

# Node Selection Toward Faster Convergence for Federated Learning on Non-IID Data

Hongda Wu, *Student Member, IEEE* and Ping Wang, *Fellow, IEEE*

**Abstract**—Federated Learning (FL) is a distributed learning paradigm that enables a large number of resource-limited nodes to collaboratively train a model without data sharing. The non-independent-and-identically-distributed (non-i.i.d.) data samples invoke discrepancies between the global and local objectives, making the FL model slow to converge. In this paper, we proposed Optimal Aggregation algorithm for better aggregation, which finds out the optimal subset of local updates of participating nodes in each global round, by identifying and excluding the adverse local updates via checking the relationship between the local gradient and the global gradient. Then, we proposed a Probabilistic Node Selection framework (FedPNS) to dynamically change the probability for each node to be selected based on the output of Optimal Aggregation. FedPNS can preferentially select nodes that propel faster model convergence. The convergence rate improvement of FedPNS over the commonly adopted Federated Averaging (FedAvg) algorithm is analyzed theoretically. Experimental results demonstrate the effectiveness of FedPNS in accelerating the FL convergence rate, as compared to FedAvg with random node selection.

**Index Terms**—Federated Learning, Mobile Edge Computing, Fast Convergence, Node Selection.

## I. INTRODUCTION

WITH the rapid growth of computational capability at mobile edge sides, next-generation computing network is experiencing a paradigm shift from traditional cloud computing to Mobile Edge Computing (MEC) systems [1] [2]. With the deployed computational power and storage capability, MEC systems construct the node-edge-cloud architecture in supporting the applications at resource-constrained nodes that require low latency communication (e.g., autonomous driving) or high throughput (e.g., content delivery network) [3]. Edge nodes such as sensors, mobile devices, and connected vehicles are generating an unprecedented amount of data consistently and coupled with cutting-edge Machine Learning (ML) / Deep Learning (DL) techniques, the MEC system is able to conduct intelligent inference (e.g., road congestion prediction [4]) and perceptive control (e.g., unmanned aerial vehicles (UAVs) swarm navigation [5]).

In traditional ML fashion, in order to train a complex DL model with millions of model parameters, a tremendous amount of data aggregated from multiple edge nodes is typically needed, which is offloaded via wireless network to edge server. However, collecting data for model training is unrealistic from privacy, security, regulatory, or necessity perspectives.

Hongda Wu and Ping Wang are with the Department of Electrical Engineering and Computer Science, Lassonde School of Engineering, York University, Toronto, ON M3J 1P3, Canada (e-mail: hwu1226@eecs.yorku.ca; pingw@yorku.ca)

With the ever-increasing computational capability on edge nodes, it becomes more attractive to perform model training on the edge node side instead of sending raw data to the edge server. To this end, Federated Learning (FL) has emerged as a variant of the previous Distributed ML (DML) manner, which decouples the data acquisition and model training at the edge server [6], [7]. In general, FL systems aim to optimize a global model under the orchestration of an edge server, which allows the collaboration of multiple edge nodes for data augmentation while keeping training data locally. FL involves several communication rounds, each of which includes local model training, model update transmission, and global model aggregation. Along the iterative process, the edge server is able to train a statistical model that is suitable for all participating nodes without accessing user-sensitive data. The improved data confidentiality and reduced volume of communication cost make FL one of the most promising technologies for future network intelligence [8].

Nonetheless, a fundamental challenge for FL in comparison with the optimization in DML, where algorithms run on independent and identically distributed (i.i.d) data samples partitioned from a large dataset, is the data heterogeneity [9]–[11]. To be more specific, the model is updated via feature learning on local data samples, which are user-specific and reveal a different pattern. The data samples across participating nodes may not be independent and identically distributed (non-i.i.d.). Since participating nodes in each iteration are selected randomly, data distribution on nodes cannot represent the global data distribution. Training on nodes with non-i.i.d. datasets will lead to the biased model update, which stagnates model convergence and reduces the model accuracy substantially, and consequently invokes additional communication rounds to resource-constrained edge nodes [10], [11]. Though a relatively small amount of data is sent (i.e., in general, model parameters have a smaller size than the raw training data), communication time in FL is proven to be the critical bottleneck for FL due to the network uncertainty, bandwidth limitation, and straggler effect, etc. [12].

In this paper, we design a node selection<sup>1</sup> scheme to improve the convergence rate of FL with non-i.i.d nodes, called FedPNS, which is a Probabilistic Node Selection framework with contribution-related criteria. We find out global model aggregation over all participating nodes is not of necessity, whereas excluding some adverse local updates may lead to a better global model in terms of model accuracy. In or-

<sup>1</sup>A critical property that differentiates FL from a typical distributed optimization problem is the massively distributed nodes [6]. Therefore, in each round, a small fraction of nodes is selected for participation.

der to improve the expected decrement of FL loss in each round, we propose an Optimal Aggregation algorithm to determine the optimal subset of local updates (from the participating nodes) for global model aggregation, which utilizes the inner product between local gradients and the global gradient<sup>2</sup> as the indicator. By applying the result from Optimal Aggregation, the data heterogeneity can be profiled, which is used to adjust the probability for each node to be selected in the subsequent global rounds. Consequently, the server can preferentially select nodes that propel faster model convergence. Note that our probabilistic node selection is conducted on the server-side, which does not impose additional communication costs. Our main contributions in this paper are as follows

- We analyze the convergence bound of the commonly adopted Federated Averaging (FedAvg) algorithm [6] from a theoretical perspective and derive the expected decrease of FL global loss, considering the data heterogeneity and the way to aggregate local updates.
- We challenge the necessity of global model aggregation over local updates of all participating nodes and propose Optimal Aggregation to identify and exclude the potential adverse local updates, which enlarges the expected decrease of global loss in each round.
- We design FedPNS, a Probabilistic Node Selection scheme that enables server to dynamically adjust the probability for each node to be selected in each round, based on the result of Optimal Aggregation. FedPNS tendentially selects nodes that boost model convergence. The convergence rate improvement of FedPNS over FedAvg is illustrated theoretically and the imposed computational complexity of FedPNS is discussed.
- We empirically evaluate the performance of FedPNS via extensive experiments using the synthetic dataset and real datasets with different learning objectives. The experimental results show the effectiveness of FedPNS in improving the convergence rate of the FL model compared with the commonly adopted FedAvg algorithm.

## II. RELATED WORK

Some existing works on FL focus on the communication cost reduction, with the aim of directly reducing communication cost on the wireless link, where typical methods range from important-based updating [13], [14], model quantization [15], and analog aggregation [16]. In particular, Wang *et al.* [13] proposed identifying the irrelevant update at the node side caused by different data distribution. Communication cost is reduced by excluding irrelevant updates from local nodes. Similarly, authors in [14] introduced a concept, namely important gradient, where communication reduction is achieved by sending the gradient with a larger magnitude. Different from [13], [14], Seide *et al.* [15] proposed 1-bit stochastic gradient descent to reduce model transfer data size and achieved 10× speed up in speech applications. Zhu *et al.* [16] proposed to utilize analog aggregation rather than

digital aggregation. By exploiting the waveform-superposition property of a multi-access channel, model transmission and aggregation are realized over wireless links simultaneously.

Another series of studies concentrate on the algorithmic perspective via handling the inherent non-i.i.d. data distribution across participating nodes, aiming to reduce communication rounds in FL. These studies include adaptive tuning local training [9], weighting design for model aggregation [17], and node selection strategies [18]–[25]. The algorithm FedProx by Li *et al.* [9] uses a regularization term to balance the optimizing discrepancy between the global and local objectives, and allowing participating nodes to perform a variable number of local updates, to consequently overcome the non-i.i.d. data distribution and resource heterogeneity. Authors in [17] exhibited a contribution-related weighting design, namely FedAdp to boost FL convergence rate in the presence of nodes with non-i.i.d. data samples, which assigns distinguished weights for participating nodes according to the correlation between local objective and global objective revealed by gradient information.

In general, to avoid long-tail waiting time in synchronous aggregation protocol, FL algorithm randomly selects a subset of nodes (i.e., partial node participation) in each round to participate in local training (e.g., FedAvg [6], FedProx [9], CMFL [13]). Compared with DML, candidate nodes in FL are more heterogeneous with regard to computation/communication capability, wireless connection, and data quality. Therefore, a carefully designed node selection is beneficial for performance improvement. Several works are carried out focusing on node selection design to improve the FL convergence rate, taking the system heterogeneity and uncertainty of wireless medium into consideration [18], [19], [21]–[23]. Specifically, Nishio *et al.* [18] proposed to select nodes intentionally based on the resource condition on nodes. Amiria *et al.* [19] designed a node scheduling algorithm by considering the significance of local update measured by  $\ell_2$ -norm and channel condition separately or jointly. For example, in BN2 algorithm [19], the server first selects a macro set of nodes to participate in local training. Then a subset of the macro set is finally chosen for model aggregation by ordering the norm of gradient transmitted from nodes of the macro set. In [20], the authors proposed biased client selection strategies, that is, preferentially choosing the node with higher local loss. Though the contribution-related loss measurement leads to a faster convergence, the selection skewness imposes potential error, and the local loss measurement results in additional communication and computation cost. Differently, references [21]–[25] focus on probabilistic node selection strategy where each node is eligible to contribute to the global model. In particular, Chen *et al.* [21] considered the limited bandwidth resource for model transmission where node selection for global model aggregation is of importance. The proposed method measures the node contribution according to the norm of local updates, by which the probability for each node to be selected is calculated so as to execute the node selection procedure. The nodes with higher norm of local updates are chosen with higher probability, thus boosting the convergence rate when limited bandwidth resource is provided. Along with [21],

<sup>2</sup>We use local/global gradient and local/global update interchangeably.

authors in [22] proposed to use Artificial Neural Networks (ANNs) as a predictor to estimate the model updates of nodes that are not allocated the bandwidth for transmission, based on the model updates that are successfully transmitted using limited bandwidth resource. The additionally included model updates further accelerate the model convergence. Authors in [23] proposed a probabilistic design by considering the importance of local update and transmission latency, where the importance of local update is evaluated by gradient divergence between local gradients and the ground truth global gradient. The probability for node selection is finally determined by the local gradient norm and transmission latency. Chen *et al.* [24] designed an importance sampling scheme that selects more informative nodes. The node sampling procedure minimizes the variance of local gradients for aggregation, while the probability for each node to be chosen is proportional to the norm of local updates. In addition, authors in [25] applied importance sampling for node selection on the server level and data selection on the node level. Similar to [24], the optimal node selection is achieved by minimizing the bound on the variance of gradient noise, i.e., the estimation error of the global gradient because of the partial node participation. The probability for each node to be chosen is proportional to the norm of its local updates.

None of the aforementioned node selection designs analyzed the impact of data heterogeneity on node selection. Given the heterogeneous training samples across nodes, the magnitude of local gradient norm is deficient in reflecting the contribution from each of those nodes, which is empirically shown in Section V-D, since local gradients may not align with the global gradient. In contrast to the above research, our work in this paper builds on the data heterogeneity perspective and designs a probabilistic model to choose participating nodes. The proposed method scrutinizes the relationship between local gradients and the global gradient so as to adjust the probability for each node to be selected, which is different from the criteria (i.e., the norm of local gradient/update) adopted in [19], [21]–[25]. Since FedAvg algorithm may struggle to converge on non-i.i.d. data, it is not trivial to profile the distribution of data samples across nodes. Upon identifying the contribution difference of nodes, it is profitable to accelerate model convergence by choosing nodes tendentiously, as compared with random selection.

The rest of the paper is organized as follows. Section III provides the preliminary and implementation of federated learning and the challenge of non-i.i.d. data on FL. In Section IV, the convergence analysis, the proposed aggregation scheme and probabilistic node selection, and complexity analysis are presented. Experimental results are shown in Section V, and the conclusion is presented in Section VI.

### III. PRELIMINARIES

In this section, we first introduce the key ingredients behind federated learning, including the system model (Section III-A) and the practical algorithm design to solve federated learning problem (Section III-B). Then, the challenge of FL on heterogeneous data is analyzed (Section III-C).

