Performative Federated Learning: A Solution to Model-Dependent and Heterogeneous Distribution Shifts

A PREPRINT

Kun Jin *
University of Michigan
kunj@umich.edu

Tongxin Yin *
University of Michigan
tyin@umich.edu

Zhongzhu Chen *
University of Michigan
zhongzhc@umich.edu

Zeyu Sun
University of Michigan
zeyusun@umich.edu

Xueru Zhang

The Ohio State University zhang.12807@osu.edu

Yang Liu
University of California, Santa Cruz
ByteDance AI Lab
yangliu@ucsc.edu

Mingyan Liu University of Michigan mingyan@umich.edu

May 10, 2023

ABSTRACT

We consider a federated learning (FL) system consisting of multiple clients and a server, where the clients aim to collaboratively learn a common decision model from their distributed data. Unlike the conventional FL framework that assumes the client's data is static, we consider scenarios where the clients' data distributions may be reshaped by the deployed decision model. In this work, we leverage the idea of distribution shift mappings in *performative prediction* to formalize this model-dependent data distribution shift and propose a performative federated learning framework. We first introduce necessary and sufficient conditions for the existence of a unique performative stable solution and characterize its distance to the performative optimal solution. Then we propose the performative FedAvg algorithm and show that it converges to the performative stable solution at a rate of $\mathcal{O}(1/T)$ under both full and partial participation schemes. In particular, we use novel proof techniques and show how the clients' heterogeneity influences the convergence. Numerical results validate our analysis and provide valuable insights into real-world applications.

1 Introduction

Traditional learning problems typically assume data distributions to be static. For applications such as face recognition, this is largely true and designing algorithms under such an assumption in general does not impact learning efficacy. This, however, is not true in many other domains. In some cases, there may be a natural evolution and shift in the distribution, e.g., in weather and climate data, in which case new data need to be acquired periodically and the algorithm re-trained to remain up to date. In other cases, the distribution shift is the result of the very learning outcome, when individuals respond to the algorithmic decisions they are subjected to. For instance, when users with certain accents perceive larger-than-acceptable errors from a speech recognition software and therefore stop using it, this can directly impact the type of speech samples collected by the software used for training the next generation of the product. Another example is "gaming the algorithm", where users through honest or dishonest means attempt to improve critical features so as to obtain a favorable decision by the algorithm (e.g., in loan approvals or job applications). This again can directly lead to the distributional change in features and label that the algorithm relies on for decision making.

^{*}These authors contributed equally to this work.

This latter type of distribution shifts, one of particular interest, where the deployed model itself can trigger changes in the data distribution and influence the objective, said to be *performative*. Performing prediction in the presence of such distribution shift is called *performative prediction* Perdomo et al. [2020]. Typical scenarios of performative prediction include *strategic learning* Hardt et al. [2016], Dong et al. [2018], Milli et al. [2019], Hu et al. [2019], Braverman and Garg [2020], Chen et al. [2020], Miller et al. [2020], Shavit et al. [2020], Haghtalab et al. [2020], Kleinberg and Raghavan [2020], Zrnic et al. [2021].

Performative prediction has been primarily studied in a centralized setting, with fruitful literature including the convergence analysis Mendler-Dünner et al. [2020], Drusvyatskiy and Xiao [2020], Brown et al. [2020], Li and Wai [2022], Wood et al. [2022] and algorithm development Izzo et al. [2021, 2022], Miller et al. [2021], Ray et al. [2022].

In modern large-scale machine learning, distributed learning offers greater privacy protection and better avoids the computational resource bottlenecks compared to centralized learning, and federated learning (FL) is one of the most popular examples. Here the issue of distribution shift is further compounded due to data heterogeneity in a distributed setting. Specifically, the distributed data sources can be heterogeneous in nature, and their respective distribution shifts can also be different. Prior works in FL systems that address data distribution shifts, such as Guo et al. [2021], Casado et al. [2022], Rizk et al. [2020], Hosseinalipour et al. [2022], Zhu et al. [2021a], Eichner et al. [2019], Ding et al. [2020], typically do not consider shifts in local distributions at the client end induced by the model. In this work, we propose the *performative federated learning* framework to study and handle such data shifts in FL.

Extending the current results in performative prediction to the decentralized FL has a number of challenges. To highlight a few: 1) *Data heterogeneity:* As already one of the major difficulties in FL, tackling data heterogeneity faces additional challenges when taking the disparity of client distribution shift into consideration. 2) *Central* \rightleftharpoons *Local:* During training, clients receive the aggregated model at certain steps and train from it. While fitting better as an entity, such aggregation may fail to fit well on each client, which may lead to more severe shifting issues. 3) *heterogeneity in shift:* some clients may be more sensitive to the deployed decisions and have more drastic data shifts than other clients, e.g., due to different manipulation costs in strategic learning.

Toward this end, we formally introduce the *performative FedAvg* algorithm, or P-FedAvg, and establish its convergence. Our main findings are as follows.

- We prove the uniqueness of the performative stable (PS) solution reached by the algorithm, and show that it is a provable approximation to the performative optimal (PO) solution under mild conditions. Both solutions will be formally defined in Section 2.1. More interestingly, we show that the stable solution has its own game-theoretic interpretation as the fixed point of the best response dynamics.
- We show in Section 3.3 that the P-FedAvg algorithm converges to the performative stable solution and has a $\mathcal{O}(1/T)$ convergence rate with both the full and partial participation schemes under mild assumptions similar to those in prior works.
- In doing so we also introduce some novel proof techniques: we prove convergence without a bounded gradient assumption. This technique can be directly applied to conventional FL, which is a special case of the performative setting.

1.1 Related Works

Federated Learning. Our work is strongly related to the literature on federated learning (FL). Although many studies have tried to address client heterogeneity in FL through constrained gradient optimization and knowledge distillation Li et al. [2020a], Karimireddy et al. [2020], Wang et al. [2020], Haddadpour et al. [2021], Zhu et al. [2021b], Li and Wang [2019], Lin et al. [2020], most of them still assume the data is static without considering the distribution shifts. To the best of our knowledge, only a few recent works consider distribution shifts in FL Guo et al. [2021], Casado et al. [2022], Rizk et al. [2020], Hosseinalipour et al. [2022], Zhu et al. [2021a], Eichner et al. [2019], Ding et al. [2020]. For example, Guo et al. [2021] considered FL with time-evolving clients where the time-drift of each client is modeled as a time-independent additive noise with zero-mean and bounded variance. Casado et al. [2022] proposed an FL algorithm adaptable to distribution drifts; it monitors the confidence scores of the model prediction throughout the learning process and assumes the drift happens whenever there is a substantial drop in confidence scores. Rizk et al. [2020] also studied dynamic FL and assumed the true model under time-evolving data follows a random walk. Hosseinalipour et al. [2022] considered FL with dynamic clients and modeled the drift using the variation in local loss over two consecutive time steps. Zhu et al. [2021a], Eichner et al. [2019], Ding et al. [2020] considered the periodical distribution shift of client

population in FL; they assume the *block-cyclic* structure where the clients from two different time zones alternately participate in training.

Performative Prediction. In addition to the one we discussed in the introduction that focuses on the centralized setting for performative prediction, more recently, Li et al. [2022] formalize the multi-agent/player performative predictions where agents try to learn a common decision rule but have heterogeneous distribution shifts (responses) to the model, and study the convergence of decentralized algorithms to the PS solution. The decentralized performative predictions capture the heterogeneity in agents'/clients' responses to the decision model and avoid centralized data collection for training. This work provides inspiration for our formulation of the performative federated learning framework, and our proposed P-FedAvg can be viewed as a substantial algorithmic extension that supports unbalanced data, much less frequent synchronizations, and partial device participation. Narang et al. [2022] propose a decentralized multi-player performative prediction framework where the players react to competing institutions' actions. Raab and Liu [2021] proposes a replicator dynamics model with label shift.

Strategic Classification and Regression. As discussed in Perdomo et al. [2020], performative prediction can be used to solve repeated strategic classification and regression problems. We can use Stackelberg games to model these problems, where the decision maker moves in the first stage by designing, publishing, and committing to a decision rule, then the agents move in the second stage, best responding to the decision rule by manipulating their features to get more desirable decision outcomes, and such manipulation can be modeled by the distribution shift mappings. Conventional strategic learning literature focus on finding the Stackelberg equilibrium Hardt et al. [2016], Kleinberg and Raghavan [2020], Shavit et al. [2020], Haghtalab et al. [2020], i.e., the PO solution where the decision maker and the agents know each others' utilities, whereas performative prediction can find the PS solution in repeated strategic learning problems regardless of the knowledge on the utilities.

2 Problem Formulation

In this section, we formulate the performative federated learning problem, define the learning objective, and introduce our performative federated learning algorithm to optimize the objective function.

To help with the understanding of performative federated learning, we first recall the performative prediction problem Perdomo et al. [2020]. Consider a typical loss minimization problem where the data distribution experiences a shift induced by the model parameter, expressed as a mapping $\mathcal{D}(\theta)$. The objective function is thus given by

$$f(\boldsymbol{\theta}) := \mathbb{E}_{Z \sim \mathcal{D}(\boldsymbol{\theta})} [\ell(\boldsymbol{\theta}; Z)],$$

where ℓ denotes the loss function. Then the performative optimal (PO) solution is $\boldsymbol{\theta}^{PO} := \arg\min_{\boldsymbol{\theta}} f(\boldsymbol{\theta})$. Perdomo et al. [2020] also introduces a second, decoupled objective function, also called the performatively stable (PS) model, which separates decision parameters $(\boldsymbol{\theta})$ from deployed parameters $(\tilde{\boldsymbol{\theta}})$:

$$f(\boldsymbol{\theta}; \tilde{\boldsymbol{\theta}}) := \mathbb{E}_{Z \sim \mathcal{D}(\tilde{\boldsymbol{\theta}})}[\ell(\boldsymbol{\theta}; Z)].$$

Minimizing this objective achieves minimal risk for the distribution induced by the deployed parameters, eliminating the need for retraining, which makes it more practical. The PS solution is defined as $\theta^{PS} := \arg\min_{\theta} f(\theta; \theta^{PS})$. Perdomo et al. [2020] showed that $\theta^{PS} \neq \theta^{PO}$ in general. We next consider a distributed setting and introduce performative federated learning.

2.1 System Settings and Objectives

Consider a system with N clients and a server, where the clients have feature distributions as $\mathcal{D}_i(\boldsymbol{\theta})$, supported on $\mathcal{Z} \subseteq \mathbb{R}^M$, and $\boldsymbol{\theta} \in \mathbb{R}^m$ denotes the decision (model) parameters deployed on the i-th client. We consider the general case where clients can have heterogeneous distributions $\mathcal{D}_i(\boldsymbol{\theta}) \neq \mathcal{D}_j(\boldsymbol{\theta})$, and each client represents a $p_i > 0$ fraction of the total data population, $\sum_{i=1}^N p_i = 1$.

