

Semi-Supervised Federated Learning with non-IID Data: Algorithm and System Design

Zhe Zhang^{1,2}, Shiyao Ma³, Jiangtian Nie⁴, Yi Wu^{1,2,*}, Qiang Yan⁵, Xiaoke Xu³, Dusit Niyato⁴, *Fellow, IEEE*

¹*School of Data Science and Technology, Heilongjiang University, Harbin 150080, China*

²*Cryptology and Network Security Laboratory of Heilongjiang University, Harbin 150080, China*

³*College of Information and Communication Engineering, Dalian Minzu University*

⁴*School of Computer Science and Engineering, Nanyang Technological University*

⁵*WeBank Co. Ltd., China*

*1995050@hlju.edu.cn

Abstract—Federated Learning (FL) allows edge devices (or *clients*) to keep data locally while simultaneously training a shared high-quality global model. However, current research is generally based on an assumption that the training data of local clients have ground-truth. Furthermore, FL faces the challenge of statistical heterogeneity, i.e., the distribution of the client's local training data is non-independent identically distributed (non-IID). In this paper, we present a robust semi-supervised FL system design, where the system aims to solve the problem of data availability and non-IID in FL. In particular, this paper focuses on studying the labels-at-server scenario where there is only a limited amount of labeled data on the server and only unlabeled data on the clients. In our system design, we propose a novel method to tackle the problems, which we refer to as Federated Mixing (FedMix). FedMix improves the naive combination of FL and semi-supervised learning methods and designs parameter decomposition strategies for disjointed learning of labeled, unlabeled data, and global models. To alleviate the non-IID problem, we propose a novel aggregation rule based on the frequency of the client's participation in training, namely the FedFreq aggregation algorithm, which can adjust the weight of the corresponding local model according to this frequency. Extensive evaluations conducted on CIFAR-10 dataset show that the performance of our proposed method is significantly better than those of the current baseline. It is worth noting that our system is robust to different non-IID levels of client data.

Index Terms—Federated Learning, Semi-supervised Learning, non-IID, Aggregation Algorithm

I. INTRODUCTION

Federated Learning (FL) [1], [2] is a distributed machine learning paradigm that allows multiple edge devices (or *clients*) to cooperatively train a shared global model [3]–[5]. The most obvious difference between FL and traditional distributed machine learning is that clients can privately access local training data without sharing data with cloud centers [6]–[8]. However, the current mainstream work is based on an unrealistic assumption: *the training data of local clients have ground-truth* [9]. In our daily lives, it is not common for each client to have rich labeled data. For example, in the early stage of COVID-19 epidemic, community hospitals without enough labeled data may not be able to train a high-precision pathophoresis prediction model. On the other hand, in most cases, putting together a properly labeled dataset for a given FL task is a time-consuming, expensive, and complicated endeavor

[9]. Therefore, it is challenging to train a high-quality global model in the real scenario of a lack of labeled data.

In the face of the above challenges, recent works [9]–[14] study how to design a semi-supervised FL (SSFL) system that can efficiently integrate semi-supervised learning into FL techniques. For example, Jeong *et al.* in [12] proposed an SSFL system with a new inter-client consistency loss to achieve this goal. In fact, consistency regularization technique is widely used in semi-supervised learning which keeps the same output that the same data inject two different noises [15], [16]. Furthermore, pseudo-label methods are important for SSFL, where they mainly utilize pseudo-labels whose predicted value is higher than the confidence threshold to achieve high-precision SSFL [9], [10], [17]. However, it is worth noting that there still remain gaps when deploying SSFL in practice.

First, traditional SSFL methods generally introduce semi-supervised techniques (such as consistency loss and pseudo-label) directly into the FL system, which ignores the implicit contribution between iterative updates of the global model. Previous work only focused on how to set pseudo-labels or how to decompose the parameters of labeled and unlabeled data for disjoint learning. In this way, the learned global model will be biased towards labeled data (supervised model) or unlabeled data (unsupervised model) instead of the global model [12]. This implies that we need to observe the implicit effects between iterations of the global model at a fine-grained level.

Second, the non-independent identically distributed (non-IID) of data between clients has always been a key and challenging issue in FL. The reason is that there are too many differences in data distribution, features, and the number of labels between clients, which is not conducive to the convergence of the global model. Currently, many efforts have effectively alleviated the non-IID problem, such as FedBN [18] utilized local batch normalization to alleviate the feature shift before average aggregating local models. However, methods such as these add additional computational and communication overhead to the server or client.

In this paper, to address the first issue, we propose the Federated Mixing (FedMix) algorithm, which performs param-

eter decomposition of disjointed learning for supervised model (learned on labeled data), unsupervised model (learned on unlabeled data), and global model. In particular, this algorithm analyzes the implicit effects between iterations of the global model in a fine-grained manner. To address the second issue, we propose a novel aggregation rule called Federated Frequency (FedFreq), which dynamically adjusts the weight of the corresponding local model by recording the training frequency of the client to alleviate the non-IID problem. Furthermore, we introduce the Dirchlet distribution function to simulate the different non-IID level scenario in our experiment. The main contributions of this paper are as follows:

- We present a robust Semi-supervised Federated Learning system design, where the system aims to solve the problem of data availability and non-IID in FL. In our system design, we propose FedMix algorithm to improve the naive combination of FL and semi-supervised learning methods.
- We propose a novel aggregation rule called FedFreq, which dynamically adjusts the weight of the corresponding local model by recording the training frequency of the client to alleviate the non-IID problem.
- We conduct extensive evaluations on CIFAR-10 dataset, which show that the performance of our designed system is 3% higher than the baseline.

