

面向全局不平衡问题的基于贡献度的联邦学习方法

吴飞, 宋一波, 季一木, 胥熙, 王木森, 荆晓远

引用本文

吴飞, 宋一波, 季一木, 胥熙, 王木森, 荆晓远. 面向全局不平衡问题的基于贡献度的联邦学习方法[J]. 计算机科学, 2023, 50(12): 343-348.

WU Fei, SONG Yibo, JI Yimu, XU Xi, WANG Musen, JING Xiaoyuan. [Contribution-based Federated Learning Approach for Global Imbalanced Problem](#) [J]. Computer Science, 2023, 50(12): 343-348.

相似文章推荐 (请使用火狐或 IE 浏览器查看文章)

Similar articles recommended (Please use Firefox or IE to view the article)

[Transformer在计算机视觉场景下的研究综述](#)

Review of Transformer in Computer Vision

计算机科学, 2023, 50(12): 130-147. <https://doi.org/10.11896/jsjcx.221100076>

[一种面向多模态医疗数据的联邦学习隐私保护方法](#)

Federated Learning Privacy-preserving Approach for Multimodal Medical Data

计算机科学, 2023, 50(11A): 230800021-8. <https://doi.org/10.11896/jsjcx.230800021>

[一种基于CutMix的增强联邦学习框架](#)

Enhanced Federated Learning Frameworks Based on CutMix

计算机科学, 2023, 50(11A): 220800021-8. <https://doi.org/10.11896/jsjcx.220800021>

[聚类联邦学习簇间优化](#)

Inter-cluster Optimization for Cluster Federated Learning

计算机科学, 2023, 50(11A): 221000243-5. <https://doi.org/10.11896/jsjcx.221000243>

[基于任务关联特征解耦网络的无监督领域自适应图像分类](#)

Image Classification for Unsupervised Domain Adaptation Based on Task Relevant Feature Separation Network

计算机科学, 2023, 50(11A): 230100068-8. <https://doi.org/10.11896/jsjcx.230100068>

面向全局不平衡问题的基于贡献度的联邦学习方法

吴飞¹ 宋一波² 季一木² 胥熙² 王木森² 荆晓远³

1 南京邮电大学自动化学院 南京 210003

2 南京邮电大学计算机学院 南京 210003

3 武汉大学计算机学院 武汉 430072

摘要 联邦学习在保护各方数据隐私的前提下,协同多方共同训练,提高了全局模型的精度。数据的类不平衡问题是联邦学习范式中具有挑战的问题,联邦学习中的数据不平衡问题可分为局部数据不平衡和全局数据不平衡,目前针对全局数据不平衡问题的相关研究较少。文中提出了一种面向全局不平衡问题的基于贡献度的联邦学习方法(CGIFL)。首先,设计了一种基于贡献度的全局判别损失函数,用于调整本地训练过程中的模型优化方向,使模型在训练中给予全局少数类更多的关注,以提高模型的泛化能力;然后,在全局模型更新阶段,设计了一种基于贡献度的动态联邦汇聚策略,优化了各节点的参与权重,更好地平衡了全局模型的更新方向。在MNIST,CIFAR10和CIFAR100这3个数据集上进行实验,实验结果表明了CGIFL在解决全局数据不平衡问题上的有效性。

关键词 联邦学习;数据不平衡;多方协同;图像分类

中图法分类号 TP391

Contribution-based Federated Learning Approach for Global Imbalanced Problem

WU Fei¹, SONG Yibo², JI Yimu², XU Xi², WANG Musen² and JING Xiaoyuan³

1 College of Automation, Nanjing University of Posts and Telecommunications, Nanjing 210003, China

2 School of Computer Science, Nanjing University of Posts and Telecommunications, Nanjing 210003, China

3 School of Computer Science, Wuhan University, Wuhan 430072, China

Abstract Under the premise of protecting the data privacy, federated learning unites multiple parties to train together to improve the accuracy of the global model. Class imbalance of data is a challenging problem in the federated learning paradigm. Data imbalance in federated learning can be divided into local data imbalance and global data imbalance. At present, there are few researches on global data imbalance. This paper proposes a contribution-based federated learning approach for global imbalance problem (CGIFL). First, a contribution-based global discriminant loss is designed to adjust the model optimization direction in the local training process and make models give more attention to the global minority classes in training to improve the generalization ability of models. And a contribution-based dynamic federated aggregation algorithm is designed to optimize the participation weight of each node and better balance the updating direction of the global model. Experimental results on MNIST, CIFAR10 and CIFAR100 datasets demonstrate the effectiveness of CGIFL in solving the problem of global data imbalance.