The system aims to minimize the weighted average loss across all agents, which is given by the performative optimal objective as follows

$$\boldsymbol{\theta}^{PO} := \arg\min_{\boldsymbol{\theta} \in \mathbb{R}^m} \sum_{i=1}^{N} p_i \mathbb{E}_{Z_i \sim \mathcal{D}_i(\boldsymbol{\theta})} [\ell(\boldsymbol{\theta}; Z_i)]. \tag{1}$$

We note that our objective reduces to that in Li et al. [2022] when $p_i = \frac{1}{N}$ for all i. This objective can typically model the strategic learning problem with different sub-populations in the system, where each client corresponds to a sub-population. Each sub-population may differ in some attributes so that they respond to the decision parameters differently, e.g., due to different action costs Milli et al. [2019], Hu et al. [2019], Braverman and Garg [2020], Zhang et al. [2022], Jin et al. [2022]. The decision maker uses a common decision rule for the entire population and aims to minimize the expected loss, and p_i represents the population fraction of each sub-population. Correspondingly, the decoupled/performative stable objective is

$$f_i(\boldsymbol{\theta}; \tilde{\boldsymbol{\theta}}) := \mathbb{E}_{Z_i \sim \mathcal{D}_i(\tilde{\boldsymbol{\theta}})}[\ell(\boldsymbol{\theta}; Z_i)], f(\boldsymbol{\theta}; \tilde{\boldsymbol{\theta}}) := \sum_{i=1}^N p_i f_i(\boldsymbol{\theta}; \tilde{\boldsymbol{\theta}}),$$

where the first argument denotes the client's decision parameter, and the second argument the deployed parameters, which determine the distribution of the samples together with $\mathcal{D}_i(\cdot)$. The PS solution is

$$\boldsymbol{\theta}^{PS} := \arg\min_{\boldsymbol{\theta}} \sum_{i=1}^{N} p_{i} \mathbb{E}_{Z_{i} \sim \mathcal{D}_{i}(\boldsymbol{\theta}^{PS})} [\ell(\boldsymbol{\theta}; Z_{i})]$$

$$= \arg\min_{\boldsymbol{\theta}} f(\boldsymbol{\theta}; \boldsymbol{\theta}^{PS}).$$
(2)

Note that this is a fixed point equation with θ^{PS} as a fixed point.

2.2 Key Assumptions

We make the following assumptions similar to Li et al. [2022], Perdomo et al. [2020]:

Assumption 2.1 (Strong Convexity). Given any $\tilde{\boldsymbol{\theta}} \in \mathbb{R}^m$, $f(\cdot, \tilde{\boldsymbol{\theta}})$ is μ -strongly convex in $\boldsymbol{\theta}$, i.e., $f(\boldsymbol{\theta}'; \tilde{\boldsymbol{\theta}}) \geq f(\boldsymbol{\theta}; \tilde{\boldsymbol{\theta}}) + \langle \nabla f(\boldsymbol{\theta}; \tilde{\boldsymbol{\theta}}), \boldsymbol{\theta}' - \boldsymbol{\theta} \rangle + \frac{\mu}{2} \|\boldsymbol{\theta}' - \boldsymbol{\theta}\|_2^2, \forall \boldsymbol{\theta}', \boldsymbol{\theta} \in \mathbb{R}^K$.

In Assumption 2.1, we do not require strong convexity for every single f_i but only the weighted average f.

Assumption 2.2 (Smoothness). The loss function $\ell(\theta; z)$ is L-smooth, i.e.,

$$\|\nabla \ell(\boldsymbol{\theta}; \boldsymbol{z}) - \nabla \ell(\boldsymbol{\theta}'; \boldsymbol{z}')\|_2 \le L(\|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_2 + \|\boldsymbol{z} - \boldsymbol{z}'\|_2).$$

Assumption 2.3 (Distribution Mapping Sensitivity). For any $i = 1, \ldots, n$ there exists $\epsilon_i > 0$ such that

$$W_1(\mathcal{D}_i(\boldsymbol{\theta}), \mathcal{D}_i(\boldsymbol{\theta}')) < \epsilon_i \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_2, \ \forall \boldsymbol{\theta}', \boldsymbol{\theta} \in \mathbb{R}^m,$$

where $W_1(\mathcal{D}, \mathcal{D}')$ is the 1-Wasserstein distance under L_2 norm between the distributions $\mathcal{D}, \mathcal{D}'$.

Assumption 2.2 and 2.3 together induce the smoothness of $f_i(\cdot, \cdot)$, which is a result of Lemma 2.1 in Drusvyatskiy and Xiao [2022] and will be used in the later proofs.

Lemma 2.4 (Continuity of ∇f_i). Under Assumption 2.2 and 2.3, for any $\theta_0, \theta_1, \theta, \hat{\theta} \in \mathbb{R}^m$,

$$\|\nabla f_i(\boldsymbol{\theta}_0; \boldsymbol{\theta}) - \nabla f_i(\boldsymbol{\theta}_1; \hat{\boldsymbol{\theta}})\|_2 \le L\|\boldsymbol{\theta}_0 - \boldsymbol{\theta}_1\|_2 + L\epsilon_i\|\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}\|_2.$$

We introduce the following assumptions specifically made in decentralized performative predictions Li et al. [2022].

Assumption 2.5 (Stochastic Gradient Variance Bound). For any i = 1, ..., N and $\theta \in \mathbb{R}^m$, there exists $\sigma \ge 0$ such that

$$\mathbb{E}_{Z_i \sim \mathcal{D}_i(\boldsymbol{\theta})} \|\nabla \ell(\boldsymbol{\theta}; Z_i) - \nabla f_i(\boldsymbol{\theta}; \boldsymbol{\theta})\|_2^2 \le \sigma^2 (1 + \|\boldsymbol{\theta} - \boldsymbol{\theta}^{PS}\|_2^2).$$

Assumption 2.6 (Local Gradient Bound). For any i = 1, ..., N and $\theta \in \mathbb{R}^m$, there exists $\varsigma \ge 0$ such that

$$\|\nabla f(\boldsymbol{\theta}; \boldsymbol{\theta}) - \nabla f_i(\boldsymbol{\theta}; \boldsymbol{\theta})\|_2^2 \le \varsigma^2 (1 + \|\boldsymbol{\theta} - \boldsymbol{\theta}^{PS}\|_2^2).$$

Here we elaborate on Assumption 2.6, and explain reasons for using it over another commonly used assumption in federated learning Li et al. [2020b], which is

$$\mathbb{E}_{Z_i \sim \mathcal{D}_i(\boldsymbol{\theta})}[\|\nabla \ell(\boldsymbol{\theta}; Z_i)\|_2^2] \le G^2. \tag{3}$$

First, it can be shown that (3) implies Assumption 2.6, thus Assumption 2.6 is weaker than (3). To see this: when (3) holds, let $\varsigma^2 = 4G^2$, then $\|\nabla f(\boldsymbol{\theta}; \boldsymbol{\theta}) - \nabla f_i(\boldsymbol{\theta}; \boldsymbol{\theta})\|_2^2 \le 2 \|\nabla f(\boldsymbol{\theta}; \boldsymbol{\theta})\|_2^2 + 2\|\nabla f_i(\boldsymbol{\theta}; \boldsymbol{\theta})\|_2^2 \le 4G^2 = \varsigma^2$.

We further give a concrete example where (3) does not hold but Assumption 2.6 holds.

Example 2.7. Suppose we have a two-client Gaussian mean estimation problem $\ell(\theta, Z) = \frac{1}{2}(\theta - Z)^2$ where $\theta, Z \in \mathbb{R}$, $\mathcal{D}_1(\theta) = \mathcal{N}(\frac{1}{2}\theta, \sigma^2)$, $\mathcal{D}_2(\theta) = \mathcal{N}(-\frac{1}{2}\theta, \sigma^2)$, and $p_1 = p_2 = \frac{1}{2}$. Then $\mathbb{E}_{Z_1 \sim \mathcal{D}_1(\theta)}[\|\nabla \ell(\theta; Z_1)\|_2^2] = \mathbb{E}_{Z_1 \sim \mathcal{D}_1(\theta)}[(\theta - Z_1)^2] = \sigma^2 + (\mathbb{E}_{Z_1 \sim \mathcal{D}_1(\theta)}[\theta - Z_1])^2 = \frac{1}{4}\theta^2 + \sigma^2$ and $\mathbb{E}_{Z_2 \sim \mathcal{D}_2(\theta)}[\|\nabla \ell(\theta; Z_2)\|_2^2] = \frac{9}{4}\theta^2 + \sigma^2$ which all go to infinity when θ goes to infinity. Thus (3) does not hold. On the other hand, $\nabla f_1(\theta; \theta) = \frac{1}{8}\theta$, $\nabla f_2(\theta; \theta) = \frac{9}{8}\theta$, and $\nabla f(\theta; \theta) = \frac{5}{8}\theta$, $\theta^{PS} = 0$, then by taking $\varsigma = \frac{1}{2}$, we can verify Assumption 2.6 holds.

Secondly, (3) also implies Assumption 2.5: when (3) holds, letting $\sigma^2 = G^2$ leads to $\mathbb{E}[\|\nabla l(\boldsymbol{\theta}; Z_i) - \nabla f_i(\boldsymbol{\theta}, \boldsymbol{\theta})\|_2^2] \leq \mathbb{E}[\|\nabla l(\boldsymbol{\theta}; Z_i)\|_2^2] \leq G^2 = \sigma^2$. On the other hand, Assumption 2.6 does not imply Assumption 2.5.

It turns out that Assumption 2.6 better characterizes the system heterogeneity, as we show how the heterogeneity impacts convergence (more details are in Theorem 3.1, 3.2, and 3.3).

2.3 Properties of the PS Solution

Define the average sensitivity as $\bar{\epsilon} := \sum_{i=1}^N p_i \epsilon_i$, and the mapping $\Phi(\boldsymbol{\theta}) := \arg\min_{\boldsymbol{\theta}' \in \mathbb{R}^m} f(\boldsymbol{\theta}', \boldsymbol{\theta})$. Then we can establish the existence and uniqueness of the PS solution.

Proposition 2.8 (Uniqueness of θ^{PS}). Under Assumptions 2.1, 2.2 and 2.3, if $\overline{\epsilon} < \mu/L$, then $\Phi(\cdot)$ is a contraction mapping with the unique fixed point $\theta^{PS} = \Phi(\theta^{PS})$; if $\overline{\epsilon} \ge \mu/L$, then there is an instance where any sequence generated by $\Phi(\cdot)$ will diverge.

Proposition 2.8 establishes a sufficient and necessary condition for the existence of θ^{PS} , similar to Li et al. [2022]. This condition only depends on the average sensitivity $\bar{\epsilon}$, which implies that we may still have a unique performative stable solution θ^{PS} for the whole system even if certain clients do not. The following proposition further validates the quality of θ^{PS} in terms of its distance to θ^{PO} .

Proposition 2.9 (Distance $\|\boldsymbol{\theta}^{PO} - \boldsymbol{\theta}^{PS}\|_2$ Bound). Under Assumption 2.1 and 2.3, suppose that the loss $\ell(\boldsymbol{\theta}; Z)$ is L_z -Lipschitz in Z, then for every performative stable solution $\boldsymbol{\theta}^{PS}$ and every performative optimal solution $\boldsymbol{\theta}^{PO}$, we have $\|\boldsymbol{\theta}^{PS} - \boldsymbol{\theta}^{PO}\|_2 \leq (2L_z\overline{\epsilon})/\mu$.

The proofs of Proposition 2.8 and 2.9 are in Appendix A.

