Emerging Trends in Federated Learning: From Model Fusion to Federated X Learning

Shaoxiong Ji¹, Teemu Saravirta¹, Shirui Pan², Guodong Long³, and Anwar Walid^{4,5}

¹ Aalto University
² Monash University
³ University of Technology Sydney
⁴ Nokia Bell Labs
⁵ Columbia University

Email: {shaoxiong.ji, teemu.saravirta}@aalto.fi, shirui.pan@monash.edu, guodong.long@uts.edu.au, anwar.walid@nokia-bell-labs.com

Abstract

Federated learning is a new learning paradigm that decouples data collection and model training via multi-party computation and model aggregation. As a flexible learning setting, federated learning has the potential to integrate with other learning frameworks. We conduct a focused survey of federated learning in conjunction with other learning algorithms. Specifically, we explore various learning algorithms to improve the vanilla federated averaging algorithm and review model fusion methods such as adaptive aggregation, regularization, clustered methods, and Bayesian methods. Following the emerging trends, we also discuss federated learning in the intersection with other learning paradigms, termed as federated x learning, where x includes multitask learning, meta-learning, transfer learning, unsupervised learning, and reinforcement learning. This survey reviews the state of the art, challenges, and future directions.

Keywords— Federated Learning, Model Fusion, Learning Algorithms

1 Introduction

Vast quantities of data are required for state-of-the-art machine learning algorithms. However, the data cannot be uploaded to a central server or cloud due

to sheer volume, privacy, or legislative reasons. Federated learning (FL) [1], also known as collaborative learning, has been a subject of many studies. FL adopts a distributed machine learning architecture with a central server for model aggregation, where clients themselves update the machine learning model. Clients can maintain ownership of their data, i.e., upload only the updated model to the central server and not expose any of their private data.

The federated learning paradigm addresses several challenges. The first challenge is privacy. Local data ownership inherits a basic level of privacy. However, federated learning systems can be vulnerable to model poisoning [2]. The second challenge is the communication cost for model uploading and downloading. Improving communication efficiency is a critical issue [3, 4]. Centralized network architecture also makes the central server suffer from heavy communication workload, calling for a decentralized server architecture. [5]. The third challenge is statistical heterogeneity. Aggregating clients' models together can result in a nonoptimal combined model as client data is often non-IID (independent and identically distributed). Statistical heterogeneity introduces a degree of uncertainty into the learning model. Therefore, adopting the right aggregation and learning techniques is vital for robust implementation. This survey gives a particular focus on how different federated learning solutions address statistical heterogeneity.

The robust model aggregation has recently garnered considerable attention. Traditionally, client contributions are weighed according to their sample quantity, while recent research has introduced adaptive weighting [6, 7], attentive aggregation [8], regularization [9], clustering [10], and Bayesian methods [11]. Many methods generally attempt to derive client characteristics by adjusting the relative weights better. Aggregation in the federated setting has also addressed fairness [12] in taking underrepresented clients and classes better into account.

Statistical heterogeneity, or *non-IID data*, leads to the difficulties of choosing models and performing hyperparameter tuning, as the data resides at clients, out of the reach of a preliminary analysis. The edge clients provide the supervision signal for supervised machine learning models. However, the lack of human annotation or interaction between humans and learning systems induces the *label scarcity* and leads to a more restricted application domain.

Label scarcity is one of the problems emblematic to the federated setting. The inability to access client data and the resulting black-box updates are tack-led by careful selection of the aggregation method and supplementary learning paradigms to fit specific real-world scenarios. As a result of label scarcity, the semi-supervised and unsupervised learning paradigms introduce essential techniques to deal with the uncertainty arising from unlabeled data. Faced with the problem that clients' local models can diverge during multiple epochs of local training, the server can be tasked with selecting the *most reliable* client models of the preceding round, regularizing the aggregation for achieving consistency. Fully unsupervised data can be enhanced via domain adaption, where the aim is to transfer knowledge from a labeled domain to an unlabeled one.

Taxonomy. To establish critical solutions for problems arising from private, non-IID data, we assess the current leading solutions in model fusion and how other learning paradigms are incorporated into the federated learning scenario. We propose a novel taxonomy of federated learning according to the model fusion principle and the connection to other learning paradigms. The taxonomy scheme as illustrated in Table 1 with some representative instantiations is organized as below.

- Federated Model Fusion. We categorize the major improvements to the pioneering FedAvg model aggregation algorithm into four subclasses (i.e., adaptive/attentive methods, regularization methods, clustered methods, and Bayesian methods), together with a special focus on fairness (Section 3).
- Federated Learning Paradigms. We investigate how the various learning paradigms are fitted into the federated learning setting (Section 4). The learning paradigms include some key supervised learning scenarios such as transfer learning, multitask and meta-learning, and learning algorithms beyond supervised learning such as semi-supervised learning, unsupervised learning, and reinforcement learning.

Contributions. This survey starts from a novel viewpoint of federated learning by coupling federated learning with different learning algorithms. We propose a new taxonomy and conduct a timely and focused survey of recent advances on solving the heterogeneity challenge. Our survey's distinction compared with other comprehensive surveys is that we focused on the emerging trends of federated model fusion and learning paradigms, which are not intensively discussed in previous surveys. Besides, we connect these recent advances with real-word applications and discuss limitations and future directions in this focused context.

