

# POSITIVE AND UNLABELED FEDERATED LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

We study the problem of learning from positive and unlabeled (PU) data in the federated setting, where each client only labels a little part of their dataset due to the limitation of resources and time. Different from the settings in traditional PU learning where the negative class consists of a single class, the negative samples which cannot be identified by a client in the federated setting may come from multiple classes which are unknown to the client. Therefore, existing PU learning methods can be hardly applied in this situation. To address this problem, we propose a novel framework, namely Federated learning with Positive and Unlabeled data (FedPU), to minimize the expected risk of multiple negative classes by leveraging the labeled data in other clients. We theoretically prove that the proposed FedPU can achieve a generalization bound which is no worse than  $C\sqrt{C}$  times (where  $C$  denotes the number of classes) of the fully-supervised model. Empirical experiments show that the FedPU can achieve much better performance than conventional learning methods which can only use positive data.

## 1 INTRODUCTION

With the development of edge devices (*e.g.*, cameras, microphones, and GPS), more and more decentralized data are collected and locally stored by different users. Due to the privacy and transmission concerns, users are unwilling or not allowed to share the data with each other. In this case, classical machine learning scheme can hardly learn a globally effective model for all the users. Therefore, federated learning (McMahan et al. (2017)) is proposed to derive a model with high performance in the central server by leveraging multiple local models trained by users (clients) themselves, which ensures the privacy of the local data.

Typically, there is a common assumption in federated learning that the local data (private data) stored on user devices is well refined (*i.e.*, all of the local data is labeled with ground truth). However, considering the limitation of time and resources, only part of the private data in each client are labeled in reality. To this end, some of the previous works were proposed to address this federated learning problem following a semi-supervised scheme. Jeong et al. (2021) proposed the FedMatch algorithm which introduced a new inter-client consistency loss and decomposed the parameters for labeled and unlabeled data. Zhang et al. (2020) managed to solve this problem by conducting a novel grouping-based model average method and improved the convergence efficiency. Itahara et al. (2020) proposed a distillation-based algorithm to exchange the local models among each client and learned the unlabeled data by pseudo labels. Although these methods can successfully address the semi-supervised learning problem for federated learning, they assume that there are labeled samples in the central server in order to help the learning of unlabeled data in each client. However, such assumption violates the universal federated learning setting that data only exists on local clients.

We consider the most general setting of federated learning with unlabeled data: 1) each client only labels *part of* their own data which comes from *part of* the classes; 2) there are no data in the central server; 3) nothing except parameters of models can be exchanged between clients and the central server. Note that the first constraint of our setting meets the problem of learning from positive and unlabeled (PU) data. Existing PU methods (Liu et al. (2003); Liu & Tao (2015); Xu et al. (2017)) focused on solving the PU problem regard the negative class (class that contains no labeled samples) as a single class. However, since negative class in one client may consist of multiple positive classes in other clients, there are multiple negative classes in one client in federated learning, which results in a multiple-positive-multiple-negative PU (MPMN-PU) learning problem.

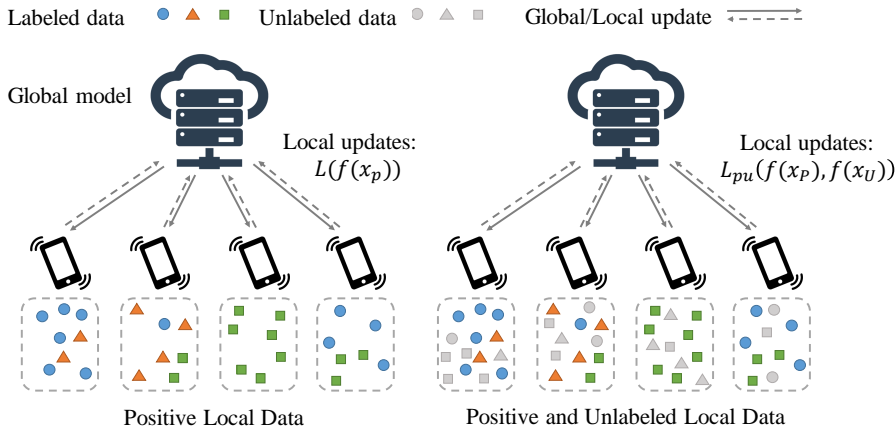


Figure 1: Illustration of the conventional federated learning (left) and the proposed method (right). Conventional federated learning method only learns from labeled data. In contrast, we propose the federated learning with positive and unlabeled data to fully inherit the information from the unlabeled data.

In this paper, we propose to solve the aforementioned problem with the Federated Averaging with Positive and Unlabeled data (FedPU) algorithm, where the local model in each client is trained with MPMN-PU data. We first analyze the expected risk of each class in each client and show that the risks of multiple negative classes can be successfully minimized by leveraging unlabeled data in this client and labeled data in other clients, which is shown in Figure 1. Moreover, we present a generalization bound of proposed FedPU and show that the FedPU algorithm is no worse than  $C\sqrt{C}$  times (where  $C$  denotes the number of classes) of the fully-supervised model in federated setting. Experiments on MNIST and CIFAR-10 datasets empirically show that the proposed method can achieve better performance than existing federated learning algorithms.

## 2 RELATED WORKS

In this section, we briefly review the related works about the federated learning and positive-unlabeled learning.

### 2.1 FEDERATED LEARNING

Federated learning was firstly proposed by McMahan et al. (2017) to collaboratively learn a model without collecting data from the participants. A lot of recent works (Li et al. (2018); Karimireddy et al. (2020); Sattler et al. (2019)) focused on improving model performance on non-iid data and proposed advanced versions of McMahan et al. (2017). Recently, Yu et al. (2020) proposed FedAwG for the setting of distributed learning with only positive data. To the best of our knowledge, there is no existing work on federated learning with positive and unlabeled data.

### 2.2 POSITIVE AND UNLABELED LEARNING

Various effective algorithms have been developed to solve the PU learning problem. The two-step technique was proposed by Liu et al. (2003) based on the assumption that all the positive samples are similar to the labeled examples and the negative samples are very different from them. Biased PU learning methods Liu & Tao (2015) treated the unlabeled samples as negative ones with label noise. Lee & Liu (2003) regarded unlabeled data as negative data with smaller weights, then performed logistic regression after weighting the samples to handle the situation that noise rate is greater than a half. In order to avoid tuning the weights, Elkan & Noto (2008) regarded unlabeled data as weighted positive and negative data simultaneously. du Plessis et al. (2014) proposed the first unbiased risk estimator, then Kiryo et al. (2017) made a progress by proposing a non-negative risk estimator for PU learning to mitigate the overfitting problem when using a flexible model. Xu et al. (2017) adapted PU learning to the setting with multi-class data input, in which labeled data are from  $k - 1$  positive

classes and unlabeled data are from all the  $k$  classes. These methods regard the negative class as a single class, which is reasonable when there is only a single dataset. However, in federated learning the datasets are distributed in different clients, where samples from the negative classes in one client may become positive in another client since different clients are free to label their data. To this end, an effective PU learning algorithm for the federated setting is urgently required.