Keywords Federated learning, Data imbalance, Multi-party coordination, Image classification

1 引言

随着近年来科技水平的不断提高,深度学习的相关研究得到了爆发式的增长,并在计算机视觉和自然语言处理等多个领域都取得了巨大的成果^[1]。深度神经网络模型取得优异性能的一个关键因素是海量优质的数据样本,如在图像分类

领域中广泛使用的 ImageNet 数据集。然而,在实践中,数据通常分散在不同的机构与组织中,集中管理的成本高、难度大、风险高,同时近年来隐私数据泄露事件频发以及相关数据保护法规的颁布与完善,传统的深度学习方案变得难以实施,各方不愿再将其私有数据发送到中心服务器进行模型训练。

为了解决使用多方数据的成本与安全问题,联邦学习^[2]

到稿日期:2022-11-11 返修日期:2023-02-17

基金项目:国家自然科学基金(62076139);之江实验室开放课题(2021KF0AB05);未来网络科研基金项目(FNSRFP-2021-YB-15);南京邮电大学1311人才计划

This work was supported by the National Natural Science Foundation of China (62076139), Open Research Project of Zhejiang Lab (2021KF0AB05), Future Network Scientific Research Fund Project (FNSRFP-2021-YB-15) and 1311 Talent Program of Nanjing University of Posts and Telecommunications.

通信作者:吴飞(wufei_8888@126.com)

(Federated Learning, FL)作为一种新的分布式机器学习架构被提出。联邦学习框架可以使多方机构与组织在无须交换本地隐私数据的前提下,完成联合训练机器学习模型的任务。联邦学习的流程主要分为本地模型训练与全局模型聚合。在本地模型训练流程中,各方按照约定模型结构,使用本地隐私数据进行训练得到本地模型;之后将本地模型交由服务器进行全局模型汇聚,得到性能均衡的全局模型。在训练过程中,不交换原始的隐私数据,只交换模型参数。目前联邦学习已经成为机器学习中的一个重要领域,吸引了许多研究学者,并被应用于医学成像、目标检测和地标分类等许多场景。

联邦学习在实际中面临一个关键性的挑战,即数据样本分布的不均衡。ACIFL^[3](Address Class Imbalance in Federated Learning)提出了数据分布的局部不平衡与全局不平衡的概念。局部不平衡,指各方数据分布呈现非独立同分布(non-IID)^[4]的特性,各本地数据集的各类样本分布不均衡,而全局不平衡指若将各方数据汇总,数据则整体呈现出不平衡的性质,例如长尾分布。这些情况可能会降低联邦学习的性能,本地数据的异质分布使得本地模型在不同的更新方向上远离全局目标,从而使全局模型在汇聚更新时远离全局目标,而整体数据的长尾分布则注定了全局模型的性能将会偏向于多数类样本^[5]。目前针对全局数据不平衡的研究工作较少,已有工作的性能仍有待提高。

本文从本地训练算法以及全局模型汇聚两个角度出发,提出了一种面向全局不平衡问题的基于贡献度的联邦学习方法(Contribution-based Federated Learning Approach for Global Imbalance Problem, CGIFL),动态调整本地模型优化方向和全局模型更新偏向。具体来说,CGIFL根据当前全局模型的性能,推测当前全局数据分布态势,并动态调整本地模型的优化方向,给予少数类样本更多的关注。在全局模型汇聚过程中,为了体现各本地模型的更新贡献的不同,本文依据各模型的性能计算每轮训练流程的各方更新贡献度,优化了更新汇聚流程。本文的具体贡献如下:

1)引入 PolyLoss^[6]作为基本的损失函数,用于学习样本特征, PolyLoss 的多项式系数可以灵活调整以适用于多种分类场景。此外,本文设计了一种基于贡献度的全局判别损失函数,根据模型判别性能和全局数据分布预测,调整模型的优化方向,提升对少数类样本的判别能力。

2)提出了一种基于贡献度的动态联邦汇聚策略,计算各节点模型对联邦学习任务的贡献价值,提高优秀模型对全局模型的影响力。

3)进行了充足的实验来评估本文方法的有效性。在3个广泛使用的图像数据集 MNIST^[7], CIFAR10^[8]和 CIFAR100^[8]上进行对比实验,验证了本文方法相比目前最新工作更具性能优势。

