

联邦学习中的安全威胁与防御措施综述

陈学斌^{1,2,3}, 任志强^{1,2,3*}, 张宏扬^{1,2,3}

(1. 华北理工大学 理学院, 河北 唐山 063210; 2. 河北省数据科学与应用重点实验室(华北理工大学), 河北 唐山 063210;

3. 唐山市数据科学重点实验室(华北理工大学), 河北 唐山 063210)

(* 通信作者电子邮箱 psp1274632466@qq.com)

摘要: 联邦学习是一种用于解决机器学习中数据共享问题和隐私保护问题的分布式学习方法,旨在多方共同训练一个机器学习模型并保护数据的隐私;但是,联邦学习本身存在安全威胁,这使得联邦学习在实际应用中面临巨大的挑战,因此,分析联邦学习面临的攻击和相应的防御措施对联邦学习的发展和应用至关重要。首先,介绍联邦学习的定义、流程和分类,联邦学习中的攻击者模型;其次,从联邦学习系统的鲁棒性和隐私性两方面介绍可能遭受的攻击,并介绍不同攻击相应的防御措施,同时也指出防御方案的不足;最后,展望安全的联邦学习系统。

关键词: 联邦学习;隐私保护;攻击与防御;机器学习;鲁棒性与隐私性

中图分类号: TP309.2 **文献标志码:** A

Review on security threats and defense measures in federated learning

CHEN Xuebin^{1,2,3}, REN Zhiqiang^{1,2,3*}, ZHANG Hongyang^{1,2,3}

(1. College of Sciences, North China University of Science and Technology, Tangshan Hebei 063210, China;

2. Hebei Provincial Key Laboratory of Data Science and Application (North China University of Science and Technology),

Tangshan Hebei 063210, China;

3. Tangshan Data Science Key Laboratory (North China University of Science and Technology), Tangshan Hebei 063210, China)

Abstract: Federated learning is a distributed learning approach for solving the data sharing problem and privacy protection problem in machine learning, in which multiple parties jointly train a machine learning model and protect the privacy of data. However, there are security threats inherent in federated learning, which makes federated learning face great challenges in practical applications. Therefore, analyzing the attacks faced by federation learning and the corresponding defensive measures are crucial for the development and application of federation learning. First, the definition, process and classification of federated learning were introduced, and the attacker model in federated learning was introduced. Then, the possible attacks in terms of both robustness and privacy of federated learning systems were introduced, and the corresponding defense measures were introduced as well. Furthermore, the shortcomings of the defense schemes were also pointed out. Finally, a secure federated learning system was envisioned.

Key words: federated learning; privacy protection; attack and defense; machine learning; robustness and privacy

0 引言

机器学习通过经验自动改进计算机的能力已经得到广泛应用^[1],但是数据的孤岛现象^[2]和个人隐私保护的要求对传统中心化的机器学习方法提出了挑战。为了解决这些问题,联邦学习^[3]作为一种新兴的技术,被用于打破数据孤岛的同时保护数据隐私。联邦学习的概念最早是在2016年由谷歌提出,本质上是一种分布式机器学习技术、一种机器学习框架,能够让多个数据拥有方共同参与机器学习过程,而不必把数据集中在一个地方。这种方法可以在不暴露数据的情况下训练模型,并且可以实现多方共同受益。此外,联邦学习也符合欧盟通用数据保护条例(General Data Protection Regulation, GDPR)^[4]等隐私保护法规的要求,它可以确保个人数据不出边界,并且只有在得到用户允许的情况

下才进行模型训练;因此,联邦学习有望成为未来机器学习领域的重要技术之一。

联邦学习作为一种分布式机器学习技术,通过在本地设备训练模型,再将模型参数上传至中央服务器进行聚合,以实现在保护用户数据隐私的前提下,利用大量分散数据训练更准确的模型。然而,联邦学习本身也存在一些安全问题:

1)数据隐私泄露风险。在联邦学习中,攻击者可能通过模型训练过程中的模型参数推断本地设备上的数据信息,从而泄露用户隐私。

2)恶意攻击风险。在联邦学习中,攻击者故意上传恶意模型参数,导致聚合后的模型出现错误,从而破坏模型的准确性;或者攻击者通过篡改数据和模型参数影响模型的训练和聚合过程,从而影响模型。

3)模型逆向工程风险。攻击者通过逆向工程可能分析

收稿日期:2023-07-04;修回日期:2023-07-15;录用日期:2023-07-25。 基金项目:国家自然科学基金资助项目(U20A20179)。

作者简介:陈学斌(1970—),男,河北唐山人,教授,博士,CCF杰出会员,主要研究方向:大数据安全、物联网安全、网络安全;任志强(2000—),男,四川广元人,硕士研究生,CCF会员,主要研究方向:数据安全、隐私保护;张宏扬(1999—),男,江苏淮安人,硕士研究生,主要研究方向:数据安全、隐私保护。

出联邦学习模型的参数,获取模型的机密信息。

本文介绍了联邦学习系统存在的安全威胁,主要分为两个方面:对联邦学习系统鲁棒性的威胁和对联邦学习系统隐私性的威胁(见表1)。对联邦学习系统鲁棒性的威胁包括数据投毒^[5-8]、模型投毒^[9-13]和后门攻击^[12,14-19]。为了讨论方便,本文将参与方训练的模型称为局部模型,将服务端的模型称为全局模型。数据投毒是指攻击者恶意篡改本地数据,影响全局模型的训练和性能;模型投毒是指攻击者篡改全局模型的更新或训练过程,从而影响模型的性能和可靠性;后门攻击是指攻击者在模型中植入后门,从而可以随时控制模

型的行为和结果。针对这些鲁棒性威胁,本文介绍了鲁棒性威胁的防御措施^[7,9-11,15-17,20-24],这些防御措施主要通过检查更新的可信性进行防御。对联邦学习系统隐私性的威胁主要涉及推理攻击^[25-29]、重构攻击^[30-34]和窃取攻击^[35-38]。推理攻击是指攻击者通过观察全局模型的输出,推断某些敏感数据或信息;重构攻击是指攻击者通过分析全局模型的参数或更新,重构原始数据或信息;窃取攻击是指攻击者通过窃取全局模型或本地数据,获取敏感信息或知识。针对这些隐私性威胁,本文介绍了隐私性威胁的防御措施^[39-52],这些防御措施主要采用加噪机制和加密机制(见表1)。

表1 联邦学习系统安全威胁和防御措施

系统安全威胁	攻击者模型	攻击者来源	攻击类型	文献序号(攻击)	防御措施	文献序号(防御)
鲁棒性威胁	恶意攻击	参与方	数据投毒 后门攻击 模型攻击	[5-19]	检查更新的可信性	[7,9-11,15-17,20-24]
隐私性威胁	诚实但好奇攻击	参与方/服务端 系统外部	推理攻击 重构攻击 窃取攻击	[25-38]	加密机制和加噪机制	[39-52]

1 理论知识

1.1 联邦学习

联邦学习的定义^[53]为: n 个数据拥有者 $\{F_1, F_2, \dots, F_n\}$ 希望利用本地数据 $\{D_1, D_2, \dots, D_n\}$ 共同训练一个机器学习模型。传统方式是将所有数据收集到一个中心,使用 $D = D_1 \cup D_2 \cup \dots \cup D_n$ 训练一个模型 M_{sum} 。联邦学习是一个多方协作学习的过程,数据的拥有者 F_i 利用本地数据 D_i 协同训练全局模型 M_{fed} ,且 F_i 不会将 D_i 暴露给其他数据拥有者 $F_j(j \neq i)$ 。将 M_{fed} 的准确率记为 V_{fed} ,它应当非常接近 M_{sum} 对应的 V_{sum} 。形式上,设 δ 为非负实数(δ 是一个很小的正数),如果满足:

$$|V_{fed} - V_{sum}| < \delta$$

(1)