### 3 METHOD

In this section, we study federated learning problem under the MPMN-PU learning setting for each client.

#### 3.1 PROBLEM SETUP

Here we first introduce the definition in federated learning, which consists of  $K$  different clients and one central server. Given the data space  $\mathcal{S}$  and the hypothesis space of parameters  $\mathcal{W}$ , the training data is distributed on  $K$  different clients and is generated from the data space  $\mathcal{S}$ , which is denoted as  $\{\mathbf{S}_k\}_{k=1}^K \in \mathcal{S}$ . Denote  $T$  as the number of communication rounds and  $w_t \in \mathcal{W}$  as the weight matrix in the central server in time  $t \in \{1, \dots, T\}$ , the weights  $w_t$  is first transferred from the central server to each client, and then updated using the training data in each client respectively and derive  $K$  different weights  $w_{t+1}^k$ :

$$w_{t+1}^k \leftarrow \text{ClientUpdate}(k, w_t), \quad (1)$$

where the client update stage is a conventional training method for updating the gradient. After that, the updated weights are then transferred back to renew the weight matrix in central server:

$$w_{t+1} \leftarrow \sum_{k=1}^K \frac{n^k}{n} w_{t+1}^k, \quad (2)$$

where  $n^k$  is the number of training samples in client  $k$  and  $n = \sum_{k=1}^K n^k$ .

In the traditional federated learning setting, the training data in each client is fully labeled. Nevertheless, samples are not always fully labeled in many real world scenario because of the time and resources limitation in each client. Specifically, the training data  $\mathbf{S}_k$  in client  $k$  consists of positive data  $\mathbf{P}_k$  and unlabeled data  $\mathbf{U}_k$ , which can be formulated as:

$$\mathbf{S}_k = \mathbf{P}_k \cup \mathbf{U}_k, \quad k = 1, \dots, K. \quad (3)$$

Given the set of classes as  $\mathbf{C} = 1, \dots, C$  in which  $C$  is the total number of classes, the classes of positive data in client  $k$  is denoted as  $\mathbf{C}_{\mathbf{P}_k}$ , which is the subset of  $\mathbf{C}$ , *i.e.*,  $\mathbf{C}_{\mathbf{P}_k} \subset \mathbf{C}$ . In other words, each client can only identify part of the classes from the dataset  $\mathbf{S}_k$ . Besides, only a portion of the data in the positive classes can be labeled since the data is too much to be fully labeled. Therefore, there exists unlabeled data from not only the negative classes but also the positive classes, *i.e.*,  $\mathbf{C}_{\mathbf{U}_k} = \mathbf{C}$ . Specifically, we have:

$$\forall x \in \mathbf{P}_k, \mathbf{Class}(x) \in \mathbf{C}_{\mathbf{P}_k}; \forall x \in \mathbf{U}_k, \mathbf{Class}(x) \in \mathbf{C}. \quad (4)$$

Note that different clients have different set of positive classes, and all of the positive classes should cover the whole classes in the dataset, *i.e.*,  $\bigcup_{\mathbf{P}_k} \mathbf{C}_{\mathbf{P}_k} = \mathbf{C}$ .

In this setting, the conventional federated learning algorithms cannot be directly applied. Fortunately, positive and unlabeled learning Liu et al. (2003) has been proposed to solve this problem. However, we meet a MPMN-PU learning problem, which cannot be directly handled by existing methods.

#### 3.2 FEDERATED LEARNING WITH POSITIVE AND UNLABELED DATA

In this section, we propose our Federated learning with Positive and Unlabeled data (FedPU) method.

We first present our MPMN-PU learning scheme. Denote the training samples as  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n \in \mathbf{S}$ . In classical multi-class classification, given the class prior  $\pi_i = p(y = i), i = 1, 2, \dots, C$ , the classifier

$f(\mathbf{x}; w)$  (short as  $f(\mathbf{x})$ ), in which  $w$  is the parameter of the classifier, can be learned by minimizing the expected misclassification rate:

$$R(f) = \sum_{i=1}^C \pi_i R_i(f), \quad (5)$$

where  $\sum_{i=1}^C \pi_i = 1$  and  $R_i(f) = P_i(f(\mathbf{x}) \neq i)$  denotes the expected misclassification rate on  $i$ -th class.

However, in MPMN-PU setting, only samples in a few classes are labeled in the training set for each client. Some of classes in Eq. 5 is unlabeled and the expected risk cannot be calculated in each client. Therefore, we introduce  $R_U(f)$  to denote the probability that the unlabeled samples have not been classified to the set of unlabeled classes:

$$\begin{aligned} R_U(f) &= \sum_{m \notin \mathbf{C}_{\mathbf{P}}} P_U(f(\mathbf{x}) \neq m) \\ &= \sum_{i \in \mathbf{C}_{\mathbf{P}}} \sum_{m \notin \mathbf{C}_{\mathbf{P}}} \pi_i P_i(f(\mathbf{x}) \neq m) + \sum_{j \notin \mathbf{C}_{\mathbf{P}}} \sum_{m \notin \mathbf{C}_{\mathbf{P}}} \pi_j P_j(f(\mathbf{x}) \neq m) \\ &= \sum_{i \in \mathbf{C}_{\mathbf{P}}} \sum_{m \notin \mathbf{C}_{\mathbf{P}}} \pi_i P_i(f(\mathbf{x}) \neq m) + \sum_{j \notin \mathbf{C}_{\mathbf{P}}} \pi_j P_j(f(\mathbf{x}) \neq j) + \sum_{j, m \notin \mathbf{C}_{\mathbf{P}}, j \neq m} \pi_j P_j(f(\mathbf{x}) \neq m). \end{aligned} \quad (6)$$