2 相关工作

2.1 传统不平衡学习

虽然机器学习领域中关于数据不平衡问题已经有许多学者进行了研究,解决方法主要分为两类,分别为数据层面和

算法层面。数据层面的方法通过平衡实际数据的分布情况来降低类不平衡对模型泛化性的影响,包括过采样、欠采样等方法。Bennin 等^[9]提出新的过采样方法,在少数类中生成尽可能多的不同的合成数据,平衡数据分布。Liu 等^[10]加大了对欠采样过程中被忽略的数据的关注,挖掘潜在的有效信息。而算法层面吸引着更多学者的研究,通过改变模型训练策略,来优化更新梯度,从而增强模型的泛化性。经典的 focal loss^[11]基于代价敏感的思想,提升模型对少数类样本和难分类样本的关注度。Zhong 等^[5]在两个训练流中分别训练多数类样本和少数类样本,为模型训练提供互补信息。Center invariant loss^[12]通过对齐每个恒等式的中心,使少数类与多数类的特征中心保持一致。

2.2 联邦不平衡学习

随着联邦学习越来越受到学者的关注,针对联邦学习中的数据不平衡问题也有了许多研究成果。在局部数据不平衡场景下, FedProx^[13]是通过降低欧几里得范数来限制局部更新的,减小了局部不平衡数据的影响。MOON^[14]引入了对比学习的思想,缩小了本地模型与全局模型的差距,提升了汇聚后的全局模型性能。Zhao 等^[15]提出了一种策略,通过在边缘设备之间全局共享少量数据来改进对 non-IID 数据的训练。在全局不平衡问题方面,目前相关研究工作则相对较少, ACIFL^[3]设计了一种监视器,用于检测全局样本的分布情况,提出 ratio loss 来动态调整训练权重,进而调整局部更新偏向, Fed-focal^[16]将 focal loss 应用于联邦学习中,试图提高全局模型对少数类样本的识别能力; Astraea^[17]增加了中介层,通过组合节点重新平衡训练。然而,目前这些方法在联邦学习中的全局数据不平衡场景下并未取得很好的效果,性能仍有待提高。

3 CGIFL 方法

3.1 CGIFL 的整体架构

假定共有 N 个参与节点 P_1, \dots, P_N 加入联邦学习任务中共同完成深度神经网络模型的训练工作。参与节点 P_n 所拥有的本地样本集为 $D_n = \{x_i^n\}_{i=1}^{M_n}$, 其中 M_n 为 D_n 的总样本数。假设共有 C 类标签, 则 D_n 可表示为 $D_n = \{D_n^1, \dots, D_n^C\}$, 其中 D_n^c 为数据集 D_n 中第 c 类标签的样本集合, D_n^c 的样本数为 M_n^c 。全局样本集合 D_{global} 由参与节点 P_1, \dots, P_N 的样本集组成, 表示为 $D_{\text{global}} = \{D_{\text{global}}^1, \dots, D_{\text{global}}^C\}$, 集合大小为 $M_{\text{global}} = M_{\text{global}}^1 + \dots + M_{\text{global}}^C$, 其中 D_{global}^c 为属于第 c 类的全局样本集合, 表示为 $D_{\text{global}}^c = \{D_1^c, \dots, D_N^c\}$, M_{global}^c 为属于第 c 类的全局样本集合的大小。全局不平衡问题可被认为是全局各类样本数 $M_{\text{global}}^1, \dots, M_{\text{global}}^C$ 之间存在不平衡的问题。

图1为CGIFL的总体框架,CGIFL总体分为模型训练与验证汇聚两部分。在本地训练过程中,节点 P_n 首先从联邦中心服务器中下载最新的全局模型,在此基础上使用数据集 D_n 训练得到本地模型。在验证汇聚过程中,所有训练节点将其本地模型上传至中心服务器中,由验证汇聚器进行性能验证并汇聚更新全局模型。本文方法首先计算各模型的贡献度,之后在本地训练过程中利用本文设计的基于贡献度的

