

联邦学习浅析

王佳¹, 苗璐²

(1. 山西金融职业学院, 太原 030008; 2. 中国科学技术大学, 合肥 230000)

摘要:

机器学习的成功依赖于大量可用的训练数据。然而,现实中通常面临的情况是,数据分布在多个不同的数据源中,而要训练一个强大的模型需要使用所有数据源中的数据。因此一个直接的想法是,将多个数据源合并到一起。然而,由于通信成本,想要将分散在各个组织、机构的数据进行整合所需的代价是巨大的。而且本地数据的分享会泄露用户隐私。近年来,世界各地都在加强对数据隐私的保护,许多新法案对个人数据的存储和共享施加了严格限制。针对从多个数据源汇集数据的困难,联邦学习的概念被提出,目的是在满足数据隐私和监管要求的前提下,设计一个机器学习框架,使人工智能系统能够更加高效、准确地使用各自的数据。对联邦学习的研究背景、提出过程和定义进行简要的阐述,并且介绍近期的研究进展,以及联邦学习的相关概念。

关键词:

人工智能; 机器学习; 联邦学习

1 研究的背景

2016 年是人工智能(AI)成熟的一年,随着 AlphaGo 击败了顶尖的人类围棋选手,我们真正见证了人工智能的巨大潜力,并开始期待在许多应用领域出现更复杂、更尖端的人工智能技术,例如无人驾驶。目前,人工智能的成就就依赖于大量可用的标记数据。例如,AlphaGo 使用了 160000 个实际游戏的 3000 万步作为训练数据,ImageNet 数据集中则包含了超过 1400 万张图片^[1]。

随着 AlphaGo 的成功,人们自然希望大数据驱动的人工智能能很快在我们生活的方方面面实现。然而,现实情况却有些令人失望:在各个行业中,更多的应用领域中只有少量或低质量的数据,并且标记数据非常昂贵,特别是在需要人类专业知识的领域。因此,人工智能技术在这些行业的发展比较缓慢。

此外,特定任务所需的数据可能不会保存在一个地方。许多组织可能只有未标记的数据,而其他一些组织的标签数量可能非常有限。那么能否通过跨组织

传输数据,将来自多个站点的数据融合在一个共同的站点上? 事实上,在许多情况下,要打破数据源之间的障碍是非常困难的,甚至是不可能的。一般来说,任何人工智能项目所需的数据都涉及多种类型。例如,在人工智能驱动的产品推荐服务中,产品销售商有关于产品的信息和用户购买的数据,但没有描述用户购买能力和支付习惯的数据,这些数据可能存在于银行中。在大多数行业中,数据以孤岛的形式存在。由于行业竞争、隐私安全和复杂的管理程序,即使是同一公司不同部门之间的数据集成也面临着很大的阻力。几乎不可能将分散在全国各地的数据和机构整合起来。

此外,随着公众的数据安全和用户隐私的意识日益增强,数据隐私和安全已成为世界性的重大课题。有关公共数据泄露的新闻引起了媒体和公众的高度关注。例如,Facebook 最近的数据泄露引起了广泛的抗议。作为回应,世界各国都在加强保护数据安全和隐私的法律。例如,欧盟于 2018 年 5 月 25 日实施的《通用数据保护条例》(GDPR),旨在保护用户的个人隐私并提供数据安全。它要求企业在用户协议中使用清晰

明了的语言,并授予用户“被遗忘的权利”,即用户可以删除或撤回其个人数据。违反该法案的公司将面临严厉的罚款。中国 2017 年颁布的《网络安全法》和《民法通则》要求,互联网企业不得泄露或篡改其收集的个人信息,在与第三方进行数据交易时,它们需要确保拟议的合同遵守法律规定的保护义务。这些法规的建立将有助于建立一个更为安全的社会,然而,如何在满足数据隐私的前提下,为缺少相关数据的企业和机构建立有效、准确的人工智能模型,是一个重大挑战。

更具体地说,人工智能中传统的数据处理模型往往涉及到简单的数据事务模型,一方收集数据并将数据传输给另一方,另一方负责数据的清洗和融合。最后,第三方将采用集成的数据,并构建模型供其他方使用。模型通常是作为服务销售的最终产品。这一传统过程面临着上述新的数据法律法规的挑战。因此,我们面临着一个两难的境地,即我们的数据是以孤岛的形式存在的,但在许多情况下,我们被禁止在不同的地方收集、融合和使用这些数据进行处理。如何解决数据碎片化和隔离问题是当今人工智能研究者和实践者面临的一大挑战。