2.4 The P-FedAvg Algorithm

In P-FedAvg, the clients communicate with the server every E local updates. Denote $\mathcal{I}_E := \{nE | n=1,2,\dots\}$ as the set of aggregation steps. Next, we formalize the full and partial participation schemes of the proposed **P-FedAvg**.

Full client participation. All clients communicate with the server at every aggregation step and update the local models θ_i^{t+1} based on the following: let $Z_i^{t+1} \sim \mathcal{D}_i(\theta_i^t)$, then

$$\begin{aligned} \boldsymbol{w}_i^{t+1} &= \boldsymbol{\theta}_i^t - \eta_t \nabla \ell(\boldsymbol{\theta}_i^t; Z_i^{t+1}); \\ \boldsymbol{\theta}_i^{t+1} &= \left\{ \begin{array}{cc} \sum_{j=1}^N p_j \boldsymbol{w}_j^{t+1} & \text{if } t+1 \in \mathcal{I}_E \\ \boldsymbol{w}_i^{t+1} & \text{o.w.} \end{array} \right. \end{aligned}$$

Partial client participation. A more realistic setting that does not require the response of all clients' output at every aggregation step. In this case, the central server only collects the outputs of the first K < N responded clients at the aggregation step. Denote the first K < N responded clients in t-th step as a size-K set $\mathcal{S}_t := \{i_1, \dots, i_K\} \in [N]$. Let

$$Z_i^{t+1} \sim \mathcal{D}_i(\boldsymbol{\theta}_i^t)$$
, then

$$\begin{split} \boldsymbol{w}_i^{t+1} &= \boldsymbol{\theta}_i^t - \eta_t \nabla \ell(\boldsymbol{\theta}_i^t; \boldsymbol{Z}_i^{t+1}); \\ \boldsymbol{\theta}_i^{t+1} &= \left\{ \begin{array}{c} \left(\text{samples } \mathcal{S}_{t+1}, \text{and} \\ \text{average } \left\{ \boldsymbol{w}_{t+1}^k \right\}_{k \in \mathcal{S}_{t+1}} \end{array} \right) & \text{if } t+1 \in \mathcal{I}_E \\ \boldsymbol{w}_i^{t+1} & \text{o.w.} \end{split} \end{split}$$

We further consider two schemes of partial participation:

- 1. (Scheme I) The server establishes S_{t+1} by *i.i.d. with replacement* sampling an index $k \in \{1, \dots, N\}$ with probabilities p_1, \dots, p_N for K times. Hence S_{t+1} is a multiset that allows an element to occur more than once. Then the server averages the parameters by $\boldsymbol{\theta}_i^{t+1} = \frac{1}{K} \sum_{k \in S_{t+1}} \boldsymbol{w}_k^{t+1}$. This sampling scheme is first proposed in Sahu et al. [2018] but the theoretical analysis was first done in Li et al. [2020b].
- 2. (Scheme II) The server samples S_{t+1} uniformly without replacement. Hence each element in S_{t+1} only occurs once. Then the server averages the parameters by $\theta_i^{t+1} = \sum_{k \in S_{t+1}} p_k \frac{N}{K} w_k^{t+1}$. Note that when the probabilities $\{p_k\}$ are not the same, one cannot ensure $\sum_{k \in S_{t+1}} p_k \frac{N}{K} = 1$ Li et al. [2020b].

Pseudo-codes of P-FedAvg are in Appendix ??.

Communication cost. The P-FedAvg requires two rounds of communications, aggregation, and broadcast for every E iterations. So at time step T, the system completes $2\lfloor T/E \rfloor$ communications. We follow the setting in Li et al. [2020b] where the server aggregates based on the chosen scheme and broadcasts the aggregated parameters to all clients.

Next, we will prove that P-FedAvg has $\mathcal{O}(1/T)$ convergence rate under the above assumption. As a supplement, we prove in Appendix D that P-FedAvg also has $\mathcal{O}(1/T)$ convergence rate if we replace Assumption 2.6 with (3).

3 Convergence Analysis

In this section, we show that the P-FedAvg converges to the unique θ^{PS} at a rate of $\mathcal{O}(1/T)$ under the assumptions made in Section 2, which holds for all above-introduced schemes. The key observation is that for sufficiently small and decaying learning rates, the effect of E steps is similar to a one-step update with a larger learning rate in the static case, as stated in Li et al. [2020a] without the performative setting. Therefore, given appropriate sampling and updating schemes that satisfy the above assumptions, the global update behaves similarly to the repeated performative SGD in Perdomo et al. [2020]. We also show that partial device participation makes the averaged parameter sequence $\{\overline{\theta}^t\}$ have the same mean as but a larger variance than the full participation, where the variance can be controlled with carefully chosen learning rates. It's worth noting that the heterogeneity of clients plays a key role in the convergence analysis, which we elaborate on below.

Quantifying the heterogeneity. The client heterogeneity can be quantified by the consensus error $\sum_{i=1}^{N} p_i \|\boldsymbol{\theta}_i^t - \overline{\boldsymbol{\theta}}^t\|_2^2$, which is dynamic due to the nature of performative prediction. It depends on both the shift mappings \mathcal{D}_i and the decision parameters. After every broadcast, the heterogeneity leads to heterogeneous distribution shifts, causing heterogeneous local updates and eventually resulting in the consensus error. The ς value in Assumption 2.6 is also a good indicator for heterogeneity.

Next, we will first present the convergence analysis of the full participation scheme and later extend the analysis to partial participation schemes. Due to the complexity of analysis in the performative setting, we define the following constants for ease of analysis and clarity of presentation.

Constants independent of system design.

$$\overline{\epsilon} := \sum_{i=1}^{N} p_i \epsilon_i, \quad \epsilon_{max} := \max_i \epsilon_i,
\widetilde{\mu} := \mu - (1 + \delta)\overline{\epsilon}L, \quad \delta > 0,
c_1 := \left(L(1 + \epsilon_{max})^2\right)/(2\delta\overline{\epsilon}),
c_2 := 4\left[\sigma^2 + L^2(1 + \epsilon_{max})^2\right],
c_3 := 6\left[2\sigma^2 + 3L^2(1 + \epsilon_{max})^2\right],
c_4 := 16\sigma^2 + 12\varsigma^2 + (8\sigma^2 + 12\varsigma^2)/\mathbb{E}\|\overline{\theta}^0 - \theta^{PS}\|_2^2,
c_5 := (48\sigma^2 + 36\varsigma^2)\mathbb{E}\|\overline{\theta}^0 - \theta^{PS}\|_2^2 + (24\sigma^2 + 36\varsigma^2).$$

Constants related to system design (e.g., E, K).

$$c_{6} := (2E^{2} + 3E + 1) \log(E + 1),$$

$$\tilde{\eta}_{0} := \tilde{\mu}/(2\sigma^{2} + (c_{1}c_{3} + c_{2}/6)c_{4}c_{6}),$$

$$\hat{\eta}_{0} := \tilde{\mu}/(2\sigma^{2} + (c_{1}c_{3} + c_{2}/6)c_{4}(2E^{2} - E) \log E),$$

$$B := 2\sigma^{2} + (4c_{1}\hat{\eta}_{0} + 4c_{2}\hat{\eta}_{0}^{2})c_{5}(2E^{2} - E) \log E,$$

$$B_{1} := 2\sigma^{2} + (4c_{1}\tilde{\eta}_{0} + 4c_{2}\tilde{\eta}_{0}^{2} + 1/K)c_{5}c_{6},$$

$$B_{2} := 2\sigma^{2} + (4c_{1}\tilde{\eta}_{0} + 4c_{2}\tilde{\eta}_{0}^{2} + \frac{N - K}{K(N - 1)})c_{5}c_{6}.$$

3.1 Convergence of Full Participation

Theorem 3.1 (Full Participation). Consider P-FedAvg with full participation and diminishing step size $\eta_t = \frac{2}{\tilde{\mu}(t+\gamma)}$, where $\gamma = \max\left\{\frac{2}{\tilde{\mu}\hat{\eta}_0}, E, \frac{2}{\tilde{\mu}}\sqrt{(4E^2+2E)c_3}\right\}$. Under Assumption 2.1, 2.2, 2.3, 2.5, 2.6, the following holds

$$\mathbb{E}[\|\overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}\|_2^2] \le \frac{v}{\gamma + t}, \ \forall t$$

where
$$\upsilon = \max \left\{ \frac{4B}{\overline{\mu}^2}, \gamma \mathbb{E} \| \overline{\boldsymbol{\theta}}^0 - \boldsymbol{\theta}^{PS} \|_2^2 \right\}$$
.

The proof of Theorem 3.1 can be found in Appendix B. The key to the proof is that expected distance $\mathbb{E}\|\overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}\|_2^2$ and expected consensus error $\sum_{i=1}^N p_i \mathbb{E}\|\boldsymbol{\theta}_i^{t-1} - \overline{\boldsymbol{\theta}}^t\|_2^2$ all depend on expected distance $\mathbb{E}\|\overline{\boldsymbol{\theta}}^{t-1} - \boldsymbol{\theta}^{PS}\|_2^2$ and expected consensus error $\sum_{i=1}^N p_i \mathbb{E}\|\boldsymbol{\theta}_i^{t-1} - \overline{\boldsymbol{\theta}}^{t-1}\|_2^2$ in the previous step. While we can establish a descent lemma for expected distance including the expected consensus error, it is impossible to establish one for expected consensus error, which makes it impossible to establish a joint descent lemma for expected distance and expected consensus error as in Li et al. [2022]. Fortunately, consensus error will become zero at every aggregation step, which enables us to control expected consensus error at every step within a constant with a novel double-iteration technique under small enough step sizes. Then by relaxing the expected consensus error to the constant, we can establish a standard descent lemma in SGD analysis for expected distance.

3.2 Convergence of Partial Participation

As mentioned in Section 2, the partial participation scheme is more realistic in federated learning Li et al. [2020b] and is of more interest since it reduces the stragglers' effect.

We first present the convergence result of Scheme I.

Theorem 3.2 (Partial Participation, Scheme I). Consider P-FedAvg with partial participation (scheme I) and a diminishing step size $\eta_t = \frac{2}{\bar{\mu}(t+\gamma)}$, where $\gamma = \max\left\{\frac{2}{\bar{\mu}\tilde{\eta}_0}, E, \frac{2}{\bar{\mu}}\sqrt{(4E^2+10E+6)c_3}\right\}$. Under Assumption 2.1, 2.2, 2.3, 2.5, 2.6, the following holds

$$\mathbb{E}[\|\overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}\|_2^2] \le \frac{\upsilon}{\gamma + t}, \ \forall t$$

where
$$v = \max \left\{ \frac{4B_1}{\overline{\mu}^2}, \gamma \mathbb{E} \| \overline{\boldsymbol{\theta}}^0 - \boldsymbol{\theta}^{PS} \|_2^2 \right\}$$
.

Then we present the convergence result of Scheme II. As discussed in Section 2, we need probabilities $p_i = \frac{1}{N}, \forall i$ to ensure $\sum_{i \in \mathcal{S}_t} p_k \frac{N}{K} = 1$.