II. RELATED WORK

A. Semi-supervised Federated Learning

Semi-supervised federated learning attempts to use semi-supervised learning techniques [19]–[23] to further improve the performance of the FL model in scenarios where there is unlabeled data on the client side [11]. For example, Long *et al.* in [10] proposed a semi-supervised federated learning (SSFL) system, *FedSemi*, which unifies the consistency-based semi-supervised learning model [24], dual model [15], and average teacher model [25] to achieve SSFL. The DS-FL system [26] was proposed to solve the communication overhead problem in SSFL. Reference [27] proposes a method to study the distribution of non-IID data, which introduces a probability distance metric to evaluate the difference in client data distribution in SSFL. Different from the literature [9], [10], [26], in this paper, we focus on labels-at-server scenario and also solve the problem of data availability and data heterogeneity in SSFL.

B. Robust Federated Learning

If the local data set distribution of each client is inconsistent (i.e., non-IID problem) [18], [28]–[33], the local objective loss function of the client will be inconsistent with the global objective [34]. In particular, when the model of the local client is updated larger, such a difference will be more obvious. Therefore, we need to design some robust FL systems to solve the above problems. Some studies try to design a robust federated learning algorithm to solve the non-IID problem. For example, *FedProx* [6] limits the distance between the local model and the global model by introducing an additional \mathcal{L}_2

regularization term in the local target function to limiting the size of the local model update. However, this method has a disadvantage that each client needs to individually adjust the local regularization term to obtain good model performance. *FedNova* [35] improved FedAvg in the aggregation phase, which normalized and scaled the model update according to the local training batch of the client. Although previous studies have alleviated the problem of non-IID to some extent, they only evaluated the data distribution at specific non-IID levels and lacked extensive experimental verification for different non-IID scenarios. Therefore, we propose a more comprehensive data distribution and data partition strategy, i.e., we introduce the Dirichlet distribution function to simulate different non-IID levels of client data.

III. PRELIMINARIES

A. Federated Learning

Federated learning solves the problem of data island on the premise of privacy protection. In particular, the FL is a distributed machine learning framework, which requires clients to hold data locally, where these clients coordinate to train a shared global model ω^* . In FL, there is a server \mathcal{S} and k clients, each of which holds an IID or non-IID datasets \mathcal{D}_k . Specifically, for a training sample x on the client side, let $\ell(\omega; x)$ be the loss function at the client, where $\omega \in \mathbb{R}^d$ denotes the model's trainable parameters. Therefore, we let $\mathcal{L}(\omega) = \mathbb{E}_{x \sim \mathcal{D}}[\ell(\omega; x)]$ be the loss function at the server. Thus, FL needs to optimize the following objective function at the server:

$$\min_{\omega} \mathcal{L}(\omega), \text{ where } \mathcal{L}(\omega) := \sum_{k=1}^K p_k \mathcal{L}_k(\omega), \quad (1)$$

where $p_k \geq 0, \sum_k p_k = 1$ indicates the relative influence of k -th client on the global model. In FL, to minimize the above objective function, the server and clients execute the following steps:

- **Step 1, Initialization:** The server sends the initialized global model ω_0 to the selected clients.
- **Step 2, Local training:** The client uses the local optimizer (e.g., SGD, Adam) on the local dataset \mathcal{D}_k to train the received initialization model. Then, each client uploads the local model ω_t^k to the server.
- **Step 3, Aggregation:** The server collects and uses a certain algorithm (e.g., FedAvg [1]) to aggregate the model updates uploaded by these clients to obtain a new global model, i.e., $\omega_{t+1} = \omega_t + \sum_{k=1}^K \frac{D_k}{D} \omega_t^k$. Then, the server sends the updated global model ω_{t+1} to all selected clients.

Note that FL repeats the above steps until the global model converges.

B. Semi-supervised Learning

In the real world, (e.g., financial and medical fields), unlabeled data is easy to gain, while labeled data is often difficult to obtain. Meanwhile, annotating data requires a lot of manpower

and material resources. To this end, the researchers proposed a machine learning paradigm namely Semi-supervised Learning [36], [37] to learn a high-precision model on a mixed dataset (part of the data is labeled, and some of the data is unlabeled). Thus, in recent years, semi-supervised learning has become a hot research direction in the field of deep learning. In this section, we introduce a basic assumption and two methods of semi-supervised learning.

Assumption 1: In machine learning, there is a basic assumption that if the feature of two unlabeled samples u_1 and u_2 are similar, the corresponding model prediction results y_1 and y_2 are the same [38], i.e., $f(u_1) = f(u_2)$, where $f(\cdot)$ is the prediction function.

According to the above assumption, we adopt two common semi-supervised learning methods as follows:

Consistency Regularization: The main idea of this method is that the model prediction results should be the same whether noise is added or not on an unlabeled training sample [15], [36]. We generally use data augmentation (such as image flipping and shifting) methods to add noise to increase the diversity of the dataset. Specifically, for an unlabeled sample u_i in unlabeled dataset $\mathbf{u} = \{u_i\}_{i=1}^m$ and its perturbation form \hat{u}_i , our goal is to minimize the distance $d(f_\theta(u_i), f_\theta(\hat{u}_i))$, where $f_\theta(u_i)$ is the output of sample u_i on model θ . The common distance measurement method is Kullback-Leiber (KL) divergence. Thus, we can calculate the consistency loss as follows:

$$d(f_\theta; \mathbf{u})_{KL} = \frac{1}{m} \sum_{i=1}^m f_\theta(u_i) \log \frac{f_\theta(u_i)}{f_\theta(\hat{u}_i)}, \quad (2)$$

where m is the total number of unlabeled samples and $f_\theta(u_i)$ indicates the model output of unlabeled sample u_i .

Pseudo-label: The pseudo-label method [24] is to utilize some labeled samples to train a model that to set the pseudo labels for unlabeled samples. Previous work generally used sharpening [39] and argmax [24] methods to set pseudo-label, where the former make the distribution of model output extreme of unlabeled samples and the latter change the model output to one-hot of unlabeled samples.