Therefore, Eq. 5 can be reformulated as:

$$\begin{aligned} R(f) &= \sum_{i \in \mathbf{C}_{\mathbf{P}}} \pi_i R_i(f) + \sum_{j \notin \mathbf{C}_{\mathbf{P}}} \pi_j R_j(f) \\ &= \sum_{i \in \mathbf{C}_{\mathbf{P}}} \pi_i R_i(f) + R_U(f) - \sum_{i \in \mathbf{C}_{\mathbf{P}}} \sum_{m \notin \mathbf{C}_{\mathbf{P}}} \pi_i P_i(f(\mathbf{x}) \neq m) - \sum_{j, m \notin \mathbf{C}_{\mathbf{P}}, j \neq m} \pi_j P_j(f(\mathbf{x}) \neq m) \\ &= \sum_{i \in \mathbf{C}_{\mathbf{P}}} \pi_i [P_i(f(\mathbf{x}) \neq i) - \sum_{m \notin \mathbf{C}_{\mathbf{P}}} P_i(f(\mathbf{x}) \neq m)] \\ &\quad + \sum_{m \notin \mathbf{C}_{\mathbf{P}}} P_U(f(\mathbf{x}) \neq m) - \sum_{j, m \notin \mathbf{C}_{\mathbf{P}}, j \neq m} \pi_j P_j(f(\mathbf{x}) \neq m). \end{aligned} \quad (7)$$

Now we are ready to solve the federated learning problem with MPMN-PU data. In the federated learning setting, the expected risk can be formulated as:

$$R(f) = \sum_{k=1}^K R^k(f), \quad (8)$$

where  $R^k(f)$  denote the expected risk in client  $k$ . Given Eq. 7, the corresponding expectation of the expected risk using a surrogate loss function  $l$  in PU setting can be reformulated as:

$$\begin{aligned} \mathbb{E}[l(R^k(f))] &= \sum_{i \in \mathbf{C}_{\mathbf{P}_k}} \pi_i \mathbb{E}_i^k \left[ l(f(\mathbf{x}) \neq i) - \sum_{m \notin \mathbf{C}_{\mathbf{P}_k}} l(f(\mathbf{x}) \neq m) \right] \\ &\quad + \sum_{m \notin \mathbf{C}_{\mathbf{P}_k}} \mathbb{E}_U^k [l(f(\mathbf{x}) \neq m)] - \sum_{j, m \notin \mathbf{C}_{\mathbf{P}_k}, j \neq m} \pi_j \mathbb{E}_j^k [l(f(\mathbf{x}) \neq m)], \end{aligned} \quad (9)$$

where  $\mathbb{E}_i^k$  means the expectation for the labeled data of  $i$ th class in client  $k$ , and  $\mathbb{E}_U^k$  means the expected risk for unlabeled data in client  $k$ .

Note that the federated MPMN-PU learning problem has several negative classes, which is fundamentally different with conventional PU learning problem Liu et al. (2003); Xu et al. (2017) whose negative class is a single class. We have an additional term  $\sum_{j, m \notin \mathbf{C}_{\mathbf{P}_k}, j \neq m} \pi_j \mathbb{E}_j^k [l(f(\mathbf{x}) \neq m)]$  in Eq. 9, which cannot be directly calculated in client  $k$  since there is no labeled data for these classes. Thus, existing PU learning methods failed to address this situation.

Fortunately, we have  $\bigcup_{\mathbf{P}_k} \mathbf{C}_{\mathbf{P}_k} = \mathbf{C}$ , which means that although we have no information for the negative class in one client, there exists labeled data for these classes in other clients. Since the weights in central server is derived from the combination of each client, we can use the labeled

data in other clients to handle the expected risk that cannot be directly calculated, which can be formulated as:

$$\begin{aligned} \mathbb{E}_j^{k_1} [l(f(\mathbf{x}) \neq m)] &\leftarrow \mathbb{E}_j^{k_2} [l(f(\mathbf{x}) \neq m)], \\ \text{where } j &\notin \mathbf{C}_{\mathbf{P}_{k_1}}, j \in \mathbf{C}_{\mathbf{P}_{k_2}}. \end{aligned} \quad (10)$$

Specifically, we use labeled samples from the  $j$ -th class in client  $k_2$  to help calculating the risk of  $j$ -th class which is treated as negative in client  $k_1$ . According to the Eq. 8, the overall risk  $R(f)$  remains the same after applying this approximation.

By applying Eq. 10 to Eq. 9, we can successfully formulated the PU learning risk as:

$$\begin{aligned} \mathbb{E}[l(R^k(f))] &= \sum_{i \in \mathbf{C}_{\mathbf{P}_k}} \pi_i \mathbb{E}_i^k \left[ l(f(\mathbf{x}) \neq i) - \sum_{m \notin \mathbf{C}_{\mathbf{P}_k}} l(f(\mathbf{x}) \neq m) \right] \\ &+ \sum_{m \notin \mathbf{C}_{\mathbf{P}_k}} \mathbb{E}_U^k [l(f(\mathbf{x}) \neq m)] - \sum_{k_q \neq k} \sum_{i \in \mathbf{C}_{\mathbf{P}_k}, i, m \notin \mathbf{C}_{\mathbf{P}_{k_q}}, i \neq m} \pi_i \mathbb{E}_i^k [l(f(\mathbf{x}) \neq m)], \end{aligned} \quad (11)$$

Different with Eq. 8 that contains risk of negative classes, the above equation can be easily minimized since it only consists of the risk of positive data and unlabeled data. Therefore, the final loss function for updating the weights  $w$  in each client  $k$  using its data  $\{\mathbf{x}^k\} \in \mathbf{S}_k$  can be formulated as:

$$\begin{aligned} L(f(w)) &= \sum_{i \in \mathbf{C}_{\mathbf{P}_k}} \pi_i \sum_{j=1}^{n_i^k} \left[ l(f(\mathbf{x}_j^k) \neq i) - \sum_{m \notin \mathbf{C}_{\mathbf{P}_k}} l(f(\mathbf{x}_j^k) \neq m) \right] \\ &- \sum_{k_q \neq k} \sum_{i \in \mathbf{C}_{\mathbf{P}_k}, i, m \notin \mathbf{C}_{\mathbf{P}_{k_q}}, i \neq m} \pi_i \sum_{j=1}^{n_i^k} [l(f(\mathbf{x}_j^k) \neq m)] + \sum_{m \notin \mathbf{C}_{\mathbf{P}_k}} \sum_{j=1}^{n_U^k} [l(f(\mathbf{x}_j^k) \neq m)], \end{aligned} \quad (12)$$

where  $n_i^k$  and  $n_U^k$  denotes the number of  $i$ -th samples and unlabeled samples in client  $k$ , respectively.

By minimizing the above loss function in each client, the overall expected risk in Eq. 5 can be minimized.

### 3.3 THEORETICAL ANALYSIS

We use the loss function in Eq. 12 to minimize the risk in Eq. 11. In this section, we analyze the generation bound of the proposed FedPU. We first evaluate the bound in each client. Then the overall bound can be derived by summing these bounds. Note that the proof of theorems and lemma can be found in the supplementary materials.

Since the Eq. 11 has three terms, we begin with the first and second terms.