Theorem 3.3 (Partial Participation, Scheme II). Consider P-FedAvg with partial participation (scheme II) and a diminishing step size $\frac{2}{\tilde{\mu}(t+\gamma)}$, where $\gamma = \max\left\{\frac{2}{\tilde{\mu}\tilde{\eta}_0}, E, \frac{2}{\tilde{\mu}}\sqrt{(4E^2+10E+6)c_5}\right\}$. Under Assumption 2.1, 2.2, 2.3, 2.5, 2.6, the following holds

$$\mathbb{E}[\|\overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}\|_2^2] \le \frac{\upsilon}{\gamma + t}, \ \forall t$$

where
$$v = \max \left\{ \frac{4B_2}{\tilde{\mu}^2}, \gamma \mathbb{E} \| \overline{\boldsymbol{\theta}}^0 - \boldsymbol{\theta}^{PS} \|_2^2 \right\}$$
.

The proofs of Theorem 3.2 and 3.3 can be found in Appendix C. Besides the technical difficulty as that of Theorem 3.1, we also need to bound the variance of $\overline{\theta}^t$ at the aggregation step. Fortunately, it can be bounded by the consensus error. Then by similar techniques as in the proof of Theorem 3.1, we can establish a standard descent lemma in SGD analysis for the expected distance.

Scheme II requires $p_i = \frac{1}{N}, \forall i$, which violates the unbalanced nature of FL. One solution in Li et al. [2020b] is scaling the local objectives to $g_i(\theta; \tilde{\theta}) = p_i N f_i(\theta; \tilde{\theta})$, and then the global objective is a simple average of the scaled local objectives

$$f(\boldsymbol{\theta}; \tilde{\boldsymbol{\theta}}) := \sum_{i=1}^{N} p_i f_i(\boldsymbol{\theta}; \tilde{\boldsymbol{\theta}}) = \frac{1}{N} \sum_{i=1}^{N} g_i(\boldsymbol{\theta}; \tilde{\boldsymbol{\theta}}).$$

We need to be careful with the Assumptions in Section 2 since scaling the objective will change those properties. The convergence theorems still hold if we replace $L, \mu, \sigma, \varsigma$ with $L' := q_{max}L, \mu' := q_{min}\mu, \sigma' := \sqrt{q_{max}}\sigma, \varsigma' := \sqrt{q_{max}}\varsigma$, where $q_{max} := N \cdot \max_i p_i, q_{min} := N \cdot \min_i p_i$.

3.3 Discussions on the Algorithm and Solution

We will only discuss with respect to the aggregation step in this sub-section for convenience, denoted as $T \in \mathcal{I}_E$, then we can simply use $\frac{T}{E}$ when dividing E. Note for a general step t, we only need to use $\lfloor \frac{t}{E} \rfloor$ to obtain an integer.

Choice of E. We are interested in the total time we need to achieve an ϵ accuracy, and how this total time changes with E. We use our results in Theorem 3.1, 3.2, and 3.3, and denote $T_{\epsilon} := \frac{v}{\epsilon} - \gamma$ as the number of computation steps that is sufficient to guarantee an ϵ -accuracy. To connect T_{ϵ} to the total time needed, suppose the expected time for each communication step is C times the expected time of each computation step, then the total time required for ϵ -accuracy is linear in $T_{\epsilon} + C \cdot \frac{T_{\epsilon}}{E}$. Below we separately analyze the influence of E on $\frac{T_{\epsilon}}{E}$ and T_{ϵ} , and then discuss how to choose the optimal E for different C values.

Let $B_0 := B$ in Theorem 3.1 for full participation and γ_i (i=0,1,2) denotes the γ in Theorem 3.1, 3.2, and 3.3 respectively. Then in Theorem 3.1, 3.2, and 3.3, T_ϵ is dominated by $\mathcal{O}\left(4B_i/\tilde{\mu}^2 + \gamma_i \mathbb{E}[\|\overline{\boldsymbol{\theta}}^0 - \boldsymbol{\theta}^{PS}\|_2^2]\right)$ where i=0,1,2. From the definition, we know that B_i (i=0,1,2) is almost a constant w.r.t. E and γ_i is of $\mathcal{O}(E^2 \log E)$. This means that when E grows, the total update steps to reach ϵ -accuracy, T_ϵ will grow, while the number of aggregation steps needed, $\frac{T_\epsilon}{E}$ will first grow and then decrease.

Now we consider $T_{\epsilon} + C \cdot \frac{T_{\epsilon}}{E}$, the total time needed to reach ϵ -accuracy. From the above analysis, we know it is of order $\mathcal{O}(E^2 \log E) + C \cdot \mathcal{O}(E \log E) + C \cdot \mathcal{O}(\log E/E)$. When communication is fast, i.e., C is small, $\mathcal{O}(E^2 \log E)$ is the dominating term, and we can focus more on the number of computation iterations T_{ϵ} , and smaller E values are preferable. However, when C is large, $C \cdot \mathcal{O}(E \log E) + C \cdot \mathcal{O}(\log E/E)$ becomes the dominating term, and we should focus more on the number of communication rounds $\frac{T_{\epsilon}}{E}$ and some middle E values are preferable.

Choice of K. Again T_{ϵ} is dominated by $\mathcal{O}\left(4B_i/\tilde{\mu}^2 + \gamma_i \mathbb{E}[\|\overline{\theta}^0 - \theta^{PS}\|_2^2]\right)$ where i=1,2. Then by the formulae of B_i (i=1,2), we know T_{ϵ} monotonically decreases with K, but the total communication time increases with K due to more severe stragglers' effect. Generally, as we show in Theorem 3.2 and 3.3, the convergence rate has a weak dependence on K. We have empirically observed this phenomenon in Figure 3(a). Therefore, we can set $\frac{K}{N}$ to an appropriate small value to reduce the straggler's effect while keeping the convergence rate.

Choice of sampling schemes. We formalize the two sampling schemes in Section 2 and show their convergence properties in Theorem 3.2 and 3.3. We note that Scheme I has a desirable property that it naturally supports unbalanced clients, so if the server has control over the sampling, Scheme I should be chosen.

But as discussed in Li et al. [2020b], sometimes the server may have no control over the sampling and simply use the first K received results for update. In this case, if the reception times from each client are IID random variables, we can treat this process as uniformly sampling K out of N at random without replacement. Theorem 3.3 showed the convergence, and the discussion on scaling the objectives provides instructions on how to make the system work with arbitrary initial p_1, \ldots, p_N values. However, it's worth noting that when p_1, \ldots, p_N are highly non-uniform, the corresponding L', σ', ς' values will be much larger than from L, σ, ς , and μ' will be much smaller than μ . Then by the formula of $\hat{\eta}_0$ and $\tilde{\eta}_0$, we have to use much smaller starting learning rates and thus slower convergence. However, such a small learning rate may cause the model to fail to train at all. We also empirically show this in Figure 7.

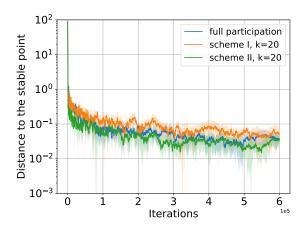


Figure 1: Distance to the performative stable solution vs. the number of iterations for full participation, Scheme I, and Scheme II.

However, an interesting observation is that when $p_i = \frac{1}{N}$, we empirically show in Figure 4 and 5, Scheme II slightly outperforms Scheme I.

Learning rate decay. The learning rate decay is a necessity for stochastic gradient descent (SGD) to converge, even when clients have static, independent and individually distributed (IID) data. The decay is used in Li et al. [2022] in decentralized performative prediction and the necessity for such decay is proved in FedAvg with static, non-IID clients. We also empirically show that constant learning rates fail to converge in Figure 6.

 θ^{PS} and θ^{PO} . Here we discuss the relationship between the θ^{PS} and θ^{PO} solutions more in depth. In the strategic learning setting, Perdomo et al. [2020] showed that θ^{PO} is the Stackelberg equilibrium. It's worth noting that θ^{PS} is not merely an approximation to θ^{PO} , but a natural convergence point of the best response dynamics (BRD). More specifically, when the clients and the decision maker have no information about others' utilities, backward induction is unavailable, and playing the Stackelberg equilibrium is unrealistic. In this case, treating others' strategies in the previous time step as constants, and optimizing one's own strategy accordingly is a rational strategy. Such an optimization step is a best response, and in multi-round sequential strategic learning problems Zrnic et al. [2021], the best responses can form the BRD, and θ^{PS} is the convergence point of the BRD. Although the decision maker's natural best response step is a risk minimization step, the gradient-based P-FedAvg can find the same θ^{PS} . Another interesting observation of θ^{PS} is that if we remove the sequential decision nature, then $\{\theta^{PS}, \mathcal{D}_1(\theta^{PS}), \dots, \mathcal{D}_N(\theta^{PS})\}$ is a Nash equilibrium since no participant has an incentive to unilaterally deviate.

4 Numerical Experiments

4.1 Weighted Gaussian mean performative prediction

As a numerical simulation, we perform P-FedAvg to estimate the mean of heterogeneous Gaussian data under performative effects and examine the impact of the hyperparameters, the sampling schemes, and client heterogeneity. We consider N=25 clients, with the i-th client minimizing the loss function $\ell(\theta;Z_i):=(\theta-Z)^2/2, \, \theta, \, Z\in\mathbb{R}$ on data $Z_i\sim \mathcal{D}_i(\theta):=\mathcal{N}(m_i+\epsilon_i\theta,\sigma^2)$. For this loss function, we have $\mu=1,\, L=1$. For $\bar{\epsilon}\in[0,1)$, the PS solution is $\theta^{PS}=\frac{\sum_{i=1}^N p_i m_i}{1-\bar{\epsilon}}$; while θ^{PS} does not exist when $\bar{\epsilon}\geq 1$. Denote the weighted average of m_i as $\overline{m}=\sum_{i=1}^N p_i m_i$ and the variance as $\mathrm{Var}(m)=\sum_{i=1}^N p_i (m_i-\overline{m})^2$. In experiment, we set $\bar{\epsilon}=0.9,\, \overline{m}=10$.

Figure 1 shows P-FedAvg converges to the performative stable solution in all three communication settings: full participation and the two schemes for the partial participation. Interestingly, partial participation with scheme II converges the fastest in this experiment. Despite the full participation scheme having the lowest upper bound on the number of iterations sufficient to convergence, our experimental results show that the actual convergence behaviors of all three schemes are very similar and weakly depend on K, especially when $p_i = \frac{1}{N}$.

Impact of E. We conduct an experiment to compare the performance of our algorithm with a variety of E values, under a homogeneous system. Figure 2 shows the result on both sampling schemes, with K=20. A slightly larger E leads to faster convergence. However, an extremely large E(E=50) in the experiment can also cause slower convergence. Since at this case, the clients deviate too much at each aggregation, which causes low efficiency issues. In real world scenarios, as the communication cost changes, E should be carefully chosen.

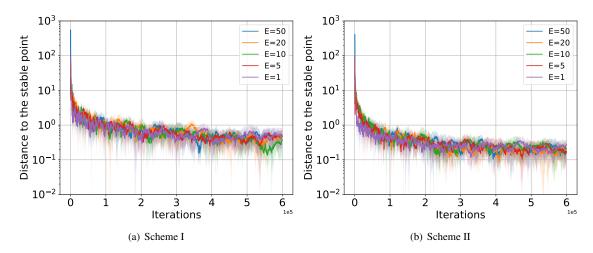


Figure 2: Impact of E on Scheme I and Scheme II. K=20, Var(m)=0.6, $Var(\epsilon)=0.1$ for both (a) and (b). For (b), $p_i=\frac{1}{25}$.