则称联邦学习有 δ -accuracy损失。

联邦学习系统的体系结构如图1所示。

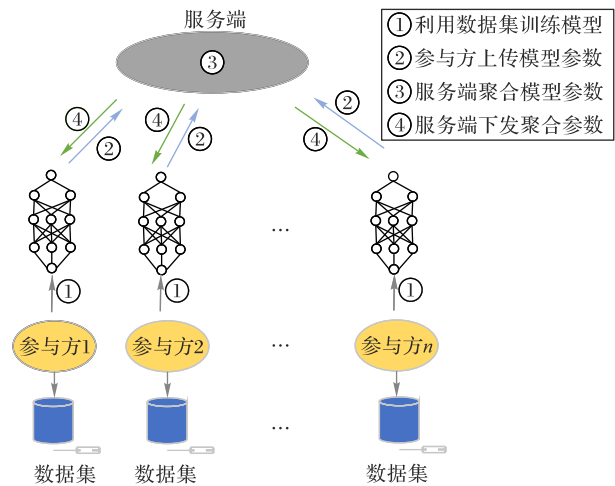


图1 联邦学习系统的体系结构

Fig. 1 Architecture of federated learning system

此系统的训练过程通常包含以下步骤:

- 1)参与方利用本地数据集训练局部模型。
- 2)参与方上传更新后的模型参数。
- 3)服务端聚合各个参与方的模型参数。

4)服务端广播聚合后的模型参数。

不断地迭代步骤1)~4)直至损失函数收敛,从而完成整个训练过程。

1.2 联邦学习分类

联邦学习可以根据参与方数和数据量的不同,分类为跨孤岛联邦学习(Cross-Silo Federated Learning)和跨设备联邦学习(Cross-Device Federated Learning)^[54]。跨孤岛联邦学习参与方数较少但拥有较大数据量,例如某银行和借贷公司利用各自本地的数据共同训练一个判断是否向客户借贷的模型;跨设备联邦学习参与方数较多但拥有较小的数据量,例如众多手机用户使用本地数据共同训练键盘输入词的预测模型。此外,联邦学习可以根据参与方的数据包含的特征,分类为横向联邦学习(Horizontal Federated Learning)、纵向联邦学习(Vertical Federated Learning)和联邦迁移学习(Federated Transfer Learning)^[53]。

1.3 攻击者模型

本文讨论联邦学习中的2种攻击者模型:恶意攻击(Malicious attack)和诚实但好奇攻击(Honest-but-Curious attack)。恶意攻击者的主要目的是攻击和破坏系统,而不是获取系统中的信息或数据。在联邦学习中,这种攻击行为会导致模型向偏离正常方向的方向学习。与之相比,诚实但好奇的攻击者会遵守协议规则和保密性,但仍然会尝试了解关于协议所传输数据的更多信息的模型。在联邦学习中,诚实但好奇的攻击者不会破坏模型的学习,但会尝试从传递的信息中推断敏感信息。

2 安全威胁与防御措施

在联邦学习场景中,攻击者的主要目标通常是破坏模型或者推断隐私信息。这些攻击者可能来自参与方或者服务端(见图2,以参与方 n 为例)。来自参与方的攻击者可能是恶意的,也可能是诚实但好奇的,无论哪种情况都可能威胁系统的鲁棒性和隐私性;来自服务端的攻击者通常被假设为诚实但好奇的,可能会威胁系统的隐私性,这些攻击者可以在不同的时间点发动攻击。来自参与方的攻击者可以在本地训练阶段或与服务端进行信息交互的阶段发动攻击,而来

自服务端的攻击者则只能在与客户端进行信息交互的阶段发动攻击。此外,当最终的模型通过应用程序接口(Application Programming Interface, API)提供服务时,非系统内部的成员也可能发动窃取攻击。

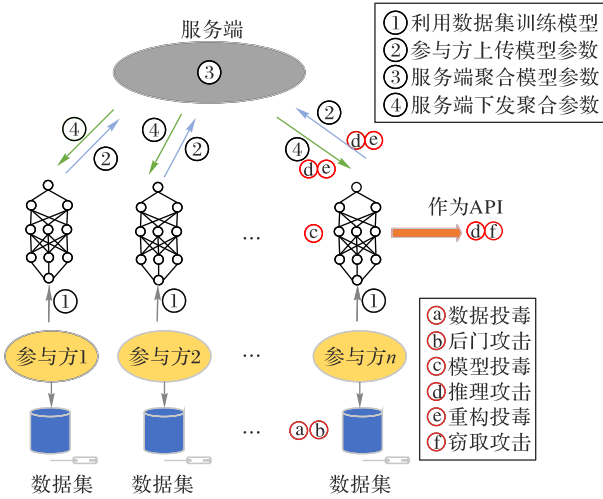


图2 联邦学习系统各阶段存在的攻击

Fig. 2 Attacks at all stages of federated learning system

2.1 针对系统鲁棒性攻击与防御

机器学习中的鲁棒性是指模型在面对输入数据中的干扰或噪声时的表现。具有良好鲁棒性的模型即使在输入数据中存在扰动或异常情况下也能保持良好的预测能力。在联邦学习中,破坏系统的鲁棒性意味着破坏联邦学习中的全局模型。攻击者可能使用各种方法破坏系统的鲁棒性,包括数据中毒、后门攻击和模型中毒等。本节讨论对系统鲁棒性造成威胁的攻击者不同情况下采用的攻击方法(见表2)、攻击效果(见表3)和相应的防御措施(见表4)、防御效果(见表5)。显然,后门攻击属于数据投毒,但考虑后门攻击的独特性,本文将分别介绍后门攻击和数据投毒。

表2 联邦学习系统鲁棒性威胁的攻击

Tab. 2 Robustness threat attacks in federated learning system

鲁棒性威胁	数据分布	攻击类型	文献序号	攻击方法论	补充说明
数据投毒	非独立同分布	无目标投毒	[6]	利用投影随机梯度上升算法最大限度地增加了目标节点的经验损失	
	独立同分布		[5]	预测由恶意输入引起的SVM决策函数的变化,并利用这种能力构造恶意数据	
	独立同分布	有目标投毒	[7]	力求最后几轮中成功投毒,并选择合适的标签翻转对象	未考虑参与方数据是非独立同分布的情况
	非独立同分布		[8]	利用GAN技术生成数据并实施标签翻转	攻击效果依赖攻击者选择的攻击时机
模型投毒	独立同分布或非独立同分布	无目标模型投毒	[9]	向局部模型添加随机噪声	
			[9]	根据安全聚合算法,修改局部模型参数使全局模型向相反方向更新	
	非独立同分布	有目标模型投毒	[10]	最大化恶意模型的更新,同时限制更新避免被检测	
			[12]	放大恶意模型更新,使得一次攻击留下足够强的后门	需要人工估计局部模型的放大因子
后门攻击	非独立同分布	标记后门	[13]	利用全局模型的历史更新推出反向更新的方向,并提交大量放大后的攻击模型更新	
			[14]	多个攻击者分布式地植入后门	
		语义后门	[12]	选择诚实参与者数据集中较少的特征作为后门特征	假设攻击者已知异常检测策略
			[17]	增大攻击者比例,降低后门任务的复杂性,限制恶意更新	
		后门攻击	[18]	利用数据集的分布特点,挑选边缘数据实施后门攻击,并限制恶意更新	
			[19]	后门模型和主任务模型分开训练,再优化结合为一个模型	

2.1.1 数据投毒

在机器学习中,数据投毒是攻击者通过恶意更改本地源数据影响模型的学习和预测结果,以达到攻击目的的一种方法。在联邦学习中,数据投毒攻击可以分为无目标的投毒攻击^[5-6]和有目标的投毒攻击^[7-8]两种。无目标的投毒攻击注重完全破坏模型性能,使全局模型的准确率越低越好,甚至使模型不能收敛,这样才算攻击成功。而有目标的投毒攻击则注重特定预测任务,攻击者不希望破坏模型对其他任务的预测,而是通过将原标签值更改为另一个标签值等方式干扰模型对特定任务的预测^[7]。