全局判别损失函数,来提升本地模型对全局少数类的识别能力;在验证汇聚模块中,使用基于贡献度的联邦汇聚算法动态

调整各模型的汇聚权重,提升优质本地模型在全局更新中的影响力,以有效提升全局模型的性能。

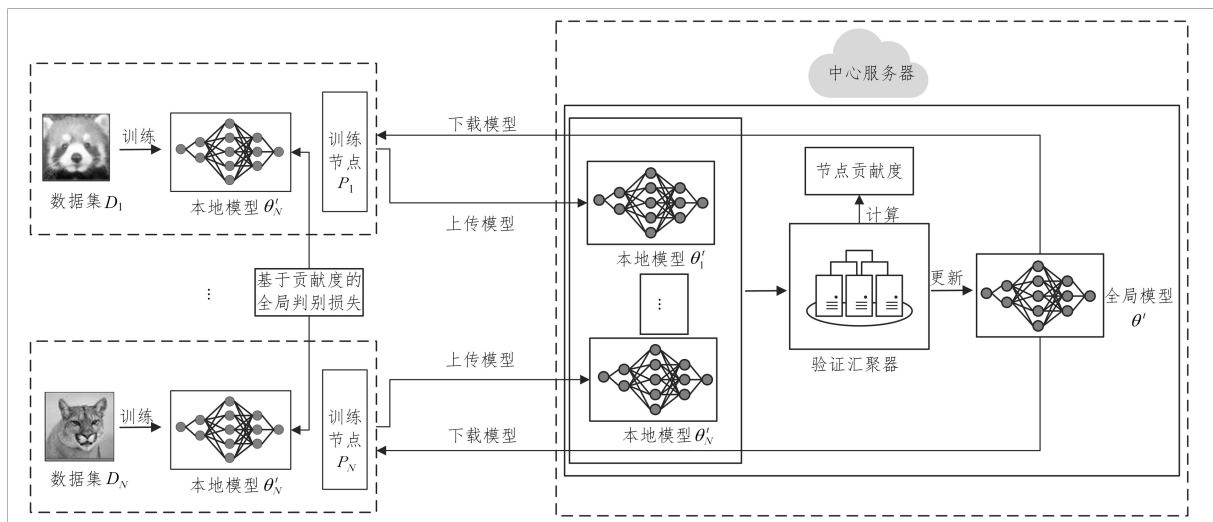


图1 CGIFL的总体框架

Fig.1 Overall framework of CGIFL

3.2 验证汇聚器

由于样本分布的不同,各本地模型的性能也不同,其对全局模型更新的贡献价值也是不同的,性能较差的本地模型对全局模型的更新价值较小,甚至在更新过程中会降低全局模型的性能。对此,本文设计了验证汇聚器,用于在更新全局模型前对所有训练节点的本地模型进行验证并评估其贡献度。

验证汇聚器部署在中心服务器中,验证汇聚器持有一个验证样本集 \$D^{val}\$, \$D^{val}\$ 由各节点提供的极少量样本组成, \$D^{val}\$ 的各类样本数是均衡的(本文方法所需的验证集极小,不会造成大量隐私数据泄露的问题)。对于训练轮次 \$t\$, 当全部训练节点将本地模型上传完毕后,验证汇聚器在验证集 \$D^{val}\$ 上对本地模型的性能进行验证,利用交叉熵损失函数计算得到各模型在验证时的损失值,记为 \$L' = \{L'_1, \dots, L'_N\}\$。使用各节点的验证损失计算各节点的更新贡献度,本文期望验证损失较小、性能较好的节点的贡献度相对更大。训练节点 \$P_n\$ 在第 \$t\$ 轮全局更新中的贡献度如下:

$$Con_n^t = \frac{1}{L'_n} \frac{1}{\sum_{i=1}^N \frac{1}{L'_i}} \quad (1)$$

完成各节点的更新贡献度计算之后,验证汇聚器对所有本地模型进行汇聚,更新全局模型。传统的 FedAvg 算法的更新式为:

$$\theta^t = \sum_{n=1}^N \frac{M_n}{M} \theta'_n \quad (2)$$

其中, \$\theta_n\$ 为节点 \$P_n\$ 在第 \$t\$ 轮全局训练的本地模型参数, \$\theta^t\$ 为第 \$t\$ 轮更新后的全局模型参数。

考虑到全局少数类样本可能集中分布在少部分节点中,且节点样本总数占全局样本总数的极小部分,而此类节点的更新权重显然是不能忽视的,因此使用基于节点样本数的联邦汇聚策略是不合适的。本文提出基于贡献度的联邦汇聚策略,使用节点贡献度作为更新权重,增大全局少数类样本敏感

的本地模型的影响力,更好地提升全局模型性能,具体定义如下:

$$\theta^t = \sum_{n=1}^N \frac{Con_n^t}{\sum_{n=1}^N Con_n^t} \theta'_n = \sum_{n=1}^N Con_n^t \theta'_n \quad (3)$$

3.3 基于贡献度的全局判别损失函数

联邦学习的目的是共同学习一个优质的全局模型。为了使全局模型具有更好的语义区分能力,每个本地模型也需要具备很好的语义区分能力。PolyLoss^[6] 可以通过调整多项式系数灵活地适用于不同任务,在二维图像分类、实例分割和目标检测等多个领域,相比传统的交叉熵损失表现都更加优秀。因此,本文使用 PolyLoss 作为基本的语义辨析损失函数。给定训练节点 \$P_n\$ 和样本数据集 \$D_n = \{x_i^n\}_{i=1}^{M_n}\$, PolyLoss 损失函数的定义如下:

$$L_{Poly-J} = \frac{1}{M_n} \sum_{i=1}^{M_n} (-\log(p_i^{n,t}) + \sum_{j=1}^J \epsilon_j (1 - p_i^{n,t})) \quad (4)$$