#### A. Federated Learning Model

In general, federated learning methods [6], [9], [26], are designed to handle the consensus learning task in a decentralized manner, where a central server coordinates the global learning objective and multiple devices train the model with locally collected data. Consider a network with  $\mathcal{K}$  local nodes (i.e.,  $i \in \{1, 2, \dots, |\mathcal{K}|\}$ ), where each node  $i$  possesses a local (private) dataset  $\mathcal{D}_i$  with size  $D_i$ . The nodes are connected with a central server and seek to collaboratively find a global model parameterized by  $\mathbf{w}$  that minimizes the empirical risk

$$F(\mathbf{w}) = \frac{1}{\sum_{i=1}^{|\mathcal{K}|} D_i} \sum_{i=1}^{|\mathcal{K}|} \sum_{\{\mathbf{x}, y\} \in \mathcal{D}_i} f(\mathbf{w}, \mathbf{x}, y), \quad (1)$$

where  $f(\mathbf{w}, \mathbf{x}, y)$  is the composite loss for training sample  $\{\mathbf{x}, y\}$ . Specifically, in the context of  $C$ -class classification problem hereinafter, each training sample  $\{\mathbf{x}, y\} \in \mathcal{D}_i$  is assumed to contain a feature vector  $\mathbf{x}$  and label  $y$  over feature space  $\mathbb{X}$  and label space  $\mathbb{Y}$  (i.e.,  $\mathbb{Y} = [C]$ , where  $[C] = \{1, \dots, C\}$ ). For each available training sample  $\{\mathbf{x}, y\} \in \bigcup_i \mathcal{D}_i$  in FL problem, the federated learning model parameterized by  $\mathbf{w}$  is considered to learn the predicted probability vector  $\bar{\mathbf{y}}$ , i.e.,  $\bar{\mathbf{y}} \sum_{j=1}^C \bar{y}_j = 1, \bar{y}_j \geq 0, \forall j \in [C]$ , with empirical risk.

From a federation perspective, the global objective  $F(\mathbf{w})$  in (1) is surrogated by local objectives  $F_i(\mathbf{w})$  and can be further represented as follows

$$F(\mathbf{w}) = \sum_{i=1}^{|\mathcal{K}|} \frac{D_i}{\sum_{i=1}^{|\mathcal{K}|} D_i} F_i(\mathbf{w}), \quad (2)$$

For node  $i \in \mathcal{K}$ ,  $F_i(\mathbf{w})$  commonly measures the local empirical risk (e.g., cross entropy loss) over the dataset  $\mathcal{D}_i$  with possibly differing data distribution  $q^{(i)}$ , which is defined as follows

$$\begin{aligned} F_i(\mathbf{w}) &= \mathbb{E}_{\mathbf{x}, y \sim q^{(i)}} \left[ - \sum_{j=1}^C \mathbb{1}_{y=j} \log l_j(\mathbf{w}, \mathbf{x}, y) \right] \\ &= - \sum_{j=1}^C q^{(i)}(y=j) \mathbb{E}_{\mathbf{x}|y=j} [\log l_j(\mathbf{w}, \mathbf{x}, y)], \end{aligned} \quad (3)$$

where  $l_j(\mathbf{w}, \mathbf{x}, y)$  denotes the probability that the data sample  $\{\mathbf{x}, y\}$  is classified as the  $j$ -th class given model  $\mathbf{w}$ .  $q^{(i)}(y=j)$  denotes the data distribution on node  $i$  over class  $j \in [C]$ .

#### B. FedAvg with Partial Node Participation

The most commonly used algorithm to solve (2) is Federated Averaging (FedAvg) [6], [26], where the training consists of multiple communication rounds. At each communication round  $t$ , the server selects a fraction  $c$  of nodes  $|\mathcal{S}_t| = c|\mathcal{K}|$  to participate in the training. Taking the global model  $\mathbf{w}^{t-1}$  in previous round as the reference, each participating node  $i \in \mathcal{S}_t$  performs local Stochastic Gradient Descent (SGD) to optimize its objective

$$\mathbf{w}_i^t = \mathbf{w}^{t-1} - \eta \nabla F_i(\mathbf{w}^{t-1}), \quad (4)$$

where  $\eta$  is the learning rate and  $\nabla F_i(\cdot)$  is the gradient<sup>3</sup> at node  $i$ . (4) gives a general principle of SGD optimization, where  $\mathbf{w}_i^t$  is the result after  $\tau$  local updates of mini-batch SGD (i.e.,  $\tau = \frac{D_i}{B}E$ , where  $E$  is the number of local training epochs,  $B$  is the batch size of mini-batch training samples).

The participating nodes then communicate their model update  $\Delta_i^t = \mathbf{w}_i^t - \mathbf{w}^{t-1}$  back to the server, which aggregates them and updates the global model<sup>4</sup> as follows

$$\begin{aligned} \Delta^t &= \frac{1}{|\mathcal{S}_t|} \sum_{i \in \mathcal{S}_t} \Delta_i^t \\ \mathbf{w}^t &= \mathbf{w}^{t-1} + \Delta^t. \end{aligned} \quad (5)$$

### C. The Challenges of Non-i.i.d. Data Distribution

Though FedAvg can achieve a decent convergence rate with random node selection policy and simple averaging design, partial node participation and non-i.i.d. training data slows the convergence rate [26], which is also observed in [10]–[12], [17], [20]. Since communication cost becomes a critical bottleneck in FL, one can increase the local computing (i.e., more local updates), which is shown to be beneficial to save communication rounds and improve the convergence rate [6], [12].

However, model performance on non-i.i.d. dataset is not satisfactory, even with increased local computing [26], [27]. This is because the local objective  $F_i(\mathbf{w})$ , which the local optimizer minimizes, is closely related to data distribution  $q^{(i)}$ . In trivial node selection policy (e.g., random selection in FedAvg), the distribution of data samples on selected nodes differs from each other. Local updates lead the model towards optima to its local objective, which is deviated from the global objective in a non-i.i.d. setting, causing training instability that makes the FL model struggle to converge.

It is crucial to understand and analyze the non-trivial node selection policy from the data heterogeneity perspective, identifying and choosing the nodes that contribute better to model convergence. By taking the inner product between the local update and global update as the criterion, which implicitly profiles the difference between data distribution on nodes and population distribution, we first identify the nodes whose updates adversely contribute to the global update. By excluding the potential adverse local updates and reducing the probability for those nodes to be selected, one can ensure that the node with a higher contribution to the decrease of global loss enjoys a higher probability of being chosen. Consequently, the non-trivial node selection accelerates model convergence compared with FedAvg.

<sup>3</sup>Through this paper, the gradient refers to the stochastic version instead of the actual gradient calculated from the entire dataset.

<sup>4</sup>It is worth to mention that the aggregation scheme is applied over all nodes in vanilla FedAvg [6], i.e.,  $\Delta^t = \sum_{i \in \mathcal{S}_t} \psi_i \Delta_i^t + \sum_{i \in \mathcal{K} - \mathcal{S}_t} \psi_i \mathbf{w}^{t-1}$ , where  $\psi_i = \frac{D_i}{\sum_{i=1}^K D_i}$ . The subsequent work [9] proposed a variant of aggregation over participating nodes as in (5). Hereinafter, FedAvg denotes the algorithm that involves random selection and partial aggregation of nodes with equal data size [9].

## IV. CONTRIBUTION-BASED NODE SELECTION

In this section, we design a probabilistic node selection scheme to improve the convergence rate of federated learning. For FL with the heterogeneous dataset, we analyze the convergence property of FedAvg theoretically (Section IV-A). In Section IV-B, we challenge the necessity of global model aggregation over all participating nodes. Then, the Optimal Aggregation algorithm is proposed, which can identify and exclude the adverse local updates to make greater progress on reducing the expected decrement of global loss in each round. The FL with Probabilistic Node Selection (FedPNS) is proposed based on the result of Optimal Aggregation. FedPNS adjusts the probability for each node to be selected, and the server is able to preferentially select nodes that propel a faster model convergence (Section IV-C). The convergence rate improvement of FedPNS over FedAvg is analyzed theoretically (Section IV-D) and the computation complexity of FedPNS is discussed in Section IV-E.

### A. Convergence Analysis

For theoretical analysis purposes, we employ the following assumptions to the loss function, which have also been commonly made in the literature [9], [26]–[28].

#### Assumption 1. Convex, $\zeta$ -Lipschitz, and $L$ -smooth.

$F_i(\mathbf{w})$  is convex,  $\zeta$ -Lipschitz, and  $L$ -smooth for all node  $i$ , i.e.,  $\|F_i(\mathbf{w}) - F_i(\mathbf{w}')\| \leq \zeta \|\mathbf{w} - \mathbf{w}'\|$ ,  $\|\nabla F_i(\mathbf{w}) - \nabla F_i(\mathbf{w}')\| \leq L \|\mathbf{w} - \mathbf{w}'\|$ , for any  $\mathbf{w}, \mathbf{w}'$ .

Based on Assumption 1, the definition of  $F(\mathbf{w})$ , and triangle inequality, we can easily get that  $F(\mathbf{w})$  is convex,  $\zeta$ -Lipschitz, and  $L$ -smooth.

#### Assumption 2. $\delta$ -local dissimilarity.

Local loss functions  $F_i(\mathbf{w}^t)$  are  $\delta$ -local dissimilar at  $\mathbf{w}^t$ , i.e.,  $\mathbb{E}_{i \sim \mathcal{S}_t} [\|\nabla F_i(\mathbf{w}^t)\|^2] \leq \|\nabla F(\mathbf{w}^t)\|^2 \delta^2$  for  $i \in \mathcal{S}_t$  and  $t = 1, \dots, T$ , where  $T$  is the number of global rounds.  $\mathbb{E}_{i \sim \mathcal{S}_t}[\cdot]$  denotes the expectation over participating nodes  $\mathcal{S}_t$  with weight  $\frac{1}{|\mathcal{S}_t|}$  (as in (5)).  $\nabla F(\mathbf{w}^t)$  is the global gradient at the  $t$ -th global round defined as  $\nabla F(\mathbf{w}^t) = \frac{1}{|\mathcal{S}_t|} \sum_{i \in \mathcal{S}_t} \nabla F_i(\mathbf{w}^t)$ .

#### Assumption 3. Bounded gradient.

The norm of gradient in each node is bounded, i.e.,  $\|\nabla F_i(\mathbf{w}^t)\| \leq \gamma_i$  for all  $i \in \mathcal{K}$  and  $t = 1, \dots, T$ .

Assumption 1 is standard, which can be satisfied when the logistic regression with cross entropy loss is adopted<sup>5</sup>. The discrepancy between the local objective and global objective caused by the data heterogeneity is captured by Assumption 2, which has been made in previous work [26], [28]. As the data distribution across participating nodes becomes more heterogeneous, the local updates (i.e., gradient) will diverge from each other, and  $\delta$  will increase. On the other hand, if the data samples on participating nodes follow the same data distribution, the local gradients become more similar and  $\delta$  goes to 1. Assumption 3 has been made in different forms by previous works [26], [27]. Besides, with  $\mathbf{w}$  trained by

<sup>5</sup>More examples include  $\ell_2$ -norm regularized linear regression with mean square error, and the support vector machine with hinge loss.

heterogeneous data,  $\gamma$  is different for different nodes, which is closely related to the data distribution on each node. If the data distribution on node  $i$  is more similar to the population distribution over all nodes,  $\gamma_i$  is lower, and vice versa. This observation is empirically illustrated in Section V-D.

**Lemma 1.** *Let assumptions 1 and 2 hold. Suppose that  $\mathbf{w}^t$  is not a stationary solution, the expected decrement on the global loss of FedAvg between two consecutive rounds satisfies*

$$F(\mathbf{w}^{t+1}) \leq F(\mathbf{w}^t) - \eta \mathbb{E}_{i \sim \mathcal{S}_t} [\langle \nabla F(\mathbf{w}^t), \nabla F_i(\mathbf{w}^t) \rangle] + \frac{L\eta^2}{2} \|\nabla F(\mathbf{w}^t)\|^2 \delta^2, \quad (6)$$

where  $\eta$  is the learning rate of SGD,  $\langle \cdot \rangle$  is the inner product operation, and  $\|\cdot\|$  denotes the  $\ell_2$ -norm of a vector.