**Theorem 1** Fix  $f \in \mathcal{F}$ , for any  $0 < \delta < 1$ , with probability at least  $1 - \delta$ , the generalization bound holds:

$$\begin{aligned} &\mathbb{E}_i^k \left[ l(f(\mathbf{x}) \neq i) - \sum_{m \notin \mathbf{C}_{\mathbf{P}_k}} l(f(\mathbf{x}) \neq m) \right] - \frac{1}{n_i^k} \sum_{j=1}^{n_i^k} \left[ l(f(\mathbf{x}_j) \neq i) - \sum_{m \notin \mathbf{C}_{\mathbf{P}_k}} l(f(\mathbf{x}_j) \neq m) \right] \\ &\leq 2CV \left( \sum_{s \in \mathbf{C}_{\mathbf{P}_k}} \frac{1}{\sqrt{n_s^k}} + \frac{1}{\sqrt{n_U^k}} \right) + \sqrt{\frac{\log \frac{1}{\delta}}{2n_i^k}}, \end{aligned} \quad (13)$$

where  $i \in \mathbf{C}_{\mathbf{P}_k}$ ,  $V$  is a constant related to the VC-dimension of  $f$  and the bound of loss function  $l$ ,  $n_i^k$  and  $n_U^k$  denotes the number of samples in  $i$ -class and unlabeled classes in  $k$ -th client, respectively.

**Theorem 2** Fix  $f \in \mathcal{F}$ , for any  $0 < \delta < 1$ , with probability at least  $1 - \delta$ , the generalization bound holds:

$$\mathbb{E}_i^k [l(f(\mathbf{x}) \neq m)] - \frac{1}{n_i^k} \sum_{j=1}^{n_i^k} l(f(\mathbf{x}_j) \neq m) \leq CV \left( \sum_{s \in \mathbf{C}_{\mathbf{P}_k}} \frac{1}{\sqrt{n_s^k}} + \frac{1}{\sqrt{n_U^k}} \right) + \sqrt{\frac{\log \frac{1}{\delta}}{2n_i^k}}. \quad (14)$$

Theorem 1 and 2 presents the generalization bound for the labeled data in each class, which can be summarized to get the error bound for the first two terms in Eq. 11.

However, it is difficult to derive the error bound of the last term in Eq. 11 since the expectation is calculated on unlabeled data, so we decompose this term using the following lemma.

**Lemma 3** Define

$$l'(f(\mathbf{x}) \neq m) = \frac{k^{C_{U_k}}}{k^{C_{U_k}} + \prod_{i \notin \mathbf{C}_{\mathbf{P}_k}} |k - i|} l(f(\mathbf{x}) \neq m), \quad (15)$$

where  $C_{U_k}$  denotes the number of unlabeled classes in client  $k$ . The last term in Eq. 11 can be decomposed as:

$$\begin{aligned} \sum_{m \notin \mathbf{C}_{\mathbf{P}_k}} \mathbb{E}_U^k [l(f(\mathbf{x}) \neq m)] &= \sum_{i \in \mathbf{C}_{\mathbf{P}_k}} \pi_i \left( \frac{\prod_{i \notin \mathbf{C}_{\mathbf{P}_k}} |k - i|}{k^{C_{U_k}}} \right) \sum_{m \notin \mathbf{C}_{\mathbf{P}_k}} \mathbb{E}_U^k [l'(f(\mathbf{x}) \neq m)] \\ &+ \sum_{m \notin \mathbf{C}_{\mathbf{P}_k}} \mathbb{E}_U^k [l'(f(\mathbf{x}) \neq m)]. \end{aligned} \quad (16)$$

Therefore, based on Lemma 3, we can present the generalization bound for the unlabeled data with the following theorem.

**Theorem 4** Fix  $f \in \mathcal{F}$ , for any  $0 < \delta < 1$ , with probability at least  $1 - \delta$ , the generalization bound holds:

$$\begin{aligned} &\sum_{m \notin \mathbf{C}_{\mathbf{P}_k}} \mathbb{E}_U^k [l(f(\mathbf{x}) \neq m)] - \frac{1}{n_U^k} \sum_{j=1}^{n_U^k} \sum_{m \notin \mathbf{C}_{\mathbf{P}_k}} l'(f(\mathbf{x}_j) \neq m) \\ &\leq \sum_{i \in \mathbf{C}_{\mathbf{P}_k}} \frac{\pi_i}{n_i^k} \left( 1 + \frac{\prod_{i \notin \mathbf{C}_{\mathbf{P}_k}} |k - i|}{k^{C_{U_k}}} \right) \sum_{j=1}^{n_i^k} \sum_{m \notin \mathbf{C}_{\mathbf{P}_k}} l'(f(\mathbf{x}_j) \neq m) \\ &+ \left( \sum_{i \in \mathbf{C}_{\mathbf{P}_k}} \pi_i + 1 \right) CV \left( \sum_{s \in \mathbf{C}_{\mathbf{P}_k}} \frac{1}{\sqrt{n_s^k}} + \frac{1}{\sqrt{n_U^k}} \right) + \sum_{i \in \mathbf{C}_{\mathbf{P}_k}} \pi_i \left( 1 + \frac{\prod_{i \notin \mathbf{C}_{\mathbf{P}_k}} |k - i|}{k^{C_{U_k}}} \right) \sqrt{\frac{\log \frac{1}{\delta}}{2n_i^k}} + \sqrt{\frac{\log \frac{1}{\delta}}{2n_U^k}}. \end{aligned} \quad (17)$$

Now we are ready to present the generalization bound for Eq. 11.

**Theorem 5** As  $n_i^k, n_U^k \rightarrow \infty$ ,  $i \in \mathbf{C}_{\mathbf{P}_k}, k \in \{1, \dots, K\}$ , the generalization bound of the proposed FedPU is of order:

$$\mathcal{O} \left( \sum_{k=1}^K C^2 \left( \sum_{i \in \mathbf{C}_{\mathbf{P}_k}} \frac{1}{\sqrt{n_i^k}} + \frac{1}{\sqrt{n_U^k}} \right) \right). \quad (18)$$

It should be noted that for fully labeled data, the generalization bound using federated learning should be of order  $\mathcal{O} \left( \sum_{k=1}^K \left( \frac{C^2}{\sqrt{\sum_{i \in \mathbf{C}_{\mathbf{P}_k}} n_i^k + n_U^k}} \right) \right)$ . As a result, the proposed method is no worse

than  $C\sqrt{C}$  times (assuming that each class has the same order of samples) of the fully-supervised models. Moreover, for the classical learning with fully labeled data (without federated learning), the

generalization bound would be of order  $\mathcal{O} \left( \frac{C^2}{\sqrt{\sum_{k=1}^K (\sum_{i \in \mathbf{C}_{\mathbf{P}_k}} n_i^k + n_U^k)}} \right)$ . Therefore, the proposed

method is no worse than  $CK\sqrt{CK}$  times of the fully-supervised models without federated learning.