This survey is organized as follows. In Section 3, we assess in detail the significant improvements recent research has proposed on top of the pioneering FedAvg model aggregation algorithm [1]. In Section 4, we analyze how the various learning paradigms are fitted into the federated learning setting. In Section 5, we highlight recent successes in applied federated learning. Finally, in Section 6, we outline future research directions specifically from the viewpoint of model fusion and complementary learning paradigms. This paper is a focused survey, assessing only the aforementioned coupled subfields, of which learning paradigms makes the learned models more robust, and model fusion brings those models together. For a more wide-ranging survey into federated learning, we recommend readers to refer [47, 48, 49].

2 Related Survey

Several related surveys have been published in recent years as summarized in Table 2. This section introduces the existing surveys and highlight our survey's contributions to the literature.

Main area	Subarea	Study
Federated Model Fusion	Adaptive/Attentive Aggregation	IDA [6], ASTW [7], FedAtt [8] FedAttOpt [13], FedMed [14], FedAMP [15].
	Regularization Methods	FedAwS [16], FedProx [17] Mime [18]
	Clustered Methods	FL+HC [10], IFCA [19] FeSEM [20], FedFast [21]
	Bayesian Methods	FedMA [22], PFNM [11]
	Fairness	q-FFL [12], AFL [23]
Learning Paradigms	Transfer Learning	FTL [24], * UFDA [25], FedSteg [26]
	Multitask and Meta Learning	Mocha [27], Kernelized FMTL [28] CFL [29], FedMeta [30], Per-FedAvg [31]
	Knowledge	FedMD [32], FedGKT [33], FedDF [34]
	Distillation	PATE [35]
	Semi-Supervised Learning	FedMatch [36]
		PATE-G [35]
	Generative Adversarial Learning	Sync. Strategies [37], FedGAN [38] DP-FedAvg-GAN [39]
	Unsupervised	* UFDA [25]
	Learning	FURL [40], FPCA [41], FedCA [42]
	Reinforcement	FedRL [43], DRL-based Aggregator [44]
	Learning	Favor [45], FRD and MixFRD [46]

Table 1: Federated learning with other learning algorithms: categorization, conjunctions and representative methods.

General Survey of Federated Learning Yang et al. [47] firstly defined the concepts of federated learning, introduced federated applications, and discussed data privacy and security aspects. Li et al. [48] systematically reviewed the federated learning building blocks, including data partitioning, machine learning model, privacy mechanism, communication architecture, the scale of the federation, and motivation of federation. Kairouz et al. [49] detailed definitions of federated learning systems variations. Li et al. [9] discussed core challenges of federated learning in communication efficiency, privacy, and some future research directions

Domain-specific Survey Other surveys review a specific domain. Xu et al. [50] surveyed the healthcare and medical informatics domain. Lyu et al. [51] discussed

about the security threats and vulnerability challenges dealing with adversaries in federated learning systems Lim et al. [52] focused on the mobile edge networks. Niknam et al. [53] reviewed federated learning in the context of wireless communications, covering the data security and privacy challenges, algorithm challenges, and wireless setting challenges Jin et al. [54] conducted a review on federated semi-supervised learning. Jin et al.'s survey is the most related work to our paper. However, it only concentrates on semi-supervised learning. Our paper fills in its gap by including a wider range of model fusion and learning algorithms.

Table 2: Comparison of related survey articles about federated learning

Publication	Scope
This survey	Learning algorithms
Jin et al. [54]	Semi-supervised learning
Xu et al. [50]	Healthcare informatics
Lo et al. [55]	Software engineering
Lim et al. [52]	Mobile edge networks
Lyu et al. [51]	Threats
Niknam et al. [53]	Wireless communication
Yang et al. [47]	General
Li et al. [48]	General
Kairouz et al. [49]	General
Li et al. [9]	General

Distinction of Our Survey Our paper reviews the emerging trends of federated learning from a unique and novel angle, i.e., the learning algorithms used in the federated learning paradigms, including the model fusion algorithms (Sec. 3) and the conjunction of federated learning and other learning paradigms (named as Federated X Learning in Sec. 4). This unique perspective has not been well-discussed in any of the aforementioned surveys. Our survey fills in this gap by reviewing recent publications. Besides, we point out challenges and outlook future directions in this specific category of research on federated learning.

3 Federated Model Fusion

3.1 Overview

The goal of federated learning is to minimize the empirical risks over local data as

$$\min_{\theta} f(\theta) = \sum_{k=1}^{m} p_k \mathcal{L}_k(\theta)$$
 (1)

where \mathcal{L}_k is the local objective of the *k*-th client and $\sum_k p_k = 1$. The widely applied federated learning algorithm, i.e., Federated Averaging (FedAvg) [1], starts with a random initialization or warmed-up model of clients followed by local training, uploading, server aggregation, and redistribution. The learning objective is configured by setting p_k to be $\frac{n_k}{\sum_k n_k}$. Federated averaging assumes a regularization effect, similar to dropout in neural networks, by randomly selecting a fraction of clients on each communication round. Sampling on each round leads to faster training without a significant drop in accuracy. Li et al. [56] conducted a theoretical analysis on the convergence of FedAvg without strong assumptions and found that sampling and averaging scheme affects the convergence. Recent studies investigate some significant while less considered problems and explore different possibilities of improving vanilla averaging. To mitigate the client drift caused by heterogeneity in FedAvg, the SCAFFOLD algorithm [57] estimates the client drift as the difference between the update directions of the server model and each client model and adopt stochastically controlled averaging the correct client drift. Reddi et al. [58] proposed adaptive optimization algorithms such as Adagrad and Adam to improve the standard federated averaging-based optimization with convergence guarantees. Singh et al. [59] adopted optimal transport, which minimizes the transportation cost of neurons, to conduct layer-wise model fusion.