Impact of K. Figure 3 shows the convergence of FedAvg under different k values, For scheme I, larger k leads to faster convergence. While for scheme II, as k increasing, the convergence rate will first increase and then decrease.

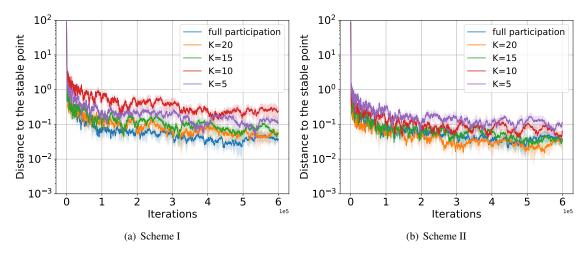


Figure 3: Impact of K on Scheme I and Scheme II . E=5, Var(m)=0.6, $Var(\epsilon)=0.1$ for both (a) and (b). For (b), $p_i=\frac{1}{25}$.

Impact of sampling schemes. Figure 1 also compares different schemes. We can see if the clients' data are uniformly sampled $(p_i = \frac{1}{N})$, then scheme II achieves a better convergence rate, which conforms to our theoretical result because $B_1 > B_2$.

Data heterogeneity and shifting heterogeneity. In Figure 4 we test our algorithm under data heterogeneity. Specifically, we set m and ϵ to have large variances, respectively. In this example, m mainly captures the data heterogeneity and ϵ capture the shifting heterogeneity. This experiment shows our algorithm still converges under a certain amount of

heterogeneity. Comparing the performance of our algorithm on both figures, we can see shifting heterogeneity is the main factor in performative federated learning.

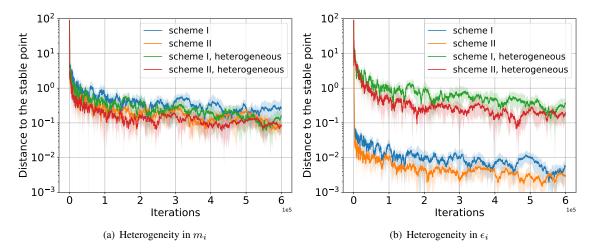


Figure 4: Impact of heterogeneity on the two schemes of partial participation. The impact of m_i is shown in (a) and the impact of ϵ_i is shown in (b). $K=20, p_i=\frac{1}{25}$ and. In (a), Var(m)=6 for hetergeneous case and 0 for homogeneous case, $Var(\epsilon)=0.1$. In (b), $Var(\epsilon)=0.6$ for hetergeneous case and 0.1 for homogeneous case, Var(m)=0.6.

4.2 Credit Score Strategic Classification

To show the performance of P-FedAvg on a real world dataset, we follow Perdomo et al. [2020] and use the same Kaggle dataset *, where a bank predicts whether loan applicants are creditworthy. The features consist of the information about an individual, and the target is 1 if the individual defaulted on a loan, and 0 otherwise. We use the same strategic setting as in Perdomo et al. [2020] where the applicants can manipulate their features in (1) revolving utilization of unsecured lines, (2) number of open credit lines and loans, and (3) number real estate loans or lines. The strength of manipulation for the *i*-th population is controlled by ϵ_i . We equally partition the training set into 10 subsets and distributed it to 10 clients, and thus $p_i = 0.1, \forall i$. The sensitivities ϵ_i for the 10 clients are independently and uniformly sampled from [0.9, 1.1]. We set K = 5 in partial participation. We train a logistic regression binary classifier. In each round of P-FedAvg, we perform E = 5 gradient descent steps on a random minibatch of size 4. A discussion on the effect of the batch size can be found in Appendix E.2.

Figure 5 shows the loss function and the distance to the PS solution as the number of deployment rounds increases. The mean and 1 standard deviation error bar are generated from 5 experiments with different random seeds. Similar to the numerical simulation, the actual convergence behaviors of all three schemes are very similar.

5 Conclusion

In this work, we leveraged the idea of distribution shift mappings in performative predictions to study federated learning problems where data shifts exists and such shifts are induced by the decision parameters. We formulated the performative federated learning problem and showed that a unique performative stable solution exists, which is a natural equilibrium of the iterative updating process between the server and the clients. Then we formalized the P-FedAvg algorithm and proved that both the full device participation and the partial device participation schemes have $\mathcal{O}(1/T)$ convergence rate to the performative stable solution. We also thoroughly discussed the impact of how some of the key system parameters influence the convergence, including the aggregation interval size, the number of sampled devices in partial participation, the sampling schemes, and the heterogeneity among clients. Our numerical results validate our theory and discussion, and provide valuable insights into the real-world applications of performative federated learning.

^{*}www.kaggle.com/competitions/GiveMeSomeCredit/data

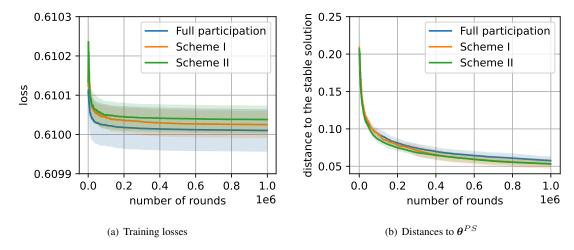


Figure 5: The losses (a) and the distances to the PS solution (b) for the full participation, Scheme I and Scheme II.

References

Mark Braverman and Sumegha Garg. The role of randomness and noise in strategic classification. In 1st Symposium on Foundations of Responsible Computing, 2020.

Gavin Brown, Shlomi Hod, and Iden Kalemaj. Performative prediction in a stateful world. *CoRR*, abs/2011.03885, 2020. URL https://arxiv.org/abs/2011.03885.

Fernando E Casado, Dylan Lema, Marcos F Criado, Roberto Iglesias, Carlos V Regueiro, and Senén Barro. Concept drift detection and adaptation for federated and continual learning. *Multimedia Tools and Applications*, 81(3): 3397–3419, 2022.

Yatong Chen, Jialu Wang, and Yang Liu. Strategic recourse in linear classification. *arXiv preprint arXiv:2011.00355*, 2020.

Yucheng Ding, Chaoyue Niu, Yikai Yan, Zhenzhe Zheng, Fan Wu, Guihai Chen, Shaojie Tang, and Rongfei Jia. Distributed optimization over block-cyclic data. *arXiv preprint arXiv:2002.07454*, 2020.

Jinshuo Dong, Aaron Roth, Zachary Schutzman, Bo Waggoner, and Zhiwei Steven Wu. Strategic classification from revealed preferences. In *Proceedings of the 2018 ACM Conference on Economics and Computation*, pages 55–70, 2018.

Dmitriy Drusvyatskiy and Lin Xiao. Stochastic optimization with decision-dependent distributions, 2020. URL https://arxiv.org/abs/2011.11173.

Dmitriy Drusvyatskiy and Lin Xiao. Stochastic optimization with decision-dependent distributions. *Mathematics of Operations Research*, 2022.

Hubert Eichner, Tomer Koren, Brendan McMahan, Nathan Srebro, and Kunal Talwar. Semi-cyclic stochastic gradient descent. In *International Conference on Machine Learning*, pages 1764–1773. PMLR, 2019.

Yongxin Guo, Tao Lin, and Xiaoying Tang. Towards federated learning on time-evolving heterogeneous data. *arXiv* preprint arXiv:2112.13246, 2021.

Farzin Haddadpour, Mohammad Mahdi Kamani, Aryan Mokhtari, and Mehrdad Mahdavi. Federated learning with compression: Unified analysis and sharp guarantees. In *International Conference on Artificial Intelligence and Statistics*, pages 2350–2358. PMLR, 2021.

Nika Haghtalab, Nicole Immorlica, Brendan Lucier, and Jack Wang. Maximizing welfare with incentive-aware evaluation mechanisms. pages 160–166, 07 2020. doi: 10.24963/ijcai.2020/23.

Moritz Hardt, Nimrod Megiddo, Christos Papadimitriou, and Mary Wootters. Strategic classification. pages 111–122, 01 2016. doi: 10.1145/2840728.2840730.

- Seyyedali Hosseinalipour, Su Wang, Nicolo Michelusi, Vaneet Aggarwal, Christopher G Brinton, David J Love, and Mung Chiang. Parallel successive learning for dynamic distributed model training over heterogeneous wireless networks. *arXiv* preprint arXiv:2202.02947, 2022.
- Lily Hu, Nicole Immorlica, and Jennifer Vaughan. The disparate effects of strategic manipulation. pages 259–268, 01 2019. doi: 10.1145/3287560.3287597.
- Zachary Izzo, Lexing Ying, and James Zou. How to learn when data reacts to your model: Performative gradient descent. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 4641–4650. PMLR, 18–24 Jul 2021. URL https://proceedings.mlr.press/v139/izzo21a.html.
- Zachary Izzo, James Zou, and Lexing Ying. How to learn when data gradually reacts to your model. In Gustau Camps-Valls, Francisco J. R. Ruiz, and Isabel Valera, editors, *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151 of *Proceedings of Machine Learning Research*, pages 3998–4035. PMLR, 28–30 Mar 2022. URL https://proceedings.mlr.press/v151/izzo22a.html.
- Kun Jin, Xueru Zhang, Mohammad Mahdi Khalili, Parinaz Naghizadeh, and Mingyan Liu. Incentive mechanisms for strategic classification and regression problems. In David M. Pennock, Ilya Segal, and Sven Seuken, editors, *EC '22: The 23rd ACM Conference on Economics and Computation, Boulder, CO, USA, July 11 15, 2022*, pages 760–790. ACM, 2022. doi: 10.1145/3490486.3538300. URL https://doi.org/10.1145/3490486.3538300.
- Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In *International Conference on Machine Learning*, pages 5132–5143. PMLR, 2020.
- Jon Kleinberg and Manish Raghavan. How do classifiers induce agents to invest effort strategically? *ACM Transactions on Economics and Computation*, 8:1–23, 11 2020. doi: 10.1145/3417742.
- Daliang Li and Junpu Wang. Fedmd: Heterogenous federated learning via model distillation. *arXiv preprint* arXiv:1910.03581, 2019.
- Qiang Li and Hoi-To Wai. State dependent performative prediction with stochastic approximation. In Gustau Camps-Valls, Francisco J. R. Ruiz, and Isabel Valera, editors, *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151 of *Proceedings of Machine Learning Research*, pages 3164–3186. PMLR, 28–30 Mar 2022. URL https://proceedings.mlr.press/v151/li22c.html.
- Qiang Li, Chung-Yiu Yau, and Hoi-To Wai. Multi-agent performative prediction with greedy deployment and consensus seeking agents, 2022. URL https://arxiv.org/abs/2209.03811.
- Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. *Proceedings of Machine Learning and Systems*, 2:429–450, 2020a.
- Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence of fedavg on non-iid data. In *International Conference on Learning Representations*, 2020b. URL https://openreview.net/forum?id=HJxNAnVtDS.
- Tao Lin, Lingjing Kong, Sebastian U Stich, and Martin Jaggi. Ensemble distillation for robust model fusion in federated learning. *Advances in Neural Information Processing Systems*, 33:2351–2363, 2020.
- Celestine Mendler-Dünner, Juan Perdomo, Tijana Zrnic, and Moritz Hardt. Stochastic optimization for performative prediction. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 4929–4939. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper/2020/file/33e75ff09dd601bbe69f351039152189-Paper.pdf.
- John Miller, Smitha Milli, and Moritz Hardt. Strategic classification is causal modeling in disguise. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 6917–6926. PMLR, 13–18 Jul 2020.
- John P Miller, Juan C Perdomo, and Tijana Zrnic. Outside the echo chamber: Optimizing the performative risk. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 7710–7720. PMLR, 18–24 Jul 2021. URL https://proceedings.mlr.press/v139/miller21a.html.