The proof of Lemma 1 is presented in Appendix-A. Lemma 1 provides a bound on how rapid the decrease of the global FL loss can be expected. The decrease of global FL loss between two consecutive rounds shows a dependency on  $\delta$ , which represents the variance between local data distributions, and the aggregation strategy  $\mathbb{E}_{i \sim \mathcal{S}_t}[\cdot]$ , where  $\nabla F(\mathbf{w}^t)$  is obtained by aggregating over local updates from all participating nodes, i.e.,  $\nabla F_i(\mathbf{w}^t), i \in \mathcal{S}_t$  with weight  $1/|\mathcal{S}_t|$ .

### B. Aggregation with Gradient Information

In vanilla FedAvg [6] and the subsequent work [9]–[12], the averaging technique is used for global update aggregation due to its simplicity. One can challenge the inherent rule that the global update is aggregated over local updates of all participating nodes since the local updates may contribute global model in an adverse way. As a sanity check, at any communication round  $t$ , the local updates from the participating nodes whose inner product between their gradients and the global gradient is negative i.e.,  $\langle \nabla F(\mathbf{w}^t), \nabla F_i(\mathbf{w}^t) \rangle < 0$ , will slow the model convergence because of the reduced expected loss decrement (i.e., a lower expectation value as in (6)) in this round. As such, it is not trivial to exclude the adverse local updates, which is realized by examining the value of expectation term in Lemma 1, as illustrated later. Excluding adverse local updates gives an impact on the reduction of overall data heterogeneity, which, in the meanwhile, changes the relationship between local gradients and the global gradient  $\langle \nabla \bar{F}(\mathbf{w}^t), \nabla F_i(\mathbf{w}^t) \rangle$ , where  $\nabla \bar{F}(\mathbf{w}^t) = \frac{1}{|\mathcal{S}_t^*|} \sum_{i \in \mathcal{S}_t^*} \nabla F_i(\mathbf{w}^t)$  is defined over  $\mathcal{S}_t^*$ , i.e., the subset of participating nodes  $\mathcal{S}_t$  after successfully excluding the nodes with adverse local updates.

To find the optimal subset of local updates to aggregate, we first check the expectation term  $\mathbb{E}_{i \sim \mathcal{S}_t} [\langle \nabla F(\mathbf{w}^t), \nabla F_i(\mathbf{w}^t) \rangle]$  in Lemma 1 and exclude the local updates from participating nodes  $k$ , i.e.,  $k \in \mathcal{S}_t - \bar{\mathcal{S}}_t$  if  $\mathbb{E}_{i \sim \bar{\mathcal{S}}_t} [\langle \nabla \bar{F}(\mathbf{w}^t), \nabla F_i(\mathbf{w}^t) \rangle] > \mathbb{E}_{i \sim \mathcal{S}_t} [\langle \nabla F(\mathbf{w}^t), \nabla F_i(\mathbf{w}^t) \rangle]$  is satisfied. However, excluding local updates gives an impact on the global update and overall data heterogeneity, i.e.,  $\|\nabla F(\mathbf{w}^t)\|^2 \delta^2$ , the last term on the right hand side of (6), which makes the expected decrement of global loss, i.e.,  $\Delta F(\mathbf{w}^t) = \frac{L\eta^2}{2} \|\nabla F(\mathbf{w}^t)\|^2 \delta^2 - \eta \mathbb{E}_{i \sim \mathcal{S}_t} [\langle \nabla F(\mathbf{w}^t), \nabla F_i(\mathbf{w}^t) \rangle]$ , difficult to be analyzed quantitatively given  $L$  and  $\delta$ . Therefore, in the second step, test loss is adopted to ensure that excluding local updates makes

---

### Algorithm 1 Optimal Local Updates for Aggregation

---

#### Procedure OPTIMAL AGGREGATION

**Input:**  $\mathcal{S}_t, \Delta_t^t, v, \text{temp} = \{\}$

```

1:  $\nabla F(\mathbf{w}^t) = -\Delta_t^t / \eta$ 
2:  $\nabla F(\mathbf{w}^t) = \frac{1}{|\mathcal{S}_t|} \sum_{i \in \mathcal{S}_t} \nabla F_i(\mathbf{w}^t)$ 
3:  $\max = \mathbb{E}_{i \sim \mathcal{S}_t} [\langle \nabla F(\mathbf{w}^t), \nabla F_i(\mathbf{w}^t) \rangle]$ 
4: while  $|\mathcal{S}_t| \geq v$  do
5:    $\text{temp} \leftarrow \text{CHECK EXPECTATION}(\nabla F_i(\mathbf{w}^t), \mathcal{S}_t, \text{temp})$ 
6:   if  $\max(\text{temp}).\text{value} < \max$  do
7:     break with  $\mathcal{S}_t^* = \mathcal{S}_t$ 
8:   else
9:      $\text{key} = \max(\text{temp}).\text{key}$ 
10:     $\text{ls}(\mathbf{w}), \text{ls}(\bar{\mathbf{w}}), \bar{\mathcal{S}}_t \leftarrow \text{CHECK LOSS}(\nabla F_i(\mathbf{w}^t), \mathcal{S}_t, \text{key})$ 
11:    if  $\text{ls}(\mathbf{w}) > \text{ls}(\bar{\mathbf{w}})$  do
12:      break with  $\bar{\mathcal{S}}_t, \mathcal{S}_t^* = \mathcal{S}_t$ 
13:    else
14:       $\mathcal{S}_t, \mathcal{S}_t^* \leftarrow \mathcal{S}_t.\text{pop}(\text{key})$ 
15:       $\max \leftarrow \text{temp}(\text{key}).\text{value}$ 
16: return  $\mathcal{S}_t^*, \bar{\mathcal{S}}_t$ 
17:  $\mathbf{w}^{t+1} \leftarrow \text{GLOBAL UPDATE}(\nabla F_i(\mathbf{w}^t), \mathcal{S}_t^*)$ 
```

#### Procedure CHECK EXPECTATION

**Input:**  $\nabla F_i(\mathbf{w}^t), \mathcal{S}_t, \text{temp}$

```

18: for  $i = 1, \dots, |\mathcal{S}_t|$  do
19:    $\bar{\mathcal{S}}_t \leftarrow \mathcal{S}_t.\text{pop}(\mathcal{S}_t[i])$ 
20:    $\nabla \bar{F}(\mathbf{w}^t) = \frac{1}{|\bar{\mathcal{S}}_t|} \sum_{i \in \bar{\mathcal{S}}_t} \nabla F_i(\mathbf{w}^t)$ 
21:    $\text{temp}(\mathcal{S}_t[i]) = \mathbb{E}_{i \sim \bar{\mathcal{S}}_t} [\langle \nabla \bar{F}(\mathbf{w}^t), \nabla F_i(\mathbf{w}^t) \rangle]$ 
```

#### Procedure CHECK LOSS

**Input:**  $\nabla F_i(\mathbf{w}^t), \mathcal{S}_t, \text{key}$

```

22:  $\bar{\mathcal{S}}_t \leftarrow \mathcal{S}_t.\text{pop}(\text{key})$ 
23: Generate global model  $\mathbf{w}^{t+1}$  by  $\nabla F_i(\mathbf{w}^t), i \in \mathcal{S}_t$  and  $\bar{\mathbf{w}}^{t+1}$ 
    by  $\nabla F_i(\mathbf{w}^t), i \in \bar{\mathcal{S}}_t$ , respectively
24: Evaluate  $\mathbf{w}^{t+1}, \bar{\mathbf{w}}^{t+1}$  by using batch samples (with size  $\bar{B}$ )
    from  $\mathcal{D}_{\text{test}}$  and get the loss  $\text{ls}(\mathbf{w})$  and  $\text{ls}(\bar{\mathbf{w}})$ , respectively
25: return  $\text{ls}(\mathbf{w}), \text{ls}(\bar{\mathbf{w}}), \bar{\mathcal{S}}_t$ 
```

#### Procedure GLOBAL UPDATE

**Input:**  $\nabla F_i(\mathbf{w}^t), \mathcal{S}_t^*$

```

26: Generate  $\mathbf{w}^{t+1}$  by  $\nabla F_i(\mathbf{w}^t), i \in \mathcal{S}_t^*$  via (4) and (5)
27: return  $\mathbf{w}^{t+1}$ 
```

---

global update better in terms of model convergence, as in [20]. In particular, the global model  $\mathbf{w}^{t+1}$  and  $\bar{\mathbf{w}}^{t+1}$  generated by  $\nabla F_i(\mathbf{w}^t), i \in \mathcal{S}_t$  and  $\nabla F_i(\mathbf{w}^t), i \in \bar{\mathcal{S}}_t$ , respectively, are evaluated using mini-batch of samples with size  $\bar{B}$  that are sampled uniformly at random from  $\mathcal{D}_{\text{test}}$  (e.g., test dataset in MNIST).

An iterative algorithm called Optimal Aggregation is proposed for a better local update aggregation in each round, which finds the *optimal* subset of local updates  $\Delta_t, i \in \mathcal{S}_t^* \subseteq \mathcal{S}_t$  by excluding the adverse local updates  $\Delta_k, k \in \mathcal{S}_t - \mathcal{S}_t^*$ , as in Algorithms 1. Specifically, for a given set of participating nodes  $\mathcal{S}_t$  in each global round  $t$ , the server iteratively removes one of the local updates  $\nabla F_i(\mathbf{w}^t), i \in \mathcal{S}_t$ , generates the potential global gradient, and calculates the expectation term in (6) (i.e., CHECK EXPECTATION, line 18-21). If excluding one local update gives a higher expectation value, compared with the case that includes all local updates retained in  $\mathcal{S}_t$ ,

that local update will be labeled, and loss comparison will be performed to check the loss criterion (CHECK LOSS, line 22-25), otherwise the server keeps all local updates (line 6). If the loss criterion is satisfied (line 13), the labeled local update is eventually removed from set  $\mathcal{S}_t$  (line 14). Otherwise, the server keeps that local update retained in  $\mathcal{S}_t$  (line 12). The process repeats until no adverse local update can be found or the number of remaining local updates is below a threshold  $v$  (line 4). In Algorithm 1, the function `pop` is defined as removing element (line 14). The introduced “temp” is a dictionary with key-value pairs (line 5) and the function `max` returns the maximum value (line 6) or the key (i.e., the node index  $i$ ) corresponding to that value (line 9), respectively.

Given a set of participating nodes  $\mathcal{S}_t$ , the benefits of finding optimal local updates are twofold: (i) Excluding the potential local updates that contribute to the global model adversely results in a larger decrement of the expected loss in each round. (ii) By CHECK EXPECTATION, the potential adverse nodes  $k, k \in \mathcal{S}_t - \bar{\mathcal{S}}_t$  (nodes with non-i.i.d. dataset normally) are identified. This identification can be used for consequent probabilistic node selection, as illustrated in Section IV-C.

### C. FL with Probabilistic Node Selection (FedPNS)

Providing the variety of different nodes on contributing global model, to improve the convergence rate, one can seek to preferentially select the nodes with higher contribution (i.e., the nodes with i.i.d. dataset, as observed in [6], [17]). As such, we propose a probabilistic node selection design that dynamically changes the probability for each node to be selected in each communication round, based on their data distribution-related contribution, which can be distinguished by the procedure CHECK EXPECTATION in Optimal Aggregation.

