# FedCor: Correlation-Based Active Client Selection Strategy for Heterogeneous Federated Learning

Minxue Tang[1], Xuefei Ning[2], Yitu Wang[1], Jingwei Sun[1], Yu Wang[2], Hai Li[1], Yiran Chen[1]*
[1]Department of Electrical and Computer Engineering, Duke University
[2]Department of Electronic Engineering, Tsinghua University
[1]{minxue.tang,yitu.wang,jingwei.sun,hai.li,yiran.chen}@duke.edu
[2]foxdoraame@gmail.com   [2]yu-wang@tsinghua.edu.cn

## Abstract

*Client-wise data heterogeneity is one of the major issues that hinder effective training in federated learning (FL). Since the data distribution on each client may vary dramatically, the client selection strategy can significantly influence the convergence rate of the FL process. Active client selection strategies are popularly proposed in recent studies. However, they neglect the loss correlations between the clients and achieve only marginal improvement compared to the uniform selection strategy. In this work, we propose FedCor—an FL framework built on a correlation-based client selection strategy, to boost the convergence rate of FL. Specifically, we first model the loss correlations between the clients with a Gaussian Process (GP). Based on the GP model, we derive a client selection strategy with a significant reduction of expected global loss in each round. Besides, we develop an efficient GP training method with a low communication overhead in the FL scenario by utilizing the covariance stationarity. Our experimental results show that compared to the state-of-the-art method, FedCorr can improve the convergence rates by $34\% \sim 99\%$ and $26\% \sim 51\%$ on FMNIST and CIFAR-10, respectively.*

## 1. Introduction

As a newly emerging distributed learning paradigm, federated learning (FL) [9,12,13,17,23] has recently attracted attention because of the offered data privacy. FL aims at dealing with scenarios where training data is distributed across a number of clients. Considering limited communication bandwidth and the privacy requirement, in each communication round, FL usually selects only a fraction of clients, and the selected clients will perform multiple iterations of local updating without exposing their own datasets [23]. This special

---

*Corresponding Author

scenario also introduces other challenges that distinguish FL from the conventional distributed learning [2,35].

One major challenge in FL is the high degree of client-wise data heterogeneity [17], which is the inherent characteristic of a large number of clients. There have been many studies [10,15,16,18,20,25,27,32] trying to tackle non-IID (independent and identically distributed) and unbalanced data of the clients in FL. Most of these studies [10,18,20,32] focus on amending the local model updates or the central aggregation based on FedAvg [23].

Recently, active client selection arises as a complement of the aforementioned studies, aiming at accelerating the convergence of FL with non-IID data. Some recent studies propose to assign higher probability of being selected to the clients with larger training loss value [4,6]. However, they neglect the correlations between the clients and consider their losses independently, which leads to only marginal performance improvement. In this paper, we propose a correlation-based active client selection strategy that can effectively alleviate the accuracy degradation caused by data heterogeneity and significantly boost the convergence of FL. Our key idea is mainly based on the following intuitions:

1. Clients do not contribute **equivalently**. For example, training with a large and balanced dataset on a "good" client can reduce the losses of most clients, while training with a small and extremely biased dataset on a "bad" client may increase the losses of other clients.

2. Clients do not contribute **independently**. The influence of selecting one client depends on the other selected clients because their local updates will be aggregated.

A toy experiment shown in Fig. 1 also illustrates the necessity of considering the correlations for client selection. In this experiment, each client has only one data sample, and thus each data point in the figure represents a client. The task is to select two clients (different markers represent the client selections of different strategies) for training a binary
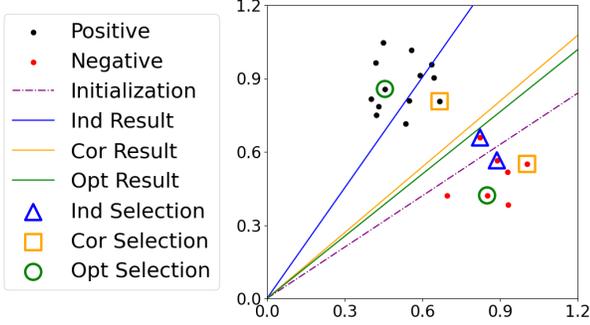
Figure 1. A toy experiment of different client selection strategies.

classifier (shown as the lines). The selection strategy that independently selects two clients with the highest local losses ("Ind Result") fails to reduce the global loss. In contrast, our method considers the correlations between the clients ("Cor Result") and derives a client selection that can achieve an almost lowest global loss ("Opt Result").

Based on the above intuitions, this work proposes Fed-Cor, an FL framework built on a correlation-based client selection strategy, to boost the convergence of FL. Our main contributions are summarized as follows:

1. We model the client loss changes with a Gaussian Process (GP) and propose an interpretable client selection strategy with a significant reduction of the expected global loss in each communication round.

2. We propose a GP training method that utilizes the co-variance stationarity to reduce the communication cost. Experiments show that the GP trained with our method can capture the client correlations well.

3. Experimental results demonstrate that FedCor stabilizes the training convergence and significantly improves the convergence rates by $34\% \sim 99\%$ and $26\% \sim 51\%$ on FMNIST and CIFAR-10, respectively.

## 2. Related Work

An important characteristic of FL [12, 13, 23] is the heterogeneity of clients, which raises new challenges of the training [9, 17, 31]. There are two kinds of heterogeneity in FL: systemic heterogeneity (computation ability, communication bandwidth, etc.) and statistical heterogeneity (non-IID, imbalanced data distribution) [17, 18]. This work mainly focuses on the latter one. A number of methods have been proposed to improve the basic FL algorithm, FedAvg [23], in heterogeneous settings. Some of them manipulate the local training loss like adding regularization terms to stabilized the training [8, 10, 18, 20, 26], while some other works amend the aggregation method to reduce the variance [24, 32].

Complementary to such methods, another way to improve the convergence of FL in non-IID settings is active client selection, which tries to strategically select clients for training in each round in stead of uniformly selecting. Goetz et al. [6] first propose to assign a high selection probability to the clients with large local loss. Cho et al. [4] select $C$ clients with the largest loss among a randomly sampled subset $\mathbb{A} \subseteq \mathbb{U}$ with size $d > C$ to reduce the selection bias. However, neither of them consider the correlations between clients while making the client selection.

## 3. Preliminary

FL seeks for a global model $\boldsymbol{w}$ that achieves the best performance (e.g., the highest classification accuracy) on all $N$ clients. The global loss function in FL is defined as:

$$L(\boldsymbol{w}) = \sum_{k=1}^{N} \frac{|\mathbb{D}_k|}{\sum_j |\mathbb{D}_j|} l(\boldsymbol{w}; \mathbb{D}_k) = \sum_{k=1}^{N} p_k l_k(\boldsymbol{w}), \quad (1)$$

$$l_k(\boldsymbol{w}) = l(\boldsymbol{w}; \mathbb{D}_k) = \frac{1}{|\mathbb{D}_k|} \sum_{\xi \in \mathbb{D}_k} l(\boldsymbol{w}; \xi), \quad (2)$$

where $l(\boldsymbol{w}; \xi)$ is the objective loss of data sample $\xi$ evaluated on model $\boldsymbol{w}$. We refer to $l_k(\boldsymbol{w})$ as the local loss of client $k$, which is evaluated with the local dataset $\mathbb{D}_k$ (of size $|\mathbb{D}_k|$) on client $k$. The weight $p_k = |\mathbb{D}_k| / \sum_j |\mathbb{D}_j|$ of the client $k$ is proportional to the size of its local dataset.

In consideration of the privacy and communication constraints, FL algorithms usually assume partial client participation and perform local model updates. In particular, in communication round $t$, only a subset $\mathbb{K}_t$ with size $|\mathbb{K}_t| = C \leq N$ of the overall client set $\mathbb{U}$ is selected to receive the global model $\boldsymbol{w}^t$ and conduct training with their local dataset for several iterations independently. After the local training, the server collects the trained models from these selected clients and aggregates them (usually by averaging [23]) to produce a new global model $\boldsymbol{w}^{t+1}$. We formulate this procedure as follows:

$$\boldsymbol{w}_k^{t+1} = \boldsymbol{w}^t - \eta_t \tilde{\nabla} l_k(\boldsymbol{w}^t), \quad (3)$$

$$\boldsymbol{w}^{t+1}(\mathbb{K}_t) = \frac{1}{C} \sum_{k \in \mathbb{K}_t} \boldsymbol{w}_k^{t+1} \quad (4)$$

$$= \boldsymbol{w}^t - \frac{\eta_t}{C} \sum_{k \in \mathbb{K}_t} \tilde{\nabla} l_k(\boldsymbol{w}^t), \quad (5)$$

where $\eta_t$ is the learning rate and $\tilde{\nabla} l_k(\boldsymbol{w}^t)$ is the equivalent cumulative gradient [32] in the $t$-th communication round. More specifically, for an arbitrary optimizer on the client $k$, it produces $\Delta \boldsymbol{w}_k^{t,\tau} = -\eta \boldsymbol{d}_k^{t,\tau}$ as the local model update at the $\tau$-th iteration in this round, and the cumulative gradient is calculated as $\tilde{\nabla} l_k(\boldsymbol{w}^t) = \sum_\tau \boldsymbol{d}_k^{t,\tau}$.

## 4. Methodology

In this section, we elaborate our proposed method, i.e., FedCor, that can effectively boost the convergence of FL.

We first formulate our goal of accelerating the convergence of FL as optimization problems that maximize the posterior expectation of loss decrease in Sec. 4.1. Then, Sec. 4.2 demonstrates empirical evidence that the prior distribution of loss changes in each communication round can be modeled as Gaussian Processes (GP). Based on this observation, we utilize GP to solve the optimization problems and obtain an effective client selection strategy for heterogeneous FL in Sec. 4.3. We further analyze the selection criterion of our client selection strategy and give out its intuitive interpretation in Sec. 4.4. Finally, in Sec. 4.5, we describe how we train the GP parameters in communication-constrained FL.

## 4.1. Problem Formulation

To achieve a fast convergence, we hope to find the client selection strategy which can lead to the maximal global loss decrease after each communication round. Accordingly, we define our target as solving a series of optimization problems, one for each communication round $t$:

$$\min_{\mathbb{K}_t} \quad \Delta L^t(\mathbb{K}_t) = L(\boldsymbol{w}^{t+1}(\mathbb{K}_t)) - L(\boldsymbol{w}^t)$$
$$\text{subject to} \quad \boldsymbol{w}^{t+1}(\mathbb{K}_t) = \boldsymbol{w}^t - \frac{\eta_t}{C} \sum_{k \in \mathbb{K}_t} \tilde{\nabla} l_k(\boldsymbol{w}^t). \quad (6)$$

It is impractical in FL to search for the best client selection with multiple trials of different client selections since it introduces large communication and computation overhead. Therefore, we need an efficient way to predict the global loss decreases for different client selections and make a decision with very limited trials. To achieve this goal, we first reformulate the optimization problem in Eq. (6) with the following lemma. The proof of this lemma is in Appendix A.2.

**Lemma 1.** *The optimization problem in Eq. (6) is approximately equivalent to the following probabilistic form.*

$$\min_{\mathbb{K}_t} \quad \mathbb{E}_{\Delta \mathbf{l}^t | \Delta l_{\mathbb{K}_t}^t(\mathbb{K}_t)} \Big[ \sum_i p_i \Delta l_i^t \Big] = \sum_i p_i \tilde{\mu}_i^t(\Delta \boldsymbol{l}_{\mathbb{K}_t}^t(\mathbb{K}_t)),$$
$$(7)$$

*where $\Delta \mathbf{l}^t = [\Delta l_1^t, \cdots, \Delta l_N^t]$ is the loss changes of all clients in round t, which is a random variable w.r.t random client selection in round t. $\tilde{\boldsymbol{\mu}}^t(\Delta \boldsymbol{l}_{\mathbb{K}_t}^t(\mathbb{K}_t))$ is the posterior mean of $\Delta \mathbf{l}^t$ conditioned on $\Delta \boldsymbol{l}_{\mathbb{K}_t}^t(\mathbb{K}_t) = [\Delta l_i^t(\mathbb{K}_t)]_{i \in \mathbb{K}_t}$.*

The reformulated objective in Eq. (7) tells that if we can predict the loss changes of those clients selected for training $(\Delta \boldsymbol{l}_{\mathbb{K}_t}^t(\mathbb{K}_t))$, we can predict the global loss change with its posterior mean and make decision according to it. Now what we need is a probabilistic model of the loss changes $\Delta \mathbf{l}^t$ to make the prediction and calculate the posterior.

## 4.2. Modeling Loss Changes with GP

It is a common practice to assume a GP prior over an unknown objective function in Bayesian Optimization [3,30].
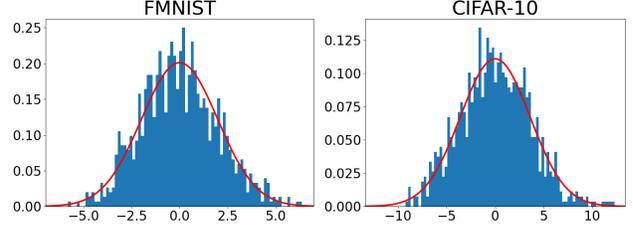


FMNIST          CIFAR-10

Figure 2. Histograms of the first principle component in Non-IID FL [23]. More details and full results can be found in Appendix D.2.

Our preliminary investigation (partly) shown in Fig. 2 also indicates that the prior distribution of the loss changes in one communication round follow a GP. Specifically, we randomly sample a number of client selections and perform one round of training to get samples of the loss changes. Then, we conduct PCA on these loss change samples and plot histograms of the first several principle components. The red line in the Fig. 2 is the Gaussian PDF with the sample mean and sample variance. And we can see that this Gaussian distribution can approximate the distribution of the samples well. A mathematical explanation of this observation is also given out in Appendix A.1.

Accordingly, we propose to model the loss changes in one communication round $t$ with a GP prior as follows:

$$\Delta \mathbf{l}^t = [\Delta l_1^t, \cdots, \Delta l_N^t] \sim \mathcal{N}(\Delta \boldsymbol{l}^t; \boldsymbol{\mu}^t, \boldsymbol{\Sigma}^t). \quad (8)$$

**Remark.** In order to efficiently learn the covariance in FL, rather than directly working with the covariance matrix, we embed all clients into a continuous vector space and use a kernel function to calculate the covariance (see Sec. 4.5). Thus, we still use the term GP instead of Multivariate Gaussian Distribution, though the dimension of $\Delta \mathbf{l}^t$ is finite.

A good property of GP is that we can get a closed form of the posterior expectation in Eq. (7), which makes our client selection strategy interpretable. In the next sections, we will propose our client selection strategy based on the GP model, and then give an interpretation of it. We leave the training method for the parameters $(\boldsymbol{\mu}^t, \boldsymbol{\Sigma}^t)$ in GP to Sec. 4.5.