Table 1: Classification result on iid data.

Num of Clients	Num of P-class	Overlap	Baseline-1	Proposed Method	Baseline-2
10	2	✓	85.47%	92.50%	97.95%
4	6	✓	92.10%	95.08%	98.05%
2	9	✓	93.15%	95.37%	98.20%
10	1	✗	37.13%	84.15%	97.95%
5	2	✗	73.41%	93.45%	98.03%
2	5	✗	74.00%	93.73%	98.20%

## 4 EXPERIMENTS

In this section, we show the experimental results of the proposed method in both iid data and non-iid data on the MNIST and CIFAR-10 dataset. We also conduct ablation study to verify the effectiveness of the proposed method in different settings.

We first detail the training strategy used in the following experiments. The SGD optimizer is used to train the network with momentum 0.5. For federated learning, we set the communication round as 200. For each client, the local epoch and local batchsize for training the network in each round is set as 1 and 100. The learning rate is initialized as 0.01 and exponentially decayed by 0.995 over communication rounds on the MNIST dataset. To show the effectiveness of the proposed method, we compare the proposed method with two different baselines. **Baseline-1** denotes that the network is trained using only positive data and FedAvg( McMahan et al. (2017)). **Baseline-2** denotes that the network is trained using fully-supervised data and FedAvg.

### 4.1 PERFORMANCE ON IID DATA WITH BALANCED POSITIVE CLASSES

We evaluate our method in iid setting of federated learning, where the training data in each client is uniformly sampled from the original dataset. Specifically, we uniformly divide the training set into  $K$  parts, where each part of data is class-imbalanced. Since the ability of each client is limited, only a few classes can be labeled. Moreover, only part of data in these classes is labeled. To fully investigate the ability of the proposed method, we conduct different settings as shown in Table 1, including using different number of clients ( $\{2, 4, 5, 10\}$ ) and different number of positive classes ( $\{1, 2, 5, 6, 9\}$ ). We also investigate the influence of overlap of positive classes between different clients. Only half of data in each positive class is labeled.

We conduct experiments on the MNIST dataset, which is composed of images with  $28 \times 28$  pixels from 10 categories. The MNIST dataset consists of 60,000 training images and 10,000 testing images. The results are shown in Table 1. We first investigate the setting that each client has overlap in positive classes and the number of clients varies from 2 to 10. The Baseline-1 trained with positive data can only achieve 85.47%, 92.10% and 93.15% accuracies for 10, 4 and 2 clients, respectively. It can be seen that as the number of clients increases, the data is more discrete, which makes the accuracies of learned networks lower. Although the Baseline-2 can achieve higher performance (97.95%, 98.05% and 98.20%), the networks should be trained with fully supervised data, which is usually unavailable in real-world applications. In contrast, the proposed method can achieve 92.50%, 95.08% and 95.37% accuracies, respectively, which is consistently higher than those of the Baseline-1 and comparable to Baseline-2.

We further investigate the non-overlap setting, where positive classes in each client are not overlapped. This setting is challenging, since the information of every class is contained in only one client. As a result, the Baseline-1 trained with positive data achieves only 37.13%, 73.41% and 74.00% accuracies for 10, 5 and 2 clients, respectively. The proposed FedPU can still achieve 84.15%, 93.45% and 93.73% accuracies by fully inheriting the information from the unlabeled data. These experiments show that the proposed method can perform well with iid data in federated setting.

Table 2: Classification results with different number of positive classes in each client.

Division of P-class	Overlap	Baseline-1	Proposed Method	Baseline-2
[2,3,4,6,7,8]	✓	93.84%	95.32%	97.91%
[1,2,4,6,7]	✓	93.81%	95.01%	98.03%
[2,4,6,8]	✓	92.27%	95.28%	98.05%
[3,7]	✗	89.68%	94.68%	98.20%
[2,3,5]	✗	71.46%	93.65%	98.16%
[1,2,3,4]	✗	74.27%	94.48%	98.05%

Table 3: Classification result on non-iid data.

Num of Partitions	Division of P-class	Overlap	Baseline-1	Proposed Method	Baseline-2
5	[2,2,...,2]	✓	25.47%	91.67%	97.47%
5	[1,1,...,1]	✗	24.37%	91.24%	97.47%
5	[4,4,3,3,2,2,1,1,1,1]	✓	76.92%	92.16%	97.47%
2	[1,1,...,1]	✗	69.24%	91.29%	96.19%

#### 4.2 PERFORMANCE ON IID DATA WITH IMBALANCED POSITIVE CLASSES

To further investigate the effectiveness of the proposed method, we study a more complicated setting that the number of positive classes is different in each client. The results are in shown in Table 2. For example, the division of P-class is [2, 3, 4, 6, 7, 8] means there are 6 clients consists of 2, 3, 4, 6, 7 and 8 positive classes, respectively. We also study both the overlap and non-overlap settings. The Baseline-1 achieves 93.84%, 93.81% and 92.27% accuracies for different divisions of positive classes in the overlap setting, while the proposed method can achieve 95.32%, 95.01% and 92.28% accuracies, respectively, which is much higher than Baseline-1. The results in non-overlap setting is worse than those in overlap setting, which is consistent with the results in Table 1 where the number of positive classes is the same in each client. The Baseline-1 achieves only 89.68%, 71.46% and 74.27% accuracies. When the number of clients grows, the performance of Baseline-1 drops dramatically. In contrast, the proposed method can achieve 94.68%, 93.65% and 94.48% accuracies, which surpasses those of the Baseline-1 and is stable with different numbers of clients. Note that although the Baseline-2 can achieve a  $\sim 98\%$  accuracy in all settings, it should trained with fully supervised data and violate most of the scenarios in real-world applications.

#### 4.3 PERFORMANCE ON NON-IID DATA

Another important setting for federated learning is that the data in different client is under the non-iid distribution. Therefore, we follow the settings in Li et al. (2018) to construct the non-iid data, where the data is sorted by class and divided to create two extreme cases: (a) 5-class non-iid, where the sorted data is divided into 50 partitions and each client is randomly assigned 5 partitions from 5 classes. (b) 2-class non-iid, where the sorted data is divided into 20 partitions and each client is randomly assigned 2 partitions from 2 classes.