2 联邦学习的定义

为了克服这些挑战,Google 首先引入了联邦学习 (FL) 系统^[2]。谷歌的主要想法是基于分布在多个设备上的数据集构建机器学习模型,同时防止数据泄漏。最近的改进集中在克服统计数据挑战^[3]和提高联邦学习的安全性^[4]。也有一些研究致力于使联邦学习更加个性化^[5]。以上工作都集中在涉及分布式移动用户交互的设备的联邦学习上,其中大规模分布的通信成本、不平衡的数据分布以及设备可靠性是优化的主要因素。

此外,数据是按用户 ID 或设备 ID 进行分区的,因此在数据空间中是水平的。联邦学习与保护隐私的机器学习密切相关,因为它还考虑了分散协作学习环境中的数据隐私。为了将联邦学习的概念扩展到组织间的协作学习场景,我们将原来的“联邦学习”推广到一个通用概念,即所有隐私保护的分散协作机器学习技术。

假设有 N 个数据所有者 $\{F_1, \dots, F_N\}$, 他们拥有的数据分别是 $\{D_1, \dots, D_N\}$, 每个数据所有者都希望通过

整合各自的数据来训练一个机器学习模型。传统的方法是把所有的数据放在一起,使用 $D = D_1 \cup \dots \cup D_N$ 来训练模型 M_{SUM} 。而联邦学习系统是一个协作过程,在这个过程中,数据所有者协同训练一个 M_{FED} , 任何数据所有者 F_i 都不会将其数据 D_i 公开给其他人。此外, M_{FED} 的精度 V_{FED} 应该非常接近 M_{SUM} 的精度 V_{SUM} 的性能。形式上,设 δ 为非负实数;如果 $|V_{FED} - V_{SUM}| < \delta$ 则称该联邦学习算法有 δ 精度损失。

3 联邦学习研究的进展

联邦学习是人工智能当中发展较快的领域,研究成果层出不穷。接下来,本文将介绍近期的两项研究进展:联邦迁移学习和基于概率的联邦学习。

3.1 联邦迁移学习

迁移学习 (TL)^[6] 是一种为数据集较小或只有部分标签的应用提供解决方案的强大技术。近年来,将迁移学习技术应用于各个领域的研究工作已经取得了很大的进展,比如图像分类以及情绪分析。迁移学习的性能取决于领域之间的关联程度。直观地说,同一个数据联邦中的各方通常是来自同一行业或相关行业的组织,因此更容易进行知识传播。

联邦迁移学习 (FTL) 适用于两个数据集不仅在样本上不同,而且在特征空间上也不同的场景。假设有两个机构,一个是位于中国的银行,另一个是位于美国的电子商务公司。由于地域限制,两家机构的用户群有只一个小的交集。另一方面,由于业务的不同,双方的特征空间只有一小部分重叠。在这种情况下,迁移学习技术为联邦内的整个样本和特征空间提供解决方案。具体来说,就是使用有限的公共样本集,学习到两个特征空间的公共表示,随后用这个公共表示获得只有单侧特征的样本的预测。联邦迁移学习是对现有联邦学习系统的重要扩展,因为它处理的问题超出了现有的联邦学习算法的范围。

假设源域数据集为 $D_A := \{(x_i^A, y_i^A)\}_{i=1}^{N_A}$, 其中 $x_i^A \in R^a$, $y_i^A \in \{+1, -1\}$ 是第 i 个数据的标签。目标域数据集为 $D_B := \{x_j^B\}_{j=1}^{N_B}$, 其中 $x_j^B \in R^a$ 。 D_A 和 D_B 分别被两个私有方所拥有,不能泄露给对方。同时,假设存在一个有限的共享数据集 $D_{AB} := \{(x_i^A, x_i^B)\}_{i=1}^{N_{AB}}$, 并且 A 中