## 4.3. Client Selection Strategy

While we have get the probabilistic model to calculate the posterior expectation, it is still not determined how to predict the loss changes of the clients selected for training, namely $\Delta \boldsymbol{l}_{\mathbb{K}_t}^t(\mathbb{K}_t)$. Inspired by UCB methods [1,5,28], we develop an iterative method that predict the loss change and select one client in each iteration, as shown in Algorithm 1. There are three steps in one iteration:

**(i) Prediction.** In each iteration, we first make an prediction $\Delta \hat{l}_k^t$ for each client $k$ if it is selected. Generally, the selected client would have a large loss decrease since it directly participate in the model update. Thus, we propose to use the

**Algorithm 1** Client Selection Strategy with GP

**Require:** $\boldsymbol{\mu}^t$ and $\boldsymbol{\Sigma}^t$ of the GP, scale factor $\boldsymbol{\alpha}^t$
**Ensure:** Client Selection $\mathbb{K}_t$
1: Initialize $\mathbb{K}_t \leftarrow \emptyset, \mathbb{P} \leftarrow \mathbb{U}$.
2: **while** $|\mathbb{K}_t| < C$ **do**
3:     **for** each client $k \in \mathbb{P}$ **do**
4:         Predict its loss change if select it: $\Delta \hat{l}_k^t = \mu_k^t - \alpha_k^t \sigma_k^t$.
5:         Calculate the posterior mean of the loss changes $\tilde{\boldsymbol{\mu}}^t(\Delta \hat{l}_k^t)$.
6:     **end for**
7:     Select the client by $k^* = \arg\min_k \sum_i p_i \tilde{\mu}_i^t(\Delta \hat{l}_k^t)$.
8:     Add $k^*$ into $\mathbb{K}_t$ and remove it from $\mathbb{P}$.
9:     $\boldsymbol{\mu}^t \leftarrow \tilde{\boldsymbol{\mu}}^t(\Delta \hat{l}_{k^*}^t), \boldsymbol{\Sigma}^t \leftarrow \tilde{\boldsymbol{\Sigma}}^t(\Delta \hat{l}_{k^*}^t)$.
10: **end while**

---

**Algorithm 2** FedCor

1: Initialize $\boldsymbol{X}_0$ and Global Model $\boldsymbol{w}_0$.
2: **for** each round $t = 0, 1, ...$ **do**
3:     **if** $t \% \Delta t == 0$ **then**
4:         Uniformly sample $S$ client selections $\mathbb{S}_{t,i}, i = 1, 2, ..., S$.
5:         **for** $i = 1, 2, ..., S$ **do**
6:             $\boldsymbol{w}^{t+1}(\mathbb{S}_{t,i}) \leftarrow \boldsymbol{w}^t - \frac{\eta_t}{C}\sum_{k \in \mathbb{S}_{t,i}} \tilde{\nabla} l_k(\boldsymbol{w}^t)$.
7:             Collect $\Delta \boldsymbol{l}^t(\mathbb{S}_{t,i}) \leftarrow \boldsymbol{l}(\boldsymbol{w}^{t+1}(\mathbb{S}_{t,i})) - \boldsymbol{l}(\boldsymbol{w}^t)$.
8:         **end for**
9:         Reset $\alpha_k \leftarrow 1, \forall k \in \mathbb{U}$.
10:     **end if**
11:     Update $\boldsymbol{X}_t$ with Eq. (16).
12:     Select clients $\mathbb{K}_t$ with Algorithm 1 ($\boldsymbol{\mu}^t = \boldsymbol{0}, \boldsymbol{\Sigma}^t = \boldsymbol{X}^{t^T}\boldsymbol{X}^t, \boldsymbol{\alpha}^t = \boldsymbol{\alpha}$).
13:     $\boldsymbol{w}^{t+1} \leftarrow \boldsymbol{w}^{t+1}(\mathbb{K}_t) = \boldsymbol{w}^t - \frac{\eta_t}{C}\sum_{k \in \mathbb{K}_t}\tilde{\nabla} l_k(\boldsymbol{w}^t)$.
14:     Update $\boldsymbol{\alpha}_{\mathbb{K}_t} \leftarrow \beta \boldsymbol{\alpha}_{\mathbb{K}_t}$.
15: **end for**

---

lower confidence bound as the prediction:

$$\Delta \hat{l}_k^t = \mu_k^t - \alpha_k^t \sigma_k^t; \quad \alpha_k^t = a\beta^{\tau_k^t}, \tag{9}$$

where $\sigma_k^t = \sqrt{\Sigma_{k,k}^t}$, and $a$ is a scale constant. $\beta \in (0,1)$ is an annealing coefficient, and its index $\tau_k^t$ denotes how many times client $k$ has been selected. We will discuss this annealing coefficient more in Sec. 4.5.

**(ii) Selection.** The client $k^*$ is selected to minimize the posterior expectation of the overall loss conditioned on its loss change prediction made in the last step:

$$k^* = \arg\min_k \sum_i p_i \tilde{\mu}_i^t(\Delta \hat{l}_k^t) \tag{10}$$

**(iii) Posterior.** After selecting the client $k^*$, we update the GP for the next iteration with the posterior conditioned on the loss change prediction of $k^*$:

$$\boldsymbol{\mu}^t \leftarrow \tilde{\boldsymbol{\mu}}^t(\Delta \hat{l}_{k^*}^t), \quad \boldsymbol{\Sigma}^t \leftarrow \tilde{\boldsymbol{\Sigma}}^t(\Delta \hat{l}_{k^*}^t). \tag{11}$$

By updating the GP with its posterior, we iteratively add conditions into the probabilistic model to approach the fully conditioned distribution $p(\Delta \boldsymbol{l}^t | \Delta \boldsymbol{l}_{\mathbb{K}_t}^t(\mathbb{K}_t))$, and make the next prediction of the loss change more accurate.

There are some similarities between our method and traditional Bayesian Optimization: Using GP as a prior of the objective function, and using UCB as well as posterior distribution for iterative selection [3, 5, 28]. However, there is a key difference: In each communication round, we determine the client selection with only predictions instead of measurements of the global loss changes, while traditional Bayesian Optimization requires a sequence of measurements as new information to make decisions. The measurements of global loss changes will introduce large communication overhead and are unfeasible in FL.

## 4.4. Insights into Our Selection Strategy

In this section, we give an intuitive interpretation of our selection strategy and show the benefits of it within a simple case. A more detailed analysis of the selection criterion and convergence of FedCor can be found in Appendix B.

For simplicity, we omit all superscript $t$ in this section. Lemma 2 gives the selection criterion of FedCor in a simple case where we only select two clients, and the proof can be found in Appendix A.3.

**Lemma 2.** *The selection criterion of FedCor when selecting two clients $k_1$ and $k_2$ can be written as*

$$k_1 = \arg\max_k \quad \beta^{\tau_k} \sum_i p_i \sigma_i r_{ik}, \tag{12}$$

$$k_2 = \arg\max_{k'} \quad \frac{\beta^{\tau_{k'}}\left[\overbrace{\sum_i p_i \sigma_i r_{ik'}}^{(A)} - r_{k_1 k'}\overbrace{\sum_i p_i \sigma_i r_{ik_1}}^{(B)}\right]}{\sqrt{1 - r_{k'k_1}^2}}, \tag{13}$$

*where $r_{ij} = \Sigma_{i,j}/\sigma_i \sigma_j$ is the Pearson correlation coefficient.*

**(i) Single-Iteration.** Eq. (12) has a clear interpretation to select the client that has large correlations with other clients ($r_{ik}$), so that other clients can benefit more from training on the selected client. Our selection criterion takes the correlations between the clients into consideration, and can conduct better selection compared with those algorithms that only consider the loss of each client independently [4, 6].

**(ii) Multi-Iteration.** In Eq. (13), term (A) and (B) are the single-iteration selection criterion in Eq. (12) of client $k'$ and $k_1$, respectively. Since we have maximized (B) when

selecting client $k_1$, term (B) is usually positive. Therefore, the selection of $k'$ does not only consider its correlations with other clients ($r_{ik'}$), but also prefers the clients that have small correlations $r_{k_1 k'}$ with the previous selected client $k_1$. This criterion penalizes selection redundancy and leads to a client selection with diverse data, which reduces the variance and makes the training process more stable. Since clients with similar data generate similar local updates, selecting redundant clients only brings marginal gains to the global performance or would even drive the optimization into bad local optimum. This selection preference is also demonstrated in Fig. 1, where FedCor chooses one positive and one negative point as the optimal selection does.

## 4.5. Training GP in FL

As a classical machine learning model, GP has been widely discussed and well studied [33]. There have been many methods to train the parameters in GP, namely, the covariance $\mathbf{\Sigma}^t$ in Eq. (8) [1]. Nevertheless, to make the GP training feasible in the communication-constrained FL procedure, we should revise the GP training method to reduce the number of samples and better utilize historical information.

In GP, a kernel function $K(\boldsymbol{x}_i, \boldsymbol{x}_j)$ is used to calculate the covariance [33] as $\Sigma^t_{i,j} = K(\boldsymbol{x}^t_i, \boldsymbol{x}^t_j)$, where $\boldsymbol{x}^t_i, \boldsymbol{x}^t_j$ are the features of the data points $i$ and $j$, respectively. Following this, we assign a trainable embedding in a latent space to each client. The embedding of the $k$-th client is noted as $\boldsymbol{x}^t_k \in \mathbb{R}^d$ ($d < N$), and we choose the kernel function as

$$K(\boldsymbol{x}^t_i, \boldsymbol{x}^t_j) = \boldsymbol{x}^{t^T}_i \boldsymbol{x}^t_j, \qquad (14)$$

which is a homogeneous linear kernel [33]. This low-rank formulation reduces the number of parameters we need to learn, thus making the GP training more data-efficient.

A commonly used GP training method is maximum likelihood evaluation, where we uniformly sample $S$ client selection $\{\mathbb{S}_{t,i} : i = 1, \cdots, S\}$, and maximize the likelihood of the corresponding loss changes $\{\Delta \boldsymbol{l}^t(\mathbb{S}_{t,i}) : i = 1, \cdots, S\}$ to learn the embedding matrix $\boldsymbol{X}^t = [\boldsymbol{x}^t_1, \cdots, \boldsymbol{x}^t_N]$:

$$\boldsymbol{X}^t = \arg\max_{\boldsymbol{X}} \sum_{i=1}^{S} \log p(\Delta \boldsymbol{l}^t(\mathbb{S}_{t,i})|\boldsymbol{X}). \qquad (15)$$

However, to collect each sample $\Delta \boldsymbol{l}^t(\mathbb{S}_{t,i})$, we have to broadcast $\boldsymbol{w}^{t+1}(\mathbb{S}_{t,i})$ to all the clients. And since a large $S$ is usually required for an unbiased estimation in each communication round $t$, the vanilla training procedure in Eq. (15) introduces a high communication overhead.

Actually, the correlations between loss changes of different clients mainly arise from similarities between their datasets, which are invariant during the FL process. Thus,

---

[1]We do not train $\boldsymbol{\mu}$ and set it to $\boldsymbol{0}$, since it does not affect the selection strategy as we can see in Lemma 2 and Appendix B.
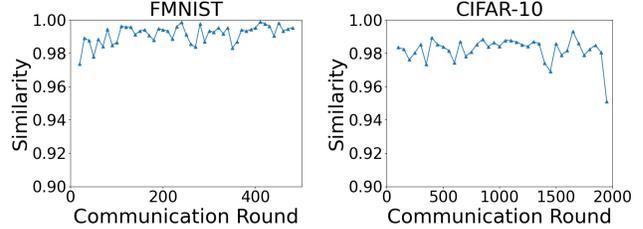


Figure 3. Covariance Stationarity in Non-IID FL [23]. Full experiment results and more details can be found in Appendix D.3.

we hypothesise that the covariance also changes slowly in the concerned time range. To verify this, we use a large number of samples to evaluate the covariance $\mathbf{\Sigma}^t$ in each communication round, and calculate the cosine similarity between $\mathbf{\Sigma}^t$ and $\mathbf{\Sigma}^{t+\Delta t}$. We set $\Delta t = 10$ for FMNIST and $\Delta t = 50$ for CIFAR-10. As shown in Fig. 3, we can see that the similarity keeps very high ($> 0.97$ for FMNIST and $> 0.95$ for CIFAR-10) during the whole FL training process.

Accordingly, we do not need to update $\boldsymbol{X}^t$ in every round but inherit the embedding matrix $\boldsymbol{X}^{t-1}$ from the last round and train it only every $\Delta t$ rounds. Furthermore, we can reuse historical samples for GP training to reduce the number of samples $S$ that we need to collect in each GP training round. We summarize our update rule of $\boldsymbol{X}^t$ as follows:

$$\boldsymbol{X}^t = \begin{cases} \boldsymbol{X}^{t-1}, & t\%\Delta t \neq 0; \\ \arg\max_{\boldsymbol{X}} \Phi_t(\boldsymbol{X}), & t\%\Delta t = 0, \end{cases} \qquad (16)$$

where

$$\Phi_t(\boldsymbol{X}) = \sum_{m=0}^{M} \sum_{i=1}^{S} \gamma^m \log p(\Delta \boldsymbol{l}^{t-m\Delta t}(\mathbb{S}_{t-m\Delta t,i})|\boldsymbol{X}). \qquad (17)$$

$M$ is the number of reused historical samples, and $\gamma < 1$ is the discount factor to weight the historical samples. Our method is able to reduce the communication overhead with a large $\Delta t$ and $S = 1$, while guaranteeing the performance.

As we only update the covariance $\mathbf{\Sigma}$ every $\Delta t$ rounds, the annealing factor $\beta^{\tau_k}$ can prevent us from making the same selection during the $\Delta t$ rounds. Repeatedly training with the same group of clients would cause the global model to overfit on their data, which may hinder the convergence of FL. In practice, we reset $\tau_k$ to 0 after each GP training round to achieve the fastest convergence while avoiding overfitting on some clients.

We summarize our overall framework FedCor in Algorithm 2. It is noteworthy that our method is orthogonal to existing FL optimizers that amend the training loss or the aggregation scheme, e.g., FedAvg [23] and FedProx [18]. So our method can be combined with any of them.
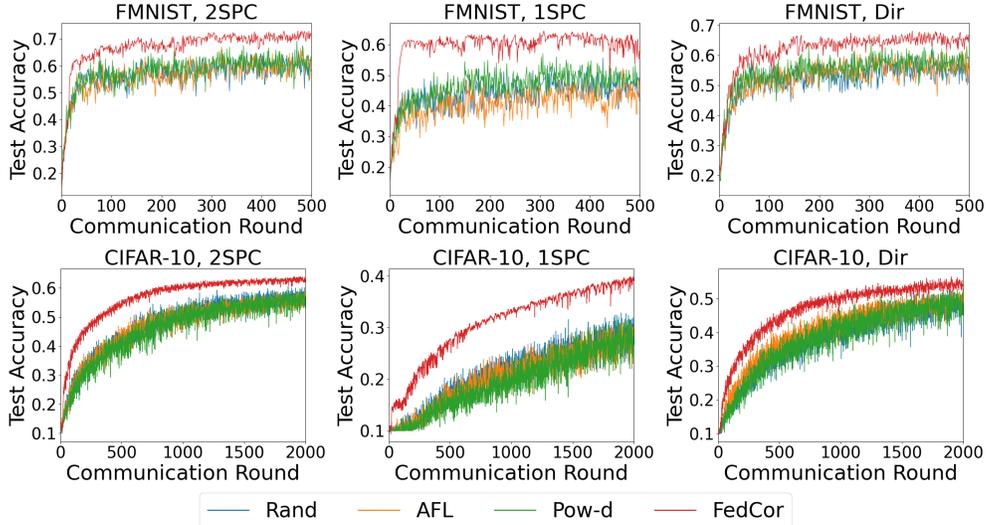
Figure 4. Test accuracy on FMNIST and CIFAR-10 under three heterogeneous settings (2SPC, 1SPC and Dir). All experiments in one figure share the same hyperparameters except for the client selection strategy.

| Method | FMNIST | | | CIFAR-10 | | |
|---|---|---|---|---|---|---|
| | 2SPC(69%) | 1SPC(62%) | Dir(64%) | 2SPC(62%) | 1SPC(36%) | Dir(54%) |
| Rand | $295.8 \pm 92.0$ | N/A | $141.0 \pm 73.0$ | $1561.2 \pm 236.2$ | $1750.4 \pm 190.3$ | N/A |
| AFL | $218.6 \pm 117.3$ | N/A | $169.0 \pm 166.1$ | N/A | $1845.2 \pm 28.8$ | $1524.4 \pm 267.9$ |
| Pow-d | $126.6 \pm 78.2$ | $167.2 \pm 72.3$ | $123.0 \pm 101.0$ | $1558.2 \pm 227.0$ | $1752.2 \pm 186.2$ | $1355.2 \pm 151.3$ |
| FedCor (Ours) | $\mathbf{94.8 \pm 18.4}$ | $\mathbf{84.0 \pm 53.1}$ | $\mathbf{68.8 \pm 27.5}$ | $\mathbf{1033.4 \pm 123.7}$ | $\mathbf{1269.2 \pm 70.6}$ | $\mathbf{1076.8 \pm 262.8}$ |

Table 1. The number of communication rounds for each selection strategy to achieve target test accuracies (specified in parentheses) under three heterogeneous settings (2SPC, 1SPC and Dir). The results consist of the mean and the standard deviation over 5 random seeds. N/A means that the corresponding selection strategy cannot achieve the target accuracy with some random seeds within the maximal number of communication rounds (500 for FMNIST and 2000 for CIFAR-10).