对于投毒攻击,攻击者的目标是通过篡改本地数据影响联邦学习中的局部模型,从而影响整个模型的准确性和鲁棒性。攻击者可以通过参与最后几轮迭代学习,增大攻击的威胁,因为在这些轮次中,全局模型已经接近收敛,而攻击者能够让它偏离正确方向。此外,攻击者参与的轮数和操作的数据量也会影响攻击的效果。如果攻击者参与的轮数较少,那么它篡改的数据可能无法影响整个模型;如果攻击者操作的数据量很少,那么它的影响也会被削弱。因此,攻击者的人数、参与的轮数和操作的数据量越多,攻击的效果就越强。相反地,诚实参与者的参与会削弱甚至消除投毒攻击的影响^[7]。

针对无目标的投毒攻击,文献[5]中提出了一种针对支持向量机(Support Vector Machine, SVM)的攻击方式,利用梯度上升策略,根据SVM最优解的性质计算梯度,并通过注入精心构建的训练数据,从而增加SVM的测试错误;但是,这种攻击仅适用于SVM模型。文献[6]中提出了AT2FL(ATTack on Federated Learning)算法,设计了一种基于投影随机梯度上升的方法,能够有效地推导中毒数据的隐式梯度,并利用它计算最优攻击策略。文献[6]中还证明了在多任务联邦学习情景下,联邦多任务学习模型容易受数据投毒攻击的影响,且随着恶意数据量的增加,模型的性能将变得越来越差。

表 3 联邦学习系统鲁棒性威胁的攻击效果
Tab. 3 Attack effects of robustness threats in federated learning system

文献 序号	攻击评价 指标	数据集	训练设置	结果/%	说明
[6]	模型错误率	EndAD	6;*:no-iid	base:6.881±0.52 result:28.588±3.74	
		Human Activity	30;*:no-iid	base:2.586±0.84 result:29.422±2.96	
		Landmine	29;*:no-iid	base:5.682±0.28 result:13.648±0.54	
[5]	模型错误率	MNIST	*;*;*	base:2-5 result:15-20	
[7]	最大召回率 损失	CIFAR-10	50;*:iid	base:0 result:2;1.42; 20:25.4	随不同恶意参与者百分比 (a:b,a表示百分比,b表示损失)
		Fashion-MNIST	50;*:iid	base:0 result:2;0.61; 20:29.2	
[8]	中毒任务 准确率	MNIST	10;*:no-iid	base:0 result:20:60±;40:80±;60:85±	随不同放大因子 (a:b,a表示放大因子,b表示准确率)
		AT&T	10;*:no-iid	base:0 result:20:70±;40:85±;60:90±	
[10]	最大准确度 损失	CIFAR10 (Alexnet)	50;1 000;*	base:0 result:Krum:43.6; Mkrum:36.8; Bulyan:45.6; TrMean:45.8; Median:40.9; AFA:47.0; FangTrmean:56.3	在不同聚合算法下 (a:b,a表示聚合算法,b表示损失)
[12]	后门任务 准确率	CIFAR10	100;*:no-iid	base:0 result:95:30±,75±	一次攻击、随全局迭代 (a:b,c,a表示轮次,b代表最低 准确度,c代表最高准确率)
		Reddit dataset	83 293;247;no-iid	base:0 result:95:0±,60±	
		CIFAR10	100;*:no-iid	base:0 result:1:30±,50±;5:80±,80±	
		Reddit dataset	83 293;247;no-iid	base:0 result:0.01:30±,65±;0.1:80±,90±	
[13]	模型准确率	MNIST	1 000;*:no-iid	base:99± result:FedAvg:1-25:10±; Median:1:90±; Median:25:60±; Trimmed-mean:1:95±; Trimmed-mean:1:50±	随恶意参与者百分比 (a:b;c,a表示聚合算法,b表示恶意 参与者占比,c表示准确率)
[14]	攻击成功率	LOAN	51;*:no-iid	base:0 result:8:99±	随全局迭代 (a:b,a表示迭代数,b表示成功率)
		MNIST	100;*:no-iid	base:0 result:20:99±	
		CIFAR	100;*:no-iid	base:0 result:350:80±	
		Tiny-imagenet	100;*:no-iid	base:0 result:80:80±	
[17]	后门任务 准确率	MNIST	*;*:no-iid	base:0 result:10:120:95±;50:300:70±	固定任务数、随全局迭代 (a:b;c,a表示任务数,b表示迭代数, c表示准确率)
[18]	后门任务 准确率	CIFAR-10	200;*:no-iid	base:0 result:100:400:55±;10:400:15±	(a:b;c,a表示攻击资源占比,b表示迭 代数,c表示准确率)
		MNIST	20;*;*	base:0 result:10:80±;20:90±	
		CIFAR-10	50;*;*	base:0 result:300:40±;400:90±	

注:在“训练设置列”采用x;y;z格式,x代表参与者数,y代表每个参与方拥有数据的数据量(条数),z代表数据集的划分
(iid或no-iid,iid表示以独立同分布划分数据,no-iid表示以非独立同分布划分数据),
“结果列”中“base”表示未攻击的情况下的基准值,“result”表示攻击后,“*”代表未知。

表 4 防御鲁棒性威胁的措施
Tab. 4 Measures to defend against robustness threats

防御 类型	数据分布	文献 序号	针对攻击类型	防御思想	防御方式	补充说明
数据 投毒	独立同分布	[20-21]	无目标投毒或有目标投毒	基于行为	利用鲁棒性的分布式梯度下降算法聚合模型	存在超参数
		[7,22]	有目标投毒	基于聚类	利用聚类算法鉴别恶意模型	
	独立同分布或非独立同分布	[11]	有目标投毒	基于行为	根据局部模型与全局模型的余弦相似度判断恶意模型	
模型 投毒	独立同分布或非独立同分布	[23]	有目标投毒		根据局部模型与全局模型的余弦相似度并结合信誉机制共同判断恶意模型	存在超参数
		[9]	无目标模型投毒	基于行为	利用拜占庭鲁棒性算法	
	独立同分布或非独立同分布	[11]			基于局部更新与全局更新的余弦相似度去除恶意梯度	
		[9]	有目标模型投毒	基于行为	基于错误率和基于损失函数的评价指标并结合拜占庭鲁棒性聚合算法防御	
后门 攻击	非独立同分布	[10]			基于奇异值分解(SVD)的谱方法检测和去除异常值	存在超参数
		[15]	标记后门攻击	基于行为	设置一个阈值,对每一轮中更新的每个维度进行投票,根据投票值是否超过阈值,动态调节该维度的学习率	
	独立同分布或非独立同分布	[16]	标记后门攻击	针对机器学习模型本身	剪裁异常神经元,约束神经元权重,微调模型	
	独立同分布或部分非独立同分布	[24]	语义后门攻击	混合策略 (基于聚类、 基于行为)	结合聚类和分类综合判定一个模型是否有害,通过剪裁策略削弱绕过检测的有毒模型的影响	微调可能导致后门加深

表5 防御鲁棒性威胁措施的效果

Tab. 5 Effectiveness of measures to defend against robustness threats

防御指标	文献序号	防御结果/%	说明
模型错误率	[9]	base:0.12 attack:* after:0.12	随着全局迭代的最终结果
	[11]	base:2.80±0.12 attack:* after:2.99±0.12±2.96±0.15/3.04±0.14	随着全局迭代的最终结果,“after”中分别对应拜占庭、标签翻转和噪声攻击下的结果
	[20]	base:10± attack:60± after:10±	随着全局迭代的最终结果
模型准确率	[21]	base:94.3± attack:77.3± after:90.7±	随着全局迭代的最终结果
	[22]	base:78± attack:76±/74.5± after:78±/77.5±	随着全局迭代的最终结果,“attack”和“after”中分别对应20%和30%的恶意参与者占比
	[23]	base:* attack:* after:83.11/81.23	随着全局迭代的最终结果,“after”中分别对应5%和50%的恶意参与者占比
模型准确率损失	[10]	base:0 attack:* after:4.3	随着全局迭代的最终结果
后门任务准确率	[15]	base:6.6 attack:88.6 after:9.0	随着全局迭代的最终结果
	[16]	base:* attack:85.5 after:4.8	
	[24]	base:* attack:100 after:0	