- Smitha Milli, John Miller, Anca Dragan, and Moritz Hardt. The social cost of strategic classification. pages 230–239, 01 2019. doi: 10.1145/3287560.3287576.
- Adhyyan Narang, Evan Faulkner, Dmitriy Drusvyatskiy, Maryam Fazel, and Lillian J. Ratliff. Multiplayer performative prediction: Learning in decision-dependent games. *CoRR*, abs/2201.03398, 2022. URL https://arxiv.org/abs/2201.03398.
- Juan Perdomo, Tijana Zrnic, Celestine Mendler-Dünner, and Moritz Hardt. Performative prediction. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 7599–7609. PMLR, 13–18 Jul 2020. URL https://proceedings.mlr.press/v119/perdomo20a.html.
- Reilly Raab and Yang Liu. Unintended selection: Persistent qualification rate disparities and interventions. *Advances in Neural Information Processing Systems*, 34:26053–26065, 2021.
- Mitas Ray, Lillian J. Ratliff, Dmitriy Drusvyatskiy, and Maryam Fazel. Decision-dependent risk minimization in geometrically decaying dynamic environments. In *Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022, Thirty-Fourth Conference on Innovative Applications of Artificial Intelligence, IAAI 2022, The Twelveth Symposium on Educational Advances in Artificial Intelligence, EAAI 2022 Virtual Event, February 22 March 1, 2022*, pages 8081–8088. AAAI Press, 2022. URL https://ojs.aaai.org/index.php/AAAI/article/view/20780.
- Elsa Rizk, Stefan Vlaski, and Ali H. Sayed. Dynamic federated learning. *CoRR*, abs/2002.08782, 2020. URL https://arxiv.org/abs/2002.08782.
- Anit Kumar Sahu, Tian Li, Maziar Sanjabi, Manzil Zaheer, Ameet Talwalkar, and Virginia Smith. On the convergence of federated optimization in heterogeneous networks. *CoRR*, abs/1812.06127, 2018. URL http://arxiv.org/abs/1812.06127.
- Yonadav Shavit, Benjamin Edelman, and Brian Axelrod. Causal strategic linear regression, 2020.
- Jianyu Wang, Qinghua Liu, Hao Liang, Gauri Joshi, and H Vincent Poor. Tackling the objective inconsistency problem in heterogeneous federated optimization. *Advances in neural information processing systems*, 33:7611–7623, 2020.
- Killian Wood, Gianluca Bianchin, and Emiliano Dall'Anese. Online projected gradient descent for stochastic optimization with decision-dependent distributions. *IEEE Control Systems Letters*, 6:1646–1651, 2022. doi: 10.1109/LCSYS.2021.3124187.
- Xueru Zhang, Mohammad Mahdi Khalili, Kun Jin, Parinaz Naghizadeh, and Mingyan Liu. Fairness interventions as (Dis)Incentives for strategic manipulation. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 26239–26264. PMLR, 17–23 Jul 2022. URL https://proceedings.mlr.press/v162/zhang221.html.
- Chen Zhu, Zheng Xu, Mingqing Chen, Jakub Konečný, Andrew Hard, and Tom Goldstein. Diurnal or nocturnal? federated learning of multi-branch networks from periodically shifting distributions. In *International Conference on Learning Representations*, 2021a.
- Zhuangdi Zhu, Junyuan Hong, and Jiayu Zhou. Data-free knowledge distillation for heterogeneous federated learning. In *International Conference on Machine Learning*, pages 12878–12889. PMLR, 2021b.
- Tijana Zrnic, Eric Mazumdar, Shankar Sastry, and Michael Jordan. Who leads and who follows in strategic classification? In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 15257–15269. Curran Associates, Inc., 2021. URL https://proceedings.neurips.cc/paper/2021/file/812214fb8e7066bfa6e32c626c62c688b-Paper.pdf.

A Proof of Proposition 2.8 and 2.9

Proposition 2.5. (Uniqueness of θ^{PS}) Under Assumptions 2.1, 2.2 and 2.3, define the map $\Phi: \mathbb{R}^m \to \mathbb{R}^m$

$$\Phi(\boldsymbol{\theta}) := arg \min_{\boldsymbol{\theta}' \in \mathbb{R}^M} f(\boldsymbol{\theta}', \boldsymbol{\theta})$$

If $\overline{\epsilon} := \sum_{i=1}^N p_i \epsilon_i < \mu/L$, then $\Phi(\cdot)$ is a contraction mapping with the unique fixed point $\boldsymbol{\theta}^{PS} = \Phi(\boldsymbol{\theta}^{PS})$. On the contrary, if $\overline{\epsilon} \ge \mu/L$, then there is an instance where any sequence generated by $\Phi(\cdot)$ will diverge.

Proof. This proof simulates the proof of Proposition 1 in Li et al. [2022].

Fix θ' , $\theta \in \mathbb{R}^m$, the optimality condition implies that

$$\sum_{i=1}^{N} p_i \nabla f_i(\Phi(\boldsymbol{\theta}); \boldsymbol{\theta}) = \mathbf{0}, \quad \sum_{i=1}^{N} p_i \nabla f_i(\Phi(\boldsymbol{\theta}'); \boldsymbol{\theta}') = \mathbf{0}$$

where the gradients are taken w.r.t the first argument in f_i . Then we have

$$0 = \langle \mathbf{0}, \Phi(\boldsymbol{\theta}) - \Phi(\boldsymbol{\theta}') \rangle$$
$$= \langle \sum_{i=1}^{N} p_i (\nabla f_i(\Phi(\boldsymbol{\theta}); \boldsymbol{\theta}) - \nabla f_i(\Phi(\boldsymbol{\theta}'); \boldsymbol{\theta}')), \Phi(\boldsymbol{\theta}) - \Phi(\boldsymbol{\theta}') \rangle.$$

Rearranging the above equation and adding $\sum_{i=1}^{N} p_i f_i(\Phi(\theta), \theta')$ to both hand sides leads to

$$\langle \sum_{i=1}^{N} p_i (\nabla f_i(\Phi(\boldsymbol{\theta}); \boldsymbol{\theta}') - \nabla f_i(\Phi(\boldsymbol{\theta}); \boldsymbol{\theta})), \Phi(\boldsymbol{\theta}) - \Phi(\boldsymbol{\theta}') \rangle$$

$$= \langle \sum_{i=1}^{N} p_i (\nabla f_i(\Phi(\boldsymbol{\theta}); \boldsymbol{\theta}') - \nabla f_i(\Phi(\boldsymbol{\theta}'); \boldsymbol{\theta}')), \Phi(\boldsymbol{\theta}) - \Phi(\boldsymbol{\theta}') \rangle.$$

By strong convexity in assumption 2.1, we have

$$f(\Phi(\boldsymbol{\theta}); \boldsymbol{\theta}') \geq f(\Phi(\boldsymbol{\theta}'); \boldsymbol{\theta}') + \langle \nabla f(\Phi(\boldsymbol{\theta}'); \boldsymbol{\theta}'), \Phi(\boldsymbol{\theta}) - \Phi(\boldsymbol{\theta}') \rangle + \frac{\mu}{2} \|\Phi(\boldsymbol{\theta}) - \Phi(\boldsymbol{\theta}')\|_{2}^{2},$$

$$f(\Phi(\boldsymbol{\theta})';\boldsymbol{\theta}') \geq f(\Phi(\boldsymbol{\theta});\boldsymbol{\theta}') + \langle \nabla f(\Phi(\boldsymbol{\theta});\boldsymbol{\theta}'), \Phi(\boldsymbol{\theta}) - \Phi(\boldsymbol{\theta}') \rangle + \frac{\mu}{2} \|\Phi(\boldsymbol{\theta}) - \Phi(\boldsymbol{\theta}')\|_{2}^{2},$$

and thus

$$\langle \nabla f(\Phi(\boldsymbol{\theta}); \boldsymbol{\theta}') - \nabla f(\Phi(\boldsymbol{\theta}'); \boldsymbol{\theta}'), \Phi(\boldsymbol{\theta}) - \Phi(\boldsymbol{\theta}') \rangle \ge \mu \|\Phi(\boldsymbol{\theta}) - \Phi(\boldsymbol{\theta}')\|_2^2. \tag{4}$$

Applying Lemma 2.4, we have

$$\sum_{i=1}^{N} p_i \langle \nabla f_i(\Phi(\boldsymbol{\theta}); \boldsymbol{\theta}') - \nabla f_i(\Phi(\boldsymbol{\theta}'); \boldsymbol{\theta}'), \Phi(\boldsymbol{\theta}) - \Phi(\boldsymbol{\theta}') \rangle \leq \sum_{i=1}^{N} p_i L \epsilon_i \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_2 \cdot \|\Phi(\boldsymbol{\theta}) - \Phi(\boldsymbol{\theta})\|_2.$$
 (5)

Combine (4) and (5), we have

$$\|\Phi(\boldsymbol{\theta}) - \Phi(\boldsymbol{\theta}')\|_{2} \leq \frac{\sum_{i=1}^{N} p_{i} \epsilon_{i} L}{\mu} \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_{2} = \frac{\overline{\epsilon} L}{\mu} \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_{2}.$$
 (6)

Therefore, if $\bar{\epsilon} < \frac{L}{\mu}$, $\Phi(\cdot)$ is a contraction mapping by Banach fixed point theorem and admits a unique fixed point θ^{PS} .

To show the divergence when $\bar{\epsilon} \geq \frac{L}{\mu}$, we consider the following example where $\theta \in \mathbb{R}$, $\frac{L}{\mu} = 1$, $\bar{\gamma} := \sum_{i=1}^{N} p_i \gamma_i \neq 0$, and

$$\ell(\theta; Z) = \frac{1}{2}(\theta - Z)^2, Z \sim \mathcal{D}_i(\theta) = \mathcal{N}(\gamma_i + \epsilon_i \theta, 1)$$

we observe

$$f_{i}(\theta';\theta) = \mathbb{E}_{Z \sim \mathcal{D}_{i}(\theta)} \left[\frac{1}{2} (\theta - Z)^{2} \right]$$

$$= \mathbb{E}_{\tilde{Z} \sim \mathcal{N}(0,1)} \left[\frac{1}{2} (\theta' - \gamma_{i} - \epsilon_{i}\theta - \tilde{Z})^{2} \right]$$

$$= \frac{1}{2} (\theta' - \gamma_{i} - \epsilon_{i}\theta)^{2} + \frac{1}{2},$$

$$\Phi(\theta) = \operatorname{argmin}_{\theta \in \mathbb{R}} \sum_{i=1}^{N} p_{i} (\theta' - \gamma_{i} - \epsilon_{i}\theta)^{2} = \overline{\epsilon}\theta + \overline{\gamma},$$

so by applying $\Phi(\cdot)$ t times, we obtain

$$\Phi^t(\theta) = \overline{\epsilon}^t \theta + (1 + \overline{\epsilon} + \dots + \overline{\epsilon}^{(t-1)}) \overline{\gamma},$$

and since $\bar{\epsilon} \geq \frac{L}{\mu} = 1, \bar{\gamma} \neq 0$, we have $\lim_{t \to \infty} \|\Phi^t(\theta)\|_2 = \infty$.