## 5. Experiments

### 5.1. Experiment Settings

We conduct experiments on two datasets, FMNIST [34] and CIFAR-10 [14]. For FMNIST, we adopt an MLP model with two hidden layers, and this model achieves an accuracy of $85.92\%$ with centralized training. For CIFAR-10, we adopt a CNN model with three convolutional layers followed by one fully connected layer, and this model can achieve an accuracy of $73.84\%$ with centralized training. More details on the model construction and training hyperparameters can be found in Appendix C.1. For each dataset, we experiment with three different heterogeneous data partitions on $N = 100$ clients as follows.

**(i) 2 shards per client (2SPC)**: This setting is the same as the non-IID setting in [23]. We sort the data by their labels and divide them into 200 shards so that all the data in one shard share the same label. We randomly allocate these shards to clients, and each client has two shards. Since all the shards have the same size, the data partition is balanced.

That is to say, all the clients have the same dataset size. We select $C = 5$ clients in each round within this setting.

**(ii) 1 shard per client (1SPC)**: This setting is similar to the 2SPC setting, and the only difference is that each client only has one shard, i.e., each client only has the data of one label. This is the data partition with the highest heterogeneity, and it is also balanced. We select $C = 10$ clients in each round within this setting.

**(iii) Dirichlet Distribution with $\alpha = 0.2$ (Dir)**: We inherit and slightly change the setting from [7] to create an unbalanced data partition. We sample the ratio of the data with each label on one client from a Dirichlet Distribution parameterized by the concentration parameter $\alpha = 0.2$. More details can be found in the Appendix C.2. We select $C = 5$ clients in each round within this setting.

We divide the training process of FedCor into two phases: (i) Warm-up phase: We uniformly sample client selection $\mathbb{K}_t$ and collect the loss values of all the clients in $\mathbb{U}$ to train the GP in each round, i.e., $\Delta t = 1$ and $S = 1$. We set the length of the warm-up phase to 15 for FMNIST and 20 for CIFAR-10. (ii) Normal phase: After the warm-up phase, we
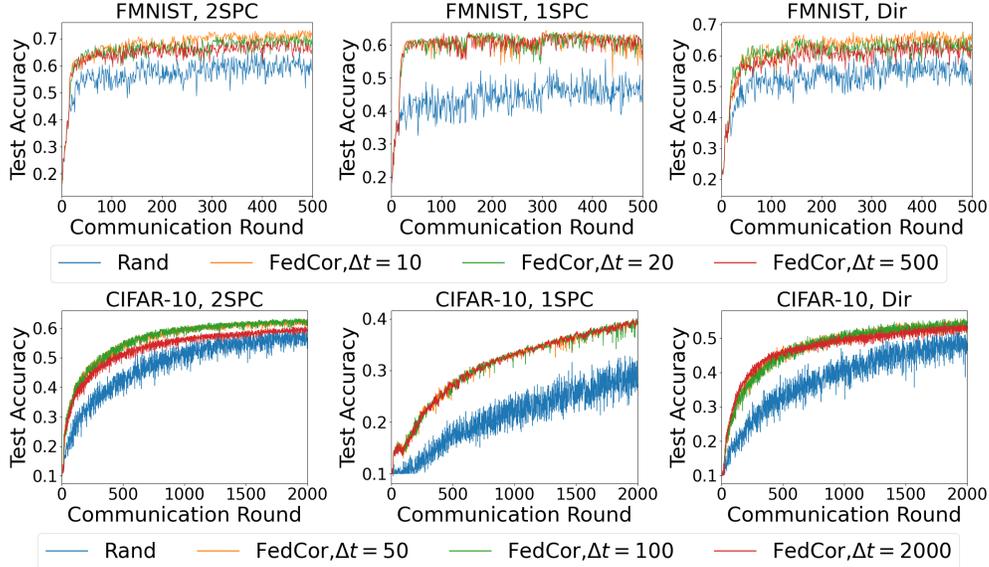
Figure 5. Test accuracy with different GP training interval $\Delta t$ on FMNIST and CIFAR-10 under 2SPC, 1SPC and Dir.

follow Algorithm 2 to select clients and update the GP.

In all the experiments, we use FedAvg [23] as the FL optimizer. We present the average results using five random seeds in all experiments. We will first show that our method can achieve faster and more stable convergence, compared with three baselines: random selection (Rand), Active FL (AFL) [6] and Power-of-choice Selection Strategy (Pow-d) [4]. Then, we will give ablation studies on the GP training interval $\Delta t$ as well as the annealing coefficient $\beta$. Finally, we visualize the client embeddings $X$ with t-SNE [21] and show that FedCor can effectively capture the correlations.

### 5.2. Convergence under Heterogeneous Settings

We compare the convergence rate of our method FedCor with the other baselines on both FMNIST and CIFAR-10, and demonstrate the results in Figure 4. We set the GP update interval $\Delta t = 10$ and the annealing coefficient $\beta = 0.95$ for FMNIST experiments, and $\Delta t = 50$ and $\beta = 0.9$ for CIFAR-10 experiments.

As shown in Figure 4, FedCor achieves the highest test accuracy and the fastest convergence in all experiments. While other active client selection strategies show only slight or even no superiority compared with the fully random strategy, our method clearly outperforms all baselines, especially under the extremely heterogeneous setting when data on each client contains only one label (1SPC). Furthermore, the learning curves of FedCor are more smooth and less noisy than those of other methods, meaning that FedCor reduces the variance and makes the federated optimization more stable.

Table 1 shows the numbers of communication rounds for each selection strategy to achieve a specified test accuracy. We can see that FedCor achieves the specified accuracy

$34\% \sim 99\%$ and $26\% \sim 51\%$ faster than Pow-d on FMNIST and CIFAR-10, respectively.

### 5.3. Results with Larger GP Training Interval

Collecting training data in the GP update rounds brings communication overhead, since we need to broadcast the model to all the clients. Thus, it is important to investigate the minimal GP update frequency. We vary the GP training interval and show the accuracy curves in Figure 5. We set $\Delta t = 10, 20, 500$ with $\beta = 0.95, 0.95, 0.99$ for the experiments on FMNIST, and $\Delta t = 50, 100, 2000$ with $\beta = 0.97, 0.97, 0.999$ for the experiments on CIFAR-10, respectively. As shown in the figures, the performance degrades very slightly with larger training intervals. It is noteworthy that even if we do not update the GP model after the warm-up phase (noted as $\Delta t = 500$ for FMNIST, and $\Delta t = 2000$ for CIFAR-10), FedCor still achieves faster convergence than the random selection strategy. These results indicate that the correlations learned by the GP model are stable, which supports our assumption in Section 4.5. In a word, one can largely reduce the communication overhead by training the GP model with a very low frequency while guaranteeing the convergence rate and accuracy under the communication-bounded FL setting.

### 5.4. Influence of Annealing Coefficient

We also conduct experiments with different annealing coefficient $\beta$ that controls how "concentrated" the client selection is. We perform FedCor with $\Delta t = 10$ and $\beta = 0.5, 0.75, 0.9$ for FMNIST, and $\Delta t = 50, \beta = 0.9, 0.95, 0.99$ for CIFAR-10. The learning curves as well as the client selection frequencies under 2SPC setting are
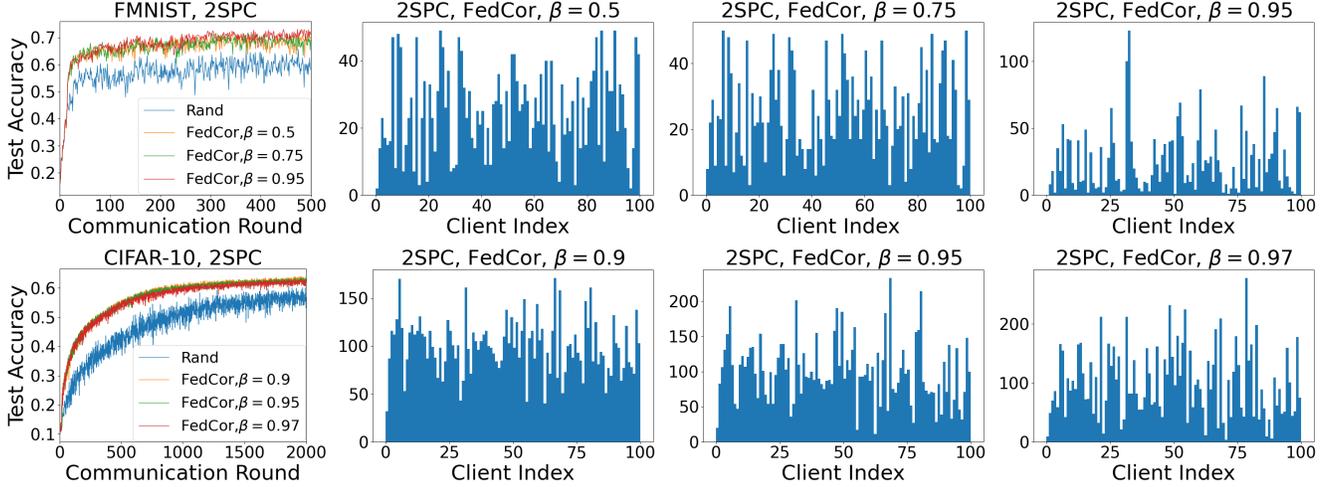
Figure 6. Test accuracy and client selection frequency with different annealing coefficient $\beta$ on FMNIST and CIFAR-10 under the 2SPC setting. The frequency is represented as the number of times each client is selected during the whole training process.
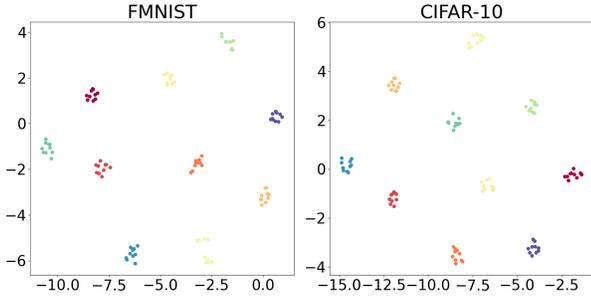


Figure 7. Visualization of client embedding under the 1SPC setting.

shown in Fig. 6, and we leave the full results under the 1SPC and Dir settings to Appendix D.1. We observe that when using a smaller $\beta$, the overall client selections appear to be more "uniform", while the learning curves are almost invariant. Notice that this does not mean that FedCor with small $\beta$ is equivalent to uniform sampling, instead, FedCor still achieves consistent improvements compared to uniform sampling. And Sec. 4.4 havs discussed the reason: FedCor not only considers the benefit that each client brings to the federation, but also considers the correlations among the clients to select the best group of clients. The experimental results here show that it is more important to select a good "group" of clients than just good individuals.

## 5.5. Visualization of Client Embedding

To obtain an insight into the correlations learned by the GP model, we show the t-SNE [21] plot of the client embeddings learned in the warm-up phase under the 1SPC setting. In Fig. 7, each embedding is labeled with the only data label on the corresponding client. We normalize the length of embedding vectors to 1 so that the distance between two

embeddings can reveal the correlation. We can see that the embeddings of clients with the same label are clustered together, which demonstrates that FedCor has captured the correlations between clients correctly in the warm-up phase.

## 6. Conclusion and Future Work

This work proposes FedCor, an FL framework with a novel client selection strategy for heterogeneous settings. FedCor is based on the intuition that it is crucial to utilize the correlations between clients to achieve a faster and more stable convergence in heterogeneous FL. Specifically, we model the client correlations with a GP, and design an effective and interpretable client selection strategy based on it. We also develop a efficient method to train the GP with a low communication overhead. Experimental results on FMNIST and CIFAR-10 show that FedCor effectively accelerates and stabilizes the training process under highly heterogeneous settings. In addition, we verify that FedCor captures the client correlation correctly using only the loss information. How to extend FedCor to the other tasks and further utilize the captured correlations is an interesting direction for future work. Besides, our method focuses on the cross-silo federated learning scenario [9], and how to extend it to the cross-device scenario is a meaningful topic.

