

联邦学习综述

牛远卓

(西北工业大学计算机科学与工程系西安 710072)

摘要 在隐私问题和深度学习愿景的驱动下,过去四年见证了机器学习适用机制的范式转变,一种称为联邦学习的新兴模式正在超越集中式系统和分布式现场学习。这是一种保护隐私的分散方法,它将原始数据保留在本机并训练本地模型,减轻了数据通信的负担,然后在中心服务器上执行学习和共享模型的联合,以在参与者之间聚合和共享构建的知识。本文首先检查和比较不同的基于机器学习的部署架构,然后对联邦学习进行深入和广泛的调查。与该领域现有的评论相比,我们在全面分析主要技术挑战和当前相关工作的基础上,对联邦学习课题和研究领域进行了新的分类。在此背景下,我们详细阐述了文献中各种具有挑战性的方面、贡献和趋势,包括核心系统模型和设计、应用领域、隐私和安全以及资源管理。此外,我们还讨论了一些重要的挑战和开放性的研究方向,以实现更强大的联邦学习系统。

关键词 人工智能,深度学习,分布智能,联邦学习应用,联邦学习,机器学习,隐私,资源管理,安全

中图法分类号 TP

DOI 号:

Federated Learning Survey

Yuanzhuo Niu

(Department of Computer Science, Northwestern Polytechnical University, 710072, China)

Driven by privacy concerns and the visions of deep learning, the last several years have witnessed a paradigm shift in the applicability mechanism of machine learning (ML). An emerging model, called federated learning (FL), is rising above both centralized systems and on-site analysis, to be a new fashioned design for ML implementation. It is a privacy-preserving decentralized approach, which keeps raw data on devices and involves local ML training while eliminating data communication overhead. A federation of the learned and shared models is then performed on a central server to aggregate and share the built knowledge among participants. This article starts by examining and comparing different ML-based deployment architectures, followed by in-depth and in-breadth investigation on FL. Compared to the existing reviews in the field, we provide in this survey a new classification of FL topics and research fields based on thorough analysis of the main technical challenges and current related work. In this context, we elaborate comprehensive taxonomies covering various challenging aspects, contributions, and trends in the literature, including core system models and designs, application areas, privacy and security, and resource management. Furthermore, we discuss important challenges and open research directions toward more robust FL systems.

Keywords Artificial intelligence, deep learning, distributed intelligence, federated learning applications, federated learning, machine learning, privacy, resource management, security

1 简介

应用场景

对于任何组织、任何人或世界上的任何人来说,数据都是最重要的东西。无论是个人还是组织,每个人都希望自己的数据不被泄露。但要训练机器学习算法,就需要数据,准确地说,需要大量高质量的数据。传统的机器学习方法是将数据保存在一台服务器中,然后训练模型。这种方法有可能造成个人数据泄露。

联邦学习 (FL) 是一种机器学习方法,它能让机器学习模型在位于不同站点的不同数据集上进行训练,而无需共享数据。它允许创建一个共享的全局模型,而无需将训练数据放在一个中心位置。它还允许个人数据保留在本地站点,降低了个人数据泄露的可能性。在联合学习中,不是数据向模型移动,而是模型向数据移动。这意味着训练是在终端设备上进行的。

联邦学习定义

目前联邦学习有很多种定义的方式,虽然说法各有不同,但是核心思想却是一样的。本文采用 2019 年 Google 发表的《Advances and Open Problems in Federated Learning》中的定义:联邦学习是一种机器学习范式,可以在一个中心服务器的协调下让多个客户端互相合作,即便在数据分散在客户端的情况下也可以得到一个完整的机器学习模型。

2 联邦学习分类

本节从数据划分、隐私机制、适用的机器学习模型、通信体系结构和解决异质问题的方法五个方面总结了联邦学习的分类。

2.1 数据划分

根据数据的样本空间和特征空间的不同分布模式,如下图所示,联邦学习可以分为三大类:横向联邦学习、纵向联邦学习和联邦迁移学习。

2.1.1 横向联邦学习

适用于两个数据集的用户特征重叠较多而用户重叠较少的情况。比如在两个不同的地区,两个当地电力局提供的服务类似,因此对用户产生了类似的数据,比如月用电量,因此两个电力局的数据集特征重叠较多。但是,这两个地区的用户分别受各自电力局的管辖,两个电力局的数据集用户重叠较少。