注:“防御结果”中的“base”表示未受到攻击时的结果,“attack”表示攻击后的结果,“after”表示使用防御措施后的结果。

针对有目标的数据投毒攻击,文献[7]中采用标签翻转攻击进行实验,该攻击证明即使恶意参与方的比例很小(低至总参与方的4%),攻击也可以显著地影响联邦学习的效果,同时,该攻击可能会因数据集和翻转对象的不同产生较大的差异;然而,文献[7]中没有考虑数据非独立同分布的情况。在此基础上,文献[8]中结合生成对抗网络(Generative Adversarial Nets, GAN)实现了标签翻转攻击,该攻击可以使主要任务和中毒任务的准确率都达到80%以上。攻击者首先伪装成正常的参与方参与联邦学习,当全局模型达到一定的准确度后,攻击者使用GAN技术生成一些类似其他正常参与方的数据,并翻转数据的标签;其次,攻击者将局部模型训练时的数据源更换为上一步得到的数据,并训练局部模型。该攻击适用于特定场景,例如在分类任务中,攻击者本地没有某类数据,但希望将该类作为攻击的目标。需要注意的是,该攻击依赖GAN生成的模拟数据,而GAN生成数据的质量与联邦学习中的全局模型质量有关,因此,攻击者选择的攻击时机非常重要。

2.1.2 模型投毒

本文联邦学习中的模型投毒攻击指直接修改模型更新的行为。这种攻击比数据投毒更直接地危害系统的鲁棒性,因为它直接作用于模型参数而非修改训练数据集。虽然模型投毒攻击和数据投毒攻击都修改了局部模型,但它们的攻击时机不同:数据投毒攻击通过修改训练数据集学习不同于正常情景下的模型参数;模型投毒攻击直接修改模型更新的参数。数据投毒攻击通常是有目标的攻击,模型投毒攻击更能直接有目标地缩放模型参数,因此,模型投毒攻击更常与其他攻击结合使用,例如结合后门攻击增加后门的效果^[12]。模型投毒^[9-13]攻击通常分为两种实施方式:第一种是无目标模型投毒,攻击者向训练的局部模型中添加随机噪声^[9,11]扰乱它的学习过程;第二种是有目标模型投毒,攻击者会更精细地修改局部模型^[12-13],使全局模型朝着某个特定的方向学习,或是为了针对已有的防御算法进行攻击^[9-10]。由于有目标的模型投毒攻击更普遍,因此本文主要讨论这种类型的攻击。

文献[9]中证明了有目标的模型投毒攻击的可行性,并指出针对无目标投毒攻击的拜占庭鲁棒性聚合防御(如中值聚合、修剪平均聚合和Krum)不能单独用于防御有目标的模型投毒攻击。文献[10]中提出了一个通用的模型投毒攻击

框架,通过在恶意方向上最大限度地扰动良性参考聚合计算恶意模型更新,并限制模型更新以避免被健壮的聚合算法检测。文献[12]中提出了模型替换攻击,通过假设并反解联邦平均聚合算法,解出使用一次攻击就能使全局模型达到局部模型的解,最终只需要确定一个局部模型的放大因子就能实现模型替换攻击;但该攻击的缺点是需要人工估计局部模型的放大因子。文献[13]中假设攻击者仅知道训练过程中的全局模型,甚至没有正常的数据集。在这种情况下,提出了一种基于伪客户端的模型中毒攻击MPAF(Model Poisoning Attack based on Fake clients),与上述攻击方式不同,它旨在将全局模型“拉向”攻击者指定的基本模型。它的攻击的关键步骤是:1)随机选择测试精度低的基本模型。2)根据全局模型的历史更新生成攻击更新,并在将它发送到服务端前放大,以扩大攻击效果。

2.1.3 后门攻击

联邦学习中的后门攻击^[12,14-19]指恶意攻击者利用特殊的数据集参与联邦学习,最终影响全局模型,并在特定的输入下激活后门输出攻击者想要的输出。后门攻击与数据投毒相似又不同:相似之处是它们都攻击参与方的源数据;不同之处是后门攻击具有隐蔽性,且不会干扰模型对正常输入的预测。将后门攻击又称为有针对性的模型投毒攻击。

后门攻击可以根据实施方式是否更改本地数据分为两种类型:标记后门攻击^[14-16]和语义后门攻击^[12,17-18]。在标记后门攻击中,攻击者向原始数据添加特殊标记,例如在图像中添加马赛克^[16],这种攻击方式的效果更难消除,因为模型已经记住了特定标记。在语义后门攻击中,攻击者利用数据集的特征,例如将数据集中所有绿色汽车的标签值修改为鸟类^[12],这种攻击方式可以通过诚实参与者拥有相似数据集的参与逐步消除后门。无论哪种攻击方式,后门攻击的效果只在特定的输入时才触发后门。多位研究者的研究^[12,14]证明了后门攻击凭借隐蔽性能够在仅发动攻击成功一次的情况下,使全局模型能在多轮的迭代中保留后门。本文主要讨论语义后门攻击。

针对标记后门攻击,文献[14]中提出了一种分布式的后门攻击方法,并证明了相较于集中式的后门攻击,联邦学习遭受分布式后门攻击的危害更大。此外,文献[14]中还证明了联邦学习中一些鲁棒性的聚合算法(例如RFA(Robust Federated Aggregation)和FoolsGold)无法防御分布式后门攻击。

针对语义后门攻击,已有多篇文献提出了不同的防御方法和策略:文献[16]中证明了基于特定标记的后门攻击在数据非独立同分布程度越高时攻击越有效。文献[12]中提出了规避防御的语义后门攻击,并通过放大本地图更新实现模型替换,同时也证明了攻击效果随诚实参与者的参与而下降;然而,假设攻击者已经熟知异常检测策略,这一假设并不太现实。为此,文献[17]中提出了范数有界后门攻击(Norm Bounded Backdoor Attack),通过约束更新规避一些防御措施,并证明了攻击的有效性;同时,量化恶意攻击者数与参与攻击的频率对后门攻击的影响,结果是攻击者比例越大、攻击者参加频率越高,后门攻击的效果越好。文献[18]中提出了边缘后门攻击(Edge-Case Backdoors Attack),利用分布在边缘的数据(出现在训练集和测试集的可能性较小)实现后门攻击,并从理论和实验证明了它难以检测和防御。文献[19]中提出了基于优化模型的后门攻击,通过将冗余神经元训练为对抗神经元实现攻击,实验结果验证了这种后门攻击不仅能实现较高的攻击成功率,还能规避一些防御措施。

2.1.4 防御破坏系统鲁棒性的措施

在联邦学习中,数据投毒、模型投毒和后门攻击都会破坏系统的鲁棒性。尽管针对数据投毒的防御措施通常直接筛选数据集并去除其中的恶意数据,但由于联邦学习不能操作参与方的数据,因此这些针对数据集的防御措施违背了联邦学习的隐私保护要求,必须从其他方面考虑防御措施。通过分析联邦学习的流程可知,防御措施只能在服务端的聚合阶段实施。数据投毒对模型的影响难以预测,模型投毒能精细缩放模型,后门攻击由于隐蔽性和常与模型投毒结合,使防御更困难;因此,后门攻击的防御难度更高,需要更精细的设计。同时,现实数据的非独立同分布特性使得在判断一个模型更新是否为恶意模型更新时需要更细致的考虑。

在联邦学习中,参与方的模型参数或梯度信息被上传到服务端,因此防御措施集中在模型质量和模型之间的差异方面。目前,主要有3种防御思想:1)聚类。通过聚类分析模型参数区分好坏模型,从而识别潜在的恶意参与方,相关研究包括文献[7,22,24]等。2)基于行为的防御。通过分析参与方上传的模型在行为方面的特征,如局部更新与全局更新之间的相似性、部分模型聚合后的错误率、模型更新的阈值等识别潜在的恶意参与方,相关研究包括文献[9-11,15,17,20-21,23]等。3)针对机器学习模型本身进行防御。这种方法主要针对机器学习模型本身的弱点(如神经网络中的神经元)进行防御,旨在防止参与方利用这些弱点破坏模型的准确性和安全性,相关研究包括文献[16]等。