## Acknowledgement

# References

[1] Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine Learning Research*, 3(Nov):397–422, 2002. 3

[2] Stephen Boyd, Neal Parikh, and Eric Chu. *Distributed optimization and statistical learning via the alternating direction method of multipliers*. Now Publishers Inc, 2011. 1

[3] Eric Brochu, Vlad M Cora, and Nando De Freitas. A tutorial on bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning. *arXiv preprint arXiv:1012.2599*, 2010. 3, 4

[4] Yae Jee Cho, Jianyu Wang, and Gauri Joshi. Client selection in federated learning: Convergence analysis and power-of-choice selection strategies. *arXiv preprint arXiv:2010.01243*, 2020. 1, 2, 4, 7, 14, 17, 18

[5] Dennis D Cox and Susan John. A statistical method for global optimization. In *[Proceedings] 1992 IEEE International Conference on Systems, Man, and Cybernetics*, pages 1241–1246. IEEE, 1992. 3, 4

[6] Jack Goetz, Kshitiz Malik, Duc Bui, Seungwhan Moon, Honglei Liu, and Anuj Kumar. Active federated learning. *arXiv preprint arXiv:1909.12641*, 2019. 1, 2, 4, 7, 18

[7] Tzu-Ming Harry Hsu, Hang Qi, and Matthew Brown. Measuring the effects of non-identical data distribution for federated visual classification. *arXiv preprint arXiv:1909.06335*, 2019. 6, 18

[8] Tzu-Ming Harry Hsu, Hang Qi, and Matthew Brown. Federated visual classification with real-world data distribution. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part X 16*, pages 76–92. Springer, 2020. 2

[9] Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Keith Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Advances and open problems in federated learning. *arXiv preprint arXiv:1912.04977*, 2019. 1, 2, 8

[10] Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank J Reddi, Sebastian U Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for on-device federated learning. *arXiv preprint arXiv:1910.06378*, 2019. 1, 2

[11] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 18

[12] Jakub Konečnỳ, Brendan McMahan, and Daniel Ramage. Federated optimization: Distributed optimization beyond the datacenter. *arXiv preprint arXiv:1511.03575*, 2015. 1, 2

[13] Jakub Konečnỳ, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*, 2016. 1, 2

[14] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 6

[15] Ang Li, Jingwei Sun, Pengcheng Li, Yu Pu, Hai Li, and Yiran Chen. Hermes: an efficient federated learning framework for heterogeneous mobile clients. In *Proceedings of the 27th Annual International Conference on Mobile Computing and Networking*, pages 420–437, 2021. 1

[16] Ang Li, Jingwei Sun, Binghui Wang, Lin Duan, Sicheng Li, Yiran Chen, and Hai Li. Lotteryfl: Personalized and communication-efficient federated learning with lottery ticket hypothesis on non-iid datasets. *arXiv preprint arXiv:2008.03371*, 2020. 1

[17] Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia Smith. Federated learning: Challenges, methods, and future directions. *IEEE Signal Processing Magazine*, 37(3):50–60, 2020. 1, 2

[18] 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. 1, 2, 5

[19] Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence of fedavg on non-iid data. *arXiv preprint arXiv:1907.02189*, 2019. 14, 15, 17

[20] Xianfeng Liang, Shuheng Shen, Jingchang Liu, Zhen Pan, Enhong Chen, and Yifei Cheng. Variance reduced local sgd with lower communication complexity. *arXiv preprint arXiv:1912.12844*, 2019. 1, 2

[21] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605, 2008. 7, 8

[22] Stephan Mandt, Matthew Hoffman, and David Blei. A variational analysis of stochastic gradient algorithms. In *International conference on machine learning*, pages 354–363, 2016. 11

[23] 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*, pages 1273–1282. PMLR, 2017. 1, 2, 3, 5, 6, 7, 18

[24] Tomoya Murata and Taiji Suzuki. Bias-variance reduced local sgd for less heterogeneous federated learning. *arXiv preprint arXiv:2102.03198*, 2021. 2

[25] Amirhossein Reisizadeh, Farzan Farnia, Ramtin Pedarsani, and Ali Jadbabaie. Robust federated learning: The case of affine distribution shifts. *arXiv preprint arXiv:2006.08907*, 2020. 1

[26] Neta Shoham, Tomer Avidor, Aviv Keren, Nadav Israel, Daniel Benditkis, Liron Mor-Yosef, and Itai Zeitak. Overcoming forgetting in federated learning on non-iid data. *arXiv preprint arXiv:1910.07796*, 2019. 2

[27] Virginia Smith, Chao-Kai Chiang, Maziar Sanjabi, and Ameet S Talwalkar. Federated multi-task learning. In *Advances in neural information processing systems*, pages 4424–4434, 2017. 1

[28] Niranjan Srinivas, Andreas Krause, Sham M Kakade, and Matthias W Seeger. Information-theoretic regret bounds for gaussian process optimization in the bandit setting. *IEEE Transactions on Information Theory*, 58(5):3250–3265, 2012. 3, 4

[29] TensorFlow team. Tensorflow convolutional neural networks tutorial. https://www.tensorflow.org/tutorials/images/cnn, 2016. 17

[30] Ngo Anh Vien, Heiko Zimmermann, and Marc Toussaint. Bayesian functional optimization. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018. 3

[31] Jianyu Wang, Zachary Charles, Zheng Xu, Gauri Joshi, H Brendan McMahan, Maruan Al-Shedivat, Galen Andrew, Salman Avestimehr, Katharine Daly, Deepesh Data, et al. A field guide to federated optimization. *arXiv preprint arXiv:2107.06917*, 2021. 2

[32] Jianyu Wang, Qinghua Liu, Hao Liang, Gauri Joshi, and H Vincent Poor. Tackling the objective inconsistency problem in heterogeneous federated optimization. *arXiv preprint arXiv:2007.07481*, 2020. 1, 2

[33] Christopher KI Williams and Carl Edward Rasmussen. *Gaussian processes for machine learning*, volume 2. MIT press Cambridge, MA, 2006. 5

[34] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017. 6

[35] Tao Yang, Xinlei Yi, Junfeng Wu, Ye Yuan, Di Wu, Ziyang Meng, Yiguang Hong, Hong Wang, Zongli Lin, and Karl H Johansson. A survey of distributed optimization. *Annual Reviews in Control*, 47:278–305, 2019. 1

## A. Theoretical Analysis

### A.1. Analytical Insight into Gaussian Processes

In this section, we give a mathematical explanation about why the loss changes obey Gaussian Distributions. Our analysis based on the following assumption where we assume that the global weight update in one communication round follow a Gaussian Distribution under uniformly client selection.

**Assumption 1.** *In any communication round $t$, if the client selection $\mathbb{K}_t$ is a random variable sampled from a uniform distribution, the global model update $\Delta \mathbf{w}^t(\mathbb{K}_t) = \mathbf{w}^{t+1}(\mathbb{K}_t) - \boldsymbol{w}^t$ follows Gaussian Distribution, i.e.,*

$$
\begin{aligned}
&\mathbb{K}_t \sim \text{Uniform}\big(\{\mathbb{K} \subseteq \mathbb{U} : |\mathbb{K}| = C\}\big) \\
&\Rightarrow \Delta \mathbf{w}^t(\mathbb{K}_t) \sim \mathcal{N}(\Delta \boldsymbol{w}^t; -\eta_t \tilde{\boldsymbol{g}}^t, \frac{\eta_t^2 \boldsymbol{B} \boldsymbol{B}^T}{C}),
\end{aligned}
\tag{18}
$$

*where $\tilde{\boldsymbol{g}}^t = \mathbb{E}_k[\tilde{\nabla} l_k(\boldsymbol{w}^t)]$ is the mean cumulative gradient of all the clients in $\mathbb{U}$, and $\boldsymbol{B}$ is a constant matrix.*

Assumption 1 is inspired by [22] who assumes the stochastic gradients in SGD are Gaussian, and therefore the parameter update after one iteration follows a Gaussian Distribution. Note that in the FL procedure, the form in Eq. 5 is very similar to that in the SGD update. The only difference is that the average gradients within one mini-batch is replaced by the average cumulative gradients of the selected clients. Therefore, it is reasonable to make this assumption similar to [22].

To make a distinction, we use $\Delta \mathbf{w}$ without parentheses to denote a random variable w.r.t. the uniformly sampled client selection, and use $\Delta \boldsymbol{w}(\mathbb{K})$ to denote a determinate value without randomness where the client selection $\mathbb{K}$ is determined. The rule for $\Delta \mathbf{l}$ and $\Delta \boldsymbol{l}(\mathbb{K})$ in the following contents is the same.

Based on this assumption, we can easily show that the loss changes in each communication round follow a Gaussian Process under first-order approximation, with the property of Gaussian Distribution.

**Corollary 1.** *In any communication round $t$, $\forall \mathbb{S} = \{i_1, \cdots, i_{|\mathbb{S}|}\} \subseteq \mathbb{U}$, the loss changes $\Delta \mathbf{l}_{\mathbb{S}}^t = [\Delta \mathbf{l}_{i_1}^t, \cdots, \Delta \mathbf{l}_{i_{|\mathbb{S}|}}^t]^T$ follow a Multivariate Gaussian Distribution (or a Gaussian Process) under first-order approximation, i.e.,*

$$
\Delta \mathbf{l}_{\mathbb{S}}^t \sim \mathcal{N}(\Delta \boldsymbol{l}_{\mathbb{S}}^t; \boldsymbol{\mu}_{\mathbb{S}}^t, \boldsymbol{\Sigma}_{\mathbb{S}}^t),
$$
*where*
$$
\begin{aligned}
&\boldsymbol{\mu}_{\mathbb{S}}^t = -\eta_t \boldsymbol{G}_{\mathbb{S}}^{t\,T} \tilde{\boldsymbol{g}}^t; \\
&\boldsymbol{\Sigma}_{\mathbb{S}}^t = \frac{\eta_t^2}{C} \boldsymbol{G}_{\mathbb{S}}^{t\,T} \boldsymbol{B} \boldsymbol{B}^T \boldsymbol{G}_{\mathbb{S}}^t; \\
&\boldsymbol{G}_{\mathbb{S}}^t = \left[ \nabla l_{i_1}(\boldsymbol{w}^t), \cdots, \nabla l_{i_{|\mathbb{S}|}}(\boldsymbol{w}^t) \right].
\end{aligned}
\tag{19}
$$

We remove the subscript $\mathbb{S}$ to simplify the corresponding representation for the client set $\mathbb{U}$ as

$$
\Delta \mathbf{l}^t \sim \mathcal{N}(\Delta \boldsymbol{l}^t; \boldsymbol{\mu}^t, \boldsymbol{\Sigma}^t),
\tag{20}
$$

which is exactly the result in Eq. 8. And we can also obtain a mathematical reason from Eq. 19 for our choice of homogeneous linear kernel in Section 4.5, where $\boldsymbol{X}^t = \boldsymbol{B}^T \boldsymbol{G}^t$.

**Remark** Although an uniformly sampled client selection is required in Assumption 1 to get the loss changes to follow a GP prior, it is not necessary for the final selection to be uniformly sampled since we are predicting its loss changes with the GP posterior conditioned on the selected clients. We can view each posterior during the iterative selection process in Section 4.3 as the distribution of the loss changes w.r.t. the client selection that consists of two parts: (i) fixed selected clients in the previous iteration and (ii) uniformly sampled clients from the rest of the clients.

### A.2. Proof of Lemma 1

To prove Lemma 1, we first introduce another assumption.

**Assumption 2.** *In any communication round $t$, for any client selection $\mathbb{K}$, we have*

$$
Pr(\mathbb{K}|\Delta \boldsymbol{l}_{\mathbb{K}}^t(\mathbb{K})) \approx 1.
\tag{21}
$$

This assumption asserts that for any client selection $\mathbb{K}$, there is unlikely another client selection other than $\mathbb{K}$ which can produce the same loss changes on $\mathbb{K}$, i.e.,

$$
\begin{aligned}
&\forall \mathbb{K}', \mathbb{K} \subseteq \mathbb{U}, |\mathbb{K}'| = |\mathbb{K}| \\
&\Rightarrow Pr(\Delta \boldsymbol{l}_{\mathbb{K}}^t(\mathbb{K}') = \Delta \boldsymbol{l}_{\mathbb{K}}^t(\mathbb{K})|\mathbb{K}' \neq \mathbb{K}) \approx 0.
\end{aligned}
\tag{22}
$$

We anticipate that this is a realistic assumption because of the heterogeneity between clients and the highly complexity of the neural network. When selecting different clients, the data used for training varies a lot under heterogeneous federated learning settings. This fact makes it almost impossible to produce the same neural network, and thus the same loss changes, with two different client selections. Furthermore, the selected clients usually have larger loss decreases than other clients who are not selected, because the model update is based on the mean cumulative gradient of these selected clients. The other client selection is unlikely to generate the same large loss decreases on all of them.

With Assumption 2, we can get the following corollary 2.

**Corollary 2.** *In any communication round $t$, for any client selection $\mathbb{K}$, we have*

$$
Pr(\Delta \boldsymbol{l}^t(\mathbb{K})|\Delta \boldsymbol{l}_{\mathbb{K}}^t(\mathbb{K})) \approx 1.
\tag{23}
$$

*Proof.* When client selection $\mathbb{K}$ is given, we get the determinate model update $\Delta \boldsymbol{w}^t(\mathbb{K})$, thus the loss changes are known without randomness. In the other word,

$$Pr(\Delta \boldsymbol{l}^t(\mathbb{K})|\mathbb{K}) = 1 \tag{24}$$

always holds. Besides, we can extend the condition in Eq. 21 to the loss changes of all the clients and get

$$Pr(\mathbb{K}|\Delta \boldsymbol{l}^t(\mathbb{K})) \approx 1. \tag{25}$$

Combining Eq. 21, Eq. 24 and Eq. 25, we have

$$
\begin{align}
Pr(\Delta \boldsymbol{l}^t(\mathbb{K})) &\approx Pr(\Delta \boldsymbol{l}^t(\mathbb{K}), \mathbb{K}) \tag{26}\\
&= Pr(\mathbb{K}) \tag{27}\\
&= Pr(\Delta \boldsymbol{l}^t_{\mathbb{K}}(\mathbb{K}), \mathbb{K}) \tag{28}\\
&\approx Pr(\Delta \boldsymbol{l}^t_{\mathbb{K}}(\mathbb{K})) \tag{29}
\end{align}
$$

By substituting Eq. 29 into the expression of $Pr(\Delta \boldsymbol{l}^t(\mathbb{K})|\Delta \boldsymbol{l}^t_{\mathbb{K}}(\mathbb{K}))$, we get

$$
\begin{align}
Pr(\Delta \boldsymbol{l}^t(\mathbb{K})|\Delta \boldsymbol{l}^t_{\mathbb{K}}(\mathbb{K})) &= \frac{Pr(\Delta \boldsymbol{l}^t(\mathbb{K}), \Delta \boldsymbol{l}^t_{\mathbb{K}}(\mathbb{K}))}{Pr(\Delta \boldsymbol{l}^t_{\mathbb{K}}(\mathbb{K}))} \tag{30}\\
&= \frac{Pr(\Delta \boldsymbol{l}^t(\mathbb{K}))}{Pr(\Delta \boldsymbol{l}^t_{\mathbb{K}}(\mathbb{K}))} \tag{31}\\
&\approx 1. \tag{32}
\end{align}
$$

$\square$

Now we are ready to prove Lemma 1.