3.2 Adaptive Weighting

The adaptive weighting approach calculates adaptive weighted averaging of model parameters as:

$$\theta_{t+1} = \sum_{k=1}^{K} \alpha_k \cdot \theta_t^{(k)},\tag{2}$$

where $\theta_t^{(k)}$ is current model parameter of k-th client, θ_{t+1} is the updated global model parameter after aggregation, and α_k is the adaptive weighting coefficient. Aiming to train a low variance global model with non-IID robustness, Yeganeh et al. [6] proposed an adaptive weighting approach called Inverse Distance Aggregation (IDA) by extracting meta information from the statistical properties of model parameters. Specifically, the weighting coefficient with inverse distance is calculated as:

$$\alpha_k = \|\theta_t - \theta_t^{(k)}\|^{-1} / (\sum_{k=1}^K \|\theta_t - \theta_t^{(k)}\|^{-1}).$$
 (3)

Considering the time effect during federated communication, Chen et al. [7] proposed temporally weighted aggregation of the local models on the server as:

$$\theta_{t+1} = \sum_{k=1}^{K} \frac{n_k}{n} (\frac{e}{2})^{-(t-t^{(k)})} \theta_t^{(k)}, \tag{4}$$

where e is the natural logarithm, t is the current update round and $t^{(k)}$ is the update round of the newest $\theta^{(k)}$.

3.3 Attentive Aggregation

The federated averaging algorithm takes the instance ratio of the client as the weight to calculate the averaged neural parameters during model fusion [1]. In attentive aggregation, the instance ratio is replaced by adaptive weights as Eq. 5:

$$\theta_{t+1} \leftarrow \theta_t - \epsilon \sum_{k=1}^m \alpha_k \nabla \mathcal{L}(\theta_t^{(k)}),$$
 (5)

where α_k is the attention scores for client model parameters. FedAtt [8] proposes a simple layer-wise attentive aggregation scheme that takes the server model parameter as the query. FedAttOpt [13] enhance the attentive aggregation of FedAtt by scaled dot product. Like attentive aggregation, FedMed [14] proposes an adaptive aggregation algorithm using Jensen-Shannon divergence as the non-parametric weight estimator. These three attentive approaches use centralized aggregation architecture with only one shared global model for client model fusion. Huang et al. [15] studied pairwise collaboration between clients and proposed FedAMP with attentive message passing among similar personalized cloud models of each client.

3.4 Regularization Methods

We summarize federated learning algorithms with additional regularization terms to client learning objectives or server aggregation formulas. One category is to add local constraints for clients. FedProx [17] adds proximal terms to clients' objectives to regularize local training and ensure convergence in the non-IID setting. After removing the proximal term, FedProx degrades to FedAvg. Another direction is to conduct federated optimization on the server side. Mime [18] adapts conventional centralized optimization algorithms into federated learning and uses momentum to reduce client drift with only global statistics as

$$\mathbf{m}_t = (1 - \beta) \nabla f_i(\mathbf{x}_{t-1}) + \beta \mathbf{m}_{t-1}$$
(6)

 \mathbf{m}_{t-1} is a moving average of unbiased gradients computed over multiple clients. Federated averaging may lead to class embedding collapse to a single point for embedding-based classifiers. To tackle the embedding collapse, Yu et al. [16] studied the federated setting where users only have access to a single class, for example, face recognition in the mobile phone. They proposed the FedAwS framework with a geometric regularization and stochastic negative mining over the server optimization to spread class embedding space.

3.5 Clustered Methods

We formulate clustered methods as algorithms that take additional steps with client clustering before federated aggregation or optimization to improve model fusion. One straightforward strategy is the two-stage approach, for example, the

clustering then aggregation scheme. Briggs et al. [10] propose to take an additional hierarchical clustering for client model updates and apply federated averaging for each cluster. Diverting client updates to multiple global models from user groups can help better capture the heterogeneity of non-IID data. Xie et al. [20] proposed multi-center federated learning, where clients belong to a specific cluster, clusters update along with the local model updates, and clients also update their belongings to different clusters. The authors formulated a joint optimization problem with distance-based multi-center loss and proposed the FeSEM algorithm with stochastic expectation maximization (SEM) to solve the optimization. Muhammad et al. [21] proposed an active aggregation method with several update steps in their FedFast framework going beyond average. The authors worked on recommendation systems and improved the conventional federated averaging by maintaining user-embedding clusters. They designed a pipelined updating scheme for item embeddings, client delegate embeddings, and subordinate user embeddings to propagate client updates in the cluster with similar clients. Ghosh et al. [19] formulated clustered federated learning by partitioning different user groups with the same learning tasks and conducting aggregation within the cluster partition. The authors proposed an Iterative Federated Clustering Algorithm (IFCA) with alternate cluster identity estimation and model optimization to capture the non-IID nature.

3.6 Bayesian Methods

Bayesian non-parametric machinery is applied to federated deep learning by matching and combining neurons for model fusion. Yurochkin et al. [11] proposed probabilistic federated neural matching (PFNM) using a Beta Bernoulli Process to model the multi-layer perceptron (MLP) weight parameters. Observing the permutation invariance of fully connected layers, the proposed FGNM algorithm first matches the neurons of neural models of clients to the global neurons. It then aggregates via maximum a posteriori estimation of global neurons. However, the authors only considered simple MLP architectures. FedMA [22] extends PFNM to convolutional and recurrent neural networks by matching and averaging hidden elements, specifically, channels for CNNs and hidden units for RNNs. It solves the matched averaging objective by iterative optimization.