IV. PROBLEM DEFINITION

In the SSFL system, there are two essential scenarios of SSFL based on the location of the labeled data. The first scenario considers a conventional case where clients have both labeled and unlabeled data (*labels-at-client*), and the second scenario considers a more challenging case, where the labeled data is only available at the server (*labels-at-server*). In particular, in this paper, we consider only the labels-at-server scenario. Next, we give the definition of the problem studied in this paper as follows:

Labels-at-server Scenario: In SSFL, we assume that there is a server \mathcal{S} and K clients, where the server holds a labeled dataset $\mathcal{D}_s = \{(x_i, y_i)\}_{i=1}^n$ and each client holds a local unlabeled dataset $\mathcal{D}_k = \{u_i\}_{i=1}^m$. Thus, in this scenario, for

unlabeled training sample u_i , let \mathcal{L}_u^k be the loss function at the client side:

$$\mathcal{L}_u^k = \frac{1}{m} \sum_{u_i \in \mathcal{D}_k} CE(\hat{y}_i, f_{\theta_k}(u_i)) + \frac{1}{m} \sum_{u_i \in \mathcal{D}_k} \|f_{\theta_k}(u_i) - f_{\theta_k}(\pi(u_i))\|^2, \quad (3)$$

where m is the number of unlabeled samples, $\pi(\cdot)$ is the data augmentation function (e.g., flip and shift of the unlabeled samples), \hat{y}_i is the pseudo label of unlabeled sample u_i , and $f_{\theta_k}(u_i)$ indicates the output of unlabeled sample u_i on model θ_k of the k -th client. For labeled sample x_i , let \mathcal{L}_s be the loss function at the server side:

$$\mathcal{L}_s = \frac{1}{n} \sum_{x_i, y_i \in \mathcal{D}_s} CE(y_i, f_\theta(x_i)), \quad (4)$$

where n is the number of labeled samples, and $f_\theta(x_i)$ indicates the output of labeled sample x_i on model θ .

Therefore, the objective function of this scenario in SSFL system is to minimize the following loss function:

$$\min \mathcal{L}, \text{ where } \mathcal{L} \doteq \sum_{k=1}^K \mathcal{L}_u^k + \mathcal{L}_s. \quad (5)$$

Note that the whole learning process is similar to the traditional FL system, except that the server not only aggregates the client model parameters but also trains the model with labeled data.

V. ALGORITHM AND SYSTEM DESIGN

A. Semi-supervised Federated Learning System Design

In our system setting, the server \mathcal{S} holds a labeled dataset $\mathcal{D}_s = \{(x_i, y_i)\}_{i=1}^n$, where n indicates the number of labeled samples. For K clients, we assume that k -th client holds a local unlabeled dataset $\mathcal{D}_k = \{u_i\}_{i=1}^m$, where m denotes the number of unlabeled samples in the local client. Similar to the traditional FL system, the server and clients in the SSFL are cooperative to train a high-performance global model ω^* . The goal of previous work is to optimize the objective function mentioned above, i.e., Equation (5). However, they ignore the implicit contribution between iterations of the global model, which results in the learned global model not being optimal. Inspired by the above facts, we propose an SSFL algorithm called FedMix that focuses on the implicit contributions between iterations of the global model in a fine-grained manner. We define the supervised model trained on the labeled dataset as σ , the unsupervised model trained on the unlabeled dataset as ψ , and the aggregated global model as ω . Specifically, we design a strategy that assigns three weights α, β , and γ to the unsupervised model ψ , supervised model σ , and the previous round of global model, respectively. The designed algorithm can capture the implicit relationship between each iteration of the global model in a fine-grained manner. Thus, the steps of our proposed FedMix algorithm are as follows:

- **Step 1, Initialization:** The server randomly selects a certain proportion of F ($0 < F < 1$) clients from all

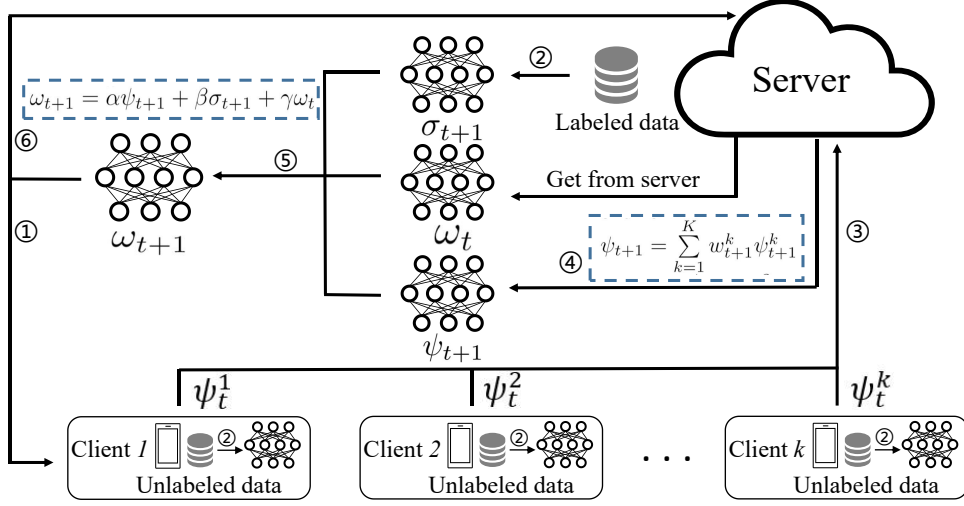


Fig. 1. Overview of semi-supervised federated learning system.

local clients to send the initialized global model ω_0 . Note that the global model ω_0 also remains on the server-side.

- **Step 2, Server Training:** Unlike FL, in our SSFL system, the server not only aggregates the model uploaded by the clients, but also trains the supervised model σ (i.e., $\sigma_t \leftarrow \omega_t$) on the labeled dataset \mathcal{D}_s . Thus, the server uses the local optimizer on the labeled dataset \mathcal{D}_s to train the supervised model σ . The minimization of the objective function is defined as follows:

$$\min_{\sigma \in \mathbb{R}^d} \mathcal{L}_s(\sigma_t), \text{ where } \mathcal{L}_s(\sigma_t) \doteq \lambda_s CE(y, f_{\sigma_t}(x)), \quad (6)$$

where λ_s is the hyperparameter, x and y are from labeled dataset \mathcal{D}_s , and $f_{\sigma_t}(x)$ means the output of labeled samples on supervised model σ at t -th training round.