2.1.2 纵向联邦学习

适用于两个数据集用户特征重叠较少,但用户重叠较多的情况。例如同一个地区有两个机构:银

行和电力局。银行产生用户的收入和支出行为等数据,电力局产生用户用电数据,两个数据集的特征几乎没有重叠,但两个数据集中的用户却几乎是同一批用户(有较大重叠)。

2.1.3 联邦迁移学习

适用于两个数据集的用户特征和用户都重叠较少的情况。在这种情况下,我们不分割数据,但可以使用迁移学习克服数据或标签的缺乏。比如地区 1 的银行和地区 2 的电力局,用户重叠少(2 个地区),用户特征重叠也较少(银行数据和电力数据)。

2.2 隐私保护机制

联邦学习最重要的特点是各个客户端可以将自己的数据保存在本地,只是需要共享模型信息来训练目标模型。但是,模型信息也会泄露一些私有信息。保护联邦隐私的常用方法是模型聚合、同态加密和差分隐私。这些机制的论文可以在综述中找到。

2.2.1 模型聚合

模型聚合通过总结各方的模型参数来训练全局模型,从而避免在训练过程中传输原始数据。

此外,联邦学习和多任务的结合允许多个用户局部训练不同任务的模型,这也是典型的模型聚合方法。在一些论文中,联邦学习和区块链相结合,基于区块链交换和更新每个设备的模型数据,最后,在区块链协议的保证下,对模型参数进行安全聚合。

2.2.2 同态加密

一般的加密方案关注的是数据存储的安全性,没有密钥的用户不可能从加密结果中获得原始数据的任何信息,也不能对加密后的数据进行任何计算操作。与一般加密不同,同态加密最重要的特点是用户可以对加密后的数据进行计算和处理,但在处理过程中不会泄露原始数据,计算结束后,用户用密钥对处理后的数据进行解密。

在使用联邦学习时,用户与服务器之间的梯度交换可能会泄露用户的隐私信息。同态加密可以很好地解决这一问题,它可以在不影响模型训练结果的前提下对加密模型进行处理。

2.2.3 差分隐私

在差分隐私定义下,数据库的计算结果对特定记录的变化不敏感,数据集中是否存在一条记录对计算结果的影响很小。因此,向数据集添加一条记录所导致的隐私泄露风险被控制在一个非常小的可接受的范围内,攻击者无法通过加入或减少一个记录,观察计算结果来获得准确的个人信息。

在传统的机器学习和深度学习的模型训练过程中,流行的做法是在输出中加入噪声,然后在梯度

迭代的过程中应用差分隐私,从而达到保护用户隐私的目的。在实际应用中,通常采用拉普拉斯机制和指数机制来实现差异化隐私保护。其中拉普拉斯机制用于保护数值型的结果,指数机制用于保护离散型的结果。在差分隐私中有两种敏感度,全局敏感度和局部敏感度。前者由函数本身决定;后者由有数据集 \mathcal{D} 数据分布决定。然而,添加更多的噪声将不可避免地影响有效性,因此如何在隐私和有效性之间取得平衡是目前较为热门的研究方向。

2.3 可应用的机器学习模型

联邦学习逐渐渗透到当前流行的机器学习模型中,其目的是保证模型的私密性和效率。我们主要考虑联邦学习支持的三种模型:线性模型、决策树和神经网络。这些机制的论文可以在综述中找到。

2.3.1 线性模型

线性模型主要分为三类:线性回归、ridge 回归和 Lasso 回归。与其他模型相比,线性模型简单,易于实现,是实现联邦学习的有效模型。

2.3.2 树模型

联邦学习可用于训练单个或多个决策树,如梯度增强决策树和随机森林。梯度增强决策树(Gradient Boosting Decision Tree, GBDT)算法(将普通残差树转换成梯度下降)是近年来被广泛提及的一种算法,主要原因是它在许多分类和回归任务中都具有良好的性能。