含有 B 的一部分标签 $D_C := \{(x_i^B, y_i^A)\}_{i=1}^{N_C}$, 其中 N_C 是可用的目标标签数量。不失一般性, 我们假设所有标签都在 A 方, 但这里的所有推导都可以适用于标签存在于 B 方的情况。通过使用 RSA 等加密技术屏蔽数据 ID, 可以在隐私保护设置中找到双方共享的样本 ID 集。我们假设 A 和 B 已经找到它们共享的样本 ID。在以上假设条件下, 目标是让双方建立一个迁移学习模型, 以便在不暴露私有数据的情况下, 尽可能准确地预测目标域方的标签。

近年来, 深层神经网络被广泛应用于迁移学习中, 来寻找隐含的迁移机制。在一般的场景中, A、B 双方通过两个神经网络产生各自的隐层表示: $u_i^A = \text{Net}^A(x_i^A)$, $u_i^B = \text{Net}^B(x_i^B)$, 其中 $u_i^A \in R^{N_A \times d}$, $u_i^B \in R^{N_B \times d}$, d 是隐层表示的维度。为了标记目标域, 通用的方法是引入一个预测函数 $\varphi(u_j^B) = \varphi(u_1^A, y_1^A, \dots, u_{N_A}^A, y_{N_A}^A, u_j^B)$ 。不失一般性, 假设 $\varphi(u_j^B)$ 线性可分, 即 $\varphi(u_j^B) = \Phi^A G(u_j^B)$ 。于是, 训练目标函数可写为:

$$\arg \min_{\Theta^A, \Theta^B} L_1 = \sum_i^{N_C} l_1(y_i^A, \varphi(u_i^B))$$

其中, Θ^A, Θ^B 分别是 $\text{Net}^A, \text{Net}^B$ 的训练参数。设 L_A, L_B 分别是 $\text{Net}^A, \text{Net}^B$ 的层数, 那么 $\Theta^A = \{\theta_l^A\}_{l=1}^{L_A}$, $\Theta^B = \{\theta_l^B\}_{l=1}^{L_B}$, 其中 θ_l^A, θ_l^B 是第 l 层的训练参数。 l_1 表示损失函数, 对于 logistic 损失, $l_1(y, \varphi) = \log(1 + \exp(-y\varphi))$ 。

另外, 我们希望最小化 A 和 B 的对齐误差:

$$\arg \min_{\Theta^A, \Theta^B} L_2 = - \sum_i^{N_{AB}} l_2(u_i^A, u_i^B)$$

其中, l_2 表示对齐误差。典型的对齐误差可以是 $-u_i^A(u_i^B)'$ 。

最终的目标函数为:

$$\arg \min_{\Theta^A, \Theta^B} L = L_1 + \gamma L_2 + \frac{\lambda}{2} (L_3^A + L_3^B)$$

其中, γ, λ 是权重参数, $L_3^A = \sum_l^{L_A} \|\theta_l^A\|_F^2$, $L_3^B = \sum_l^{L_B} \|\theta_l^B\|_F^2$ 是正则化项。

接下来要获取反向传播过程中更新 Θ^A, Θ^B 所需的梯度:

$$\frac{\partial L}{\partial \theta_i^j} = \frac{\partial L_1}{\partial \theta_i^j} + \gamma \frac{\partial L_2}{\partial \theta_i^j} + \lambda \theta_i^j$$

其中 $i \in \{A, B\}$ 。联邦迁移学习要保证 A 和 B 不泄露自己的数据, 因此需要隐私保护算法来计算损失函数和梯度。文献[7]中提供了一种新颖的方法, 将加性同态加密(HE)应用于神经网络的多方计算(MPC), 从而仅需要对神经网络进行最小的修改, 并且准确性几乎是无损的, 而大多数现有的安全深度学习框架在采用隐私保护技术时会失去一定的准确性。联邦迁移学习的未来工作可能包括采用该方法到其他需要隐私保护数据协作的深度学习系统, 并通过使用分布式计算技术继续提高算法的效率, 以及寻找成本较低的加密方案。

3.2 基于概率的联邦学习

联邦学习中的每个数据源是隔离的, 联邦学习算法在训练每个数据源上的本地模型和将它们提取为全局联邦模型之间进行迭代, 而无需显式地组合来自不同数据源的数据。典型的联邦学习算法需要访问本地存储的数据进行学习, 更极端的情况是访问本地数据预先训练的模型, 而不是数据本身。文献[8]解决的问题是, 将根据不同来源的数据独立训练的“遗留”模型组合成一个改进的联邦模型。