**Lemma 1.** *The optimization problem in Eq. (6) is approximately equivalent to the following probabilistic form.*

$$\min_{\mathbb{K}_t} \quad \mathbb{E}_{\Delta \mathbf{l}^t|\Delta \boldsymbol{l}^t_{\mathbb{K}_t}(\mathbb{K}_t)}\Big[\sum_i p_i \Delta l^t_i\Big] = \sum_i p_i \tilde{\mu}^t_i(\Delta \boldsymbol{l}^t_{\mathbb{K}_t}(\mathbb{K}_t)), \tag{33}$$

*where $\Delta \mathbf{l}^t = [\Delta \mathrm{l}^t_1, \cdots, \Delta \mathrm{l}^t_N]$ is the loss changes of all clients in round t, which is a random variable w.r.t random client selection in round t. $\tilde{\boldsymbol{\mu}}^t(\Delta \boldsymbol{l}^t_{\mathbb{K}_t}(\mathbb{K}_t))$ is the posterior mean of $\Delta \mathrm{l}^t$ conditioned on $\Delta \boldsymbol{l}^t_{\mathbb{K}_t}(\mathbb{K}_t) = [\Delta l^t_i(\mathbb{K}_t)]_{i \in \mathbb{K}_t}$.*

*Proof.* According to Corollary 2, we can transform the optimization problem in Eq. 6 into the form in Eq. 33.

$$
\begin{align}
& \min_{\mathbb{K}_t} \quad \Delta L^t(\mathbb{K}_t) \tag{34}\\
=& \min_{\mathbb{K}_t} \quad \sum_i p_i \Delta l^t_i(\mathbb{K}_t) \tag{35}\\
\approx& \min_{\mathbb{K}_t} \quad Pr(\Delta \boldsymbol{l}^t(\mathbb{K}_t)|\Delta \boldsymbol{l}^t_{\mathbb{K}_t}(\mathbb{K}_t)) \sum_i p_i \Delta l^t_i(\mathbb{K}_t) \tag{36}\\
\approx& \min_{\mathbb{K}_t} \quad \mathbb{E}_{\Delta \mathbf{l}^t|\Delta \boldsymbol{l}^t_{\mathbb{K}_t}(\mathbb{K}_t)}\Big[\sum_i p_i \Delta l^t_i\Big] \tag{37}\\
=& \min_{\mathbb{K}_t} \quad \sum_i p_i \tilde{\mu}^t_i(\Delta \boldsymbol{l}^t_{\mathbb{K}_t}(\mathbb{K}_t)). \tag{38}
\end{align}
$$

$\square$

## A.3. Proof of Lemma 2

**Lemma 2.** *The selection criterion of FedCor when selecting two clients $k_1$ and $k_2$ can be written as*

$$k_1 = \arg\max_k \quad \beta^{\tau_k} \sum_i p_i \sigma_i r_{ik}, \tag{39}$$

$$k_2 = \arg\max_{k'} \quad \frac{\beta^{\tau_{k'}}\left[\sum_i p_i \sigma_i r_{ik'} - r_{k_1 k'}\sum_i p_i \sigma_i r_{ik_1}\right]}{\sqrt{1 - r^2_{k' k_1}}}, \tag{40}$$

*where $r_{ij} = \Sigma_{i,j}/\sigma_i \sigma_j$ is the Pearson correlation coefficient.*

*Proof.* We first deduce Eq. 39 for the first client $k_1$. By substituting the loss change estimation $\Delta \hat{l}_k$ from Eq. 9 into the criterion in Eq. 10, we can calculate the weighted sum of the posterior mean as

$$
\begin{align}
& \sum_i p_i \tilde{\mu}_i(\Delta \hat{l}_k) \tag{41}\\
=& \sum_i p_i \mu_i + \sum_i p_i \frac{\Sigma_{i,k}}{\sigma^2_k}(\Delta \hat{l}_k - \mu_k) \tag{42}\\
=& \sum_i p_i \mu_i - a\beta^{\tau_k}\sum_i p_i \sigma_i r_{ik}, \tag{43}
\end{align}
$$

where $r_{ik}$ is the Pearson correlation coefficient. The first item in Eq. 43 and the factor $a$ are constant for all $k$, thus the selection strategy becomes

$$k_1 = \arg\max_k \quad \beta^{\tau_k}\sum_i p_i \sigma_i r_{ik}, \tag{44}$$

which is Eq. 39.

Then we deduce Eq. 40 for selecting $k_2$. We can calculate the posterior covariance conditioned on $\Delta \hat{l}_{k_1}$ as

$$
\begin{align}
\tilde{\Sigma}_{i,j}(\Delta \hat{l}_{k_1}) &= \Sigma_{i,j} - \frac{\Sigma_{i,k_1}\Sigma_{k_1,j}}{\sigma^2_{k_1}} \tag{45}\\
&= \sigma_i \sigma_j(r_{ij} - r_{ik_1}r_{k_1 j}) \tag{46}\\
\tilde{\sigma}_i(\Delta \hat{l}_{k_1}) &= \sqrt{\tilde{\Sigma}_{i,i}(\Delta \hat{l}_{k_1})} \tag{47}\\
&= \sigma_i\sqrt{1 - r^2_{ik_1}}. \tag{48}
\end{align}
$$

We substitute the posterior covariance into the simplified selection criterion in Eq. 44 and get

$$
\begin{align}
& \beta^{\tau_{k'}}\sum_i p_i \frac{\tilde{\Sigma}_{i,k'}(\Delta \hat{l}_{k_1})}{\tilde{\sigma}_{k'}(\Delta \hat{l}_{k_1})} \\
=& \frac{\beta^{\tau_{k'}}\left[\sum_i p_i \sigma_i r_{ik'} - r_{k_1 k'}\sum_i p_i \sigma_i r_{ik_1}\right]}{\sqrt{1 - r^2_{k' k_1}}}. \tag{49}
\end{align}
$$

So we have Eq. 40:

$$k_2 = \arg\max_{k'} \frac{\beta^{\tau_{k'}} \left[ \sum_i p_i \sigma_i r_{ik'} - r_{k_1 k'} \sum_i p_i \sigma_i r_{ik_1} \right]}{\sqrt{1 - r_{k' k_1}^2}}.$$

(50)

□

## B. Selection Criterion and Convergence Analysis

In this section, we will analyse FedCor when selecting arbitrary number of clients. While the iterative client selection makes it obscure to analyse the convergence, we will show that we can construct a simpler proxy algorithm who can approximate the selection strategy of FedCor and there for share similar convergence characteristic. We will prove the convergence of this proxy algorithm.

### B.1. Definitions

We first introduce some important definitions. In the following analysis, We denote the client selection sampled from FedCor as $\mathbb{K}_t \sim \pi$ and client selection sampled uniformly as $\mathbb{K}_t \sim \mathcal{U}$.

In the $j$-th iteration of FedCor, we select a client $k_j$ to minimize the posterior mean of the loss change. Since the prior mean in each iteration is fixed, we can say that we are maximizing the decrease from prior mean $\boldsymbol{\mu}^{t,j}$ to posterior mean $\tilde{\boldsymbol{\mu}}^{t,j}$. We define the posterior gain of this iteration as the decrease from prior mean to posterior mean, namely,

$$g^{t,j}(k_j) = \sum_i p_i(\mu_i^{t,j} - \tilde{\mu}_i^{t,j}(\Delta \hat{l}_{k_j}^t)) \quad (51)$$

$$= \alpha_{k_j}^t \sum_i p_i \sigma_i^{t,j} r_{ik_j}^{t,j}. \quad (52)$$

We define $\boldsymbol{\mu}^{t,1} = \boldsymbol{\mu}^t$ and $\boldsymbol{\Sigma}^{t,1} = \boldsymbol{\Sigma}^t$. And for $j > 1$ we have

$$\boldsymbol{\mu}^{t,j} = \tilde{\boldsymbol{\mu}}^{t,j-1}(\Delta \hat{l}_{k_{j-1}}^t), \qquad \boldsymbol{\Sigma}^{t,j} = \tilde{\boldsymbol{\Sigma}}^{t,j-1}(\Delta \hat{l}_{k_{j-1}}^t). \quad (53)$$

With Lemma 2, we get

$$g^{t,j}(k_j) = \frac{g^{t,j-1}(k_j) - \frac{\alpha_{k_j}^t}{\alpha_{k_{j-1}}^t} r_{k_{j-1} k_j}^{t,j-1} g^{t,j-1}(k_{j-1})}{\sqrt{1 - r_{k_{j-1} k_j}^{t,j-1}{}^2}}. \quad (54)$$

With this notation, we can simplify our selection strategy as follows.

$$k_j^* = \arg\max_{k_j} g^{t,j}(k_j). \quad (55)$$

We further define the one-round advantage of FedCor compared with uniform sampling as follows.

$$A^t = \mathbb{E}_{\mathbb{K}_t \sim \mathcal{U}}[L(\boldsymbol{w}^{t+1}) - L(\boldsymbol{w}^t)] - \mathbb{E}_{\mathbb{K}_t \sim \pi}[L(\boldsymbol{w}^{t+1}) - L(\boldsymbol{w}^t)] \quad (56)$$

$$= \sum_{j=1}^C g^{t,j}(k_j^*). \quad (57)$$

The second equation directly arises from the definition of our prior distribution where $\mathbb{E}_{\mathbb{K}_t \sim \mathcal{U}}[L(\boldsymbol{w}^{t+1}) - L(\boldsymbol{w}^t)] = \sum_i \mu_i^t$.

Unfortunately, because of the iterative selection, the selection criterion of $k_j$ depends on the previous selected clients, which makes a quantitatively analysis complicated. To bypass this difficulty, we will first point out that $A^t$ has a lower bound that is tight in some special cases. We find that a proxy client selection strategy that maximizes this lower bound has a similar but simpler behaviour compared with FedCor, and we will also give a convergence guarantee of the proxy algorithm.

### B.2. Approximation of FedCor

An important property of FedCor is that it prefers clients who have lower correlations with those selected in the previous iteration, since

$$\forall r_{k_{j-1} k_j}^{t,j-1} \in (-1,1), \frac{\partial g^{t,j}(k_j)}{\partial r_{k_{j-1} k_j}^{t,j-1}} < 0. \quad (58)$$

We further predict that FedCor tends to select clients that with $r_{k_{j-1} k_j}^{t,j-1}$ close to 0 instead of $r_{k_{j-1} k_j}^{t,j-1} < 0$ because if $r_{k_{j-1} k_j}^{t,j-1} < 0$, $k_j$ should be far away from $k_{j-1}$ who is closed to other clients in the embedding space, which makes $k_j$ has low correlation with the other clients and not be selected. Therefore, we can infer that FedCor will select a group of clients who have nearly zero correlations with each other, which simplifies the expression of $g^{t,j}(k_j)$ to $g^{t,1}(k_j)$.

Based on the analysis above, we define a proxy algorithm $\tilde{\pi}$ who maximize the following objective.

$$\tilde{A}^t = \sum_{k_j \in \mathbb{K}_t} g^{t,1}(k_j) \approx \sum_{k \in \mathbb{K}_t} \sum_i p_i \Sigma_{i,k}^t, \quad (59)$$

where we further omit the difference of $\alpha_k^t$ and $\sigma_k^t$ for different client $k$. We can use the client selection generated by this proxy algorithm to approximate the client selection of FedCor, and thus they share similar convergence characteristic.

In the following section, we will show that this proxy algorithm has a good property that enable it to converge to the optimal solution of the global loss $L$ without gap, even it is a biased selection strategy.

## B.3. Convergence Analysis of the Proxy Algorithm

In the following section, we denote the client selection sampled from the proxy client selection strategy as $\mathbb{K}_t \sim \tilde{\pi}$. We use $\mathbb{E}[\cdot]$ as the expectation over the mini-batch and $\mathbb{E}_{\mathbb{K}_t}[\cdot]$ as the expectation over the client selection strategy. We first give the common assumptions used in Federated Learning [4, 19].

**Assumption 3.** $l_1, l_2, \cdots, l_N$ are all $M$-smooth: for all $\boldsymbol{v}$ and $\boldsymbol{w}$, $l_k(\boldsymbol{v}) \leq l_k(\boldsymbol{w}) + (\boldsymbol{v} - \boldsymbol{w})^T \nabla l_k(\boldsymbol{w}) + \frac{M}{2}\|\boldsymbol{v} - \boldsymbol{w}\|_2^2$.

**Assumption 4.** $l_1, l_2, \cdots, l_N$ are all $m$-strongly convex: for all $\boldsymbol{v}$ and $\boldsymbol{w}$, $l_k(\boldsymbol{v}) \geq l_k(\boldsymbol{w}) + (\boldsymbol{v} - \boldsymbol{w})^T \nabla l_k(\boldsymbol{w}) + \frac{m}{2}\|\boldsymbol{v} - \boldsymbol{w}\|_2^2$.

**Assumption 5.** For the mini-batch $\xi_k \in \mathbb{D}_k$ sampled uniformly on each client $k \in \mathbb{U}$, the variance of stochastic gradients is bounded: $\mathbb{E}\|\nabla l_k(\boldsymbol{w}_k, \xi_k) - \nabla l_k(\boldsymbol{w}_k)\|^2 \leq s_k^2$.

**Assumption 6.** For each client $k \in \mathbb{U}$ and any communication round $t$, the expected squared norm of stochastic gradients is uniformly bounded: $\mathbb{E}\|\nabla l_k(\boldsymbol{w}_k, \xi_k)\|^2 \leq G^2$.

For concision, we omit $\mathbb{E}$ in the following content and apply an expectation over the mini-batch by default.

Now we give an important property of the proxy algorithm that will be used for proving the convergence.

**Lemma 3.** In any communication round $t$, with Assumption 1 and Assumption 2 holds, we have

$$\mathbb{K}_t \sim \tilde{\pi} = \arg\max_{\mathbb{K}} (\boldsymbol{B}^T \nabla L(\boldsymbol{w}^t))^T \sum_{k \in \mathbb{K}} \boldsymbol{B}^T \nabla l_k(\boldsymbol{w}^t). \tag{60}$$

*Proof.* In the proxy algorithm, we have

$$\mathbb{K}_t = \arg\max_{\mathbb{K}} \sum_{k \in \mathbb{K}_t} \sum_i p_i \Sigma_{i,k}^t \tag{61}$$

$$= \arg\max_{\mathbb{K}} \frac{\eta_t^2}{C} \sum_{k \in \mathbb{K}} \sum_i p_i \nabla l_i(\boldsymbol{w}^t) \boldsymbol{B}\boldsymbol{B}^T \nabla l_k(\boldsymbol{w}^t) \tag{62}$$

$$= \arg\max_{\mathbb{K}} (\boldsymbol{B}^T \nabla L(\boldsymbol{w}^t))^T \sum_{k \in \mathbb{K}} \boldsymbol{B}^T \nabla l_k(\boldsymbol{w}^t). \tag{63}$$

Eq. 62 comes from the expression of $\boldsymbol{\Sigma}^t$ in Corollary 1, and Eq. (63) arises from $L(\boldsymbol{w}^t) = \sum_i p_i l_i(\boldsymbol{w}^t)$. □

To connect this property with the convergence of the algorithm, we first define a sequence and show that the convergence of this sequence is equivalent to the convergence of the algorithm with this property. We define Sequence $\Delta_t$ as follows.

$$\Delta_t = \mathbb{E}_{\mathbb{K}_t \sim \tilde{\pi}} \|\boldsymbol{w}^t - \boldsymbol{w}^*_{\mathbb{K}_t}\|^2, \tag{64}$$

where

$$\boldsymbol{w}^*_{\mathbb{K}_t} = \arg\min_{\boldsymbol{w}} \sum_{k \in \mathbb{K}_t} l_k(\boldsymbol{w}). \tag{65}$$

We now show that if $\Delta_t \to 0$, we have $\boldsymbol{w} \to \boldsymbol{w}^*$.