2.3.3 神经网络模型

这模型当前在解决复杂问题时很流行,在联邦学习中也一样。

2.4 通信体系结构

联邦学习应用场景通常面临一下问题:客户数据分布不均、客户当地训练资源限制或不均,通信成本等问题。

1. 联邦平均 (FedAvg) 是联邦学习中最常用的模型优化方法。该方法对本地上传的梯度数据进行平均,然后更新并分发回本地。在多任务学习中,证明了 FedAvg 模型优化方法具有良好的性能。按我的理解, FedAvg 让客户在本地训练多轮,提高当地正确率的同时减少通信次数,减少通信成本。

2. 一些论文提出了一个名为 FedProx 的模型,该模型结合边缘设备数据进行分布式训练,并使用联邦平均模型优化方法来保证目标任务的鲁棒性和稳定性。按我的理解, FedProx 相比于 FedAvg 提高了对客户当地训练效果的容忍度,让他们按照自己设备所能承受的训练量来训练,而不是粗暴的剔除训练失败的模型。

3. 为了解决联邦学习中模型更新的通信代价过高的关键问题, Konecny 等人通过量化、随机旋转、二次采样等方法对模型数据进行压缩,以降低中心服务器与所有用户之间的通信压力。

4. Caldas 等人采用有损压缩和 Federated Dropout 来减少服务器到设备的通信。

5. Sattler 等人提出了一种稀疏三元压缩协议,该协议在对非独立同分布数据进行联邦训练时收敛速度比联邦平均算法快。

6. 为了解决 NonIID 数据的不平衡问题, Yang 等人提出了一种新的联邦平均算法,该算法通过计算不同设备的模型加权平均来聚合全局模型。

2.5 解决异质性的办法

在联邦学习的应用场景中,设备的差异会使整个训练过程的效率低下。为了解决系统异质问题,有四种方式:异步通信、设备采样、容错机制和模型异质。

2.5.1 异步通信

在传统的数据中心建设中,基于并行迭代优化算法有两种常见的方案:同步通信和异步通信。但是,面对设备的多样性,同步方案容易受到干扰,所以在联邦学习多设备环境中,异步通信方案可以更好地解决设备分散的问题。

2.5.2 设备采样

在联邦学习中,并不是每个设备都需要参与到每个迭代训练过程中。

1. Nishio 等人提出了一种新的协议 FedCS,解决了资源受限的客户端的选择问题,在训练过程中增加了更多的客户端,提高了模型的性能。

2. Kang 等人基于契约理论设计了一种激励机制,鼓励拥有高质量数据的本地设备积极参与有效的联邦学习过程,提高学习的准确性。

3. Qi 等人设计了一种基于联邦学习的新闻推荐模型,该模型也随机选择用户的局部梯度上传到服务器上,以训练全局模型。

4. Wang 等人提出了一种新的基于联邦学习的带局部补偿的 PRLC 方法来实现端到端通信。PRLC 的主要思想是:在每次迭代中,只有部分设备参与模型更新,没有参与的设备通过 PRLC 方法进行局部更新,以缩小与全局模型的差距。最后证明,在强凸性和非凸性情况下, PRLC 方法的收敛速度与未压缩方法相同,具有更好的可扩展性。

2.5.3 容错机制

在不稳定的网络环境中,容错机制可以防止系统崩溃。特别是在分布式环境中,当多台设备同时工作时,一旦其中一台设备出现故障,将会影响到

其他设备。

1.Wang 等人提出了一种确定局部更新和全局参数聚合的最佳权衡的控制算法,以适应设备资源的限制。

2.Yu 等人通过减少通信,改善了分布式随机梯度下降算法的线性加速特性。

3.也有一些研究直接忽略了设备的参与,在多任务学习中不影响联邦学习的效率。

4.容忍设备故障的另一种方法是通过编码计算引入算法冗余。移动设备上不正确的数据可能导致联邦学习欺诈。Kang 等人通过引入声誉作为度量,引入区块链作为声誉管理方案,提出了一种基于可靠员工选择的联邦学习方案,可以有效防止恶意攻击和篡改。

2.5.4 模型异质性

当从多方设备中收集分布不均匀的数据来训练联邦模型时,会严重影响模型的最终效率。合理处理来自不同设备的数据对联邦学习有着至关重要的影响。为了解决统计数据异质的问题,联邦学习网络主要分为三种建模方法:单个设备有自己的模型;培训适用于所有设备的全局模型;为任务训练相关的学习模式。