- **Step 3, Local Training:** The k -th client utilizes the local unlabeled data to train the received global model ω_t (i.e., $\psi_t^k \leftarrow \omega_t$) and then obtains the unsupervised model ψ_{t+1}^k . Thus, we define the following objective function:

$$\min_{\psi \in \mathbb{R}^d} \mathcal{L}_u(\psi_t^k), \mathcal{L}_u(\psi_t^k) \doteq \lambda_2 \|f_{\psi_t^k}(\pi_1(u)) - f_{\psi_t^k}(\pi_2(u))\|^2 + \lambda_1 CE(\hat{y}, f_{\psi_t^k}(u)) + \lambda_{L1} \|\sigma_t - \psi_t^k\|^2, \quad (7)$$

where λ_1, λ_2 , and λ_{L1} are hyperparameters to control the ratio between the loss terms, ψ_t^k is the unsupervised model of the k -th client at t -th training round, u is from unlabeled dataset \mathcal{D}_k , $\pi(\cdot)$ is the form of perturbation, i.e., π_1 is the shift augmentation, π_2 is the flip augmentation, $\|\sigma_t - \psi_t^k\|^2$ is a penalty term that aims to let the k -th client unsupervised model ψ_t^k learn the knowledge of the supervised model σ_t (note that σ_t is inferred from Equation (9)), and \hat{y} is pseudo label obtained by using our proposed argmax method. The argmax method is defined as follows:

$$\hat{y} = \mathbf{1}(\text{Max}(\sum_{i=1}^A f_{\psi_t^k}(\pi_i(u)))), \quad (8)$$

where $\text{Max}(\cdot)$ is a function that can output the maximum probability that unlabeled data belongs to a certain class, $\mathbf{1}(\cdot)$ is the one-hot function that can change the numerical value to 1, A represents the number of unlabeled data after data augmentation, and u is from unlabeled dataset \mathcal{D}_k . Specifically, we discard low-confident predictions below confidence threshold $\tau = 0.80$ when generating pseudo-labels.

- **Step 4, Aggregation:** The server uses the proposed FedFreq (see Section IV-B) aggregation algorithm to aggregate the unsupervised models uploaded by the clients to obtain the global unsupervised model, i.e., $\psi_{t+1} = \sum_{k=1}^K w_{t+1}^k \psi_{t+1}^k$, where ψ_{t+1}^k is the unsupervised model of the k -th client at $t+1$ -th training round and w_{t+1}^k is the weight of the k -th client. The server then aggregates the global unsupervised model ψ_{t+1} , the supervised model σ_{t+1} , and the global model ω_t from the previous round t to obtain a new global model ω_{t+1} :

$$\omega_{t+1} = \alpha \psi_{t+1} + \beta \sigma_{t+1} + \gamma \omega_t, \quad (9)$$

where α, β , and γ are the corresponding weights of the three models and $(\alpha, \beta, \gamma) \in \{\alpha + \beta + \gamma = 1 \wedge \alpha, \beta, \gamma \geq 0\}$.

Repeat all the above steps until the global model converges. The proposed FedMix algorithm is shown in Algorithm 1.

B. FedFreq Aggregation Algorithm

In this section, we present the designed FedFreq aggregation algorithm, which can dynamically adjust the weight of the corresponding local model according to the training frequency of the client to alleviate the non-IID problem. We observe that the parameter distribution of the global model will be biased towards clients that often participate in federated training, which is obviously not friendly to the robustness of the global model. Therefore, our insight is to reduce the influence of clients with high training frequency on the global model to improve the

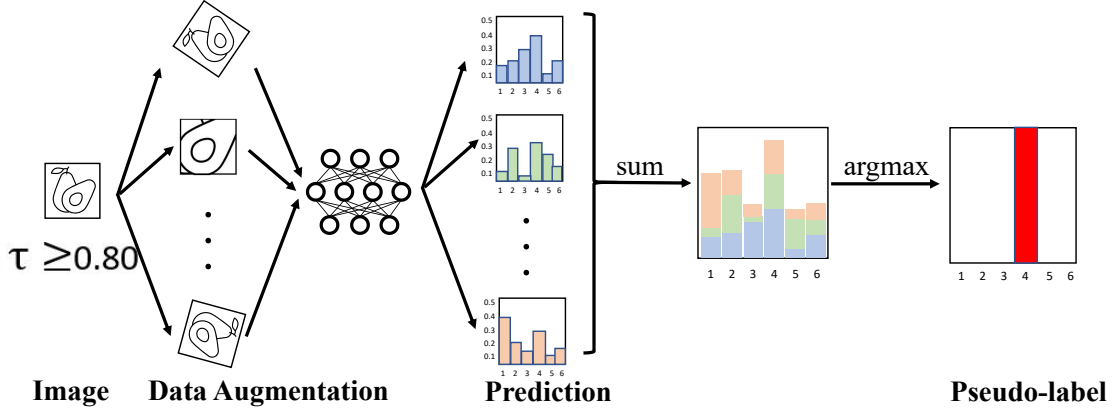


Fig. 2. Overview of the proposed argmax method.

Algorithm 1 FedMix algorithm on labels-at-server scenario.

Input: The client set \mathcal{K} , B_{server} is the mini-batch size at the server side, E_{server} is the number of epochs at the server side, B_{client} is the local mini-batch size at the client side, E_{client} is the number of local epochs at the client side, and η is the learning rate.