As we know in each round of FL, a number of nodes are selected to participate in the local training and global aggregation. It is natural to lower the node selection probabilities for those nodes whose local updates slow model convergence. Therefore, on the server-side, we propose to dynamically change the probability for each node to be selected via using the output of Optimal Aggregation (i.e.,  $\bar{\mathcal{S}}_t$ ). In particular, the probabilities for those nodes that are labeled by the procedure CHECK EXPECTATION (i.e.,  $i \in \mathcal{S}_t - \bar{\mathcal{S}}_t$ ) are decreased according to the parameter  $x$  in (7), and the probabilities for all the rest nodes will be increased.

$$\Delta p_i^t = p_i^t \cdot \min[(x + \beta)^\alpha, 1], \quad i \in \mathcal{S}_t - \bar{\mathcal{S}}_t, \quad (7)$$

where  $p_i^t$  and  $\Delta p_i^t$  denote the probability for node  $i$  to be selected in the  $t$ -th global round, and its probability decrement in next round, respectively.  $\min$  function returns the minimum value among all arguments,  $x \in (0, 1]$  is defined as the ratio between the accumulated times that a node is labeled by the procedure CHECK EXPECTATION and the accumulated times that the node is selected,  $\alpha \in \mathbb{Z}^+, \beta \in [0, 1]$  are coefficients as explained in the following

- $\lim_{x \rightarrow \epsilon} (x + \beta)^\alpha \approx 1$ , where  $\epsilon \propto \alpha$  is constant.
- $\lim_{0 \rightarrow x \rightarrow v} (x + \beta)^\alpha \approx \beta$ , where  $v \propto \alpha$  is a constant.

---

### Algorithm 2 FL with Probabilistic Node Selection

---

#### Procedure FEDERATED OPTIMIZATION

**Input:**  $E, B, \eta, \mathcal{K}, T, p_i^t, i = 1, \dots, |\mathcal{K}|$

- 1: Server initializes  $\mathbf{w}^0, p_i^0 = 1/|\mathcal{K}|$
- 2: **for**  $t = 1, \dots, T$  **do**
- 3:   Server samples a subset  $\mathcal{S}_t$  of nodes according to  $p_i^{t-1}$
- 4:   Server sends  $\mathbf{w}^t$  to nodes  $i \in \mathcal{S}_t$
- 5:   Each node  $i \in \mathcal{S}_t$  finds  $\mathbf{w}_i^t$  to optimize  $F_i(\mathbf{w}^t)$  using SGD, as in (4), and sends back  $\Delta_i^t$  to the server
- 6:    $\mathbf{w}^{t+1}, \bar{\mathcal{S}}_t \leftarrow \text{OPTIMAL AGGREGATION}$
- 7:   Server updates the probability  $p_i^t, i = 1, \dots, |\mathcal{K}|$  by (7) and (8) for next round's usage
- 8: **return**  $\mathbf{w}^T$

#### Procedure OPTIMAL AGGREGATION

**Input:**  $\mathcal{S}_t, \Delta_i^t, v, \text{temp} = \{\}$

- 9: Direct to Algorithm 1
  - 10: **return**  $\mathbf{w}^{t+1}, \bar{\mathcal{S}}_t$
- 

$\alpha$  controls how big the probability decrement is achieved by  $(x + \beta)^\alpha$  given a ratio  $x$ . For example, a large value of  $\alpha$  brings an aggressive decrement since the probability decrement happens in a wide range  $(\beta, 1)$  as  $x$  increases within a small range  $(v, \epsilon)$ , making the node selection probability drop very quickly when  $x$  grows. Meanwhile, the large  $\alpha$  makes node selection sensitive to the identification mistake, which may prevent i.i.d. nodes from being selected in the subsequent rounds. However, setting a small value of  $\alpha$  is not consistently effective to differentiate the nodes since the probability change is marginal.  $\beta$  is adopted to keep the rate of probability change in a visible range  $[\beta, 1]$ . From experiments, we find out  $\alpha = 2, \beta = 0.7$  is a good choice that balances the tradeoff. The choice of  $\alpha$  and  $\beta$  is empirically investigated in Section V-C.

After getting the probability change for the labeled nodes (i.e.,  $i \in \mathcal{S}_t - \bar{\mathcal{S}}_t$ ), we equally increase the probability for all the rest nodes  $i \in \mathcal{K} - (\mathcal{S}_t - \bar{\mathcal{S}}_t)$ , as shown in (8).

$$p_i^{t+1} = \begin{cases} p_i^t - \Delta p_i^t & i \in \mathcal{S}_t - \bar{\mathcal{S}}_t \\ p_i^t + \frac{\sum_{i \in \mathcal{S}_t - \bar{\mathcal{S}}_t} \Delta p_i^t}{|\mathcal{K} - (\mathcal{S}_t - \bar{\mathcal{S}}_t)|} & i \in \mathcal{K} - (\mathcal{S}_t - \bar{\mathcal{S}}_t) \end{cases}, \quad (8)$$

where  $p_i^{t+1}, i \in \mathcal{K}$  are used for the  $(t+1)$ -th round.

We summarize the proposed FL design with probabilistic node selection and optimal aggregation in Algorithm 2. Particularly, in each communication round  $t$ , after the server receives the local update from participating nodes  $i \in \mathcal{S}_t$ , the server identifies the nodes that are labeled by the procedure CHECK EXPECTATION (i.e.,  $\mathcal{S}_t - \bar{\mathcal{S}}_t$ ) and the remaining nodes for aggregation  $\bar{\mathcal{S}}_t^*$ , which are used to regulate the probability for subsequent rounds (line 7) and aggregate the global model for this round (line 6).

### D. Convergence Rate of FedPNS

To facilitate theoretical analysis, we introduce the auxiliary parameter  $\mathbf{v}^t$ , which is optimized w.r.t. the global loss function  $F(\mathbf{v})$  in the centralized setting.  $\mathbf{v}^t$  is a virtual sequence since

$F(\mathbf{v})$  is only observable when all data samples are available at a central place. We use  $\tilde{\mathbf{w}}^t$  to denote model weight with full node participation, i.e.,  $\tilde{\mathbf{w}}^t = \sum_{i=1}^{|\mathcal{K}|} \frac{1}{|\mathcal{K}|} \mathbf{w}_i^t$ . We define that  $\mathbf{v}^t$  is “synchronized” with  $\tilde{\mathbf{w}}^t$  at the beginning of each global round, i.e., at the beginning of the  $t$ -th global round, the initial value of  $\mathbf{v}^t$  is set as  $\mathbf{v}^{t-1} = \tilde{\mathbf{w}}^{t-1}$ . At the end of the  $t$ -th global round, the update rule of the centralized SGD is as follows

$$\mathbf{v}^t = \mathbf{v}^{t-1} - \eta \left( - \sum_{j=1}^C q(y=j) \mathbb{E}_{\mathbf{x}|y=j} [\log l_j(\mathbf{v}^{t-1}, \mathbf{x}, y)] \right), \quad (9)$$

where  $q(y=j)$  is the population distribution over class  $j$ .

We first quantify the weight divergence  $\mathbb{E}_{\mathcal{S}_t} \|\mathbf{w}^t - \mathbf{v}^t\|$  between  $\mathbf{w}^t$  and  $\mathbf{v}^t$ , for any global round  $t, t = 1, \dots, T$ . Then, by combing the result in [12], we obtain the convergence rate of FedPNS.

**Theorem 1.** *Consider  $\mathcal{K}$  local nodes with equal data size and the data samples on node  $i \in \mathcal{K}$  follow the data distribution  $q^{(i)}$ . Let assumptions 1 and 3 hold. Assume a fixed number of local updates  $\tau$  exists between two consecutive global rounds. Then, the weight divergence in FedPNS after the  $(t-1)$ -th synchronization satisfies*

$$\mathbb{E}_{\mathcal{S}_t} \|\mathbf{w}^t - \mathbf{v}^t\| \leq \eta \sum_{i=1}^{|\mathcal{K}|} \left( p_i \gamma_i + \frac{1}{|\mathcal{K}|} q_{dif}^{(i)} \left( \sum_{k=1}^{\tau-1} a^k g_{\max}(\mathbf{v}^{t-\tau-1-k}) \right) \right), \quad (10)$$

where  $g_{\max}(\mathbf{v}) = \max_{j=1}^C \|\nabla \mathbb{E}_{\mathbf{x}|y=j} [\log l_j(\mathbf{v}, \mathbf{x}, y)]\|$ ,  $a = 1 + \eta L$ , and  $q_{dif}^{(i)} = \sum_{j=1}^C \|(q^{(i)}(y=j) - q(y=j))\|$ .

**Remark 1.** The weight divergence between  $\mathbf{w}^t$  and  $\mathbf{v}^t$  mainly comes from two parts, the bound of the norm of local gradient from each participating node, i.e.,  $\sum_{i=1}^{|\mathcal{K}|} p_i \gamma_i$ , and the weight divergence introduced by the difference between the data distribution on node and population distribution, i.e.,  $q_{dif}^{(i)}$ . FedPNS preferentially selects nodes with a smaller bounded gradient, which results in a smaller weight divergence, compared with node selection with equal probability in FedAvg, i.e.,  $\sum_{i=1}^{|\mathcal{K}|} p_i \gamma_i \leq \sum_{i=1}^{|\mathcal{K}|} \frac{1}{|\mathcal{K}|} \gamma_i$ .

**Theorem 2.** *When  $\eta \leq \frac{1}{L}$ , compared with FedAvg, FedPNS with a smaller weight divergence achieves tighter upper bound after  $T$  global rounds, i.e.,  $F(\mathbf{w}^T) - F(\mathbf{w}^*)$ , where  $F(\mathbf{w}^*)$  denotes the optimal model parameter that minimizes  $F(\mathbf{w})$ .*

*Proof.* Theorem 2 is proven by combing the weight divergence  $\mathbb{E}_{\mathcal{S}_t} \|\mathbf{w}^t - \mathbf{v}^t\|$  in Theorem 1 with the result in [12, Theorem 2]. From Theorem 1, it is straightforward to see that the weight divergence  $\mathbb{E}_{\mathcal{S}_t} \|\mathbf{w}^t - \mathbf{v}^t\|$  in FedPNS is smaller than that in FedAvg. From [12, Theorem 2], we have  $F(\mathbf{w}^T) - F(\mathbf{w}^*) \propto \mathbb{E}_{\mathcal{S}_t} \|\mathbf{w}^t - \mathbf{v}^t\|$ , i.e., a smaller weight divergence in each global round  $t, t = 1, \dots, T$  results in a smaller gap between the global loss after  $T$  global round and the global loss with optimal model,  $F(\mathbf{w}^T) - F(\mathbf{w}^*)$ , which completes the proof.

### E. Complexity Analysis

We consider the model in *float* format (i.e., 32 bits for each parameter) and the operations in algorithms are float point

operations. For simplicity, we consider a general  $n_{layer}$  layers fully connected neural network (FCNN) with the same number of parameters  $n$ , in each layer (i.e., the total parameters of model update/gradient is  $n \cdot n_{layer}$  and  $n \gg n_{layer}$  holds typically). The output of the  $k$ -th layer in forward propagation (FP) is represented as  $a^{(k)} = g(z^{(k)})$ ,  $z^{(k)} = \mathbf{w}^{(k)} a^{(k-1)}$ , where  $g(\cdot)$  is the activation function which is evaluated elementwise,  $\mathbf{w}^{(k)}$  is the model parameter in the  $k$ -th layer, and the bias component in FCNN is omitted for simplicity. We assume the number of features for the input layer is  $n$ . The computation in each layer is viewed as a matrix-vector multiplication, and an activation function, thus the complexity for multiplications in FP and for activation function applied in FP are  $\sum_{k=1}^{n_{layer}} n^2 = n_{layer} \cdot n^2$  and  $\sum_{k=1}^{n_{layer}} n = n_{layer} \cdot n$ , respectively. Therefore, the complexity for forward propagation of FCNN (also for CHECK LOSS, line 24) is  $\mathcal{O}(n^2)$  since  $n \gg n_{layer}$ . Given the total  $n \cdot n_{layer}$  parameters of local gradient, the complexity for arithmetic addition and arithmetic multiplication are  $\mathcal{O}(n)$  and  $\mathcal{O}(n^2)$ , respectively.