对数据投毒的防御 针对数据投毒攻击,可以根据攻击是否有目标采用不同的防御措施。对于无目标攻击,一种轻量级的防御措施是Krum聚合算法^[20],它选择欧几里得距离最小的局部模型进行聚合。另外,中值聚合(Median Aggregation)和修剪平均聚合(Trimmed Mean Aggregation)也被证明在遭受无目标攻击时具有鲁棒性^[21]。这些方法的一个潜在缺点是,它们可能无法有效地抵御有目标攻击。对于有目标攻击,防御标签翻转攻击的一个算法是基于聚类的思想,它记录每个参与方的局部更新与全局更新的差值,并使用主成分分析(Principal Component Analysis, PCA)技术对数据降维,以观察正常参与方与恶意攻击者上传的更新^[7]。文献[22]中在此基础上提出了使用KPCA(Kernel Principal Component Analysis)和K-means聚类代替PCA的方法,从而获得更好的防御效果。此外,自适应联邦平均(Adaptive Federated Averaging, AFA)和CONTRA^[23]是基于行为的防御

方法,它们利用余弦相似度确定局部模型的可信度,并通过信誉方案,根据单个客户的每一轮和对全局模型的历史贡献动态提升或惩罚单个客户;然而,这些基于聚类 and 行为的防御方法在面对非独立同分布的数据时可能出现误判的情况。

对模型投毒的防御 对于模型投毒攻击,可以采取不同的防御措施提高模型的鲁棒性。针对无目标的模型投毒攻击,可以使用拜占庭鲁棒性聚合算法^[9,11]防御,通过聚合多个参与方的模型权重降低恶意攻击的影响。具体地,可以采用中值聚合、修剪平均聚合和Krum算法等实现。这些算法可以有效地防御无目标的模型投毒攻击,因为这种攻击与无目标的投毒攻击在一定程度上对模型造成的危害相似。对于更精细的有目标模型投毒攻击,则需要设计更精细的防御才能有效地抵抗。一种常见的方法是采用基于错误率和基于损失函数的评价指标,并结合拜占庭鲁棒性聚合算法^[9]实现防御,这种方法可以在一定条件下有效地提高模型的鲁棒性。此外,文献[11]中提出了AFA算法,基于局部更新与全局更新的余弦相似度去除恶意梯度,可以有效防御无目标的模型投毒攻击;但是,该算法存在一些超参数需要人工指定。文献[10]中提出了DNC(Divide-aNd-Conquer)算法,利用基于奇异值分解(Singular Value Decomposition, SVD)的光谱方法检测和去除异常值,可以提高模型的鲁棒性。

对后门攻击的防御 后门攻击是一种特殊的攻击手段,旨在通过在模型中植入后门,在特定条件下产生错误的输出。由于后门攻击的高度隐蔽性,传统的数据投毒攻击防御策略并不一定适用于此种攻击方式^[24],因此,需要更精细的设计防御后门攻击。在防御标记后门攻击方面,多种方法被提出:文献[16]中通过剪枝冗余神经元并对权重偏离正常值的神经元进行参数约束防御后门攻击;文献[15]中通过动态调整每轮中每个更新维度的学习率降低恶意参与者的更新的影响;文献[17]中提出了使用更新范数阈值约束和差分隐私(Differential Privacy, DP)^[39]技术防御后门攻击,其中,更新范数阈值约束方法是服务器简单地忽略规范高于某个阈值 M 的更新,而DP技术则通过添加噪声防御后门攻击。对于语义后门攻击,文献[24]中提出了一种名为DeepSight的模型过滤方法,结合了聚类和分类策略鉴别有毒模型,并通过剪裁策略削弱绕过检测的有毒模型的影响。

2.2 针对系统隐私性攻击与防御

机器学习面临的隐私攻击也对联邦学习构成威胁,因为联邦学习过程中涉及参与者的数据隐私和全局模型的隐私。攻击者可以通过联邦学习交互过程中的信息或模型本身的信息推断隐私信息,从而破坏系统的隐私。在联邦学习中,可能会发生的隐私攻击包括推理攻击、重构攻击和窃取攻击,这些攻击可以来自系统内部的参与者和服务端,也可以来自系统外部的API请求者(见表6)。针对系统内部的攻击,参与者和服务端可能利用模型访问权限,使用模型参数或梯度信息推断私密信息。为了防止这种攻击,联邦学习系统需要采取措施限制参与者和服务端的访问权限,并对交互数据进行加密和匿名化处理(见表7)。在系统外部,API请求者可能试图利用API请求中的信息推断私密信息,例如使用模型输出推断输入数据的敏感信息。

2.2.1 推理攻击

推理攻击^[25-29]旨在利用可获得的信息推断系统不想暴露的敏感信息,这对机器学习模型尤其危险。攻击者可以利用模型的相关信息或者模型的API进行推理攻击。在机器学习中,推理攻击的经验依据包括以下几个方面:在对自然语言文本序列进行分类或预测的神经网络模型中,对于某些

特殊的训练数据序列,存在被生成模型无意识记忆的风险^[25],这意味着攻击者可以通过这些特殊的训练数据序列识别模型的一些敏感信息。机器学习模型对训练数据样本和非训练数据样本的表现不同^[26]等。

表6 联邦学习系统隐私性威胁的攻击

Tab. 6 Privacy threat attacks on federal learning system

隐私性威胁	攻击者来源	文献序号	攻击者知识
成员推理攻击	系统内部	[26, 28-29]	无需额外知识
	系统外部	[26]	模型API,模型的训练数据的背景知识
		[28]	模型损失函数,损失范围
重构攻击	系统内部	[30-31, 33-34]	无需额外知识
	系统内部服务端	[32]	无需额外知识
窃取攻击	系统外部	[35]	输出置信度的模型API,模型架构
		[36-37]	输出标签的模型API,模型架构
		[38]	输出置信度的模型API

表7 防御隐私性威胁的措施

Tab. 7 Measures to defend against privacy threats

防御措施	防御技术	防御思想	缺点	文献序号
加噪机制	客户级	掩盖原始梯度	服务端必须可信	[40]
	差分隐私	信息		
	本地级	掩盖原始梯度	较大地影响模型性能	[41-43]
加密机制	差分隐私	信息		
	同态加密	在密文上进行计算	加密效率低、密文的膨胀率高	[48-50]
	秘密分享	将秘密信息划分多份	增加计算成本和通信成本	[51-52]

推理攻击对隐私泄露的影响非常严重,可以揭示一些敏感信息,如患者的病史、种族和性别等。例如,通过成员推理攻击,攻击者可以揭示某个患者是否患有某种疾病,因为他们知道该患者的治疗记录被用于训练某种与该疾病相关的模型^[26]。另外,利用属性推理攻击,攻击者可以在分类男女的人脸识别机器学习模型中同时预测输入是否为白种人^[27]。推理攻击可以细分为不同的类型,但在联邦学习的情景下,本节仅讨论成员推断攻击。

对于成员推理攻击,文献[26]中提出了一种攻击方法,该方法适用于已经部署为服务的分类机器学习模型。攻击方法的基本思路是训练一个推理模型判断输入数据是否是分类模型的训练数据。具体地,攻击者输入分类模型的训练数据和非训练数据,并获取分类模型的置信度值输出。然后,将输出与数据对应的标签作为推理模型的训练数据,并根据输入是否用于训练过分类模型,作为上一步训练数据的标签,以此训练一个推理模型判断输入是否为训练数据。然而,将此攻击应用于联邦学习且攻击来自系统外部的情况下存在两个难题。首先,模型的输出可能不是置信度值;其次,攻击者缺乏模型的训练数据。为了解决这些难题,文献[26]中假设攻击者拥有模型的API或模型架构及训练算法,并通过一些背景知识(如统计信息、一部分训练数据和加噪后的训练数据)获得数据集,并利用它训练出与全局模型相似的影子模型,再进行上述推理模型的训练。对于系统外部的攻击者,可能较难获得模型训练的数据的背景知识,同时训练多个影子模型可能需要更多的算力资源和数据集资源。文