3.7 Fairness

When aggregating the global shared model, FedAvg applies a weighted average concerning the number of samples that participating clients used in their training. However, the model updates can easily skew towards an over-represented subgroup of clients where super-users provide the majority of samples. Mohri et al. [23] suggested that valuing each sample without clear discrimination is inherently risky as it might result in sub-optimal performance for underrepresented clients and sought to good-intent fairness to ensure federated training not overfitting to some of the specific clients. Instead of the uniform distribution in clas-

sic federated learning, the authors proposed agnostic federated learning (AFL) with minimax fairness, which takes a mixture of distributions into account. However, the overall tradeoff between fairness and performance is still not well explored. Inspired by fair resource allocation in wireless networks, the q-fair federated learning (q-FFL) [12] proposes an optimization algorithm to ensure fair performance, i.e., a more uniform distribution of performance gained in federated clients. The optimization objective (Eq. 7) adjusts the traditional empirical risk objective by tunable performance-fairness tradeoff controlled by q.

$$\min_{\theta} f_q(\theta) = \sum_{k=1}^m \frac{p_k}{q+1} \mathcal{L}_k^{q+1}(\theta)$$
 (7)

The flexible q-FFL also generalizes well to previous methods; specifically, it reduces to FedAvg and AFL when the value of q is set to 0 and ∞ respectively.

4 Federated X Learning

The customizability of federated learning objective leads to possibilities in quickly adapting FL to adversarial, semi-supervised, or reinforcement learning settings, offering flexibility to other learning algorithms in conjunction with federated learning. We term FL's intersection with other learning algorithms as Federated X Learning.

4.1 Federated Transfer Learning and Knowledge Distillation

Transfer learning focuses on transferring knowledge from one particular problem to another, and it has also been integrated into federated learning to construct a model from two datasets with different samples and feature spaces [47]. Liu et al. [24] formulated the Federated Transfer Learning (FTL) to solve the problem that traditional federated learning falters when datasets do not share sufficient common features or samples. The authors also enhance the security with homomorphic encryption and secret sharing. In real-world applications, FedSteg [26] applies federated transfer learning for secure image steganalysis to detect hidden information. Alawad et al. [60] utilized federated transfer learning without sharing vocabulary for privacy-preserving NLP applications for cancer registries.

Knowledge Distillation Given the assumption that clients have sufficient computational capacity, federated averaging adopts the same model architecture for clients and the server. FedMD [32] couples transfer learning and knowledge distillation (KD), where the centralized server does not control the architecture of models. It introduces an additional public dataset for knowledge distillation, and each client optimizes their local models on both public and private data. Strictly speaking, transfer learning differs from knowledge distillation; however, the FedMD framework puts them under one umbrella. Many technical details are

only briefly introduced in the original paper of FedMD. Recently, He et al. [33] utilized knowledge distillation with technical solidity to train computationally affordable CNNs for edge devices via knowledge distillation. The authors proposed the Group Knowledge Transfer (FedGKT) framework that optimizes the client and the server model alternatively with knowledge distillation loss. Specifically, the larger server model takes features from the edge to minimize the gap between periodically transferred ground truth and soft label predicted by the edge model, and the small model distills knowledge from the larger server model by optimizing the KD-loss using private data and soft labels transferred back from the server. However, this framework has a potential risk of privacy breach as the server holds the ground truth, especially when ground truth labels are user's typing records in the mobile keyboard application. Lin et al. [34] applied knowledge distillation to mitigate privacy risk and cost and proposed a novel ensemble distillation for model fusion that utilizes unlabeled data.

4.2 Federated Multitask and Meta Learning

This section takes multitask learning and meta-learning under the same category coupled with federated learning, where different clients adopt different models at inference time.

Federated Multitask Learning trains separate models for each client with some shared structure between models, where learning from local datasets at different clients is regarded as a separate task. In contrast to federated transfer learning between two parties, federated multitask learning involves multiple parties and formulates similar tasks clustered with specific constraints over model weights. It exploits related tasks for more efficient learning to tackle the statistical heterogeneity challenge. The Mocha framework [27] trains separate yet related models for each client by solving a primal-duel optimization. It leverages a shared representation across multiple tasks and addresses the challenges of data and system heterogeneity. However, the Mocha framework is limited to regularized linear models. Caldas et al. [28] further studied the theoretical potential of kernelized federated multitask learning to solve the non-linearity. To solve the suboptimal results, Sattler et al. [29] studied geometric properties of the federated loss surface. They proposed a federated multitask framework with non-convex generalization to cluster the client population.

Federated Meta Learning aims to train a model that be quickly adapted into new tasks with few training data, where clients serve as a variety of learning tasks. The seminal model-agnostic meta-learning (MAML) framework [61] has been intensively applied to this learning scenario. Several studies connect FL and meta-learning, for example, model updating algorithm with average difference descent [62] inspired by the first-order meta-learning algorithm. However, this study focuses on applications in the social care domain with less consideration in practical settings. Jiang et al. [63] further provided a unified view of federated

meta-learning to compare MAML and the first-order approximation method. Inspired by the connection between federated learning and meta-learning, Fallah et al. [31] adapted MAML into the federated framework Per-FedAvg, to learn an initial shared model, leading to fast adaption and personalization for each client. FedMeta [30] proposes a two-stage optimization with a controllable meta updating scheme after model aggregation as:

$$\theta_{t+1}^{meta} = \theta_{t+1} - \eta_{meta} \nabla_{\theta_{t+1}} \mathcal{L}(\theta_{t+1}; \mathcal{D}_{meta}), \tag{8}$$

where \mathcal{D}_{meta} is a small set of meta data on the server.