文中开发了一个基于概率的联邦学习框架, 称为贝叶斯非参数的神经网络联邦学习框架。假设每个数据服务器提供本地神经网络的权重, 这些权重通过该框架进行建模。然后使用一种推理方法, 合成一个更具表现力的全局网络, 这个过程无需额外的监督和数据汇集, 而且只需一个通信轮次。假设要么是本地数据, 要么是经过本地训练的模型可用。当数据可用时, 并行地为每个数据源训练本地模型。然后匹配不同数据源估计的局部模型参数(权重向量)构建全局网络。局部参数的匹配, 由贝塔-伯努利过程(BBP)控制。BBP 是一个模型, 允许局部参数匹配现有的全局参数, 或在现有的全局参数是差的匹配时, 创建新的全局参数。

以包含单个隐层的多层感知机为例, 假设已经训练出 J 个多层感知机, 分别拥有一个隐层。对于第 j 个感知机, $V_j^{(0)} \in R^{D \times L_j}$ 是隐层的权重, $\tilde{v}_j^{(0)} \in R^{L_j}$ 是隐层的偏置项, 其中 D 是数据维度, L_j 是隐层神经元的个

数。 $V_j^{(l)} \in R^{L_j \times K}$ 是 softmax 层的权重, $\tilde{v}_j^{(l)} \in R^K$ 是 softmax 层的偏置项, 其中 K 是类别数目。在拥有 J 个 $\{V_j^{(l)}, \tilde{v}_j^{(l)}, V_j^{(l)}, \tilde{v}_j^{(l)}\}$ 的情况下, 试图学习全局模型, 它的参数为: $\Theta^{(0)} \in R^{D \times L}$, $\tilde{\theta}^{(0)} \in R^L$, $\Theta^{(l)} \in R^{L \times K}$, $\tilde{\theta}^{(l)} \in R^K$, 其中 L 是全局模型的隐层神经元个数, 由推理得出。

算法的原理如图 1 所示, 三个本地多层感知机的隐层神经元经过匹配后, 形成全局模型。图中的节点表示神经元, 相同颜色的神经元已经匹配。

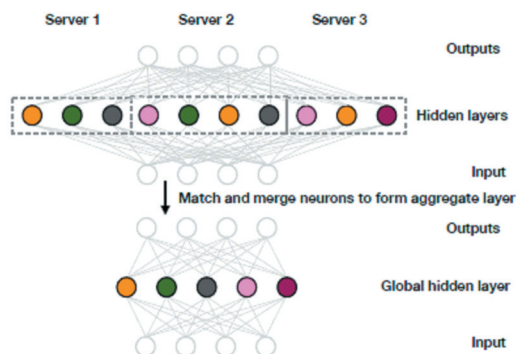


图1 单层概率联邦神经匹配算法原理示意图

文中提出的概率联邦神经匹配(PFNM)算法比现有方法有几个优点。首先, 它将局部模型的学习与局部模型合并为全局联邦模型的过程相分离。这种解耦允许我们对局部学习算法保持不可知的态度, 局部学习算法可以根据需要进行调整, 每个数据源甚至可能使用不同的学习算法。此外, 只要给定预先训练的模型, PFNM 就能够将它们组合成一个联邦全局模型, 而不需要额外的数据或关于生成预训练模型的算法的知识。而现有的神经网络联邦学习的方法需要关于局部学习的强假设, 例如, 共享相同的随机初始化, 这个假设在很多情况下是不现实的。并且, 不同于现有的方法, 文中提出的框架本质上是非参数的, 允许联邦模型灵活地增大或缩小其复杂性(即神经元的数目), 以考虑不同的数据复杂性。

4 联邦学习的相关概念

联邦学习使多方能够协同构建机器学习模型, 同时保持其私有训练数据的私有性。联邦学习作为一门新兴的技术, 有着许多独创性的思想, 其中一些思想植根于现有的领域。下面, 我们将从多个角度解释联邦

学习与其他相关概念之间的关系。

4.1 隐私保护的机器学习

联邦学习可以看作是一种隐私保护、分散协作的机器学习。过去, 许多研究工作都致力于多方、隐私保护的机器学习这一领域。例如, 文献[9]的作者提出了用于垂直分区数据的安全多方决策树的算法。Vaidya 和 Clifton 提出了安全关联挖掘规则^[10]、安全 K-means^[11]和朴素贝叶斯分类器^[12]。文献[13]的作者提出了一种基于水平分块数据的关联规则算法。文献[14]的作者提出了多方线性回归和分类的安全协议。文献[15]的作者提出了安全的多方梯度下降方法。这些作品都使用了安全多方计算(SMC)来保证隐私。