1.Yu 等人提出了一种仅使用正标签进行训练的通用框架,即 FedAwS (Federated Averaging with Spreadout)。在该框架中,服务器在每次迭代后添加一个几何正则化器,以促进类在嵌入空间中展开,这大大提高了训练效率,保证了分类任务的准确性。

2.Zhao 等人通过训练边缘设备之间的一小部分数据来建立全局模型,提高 Non-IID 数据的训练精度。

3.Khodak 等人在统计学习设置中设计并实现了一种自适应学习方法,提高了小样本学习和联邦学习的性能。

4.Eichner 等人考虑了全局模型与特定设备之间的快速数据自适应训练,以解决联邦训练时数据异质的问题。

5.Corinzia 等人提出了一种名为 VIRTUAL 的联邦学习算法,该算法将中央服务器和客户端的联邦网络视为贝叶斯网络,并采用近似变分推理在网络上进行训练,在联邦学习真实数据集上表现出了最先进的性能。

6.Liang 等人提出了一种将局部表示学习与全局模型联邦训练相结合的局部全局联邦平均 (Local Global Federated Averaging, LG-FEDAVG) 算法。理论分析表明,局部和全局模型的结合减少了数据的方差,减少了设备的方差,提高了模型在处理异质数据时的灵活性。实验表明, LG-FEDAVG 能够

降低通信成本,处理异质数据,有效学习模糊保护属性的公平表示。

3 联邦学习框架

以 FedAVG 为例,联邦学习框架如图 1 所示。假设固定的客户集,每个客户有固定的当地数据集。在每轮开始前,随机挑选固定比例的客户,然后服务器传输当前的联邦算法状态给这些客户。(由于实验表明增加客户数量达到一定程度会降低回报,所以只选择一定比例的客户)。每个被选择的客户根据联邦算法以及各自当地数据做固定次数的 SGD 当地训练,并将模型结果上传至服务器。最后,服务器根据这些模型结果更新联邦模型。整个过程持续重复。

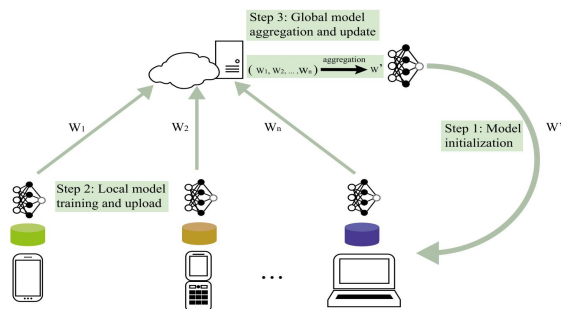


图 1 FedAVG 流程

4 挑战与未来研究方向

4.1 隐私保护

虽然目前有一些提高数据私密性的方法,但这些方法都增加了计算复杂度。为了进一步有效地保护私有数据的安全,我们需要寻找新的方法来防止私有数据在模型传输过程中被泄露。

4.2 通信成本

联邦网络可能由许多设备组成,例如数百万个远程移动设备,这意味着联邦学习模型的训练可能涉及大量的通信。此外,网络中的通信速度无法得到保证,因此联邦学习的通信代价是非常值得考虑的。

4.3 系统异质性

由于硬件和网络连接的不同,联邦网络中每个设备的计算和通信能力可能不同,网络中同时处于活动状态的设备通常只占很小的一部分。例如,数百万个设备的网络有时只能同时拥有数百个活动设备。同时,每个设备也可能是不可靠的,这些系统的异质性极大地加剧了容错的挑战。因此联邦学习

方法必须能够容忍异质硬件，并且对网络中的离线设备具有鲁棒性。

4.4 不可靠的模型上传

在联邦学习中，移动节点可能会有意或无意地误导服务器聚合全局模型。对于故意的行为，攻击者可以发送恶意的模型参数来影响全局模型的聚合，

从而造成模型训练的错误。另一方面，不稳定的移动网络环境可能会导致移动设备出现一些意想不到的行为，比如上传一些低质量的模型，这些都会对联邦学习产生不利影响。因此，抵抗这种不可靠的本地模型上传至关重要。