Output: The optimal global model ω^* .

```

1: Server executes:
2: Initialize global model  $\omega_0$ 
3: for each round  $t = 0, 1, 2, \dots$  do
4:    $\sigma_t \leftarrow \omega_t$ 
5:   for the server epoch  $e$  from 1 to  $E_{server}$  do
6:     for mini-batch  $b \in B_{server}$  do
7:        $\sigma_{t+1} = \sigma_t - \eta \nabla \mathcal{L}_s(\sigma_t, \mathcal{D}_s, b)$ 
8:     end for
9:   end for
10:   $m \leftarrow \max(F \cdot K, 1)$ 
11:   $S_t \leftarrow$  randomly select  $m$  clients from the client set  $\mathcal{K}$ 
12:  for each client  $k \in S_t$  in parallel do
13:     $\psi_t^k \leftarrow \omega_t$ 
14:     $\psi_{t+1}^k \leftarrow \text{ClientUpdate}(k, \psi_t^k)$ 
15:  end for
16:   $\psi_{t+1} = \sum_{k=1}^K w_{t+1}^k \psi_{t+1}^k$  // Refer to FedFreq algorithm
17:   $\omega_{t+1} = \alpha \psi_{t+1} + \beta \sigma_{t+1} + \gamma \omega_t$ 
18: end for
19: ClientUpdate( $k, \psi_t^k$ ): // Run on client  $k$ 
20: for each local epoch  $e$  from 1 to  $E_{client}$  do
21:   for minibatch  $b \in B_{client}$  do
22:     $\psi_{t+1}^k = \psi_t^k - \eta \nabla \mathcal{L}_u(k, \psi_t^k, \mathcal{D}_k, b)$ 
23:   end for
24: end for
25: return  $\omega^*$  to server.

```

where F is the sample proportion of the server, K is the total number of clients, $p_{t+1}^k = \frac{q_{t+1}^k}{\sum_{k \in S_{t+1}} q_{t+1}^k}$, q_{t+1}^k is the number of times that the k -th client has been trained up to the $t+1$ -th round, and S_{t+1} denotes the set of clients selected by the server in round $t+1$. Then, for the client, the FedFreq aggregation rule is expressed as follows: $\psi_{t+1} = \sum_{k=1}^K w_{t+1}^k \psi_{t+1}^k$.

C. Dirichlet Data Distribution Function

To better evaluate the robustness of the designed system to non-IID data, in this paper, we introduce the Dirichlet distribution function [40], [41], which is a popular non-IID function, to adjust the non-IID level of the local client data. Specifically, we generate data distributions of different non-IID levels by adjusting the parameter (i.e., μ) of the Dirichlet distribution function. We assume that the local dataset \mathcal{D}_k of k -th client has c classes, and thus, the definition of Dirichlet distribution function is as follows:

$$P_k(\varphi_1, \dots, \varphi_c) = \frac{\Gamma(\sum_i \mu_i)}{\prod_i \Gamma(\mu_i)} \prod_{i=1}^c \Gamma(\mu_i) \varphi_i^{\mu_i-1}, \quad (11)$$

$$p_k(\varphi_c) = \frac{P_k(\varphi_c)}{\sum_{i=1}^c P_k(\varphi_i)}, \quad (12)$$

where Θ is a set of c samples randomly selected from the Dirichlet function, i.e., $\Theta = \{\varphi_1, \dots, \varphi_c\}$ and $\Theta \sim \text{Dir}(\mu_1, \dots, \mu_c)$, μ, μ_1, \dots, μ_c are the parameters of the Dirichlet distribution function (where $\mu = \mu_1 = \mu_2 = \dots = \mu_c$), and $p_k(\varphi_c)$ denotes the proportion of the c -th class data in all data of the client. In particular, the smaller the μ , the higher the non-IID level of the data distribution of each client; otherwise, the data distribution of the client tends to the IID setting. Therefore, we adjust the parameters of the Dirichlet distribution function to simulate the different non-IID levels of the client's local dataset. For example, as shown in Fig. 3, we demonstrate the data distribution when $\mu = \{0.1, 1, 10\}$.

robustness of the model. Thus, the formal expression of the FedFreq aggregation algorithm is as follows:

$$w_{t+1}^k = \frac{1 - p_{t+1}^k}{1 - p_{t+1}^1 + \dots + 1 - p_{t+1}^k} = \frac{1 - p_{t+1}^k}{FK - 1}, \quad (10)$$

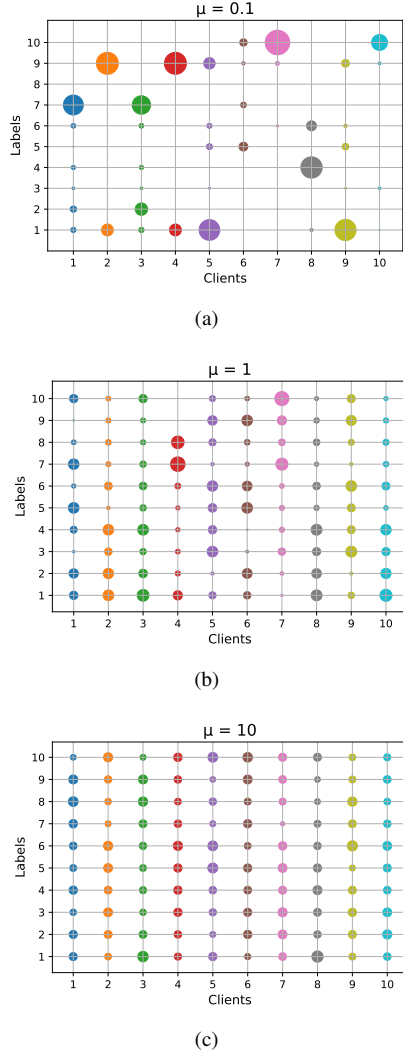


Fig. 3. Overview of the data distribution when $\mu = \{0.1, 1, 10\}$.

VI. EXPERIMENT

In the Labels-at-Server scenario, we compared our method FedMix and the baseline FedMatch [12] on the CIFAR-10 dataset. Furthermore, we simulate the federated learning setup (one server and K clients) on a commodity machine with Intel(R) Core(TM) i9-9900K CPU @ 3.60GHz and NVIDIA GeForce RTX 2080Ti GPU.