4.2 分布式机器学习

联邦学习与分布式机器学习有点相似。分布式机器学习包括很多方面, 如训练数据的分布式存储、计算任务的分布式操作、模型结果的分布式分布等。参数服务器^[16]是分布式机器学习中的一个典型元素。参数服务器作为一种加速训练过程的工具, 将数据存储在分布式工作节点上, 通过一个中心调度节点来分配数据和计算资源, 从而更有效地训练模型。对于联合学习, 工作节点表示数据所有者, 对本地数据具有完全的自主权, 可以决定何时以及如何加入联邦学习。在参数服务器中, 中心节点总是起控制作用。然而, 联邦学习面临着一个更加复杂的学习环境。此外, 在模型训练过程中, 联邦学习强调数据所有者的数据隐私保护。有效的数据隐私保护措施可以更好地应对未来日益严格的数据隐私和数据安全监管环境。

5 结语

最近, 数据的隔离和数据隐私保护成为人工智能面临的下一个挑战, 联邦学习给我们带来了新的希望。它可以在保护本地数据的同时, 为多个机构建立统一的模型, 使多个机构能够在数据安全的基础上协同工作。本文简述了联邦学习的基本定义、提出背景和研究进展, 包括联邦迁移学习和基于概率的联邦学习, 最后介绍了联邦学习的相关概念。预计在将来, 联邦学习将打破行业之间的壁垒, 使数据和知识可以安全地共享, 并根据每个参与者的贡献公平地分配收益。联邦学习的发展将会促进人工智能应用到我们生活的每个角落。

参考文献:

- [1]YANG Qiang, LIU Yang, CHEN Tianjian. Federated Machine Learning:Concept and Applications[J]. ACM Transactions on Intelligent Systems and Technology (TIST), 2019, 10(2):1-19.
- [2]Jakub Konecny, H. Brendan McMahan, Daniel Ramage, Peter Richtarik. Federated optimization:Distributed Machine Learning for on-Device Intelligence. CoRR abs/1610.02527, 2016. arxiv:1610.02527. <http://arxiv.org/abs/1610.02527>.
- [3]ZHAO Yue, LI Meng, LAI Liangzhen, Naveen Suda, Damon Civin, Vikas Chandra. Federated Learning with Non-IID Data, 2018. arxiv:cs.LG/1806.00582.
- [4]Robin C, Geyer, Tassilo Klein, Moin Nabi. Differentially Private Federated Learning:A Client Level Perspective. CoRR abs/1712.07557, 2017. arxiv:1712.07557. <http://arxiv.org/abs/1712.07557>.
- [5]CHEN Fei, DONG Zhenhua, LI Zhenguo, HE Xiuqiang. Federated Meta-Learning for Recommendation. CoRR abs/1802.07876, 2018. arxiv:1802.07876. <http://arxiv.org/abs/1802.07876>.
- [6]Sinno Jialin Pan, YANG Qiang. A Survey on Transfer Learning. IEEE Trans. Knowl. Data Eng., 2010, 22(10):1345-1359. DOI:<https://doi.org/10.1109/TKDE.2009.191>.
- [7]LIU Yang, CHEN Tianjian, YANG Qiang. Secure Federated Transfer Learning[J]. arXiv Preprint, arXiv:1812.03337, 2018.
- [8]Yurochkin M, Agarwal M, Ghosh S, et al. Bayesian Nonparametric Federated Learning of Neural Networks[J], 2019.
- [9]Jaideep Vaidya, Chris Clifton. Privacy-Preserving Decision Trees Over Vertically Partitioned Data. Data and Applications Security XIX, Sushil Jajodia and Duminda Wijesekera (Eds.). Springer, Berlin, 2005:139 - 152.
- [10]Jaideep Vaidya, Chris Clifton. Privacy Preserving Association Rule Mining in Vertically Partitioned Data. Proceedings of the 8th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'02). ACM, New York, NY, 2002:639 - 644. DOI: <https://doi.org/10.1145/775047.775142>.
- [11]Jaideep Vaidya, Chris Clifton. Privacy-Preserving K-Means Clustering over Vertically Partitioned Data. Proceedings of the 9th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'03). ACM, New York, NY, 2003:206-215. DOI: <https://doi.org/10.1145/956750.956776>.
- [12]Jaideep Vaidya, Chris Clifton. Privacy Preserving Naive Bayes Classifier for Vertically Partitioned Data. Proceedings of the 4th SIAM Conference on Data Mining, 2004:330 - 334.
- [13]Murat Kantarcioglu, Chris Clifton. 2004. Privacy-Preserving Distributed Mining of Association Rules on Horizontally Partitioned Data. IEEE Trans. on Knowl. and Data Eng., 2004, 16(9):1026 - 1037. DOI:<https://doi.org/10.1109/TKDE.2004.45>.
- [14]DU Wenliang, Yunghsiang Sam Han, CHEN Shigang. 2004. Privacy-Preserving Multivariate Statistical Analysis:Linear Regression and Classification. SDM, 2004(4):222 - 233.
- [15]LiWan, Wee Keong Ng, Shuguo Han, Vincent C S Lee. 2007. Privacy-Preservation for Gradient Descent Methods. Proceedings of the 13th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'07). ACM, New York, NY, 2007:775 - 783. DOI:<https://doi.org/10.1145/1281192.1281275>.
- [16]Qirong Ho, James Cipar, Henggang Cui, Jin Kyu Kim, Seunghak Lee, Phillip B, Gibbons, Garth A, Gibson, Gregory R, Ganger, Eric P Xing. 2013. More Effective Distributed ML via a Stale Synchronous Parallel Parameter Server. Proceedings of the 26th International Conference on Neural Information Processing Systems-Volume 1 (NIPS'13). Curran Associates Inc., 2013:1223-1231. <http://dl.acm.org/citation.cfm?id=2999611.2999748>.

