng, Ph.D., associate professor. His research interests

include blockchain and information security.

DingQingyang, Ph.D. His research interests include blockchain application and big data governance.

Background

With the development of mobile communication technology and intelligent edge devices, recent years have witnessed rapid development of Intelligent Edge Computing (IEC), where a large number of novel mobile applications integrated into our daily life, such as autonomous driving, intelligent diagnosis, smart cities and so on. In particular, federated learning (FL) is an effective way to reduce privacy disclosure risk in the process of data cooperation training, which is deeply integrated with the emerging technologies such as blockchain, cloud computing and internet of things.

However, existing cross-device federated learning still faces some challenges including privacy disclosure of intermediate parameters, malicious poisoning attack and so on. In cross-device federated learning, a centralized server is needed as a parameter server to perform model aggregation algorithm, which is also the owner of the global model. However, the centralized parameter server may exist single point failure and malicious attacks. Once the parameter server is broken, the attacker can obtain and tamper with the intermediate parameters, that may decrease the quality of aggregation model. On the other hand, it's hard to avoid a lot of remote data communication between participating nodes and parameter servers, which leads to the application limitation of federated learning.

In view of above challenges, we are motivated to construct a privacy preserving and trustworthy federated learning model based on blockchain. Existing work has adopted blockchain to reconstruct the federated learning node architecture. However, there are still privacy disclosure of intermediate parameters and trust issues among nodes, which may lead to low quality aggregation models and even affect model aggregation process. Moreover, the data redundancy and communication overhead of blockchain are not suitable for the participation of lightweight edge nodes. In this context, we constructed the blockchain and federated

YuanLiping, M.S. candidate. Her research interests include blockchain and privacy protection.

learning hierarchically, which is suitable for the lightweight edge nodes in the intelligent edge computing scene. In order to improve the quality and credibility of aggregation model, we proposed federated adaptive model aggregation algorithm, which adopted model cross entropy and historical reputation value as the metric to adjust aggregation weight. According to the model quality, Laplace random noise is dynamically adjusted to achieve the trade-off between privacy protection and model quality error. Moreover, proof of contribution (PoC) consensus algorithm is proposed to reduce computing resource overhead and incentive higher contribution of nodes.

This simulation results show that the proposed federated adaptive model aggregation algorithm can achieve higher accuracy of aggregation model when occur poisoning attack. By dynamically adjusting the Laplace random noise, the accuracy error of the aggregation model is reduced. The experiment of blockchain performance show that the computational cost, communication cost and storage cost are small, which confirmed that our scheme has good practicability. The introduction of blockchain will not reduce the credibility of federated learning process under certain conditions, and it will not cause the performance bottleneck of the whole system.

This work was supported by Nature Key Research and Development Program of China (No.2017YFB1400700), the National Natural Science Foundation of China (No.62072487), and Beijing Municipal Natural Science Foundation (No.M21036).

Blockchain, information security, and intelligent edge computing are main research topics of our research group. In the past few years, we have worked out many research papers that have been published in Chinese Journal of Computers, Journal of Software and some international journals.