With regards to the complexity of the proposed Algorithm 1, we consider the procedures CHECK EXPECTATION and CHECK LOSS. In particular, generating a global gradients/model needs  $|\mathcal{S}_t|$  additions (line 20, line 23) and  $|\mathcal{S}_t|$  multiplications are needed for calculating the expectation values (i.e., temp, line 21). As such, the complexity for CHECK EXPECTATION is  $\mathcal{O}(n^2)$  since the number of local updates  $|\mathcal{S}_t|$  and the number of iterations  $|\mathcal{S}_t| - v$  (line 4) are much smaller than  $n$ . The complexity for Algorithm 1 is  $\mathcal{O}(n^2 + n^2 + n) = \Theta(2n^2 + n)$ . In Algorithm 2, the complexity for adjusting the probability (line 7) is  $\mathcal{O}(|\mathcal{K}|)$ , which is marginal compared with the complexity of OPTIMAL AGGREGATION (line 6). Therefore, the complexity for Algorithm 2 is  $\mathcal{O}(n^2 + n^2 + n) = \Theta(2n^2 + n)$ . Compared with local training including FP and back propagation (BP) (the complexity for BP is  $\mathcal{O}(n^3)$ ) at node side, the overhead of the proposed algorithms that are conducted at server side is marginal and can be ignored.

## V. EVALUATION AND ANALYSIS

We now present empirical results for the proposed probabilistic node selection strategy. We implement FedPNS on different tasks, models, datasets, and compare with commonly used benchmark FedAvg. We first demonstrate the effectiveness of the proposed Optimal Aggregation in enlarging the expected decrement of FL global loss and in identifying the potential adverse local updates (Section V-A). Then, the superiority of the proposed FedPNS in presence of different data heterogeneity is illustrated in Section V-B. The choice  $\alpha$  and  $\beta$  for adjusting node probability in FedPNS are discussed in Section V-C. In Section V-D, by tracking the norm of gradient on different nodes, the Assumption 3 is empirically justified. In addition, we also compare the proposed FedPNS with an existing work that uses the norm of gradient  $\|\nabla F_i(\mathbf{w}^t)\|$  [19] for node selection. All code, data, and experiments are publicly available as an open-source GitHub repository at: [github.com/HongdaWu1226/FedPNS](https://github.com/HongdaWu1226/FedPNS).

We briefly describe our adopted datasets, learning model, and experiment setting as follows.



**Synthetic data.** To better characterize the data heterogeneity and study its impact on model convergence, we generate synthetic data by following the similar setup as in [9], [29]. In particular, the data samples  $\{\mathbf{x}, y\}$  on local node  $i$  are generated according to the model  $y = \arg\max(\text{softmax}(\mathbf{w}\mathbf{x} + b))$ ,  $\mathbf{x} \in \mathbb{R}^{60}$ ,  $\mathbf{w} \in \mathbb{R}^{10 \times 60}$ ,  $b \in \mathbb{R}^{10}$ . We set  $\mathbf{w}, b \sim \mathcal{N}(0, 1)$ . For the data on i.i.d. nodes,  $\mathbf{x}$  follows the same distribution  $\mathcal{N}(0, \Sigma)$ , where  $\Sigma$  is diagonal with  $\Sigma_{r,r} = r^{-1.2}$ . For the data samples on non-i.i.d. node  $i$ ,  $\mathbf{x} \sim \mathcal{N}(o_i, \Sigma)$ , each element in the mean vector  $o_i$  is drawn from  $\mathcal{N}(B_i, 1)$ ,  $B_i \sim \mathcal{N}(0, \varrho)$ . As such, a big value of  $\varrho$  denotes a more heterogeneous data scenario. The training set and testing set are randomly split with 80%–20% proportion on each node. A Multinomial Logistic Regression (MLR) model is applied to the synthetic data.

**Real data.** We explore different learning objectives on different real datasets, which are considered in prior works [6], [9]. In Section V-B, we start with a convex classification problem with MNIST [30] using MLR model. Then, for the non-convex setting, we consider two CNN models for MNIST and CIFAR-10 [31], which are referred as CNN-M<sup>6</sup> and CNN-C<sup>7</sup> hereinafter.

Through the experimental result, unless otherwise specified, we evaluate the accuracy of the trained models using the testing set from each dataset. The fraction for selecting nodes is set to be  $c = 0.2$ ,  $|\mathcal{S}_t| = c|\mathcal{K}| = 10$ ,  $D_i = 200$ ,  $B = 20$ ,  $E = 1$ ,  $T = 200$ ,  $\eta = 0.01$ , decay rate = 0.995,  $\nu = 0.7$ ,  $\bar{B} = 128$ . For real datasets, the overall data heterogeneity is measured by  $\sigma$  and the skewness of dataset on non-i.i.d. nodes is represented by  $\rho$ . For example,  $\sigma = 0.2, \rho = 2$  means that  $\sigma|\mathcal{K}| = 10$  nodes are equipped with i.i.d. dataset, where non-i.i.d. dataset lay on the rest  $(1 - \sigma)|\mathcal{K}| = 40$  nodes, and the data samples on which are evenly belong to 2 labels. As such, a small  $\sigma, \rho$  indicates a higher data heterogeneity.

#### A. Performance of Optimal Aggregation

In this part, we conduct an experiment to illustrate the performance of the proposed Optimal Aggregation algorithm. Particularly, we adopt CNN-M model on MNIST dataset where the data heterogeneity is set to be  $\sigma = 0.5, \rho = 1$ . In each global round, we randomly select  $|\mathcal{S}_t| = 10$  nodes while guaranteeing the participating nodes include half i.i.d. nodes and half non-i.i.d. nodes. To avoid the randomness of node selection, the participating nodes in each round are kept as the same for FedAvg [9] and the proposed Optimal Aggregation algorithm.

As shown in the upper part of Fig. 1, the proposed Optimal Aggregation algorithm can achieve lower training loss than FedAvg. When the global model is not robust in the several initial rounds, the local updates are more diverse due to the data heterogeneity, thus excluding adverse

<sup>6</sup>The CNN-M model has 7 layers with the following structure:  $5 \times 5 \times 10$  Convolutional  $\rightarrow 2 \times 2$  MaxPool  $\rightarrow 5 \times 5 \times 20$  Convolutional  $\rightarrow 2 \times 2$  MaxPool  $\rightarrow 320 \times 50$  Fully connected  $\rightarrow 50 \times 10$  Fully connected  $\rightarrow$  Softmax. The second convolutional layer is with 50% dropout. All Convolutional and Fully connected layers are mapped by ReLU activation.

<sup>7</sup>The CNN-C model has 8 layers as structured follows:  $5 \times 5 \times 6$  Convolutional  $\rightarrow 2 \times 2$  MaxPool  $\rightarrow 5 \times 5 \times 16$  Convolutional  $\rightarrow 2 \times 2$  MaxPool  $\rightarrow 400 \times 120$  Fully connected  $\rightarrow 120 \times 84$  Fully connected  $\rightarrow 84 \times 10$  Fully connected  $\rightarrow$  Softmax. ReLU activation is applied to all layers.

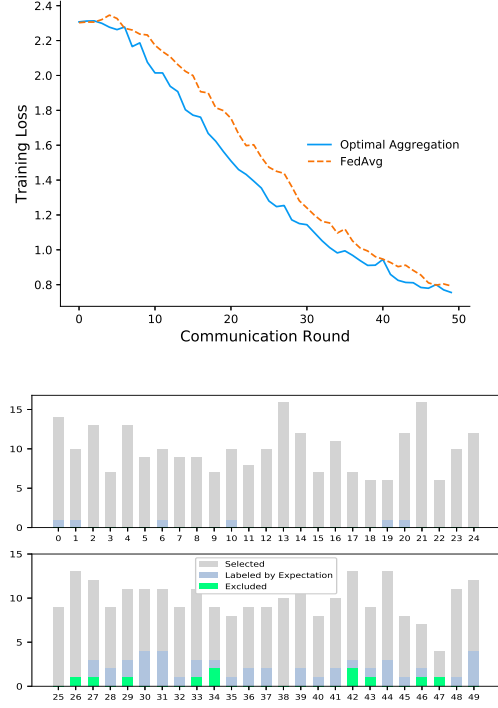


Fig. 1. Performance of the proposed Optimal Aggregation. (1) Upper: The training loss on the MNIST dataset when different aggregation strategies are adopted. Optimal Aggregation and FedAvg aggregate local updates over  $\mathcal{S}_t^*$  and  $\mathcal{S}_t$ , respectively. (2) Bottom: We use a triple to observe the result of Optimal Aggregation, which includes the accumulated times that each node is selected, labeled by CHECK EXPECTATION, and excluded eventually by CHECK LOSS during FL model training. The upper and bottom row refer to the results for i.i.d. nodes and non-i.i.d. nodes, respectively.

local updates is more effective. We count the accumulated times that each node is selected, labeled by the procedure CHECK EXPECTATION (line 7 in Algorithm 1), and finally excluded by the procedure CHECK LOSS (line 14 in Algorithm 1). As we can see from the bottom part of Fig. 1, i) the i.i.d. nodes (i.e., with index “0”,  $\dots$ , “24”) are never been excluded, yet some of the non-i.i.d. nodes (e.g., “26”, “27”, “34”, etc.) have been excluded for many times. ii) Almost all non-i.i.d. nodes are labeled at least one time, which illustrates the effectiveness of Optimal Aggregation in identifying the nodes with the skewed dataset.

#### B. Data Heterogeneity

In this section, we use different combinations of  $\sigma$  and  $\rho$  to investigate the performance of the proposed FedPNS scheme in presence of different data heterogeneity. Through all experiments,  $\alpha$  and  $\beta$  are chosen to be 2 and 0.7 respectively. The choice  $\alpha$  and  $\beta$  are discussed in Section V-C.

1) *MLR Model with Synthetic Data:* We follow the description in Section V-A to generate synthetic data samples. The ratio of i.i.d. nodes is set to be  $\sigma = 0.2, 0.3$ , and 0.5 with  $\varrho = 0.5$  and 1. For each node  $i$ , the number of data samples  $D_i = 1000$  and the number of epochs for local training is  $E = 20$ . In Fig. 2, we study how data heterogeneity affects



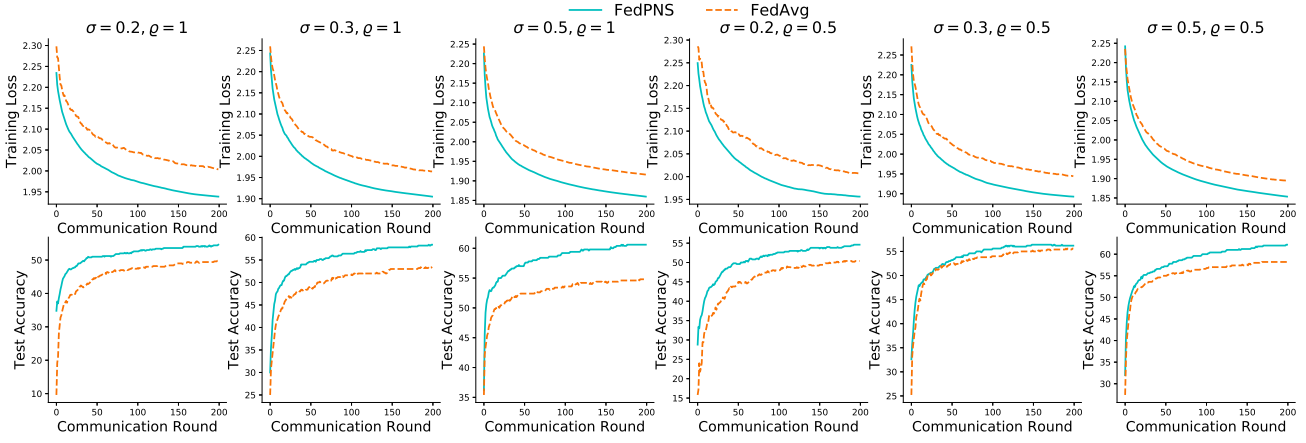


Fig. 2. Effect of data heterogeneity on convergence. (1) Top row: we show the training loss on synthetic dataset whose data heterogeneity decreases from left to right (with a fixed  $\sigma$  or  $\rho$ ). (1) Bottom row: we show the corresponding test accuracy.