A. Experiment Setup

Dataset 1) CIFAR-10 dataset under IID setting: We use the CIFAR-10 dataset including 56,000 training samples and 2,000 test samples as the validation dataset in our experiment. The training set includes 55,000 unlabeled samples and 1,000 labeled samples, where the former is used to train the unsupervised model at the local and the latter is used to train the supervised model at the server. The unlabeled samples are equally distributed to 100 clients at a ratio of 1 : 100, of which there are 550 samples for each client (i.e., 55 for each class, and 10 classes in total). Similarly, the labeled samples have a total

of 1,000 and contain 10 classes on the server, of which there are 100 samples in each class. Meanwhile, we set a participation rate $F = 0.05$ of clients, i.e., 5 clients are randomly selected for training in each round.

Dataset 2) CIFAR-10 dataset under non-IID setting: Our setting is similar to the above IID setting, except that the Dirichlet distribution function is introduced to adjust the non-IID level of the local client data. Specifically, we generate data distributions of different non-IID levels by adjusting the parameter (i.e., μ) of the Dirichlet distribution function. Meanwhile, we simulate quantity imbalance and class imbalance the local client data. In particular, we make each client hold a different number of training samples. For example, some customers have 580 samples, while some customers have less than 50 samples. Second, we make the client hold a different sample size for each category in the data. For example, some clients have ten types of data, and some clients have less than two types of data.

Baseline and training details: Our baseline is FedMatch [12] naively using unsupervised model and supervised model parameter decomposition strategy, i.e., $\omega = \psi + \sigma$. In the training process, both our model and baseline use Stochastic Gradient Descent (SGD) to optimize the ResNet-9 neural network with initial learning rate $\eta = 1e - 3$. We set training round $t = 150$, the number of labeled samples on sever is $N_s = 1000$, local client training epoch $E_{client} = 1$ and mini-batch size $B_{client} = 64$, the server training epoch $E_{server} = 1$ and mini-batch size $B_{server} = 64$. Second, we set the data augmentation number in the argmax method $A = 5$.

B. Experiment Results

As shown in Fig. 4, under IID and non-IID settings, our method FedMix is better than baseline under each different aggregation method settings. For example, under a non-IID setting, the convergence accuracy of our method is 47.5% about 3% higher than that of the baseline. In particular, the accuracy of our method increases faster and more stable in the early stage of model training. The reason is that: (1) The FedMix focuses on the implicit contributions between iterations of the global model in a fine-grained manner, while the FedMatch only naively uses model parameter decomposition. (2) Frequency-based aggregation method FedFreq is more suitable for non-IID settings. Notably, the FedFreq only requires the server to give appropriate weights in the aggregation process according to the training frequency of each client, which does not bring additional computational overhead to the server and local clients.

Fig. 5 shows the performance comparison of the proposed method under different hyperparameter settings. To be specific, the three hyperparameters are the weights of the global unsupervised model, the supervised model, and the previous round of global model. From Fig. 5, we can find that under the non-IID setting, as α decreases, the accuracy curve of the proposed method becomes unstable. The reason for this phenomenon is that: with the decrease of global unsupervised model weight,

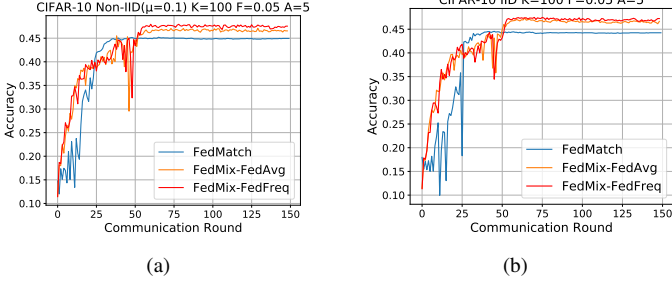


Fig. 4. Test accuracy curves on IID and non-IID setting.

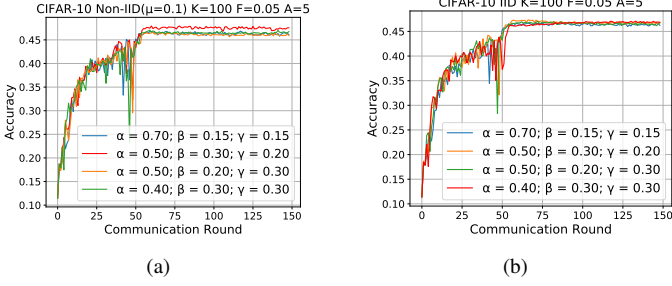


Fig. 5. Test accuracy curves under different hyperparameter settings.

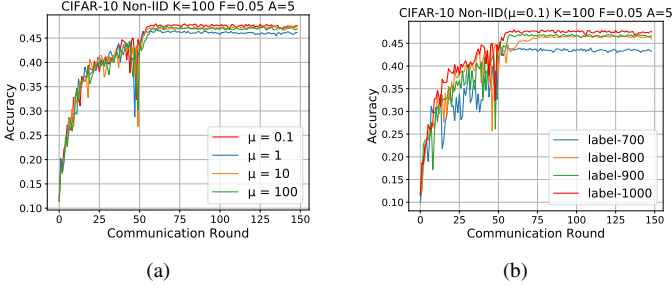


Fig. 6. (a) Performance comparison on different non-IID levels; (b) Performance comparison on different numbers of labeled samples.

FedFreq is losing its effect. In particular, when $\alpha = 0.5$, $\beta = 0.3$, and $\gamma = 0.2$, our global model is the aggregation of the optimal weights of the three models. In addition, we find that the proposed method is easy to achieve better performance when these three parameters are relatively close under the IID setting, as shown in Fig. 5(b).

Fig. 6(a) shows the performance comparison of the proposed method on different non-IID levels of client data. In this experiment, we let $\mu = 0.1$ denote the highest non-IID level of the client data. In this case, as the value of μ increases, the local client data distribution is closer to the IID setting. It can be seen from Fig. 6 that for different non-IID levels, our method all can achieve stable accuracy. Meanwhile, the model convergence accuracy under $\mu = \{0.1, 1, 10, 100\}$ settings does not differ by more than 1%. Therefore, our method is not sensitive to the different levels of client data distribution, i.e., it is robust to different types of data distribution settings.