4.3 Federated Generative Adversarial Learning

Generative Adversarial Networks (GANs) consist of two competing models, i.e., a generator and a discriminator. The generator learns to produce samples approximating the underlying ground-truth distribution. The discriminator, usually a binary classifier, tries to distinguish the samples produced by the generator from the real samples. A straightforward combination with FL is to have the GAN models trained locally on clients and the global model fused with different strategies. Fan and Liu [37] studied the synchronization strategies for aggregating discriminator and generator networks on the server and conducted a series of empirical analyses. Updating clients on each round with both the generator and the discriminator models achieves the best results; however, it is twice as computationally expensive as just syncing the generator. Updating just the generator leads to almost equivalent performance than updating both, whereas updating just the discriminator leads to considerably worse performance, closer to updating neither. Rasouli et al. [38] extended the federated GAN with different applications and proposed the FedGAN framework to use an intermediary for averaging and broadcasting the parameters of generator and discriminator. Furthermore, the authors studied the convergence of distributed GANs by connecting the stochastic approximation and communication-efficient SGD optimization for GAN and federated learning. Augenstein et al. [39] proposed differentially private federated generative models to address the challenges of non-inspectable data scenario. GANs are adopted to synthesize realistic examples of the private data for data labeling inspection at inference time.

4.4 Federated Semi-supervised Learning

Private data at a client might be partly or entirely unlabeled. Semi-supervised learning sets learning from labeled and unlabeled data, with unlabeled data comprising a much larger portion than labeled data. When combined with federated learning, it leads to a new learning setup, i.e., federated semi-supervised learning (FSSL), which is more realistic as users may not annotate all the data in their devices. Similar to centralized semi-supervised learning, FSSL also utilizes a two-part loss function on device with the loss stemming from supervised learning $\mathcal{L}_s(\theta)$ and the loss from unsupervised learning $\mathcal{L}_u(\theta)$. Jeong et al. [36] proposed

a federated matching (FedMatch) framework with inter-client consistency loss to exploit the heterogenous knowledge learned by multiple client models. The authors showed that learning on both labeled and unlabeled data simultaneously may result in the model forgetting what it had learned from labeled data. To counter this, the authors decomposed the model parameters θ to two variables $\theta = \psi + \rho$ and utilized a separate updating strategy, where only ψ is updated during unsupervised learning, and similarly, ρ is updated for supervised learning. Semi-supervised learning also couples with teacher-student learning for learning from private data. Papernot et al. [35] put forward a semi-supervised approach with a private aggregation of teacher ensembles (PATE), an architecture where each client votes on the correct label. PATE was shown empirically to particularly beneficial when used in conjunction with GANs.

4.5 Federated Unsupervised Learning

It is more common that local clients host no labeled data, which naturally leads to the learning paradigm of federated unsupervised learning without supervision in the decentralized learning scenario. A straightforward solution is to pretrain unlabeled data to learn useful features and utilized pretrained features in downstream tasks of federated learning systems [40]. There exist two challenges in federated unsupervised learning, i.e., the inconsistency of representation spaces due to data distribution shift and misalignment of representations due to the lack of unified information among clients. FedCA [42] proposes a federated contrastive averaging algorithm with the dictionary and alignment modules for client representation aggregation and alignment, respectively. The local model training utilizes the contrastive loss and the server aggregates models and dictionaries from clients. Recently, many unsupervised learning methods such as Principal Component Analysis (PCA) and unsupervised domain adaptation have been adopted to combine with federated learning. Peng et al. [25] studied the federated unsupervised domain adaptation that aligns the shifted domains under federated setting with a couple of learning paradigms. Specifically, unsupervised domain adaptation is explored by transferring labeled source domain to unlabelled target domain, and adversarial adaptation techniques are also applied. Grammenos et al. [41] proposed the federated PCA algorithm with differential privacy guarantee. The proposed FPCA method is permutation invariant and robust to straggler or fault clients.

4.6 Federated Reinforcement Learning

In deep reinforcement learning (DRL), the deep learning model gets rewards for its actions and learns which actions yield higher rewards. Zhuo et al. [43] introduced reinforcement learning to federated learning framework (FedRL), assuming that distributed agents do not share their observations. The proposed FedRL architecture has two local models: a simple neural network, such as multi-layer perceptron (MLP), and a Q-network that utilizes Q-learning to compute the re-

ward for a given state and action. The authors provided algorithms on how their model works with two clients and suggest that the approach can be extended to many clients using the same approach. In the proposed architecture, the clients update the local parameters of their respective MLPs first and then share the parameters to train these q-networks. Clients work out this parameter exchange in a peer-to-peer fashion. Federated reinforcement learning can improve federated aggregation to address the non-IID challenge, and it also has real-world applications, such as in the Internet of Things (IoT). A control framework called Favor [45] improves client selection with reinforcement learning to choose the best candidate for federated aggregation. The federated reinforcement distillation (FRD) framework [46], together with its improved variant MixFRD with mixup augmentation, utilizes policy distillation for distributed reinforcement learning. In the fusion stage of FRD, only proxy experience replay memory (ProxRM) with locally averaged policies are shared across agents, aiming to preserve privacy. Facing the tradeoff between the aggregator's pricing and the efficiency of edge computing, Zhan et al. [44] investigated the design of incentive mechanism with DRL to promote edge learning.