For 5-class non-iid, we study three different divide settings for positive classes in each client, which is shown in Table 3. Compared with the iid setting, the non-iid setting is more challenging since the data distribution in each client is different and it is hard for the model to effectively learn the latent distribution on the whole dataset. Therefore, Baseline-1 can achieve only 25.47%, 24.37% and 79.62% accuracies when dealing with non-iid and unlabeled data, which is hard to optimize. In contrast, the proposed method can still achieve 91.67%, 91.24% and 92.16% accuracies, respectively, which outperforms Baseline-1 by a large margin.

For 2-class non-iid, Baseline-1 achieves a 69.24% accuracy while the proposed method achieves a 91.29% accuracy, which still shows the superiority of the proposed FedPU. It should be noted that although the baseline method can achieve  $\sim 98\%$  accuracy, it requires the fully labeled data to train the model in each client. In contrast, the proposed method requires only a small amount of labeled data and utilizes the information on unlabeled data to learn an effective model. In conclusion, the proposed method successfully learns the latent distribution from the positive and unlabeled data in the non-iid federated setting and achieve better performance than conventional federated learning methods.



Table 4: Classification results on CIFAR-10 dataset

Data Distribution	Division of P-class	Overlap	Baseline-1	Proposed Method	Baseline-2
iid	[2,2,2,2,2]	✗	65.52%	76.81%	81.13%
iid	[1,2,4,6,7]	✓	71.42%	75.41%	81.13%
non-iid	[2,2,2,2,2,2,2,2,2]	✓	52.39%	61.05%	72.61%
non-iid	[4,4,3,3,2,2,1,1,1]	✓	55.57%	65.73%	72.61%

Table 5: Comparison with semi-supervised methods on CIFAR-10 dataset

Methods	Supervised FedAVG	UDA	FixMatch	FedMatch	Ours
<b>IID Acc.</b>	80.25%	47.45%	47.20%	52.13%	58.25%
<b>Non-IID Acc.</b>	84.70%	46.31%	46.20%	52.25%	55.20%

#### 4.4 EXPERIMENTS ON CIFAR-10

After investigating the performance of the proposed FedPU on MNIST dataset, we further evaluate our method on the CIFAR-10 dataset. The CIFAR-10 dataset consists of 50,000 training images and 10,000 testing images with size  $32 \times 32 \times 3$  from 10 categories. The training strategy is the same as that on the MNIST dataset. As shown in Table 4, experiments on different settings (*e.g.*, data distribution, division of positive classes, overlap), are conducted to evaluate the effectiveness of the proposed method.

We first investigate the results on the iid data. 5 clients with 2 positive classes in each client are used to train the model. The positive classes have no overlap. The FedAvg method trained with positive data achieves only a 62.52% accuracy. The proposed method achieves a 76.71% accuracy with the help of unlabeled data, which is more close to the result (81.13%) trained with fully supervised data. Then, we turn to explore the challenging situation that each client has different number of positive class ([1,2,4,6,7]). The proposed method still achieves a 75.41% accuracy, which is much higher than that of Baseline-1 (71.42%).

We further construct the non-iid data on CIFAR-10 dataset following Li et al. (2018), where each class of the training data is randomly divided into 5 partitions (50 partitions for 10 classes) and each client is randomly assigned 5 partitions from 5 classes. We also investigate the situation that each client has the same/different number of positive classes. As shown in Table 4, the models trained by the proposed method achieve accuracies of 61.05% and 65.73% and surpass those trained with the baseline method by a large margin (8.66% and 10.16%). In conclusion, the proposed method significantly improves the performance of the existing federated learning method in different settings on CIFAR-10 dataset.

#### 4.5 COMPARISON WITH SEMI-SUPERVISED METHODS

To further show the superiority of the proposed method, we conduct comparison with the semi-supervised algorithms in federated setting. We follow the setting in Jeong et al. (2021) to use CIFAR-10 datasets. Specifically, for the IID setting, 5 labeled images are extracted in per class for each client (100 clients) and the rest of images are used as unlabeled data. For the Non-IID setting, the distribution of the number of images per class for each client is different. Table 5 shows the performance of FedAVG using supervised data, UDA (Xie et al. (2019)), FixMatch (Sohn et al. (2020)), FedMatch (Jeong et al. (2021)) and the proposed FedPU. The proposed method achieve the state-of-the-art performance among all semi-supervised methods.

## 5 CONCLUSION

We study a real-world setting in federated learning problem, where each client could only label limited number of data in part of classes. Existing federated learning algorithms can hardly achieve satisfying performance since they cannot minimize the expected risk for each class in each client. To address this problem, we propose the Federated learning with Positive and Unlabeled data (FedPU) algorithm, which can effectively learn from both labeled and unlabeled data for each client. Theoretical analysis and empirical experiments demonstrate that the proposed method can achieve better performance than the conventional federated learning method learned by the positive data.

## REFERENCES

- Olivier Bousquet, Stéphane Boucheron, and Gábor Lugosi. Introduction to statistical learning theory. In *Summer School on Machine Learning*, pp. 169–207. Springer, 2003.
- Marthinus Christoffel du Plessis, Gang Niu, and Masashi Sugiyama. Analysis of learning from positive and unlabeled data. In *NIPS*, pp. 703–711, 2014.
- Charles Elkan and Keith Noto. Learning classifiers from only positive and unlabeled data. In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 213–220, 2008.
- Sohei Itahara, Takayuki Nishio, Yusuke Koda, Masahiro Morikura, and Koji Yamamoto. Distillation-based semi-supervised federated learning for communication-efficient collaborative training with non-iid private data. *arXiv preprint arXiv:2008.06180*, 2020.
- Wonyong Jeong, Jaehong Yoon, Eunho Yang, and Sung Ju Hwang. Federated semi-supervised learning with inter-client consistency & disjoint learning. *ICLR*, 2021.
- Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In *International Conference on Machine Learning*, pp. 5132–5143. PMLR, 2020.
- Ryuichi Kiryo, Gang Niu, Marthinus C du Plessis, and Masashi Sugiyama. Positive-unlabeled learning with non-negative risk estimator. *arXiv preprint arXiv:1703.00593*, 2017.
- Vladimir Koltchinskii, Dmitry Panchenko, et al. Empirical margin distributions and bounding the generalization error of combined classifiers. *Annals of statistics*, 30(1):1–50, 2002.
- Wee Sun Lee and Bing Liu. Learning with positive and unlabeled examples using weighted logistic regression. In *ICML*, volume 3, pp. 448–455, 2003.
- Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. *arXiv preprint arXiv:1812.06127*, 2018.
- Bing Liu, Yang Dai, Xiaoli Li, Wee Sun Lee, and Philip S Yu. Building text classifiers using positive and unlabeled examples. In *Third IEEE International Conference on Data Mining*, pp. 179–186. IEEE, 2003.
- Tongliang Liu and Dacheng Tao. Classification with noisy labels by importance reweighting. *IEEE Transactions on pattern analysis and machine intelligence*, 38(3):447–461, 2015.
- Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial Intelligence and Statistics*, pp. 1273–1282. PMLR, 2017.
- Felix Sattler, Simon Wiedemann, Klaus-Robert Müller, and Wojciech Samek. Robust and communication-efficient federated learning from non-iid data. *IEEE transactions on neural networks and learning systems*, 31(9):3400–3413, 2019.
- Kihyuk Sohn, David Berthelot, Chun-Liang Li, Zizhao Zhang, Nicholas Carlini, Ekin D Cubuk, Alex Kurakin, Han Zhang, and Colin Raffel. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. *arXiv preprint arXiv:2001.07685*, 2020.
- Qizhe Xie, Zihang Dai, Eduard Hovy, Minh-Thang Luong, and Quoc V Le. Unsupervised data augmentation for consistency training. *arXiv preprint arXiv:1904.12848*, 2019.
- Yixing Xu, Chang Xu, Chao Xu, and Dacheng Tao. Multi-positive and unlabeled learning. In *IJCAI*, pp. 3182–3188, 2017.
- Felix Yu, Ankit Singh Rawat, Aditya Menon, and Sanjiv Kumar. Federated learning with only positive labels. In *International Conference on Machine Learning*, pp. 10946–10956. PMLR, 2020.
- Zhengming Zhang, Zhewei Yao, Yaoqing Yang, Yujun Yan, Joseph E Gonzalez, and Michael W Mahoney. Benchmarking semi-supervised federated learning. *arXiv preprint arXiv:2008.11364*, 2020.