其中, \$p_i^{n,t}\$ 为样本 \$x_i^n\$ 由本地模型 \$\theta'_n\$ 进行正确预测的概率, \$\epsilon_j\$ 为多项式系数, \$J\$ 为多项式长度。参照文献[6], 本文将多项式长度设置为 1, 将多项式系数设置为 1, 使用 \$L_{Poly-1}\$ 作为训练过程中的损失主体, 定义简化为:

$$L_{Poly-1} = \frac{1}{M_n} \sum_{i=1}^{M_n} (-\log(p_i^{n,t}) + (1 - p_i^{n,t})) \quad (5)$$

各节点的更新贡献度体现了其本地模型的性能, 贡献度较小, 表明其本地模型对各类样本的判别能力较差。此外, 全局数据样本的分布情况并不透明公开, 本地样本分布与全局分布没有直接关联。受到 focal loss^[11] 的启发, 本文设计了基于贡献度的全局判别损失函数作为辅助损失函数, 基于节点贡献度与模型的样本识别能力动态地调整各类样本的损失权重, 平衡少数类和多数类样本对模型梯度更新的影响, 具体定义如下:

$$L_{Con} = \frac{1}{M_n} \sum_{i=1}^{M_n} (1 - Con_n^t) (1 - p_i^{n,t-1}) \log p_i^{n,t} \quad (6)$$

其中, \$p_i^{n,t-1}\$ 为样本 \$x_i^n\$ 由上一轮全局模型 \$\theta^{t-1}\$ 进行正确预测

的概率,若 θ^{t-1} 对该样本的预测成功概率较小,则表明当前全局模型对样本 x_i^t 的识别能力较差,该样本所属的标签类为全局少数类的可能性较大。 L_{con} 作为辅助损失函数,当模型对某样本的识别能力较差时,在训练过程中给予该样本更高的关注度,从而提升该样本类在训练中的影响力,增强模型对该类样本的判别能力。当节点贡献度较小时,表明该节点的性能较差,从节点层面出发,该节点的全部样本都应被给予更多的训练关注度,从而提升整体本地模型的性能。

结合两个损失函数 PolyLoss 与基于贡献度的全局判别损失函数,得到 CGIFL 的总损失函数为:

$$L_{CGIFL} = L_{Poly-1} + \alpha L_{con} \quad (7)$$

其中, α 为超参数,用于平衡两个损失。

算法 1 简要地总结了 CGIFL 的工作流程,本地训练模块以节点 P_n 为例。

算法 1 CGIFL 工作流程

输入:数据样本集 $D_n = \{x_i^n\}_{i=1}^{M_n}$, 对应标签集 $L_n = \{l_i\}_{i=1}^{M_n}$, 最新的全局模型 θ^{t-1}

输出:更新后的全局模型 θ^t

初始化:节点 P_n 获取最新的全局模型 θ^{t-1} , 本地训练迭代次数 S , 学习率 lr

本地训练流程:

1. $\theta_n^t \leftarrow \theta^{t-1}$
2. for $s=1, 2, \dots, S$
3. for x_i in $D_n = \{x_i^n\}_{i=1}^{M_n}$
4. $f_i \leftarrow \theta_n^t(x_i)$
5. $L \leftarrow L_{CGIFL}(l_i, f_i)$
6. $\nabla L_{CGIFL} \leftarrow L$
7. end for
8. $\theta_n^t \leftarrow \theta_n^t + lr \nabla L_{CGIFL}$
9. end for

全局模型更新流程:

10. 计算节点贡献度
11. $\theta^t \leftarrow \{\theta_i^t\}_{i=1}^N, \{Con_i^t\}_{i=1}^N$

4 实验验证

4.1 数据集介绍

本文采用了 3 个广泛使用的数据集: MNIST^[7], CIFAR10^[8] 和 CIFAR100^[8]。这 3 个数据集都是由有标签的图像构成的。

MNIST 数据集中的图像都是手写数字,具有 60000 个训练样本和 10000 个测试样本,每个样本图像的像素为 28×28 。CIFAR10 数据集共有 10 类样本,每类都有 6000 张图像,分为 5000 张训练图像和 1000 张测试图像,每张图片的尺寸都是 32×32 。CIFAR100 与 CIFAR10 具有相同的图像尺寸和相似的组成结构,它有 100 个类别,每类都包含 600 个图像,被分为 500 张训练图像和 100 张测试图像。