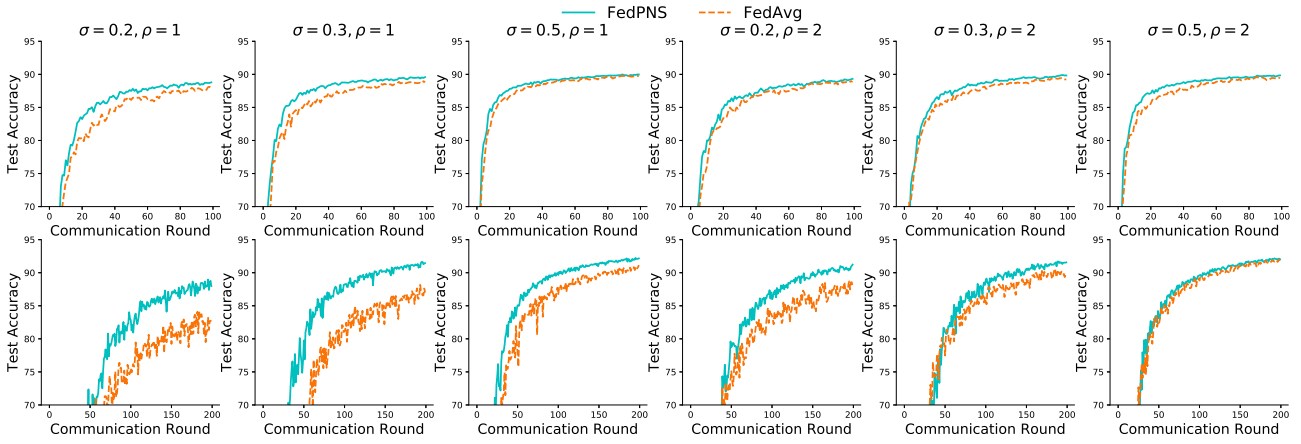


Fig. 3. Test accuracy over communication rounds of FedPNS and FedAvg with different data heterogeneity. Upper and lower subplots correspond to training performance when the MLR model and CNN-M model are adopted for MNIST, respectively. A smaller  $\alpha, \rho$  indicates a higher data heterogeneity.

model convergence using MLR model and synthetic dataset. As we can see from Fig. 2, as the data heterogeneity increases, i.e.,  $\sigma = 0.5, 0.3$  and  $0.2$  with fixed  $\rho = 1$  or  $0.5$ , FedAvg slows to converge (i.e., higher training loss) with a decreasing test accuracy in the meantime. FedPNS achieves a lower training loss and higher test accuracy, compared with FedAvg in all data setting.

2) *MLR, CNN-M Model for MNIST*: As we can tell from Fig. 3, FedPNS converges faster and achieves a higher test accuracy, compared with FedAvg for both MLR and CNN model regardless of different data heterogeneity. FedPNS achieves better improvement when the CNN model is adopted, compared with the scenario when the MLR model is utilized, which attributes to the limited learning capability of MLR. In addition, it is observable that as the data becomes more heterogeneous, the performance enhancement is enlarged (i.e.,  $\alpha$  decreases from  $0.5$  to  $0.2$  for a given  $\beta$ , or  $\beta$  changes from  $2$  to  $1$  for a given  $\alpha$ ). When the number of i.i.d. nodes is limited and the non-i.i.d. nodes are equipped with highly skewed dataset (e.g.,  $\sigma = 0.2, \rho = 1$  and  $\sigma = 0.3, \rho = 1$ ), FedPNS gains remarkable performance improvement, which verifies the effectiveness of FedPNS in identifying and selecting the

nodes that contribute global model better. For the scenario with the lowest data heterogeneity (i.e.,  $\sigma = 0.5, \rho = 2$ ), the performance gap between FedPNS and FedAvg is not obvious. This is because the impact of the non-i.i.d. nodes on the convergence is reduced when a large number of i.i.d. nodes can be selected.

3) *CNN-C Model for CIFAR-10*: For the more complex three channel image classification task, the number of local epoch is set to be  $E = 5$ . As we can see from Fig. 4, compared with FedAvg, FedPNS converges faster and leads to a higher test accuracy, especially for the high data heterogeneity scenario (i.e.,  $\sigma = 0.2$  and  $0.3, \rho = 1$ ). The performance improvement of FedPNS is not obvious when  $\sigma = 0.2, \rho = 2$ , this is because the small number of i.i.d. nodes with less heterogeneous data samples on non-i.i.d. nodes makes FedPNS hard to distinguish the node contribution.

### C. Choosing $\alpha$ and $\beta$

The choice of  $\alpha$  and  $\beta$  gives an impact on FedPNS. As discussed in Section IV-C, a large value of  $\alpha$  can help increase the model convergence rate by aggressively adjusting the node probability. On the other hand, a large value of  $\alpha$  also makes

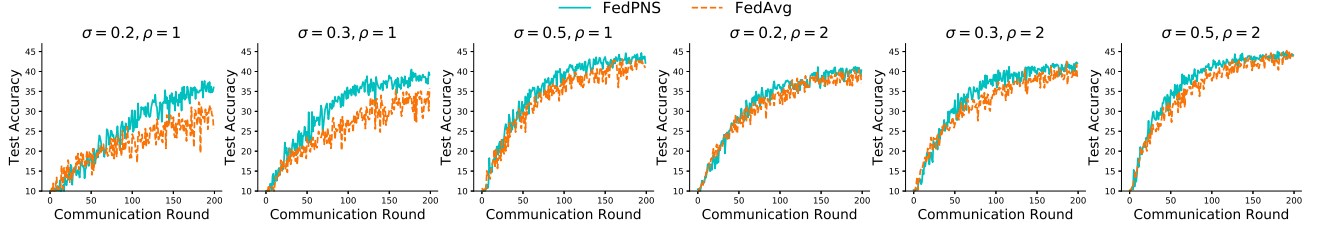


Fig. 4. Test accuracy over communication rounds of FedPNS and FedAvg with different data heterogeneity. CNN-C model is adopted for CIFAR-10.

node selection sensitive to the identification mistake, which may negatively impact the convergence. A similar effect is achieved by  $\beta$ , which keeps the rate of probability change in a range  $[\beta, 1]$ . We studied the effect of different  $\alpha$  and  $\beta$  via

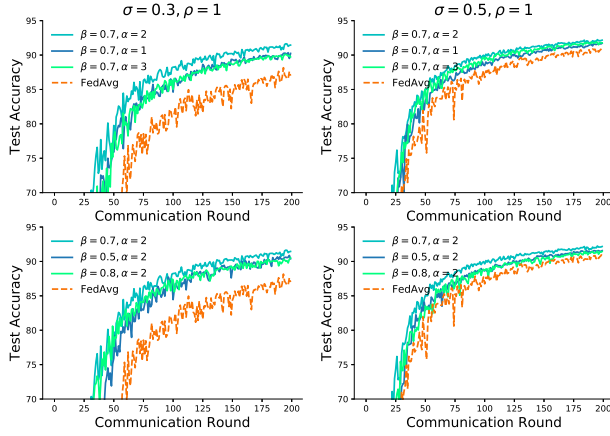


Fig. 5. Effect of adopting different  $\alpha$  and  $\beta$ . We heuristically choosing  $\alpha \in \mathbb{Z}^+$ ,  $\beta \in [0, 1]$  in ascending order. The top row and bottom row correspond to the performance with varied  $\alpha$  and  $\beta$ , respectively. CNN-M on MNIST is adopted.

heuristically choosing  $\alpha \in \mathbb{Z}^+$ ,  $\beta \in [0, 1]$  in ascending order. From the top row of Fig. 5, for a fixed  $\beta = 0.7$ , increasing  $\alpha$  from 1 to 2 boosts performance. However, keep increasing  $\alpha$  does not consistently embrace performance gain, this is because FedPNS becomes more sensitive to identification mistakes, which may prevent i.i.d. nodes from being selected in the subsequent rounds. Similarly, from the bottom plot of Fig. 5, for a fixed  $\alpha = 2$ , increasing  $\beta$  from 0.5 to 0.7 promotes model performance. However, further increasing  $\beta$  to 0.8 leads to a degraded performance. Empirically, we find  $\alpha = 2$ ,  $\beta = 0.7$  that balances the tradeoff and leads to the best performance.

#### D. Other Comparison

In this section, we take one experimental case as an example to demonstrate the bounded norm of local gradient  $\|\nabla F_i(\mathbf{w}^t)\|$ , which is related to the data distribution on each node. Besides, we compare the proposed FedPNS with another node selection scheme BN2 [19], which chooses the nodes with higher  $\|\nabla F_i(\mathbf{w}^t)\|$  for aggregation. Specifically, in each global round, BN2 first randomly selects  $|\mathcal{M}|$  nodes for local training. After that, the participating nodes send their gradient norm

$\|\nabla F_i(\mathbf{w}^t)\|$ ,  $i \in \mathcal{M}$  to the server. The server chooses the first  $|\mathcal{S}_t|$  local updates for model aggregation by sorting  $\|\nabla F_i(\mathbf{w}^t)\|$ ,  $i \in \mathcal{M}$  in descending order.

In this experiment,  $|\mathcal{M}|$  is set to be 20. We track the norm of gradient for each participating node  $i \in \mathcal{M}$  statistically in each global round. As we can see from Fig. 6, the averaged gradient norm from i.i.d. nodes is smaller than that from non-i.i.d. nodes. This is because the data distribution on i.i.d. nodes is more similar to population distribution that is defined over all nodes. As such, preferentially scheduling the nodes with higher norm of gradient would slow the convergence, as shown in the bottom of Fig. 6.

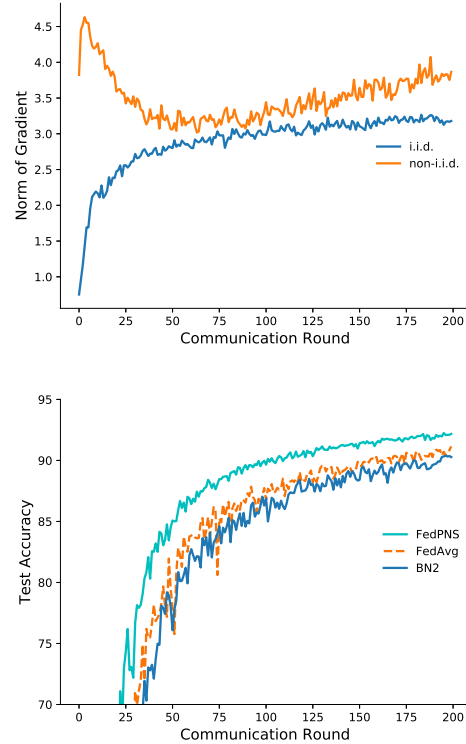


Fig. 6. Node selection design with different importance indicator. FedPNS chooses nodes by measuring the data distribution on local nodes, while BN2 selects nodes according to the norm of gradient. (1) Top plot: we track the averaged gradient norm of node  $i \in \mathcal{M}$  with different data distribution, where each node is selected from  $\mathcal{K}$  randomly. (2) Bottom plot: we compare the test accuracy for different node selection designs. CNN-M on MNIST is adopted with  $\sigma = 0.5$ ,  $\rho = 1$ .

## VI. CONCLUSION

In this paper, we have presented our design of FedPNS algorithm, a probabilistic node selection strategy that can preferentially select nodes to boost model convergence of FL with non-i.i.d. datasets. FedPNS adjusts the probability for each node to be selected in each round based on the result of the proposed Optimal Aggregation algorithm, which is able to find out the optimal subset of local updates from participating nodes and excludes the adverse local updates for a better model aggregation, by measuring the relationship between the local gradient and the global gradient from participating nodes. The convergence rate improvement of the FedPNS design over FedAvg is analyzed theoretically. Finally, experimental results on different tasks, models, and datasets have shown that FL training with FedPNS accelerates model convergences and leads to higher test accuracy, as compared to FedAvg.