Table 6: Classification result on iid data.

Num of Clients	Num of P-class	Overlap	Baseline-1	Proposed Method	Baseline-2
10	2	✓	89.59%	90.25%	97.95%
4	6	✓	94.08%	94.22%	98.05%
2	9	✓	94.36%	94.57%	98.20%
10	1	✗	74.07%	89.78%	97.95%
5	2	✗	90.58%	94.09%	98.03%
2	5	✗	94.38%	94.62%	98.20%

Table 7: Classification results with different number of positive classes in each client.

Division of P-class	Overlap	Baseline-1	Proposed Method	Baseline-2
[2,3,4,6,7,8]	✓	89.21%	93.33%	97.91%
[1,2,4,6,7]	✓	90.38%	93.53%	98.03%
[2,4,6,8]	✓	93.20%	94.74%	98.05%
[3,7]	✗	93.56%	94.08%	98.20%
[2,3,5]	✗	89.24%	93.24%	98.16%
[1,2,3,4]	✗	89.18%	93.72%	98.05%

## A RESULTS ON FEDSGD

To further demonstrate the effectiveness of the proposed method, we conduct the proposed method and baseline using FedSGD McMahan et al. (2017). The results are shown in Table 6, 7 and 8, which is consistent with those using FedAvg in the main paper.

## B ABLATION STUDY

In the above sections, we study the PU setting where there are half of data in each positive classes are labeled on the MNIST dataset. Here we make an ablation study to investigate the impact of the percentage of labeled data in positive class. We use 4 clients whose number of positive classes are all equal to 6. The data is collected with iid distributions from each client. As shown in Table 9, with the growth of the percentage of labeled data (from 1/3 to 2/3), the accuracy of the proposed method can be improved from 94.46% to 95.60%, which indicates the effectiveness of the proposed method with different percentage of labeled data.

## C PROOFS

**Theorem 6** Fix  $f \in \mathcal{F}$ , for any  $0 < \delta < 1$ , with probability at least  $1 - \delta$ , the generalization bound holds:

$$\begin{aligned} & \mathbb{E}_i^k \left[ l(f(\mathbf{x}) \neq i) - \sum_{m \notin \mathbf{C}_{\mathbf{P}_k}} l(f(\mathbf{x}) \neq m) \right] - \frac{1}{n_i^k} \sum_{j=1}^{n_i^k} \left[ l(f(\mathbf{x}_j) \neq i) - \sum_{m \notin \mathbf{C}_{\mathbf{P}_k}} l(f(\mathbf{x}_j) \neq m) \right] \\ & \leq 2CV \left( \sum_{s \in \mathbf{C}_{\mathbf{P}_k}} \frac{1}{\sqrt{n_s^k}} + \frac{1}{\sqrt{n_U^k}} \right) + \sqrt{\frac{\log \frac{1}{\delta}}{2n_i^k}}, \end{aligned} \tag{19}$$

where  $i \in \mathbf{C}_{\mathbf{P}_k}$ ,  $V$  is a constant related to the VC-dimension of  $f$  and the bound of loss function  $l$ ,  $n_i^k$  and  $n_U^k$  denotes the number of samples in  $i$ -class and unlabeled classes in  $k$ -th client, respectively.

**Proof 1** According to Koltchinskii et al. (2002), denote  $R(f)$  as the generalization error of hypothesis  $f$ ,  $\hat{R}_{S,\rho}(f)$  as its empirical margin loss with bound  $\rho$ , and  $\mathcal{R}_m(f)$  as Rademacher complexity of the family of loss functions  $f$ , with probability at least  $1 - \delta$ , we have:

$$R(f) \leq \hat{R}_{S,\rho}(f) + \frac{4C}{\rho} \mathcal{R}_n(f) + \sqrt{\frac{\log \frac{1}{\delta}}{n}}, \tag{20}$$

Table 8: Classification result on non-iid data.

Num of Partitions	Division of P-class	Overlap	Baseline-1	Proposed Method	Baseline-2
5	[2,2,...,2]	✓	90.17%	92.54%	97.47%
5	[1,1,...,1]	✗	61.70%	89.69%	97.47%
5	[4,4,3,3,2,2,1,1,1,1]	✓	81.74%	90.48%	97.47%
2	[1,1,...,1]	✗	85.07%	88.62%	96.19%

Table 9: Classification results with different percentage of positive samples.

Percentage	Baseline-1	FedPU	Baseline-2
1/3	91.22%	94.46%	98.05%
1/2	93.24%	95.31%	98.05%
2/3	94.11%	95.60%	98.05%

where  $n$  is the number of training samples and  $C$  is the number of classes.

According to Bousquet et al. (2003), we have:

$$\mathcal{R}_m(f) \leq V' \sqrt{\frac{d}{n}}, \quad (21)$$

where  $d$  is the Vapnik–Chervonenkis (VC) dimension of  $f$ ,  $V'$  is a constant. Taking  $m$  in  $n_i^k$  and  $n_U^k$ , we have:

$$\mathcal{R}_m(f) \leq V' \sqrt{d} \left( \sum_{s \in \mathbf{C}_{P_k}} \frac{1}{\sqrt{n_s^k}} + \frac{1}{\sqrt{n_U^k}} \right). \quad (22)$$

Taking  $V = 4V' \frac{\sqrt{d}}{\rho}$ , we then finish the proof.