**Corollary 3.** *(Optimal Solution Consistency)* If $\Delta_t$ converges to 0, there must be $\boldsymbol{w}^t$ converges to $\boldsymbol{w}^*$.

$$\lim_{t \to \infty} \Delta_t = 0 \Rightarrow \lim_{t \to \infty} \boldsymbol{w}^t = \boldsymbol{w}^* \tag{66}$$

*Proof.* With $\mathbb{K}_t \sim \tilde{\pi}$, we have

$$\lim_{t \to \infty} \Delta_t = 0 \tag{67}$$

$$\Rightarrow \lim_{t \to \infty} \boldsymbol{w}^t = \boldsymbol{w}^*_{\mathbb{K}_t} \tag{68}$$

$$\Rightarrow \lim_{t \to \infty} \sum_{k \in \mathbb{K}_t} \nabla l_k(\boldsymbol{w}^t) = \boldsymbol{0} \tag{69}$$

$$\Rightarrow \lim_{t \to \infty} (\boldsymbol{B}^T \nabla L(\boldsymbol{w}^t))^T \sum_{k \in \mathbb{K}_t} \boldsymbol{B}^T \nabla l_k(\boldsymbol{w}^t) = 0. \tag{70}$$

Since

$$\mathbb{K}_t = \arg\max_{\mathbb{K}} (\boldsymbol{B}^T \nabla L(\boldsymbol{w}^t))^T \sum_{k \in \mathbb{K}} \boldsymbol{B}^T \nabla l_k(\boldsymbol{w}^t), \tag{71}$$

If $\lim_{t \to \infty} \boldsymbol{B}^T \nabla L(\boldsymbol{w}^t) \neq \boldsymbol{0}$ or does not converge, we can say that

$$\forall \epsilon > 0, \exists \tau, \forall t > \tau, \forall \mathbb{K}, \tag{72}$$

$$(\boldsymbol{B}^T \nabla L(\boldsymbol{w}^t))^T \sum_{k \in \mathbb{K}} \boldsymbol{B}^T \nabla l_k(\boldsymbol{w}^t) \leq \epsilon, \tag{73}$$

which cannot be true since

$$\mathbb{E}_{\mathbb{K} \sim \mathcal{U}} \sum_{k \in \mathbb{K}} \nabla l_k(\boldsymbol{w}^t) = C \nabla L(\boldsymbol{w}^t). \tag{74}$$

Thus we conclude that

$$\lim_{t \to \infty} \boldsymbol{B}^T \nabla L(\boldsymbol{w}^t) = 0. \tag{75}$$

If the Gaussian Distribution in Assumption 1 is non-degenerate, we have

$$\lim_{t \to \infty} \nabla L(\boldsymbol{w}^t) = \boldsymbol{0} \Rightarrow \lim_{t \to \infty} \boldsymbol{w}^t = \boldsymbol{w}^* \tag{76}$$

□

We now only need to prove the convergence of $\Delta_t$, which will imply the convergence of the proxy algorithm according to Corollary 3. We first introduce one extra assumption as well as two lemmas that will be used in the proof.

For convenient, we define $L_{\mathbb{K}_t}(\boldsymbol{w}) = \frac{1}{C} \sum_{k \in \mathbb{K}_t} l_k(\boldsymbol{w})$, and thus $\boldsymbol{w}^*_{\mathbb{K}_t} = \arg\min_{\boldsymbol{w}} L_{\mathbb{K}_t}(\boldsymbol{w})$. Notice that $\mathbb{K}_t \sim \tilde{\pi}$ only depends on $\boldsymbol{\Sigma}^t$, thus we can say that $\boldsymbol{w}^*_{\mathbb{K}_t}$ is given by a function of $\boldsymbol{\Sigma}^t$, i.e., $\boldsymbol{w}^*_{\mathbb{K}_t} = \Omega(\boldsymbol{\Sigma}^t)$. We further assume the smoothness of $\Omega$:

**Assumption 7.** For any $t$, $\mathbb{E}\|\boldsymbol{w}^*_{\mathbb{K}_{t+1}} - \boldsymbol{w}^*_{\mathbb{K}_t}\|^2 = \mathbb{E}\|\Omega(\boldsymbol{\Sigma}^{t+1}) - \Omega(\boldsymbol{\Sigma}^t)\|^2 \leq \delta \mathbb{E}\|\boldsymbol{\Sigma}^{t+1} - \boldsymbol{\Sigma}^t\|_1$, where $\|\cdot\|_1$ is the $\ell_1$ norm of a vector.

Now we introduce a lemma that bounds $\mathbb{E}\|\boldsymbol{\Sigma}^{t+1} - \boldsymbol{\Sigma}^t\|_1$.

**Lemma 4.** *Assume Assumption 1, Assumption 3 and Assumption 6, if $\mathbb{E}\|\boldsymbol{w}^t - \boldsymbol{w}^{t+1}\|^2 \leq q_t^2$, we have*

$$\mathbb{E}\|\boldsymbol{\Sigma}_{t+1} - \boldsymbol{\Sigma}_t\|_1 \leq \frac{bN^2}{C}[\eta_t^2(G + Mq_t)^2 - \eta_{t+1}^2 G^2], \tag{77}$$

*where $b$ is the largest eigenvalue of $\boldsymbol{BB}^T$.*

*Proof.* According to Assumption 1, we have

$$\Sigma_{i,j} = \frac{\eta_t^2}{C}\nabla l_i^{tT}\boldsymbol{BB}^T\nabla l_j^t. \tag{78}$$

And we can calculate

$$\left|\Sigma_{i,j}^{t+1} - \Sigma_{i,j}^t\right| \tag{79}$$

$$= \left|\frac{\eta_t^2}{C}(\nabla l_i^{t+1T}\boldsymbol{BB}^T\nabla l_j^{t+1} - \nabla l_i^{tT}\boldsymbol{BB}^T\nabla l_j^t) + \frac{\eta_{t+1}^2 - \eta_t^2}{C}\nabla l_i^{t+1T}\boldsymbol{BB}^T\nabla l_j^{t+1}\right| \tag{80}$$

$$\leq \frac{\eta_t^2}{C}\left|\nabla l_i^{t+1T}\boldsymbol{BB}^T\nabla l_j^{t+1} - \nabla l_i^{tT}\boldsymbol{BB}^T\nabla l_j^t\right| + \frac{\eta_t^2 - \eta_{t+1}^2}{C}\left|\nabla l_i^{t+1T}\boldsymbol{BB}^T\nabla l_j^{t+1}\right|. \tag{81}$$

We now bound each term in Eq. (81) separately. For the first term,

$$\left|\nabla l_i^{t+1T}\boldsymbol{BB}^T\nabla l_j^{t+1} - \nabla l_i^{tT}\boldsymbol{BB}^T\nabla l_j^t\right| \tag{82}$$

$$= \left|(\nabla l_i^{t+1} - \nabla l_i^t)^T\boldsymbol{BB}^T\nabla l_j^t + \nabla l_i^{tT}\boldsymbol{BB}^T(\nabla l_j^{t+1} - \nabla l_j^t) + (\nabla l_i^{t+1} - \nabla l_i^t)^T\boldsymbol{BB}^T(\nabla l_j^{t+1} - \nabla l_j^t)\right| \tag{83}$$

$$\leq b\Big(\|\nabla l_i^{t+1} - \nabla l_i^t\|\|\nabla l_j^t\| + \|\nabla l_j^{t+1} - \nabla l_j^t\|\|\nabla l_i^t\| + \|\nabla l_i^{t+1} - \nabla l_i^t\|\|\nabla l_j^{t+1} - \nabla l_j^t\|\Big) \tag{84}$$

$$\leq b\Big[M\|\boldsymbol{w}^{t+1} - \boldsymbol{w}^t\|(\|\nabla l_j^t\| + \|\nabla l_i^t\|) + M^2\|\boldsymbol{w}^{t+1} - \boldsymbol{w}^t\|^2\Big], \tag{85}$$

where $b$ is the largest eigenvalue of $\boldsymbol{BB}^T$. For the second term,

$$\left|\nabla l_i^{t+1T}\boldsymbol{BB}^T\nabla l_j^{t+1}\right| \leq b\|\nabla l_i^{t+1}\|\|\nabla l_j^{t+1}\|. \tag{86}$$

We take the expectation over both sides and with Cauchy-Schwarz inequality, we get

$$\mathbb{E}\left|\Sigma_{i,j}^{t+1} - \Sigma_{i,j}^t\right| \tag{87}$$

$$\leq \frac{\eta_t^2}{C}b\Big[M\sqrt{\mathbb{E}\|\boldsymbol{w}^{t+1} - \boldsymbol{w}^t\|^2\mathbb{E}\|\nabla l_j^t\|^2} + M\sqrt{\mathbb{E}\|\boldsymbol{w}^{t+1} - \boldsymbol{w}^t\|^2\mathbb{E}\|\nabla l_i^t\|^2} + M^2\mathbb{E}\|\boldsymbol{w}^{t+1} - \boldsymbol{w}^t\|^2\Big] + \frac{\eta_t^2 - \eta_{t+1}^2}{C}b\sqrt{\mathbb{E}\|\nabla l_i^{t+1}\|^2\mathbb{E}\|\nabla l_j^{t+1}\|^2} \tag{88}$$

$$\leq \frac{\eta_t^2 b}{C}(G^2 + 2Mq_tG + M^2q_t^2) - \frac{\eta_{t+1}^2}{C}bG^2 \tag{89}$$

And we have

$$\mathbb{E}\|\boldsymbol{\Sigma}_{t+1} - \boldsymbol{\Sigma}_t\|_1 \tag{90}$$

$$= \sum_{i,j}^N \mathbb{E}\left|\Sigma_{i,j}^{t+1} - \Sigma_{i,j}^t\right| \tag{91}$$

$$\leq \frac{bN^2}{C}[\eta_t^2(G + Mq_t)^2 - \eta_{t+1}^2 G^2] \tag{92}$$

$\square$

We will also use the following lemma that is proved by [19].

**Lemma 5.** *Assume Assumption 3 to 6. If $\eta_t \leq \frac{1}{4M}$, with full and balanced participation in FedAvg, in any communication round $t$ and its $i$-th iteration, we have*

$$\mathbb{E}\|\bar{\boldsymbol{w}}^{t,i+1} - \boldsymbol{w}^*\|^2 \leq (1 - \eta_t m)\mathbb{E}\|\bar{\boldsymbol{w}}^{t,i} - \boldsymbol{w}^*\|^2 + \eta_t^2 F, \tag{93}$$

*where*

$$F = \frac{1}{N}\sum_{k=1}^N s_k^2 + 6M\Gamma + 8(E - 1)^2 G^2, \tag{94}$$

$$\Gamma = L^* - \frac{1}{N}\sum_{k=1}^N l_k^*. \tag{95}$$

*Here, $\bar{\boldsymbol{w}}^{t,i} = \frac{1}{N}\sum_{k=1}^N \boldsymbol{w}_k^{t,i}$, and $\boldsymbol{w}_k^{t,i}$ is the local weight at the $i$-th iteration of communication round $t$. $E$ is the total number of local training iterations. $L^* = L(\omega^*)$ and $l_k^* = l_k(\omega_k^*)$ are the optimal value of $L$ and $l_k$, respectively.*

Now we give the theorem of the convergence of $\Delta_t$ and prove it.

**Theorem 1.** *With Assumption 1 to 7 holds, with learning rate $\eta_t = \frac{\beta}{t+\gamma}$ for some $\beta > \frac{1}{m}$ and $\gamma > 0$ such that $\eta_1 \leq \min\{\frac{1}{m}, \frac{1}{4M}\} = \frac{1}{4M}$, we have*

$$\Delta_t \leq \frac{\nu}{\gamma + t}, \tag{96}$$

*where*

$$\nu = \max\{\frac{\beta^2(\tilde{F} + \tilde{D})}{\beta m - 1}, (\gamma + 1)\Delta_1\}, \qquad (97)$$

$$\tilde{F} = 2E \max_t F_t, \qquad (98)$$

$$F_t = \frac{1}{C} \sum_{k \in \mathbb{K}_t} s_k^2 + 6M\Gamma_t + 8(E-1)^2 G^2, \qquad (99)$$

$$\Gamma_t = L_{\mathbb{K}_t}^* - \frac{1}{C} \sum_{k \in \mathbb{K}_t} l_k^*, \qquad (100)$$

$$\tilde{D} = (\frac{1}{m} + \frac{1}{4M})\delta D, \qquad (101)$$

$$D = \frac{bN^2}{C}(2mG^2 + 2MEG + \frac{1}{4}ME^2 G^2). \qquad (102)$$

*Proof.* For $\mathbb{K}_t \sim \tilde{\pi}(\boldsymbol{w}^t)$ and $\mathbb{K}_{t+1} \sim \tilde{\pi}(\boldsymbol{w}^{t+1})$, we have

$$\begin{aligned}
\Delta_{t+1} =& \|\boldsymbol{w}^{t+1} - \boldsymbol{w}_{\mathbb{K}_{t+1}}^*\|^2 & (103)\\
=& \|\boldsymbol{w}^{t+1} - \boldsymbol{w}_{\mathbb{K}_t}^*\|^2 + \|\boldsymbol{w}_{\mathbb{K}_t}^* - \boldsymbol{w}_{\mathbb{K}_{t+1}}^*\|^2 + \\
& 2\langle \boldsymbol{w}^{t+1} - \boldsymbol{w}_{\mathbb{K}_t}^*, \boldsymbol{w}_{\mathbb{K}_t}^* - \boldsymbol{w}_{\mathbb{K}_{t+1}}^* \rangle & (104)\\
\leq& \|\boldsymbol{w}^{t+1} - \boldsymbol{w}_{\mathbb{K}_t}^*\|^2 + \|\boldsymbol{w}_{\mathbb{K}_t}^* - \boldsymbol{w}_{\mathbb{K}_{t+1}}^*\|^2 + \\
& \eta_t m \|\boldsymbol{w}^{t+1} - \boldsymbol{w}_{\mathbb{K}_t}^*\|^2 + \\
& \frac{1}{\eta_t m}\|\boldsymbol{w}_{\mathbb{K}_t}^* - \boldsymbol{w}_{\mathbb{K}_{t+1}}^*\|^2 & (105)\\
\leq& (1 + \eta_t m)\|\boldsymbol{w}^{t+1} - \boldsymbol{w}_{\mathbb{K}_t}^*\|^2 + \\
& (1 + \frac{1}{\eta_t m})\delta\|\boldsymbol{\Sigma}^{t+1} - \boldsymbol{\Sigma}^t\|_1, & (106)
\end{aligned}$$

where Eq. 105 arises from AM-GM inequality and Eq. 106 arises from Assumption 7.

For the first term in Eq. 106, we can bound it by Lemma 5 as follows. The key point here is that when training in one communication round $t$, we can view this round a small FL process with clients in $\mathbb{K}_t$ fully participating. In this view, the global loss and the optimal global weight becomes $L_{\mathbb{K}_t}$ and $\boldsymbol{w}_{\mathbb{K}_t}^*$ instead. Thus we can apply Lemma 5 directly to bound $\|\boldsymbol{w}_{t+1} - \boldsymbol{w}_{\mathbb{K}_t}^*\|^2$. With $\eta_t \leq \frac{1}{4M} \leq \frac{1}{m}$, we have $\eta_t m \leq 1$ and $1 + \eta_t m \leq \frac{1}{1-\eta_t m}$, and we can get

$$\begin{aligned}
& (1 + \eta_t m)\|\boldsymbol{w}^{t+1} - \boldsymbol{w}_{\mathbb{K}_t}^*\|^2 & (107)\\
=& (1 + \eta_t m)[\|\bar{\boldsymbol{w}}^{t,E} - \boldsymbol{w}_{\mathbb{K}_t}^*\|^2] & (108)\\
\leq& (1 + \eta_t m)[(1 - \eta_t m)\|\bar{\boldsymbol{w}}^{t,E-1} - \boldsymbol{w}_{\mathbb{K}_t}^*\|^2 + \eta_t^2 F_t] & (109)\\
\leq& (1 + \eta_t m)\{(1 - \eta_t m)^2\|\bar{\boldsymbol{w}}^{t,E-2} - \boldsymbol{w}_{\mathbb{K}_t}^*\|^2 + \\
& [1 + (1 - \eta_t m)]\eta_t^2 F_t\} & (110)\\
& \cdots \\
\leq& (1 + \eta_t m)\{(1 - \eta_t m)^E\|\bar{\boldsymbol{w}}^{t,0} - \boldsymbol{w}_{\mathbb{K}_t}^*\|^2 + \\
& [1 + (1 - \eta_t m) + \cdots + (1 - \eta_t m)^{E-1}]\eta_t^2 F_t\} & (111)\\
\leq& (1 - \eta_t m)^{E-1}\|\boldsymbol{w}^t - \boldsymbol{w}_{\mathbb{K}_t}^*\|^2 + \\
& (1 + \eta_t m)\frac{1 - (1 - \eta_t m)^E}{m}\eta_t F_t & (112)\\
\leq& (1 - \eta_t m)\|\boldsymbol{w}^t - \boldsymbol{w}_{\mathbb{K}_t}^*\|^2 + (1 + \eta_t m)E\eta_t^2 F_t & (113)\\
\leq& (1 - \eta_t m)\|\boldsymbol{w}^t - \boldsymbol{w}_{\mathbb{K}_t}^*\|^2 + 2E\eta_t^2 F_t, & (114)
\end{aligned}$$

where

$$F_t = \frac{1}{C} \sum_{k \in \mathbb{K}_t} s_k^2 + 6M\Gamma_t + 8(E-1)^2 G^2, \qquad (115)$$

$$\Gamma_t = L_{\mathbb{K}_t}^* - \frac{1}{C} \sum_{k \in \mathbb{K}_t} l_k^*. \qquad (116)$$

Eq. 114 arises from the inequality $1 - Ex \leq (1-x)^E \leq 1-x$ for $x \in [0, 1]$.