Fig. 6(b) shows the performance comparison of the proposed

method in the case of different numbers of labeled samples at the server. Obviously, the converged accuracy of our method is 47% with 800 labeled samples, which is 2% higher than FedMatch. However, when the number of labeled samples is reduced to 700, the accuracy of our model decreases greatly. Therefore, we regard $N_s = 800$ as the best setting for our method.

C. Analysis

In this section, we further analyze the advantages of FedMix compared to FedMatch in labels-at-server scenario.

1) The performance of FedMix training model under the CIFAR-10 dataset is better than FedMatch. This is due to FedMatch simply uses the strategy of parameter decomposition of unsupervised model and supervised model in the training process, i.e., $\omega_t = \psi_t + \sigma_t$. In this way, the learned global model will be biased towards unlabeled data (unsupervised model) or labeled data (supervised model) instead of the overall data. Thus, in order to avoid the drift problem of the global model, FedMix adds the global model from the previous round to the model parameter aggregation, i.e., $\omega_t = \alpha\psi_t + \beta\sigma_t + \gamma\omega_{t-1}$. Meanwhile, we conducted a sensitivity experiment of model performance to different hyperparameter weights to find the optimal weight combination.

2) FedMix is robust to different levels of non-IID data.

In our experiment, we introduced the Dirichlet distribution function to simulate the local client non-IID data in FL. In details, we generate data distributions of different non-IID levels by adjusting the parameters of the Dirichlet distribution function, i.e., $\mu = \{0.1, 1, 10, 100\}$ respectively correspond to different levels of non-IID. The results show that the performance difference of our model does not exceed 1% under different levels of non-IID settings. FedMatch uses a pseudo-random method to generate the non-IID data distribution of each client. However, in reality there is no such distribution, which will cause the model to lose its robustness.

VII. CONCLUSION

In this paper, we studied the labels-at-server scenario and addressed the problem of data availability and non-IID in FL. To solve the first problem, we designed a robust SSFL system that uses the FedMix algorithm to achieve high-precision semi-supervised learning. To tackle the non-IID problem, we propose a novel aggregation algorithm FedFreq, which effectively achieves the stable performance of the global model in the training process without adding additional computational overhead. Through experimental verification, our robust SSFL system is significantly better than the baseline in performance. In future work, we will further improve the algorithm to maximize the use of unlabeled data. Furthermore, we will continue to strengthen the theory of SSFL so that it can be better applied in real-world scenarios.

ACKNOWLEDGMENT

This work is supported by Heilongjiang Provincial Natural Science Foundation of China (Grant No. LH2020F044), the

2019–“Chunhui Plan” Cooperative Scientific Research Project of Ministry of Education of China (Grant No. HLJ2019015), the Fundamental Research Funds for Heilongjiang Universities, China (Grant No. 2020-KYYWF-1014), and also supported by NSFC under grant No. 6210071210.