5 Applications

Current publications yield remarkable achievements in some real-world applications, while some focus more on using synthetic data and tasks to mimic the federation. Some applications have been studied in publications reviewed, such as recommendation [21] and image steganalysis[26]. There are also many industrial applications in the Internet of Things. Applications of cross-silo federated learning, including healthcare and financial applications, have practical significance. We recommend the survey by Xu et al. [50] for more introduction about research on federated healthcare informatics. Applications of cross-device federated learning require human-device interaction to provide labels as supervision signals for federated learning systems with the widely applied supervised learning methods. Mobile keyboard suggestion [1, 8] is a typical cross-device application in which the user's typing signal acts as supervision. More efforts should be paid to implement practical applications under the federated setting.

6 Challenges and Future Directions

In recent years, federated learning has seen drastic growth in terms of the amount of research and the breadth of topics. There is still a need for comparative studies, especially when assessing which learning paradigms should be used with FL.

Statistical Heterogeneity Diverse client patterns and hardware specifications bring heterogeneity to federated learning. We consider more about statistical heterogeneity as this paper focuses on federated learning algorithms. Federated

learning coupled with many different architectures and learning paradigms widen practical applications and play an essential role in modeling heterogeneous data. With various learning and optimization algorithms such as multitask learning, meta-learning, transfer learning, and alternate optimization techniques, recent advances achieve heterogeneity-aware model fusion. Nonetheless, there is still a long way to go with heterogeneity. More work focuses on overall performance, while no performance guarantee for individual devices.

Label Scarcity Current federated learning heavily relies on supervised learning. However, in most real-world applications, clients may not have sufficient labels or lack interaction between users to provide interactive labels. The label scarcity problem makes federated learning impractical in many scenarios. The idea of keeping private data on-device is fantastic; however, taking the label deficiency into consideration is critical in a realistic situation.

On-device Personalization Conventionally, personalization is achieved by additional fine-tuning before inference. Recently, more research focuses on personalization. On-device personalization [64] brings forward multiple possible scenarios where clients would additionally benefit from personalization. Mansour et al. [65] formulated their approaches for personalization, including user clustering, data interpolation, and model interpolation. Model-agnostic meta-learning aims to learn quick adaptations and also brings the potential to personalize to individual devices. The studies of effective formulation and metrics to evaluate the personalized performance is missed. The underlying essence of personalization and the connections between global model learning and personalized on-device training should be addressed.

Unsupervised Learning Current research on federated learning utilizes supervised or semi-supervised methods. Due to the label deficiency problem mentioned above in the real-world scenario, unsupervised representation learning can be the future direction in the federated setting and other learning problems.

7 Conclusion

This paper conducts a timely and focused survey about federated learning coupled with different learning algorithms. The flexibility of FL was showcased by presenting a wide range of relevant learning paradigms that can be employed within the FL framework. In particular, the compatibility was addressed from the standpoint of how learning algorithms fit the FL architecture and how they take into account two of the critical problems in federated learning: efficient learning and statistical heterogeneity.

References

- [1] H. B. McMahan, E. Moore, D. Ramage, S. Hampson, et al., Communication-efficient learning of deep networks from decentralized data, in: International Conference on Artificial Intelligence and Statistics, 2017, pp. 1273–1282.
- [2] E. Bagdasaryan, A. Veit, Y. Hua, D. Estrin, V. Shmatikov, How to backdoor federated learning, in: International Conference on Artificial Intelligence and Statistics, PMLR, 2020, pp. 2938–2948.
- [3] J. Konečný, H. B. McMahan, F. X. Yu, P. Richtárik, A. T. Suresh, D. Bacon, Federated learning: Strategies for improving communication efficiency, arXiv preprint arXiv:1610.05492 (2016).
- [4] S. Ji, W. Jiang, A. Walid, X. Li, Dynamic sampling and selective masking for communication-efficient federated learning, arXiv preprint arXiv:2003.09603 (2020).
- [5] C. He, C. Tan, H. Tang, S. Qiu, J. Liu, Central server free federated learning over single-sided trust social networks, arXiv preprint arXiv:1910.04956 (2019).
- [6] Y. Yeganeh, A. Farshad, N. Navab, S. Albarqouni, Inverse distance aggregation for federated learning with non-iid data, in: DCL Workshop at MICCAI, 2020, pp. 150–159.
- [7] Y. Chen, X. Sun, Y. Jin, Communication-efficient federated deep learning with layerwise asynchronous model update and temporally weighted aggregation, IEEE Transactions on Neural Networks and Learning Systems (2020).
- [8] S. Ji, S. Pan, G. Long, X. Li, J. Jiang, Z. Huang, Learning private neural language modeling with attentive aggregation, in: International Joint Conference on Neural Network, 2019.
- [9] T. Li, A. K. Sahu, A. Talwalkar, V. Smith, Federated learning: Challenges, methods, and future directions, IEEE Signal Processing Magazine 37 (3) (2020) 50–60.
- [10] C. Briggs, Z. Fan, P. Andras, Federated learning with hierarchical clustering of local updates to improve training on non-iid data, in: International Joint Conference on Neural Network, 2020.
- [11] M. Yurochkin, M. Agarwal, S. Ghosh, K. Greenewald, N. Hoang, Y. Khazaeni, Bayesian nonparametric federated learning of neural networks, in: International Conference on Machine Learning, 2019, pp. 7252–7261.
- [12] T. Li, M. Sanjabi, A. Beirami, V. Smith, Fair resource allocation in federated learning, in: International Conference on Learning Representations, 2020.
- [13] J. Jiang, S. Ji, G. Long, Decentralized knowledge acquisition for mobile internet applications, World Wide Web (2020).