Proposition 2.6. Under Assumption 2.1 and 2.3, suppose that the loss $\ell(\theta; Z)$ is L_z -Lipschitz in Z, let $\overline{\epsilon} := \sum_{i=1}^N p_i \epsilon_i$, we have for every performative stable solution and every performative optimal solution θ^{PO} that

$$\|\boldsymbol{\theta}^{PS} - \boldsymbol{\theta}^{PO}\|_2 \le \frac{2L_z\overline{\epsilon}}{\mu}.$$

Proof. This proof simulates the proof of Theorem 4.3 in Perdomo et al. [2020].

First by the optimality of θ^{PO} , we have $f(\theta^{PO}; \theta^{PO}) \leq f(\theta^{PS}; \theta^{PS})$. By strong convexity in Assumption 2.1, we have

$$f(\boldsymbol{\theta}^{PO};\boldsymbol{\theta}^{PS}) \geq f(\boldsymbol{\theta}^{PS};\boldsymbol{\theta}^{PS}) + \langle \nabla f(\boldsymbol{\theta}^{PS};\boldsymbol{\theta}^{PS}), \; \boldsymbol{\theta}^{PO} - \boldsymbol{\theta}^{PS} \rangle + \frac{\mu}{2} \|\boldsymbol{\theta}^{PO} - \boldsymbol{\theta}^{PS}\|_2^2 \geq \frac{\mu}{2} \|\boldsymbol{\theta}^{PO} - \boldsymbol{\theta}^{PS}\|_2^2$$

Further by Assmption 2.3, the the loss $\ell(\theta; Z)$ is L_z -Lipschitz in Z, and Kantorovich-Rubinstein duality, we have

$$f(\boldsymbol{\theta}^{PO}; \boldsymbol{\theta}^{PS}) - f(\boldsymbol{\theta}^{PO}; \boldsymbol{\theta}^{PO})$$

$$= \sum_{i=1}^{N} p_{i} \Big(\mathbb{E}_{Z_{i} \sim \mathcal{D}_{i}(\boldsymbol{\theta}^{PS})} [\ell(\boldsymbol{\theta}^{PO}; Z_{i})] - \mathbb{E}_{Z_{i} \sim \mathcal{D}_{i}(\boldsymbol{\theta}^{PO})} [\ell(\boldsymbol{\theta}^{PO}; Z_{i})] \Big)$$

$$\leq \sum_{i=1}^{N} p_{i} L_{z} \mathcal{W}_{1}(\mathcal{D}_{i}(\boldsymbol{\theta}^{PS}), \mathcal{D}_{i}(\boldsymbol{\theta}^{PO}))$$

$$= \sum_{i=1}^{N} p_{i} L_{z} \epsilon_{i} \|\boldsymbol{\theta}^{PO} - \boldsymbol{\theta}^{PS}\|_{2} = L_{z} \overline{\epsilon} \|\boldsymbol{\theta}^{PO} - \boldsymbol{\theta}^{PS}\|_{2}. \tag{7}$$

where the inequality is a well-know conclusion in optimal transport theory. Equation 7, we have $L_z \bar{\epsilon} \epsilon_i \|\boldsymbol{\theta}^{PO} - \boldsymbol{\theta}^{PS}\|_2 \ge f(\boldsymbol{\theta}^{PO}, \boldsymbol{\theta}^{PS}) - f(\boldsymbol{\theta}^{PO}; \boldsymbol{\theta}^{PO}) \ge f(\boldsymbol{\theta}^{PO}; \boldsymbol{\theta}^{PS}) - f(\boldsymbol{\theta}^{PS}; \boldsymbol{\theta}^{PS}) \ge \frac{\mu}{2} \|\boldsymbol{\theta}^{PO} - \boldsymbol{\theta}^{PS}\|_2^2$, implying that $\|\boldsymbol{\theta}^{PO} - \boldsymbol{\theta}^{PS}\|_2 \le \frac{2L_z\bar{\epsilon}}{\mu}$.

B Proof of Theorem 3.1

B.1 Additional Notation

In our analysis, for the sake of convenience, we will define two additional sequences as $\overline{\boldsymbol{w}}^t := \sum_{i=1}^N p_i \boldsymbol{w}_i^t$ and $\overline{\boldsymbol{\theta}}^t := \sum_{i=1}^N p_i \boldsymbol{\theta}_i^t$, following that of Li et al. [2020b]. We note that $\overline{\boldsymbol{w}}^t$ results from a single step of SGD from $\overline{\boldsymbol{\theta}}^t$. When $t+1 \notin \mathcal{I}_E$, both $\overline{\boldsymbol{w}}^t$ and $\overline{\boldsymbol{\theta}}^t$ are unaccessible. When $t+1 \in \mathcal{I}_E$, we can obtain $\overline{\boldsymbol{\theta}}^t$. In addition, we also define $\overline{\boldsymbol{g}}_t := \sum_{i=1}^N p_i \nabla f_i(\boldsymbol{\theta}_i^t; \boldsymbol{\theta}_i^t), \quad \boldsymbol{g}_t := \sum_{i=1}^N p_i \nabla \ell(\boldsymbol{\theta}_i^t; \mathcal{Z}_i^{t+1})$ where $Z_i^{t+1} \sim \mathcal{D}_i(\boldsymbol{\theta}_i^t)$. It is clear that in full participation, $\overline{\boldsymbol{w}}^{t+1} = \overline{\boldsymbol{\theta}}^t - \eta_t \boldsymbol{g}_t$ and $\mathbb{E} \boldsymbol{g}_t = \overline{\boldsymbol{g}}_t$. Clearly we have $\overline{\boldsymbol{\theta}}^t = \overline{\boldsymbol{w}}^t$ for any t.

B.2 Key Lemmas

For clarity, we will present several lemmas for establishing our main theorem. In particular, we will present a descent lemma for $\mathbb{E}[\|\overline{\boldsymbol{w}}^t - \boldsymbol{\theta}^{PS}\|_2^2]$ and an upper bound for $\sum_{i=1}^N p_i \mathbb{E}\|\boldsymbol{\theta}_i^t - \overline{\boldsymbol{\theta}}^t\|_2^2$, which together gives a standard descent lemma for $\mathbb{E}[\|\overline{\boldsymbol{\theta}}^{t+1} - \boldsymbol{\theta}^{PS}\|_2^2]$ in SGD analysis and leads to $\mathcal{O}(\frac{1}{t})$ convergence. The proof of the Lemmas will be deferred.

In the following lemma, we aim to establish an upper bound for $\mathbb{E}[\|\overline{\boldsymbol{\theta}}^{t+1} - \boldsymbol{\theta}^{PS}\|_2^2]$. Because $\overline{\boldsymbol{\theta}}^{t+1} = \overline{\boldsymbol{w}}^{t+1}$ in full participation, this is equivalent to establishing an upper bound for $\mathbb{E}[\|\overline{\boldsymbol{w}}^{t+1} - \boldsymbol{\theta}^{PS}\|_2^2]$.

Lemma B.1. (Descent Lemma) Under Assumptions 2.1, 2.2, 2.3, 2.5, in full participation

$$\mathbb{E}[\|\overline{\boldsymbol{\theta}}^{t+1} - \boldsymbol{\theta}^{PS}\|_2^2] = \mathbb{E}[\|\overline{\boldsymbol{w}}^{t+1} - \boldsymbol{\theta}^{PS}\|_2^2] \leq (1 - \tilde{\mu}\eta_t)\mathbb{E}\|\overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}\|_2^2 + 2\sigma^2\eta_t^2 + (c_1\eta_t + c_2\eta_t^2)\sum_{i=1}^N p_i\mathbb{E}\|\boldsymbol{\theta}_i^t - \overline{\boldsymbol{\theta}}^t\|_2^2$$

for any t, where
$$\epsilon_{max} := \max_i \epsilon_i, \overline{\epsilon} := \sum_{i=1}^N p_i \epsilon_i, c_1 := \frac{L(1+\epsilon_{max})^2}{2\delta \overline{\epsilon}}, c_2 := 4[\sigma^2 + L^2(1+\epsilon_{max})^2], \widetilde{\mu} := \mu - (1+\delta)\overline{\epsilon}L$$

Now we are going to establish an upper bound for $\sum_{i=1}^N p_i \mathbb{E} \| \boldsymbol{\theta}_i^t - \overline{\boldsymbol{\theta}}^t \|_2^2$. Note that if $t \in \mathcal{I}_E$, the synchronization step, we have $\boldsymbol{\theta}_i^t = \overline{\boldsymbol{\theta}}^t$ for any $i \in [N]$, which implies that $\sum_{i=1}^N p_i \mathbb{E} \| \boldsymbol{\theta}_i^t - \overline{\boldsymbol{\theta}}^t \|_2^2 = 0$. If $t \notin \mathcal{I}_E$, the following lemma gives an upper bound for $\sum_{i=1}^N p_i \mathbb{E} \| \boldsymbol{\theta}_i^t - \overline{\boldsymbol{\theta}}^t \|_2^2$.

Lemma B.2. (Consensus Error) Under Assumption 2.1, 2.2, 2.3, 2.5, 2.6, if $\{\eta_t\}$ is non-increasing, $\eta_t \leq 2\eta_{t+E}$, $t \notin \mathcal{I}_E$, $\eta_t^2 \leq 1/(2c_3(t+1-t_0)(1+2(t+1-t_0)))$, and

$$\eta_0 \leq \hat{\eta}_0 := \frac{\tilde{\mu} \mathbb{E} \|\overline{\boldsymbol{\theta}}^0 - {\boldsymbol{\theta}}^{PS}\|_2^2}{2\sigma^2 + (c_1c_3 + c_2/6)(2E^2 - E)\log E((16\sigma^2 + 12\varsigma^2)\mathbb{E} \|\overline{\boldsymbol{\theta}}^0 - {\boldsymbol{\theta}}^{PS}\|_2^2 + (8\sigma^2 + 12\varsigma^2))},$$

then in full participation, we have

$$\sum_{i=1}^{N} p_{i} \mathbb{E} \|\boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t}\|_{2}^{2} \leq 4\eta_{t}^{2} (2E^{2} - E) \log E (48\sigma^{2} + 36\varsigma^{2}) \mathbb{E} \|\overline{\boldsymbol{\theta}}^{0} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + 4\eta_{t}^{2} (2E^{2} - E) \log E (24\sigma^{2} + 36\varsigma^{2}).$$

where for any t, where
$$\epsilon_{max} := \max_i \epsilon_i, \bar{\epsilon} := \sum_{i=1}^N p_i \epsilon_i, c_1 := \frac{L(1+\epsilon_{max})^2}{2\delta \bar{\epsilon}}, c_2 := 4[\sigma^2 + L^2(1+\epsilon_{max})^2], \tilde{\mu} := \mu - (1+\delta)\bar{\epsilon}L, c_3 := 12\sigma^2 + 18L^2(1+\epsilon_{max})^2.$$

(One should note that $4\eta_t^2$, $(48\sigma^2 + 36\varsigma^2)$, and $(24\sigma^2 + 36\varsigma^2)$ comes from several times of applying $\eta_{t-1} \le 2\eta_t$ and the real constants could be much smaller by choosing stepsizes carefully.)