**Theorem 7** Fix  $f \in \mathcal{F}$ , for any  $0 < \delta < 1$ , with probability at least  $1 - \delta$ , the generalization bound holds:

$$\mathbb{E}_i^k [l(f(\mathbf{x}) \neq m)] - \frac{1}{n_i^k} \sum_{j=1}^{n_i^k} l(f(\mathbf{x}_j) \neq m) \leq CV \left( \sum_{s \in \mathbf{C}_{P_k}} \frac{1}{\sqrt{n_s^k}} + \frac{1}{\sqrt{n_U^k}} \right) + \sqrt{\frac{\log \frac{1}{\delta}}{2n_i^k}}. \quad (23)$$

The proof of Theorem 2 is the same as that of Theorem 1.

**Lemma 8** Define

$$l'(f(\mathbf{x}) \neq m) = \frac{k^{C_{U_k}}}{k^{C_{U_k}} + \prod_{i \notin \mathbf{C}_{P_k}} |k - i|} l(f(\mathbf{x}) \neq m), \quad (24)$$

where  $C_{U_k}$  denotes the number of unlabeled class in client  $k$ . The decomposition is hold:

$$\begin{aligned} & \sum_{m \notin \mathbf{C}_{P_k}} \mathbb{E}_U^k [l(f(\mathbf{x}) \neq m)] \\ &= \sum_{i \in \mathbf{C}_{P_k}} \pi_i \left( \frac{\prod_{i \notin \mathbf{C}_{P_k}} |k - i|}{k^{C_{U_k}}} \right) \sum_{m \notin \mathbf{C}_{P_k}} \mathbb{E}_i^k [l'(f(\mathbf{x}) \neq m)] + \sum_{m \notin \mathbf{C}_{P_k}} \mathbb{E}_U^k [l'(f(\mathbf{x}) \neq m)]. \end{aligned} \quad (25)$$

**Proof 2** Given

$$l'(f(\mathbf{x}) \neq m) = \frac{k^{C_{U_k}}}{k^{C_{U_k}} + \prod_{i \notin \mathbf{C}_{P_k}} |k - i|} l(f(\mathbf{x}) \neq m), \quad (26)$$

we have:

$$\begin{aligned}
& \mathbb{E}_U^k [l(f(\mathbf{x}) \neq m)] \\
&= \int \sum_y \frac{k^{C_{U_k}} + \prod_{i \notin \mathbf{C}_{P_k}} |k-i|}{k^{C_{U_k}}} l'(f(\mathbf{x}) \neq m) p(\mathbf{x}, y) d\mathbf{x} \\
&= \int l'(f(\mathbf{x}) \neq m) \left[ \sum_{j=1}^K \frac{k^{C_{U_k}} + \prod_{i \notin \mathbf{C}_{P_k}} |k-i|}{k^{C_{U_k}}} p(\mathbf{x}, y = j) \right] d\mathbf{x} \\
&= \int l'(f(\mathbf{x}) \neq m) \sum_{j \in \mathbf{C}_{P_k}} \frac{\prod_{i \notin \mathbf{C}_{P_k}} |k-i|}{k^{C_{U_k}}} p(\mathbf{x}, y = j) d\mathbf{x} \\
&\quad + \int l'(f(\mathbf{x}) \neq m) \sum_{j=1}^K p(\mathbf{x}, y = j) d\mathbf{x} \\
&= \sum_{i \in \mathbf{C}_{P_k}} \pi_i \left( \frac{\prod_{i \notin \mathbf{C}_{P_k}} |k-i|}{k^{C_{U_k}}} \right) \mathbb{E}_i^k [l'(f(\mathbf{x}) \neq m)] \\
&\quad + \mathbb{E}_U^k [l'(f(\mathbf{x}) \neq m)].
\end{aligned} \tag{27}$$

**Theorem 9** Fix  $f \in \mathcal{F}$ , for any  $0 < \delta < 1$ , with probability at least  $1 - \delta$ , the generalization bound holds:

$$\begin{aligned}
& \sum_{m \notin \mathbf{C}_{P_k}} \mathbb{E}_U^k [l(f(\mathbf{x}) \neq m)] - \frac{1}{n_U^k} \sum_{j=1}^{n_U^k} \sum_{m \notin \mathbf{C}_{P_k}} l'(f(\mathbf{x}_j) \neq m) \\
&\leq \sum_{i \in \mathbf{C}_{P_k}} \frac{\pi_i}{n_i^k} \left( 1 + \frac{\prod_{i \notin \mathbf{C}_{P_k}} |k-i|}{k^{C_{U_k}}} \right) \sum_{j=1}^{n_i^k} \sum_{m \notin \mathbf{C}_{P_k}} l'(f(\mathbf{x}_j) \neq m) \\
&\quad + \left( \sum_{i \in \mathbf{C}_{P_k}} \pi_i + 1 \right) CV \left( \sum_{s \in \mathbf{C}_{P_k}} \frac{1}{\sqrt{n_s^k}} + \frac{1}{\sqrt{n_U^k}} \right) \\
&\quad + \sum_{i \in \mathbf{C}_{P_k}} \pi_i \left( 1 + \frac{\prod_{i \notin \mathbf{C}_{P_k}} |k-i|}{k^{C_{U_k}}} \right) \sqrt{\frac{\log \frac{1}{\delta}}{2n_i^k}} + \sqrt{\frac{\log \frac{1}{\delta}}{2n_U^k}}.
\end{aligned} \tag{28}$$

With the evidence of Lemma 3, the proof of Theorem 4 is the same as that of Theorem 1.

**Theorem 10** As  $n_i^k, n_U^k \rightarrow \infty$ ,  $i \in \mathbf{C}_{P_k}$ ,  $k \in \{1, \dots, K\}$ , the generalization bound of the proposed FedPU is of order:

$$\mathcal{O} \left( \sum_{k=1}^K C^2 \left( \sum_{i \in \mathbf{C}_{P_k}} \frac{1}{\sqrt{n_i^k}} + \frac{1}{\sqrt{n_U^k}} \right) \right). \tag{29}$$

By concluding the result in Theorem 1, 2 and 4. We can derive that the generalization bound in  $k$ -th client as  $\mathcal{O} \left( C^2 \left( \sum_{i \in \mathbf{C}_{P_k}} \frac{1}{\sqrt{n_i^k}} + \frac{1}{\sqrt{n_U^k}} \right) \right)$ . By summing the bound in each client, we then finish the proof.