献[28]中提出了一种成员推理攻击方法,即根据输入的损失实施攻击。具体地,对于参与过训练的数据输入,模型的函数损失在一定范围内;而对于未参与训练的数据输入,函数损失将超出这个范围。通过实验证明,这种攻击方法的推理准确率虽然略低于文献[26]算法,但不需要更多的计算资源。将此攻击应用于联邦学习,它的攻击者必须有模型损失函数的知识和损失范围。文献[29]中提出的攻击方法是对深度神经网络(Deep Neural Network, DNN)的白盒推理攻击。它利用了DNN的特性,即它具有大量的参数(数百万个)且不能正确地泛化到训练数据之外(在较多情况下,训练数据的大小要小一个数量级)。该攻击认为可以利用模型梯度信息区分成员和非成员训练数据,并利用这些信息训练攻击模型。需要注意的是,这种攻击方法需要访问目标模型的内部信息,包括模型的架构、参数和梯度信息等。因此,攻击者需要足够了解目标模型,以成功地实施这种攻击。

2.2.2 重构攻击

在机器学习中,重构攻击是指攻击者试图通过模型或模型输出恢复原始数据的一种攻击方式。这种攻击相较于推理攻击更危险,因为它试图完全重构原始数据,而不仅是笼统地判断数据。例如,攻击者可以从人脸识别模型中重构参与训练的人脸图像^[30],或者从分类模型中重构攻击者不拥有的某类数据^[31]。为了进行重构攻击,攻击者需要掌握模型或梯度等相关信息。因此,在联邦学习场景中,攻击者只能来自系统内部。

针对重构攻击,多种攻击算法被提出:文献[30]中提出了一种反演攻击算法,利用梯度下降算法最小化包含原模型的代价函数,并在优化过程中进行图像处理,以恢复原始图像。文献[31]中的协同分类机器学习攻击则利用GAN技术生成攻击者没有的某一标签对应的数据,并随着交互的进行,它的生成的数据质量越高。攻击者还可以将“尚未成熟”的生成数据加入局部模型训练,以刺激其他参与方输入更多有关此标签对应的数据。然而,这种攻击只能重构某标签对应数据的代表,而不能精确恢复训练数据。文献[32]中在文献[31]的基础上提出了mGAN-AI (multi-task GAN for Auxiliary Identification)重构攻击,旨在攻击特定的客户端,达到破坏客户端级的隐私的目的。文献[33]中的梯度的深度泄漏攻击(Deep Leakage from Gradients, DLG)证明了DLG能够仅通过分析梯度信息在计算机视觉和自然语言处理任务中重构输入,但此攻击假设梯度变化仅由一个输入造成,这个假设不切实际。为了提高攻击的准确性,文献[34]中在文献[33]的基础上提出了improved DLG (iDLG),它对任何用交叉熵损失训练的可微分模型都能提取输入的真实标签。

2.2.3 窃取攻击

在机器学习中,模型窃取攻击是指攻击者试图从一个已经训练好的模型中获取信息,以训练一个类似原始模型的新模型。在联邦学习中,讨论模型窃取攻击通常假设攻击者来自系统外部,攻击者一般需要通过模型API进行交互获取有关模型的信息。攻击者掌握的背景知识越丰富,就越可能生成与原始模型相似的新模型,这里的背景知识包括模型的架构和用于训练模型的数据集的背景知识。

模型窃取攻击的难点是需要确定被攻击模型的类型和内部结构,以及确定模型中的超参数,并构造用于模型窃取的数据集。通常情况下,攻击者需要先了解被攻击模型的类型。例如,对于输出预测置信度的模型,可以使用样本集和对应模型的输出置信度训练一个与原模型相似的模型,这种

方法在文献[35]中被提出。而对于只输出类标签的模型,也存在被攻击的风险。在DNN模型窃取攻击方面,文献[36-37]中针对只输出标签的DNN模型,提出了类似的算法,其中关键步骤包括使用原模型为生成的数据集打上标签,将打上标签的数据集作为攻击模型的训练集进行训练,并对数据进行增强处理,以提高模型的准确性。文献[37]中还提出了3种解决超参数问题的方案和2种创建训练样本的方案,并将查询预算融入攻击算法。此外,针对模型窃取攻击,还有一种方法是利用与原模型分布不同的训练样本,训练得到与原模型相似的模型。文献[38]中研究了这种方法所需的样本和模型,并证明了利用与原模型分布不同的训练样本可以训练得到与原模型相似的模型。在模型选择方面,文献[38]中还证明了原模型的输出可以用于训练具有不同模型架构的新模型,并选择较复杂的模型架构可以更有效地实现窃取攻击。

2.2.4 系统隐私性保护技术

应用于联邦学习中的隐私保护技术主要有:差分隐私(DP)^[39]、同态加密(Homomorphic Encryption, HE)^[45-46]和秘密分享(Secret Sharing)^[47]。

差分隐私 对于一个随机化算法 $M(X) \rightarrow R$,其中 X 为域, R 为值域。如果对于任意两个相邻数据集 $D, D' \in X$,以及任意子集 $S \in R$,算法 $M(X)$ 满足: $\Pr[M(D) \in S] \leq e^\epsilon \Pr[M(D') \in S] + \delta$ (δ 表示算法失败的概率上界),则称算法 $M(X)$ 满足 (ϵ, δ) -差分隐私。如果 $\delta = 0$,则称算法 $M(X)$ 满足 ϵ -差分隐私,其中 ϵ 表示隐私保护预算。

同态加密 对于原文信息 a 和 b ,如果算法 E 满足: $E(a + b) = a' \oplus b'$,则称算法 E 是加法同态加密算法,其中 a' 和 b' 分别是 a 和 b 在算法 E 下得到的密文。同态加密能细分为加法同态加密算法、乘法同态加密算法、半同态加密算法和全同态加密算法。

秘密分享 将秘密值 s 分成 n 份,分别分配给参与方 p_1, p_2, \dots, p_n ,利用公式: $f(x) = s + \sum_{i=1}^{t-1} a_i x^i$ 进行秘密分享,其中 t 是整数,且 $t \leq n$, a_i 是随机生成的系数,原秘密为常数项。对于参与方 p_i ,分配不同的 x_i 计算该参与方的秘密份额,具体可以通过公式 $s_i = f(x_i)$ 计算。当至少有 t 个参与方协作时,可以通过公式 $s = \sum_{j=1}^t s_j \prod_{l \neq j} \frac{x_l - x_i}{x_j - x_i}$ 恢复原始的秘密值 s 。

2.2.5 防御破坏系统隐私性的措施

在联邦学习中,确保参与者的数据隐私和模型隐私非常重要,因此需要采用一系列隐私保护技术。现有的隐私保护技术主要包括加噪机制^[39-43]和加密机制^[44-52]。加噪机制通过向模型更新的梯度信息中添加一些噪声,以掩盖真实的梯度信息,防止攻击者通过梯度信息推断参与者的数据隐私;但是,添加噪声会影响模型的准确性,一般地,添加的噪声越多,模型的性能就越差。加密机制加密梯度信息,以确保只有允许的所有者可以查看梯度信息,而其他人员无法访问梯度信息。这种技术不会影响最终的模型性能,但加解密过程会影响联邦学习的收敛速度,同时也会造成大量的计算消耗。因此,在将加密技术应用于联邦学习中保护梯度信息时,需要解决加密效率和通信成本的难题。

加噪机制通常在联邦学习中采用差分隐私技术保护隐私。差分隐私可以确保即使攻击者拥有除一条信息以外的所有已发布信息,也无法推断该信息。在差分隐私中,隐私预算决定了噪声添加的大小。隐私预算越小,保护级别越

高,但是添加的噪声也越多。将差分隐私技术应用于联邦学习保护隐私时,模型的收敛性能和隐私保护级别之间存在矛盾,收敛性能越好,保护级别越低。在固定隐私预算下,增加参与联邦学习的客户端数量可以提高收敛性能^[40-41]。文献[40]中提出了一种方法,即对上传至服务端的每个参与者的梯度信息添加噪声,然后聚合,从而实现客户级的隐私保护,即攻击者无法确定客户是否参与了训练;但是,这种方法无法保护参与者模型参数的隐私,因此一些研究表明,对于联邦学习中的差分隐私,参与者可以先添加噪声,再上传梯度信息,以保护模型参数的隐私^[40-42],这种方法可以实现本地级的隐私保护。例如,NbAFL(Noising before model Aggregation FL)将加噪时机定在参与方发送梯度信息前,从而达到本地级的隐私保护效果。此外,研究还表明,在固定的隐私预算下,存在一个最佳的聚合次数,以达到最佳的收敛性能。