对于 3 个数据集,本文设置了不同神经网络来提取样本特征。参照文献[18-19],对于 MNIST 数据集,本文使用 LeNet5 提取手写数字样本的特征,对于 CIFAR10 和 CIFAR100,本文使用 ResNet18 来提取特征。

为了模拟全局样本分布不平衡的场景,本文固定选择标签(个数)为 2,4,7 的样本类为少数类,其余类别为多数类。本文首先给所有节点设置所拥有的样本标签类,每个节点所拥有的标签类数被随机设置为完整数据集标签类数的 40%~70%。对于多数类样本,节点从完整数据集中随机抽取 20% 对应类样本作为本节点的多数类样本,抽取比例为 20%。对于少数类样本,为了测试本文在不同的全局不平衡程度下的性能,设置了 3 种少数类样本的抽取比例,分别为 0.4%,1% 和 2%,即与多数类的比例为 1:50,1:20,1:10。

验证集 D^{val} 的每类样本数设置为 10,所有的训练节点从本地样本集中提供每类样本各 1 个并上传至中心服务器。由于各节点的样本类各不相同,此时上传至中心服务器中的各类样本数是不均衡的,因此对于不足的样本类进行上采样补足,对数量足够的样本类进行下采样,随机删除部分样本。

4.2 对比方法和评估指标

为了评估 CGIFL 的性能,本文选取了如下方法进行对比:1)传统的 Fedavg 方案^[20],使用交叉熵损失函数,作为基本的 baseline;2)基于全局样本分布监视器的解决方法 ACIFL^[3];3)基于 focal-loss 的联邦学习解决方法 Fed-focal^[16]。为了公平起见,所有方法都采用与本文实验相同的实验设置,本文运行文献[3,16,20]提供的代码或复现原论文的运行方法以获得实验结果。

本文使用 3 个数据集对应的测试集在最终的全局模型上进行验证,将完整测试集的验证准确度 Ac 值和少数类测试集的验证准确度 $Ac.M$ 值(Accuracy of Minority)作为评估标准。为了减小实验过程中随机性的影响,本文取 3 次运行结果的平均值作为最终的实验结果。

4.3 网络设置细节

本文使用 LeNet5 和 ResNet18 神经网络作为训练网络,并添加了一层全连接层作为输出层。3 个数据集对应的全连接层的维度分别为 10,10,100。Softmax 激活函数被添加在输出层后,用于计算标签预测的概率值。

本文将参与节点数设置为 20。在训练过程中,3 个数据集上的批量大小都设置为 32,本文使用 SGD 优化器来训练模型,学习率 lr 设置为 0.01,权重衰减设置为 0.00001,动量设置为 0.9。对于 MNIST, CIFAR10 和 CIFAR100,本文将全局轮次分别设为 20,50,50。对于所有的数据集,本文将局部轮次设置为 10。在训练过程中,通过网格搜索法调整式(7)中的超参数 α ,并获得最佳值 $\alpha=1$ 。

4.4 实验评估指标与实验结果分析

本文在 MNIST, CIFAR10 和 CIFAR100 数据集上的测试结果如表 1 和表 2 所列。由表 1 可见,在全局数据不平衡的场景下,本文方法在 3 种比例的 3 个数据集上的表现都优于 Fedavg, ACIFL 和 Fed-focal。具体来说,对于 Ac 指标,CGIFL 相比 Fedavg, ACIFL 以及 Fed-focal 在所有数据集上的平均性能分别提升了 1%,0.4% 和 2.1%。对于 $Ac.M$ 指标,CGIFL 的优势更加明显,如表 2 所列,相比基本的 Fedavg, ACIFL 以及 Fed-focal, CGIFL 的平均性能分别提升了约 2.9%,1.6% 和 3.9%。

表1 各方法在3个数据集上的Ac值

Table 1 Ac of each method on three datasets

Dataset	Fedavg	ACFIL	Fed-focal	CGIFL
MNIST(1:10)	98.64	98.79	98.64	98.70
MNIST(1:20)	98.06	98.23	98.00	98.32
MNIST(1:50)	97.17	97.26	97.08	97.47
CIFAR10(1:10)	82.98	83.95	80.27	84.97
CIFAR10(1:20)	79.68	80.58	76.83	81.15
CIFAR10(1:50)	74.03	74.76	71.54	76.61
CIFAR100(1:10)	46.84	47.37	46.22	47.54
CIFAR100(1:20)	44.37	45.63	43.80	45.54
CIFAR100(1:50)	42.81	43.15	42.06	43.34
Mean	73.84	74.41	72.71	74.85

(单位:%)