We now turn to bound the second term in Eq. 106. We first find the $q_t$ in Lemma 4.

$$\begin{aligned}
\|\boldsymbol{w}^{t+1} - \boldsymbol{w}^t\|^2 =& \|\frac{1}{C} \sum_{k \in \mathbb{K}_t} \boldsymbol{w}_k^{t,E} - \boldsymbol{w}^t\|^2 & (117)\\
\leq& \frac{1}{C} \sum_{k \in \mathbb{K}_t} \|\boldsymbol{w}_k^{t,E} - \boldsymbol{w}^t\|^2 & (118)\\
=& \frac{\eta_t^2}{C} \sum_{k \in \mathbb{K}_t} \|\sum_{i=0}^{E-1} \nabla l_k(\boldsymbol{w}_k^{t,i})\|^2 & (119)\\
\leq& \frac{\eta_t^2 E}{C} \sum_{k \in \mathbb{K}_t} \sum_{i=0}^{E-1} \|\nabla l_k(\boldsymbol{w}_k^{t,i})\|^2 & (120)\\
\leq& \frac{\eta_t^2 E}{C} \sum_{k \in \mathbb{K}_t} \sum_{i=0}^{E-1} G^2 & (121)\\
\leq& \frac{\eta_t^2 E}{C} \sum_{k \in \mathbb{K}_t} EG^2 & (122)\\
=& \eta_t^2 E^2 G^2 = q_t^2, & (123)
\end{aligned}$$

where Eq. 118 and Eq. 120 comes from Jensen inequality, and Eq. 121 comes from Assumption 6. With Lemma

Lemma [4], we get

$$\|\mathbf{\Sigma}^{t+1} - \mathbf{\Sigma}^t\|_1 \le \frac{bN^2}{C}[G^2(\eta_t^2 - \eta_{t+1}^2)+ \tag{124}$$

$$2MEG\eta_t^3 + M^2E^2G^2\eta_t^4)]. \tag{125}$$

Further with a diminishing $\eta_t = \frac{\beta}{t+\gamma}$, we have

$$\eta_t^2 - \eta_{t+1}^2 = \beta^2\left(\frac{1}{(t+\gamma)^2} - \frac{1}{(t+1+\gamma)^2}\right) \tag{126}$$

$$= \beta^2 \frac{2(t+\gamma)+1}{(t+\gamma)^2(t+1+\gamma)^2} \tag{127}$$

$$\le \frac{2\beta^2}{(t+\gamma)^3} \tag{128}$$

$$= \frac{2\eta_t^3}{\beta}, \tag{129}$$

and with $\beta > \frac{1}{m}, \eta_t \le \eta_1 \le \frac{1}{4M}$, we get

$$\|\mathbf{\Sigma}^{t+1} - \mathbf{\Sigma}^t\|_1 \tag{130}$$

$$\le \frac{bN^2\eta_t^3}{C}\left(\frac{2G^2}{\beta} + 2MEG + M^2E^2G^2\eta_t\right) \tag{131}$$

$$\le \frac{bN^2\eta_t^3}{C}\left(2mG^2 + 2MEG + \frac{1}{4}ME^2G^2\right) \tag{132}$$

$$= \eta_t^3 D, \tag{133}$$

where

$$D = \frac{bN^2}{C}\left(2mG^2 + 2MEG + \frac{1}{4}ME^2G^2\right). \tag{134}$$

With Eq. [114] and Eq. [133], we have

$$\Delta_{t+1} \le (1-\eta_t m)\Delta_t + 2E\eta_t^2 F_t + \left(1 + \frac{1}{\eta_t m}\right)\eta_t^3 \delta D \tag{135}$$

$$\le (1-\eta_t m)\Delta_t + \eta_t^2\left(2EF_t + \frac{\delta}{m}D\right) + \eta_t^3 \delta D \tag{136}$$

$$\le (1-\eta_t m)\Delta_t + \eta_t^2(\tilde{F} + \tilde{D}), \tag{137}$$

where

$$\tilde{F} = 2E\max_t F_t, \tag{138}$$

$$\tilde{D} = \left(\frac{1}{m} + \frac{1}{4M}\right)\delta D. \tag{139}$$

Now we can use the same trick in [19] to finish the proof of convergence. With a diminishing learning rate, $\eta_t = \frac{\beta}{t+\gamma}$ for some $\beta > \frac{1}{m}$ and $\gamma > 0$ such that $\eta_1 \le \min\{\frac{1}{m}, \frac{1}{4M}\} = \frac{1}{4M}$, we will prove by induction that $\Delta_t \le \frac{\nu}{\gamma+t}$, where $\nu = \max\{\frac{\beta^2(\tilde{F}+\tilde{D})}{\beta m - 1}, (\gamma+1)\Delta_1\}$.

With the definition of $\nu$, we ensure that $\Delta_1 \le \frac{\nu}{\gamma+1}$. Now we assume that $\Delta_t \le \frac{\nu}{\gamma+t}$ holds for some $t$, we have

$$\Delta_{t+1} \le (1-\eta_t m)\Delta_t + \eta_t^2(\tilde{F} + \tilde{D}) \tag{140}$$

$$\le \left(1 - \frac{\beta m}{t+\gamma}\right)\frac{\nu}{t+\gamma} + \frac{\beta^2(\tilde{F}+\tilde{D})}{(t+\gamma)^2} \tag{141}$$

$$= \frac{t+\gamma-1}{(t+\gamma)^2}\nu + \left[\frac{\beta^2(\tilde{F}+\tilde{D})}{(t+\gamma)^2} - \frac{\beta m - 1}{(t+\gamma)^2}\nu\right] \tag{142}$$

$$\le \frac{t+\gamma-1}{(t+\gamma-1)^2 + 2(t+\gamma)-1}\nu \tag{143}$$

$$\le \frac{t+\gamma-1}{(t+\gamma-1)^2 + 2(t+\gamma-1)}\nu \tag{144}$$

$$\le \frac{\nu}{t+\gamma+1}. \tag{145}$$

Eq. [143] also arises from the definition of $\nu$ that $\beta^2(\tilde{F}+\tilde{D}) \le (\beta m - 1)\nu$. Accordingly, for all $t$, we have $\Delta_t \le \frac{\nu}{\gamma+t}$ holds. $\square$

With this result, we prove that $\Delta_t$ converges to 0 with convergence rate $\mathcal{O}(\frac{1}{T})$, and thus we can say that the proxy algorithm of FedCor converges to the global optimal with convergence rate $\mathcal{O}(\frac{1}{T})$ with Corollary [3].

## C. Experiment Details

We simulate the training process of federated learning on one machine. All experiments in this paper are run on one NVIDIA 2080-Ti GPU and two Intel Xeon E5-2630 v4 CPUs. The experiments on FMNIST require around 3 hours for each seed, and the experiments on CIFAR-10 require around 10 hours for each seed.

### C.1. Model Parameters

**Hyperparameters in FMNIST**   We follow [4] to construct the neural model on FMNIST: An MLP model with two hidden layers with 64 and 30 units, respectively. Under all three heterogeneous settings, we set the local batch size $B = 64$ and the number of local iterations $E = 20$. The learning rate $\eta_0$ is set to 0.005 initially, and halved at the 150-th and 300-th rounds. An SGD optimizer with a weight decay of 0.0001 and no momentum is used. We allocate data to $N = 100$ clients, and set the participation fraction $C = 10$ for the 1SPC setting, and $C = 5$ for the 2SPC and Dir settings.

**Hyperparameters in CIFAR-10**   We use a CNN with three convolutional layers [29] with 32, 64 and 64 kernels, respectively. And all convolution kernels are of size $3 \times 3$. Finally, the outputs of convolutional layers are fed into a fully-connected layer with 64 units. Under all three heterogeneous settings, we set the local batch size $B = 50$ and the

number of local iterations $E = 40$. We use a learning rate $\eta = 0.01$ without learning rate decay, and a weight decay of 0.0003 for the SGD optimizer. The total number of clients and the client participation fraction are the same as those in FMNIST.

**Hyperparameters for FedCor** We set the dimension of client embedding $d = 15$ for all experiments. In Eq. (16), we set $M = 10, S = 1$ for the warm-up phase, and $M = 1, S = 1$ for the normal phase. And we set the discount factor $\gamma = \theta^{\Delta t}$ where $\theta = 0.9$ for experiments on FMNIST and $\theta = 0.99$ for experiments on CIFAR-10. In each GP update round $t$, we use $\boldsymbol{X}^{t-1}$ as the initialization and use an Adam optimizer [11] with learning rate 0.01 to optimize for $\boldsymbol{X}^t$. Notice that although Eq. (16) has a closed form optimal solution for $\boldsymbol{X}^t$, we still learn $\boldsymbol{X}^t$ with the gradient decent method with the initialization $\boldsymbol{X}^{t-1}$ in order to utilize the covariance stationarity and reduce the evaluation bias with small number of samples.

**Hyperparameters for other baselines** We use the same parameters $\alpha_1 = 0.75, \alpha_2 = 0.01$ and $\alpha_3 = 0.1$ as those in the paper [6] for Active Federated Learning. And we set $d = 2NC$ for Power-of-choice Selection Strategy, which is empirically shown to be the best value of $d$ in a highly heterogeneous setting in the paper [4].

Note that we implement the random selection strategy as uniformly sampling clients from $\mathbb{U}$ without replacement [23], while Cho et al. [4] implement the random selection strategy as sampling clients with replacement. Thus, our implemented random selection strategy achieves better performances than their implementation.

### C.2. Dirichlet Distribution for Data Partition

We follow the idea in [7] to construct the Dir heterogeneous setting, while we make some modifications to get an unbalanced non-identical data distribution.

For each client $k$, we sample the data distribution $\boldsymbol{q}_k \in \mathbb{R}^{10}$ from a dirichlet distribution independently, which could be formulated as

$$\boldsymbol{q}_k \sim \text{Dir}(\alpha \boldsymbol{p}), \tag{146}$$

where $\boldsymbol{p}$ is the prior label distribution and $\alpha \in \mathbb{R}_+$ is the concentration parameter of the dirichlet distribution. We group $\boldsymbol{q}_k$ of all the clients together and get a fraction matrix $Q = [\boldsymbol{q}_1, \cdots, \boldsymbol{q}_n]$. We denote the size of dataset on each client as $\boldsymbol{x} = [x_1, \cdots, x_N]^T$ and we get it from a solution of a quadratic programming:

$$\min_{\boldsymbol{x}} \quad \boldsymbol{x}^T \boldsymbol{x} \tag{147}$$

$$\text{subject to} \quad Q\boldsymbol{x} = \boldsymbol{d} \tag{148}$$

$$\boldsymbol{x} \in \mathbb{R}_{++}^N, \tag{149}$$

where $\boldsymbol{d}$ is the number of data with each label. We minimize $\|\boldsymbol{x}\|_2$ to avoid the cases where data distribution is over-concentrated on a small fraction of clients. In that case, the client selection problem might become trivial, since we can always ignore those clients with a small dataset and select those with a large dataset.

## D. Extra Experimental Results

### D.1. Ablation Study: Annealing Coefficient

We conduct experiments on FMNIST and CIFAR-10 with different annealing coefficient $\beta$. We setup our experiments under three heterogeneous settings as in Section 5, with different annealing coefficient $\beta$ ($\beta = 0.95, 0.75, 0.5$ for FMNIST and $\beta = 0.97, 0.95, 0.9$ for CIFAR-10). We fix the GP training interval $\Delta t$ to 10 for FMNIST and 50 for CIFAR-10. The test accuracy curves are shown in Figure 8. We can see that within a large range, the value of annealing coefficient only slightly influence the convergence rate as well as the final accuracy. Recalling the results of different GP training intervals $\Delta t$ in Section 5.3, we can say that our method is not sensitive to the hyperparameters $\Delta t$ and $\beta$.

We present the selected frequency of each client in Figure 10 and Figure 11 for FMNIST and CIFAR-10 respectively. We can see that with a smaller $\beta$, the selected frequency tends to be more "uniform". However, this does not mean that our selection strategy is equivalent to the uniformly random selection. Our sequential selection strategy introduces dependencies between selected clients as discussed in the multi-iteration insights in Section 4.3, which makes our selection strategy prefer some combinations of selected clients to others, while the uniformly random selection treats all the combinations equally. The advantage shown in Figure 8 compared to the uniformly random strategy demonstrates that selecting a good combination of clients, not only a good individual, is important.

### D.2. Normality Verification

We setup experiments to show that Gaussian Distribution can model the loss changes w.r.t. uniformly sampled client selection. To verify this, in the last round of the warm-up phase, we perform the following procedure to examine the normality.

1. We uniformly sample 1000 different client selections $\{\mathbb{S}_{t,i} : i = 1, \cdots, 1000\}$ and collect the corresponding loss changes $\Delta \boldsymbol{l}^t(\mathbb{S}_{t,i}) = [\Delta l_1^t(\mathbb{S}_{t,i}), \cdots, \Delta l_N^t(\mathbb{S}_{t,i})]$ for each of them.

2. We perform PCA on $\{\Delta \boldsymbol{l}^t(\mathbb{S}_{t,i}) : i = 1, \cdots, 1000\}$ to extract the principle components.