作者简介:

王佳(1985-),女,山西新绛人,硕士,讲师,研究方向为计算机应用

苗璐(1997-),山西长治人,研究生在读,研究方向为计算机应用

收稿日期:2020-05-08

修稿日期:2020-07-08

(下转第 36 页)

Space-Time Modeling Method Based on Information Physics Fusion System

JIANG Yi-xun, ZHANG Li-chen

(School of Computer, Guangdong University of Technology, Guangzhou 510006)

Abstract:

Petri net is a graph-based formal modeling theory, which has been widely used in the analysis of concurrent systems. We will use Petri nets, automata and process algebra to provide mathematical formalism to verify the behavior of the space-time model, analyze its processes, conduct state space searches and state reachability analysis. The space-time modeling of our cyber-physical fusion system needs to consider both time and space specifications. As an important tool for modeling and analysis of distributed, parallel and real-time systems, time Petri nets will provide a solid foundation for CPS modeling. We will analyze the characteristics of the spatiotemporal characteristics in the model, and also compare different modeling methods of spatiotemporal characteristics, and finally decide to choose Petri nets as the modeling method, create a new space-time Petri net.

Keywords:

Cyber-Physical Systems; Petri Net; Space-Time Characteristic

(上接第 31 页)

Brief Introduction of Federated Learning

WANG JIA¹, MIAO Lu²

(1. Shanxi Vocational College of finance, Taiyuan 030008; 2. China University of Science and Technology, Hefei 230000)

Abstract:

The success of machine learning depends on the availability of a large number of training data. However, in reality, we usually face the situation that the data is distributed in many different data sources. To train a powerful model, we need to use the data from all data sources. So, a direct idea is to combine multiple data sources together. However, due to the cost of communication, integrating the data scattered in various organizations and institutions costs too much. And the sharing of local data will compromise user privacy. In recent years, the protection of data privacy has been enhanced all over the world. Many new laws impose strict restrictions on the storage and sharing of personal data. In view of the difficulty of collecting data from multiple data sources, the concept of Federated learning is proposed. The purpose is to design a machine learning framework on the premise of meeting data privacy and regulatory requirements, so that the AI system can use their data more efficiently and accurately. In this paper, the research background, process and definition of Federated learning are briefly described, and the recent progress and related concepts of Federated learning are introduced.

Keywords:

Artificial Intelligence, Machine Learning, Federated Learning