表2 各方法在3个数据集上的Ac.M值

Table 2 Ac.M of each method on three datasets

Dataset	Fedavg	ACFIL	Fed-focal	CGIFL
MNIST(1:10)	97.86	97.89	97.72	98.02
MNIST(1:20)	96.14	96.52	96.03	96.79
MNIST(1:50)	93.16	93.34	92.98	94.16
CIFAR10(1:10)	67.93	70.67	64.63	74.00
CIFAR10(1:20)	56.67	59.67	53.23	62.53
CIFAR10(1:50)	37.06	40.97	35.70	45.94
CIFAR100(1:10)	21.50	22.73	21.37	22.47
CIFAR100(1:20)	13.97	14.23	14.10	15.23
CIFAR100(1:50)	6.27	6.67	6.13	7.67
Mean	54.51	55.85	53.54	57.42

(单位:%)

4.5 消融实验结果

本小节评估了CGIFL中的组件的有效性,本文将没有使用PolyLoss的CGIFL版本称为CGIFL-pl,将没有使用基于贡献度的全局判别损失函数的CGIFL版本称为CGIFL-con,将没有使用基于贡献度的联邦汇聚策略的CGIFL版本称为CGIFL-cf。将这3种版本与完整的CGIFL方法在CIFAR10数据集上进行了实验对比,表3和表4列出了实验结果。

表3 CGIFL与其他版本在CIFAR10上的Ac值

Table 3 Ac of CGIFL and other versions on CIFAR10

Dataset	CGIFL-pl	CGIFL-con	CGIFL-cf	CGIFL
CIFAR10(1:10)	81.27	83.45	84.32	84.97
CIFAR10(1:20)	78.06	79.37	80.49	81.15
CIFAR10(1:50)	73.78	75.36	76.21	76.61

(单位:%)

表4 CGIFL与其他版本在CIFAR10上的Ac.M值

Table 4 Ac.M of CGIFL and other versions on CIFAR10

Dataset	CGIFL-pl	CGIFL-con	CGIFL-cf	CGIFL
CIFAR10(1:10)	68.53	71.37	73.23	74.00
CIFAR10(1:20)	55.61	60.29	61.73	62.53
CIFAR10(1:50)	41.15	44.64	45.26	45.94

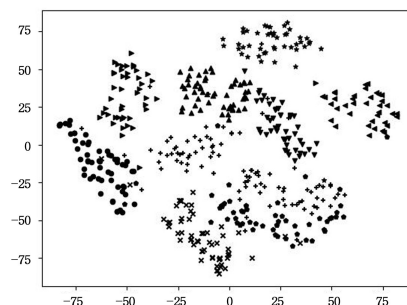
(单位:%)

由表3和表4的实验结果可以看出,CGIFL-pl,CGIFL-con和CGIFL-cf在3种比例的CIFAR10数据集上的表现都明显不如完整版本的CGIFL。这说明,本文选用的PolyLoss函数,以及本文设计的基于贡献度的全局判别损失函数和基于贡献度的联邦汇聚策略,都促进了不平衡学习的学习效果,

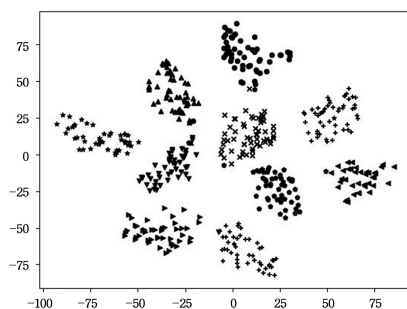
对解决联邦学习中全局数据不平衡问题很有效。

4.6 实验结果可视化分析

为了直观研究CGIFL在全局不平衡场景下对于少数类样本和多数类样本的区分能力,本文采用t分布随机邻域嵌入方法(t-Distributed Stochastic Neighbor Embedding, t-SNE),将图像样本的表示嵌入二维可视化平面中。本文以MNIST数据集为例,随机选取每类样本各50个,总计500个样本,将样本的原始特征和经过CGIFL处理后的特征嵌入二维平面中,可视化结果如图2所示,其中每种形状分别代表一种数字类别,共有10种。



(a) 原始图像特征投影



(b) 经过CGIFL处理后的特征投影

图2 图像特征可视化表示

Fig. 2 Visualization of image feature

图2(a)、图2(b)分别给出了原始图像特征和经过CGIFL处理后的图像特征的分布情况。可以看出,原始CIFAR10样本特征的分布整体呈现离散状,经过CGIFL处理后的同类图像特征在特征空间中被较好地聚合。这表明,在全局数据不平衡的联邦学习场景下,CGIFL可以有效地判别图像之间的差异,是一种有效的联邦学习方法。