REFERENCES

- [1] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, “Communication-efficient learning of deep networks from decentralized data,” in *Artificial Intelligence and Statistics*. PMLR, 2017, pp. 1273–1282.
- [2] J. Kang, Z. Xiong, D. Niyato, S. Xie, and J. Zhang, “Incentive mechanism for reliable federated learning: A joint optimization approach to combining reputation and contract theory,” *IEEE Internet of Things Journal*, vol. 6, no. 6, pp. 10700–10714, 2019.
- [3] Y. Liu, J. James, J. Kang, D. Niyato, and S. Zhang, “Privacy-preserving traffic flow prediction: A federated learning approach,” *IEEE Internet of Things Journal*, vol. 7, no. 8, pp. 7751–7763, 2020.
- [4] W. Y. B. Lim, Z. Xiong, J. Kang, D. Niyato, C. S. Leung, C. Miao, and S. Shen, “When information freshness meets service latency in federated learning: A task-aware incentive scheme for smart industries,” *IEEE Transactions on Industrial Informatics*, 2020.
- [5] Y. Liu, S. Garg, J. Nie, Y. Zhang, Z. Xiong, J. Kang, and M. S. Hossain, “Deep anomaly detection for time-series data in industrial iot: a communication-efficient on-device federated learning approach,” *IEEE Internet of Things Journal*, vol. 8, no. 8, pp. 6348–6358, 2020.
- [6] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith, “Federated optimization in heterogeneous networks,” *arXiv preprint arXiv:1812.06127*, 2018.
- [7] Y. Liu, X. Yuan, Z. Xiong, J. Kang, X. Wang, and D. Niyato, “Federated learning for 6g communications: Challenges, methods, and future directions,” *China Communications*, vol. 17, no. 9, pp. 105–118, 2020.
- [8] J. Kang, Z. Xiong, D. Niyato, Y. Zou, Y. Zhang, and M. Guizani, “Reliable federated learning for mobile networks,” *IEEE Wireless Communications*, vol. 27, no. 2, pp. 72–80, 2020.
- [9] Y. Liu, X. Yuan, R. Zhao, Y. Zheng, and Y. Zheng, “Rc-ssfl: Towards robust and communication-efficient semi-supervised federated learning system,” *arXiv preprint arXiv:2012.04432*, 2020.
- [10] Z. Long, L. Che, Y. Wang, M. Ye, J. Luo, J. Wu, H. Xiao, and F. Ma, “Fedsemi: An adaptive federated semi-supervised learning framework,” *arXiv preprint arXiv:2012.03292*, 2020.
- [11] Y. Jin, X. Wei, Y. Liu, and Q. Yang, “Towards utilizing unlabeled data in federated learning: A survey and prospective,” *arXiv e-prints*, pp. arXiv–2002, 2020.
- [12] W. Jeong, J. Yoon, E. Yang, and S. J. Hwang, “Federated semi-supervised learning with inter-client consistency & disjoint learning,” in *International Conference on Learning Representations*, 2020.
- [13] Y. Zhu, Y. Liu, J. James, and X. Yuan, “Semi-supervised federated learning for travel mode identification from gps trajectories,” *IEEE Transactions on Intelligent Transportation Systems*, 2021.
- [14] B. Wang, A. Li, H. Li, and Y. Chen, “Graphfl: A federated learning framework for semi-supervised node classification on graphs,” *arXiv preprint arXiv:2012.04187*, 2020.
- [15] L. Samuli and A. Timo, “Temporal ensembling for semi-supervised learning,” in *International Conference on Learning Representations (ICLR)*, vol. 4, no. 5, 2017, p. 6.
- [16] S. Park, J. Park, S.-J. Shin, and I.-C. Moon, “Adversarial dropout for supervised and semi-supervised learning,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018.
- [17] K. Sohn, D. Berthelot, N. Carlini, Z. Zhang, H. Zhang, C. A. Raffel, E. D. Cubuk, A. Kurakin, and C.-L. Li, “Fixmatch: Simplifying semi-supervised learning with consistency and confidence,” *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [18] X. Li, M. JIANG, X. Zhang, M. Kamp, and Q. Dou, “Fedbn: Federated learning on non-iid features via local batch normalization,” in *International Conference on Learning Representations*, 2020.
- [19] X. Zhu and A. B. Goldberg, “Introduction to semi-supervised learning,” *Synthesis lectures on artificial intelligence and machine learning*, vol. 3, no. 1, pp. 1–130, 2009.
- [20] O. Chapelle, B. Scholkopf, and A. Zien, “Semi-supervised learning (chapelle, o. et al., eds.; 2006)[book reviews],” *IEEE Transactions on Neural Networks*, vol. 20, no. 3, pp. 542–542, 2009.
- [21] D. P. Kingma, S. Mohamed, D. J. Rezende, and M. Welling, “Semi-supervised learning with deep generative models,” in *Advances in neural information processing systems*, 2014, pp. 3581–3589.
- [22] X. Zhai, A. Oliver, A. Kolesnikov, and L. Beyer, “S4l: Self-supervised semi-supervised learning,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 1476–1485.
- [23] P. K. Mallapragada, R. Jin, A. K. Jain, and Y. Liu, “Semiboost: Boosting for semi-supervised learning,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 31, no. 11, pp. 2000–2014, 2008.
- [24] D.-H. Lee et al., “Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks,” in *Workshop on challenges in representation learning, ICLR*, vol. 3, no. 2, 2013.
- [25] A. Tarvainen and H. Valpola, “Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results,” in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2017, pp. 1195–1204.
- [26] S. Itahara, T. Nishio, Y. Koda, M. Morikura, and K. Yamamoto, “Distillation-based semi-supervised federated learning for communication-efficient collaborative training with non-iid private data,” *IEEE Transactions on Mobile Computing*, pp. 1–1, 2021.
- [27] Z. Zhang, Y. Yang, Z. Yao, Y. Yan, J. E. Gonzalez, and M. W. Mahoney, “Improving semi-supervised federated learning by reducing the gradient diversity of models,” *arXiv preprint arXiv:2008.11364*, 2020.
- [28] X. Li, K. Huang, W. Yang, S. Wang, and Z. Zhang, “On the convergence of fedavg on non-iid data,” *arXiv preprint arXiv:1907.02189*, 2019.
- [29] Y. Zhao, M. Li, L. Lai, N. Suda, D. Civin, and V. Chandra, “Federated learning with non-iid data,” *arXiv preprint arXiv:1806.00582*, 2018.
- [30] F. Sattler, S. Wiedemann, K.-R. Müller, and W. Samek, “Robust and communication-efficient federated learning from non-iid data,” *IEEE transactions on neural networks and learning systems*, vol. 31, no. 9, pp. 3400–3413, 2019.
- [31] C. Briggs, Z. Fan, and P. Andras, “Federated learning with hierarchical clustering of local updates to improve training on non-iid data,” in *2020 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2020, pp. 1–9.
- [32] H. Wang, Z. Kaplan, D. Niu, and B. Li, “Optimizing federated learning on non-iid data with reinforcement learning,” in *IEEE INFOCOM 2020-IEEE Conference on Computer Communications*. IEEE, 2020, pp. 1698–1707.
- [33] Y. Chen, Y. Ning, M. Slawski, and H. Rangwala, “Asynchronous online federated learning for edge devices with non-iid data,” in *2020 IEEE International Conference on Big Data (Big Data)*. IEEE, 2020, pp. 15–24.
- [34] Q. Li, Y. Diao, Q. Chen, and B. He, “Federated learning on non-iid data silos: An experimental study,” *arXiv preprint arXiv:2102.02079*, 2021.
- [35] J. Wang, Q. Liu, H. Liang, G. Joshi, and H. V. Poor, “Tackling the objective inconsistency problem in heterogeneous federated optimization,” *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [36] M. Sajjadi, M. Javanmardi, and T. Tasdizen, “Regularization with stochastic transformations and perturbations for deep semi-supervised learning,” *arXiv preprint arXiv:1606.04586*, 2016.
- [37] Q. Xie, Z. Dai, E. Hovy, T. Luong, and Q. Le, “Unsupervised data augmentation for consistency training,” *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [38] X. Yang, Z. Song, I. King, and Z. Xu, “A survey on deep semi-supervised learning,” *arXiv preprint arXiv:2103.00550*, 2021.
- [39] D. Berthelot, N. Carlini, I. Goodfellow, N. Papernot, A. Oliver, and C. A. Raffel, “Mixmatch: A holistic approach to semi-supervised learning,” *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [40] M. Yurochkin, M. Agarwal, S. Ghosh, K. Greenewald, N. Hoang, and Y. Khazaeni, “Bayesian nonparametric federated learning of neural networks,” in *International Conference on Machine Learning*. PMLR, 2019, pp. 7252–7261.
- [41] T.-M. H. Hsu, H. Qi, and M. Brown, “Measuring the effects of non-identical data distribution for federated visual classification,” *arXiv preprint arXiv:1909.06335*, 2019.