## APPENDIX

### A. Proof of Lemma 1

From the  $L$ -smooth of  $F(\mathbf{w})$  and applying Taylor expansion, we have

$$F(\mathbf{w}^{t+1}) \leq F(\mathbf{w}^t) + \langle \nabla F(\mathbf{w}^t), \mathbf{w}^{t+1} - \mathbf{w}^t \rangle + \frac{L}{2} \|\mathbf{w}^{t+1} - \mathbf{w}^t\|^2. \quad (\text{A1})$$

• Bounding  $\|\mathbf{w}^{t+1} - \mathbf{w}^t\|^2$ : By the definition of the global aggregation in (5) and local update calculated by (4), we have

$$\begin{aligned} \|\mathbf{w}^{t+1} - \mathbf{w}^t\|^2 &= (\mathbb{E}_{i \sim \mathcal{S}_t} [\|\mathbf{w}^{t+1} - \mathbf{w}^t\|])^2 \\ &= \eta^2 (\mathbb{E}_{i \sim \mathcal{S}_t} [\|\nabla F_i(\mathbf{w}^t)\|])^2 \\ &\leq \eta^2 \mathbb{E}_{i \sim \mathcal{S}_t} [\|\nabla F_i(\mathbf{w}^t)\|^2] \\ &\leq \eta^2 \|\nabla F(\mathbf{w}^t)\|^2 \delta^2, \end{aligned} \quad (\text{A2})$$

where inequality 1 holds because of Cauchy-Schwarz inequality and the last inequality is due to the bounded dissimilarity assumption.

• Bounding  $\langle \nabla F(\mathbf{w}^t), \mathbf{w}^{t+1} - \mathbf{w}^t \rangle$ : Again, by the definition of the global aggregation for  $\mathbf{w}^{t+1}$  and SGD optimization, we have

$$\langle \nabla F(\mathbf{w}^t), \mathbf{w}^{t+1} - \mathbf{w}^t \rangle = -\eta \mathbb{E}_{i \sim \mathcal{S}_t} [\langle \nabla F(\mathbf{w}^t), \nabla F_i(\mathbf{w}^t) \rangle]. \quad (\text{A3})$$

Plugging (A2) and (A3) into (A1), we obtain

$$\begin{aligned} F(\mathbf{w}^{t+1}) - F(\mathbf{w}^t) &\leq -\eta \mathbb{E}_{i \sim \mathcal{S}_t} [\langle \nabla F(\mathbf{w}^t), \nabla F_i(\mathbf{w}^t) \rangle] \\ &\quad + \frac{L\eta^2}{2} \|\nabla F(\mathbf{w}^t)\|^2 \delta^2. \end{aligned} \quad (\text{A4})$$

### B. Proof of Theorem 1

At any global round  $t$ , the weight divergence between the model  $\mathbf{w}^t$  with partial node participation and centralized model  $\mathbf{v}^t$  is bounded as follows

$$\begin{aligned} \mathbb{E}_{\mathcal{S}_t} \|\mathbf{w}^t - \mathbf{v}^t\| &= \mathbb{E}_{\mathcal{S}_t} \|\mathbf{w}^t - \tilde{\mathbf{w}}^t + \tilde{\mathbf{w}}^t - \mathbf{v}^t\| \\ &\leq \mathbb{E}_{\mathcal{S}_t} \|\mathbf{w}^t - \tilde{\mathbf{w}}^t\| + \|\tilde{\mathbf{w}}^t - \mathbf{v}^t\|. \end{aligned} \quad (\text{B1})$$

We will separately bound the last two terms on the right-hand side of the above inequality.

• Bounding  $\|\tilde{\mathbf{w}}^t - \mathbf{v}^t\|$ : In this part, to facilitate analysis, we introduce the index of local update, e.g., the models  $\tilde{\mathbf{w}}^t$  and  $\mathbf{v}^t$  are represented by  $\tilde{\mathbf{w}}^{t\tau}$  and  $\mathbf{v}^{t\tau}$  since  $\tau$  times of local SGD are applied in each global round.

Based on the definition of  $\tilde{\mathbf{w}}^t$  and  $\mathbf{v}^t$ , we have

$$\begin{aligned} \|\tilde{\mathbf{w}}^t - \mathbf{v}^t\| &= \|\tilde{\mathbf{w}}^{t\tau} - \mathbf{v}^{t\tau}\| = \left\| \sum_{i=1}^{|\mathcal{K}|} \frac{D_i}{\sum_{i=1}^{|\mathcal{K}|} D_i} \mathbf{w}_i^{t\tau} - \mathbf{v}^{t\tau} \right\| \\ &\stackrel{1}{=} \left\| \sum_{i=1}^{|\mathcal{K}|} \frac{1}{|\mathcal{K}|} (\mathbf{w}_i^{t\tau-1} - \eta \nabla F_i(\mathbf{w}_i^{t\tau-1})) - \mathbf{v}^{t\tau-1} + \eta \nabla F(\mathbf{v}^{t\tau-1}) \right\| \\ &\stackrel{2}{\leq} \left\| \sum_{i=1}^{|\mathcal{K}|} \frac{1}{|\mathcal{K}|} \mathbf{w}_i^{t\tau-1} - \mathbf{v}^{t\tau-1} \right\| + \eta \left\| \sum_{i=1}^{|\mathcal{K}|} \frac{1}{|\mathcal{K}|} \sum_{j=1}^C q^{(i)}(y=j) \right. \\ &\quad \left. (\nabla \mathbb{E}_{\mathbf{x}|y=j} [\log l_j(\mathbf{w}_i^{t\tau-1}, \mathbf{x}, y)] - \nabla \mathbb{E}_{\mathbf{x}|y=j} [\log l_j(\mathbf{v}^{t\tau-1}, \mathbf{x}, y)]) \right\| \\ &\stackrel{3}{=} \left\| \sum_{i=1}^{|\mathcal{K}|} \frac{1}{|\mathcal{K}|} \mathbf{w}_i^{t\tau-1} - \mathbf{v}^{t\tau-1} \right\| + \eta \left\| \sum_{i=1}^{|\mathcal{K}|} \frac{1}{|\mathcal{K}|} (\nabla F_i(\mathbf{w}_i^{t\tau-1}) - \nabla F_i(\mathbf{v}^{t\tau-1})) \right\| \\ &\stackrel{4}{\leq} \sum_{i=1}^{|\mathcal{K}|} \frac{1}{|\mathcal{K}|} (1 + \eta L) \|\mathbf{w}_i^{t\tau-1} - \mathbf{v}^{t\tau-1}\|, \end{aligned} \quad (\text{B2})$$

where equality 1 holds by the updating rule of SGD and by that all nodes are with equal data size. Inequality 2 holds by applying triangle inequality and by the observation that for each class, the data distribution over all nodes is the same as the distribution over the whole data samples, i.e.,  $j \in [C], q(y=j) = \sum_{i=1}^{|\mathcal{K}|} \frac{1}{|\mathcal{K}|} q^{(i)}(y=j)$ . Equality 3 holds by (2), (3) and (9). and inequality 4 holds by Assumption 1 that the local loss function is  $L$ -smooth.

For node  $i \in \mathcal{K}$ ,  $\|\mathbf{w}_i^{t\tau-1} - \mathbf{v}^{t\tau-1}\|$  is bounded as

$$\begin{aligned} &\|\mathbf{w}_i^{t\tau-1} - \mathbf{v}^{t\tau-1}\| \\ &= \|\mathbf{w}_i^{t\tau-2} - \eta \nabla F_i(\mathbf{w}_i^{t\tau-2}) - \mathbf{v}^{t\tau-2} + \eta \nabla F(\mathbf{v}^{t\tau-2})\| \\ &\leq \|\mathbf{w}_i^{t\tau-2} - \mathbf{v}^{t\tau-2}\| + \eta \left\| \sum_{j=1}^C q^{(i)}(y=j) \nabla \mathbb{E}_{\mathbf{x}|y=j} [\log l_j(\mathbf{w}_i^{t\tau-2}, \mathbf{x}, y)] \right. \\ &\quad \left. - \sum_{j=1}^C q(y=j) \nabla \mathbb{E}_{\mathbf{x}|y=j} [\log l_j(\mathbf{v}^{t\tau-2}, \mathbf{x}, y)] \right\| \\ &\stackrel{5}{\leq} \|\mathbf{w}_i^{t\tau-2} - \mathbf{v}^{t\tau-2}\| + \eta \left\| \sum_{j=1}^C q^{(i)}(y=j) \right. \\ &\quad \left. (\nabla \mathbb{E}_{\mathbf{x}|y=j} [\log l_j(\mathbf{w}_i^{t\tau-2}, \mathbf{x}, y)] - \nabla \mathbb{E}_{\mathbf{x}|y=j} [\log l_j(\mathbf{v}^{t\tau-2}, \mathbf{x}, y)]) \right\| \\ &\quad + \eta \left\| \sum_{j=1}^C (q^{(i)}(y=j) - q(y=j)) \nabla \mathbb{E}_{\mathbf{x}|y=j} [\log l_j(\mathbf{v}^{t\tau-2}, \mathbf{x}, y)] \right\| \\ &\stackrel{6}{=} \|\mathbf{w}_i^{t\tau-2} - \mathbf{v}^{t\tau-2}\| + \eta \|\nabla F_i(\mathbf{w}_i^{t\tau-2}) - \nabla F_i(\mathbf{v}^{t\tau-2})\| \\ &\quad + \eta \left\| \sum_{j=1}^C (q^{(i)}(y=j) - q(y=j)) \nabla \mathbb{E}_{\mathbf{x}|y=j} [\log l_j(\mathbf{v}^{t\tau-2}, \mathbf{x}, y)] \right\| \\ &\stackrel{7}{\leq} (1 + \eta L) \|\mathbf{w}_i^{t\tau-2} - \mathbf{v}^{t\tau-2}\| \\ &\quad + \eta g_{\max}(\mathbf{v}^{t\tau-2}) \sum_{j=1}^C \|(q^{(i)}(y=j) - q(y=j))\|, \end{aligned} \quad (\text{B3})$$

where inequality 5 holds by introducing a term  $\sum_{j=1}^C q^{(i)}(y = j) \nabla \mathbb{E}_{\mathbf{x}|y=j} [\log l_j(\mathbf{v}^{t\tau-2}, \mathbf{x}, y)]$  and applying triangle inequality. Equality 6 holds by (2), (3) and (9). Inequality 7 holds by Assumption 1 and by defining  $g_{\max}(\mathbf{v}^{t\tau-2}) = \max_{j=1}^C \|\nabla \mathbb{E}_{\mathbf{x}|y=j} [\log l_j(\mathbf{v}^{t\tau-2}, \mathbf{x}, y)]\|$ .

Based on (B3), by mathematical induction and setting  $a = 1 + \eta L$ , we have

$$\begin{aligned}
 & \|\mathbf{w}_i^{t\tau-1} - \mathbf{v}^{t\tau-1}\| \\
 & \leq a \|\mathbf{w}_i^{t\tau-2} - \mathbf{v}^{t\tau-2}\| \\
 & \quad + \eta \sum_{j=1}^C \|(q^{(i)}(y = j) - q(y = j))\| g_{\max}(\mathbf{v}^{t\tau-2}) \\
 & \leq a^2 \|\mathbf{w}_i^{t\tau-3} - \mathbf{v}^{t\tau-3}\| + \eta \sum_{j=1}^C \|(q^{(i)}(y = j) - q(y = j))\| \\
 & \quad (g_{\max}(\mathbf{v}^{t\tau-2}) + a g_{\max}(\mathbf{v}^{t\tau-3})) \\
 & \vdots \\
 & \leq a^{\tau-1} \|\mathbf{w}_i^{(t-1)\tau} - \mathbf{v}^{(t-1)\tau}\| + \eta \sum_{j=1}^C \|(q^{(i)}(y = j) - q(y = j))\| \\
 & \quad \left( \sum_{k=0}^{\tau-2} a^k g_{\max}(\mathbf{v}^{t\tau-2-k}) \right). \tag{B4}
 \end{aligned}$$

Substituting (B4) to (B2), we obtain

$$\begin{aligned}
 \|\tilde{\mathbf{w}}^t - \mathbf{v}^t\| & \leq \sum_{i=1}^{|\mathcal{K}|} \frac{1}{|\mathcal{K}|} (a^\tau \|\mathbf{w}_i^{(t-1)\tau} - \mathbf{v}^{(t-1)\tau}\| \\
 & \quad + \eta \sum_{j=1}^C \|(q^{(i)}(y = j) - q(y = j))\| \left( \sum_{k=1}^{\tau-1} a^k g_{\max}(\mathbf{v}^{t\tau-1-k}) \right)). \tag{B5}
 \end{aligned}$$

Since  $\mathbf{v}^t$  is “synchronized” with  $\tilde{\mathbf{w}}^t$  at the beginning of each global round, we ignore the first item of the right hand side of (B5), which is the weight divergence accumulated from the previous round. Thus, the weight divergence  $\|\tilde{\mathbf{w}}^t - \mathbf{v}^t\|$  between two consecutive global round is represented as

$$\|\tilde{\mathbf{w}}^t - \mathbf{v}^t\| \leq \eta \sum_{i=1}^{|\mathcal{K}|} \frac{1}{|\mathcal{K}|} q_{dif}^{(i)} \left( \sum_{k=1}^{\tau-1} a^k g_{\max}(\mathbf{v}^{t\tau-1-k}) \right), \tag{B6}$$

where  $q_{dif}^{(i)} = \sum_{j=1}^C \|(q^{(i)}(y = j) - q(y = j))\|$ .