加密机制在联邦学习中的应用采用HE^[45-46]和安全多方计算(Secure Multi-party Computation, SMC)^[44]中的秘密分享^[47],以保护梯度信息的隐私。同态加密允许在密文状态下进行计算,避免了解密数据的风险,但需要大量的计算才能实现加密和解密操作,且加密后的数据通常比原始数据大得多。秘密分享将一个秘密信息拆分成多份小份信息,分享给不同的参与方,只有当所有参与方将份额组合在一起时才能还原原始秘密信息。在联邦学习中应用同态加密时,参与方使用公钥加密梯度信息,服务端在密文上进行聚合,参与方使用私钥解密聚合信息。为了提高训练速度和加密效率,一些研究提出了改进的同态加密算法和编码方案,如使用改进的Paillier算法^[48]和BatchCrypt编码方案^[49],以降低通信开销和提高训练速度。但是,由于参与方使用相同的公私密钥,加密的梯度信息可能被截取,导致梯度信息泄露。为此,文献[50]中提出了多密钥同态加密的隐私保护方案,确保解密需要所有参与方之间的协作。秘密分享多用于特定的联邦学习架构保护梯度隐私,如结合边缘计算和区块链的联邦学习。在这种方案中,参与者利用秘密分享将梯度信息划分多份并传给多个边缘节点,每个边缘节点聚合不同参与方上传的梯度份额并上传至区块链,区块链聚合边缘节点上传的聚合份额^[51]。同时,一些研究还利用同态加密保护梯度信息的隐私,在结合区块链的联邦学习架构中,利用秘密分享确保同态加密中密钥的安全^[52]。

3 结语

联邦学习是一种新兴的技术,能够在不共享数据的情况下训练和共享模型,从而打破数据壁垒,促进数据的共享和协作,同时保护个人数据隐私;然而,联邦学习系统本身也存在鲁棒性和隐私性两方面的安全威胁。

目前针对威胁联邦学习系统鲁棒性方面的研究更多将参与方上传的模型参数区分为恶意和非恶意两类;然而,数据的非独立同分布本身就会使模型参数相差极大,这使得鉴别恶意模型参数变得更困难。虽然有研究证明在非独立同分布数据下能正确地检测恶意模型参数,但是防御策略本身存在一些缺陷,例如计算开销大、缺乏理论依据等。针对防御威胁联邦学习系统隐私性方面的研究更多利用加噪策略和加密策略保护隐私,但加噪机制会导致模型的精度下降,加密机制会导致高额通信成本和计算成本。另外,保证联邦学习系统隐私性方面的防御策略和鲁棒性方面的防御策略存在冲突,例如,利用同态加密技术能保证联邦学习系统的隐

私性,但是无法防御针对联邦学习系统鲁棒性的攻击。

未来的研究工作可以分为以下3种情况:

1)不可信的参与方和可信的服务端。在这种情况下,需要设计更健壮的防御策略,抵抗威胁联邦学习系统鲁棒性的攻击。

2)可信的参与方和不可信的服务端。在这种情况下,需要解决为了保证联邦学习系统隐私性带来的精度下降、通信和计算成本暴增的问题。

3)不可信的参与方和不可信的服务端。兼顾联邦学习系统的隐私性和鲁棒性是一个非常严峻的挑战,但这也是目前更值得研究的一种情况。如何应对这种严峻挑战,是未来工作的重点。

综上,设计一个安全的联邦学习系统是一个长期且具有挑战性的任务。这需要不断跟进最新的研究成果,并不断探索创新的方法以确保联邦学习系统的安全性。

参考文献 (References)

- [1] JORDAN M I, MITCHELL T M. Machine learning: trends, perspectives, and prospects [J]. *Science*, 2015, 349 (6245): 255-260.
- [2] ZHANG C, XIE Y, BAI H, et al. A survey on federated learning [J]. *Knowledge-Based Systems*, 2021, 216: 106775.
- [3] 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*. New York: PMLR, 2017, 54: 1273-1282.
- [4] ALBRECHT J P. How the GDPR will change the world [J]. *European Data Protection Law Review*, 2016, 2: 287-289.
- [5] BIGGIO B, NELSON B, LASKOV P. Poisoning attacks against support vector machines [EB/OL]. (2013-03-25) [2023-07-09]. <https://arxiv.org/pdf/1206.6389.pdf>.
- [6] SUN G, CONG Y, DONG J, et al. Data poisoning attacks on federated machine learning [J]. *IEEE Internet of Things Journal*, 2021, 9(13): 11365-11375.
- [7] TOLPEGIN V, TRUOX S, GURSOY M E, et al. Data poisoning attacks against federated learning systems [C]// *Proceedings of the 2020 European Symposium on Research in Computer Security*. Cham: Springer, 2020: 480-501.
- [8] ZHANG J, CHEN J, WU D, et al. Poisoning attack in federated learning using generative adversarial nets [C]// *Proceedings of the 2019 18th IEEE International Conference on Trust, Security and Privacy in Computing and Communications/13th IEEE International Conference on Big Data Science and Engineering*. Piscataway: IEEE, 2019: 374-380.
- [9] FANG M, CAO X, JIA J, et al. Local model poisoning attacks to Byzantine-robust federated learning [C]// *Proceedings of the 29th USENIX Conference on Security Symposium*. Berkeley: USENIX Association, 2020: 1623-1640.
- [10] SHEJWALKAR V, HOUMANSADR A. Manipulating the Byzantine: optimizing model poisoning attacks and defenses for federated learning [C/OL]// *Proceedings of the 2021 Network and Distributed System Security Symposium* [2023-05-30]. <https://par.nsf.gov/servlets/purl/10286354>.
- [11] MUÑOZ-GONZÁLEZ L, CO K T, LUPU E C. Byzantine-robust federated machine learning through adaptive model averaging [EB/OL]. (2019-09-11) [2023-07-09]. <https://arxiv.org/pdf/1909.05125.pdf>.
- [12] BAGDASARYAN E, VEIT A, HUA Y, et al. How to backdoor federated learning [C]// *Proceedings of the 23rd International Conference on Artificial Intelligence and Statistics*. New York: PMLR, 2020, 108: 2938-2948.
- [13] CAO X, GONG N Z. MPAF: model poisoning attacks to federated learning based on fake clients [C]// *Proceedings of the 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition*. Piscataway: IEEE, 2022: 3395-3403.
- [14] XIE C, HUANG K, CHEN P Y, et al. DBA: distributed backdoor attacks against federated learning [C/OL]// *Proceedings of the 2020 International Conference on Learning Representations (2020-12-19)* [2023-05-30]. https://openreview.net/attachment?id=rkgyS0VFvr&name=original_pdf.
- [15] OZDAYI M S, KANTARCIOGLU M, GEL Y R. Defending against backdoors in federated learning with robust learning rate [J]. *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021, 35(10): 9268-9276.
- [16] WU C, YANG X, ZHU S, et al. Mitigating backdoor attacks in federated learning [EB/OL]. (2021-01-14) [2023-07-09]. <https://arxiv.org/pdf/2011.01767.pdf>.
- [17] SUN Z, KAIROUZ P, SURESH A T, et al. Can you really backdoor federated learning? [EB/OL]. (2019-12-02) [2023-07-09]. <https://arxiv.org/pdf/1911.07963v2.pdf>.
- [18] WANG H, SREENIVASAN K, RAJPUT S, et al. Attack of the tails: yes, you really can backdoor federated learning [J]. *Advances in Neural Information Processing Systems*, 2020, 33: 16070-16084.
- [19] ZHOU X, XU M, WU Y, et al. Deep model poisoning attack on federated learning [J]. *Future Internet*, 2021, 13(3): 73.
- [20] BLANCHARD P, EL MHAMDI E M, GUERRAOUI R, et al. Machine learning with adversaries: Byzantine tolerant gradient descent [C]// *Proceedings of the 31st International Conference on Neural Information Processing Systems*. Red Hook: Curran Associates Inc., 2017: 118-128.
- [21] YIN D, CHEN Y, KANNAN R, et al. Byzantine-robust distributed learning: towards optimal statistical rates [C]// *Proceedings of the 35th International Conference on Machine Learning*. New York: PMLR, 2018, 80: 5650-5659.
- [22] LI D, WONG W E, WANG W, et al. Detection and mitigation of label-flipping attacks in federated learning systems with KPCA and K-means [C]// *Proceedings of the 2021 8th International Conference on Dependable Systems and Their Applications*. Piscataway: IEEE, 2021: 551-559.
- [23] AWAN S, LUO B, LI F. CONTRA: defending against poisoning attacks in federated learning [C]// *Proceedings of the 26th European Symposium on Research in Computer Security*. Cham: Springer, 2021: 455-475.
- [24] RIEGER P, NGUYEN T D, MIETTINEN M, et al. DeepSight: mitigating backdoor attacks in federated learning through deep model inspection [EB/OL]. (2022-01-03) [2023-07-09]. <https://arxiv.org/pdf/2201.00763.pdf>.
- [25] CARLINI N, LIU C, ERLINGSSON Ú, et al. The secret sharer: evaluating and testing unintended memorization in neural networks [C]// *Proceedings of the 28th USENIX Security Symposium*. Berkeley: USENIX Association, 2019: 267-284.
- [26] SHOKRI R, STRONATI M, SONG C, et al. Membership inference attacks against machine learning models [C]// *Proceedings of the 2017 IEEE Symposium on Security and Privacy*. Piscataway: IEEE, 2017: 3-18.
- [27] MALEKZADEH M, BOROVYKH A, GÜNDÜZ D. Honest-but-curious nets: sensitive attributes of private inputs can be secretly