3. We plot the histogram of each principle component and compare its distribution with the Gaussian Distribution.
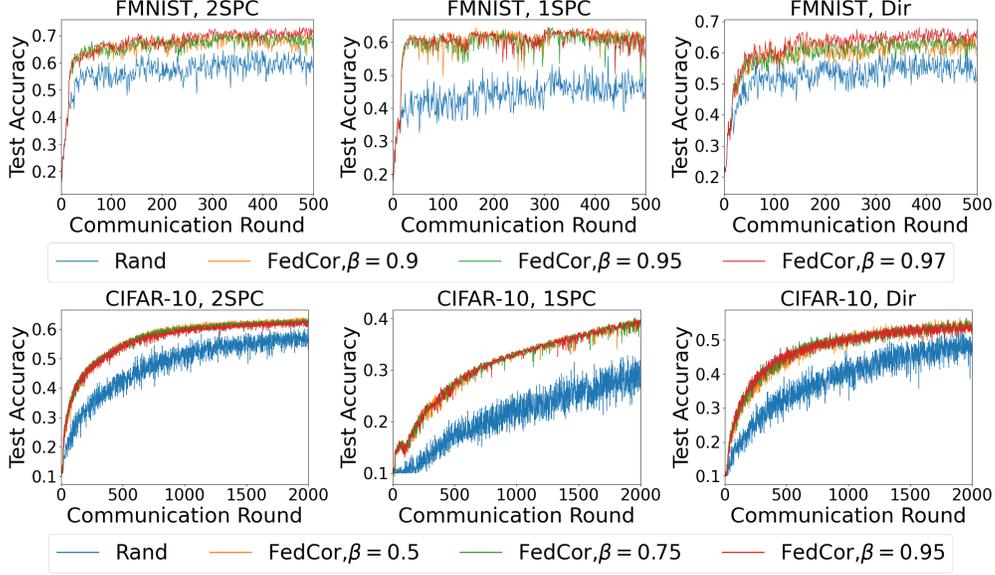
Figure 8. Test accuracy with different annealing coefficient $\beta$ on FMNIST (top) and CIFAR-10 (bottom) under three heterogeneous settings (left: 2SPC; median: 1SPC; right: Dir).
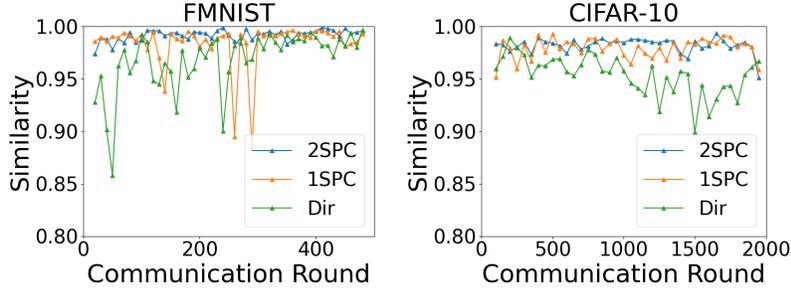


Figure 9. Verification of covariance stationarity on FMNIST and CIFAR-10.

We do not use Multivariate Normality Test directly because we find that $\mathbf{\Sigma^t}$ is always nearly singular, which makes the Multivariate Normality Test unstable. Thus, we turn to perform PCA and visualize each principle component to verify the normality.

The results of FMNIST and CIFAR-10 are shown in Figure 12 and Figure 13 respectively. The red line shows the probability density of Gaussian Distribution with the mean and variance of that principle component. We can see that in all our experiments, Gaussian Distribution can fit the distribution of the principle component well, which verifies that Lemma 1 does hold in all the experiment settings.

### D.3. Covariance Stationarity Verification

We examine that assumption in Section 4.5 that the covariance keep approximately stationary during the FL training, namely,

$$\forall t, \mathbf{\Sigma}^t \approx \mathbf{\Sigma}^{t+\Delta t}. \tag{150}$$

To verify this, every $\Delta t$ rounds ($\Delta t = 10$ for FMNIST and $\Delta t = 50$ for CIFAR-10), we randomly sample 1000 client selections $\mathbb{K}_i$ and collect the corresponding loss changes $\Delta \boldsymbol{l}^t(\mathbb{K}_i)$. We directly calculate the covariance matrix $\mathbf{\Sigma}^t$ with these samples $\{\Delta \boldsymbol{l}^t(\mathbb{K}_i) : i = 1, \cdots, 1000\}$. Then for each adjacent pair of covariance matrix, we calculate their cosine similarity as follows.

$$\text{similarity}(\Sigma^t, \Sigma^{t+\Delta t}) = \frac{\text{tr}(\Sigma^{t^T} \Sigma^{t+\Delta t})}{\text{tr}(\Sigma^{t^T} \Sigma^t) \text{tr}(\Sigma^{t+\Delta t^T} \Sigma^{t+\Delta t})} \tag{151}$$

The similarity is in range $[0, 1]$, and a larger one shows a higher similarity.

The results are shown in Figure 9. We can see that in most cases the similarity is larger than $0.9$, which verifies our claim of the covariance stationarity.
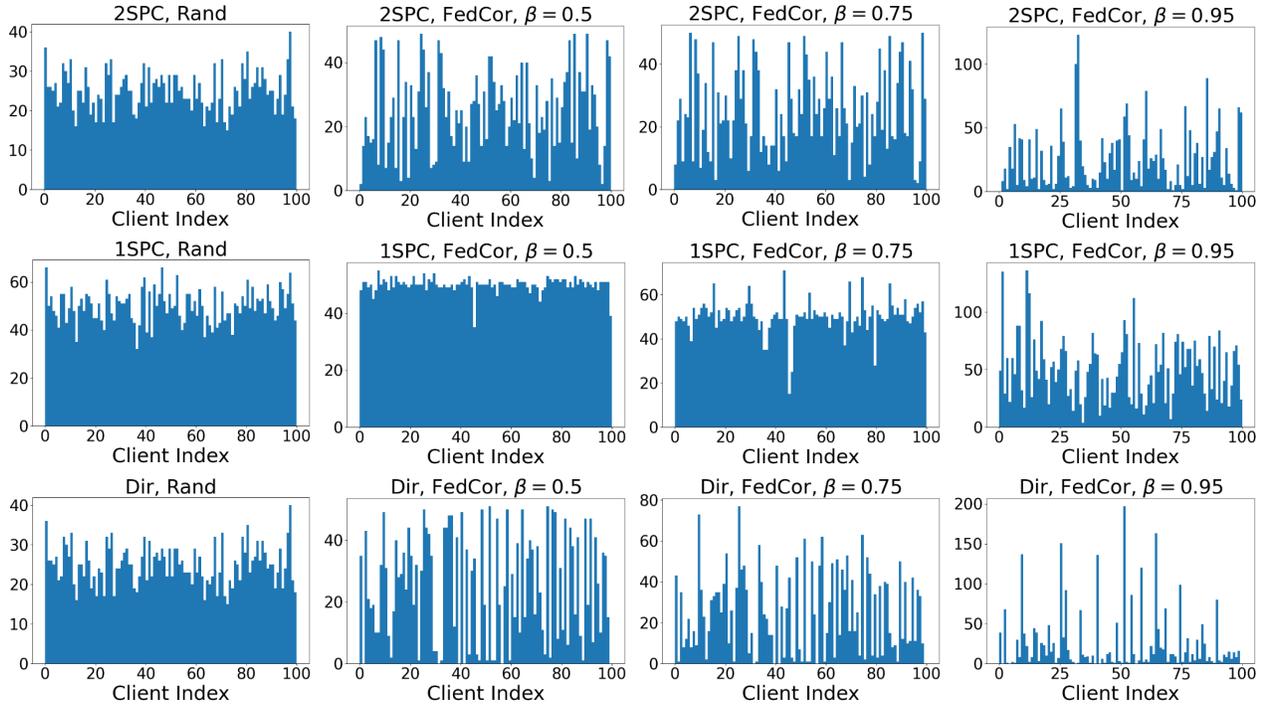
Figure 10. Selected Frequency of each client with different annealing coefficient $\beta$ on FMNIST under three heterogeneous settings (top: 2SPC; median: 1SPC; bottom: Dir).
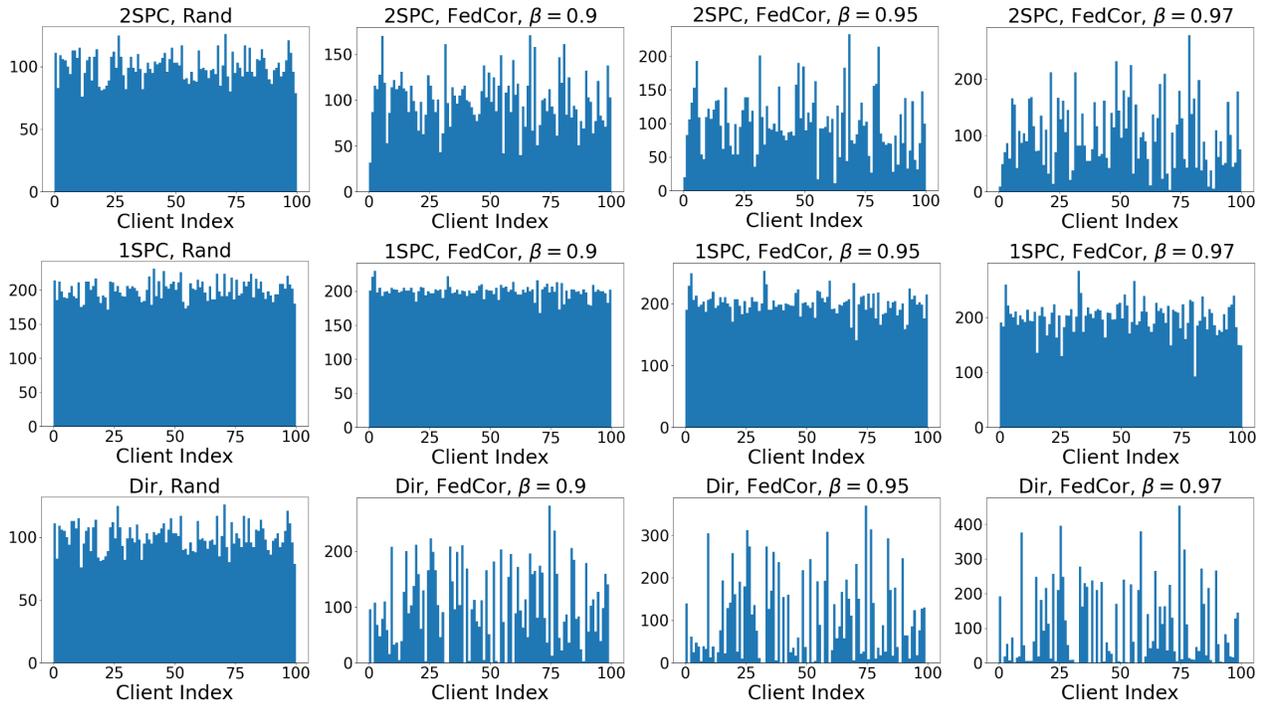


Figure 11. Selected Frequency of each client with different annealing coefficient $\beta$ on CIFAR-10 under three heterogeneous settings (top: 2SPC; median: 1SPC; bottom: Dir).
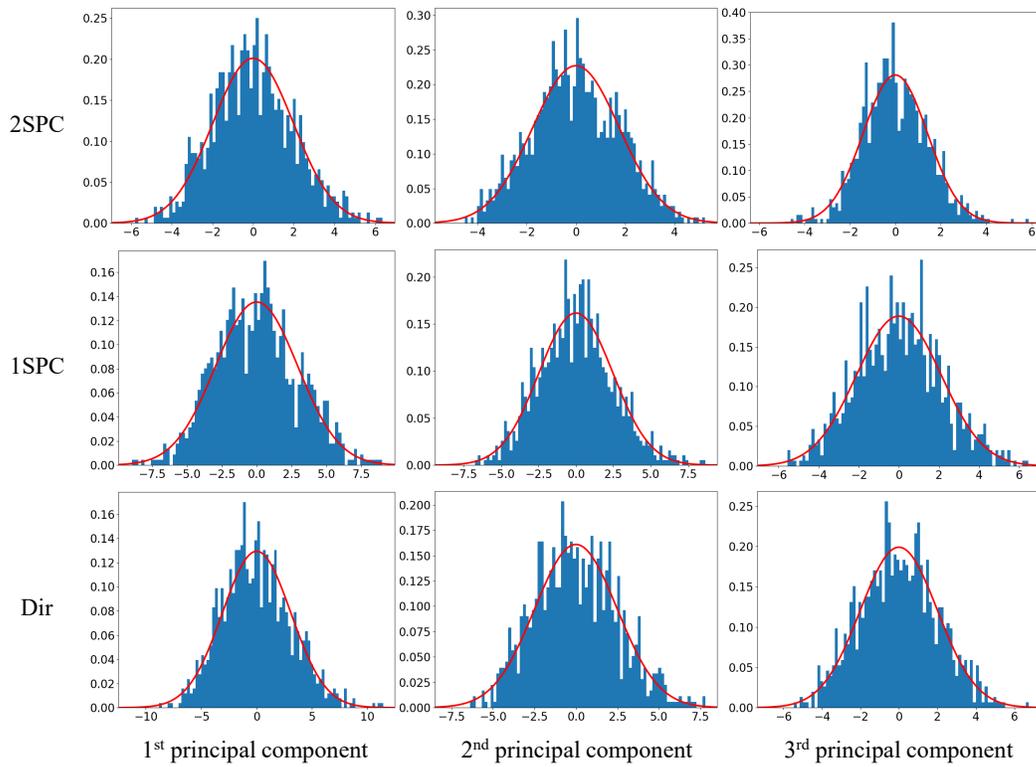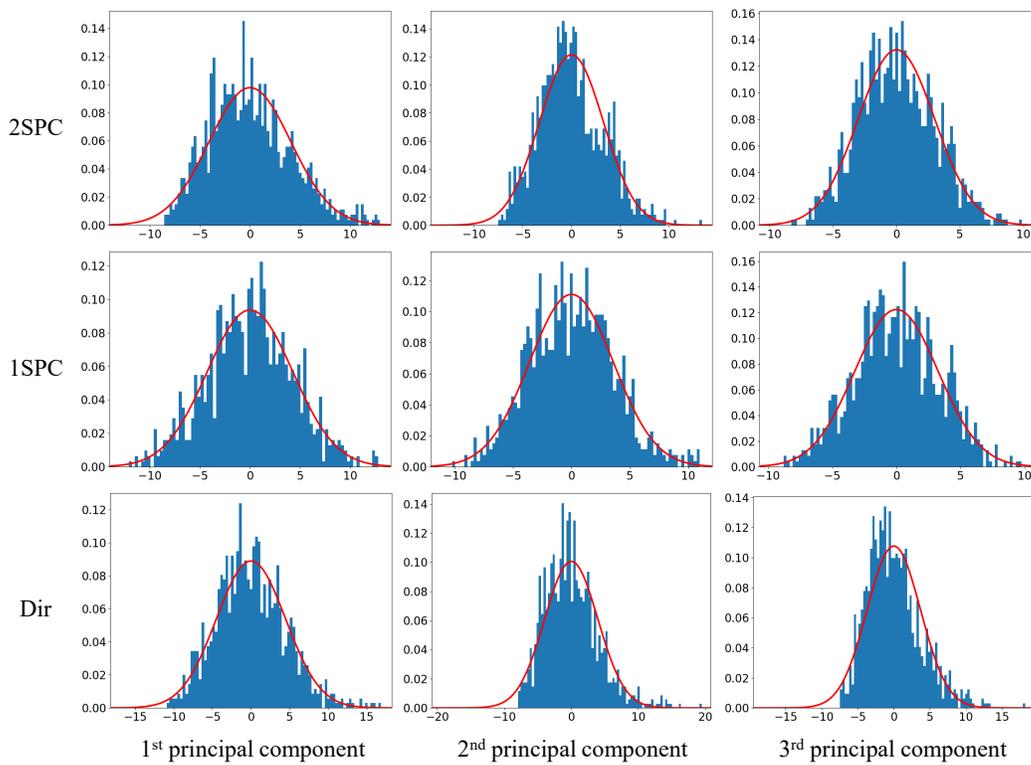
Figure 12. Normality Test on FMNIST.



Figure 13. Normality Test on CIFAR-10.