结束语 本文提出了一种面向全局不平衡问题的基于贡献度的联邦学习方法,设计了一种基于贡献度的全局判别损失函数,在训练过程中提高了本地模型对全局少数类样本的判别能力,进一步提升了全局模型的性能;设计了基于贡献度的联邦汇聚策略,降低了性能较差的本地模型的影响,增强了全局模型汇聚更新的效果。实验结果表明,相比对比方法,CGIFL可以使模型具有更好的性能,提升了模型对少数类样本的判别能力,CGIFL有效地降低了全局数据不平衡问题造成的影响。

本文以联邦学习中的图像分类任务为主要研究场景,尚未探索联邦学习下的人脸识别、目标检测等应用场景。未来计划将CGLIF应用于多种联邦任务中,以研究CGLIF的泛化能力。

参 考 文 献

- [1] MOHAMMADI M, AL-FUQAHA A, SOROUR S, et al. Deep learning for IoT big data and streaming analytics: A survey [J]. *IEEE Communications Surveys & Tutorials*, 2018, 20(4): 2923-2960.
- [2] YANG Q, LIU Y, CHEN T, et al. Federated machine learning: Concept and applications [J]. *ACM Transactions on Intelligent Systems and Technology*, 2019, 10(2): 1-19.
- [3] WANG L, XU S, WANG X, et al. Addressing class imbalance in federated learning [C] // *Proceedings of the AAAI Conference on Artificial Intelligence*. 2021, 35(11): 10165-10173.
- [4] KAIROUZ P, MCMAHAN H B, AVENT B, et al. Advances and open problems in federated learning [J]. *Foundations and Trends in Machine Learning*, 2021, 14(1/2): 214-217.
- [5] ZHONG Y, DENG W, WANG M, et al. Unequal-training for deep face recognition with long-tailed noisy data [C] // *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2019: 7812-7821.
- [6] LENG Z, TAN M, LIU C, et al. PolyLoss: A Polynomial Expansion Perspective of Classification Loss Functions [J]. *arXiv: 2204.12511*, 2022.
- [7] DENG L. The mnist database of handwritten digit images for machine learning research [best of the web] [J]. *IEEE Signal Processing Magazine*, 2012, 29(6): 141-142.
- [8] KRIZHEVSKY A, HINTON G. Learning multiple layers of features from tiny images [R]. Technical report, University of Toronto, 2009.
- [9] BENNIN K E, KEUNG J, PHANNACHITTA P, et al. Mahakil: Diversity based oversampling approach to alleviate the class imbalance issue in software defect prediction [J]. *IEEE Transactions on Software Engineering*, 2017, 44(6): 534-550.
- [10] LIU X Y, WU J, ZHOU Z H. Exploratory undersampling for class-imbalance learning [J]. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 2008, 39(2): 539-550.
- [11] LIN T Y, GOYAL P, GIRSHICK R, et al. Focal loss for dense object detection [C] // *Proceedings of the IEEE International Conference on Computer Vision*. 2017: 2980-2988.
- [12] WU Y, LIU H, LI J, et al. Deep face recognition with center invariant loss [C] // *Proceedings of the on Thematic Workshops of ACM Multimedia 2017*. 2017: 408-414.
- [13] LI T, SAHU A K, ZAHEER M, et al. Federated optimization in heterogeneous networks [J]. *Proceedings of Machine Learning and Systems*, 2020, 2: 429-450.
- [14] LI Q, HE B, SONG D. Model-contrastive federated learning [C] // *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2021: 10713-10722.
- [15] ZHAO Y, LI M, LAI L, et al. Federated learning with non-iid data [J]. *arXiv: 1806.00582*, 2018.
- [16] SARKAR D, NARANG A, RAI S. Fed-focal loss for imbalanced data classification in Federated Learning [J]. *arXiv: 2011.06283*, 2020.
- [17] DUAN M, LIU D, CHEN X, et al. Self-balancing federated learning with global imbalanced data in mobile systems [J]. *IEEE Transactions on Parallel and Distributed Systems*, 2020, 32(1): 59-71.
- [18] LECUN Y, BOTTOU L, BENGIO Y, et al. Gradient-based learning applied to document recognition [J]. *Proceedings of the IEEE*, 1998, 86(11): 2278-2324.
- [19] HE K, ZHANG X, REN S, et al. Deep residual learning for image recognition [C] // *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2016: 770-778.
- [20] MCMAHAN H B, MOORE E, RAMAGE D, et al. Federated learning of deep networks using model averaging [J]. *arXiv: 1602.05629*, 2016.



WU Fei, born in 1989, Ph.D, associate professor, is a member of China Computer Federation. His main research interests include pattern recognition and machine learning.

(责任编辑:喻黎)