• Bounding  $\|\mathbf{w}^t - \tilde{\mathbf{w}}^t\|$ : We follow the identical sampling distribution (i.e.,  $\{p_1, p_2, \dots, p_{|\mathcal{K}|}\}$ ) to select  $|\mathcal{S}_t|$  nodes from  $|\mathcal{K}|$  nodes and let  $\mathcal{S}_t = \{k_1, \dots, k_{|\mathcal{S}_t|}\}$  denote the set of indices of chosen nodes. The global model in FL with partial node participation is represented as  $\mathbf{w}^t = \frac{1}{|\mathcal{S}_t|} \sum_{i=1}^{|\mathcal{S}_t|} \mathbf{w}_{k_i}^t$ . Taking expectation over  $\mathcal{S}_t$ , we have

$$\mathbb{E}_{\mathcal{S}_t} \|\mathbf{w}^t - \tilde{\mathbf{w}}^t\| = \mathbb{E}_{\mathcal{S}_t} \frac{1}{|\mathcal{S}_t|} \sum_{i=1}^{|\mathcal{S}_t|} \|\mathbf{w}_{k_i}^t - \tilde{\mathbf{w}}^t\| = \sum_{i=1}^{|\mathcal{K}|} p_i \|\mathbf{w}_i^t - \tilde{\mathbf{w}}^t\|, \tag{B7}$$

where the last equality in (B7) is obtained by the following the observation  $\mathbb{E}_{\mathcal{S}_t} \sum_{i \in \mathcal{S}_t} x_i = \mathbb{E}_{\mathcal{S}_t} \sum_{i=1}^{|\mathcal{S}_t|} x_{k_i} = |\mathcal{S}_t| \mathbb{E}_{\mathcal{S}_t} x_{k_i} = |\mathcal{S}_t| \sum_{i=1}^{|\mathcal{K}|} p_i x_i$  given  $\mathcal{S}_t = \{x_{k_1}, \dots, x_{k_{|\mathcal{S}_t|}}\} \subset \mathcal{K}$ , and by replacing  $x_i$  with  $\mathbf{w}_i^t$  in the above observation.

We consider the model parameter in previous global round  $\mathbf{w}_i^{t-1}$ , which is identical for any  $i \in \mathcal{K}$ . As such, we have  $\sum_{i=1}^{|\mathcal{K}|} p_i (\mathbf{w}_i^t - \mathbf{w}^{t-1}) = \tilde{\mathbf{w}}^t - \tilde{\mathbf{w}}^{t-1}$ . Thus, the above equation can be bounded as

$$\begin{aligned}
 \sum_{i=1}^{|\mathcal{K}|} p_i \|\mathbf{w}_i^t - \tilde{\mathbf{w}}^t\| & = \sum_{i=1}^{|\mathcal{K}|} p_i \underbrace{\|(\mathbf{w}_i^t - \tilde{\mathbf{w}}^{t-1}) - (\tilde{\mathbf{w}}^t - \tilde{\mathbf{w}}^{t-1})\|}_{\mathbf{X}} \\
 & \leq \sum_{i=1}^{|\mathcal{K}|} p_i \|\mathbf{w}_i^t - \tilde{\mathbf{w}}^{t-1}\|, \tag{B8}
 \end{aligned}$$

where the last equality holds because  $\mathbb{E}\|\mathbf{X} - \mathbb{E}[\mathbf{X}]\| \leq \mathbb{E}\|\mathbf{X}\|$ .

Substituting (B8) into (B7), we have,

$$\begin{aligned}
 \mathbb{E}_{\mathcal{S}_t} \|\mathbf{w}^t - \tilde{\mathbf{w}}^t\| & \leq \sum_{i=1}^{|\mathcal{K}|} p_i \|\mathbf{w}_i^t - \tilde{\mathbf{w}}^{t-1}\| \\
 & \leq \sum_{i=1}^{|\mathcal{K}|} p_i \|\mathbf{w}_i^t - \mathbf{w}_i^{t-1}\| \\
 & \leq \sum_{i=1}^{|\mathcal{K}|} p_i \|\eta \nabla F_i(\mathbf{w}^{t-1})\| \\
 & \leq \eta \sum_{i=1}^{|\mathcal{K}|} p_i \gamma_i, \tag{B9}
 \end{aligned}$$

where the last inequality results from Assumption 3.

Finally, Theorem 1 is proved by substituting (B9) and (B6) into (B1).

## REFERENCES

- [1] M. Chiang and T. Zhang, “Fog and iot: An overview of research opportunities,” *IEEE Internet of Things Journal*, vol. 3, no. 6, pp. 854–864, 2016.
- [2] Z. Xiong, Y. Zhang, D. Niyato, P. Wang, and Z. Han, “When mobile blockchain meets edge computing,” *IEEE Communications Magazine*, vol. 56, no. 8, pp. 33–39, 2018.
- [3] D. Sabella, A. Vaillant, P. Kuure, U. Rauschenbach, and F. Giust, “Mobile-edge computing architecture: The role of mec in the internet of things,” *IEEE Consumer Electronics Magazine*, vol. 5, no. 4, pp. 84–91, 2016.
- [4] T. Zhang, J. Gao, T. Uehara, *et al.*, “Testing location-based function services for mobile applications,” in *Proc. the IEEE Symposium on Service-Oriented System Engineering (SOSE)*, 2015.
- [5] C. H. Liu, X. Ma, X. Gao, and J. Tang, “Distributed energy-efficient multi-uav navigation for long-term communication coverage by deep reinforcement learning,” *IEEE Transactions on Mobile Computing*, vol. 19, no. 6, pp. 1274–1285, 2020.
- [6] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, “Communication-efficient learning of deep networks from decentralized data,” in *Proc. the Artificial Intelligence and Statistics Conference (AISTATS)*, 2017.
- [7] J. Konečný, H. B. McMahan, F. X. Yu, P. Richtárik, A. T. Suresh, and D. Bacon, “Federated learning: Strategies for improving communication efficiency,” *arXiv preprint arXiv:1610.05492*, 2016.
- [8] J. Park, S. Samarakoon, M. Bennis, and M. Debbah, “Wireless network intelligence at the edge,” *Proceedings of the IEEE*, vol. 107, no. 11, pp. 2204–2239, 2019.
- [9] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith, “Federated optimization in heterogeneous networks,” in *Proc. the Machine Learning and Systems (MLSys)*, 2020.
- [10] 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.
- [11] H. Wang, Z. Kaplan, D. Niu, and B. Li, “Optimizing federated learning on non-iid data with reinforcement learning,” in *Proc. the IEEE Conference on Computer Communications (INFOCOM)*, 2020.

- [12] S. Wang, T. Tuor, T. Salonidis, K. K. Leung, C. Makaya, T. He, and K. Chan, "Adaptive federated learning in resource constrained edge computing systems," *IEEE Journal on Selected Areas in Communications*, vol. 37, no. 6, pp. 1205–1221, 2019.
- [13] L. Wang, W. Wang, and B. Li, "Cmfl: Mitigating communication overhead for federated learning," in *Proc. the IEEE International Conference on Distributed Computing Systems (ICDCS)*, 2019.
- [14] Y. Lin, S. Han, H. Mao, Y. Wang, and W. J. Dally, "Deep Gradient Compression: Reducing the communication bandwidth for distributed training," in *Proc. the International Conference on Learning Representations (ICLR)*, 2018.
- [15] F. Seide, H. Fu, J. Droppo, G. Li, and D. Yu, "1-bit stochastic gradient descent and its application to data-parallel distributed training of speech dnns," in *Proc. the Fifteenth Annual Conference of the International Speech Communication Association (Interspeech)*, 2014.
- [16] G. Zhu, Y. Wang, and K. Huang, "Broadband analog aggregation for low-latency federated edge learning," *IEEE Transactions on Wireless Communications*, vol. 19, no. 1, pp. 491–506, 2020.
- [17] H. Wu and P. Wang, "Fast-convergent federated learning with adaptive weighting," *IEEE Transactions on Cognitive Communications and Networking*, vol. 7, no. 4, pp. 1078–1088, 2021.
- [18] T. Nishio and R. Yonetani, "Client selection for federated learning with heterogeneous resources in mobile edge," in *Proc. the IEEE International Conference on Communications (ICC)*, 2019.
- [19] M. M. Amiri, D. Gündüz, S. R. Kulkarni, and H. V. Poor, "Convergence of update aware device scheduling for federated learning at the wireless edge," *IEEE Transactions on Wireless Communications*, vol. 20, no. 6, pp. 3643–3658, 2021.
- [20] Y. J. Cho, J. Wang, and G. Joshi, "Client selection in federated learning: Convergence analysis and power-of-choice selection strategies," *arXiv preprint arXiv:2010.01243*, 2020.
- [21] M. Chen, N. Shlezinger, H. V. Poor, Y. C. Eldar, and S. Cui, "Communication-efficient federated learning," *Proceedings of the National Academy of Sciences*, vol. 118, no. 17, 2021.
- [22] M. Chen, H. V. Poor, W. Saad, and S. Cui, "Convergence time optimization for federated learning over wireless networks," *IEEE Transactions on Wireless Communications*, vol. 20, no. 4, pp. 2457–2471, 2020.
- [23] J. Ren, Y. He, D. Wen, G. Yu, K. Huang, and D. Guo, "Scheduling for cellular federated edge learning with importance and channel awareness," *IEEE Transactions on Wireless Communications*, vol. 19, no. 11, pp. 7690–7703, 2020.
- [24] W. Chen, S. Horvath, and P. Richtarik, "Optimal client sampling for federated learning," *arXiv preprint arXiv:2010.13723*, 2020.
- [25] E. Rizk, S. Vlaski, and A. H. Sayed, "Optimal importance sampling for federated learning," in *Proc. the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021.
- [26] X. Li, K. Huang, W. Yang, S. Wang, and Z. Zhang, "On the convergence of fedavg on non-iid data," in *Proc. the International Conference on Learning Representations (ICLR)*, 2020.
- [27] S. U. Stich, "Local SGD converges fast and communicates little," in *Proc. the International Conference on Learning Representations (ICLR)*, 2019.
- [28] H. T. Nguyen, V. Schwag, S. Hosseinalipour, C. G. Brinton, M. Chiang, and H. V. Poor, "Fast-convergent federated learning," *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 1, pp. 201–218, 2020.
- [29] O. Shamir, N. Srebro, and T. Zhang, "Communication-efficient distributed optimization using an approximate newton-type method," in *Proc. the International Conference on Machine Learning (ICML)*, 2014.
- [30] Y. LeCun and C. Cortes, "MNIST handwritten digit database," 2010.
- [31] A. Krizhevsky, V. Nair, and G. Hinton, "Cifar-10 (canadian institute for advanced research),"