- [14] X. Wu, Z. Liang, J. Wang, FedMed: A federated learning framework for language modeling, Sensors 20 (14) (2020) 4048.
- [15] Y. Huang, L. Chu, Z. Zhou, L. Wang, J. Liu, J. Pei, Y. Zhang, Personalized cross-silo federated learning on non-iid data, in: AAAI Conference on Artificial Intelligence, 2021.
- [16] F. X. Yu, A. S. Rawat, A. K. Menon, S. Kumar, Federated learning with only positive labels, in: International Conference on Machine Learning, 2020.
- [17] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, V. Smith, Federated optimization in heterogeneous networks, in: Conference on Machine Learning and Systems, 2020.
- [18] S. P. Karimireddy, M. Jaggi, S. Kale, M. Mohri, S. J. Reddi, S. U. Stich, A. T. Suresh, Mime: Mimicking centralized stochastic algorithms in federated learning, arXiv preprint arXiv:2008.03606 (2020).
- [19] A. Ghosh, J. Chung, D. Yin, K. Ramchandran, An efficient framework for clustered federated learning, in: Advances in Neural Information Processing Systems, 2020.
- [20] M. Xie, G. Long, T. Shen, T. Zhou, X. Wang, J. Jiang, Multi-center federated learning, arXiv preprint arXiv:2005.01026 (2020).
- [21] K. Muhammad, Q. Wang, D. O'Reilly-Morgan, E. Tragos, B. Smyth, N. Hurley, J. Geraci, A. Lawlor, FedFast: Going beyond average for faster training of federated recommender systems, in: SIGKDD, 2020, pp. 1234–1242.
- [22] H. Wang, M. Yurochkin, Y. Sun, D. Papailiopoulos, Y. Khazaeni, Federated learning with matched averaging, in: International Conference on Learning Representations, 2020.
- [23] M. Mohri, G. Sivek, A. T. Suresh, Agnostic federated learning, in: International Conference on Machine Learning, 2019.
- [24] Y. Liu, Y. Kang, C. Xing, T. Chen, Q. Yang, A secure federated transfer learning framework, IEEE Intelligent Systems 35 (2020) 70–82.
- [25] X. Peng, Z. Huang, Y. Zhu, K. Saenko, Federated adversarial domain adaptation, in: International Conference on Learning Representations, 2020.
- [26] H. Yang, H. He, W. Zhang, X. Cao, FedSteg: A Federated Transfer Learning Framework for Secure Image Steganalysis, IEEE Transactions on Network Science and Engineering (2020).
- [27] V. Smith, C.-K. Chiang, M. Sanjabi, A. S. Talwalkar, Federated multi-task learning, in: Advances in Neural Information Processing Systems, 2017, pp. 4427–4437.
- [28] S. Caldas, V. Smith, A. Talwalkar, Federated kernelized multi-task learning, in: Conference on Machine Learning and Systems, 2018.

- [29] F. Sattler, K.-R. Müller, W. Samek, Clustered federated learning: Modelagnostic distributed multitask optimization under privacy constraints, IEEE Transactions on Neural Networks and Learning Systems (2020).
- [30] X. Yao, T. Huang, R.-X. Zhang, R. Li, L. Sun, Federated learning with unbiased gradient aggregation and controllable meta updating, in: Advances in Neural Information Processing Systems Workshop, 2019.
- [31] A. Fallah, A. Mokhtari, A. Ozdaglar, Personalized federated learning with theoretical guarantees: A model-agnostic meta-learning approach, in: Advances in Neural Information Processing Systems, 2020.
- [32] D. Li, J. Wang, FedMD: Heterogenous federated learning via model distillation, in: Advances in Neural Information Processing Systems Workshop, 2019.
- [33] C. He, M. Annavaram, S. Avestimehr, Group knowledge transfer: Federated learning of large cnns at the edge, Advances in Neural Information Processing Systems (2020).
- [34] T. Lin, L. Kong, S. U. Stich, M. Jaggi, Ensemble distillation for robust model fusion in federated learning, in: Advances in Neural Information Processing Systems, 2020.
- [35] N. Papernot, M. Abadi, Ú. Erlingsson, I. Goodfellow, K. Talwar, Semisupervised knowledge transfer for deep learning from private training data, in: International Conference on Learning Representations, 2017.
- [36] W. Jeong, J. Yoon, E. Yang, S. J. Hwang, Federated semi-supervised learning with inter-client consistency, in: International Conference on Machine Learning Workshop, 2020.
- [37] C. Fan, P. Liu, Federated generative adversarial learning, arXiv preprint arXiv:2005.03793 (2020).
- [38] M. Rasouli, T. Sun, R. Rajagopal, FedGAN: Federated generative adversarial networks for distributed data, arXiv preprint arXiv:2006.07228 (2020).
- [39] S. Augenstein, H. B. McMahan, D. Ramage, S. Ramaswamy, P. Kairouz, M. Chen, R. Mathews, B. A. y Arcas, Generative models for effective ml on private, decentralized datasets, in: International Conference on Learning Representations, 2020.
- [40] B. v. Bram, A. Saeed, T. Ozcelebi, Towards federated unsupervised representation learning, in: ACM EdgeSys, 2020, pp. 31–36.
- [41] A. Grammenos, R. Mendoza Smith, J. Crowcroft, C. Mascolo, Federated principal component analysis, in: Advances in Neural Information Processing Systems, 2020.