- coded into the classifiers' outputs [C]// Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security. New York: ACM, 2021: 825-844.
- [28] YEOM S, GIACOMELLI I, FREDRIKSON M, et al. Privacy risk in machine learning: analyzing the connection to overfitting [C]// Proceedings of the 2018 IEEE 31st Computer Security Foundations Symposium. Piscataway: IEEE, 2018: 268-282.
- [29] NASR M, SHOKRI R, HOUMANSADR A. Comprehensive privacy analysis of deep learning: passive and active white-box inference attacks against centralized and federated learning [C]// Proceedings of the 2019 IEEE Symposium on Security and Privacy. Piscataway: IEEE, 2019: 739-753.
- [30] FREDRIKSON M, JHA S, RISTENPART T. Model inversion attacks that exploit confidence information and basic countermeasures [C]// Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security. New York: ACM, 2015: 1322-1333.
- [31] HITAJ B, ATENIESE G, PEREZ-CRUZ F. Deep models under the GAN: information leakage from collaborative deep learning [C]// Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security. New York: ACM, 2017: 603-618.
- [32] WANG Z, SONG M, ZHANG Z, et al. Beyond inferring class representatives: user-level privacy leakage from federated learning [C]// Proceedings of the 2019 IEEE Conference on Computer Communications. Piscataway: IEEE, 2019: 2512-2520.
- [33] ZHU L, LIU Z, HAN S. Deep leakage from gradients [EB/OL]. (2019-12-19) [2023-07-09]. <https://arxiv.org/pdf/1906.08935.pdf>.
- [34] ZHAO B, MOPURI K R, BILEN H. iDLG: improved deep leakage from gradients [EB/OL]. (2020-01-08) [2023-07-09]. <https://arxiv.org/pdf/2001.02610.pdf>.
- [35] TRAMÈR F, ZHANG F, JUELS A, et al. Stealing machine learning models via prediction APIs [C]// Proceedings of the 25th USENIX Conference on Security Symposium. Berkeley: USENIX Association, 2016: 601-618.
- [36] PAPERNOT N, McDANIEL P, GOODFELLOW I, et al. Practical black-box attacks against machine learning [C]// Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security. New York: ACM, 2017: 506-519.
- [37] JUUTI M, SZYLLER S, MARCHAL S, et al. PRADA: protecting against DNN model stealing attacks [C]// Proceedings of the 2019 IEEE European Symposium on Security and Privacy. Piscataway: IEEE, 2019: 512-527.
- [38] OREKONDY T, SCHIELE B, FRITZ M. Knockoff nets: stealing functionality of black-box models [C]// Proceedings of the 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition. Piscataway: IEEE, 2019: 4954-4963.
- [39] DWORK C. Differential privacy [C]// Proceedings of the 33rd International Conference on Automata, Languages, and Programming. Berlin: Springer, 2006: 1-12.
- [40] GEYER R C, KLEIN T, NABI M. Differentially private federated learning: a client level perspective [EB/OL]. (2018-03-01) [2023-07-09]. <https://arxiv.org/pdf/1712.07557.pdf>.
- [41] WEI K, LI J, DING M, et al. Federated learning with differential privacy: algorithms and performance analysis [J]. IEEE Transactions on Information Forensics and Security, 2020, 15: 3454-3469.
- [42] TRUEX S, LIU L, CHOW K-H, et al. LDP-Fed: federated learning with local differential privacy [C]// Proceedings of the Third ACM International Workshop on Edge Systems, Analytics and Networking. New York: ACM, 2020: 61-66.
- [43] ZHAO Y, ZHAO J, YANG M, et al. Local differential privacy-based federated learning for internet of things [J]. IEEE Internet of Things Journal, 2020, 8(11): 8836-8853.
- [44] GOLDBREICH O. Secure multi-party computation [EB/OL]. (1998-06-11) [2023-07-09]. https://www.researchgate.net/publication/2934115_Secure_Multi-Party_Computation.
- [45] OGBURN M, TURNER C, DAHAL P. Homomorphic encryption [J]. Procedia Computer Science, 2013, 20: 502-509.
- [46] GENTRY C. Fully homomorphic encryption using ideal lattices [C]// Proceedings of the 41st Annual ACM Symposium on Theory of Computing. New York: ACM, 2009: 169-178.
- [47] LONGO D L, DRAZEN J M. Data sharing [J]. New England Journal of Medicine, 2016, 374(3): 276-277.
- [48] FANG H, QIAN Q. Privacy preserving machine learning with homomorphic encryption and federated learning [J]. Future Internet, 2021, 13(4): 94.
- [49] ZHANG C, LI S, XIA J, et al. BatchCrypt: efficient homomorphic encryption for cross-silo federated learning [C]// Proceedings of the 2020 USENIX Annual Technical Conference. Berkeley: USENIX Association, 2020: 493-506.
- [50] MA J, NAAS S-A, SIGG S, et al. Privacy-preserving federated learning based on multi-key homomorphic encryption [J]. International Journal of Intelligent Systems, 2022, 37(9): 5880-5901.
- [51] 陈宛桢, 张恩, 秦磊勇, 等. 边缘计算下基于区块链的隐私保护联邦学习算法 [J]. 计算机应用, 2023, 43(7): 2209-2216. (CHEN W Z, ZHANG E, QIN L Y, et al. Privacy-preserving federated learning algorithm based on blockchain in edge computing [J]. Journal of Computer Applications, 2023, 43(7): 2209-2216.)
- [52] 周炜, 王超, 徐剑, 等. 基于区块链的隐私保护去中心化联邦学习模型 [J]. 计算机研究与发展, 2022, 59(11): 2423-2436. (ZHOU W, WANG C, XU J, et al. Privacy-preserving and decentralized federated learning model based on the blockchain [J]. Journal of Computer Research and Development, 2022, 59(11): 2423-2436.)
- [53] 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): No. 12.
- [54] 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): 1-210.

This work is partially supported by National Natural Science Foundation of China (U20A20179).

CHEN Xuebin, born in 1970, Ph. D., professor. His research interests include big data security, IoT security, network security.

REN Zhiqiang, born in 2000, M. S. candidate. His research interests include data security, privacy protection.

ZHANG Hongyang, born in 1999, M. S. candidate. His research interests include data security, privacy protection.