- [42] F. Zhang, K. Kuang, Z. You, T. Shen, J. Xiao, Y. Zhang, C. Wu, Y. Zhuang, X. Li, Federated unsupervised representation learning, arXiv preprint arXiv:2010.08982 (2020).
- [43] H. H. Zhuo, W. Feng, Q. Xu, Q. Yang, Y. Lin, Federated deep reinforcement learning, arXiv preprint arXiv:1901.08277 (2019).
- [44] Y. Zhan, J. Zhang, An incentive mechanism design for efficient edge learning by deep reinforcement learning approach, in: IEEE International Conference on Computer Communications, IEEE, 2020, pp. 2489–2498.
- [45] H. Wang, Z. Kaplan, D. Niu, B. Li, Optimizing Federated Learning on Non-IID Data with Reinforcement Learning, in: IEEE International Conference on Computer Communications, IEEE, 2020, pp. 1698–1707.
- [46] H. Cha, J. Park, H. Kim, M. Bennis, S.-L. Kim, Proxy experience replay: Federated distillation for distributed reinforcement learning, IEEE Intelligent Systems (2020).
- [47] Q. Yang, Y. Liu, T. Chen, Y. Tong, Federated machine learning: Concept and applications, ACM Transactions on Intelligent Systems and Technology 10 (2) (2019) 12.
- [48] Q. Li, Z. Wen, B. He, Federated learning systems: Vision, hype and reality for data privacy and protection, arXiv preprint arXiv:1907.09693 (2019).
- [49] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings, et al., Advances and open problems in federated learning, arXiv preprint arXiv:1912.04977 (2019).
- [50] J. Xu, F. Wang, Federated learning for healthcare informatics, arXiv preprint arXiv:1911.06270 (2019).
- [51] L. Lyu, H. Yu, Q. Yang, Threats to federated learning: A survey, arXiv preprint arXiv:2003.02133 (2020).
- [52] W. Y. B. Lim, N. C. Luong, D. T. Hoang, Y. Jiao, Y.-C. Liang, Q. Yang, D. Niyato, C. Miao, Federated learning in mobile edge networks: A comprehensive survey, IEEE Communications Surveys & Tutorials (2020).
- [53] S. Niknam, H. S. Dhillon, J. H. Reed, Federated learning for wireless communications: Motivation, opportunities, and challenges, IEEE Communications Magazine 58 (6) (2020) 46–51.
- [54] Y. Jin, X. Wei, Y. Liu, Q. Yang, A survey towards federated semi-supervised learning, arXiv preprint arXiv:2002.11545 (2020).
- [55] S. K. Lo, Q. Lu, C. Wang, H. Paik, L. Zhu, A systematic literature review on federated machine learning: From a software engineering perspective, arXiv preprint arXiv:2007.11354 (2020).

- [56] X. Li, K. Huang, W. Yang, S. Wang, Z. Zhang, On the convergence of fedavg on non-iid data, in: International Conference on Learning Representations, 2020.
- [57] S. P. Karimireddy, S. Kale, M. Mohri, S. J. Reddi, S. U. Stich, A. T. Suresh, Scaffold: Stochastic controlled averaging for federated learning, in: International Conference on Machine Learning, 2020, pp. 5132–5143.
- [58] S. Reddi, Z. Charles, M. Zaheer, Z. Garrett, K. Rush, J. Konečný, S. Kumar, H. B. McMahan, Adaptive federated optimization, in: International Conference on Learning Representations, 2021.
- [59] S. P. Singh, M. Jaggi, Model fusion via optimal transport, Advances in Neural Information Processing Systems 33 (2020).
- [60] M. Alawad, H.-J. Yoon, S. Gao, B. Mumphrey, X.-C. Wu, E. B. Durbin, J. C. Jeong, I. Hands, D. Rust, L. Coyle, et al., Privacy-preserving deep learning nlp models for cancer registries, IEEE Transactions on Emerging Topics in Computing (2020).
- [61] C. Finn, P. Abbeel, S. Levine, Model-agnostic meta-learning for fast adaptation of deep networks, in: International Conference on Machine Learning, 2017, pp. 1126–1135.
- [62] S. Ji, G. Long, S. Pan, T. Zhu, J. Jiang, S. Wang, X. Li, Knowledge transferring via model aggregation for online social care, arXiv preprint arXiv:1905.07665 (2019).
- [63] Y. Jiang, J. Konečný, K. Rush, S. Kannan, Improving federated learning personalization via model agnostic meta learning, in: Advances in Neural Information Processing Systems Workshop, 2019.
- [64] K. Wang, R. Mathews, C. Kiddon, H. Eichner, F. Beaufays, D. Ramage, Federated evaluation of on-device personalization, arXiv preprint arXiv:1910.10252 (2019).
- [65] Y. Mansour, M. Mohri, J. Ro, A. T. Suresh, Three approaches for personalization with applications to federated learning, arXiv preprint arXiv:2002.10619 (2020).