The following lemma gives us a standard descent lemma in SGD analysis under technical conditions for establishing the $\mathcal{O}(\frac{1}{T})$ convergence in Theorem 3.1.

Lemma B.3. Under Assumption 2.1, 2.2, 2.3, 2.5, 2.6, if $\{\eta_t\}$ is non-increasing, $\eta_t \leq 2\eta_{t+E}$, $\eta_t^2 \leq 1/(2c_3(t+1-t_0)(1+2(t+1-t_0)))$, and

$$\eta_0 \leq \hat{\eta}_0 := \frac{\tilde{\mu} \mathbb{E} \|\overline{\boldsymbol{\theta}}^0 - \boldsymbol{\theta}^{PS}\|_2^2}{2\sigma^2 + (c_1c_3 + c_2/6)(2E^2 - E)\log E((16\sigma^2 + 12\varsigma^2)\mathbb{E} \|\overline{\boldsymbol{\theta}}^0 - \boldsymbol{\theta}^{PS}\|_2^2 + (8\sigma^2 + 12\varsigma^2))},$$

then in full participation, we have

$$\mathbb{E}[\|\overline{\boldsymbol{\theta}}^{t+1} - \boldsymbol{\theta}^{PS}\|_2^2] \le (1 - \tilde{\mu}\eta_t)\mathbb{E}\|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_2^2 + B\eta_t^2$$

where for any t, where $\epsilon_{max} := \max_i \epsilon_i, \overline{\epsilon} := \sum_{i=1}^N p_i \epsilon_i, c_1 := \frac{L(1+\epsilon_{max})^2}{2\delta\overline{\epsilon}}, c_2 := 4[\sigma^2 + L^2(1+\epsilon_{max})^2], \widetilde{\mu} := \mu - (1+\delta)\overline{\epsilon}L, c_3 := 12\sigma^2 + 18L^2(1+\epsilon_{max})^2, \text{ and } B := 2\sigma^2 + (4c_1\hat{\eta}_0 + 4c_2\hat{\eta}_0^2)(2E^2 - E)\log E\big((48\sigma^2 + 36\varsigma^2)\mathbb{E}\|\overline{\boldsymbol{\theta}}^0 - \boldsymbol{\theta}^{PS}\|_2^2 + (24\sigma^2 + 36\varsigma^2)\big).$

B.3 Completing the Proof of Theorem 3.1

We restate the definitions of all the constants here:

Constants independent of system design.

$$\epsilon_{max} := \max_{i} \epsilon_{i},$$

$$\bar{\epsilon} := \sum_{i=1}^{N} p_{i} \epsilon_{i},$$

$$\tilde{\mu} := \mu - (1 + \delta) \bar{\epsilon} L.$$

$$\mu := \mu \quad (1 + 0)$$
cL;

$$c_1 := (L(1 + \epsilon_{max})^2)/(2\delta\overline{\epsilon}),$$

$$c_2 := 4[\sigma^2 + L^2(1 + \epsilon_{max})^2],$$

$$c_3 := 6[2\sigma^2 + 3L^2(1 + \epsilon_{\max})^2],$$

$$c_4 := 16\sigma^2 + 12\varsigma^2 + (8\sigma^2 + 12\varsigma^2) / \mathbb{E} \|\overline{\boldsymbol{\theta}}^0 - {\boldsymbol{\theta}}^{PS}\|_2^2,$$

$$c_5 := (48\sigma^2 + 36\varsigma^2) \mathbb{E} \|\overline{\boldsymbol{\theta}}^0 - {\boldsymbol{\theta}}^{PS}\|_2^2 + (24\sigma^2 + 36\varsigma^2).$$

Constants related to system design (e.g., E, K).

$$\hat{\eta}_0 := \tilde{\mu} / (2\sigma^2 + (c_1c_3 + c_2/6)c_4(2E^2 - E)\log E),$$

$$B := 2\sigma^2 + (4c_1\hat{\eta}_0 + 4c_2\hat{\eta}_0^2)c_5(2E^2 - E)\log E,$$

$$c_6 := (2E^2 + 3E + 1)\log(E + 1),$$

$$\tilde{\eta}_0 := \tilde{\mu}/(2\sigma^2 + (c_1c_3 + c_2/6)c_4c_6),$$

$$B_1 := 2\sigma^2 + (4c_1\tilde{\eta}_0 + 4c_2\tilde{\eta}_0^2 + 1/K)c_5c_6,$$

$$B_2 := 2\sigma^2 + \left(4c_1\tilde{\eta}_0 + 4c_2\tilde{\eta}_0^2 + \frac{N-K}{KN(N-1)}\right)c_5c_6.$$

Instead of proving Theorem 3.1 directly, we prove a more general version of convergence results suppose that some conditions about the stepsize are satisfied. Then we will show that the stepsizes given in Theorem 3.1 satisfy the conditions.

Theorem B.4. Under Assumption 2.1, 2.2, 2.3, 2.5, 2.6, for a diminishing stepsize $\eta_t = \frac{\beta}{t+\gamma}$ where $\beta > \frac{1}{\tilde{\mu}}$, $\gamma > 0$ such that $\eta_0 \leq \hat{\eta}_0$, $\eta_t \leq 2\eta_{t+E}$, and $\eta_t^2 \leq 1/(2c_3(t+1-t_0)(1+2(t+1-t_0)))$, then in full participation, we have for any t

$$\mathbb{E}[\|\overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}\|_2^2] \le \frac{\upsilon}{\gamma + t}$$

where
$$v = \max \left\{ \frac{4B}{\tilde{\mu}^2}, \gamma \mathbb{E}[\|\overline{\boldsymbol{\theta}}^0 - \boldsymbol{\theta}^{PS}\|_2^2] \right\}$$
.

Proof. Let $\Delta_t := \mathbb{E}[\|\overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}\|_2^2]$, then from Lemma B.3, we have

$$\Delta_{t+1} \le (1 - \tilde{\mu}\eta_t) \mathbb{E} \|\overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}\|_2^2 + B\eta_t^2.$$

For a diminishing stepsize $\eta_t = \frac{\beta}{t+\gamma}$ where $\beta > \frac{1}{\bar{\mu}}$, $\gamma > 0$ such that $\eta_t^2 \leq 1/\left(2c_3(t+1-t_0)(1+2(t+1-t_0))\right)$, $\eta_0 \leq \hat{\eta}_0$, and $\eta_t \leq 2\eta_{t+E}$, we will prove that $\Delta_t \leq \frac{\upsilon}{\gamma+t}$ where $\upsilon = \max\left\{\frac{\beta^2 B}{\beta \bar{\mu}-1}, \gamma \Delta_0\right\} = \max\left\{\frac{4B}{\bar{\mu}^2}, \gamma \Delta_0\right\}$ by induction.

Firstly, $\Delta_0 \leq \frac{v}{\gamma}$ by the definition of v. Assume that for some $0 \leq t$, $\Delta_t \leq \frac{v}{\gamma + t}$, then

$$\Delta_{t+1} \leq (1 - \eta_t \tilde{\mu}) \Delta_t + \eta_t^2 B$$

$$\leq \left(1 - \frac{\beta \tilde{\mu}}{t + \gamma}\right) \frac{v}{t + \gamma} + \frac{\beta^2 B}{(t + \gamma)^2}$$

$$= \frac{t + \gamma - 1}{(t + \gamma)^2} v + \left[\frac{\beta^2 B}{(t + \gamma)^2} - \frac{\beta \tilde{\mu} - 1}{(t + \gamma)^2} v\right]$$

$$\leq \frac{v}{t + \gamma + 1}.$$

Specifically, if we choose $\beta=\frac{2}{\tilde{\mu}}, \gamma=\max\{\frac{2}{\tilde{\mu}\hat{\eta}_0}, E, \frac{2}{\tilde{\mu}}\sqrt{2E(2E+1)(12\sigma^2+18L^2(1+\epsilon_{\max})^2)}\}$, then we have

$$\eta_0 = \frac{\beta}{\gamma} \le \frac{2}{\tilde{\mu} \frac{2}{\tilde{\mu} \hat{\eta}_0}} = \hat{\eta}_0$$

and

$$\eta_t - 2\eta_{t+E} = \frac{\beta}{\gamma + t} - \frac{2\beta}{\gamma + t + E} = \frac{\beta(E - \gamma - t)}{(\gamma + t)(\gamma + t + E)} \le \frac{\beta(E - \gamma)}{(\gamma + t)(\gamma + t + E)} \le 0.$$

To prove that $\eta_t^2 \le 1/(2c_3(t+1-t_0)(1+2(t+1-t_0)))$ for any t, it suffices to prove that for $0 \le t \le E-1$ because $\{\eta_t\}$, i.e., $t_0=0$, is non-increasing and $t+1-t_0$ is periodic with period E. When $t_0=0$, we need to prove $\eta_t^2 \le 1/(2c_3(t+1)(1+2(t+1)))$ for $0 \le t \le E-1$, which is satisfied if

$$\begin{split} \max_{0 \leq t \leq E-1} \eta_t &\leq \min_{0 \leq t \leq E-1} \sqrt{1/\big(2c_3(t+1)(1+2(t+1))\big)} \\ \Longleftrightarrow \eta_0 \leq \sqrt{\frac{1}{2E(2E+1)c_3}} \\ \Longleftrightarrow \gamma \geq \beta \sqrt{2E(2E+1)c_3} &= \frac{2}{\tilde{\mu}} \sqrt{2E(2E+1)c_3} = \frac{2}{\tilde{\mu}} \sqrt{2E(2E+1)(12\sigma^2+18L^2(1+\epsilon_{\max})^2)}. \end{split}$$

B.4 Deferred Proofs of Key Lemmas

Proof of Lemma B.1. This proof follows from Lemma 3 in Li et al. [2022]. We first decompose $\|\overline{\boldsymbol{w}}^{t+1} - \boldsymbol{\theta}^{PS}\|_2^2$ as

$$\mathbb{E}\|\overline{\boldsymbol{w}}^{t+1} - \boldsymbol{\theta}^{PS}\|_{2}^{2} = \mathbb{E}\|\overline{\boldsymbol{\theta}}^{t} - \eta_{t}\boldsymbol{g}_{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2} = \mathbb{E}\|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2} - 2\eta_{t}\mathbb{E}\langle\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}, \boldsymbol{g}_{t}\rangle + \eta_{t}^{2}\mathbb{E}\|\boldsymbol{g}_{t}\|_{2}^{2}. \tag{8}$$

Next we present an upper bound for $\mathbb{E}\|\boldsymbol{g}_t\|_2^2$. By the definition of $\boldsymbol{\theta}^{PS}$, we have $\sum_{i=1}^N p_i \nabla f_i(\boldsymbol{\theta}^{PS}; \boldsymbol{\theta}^{PS}) = \mathbf{0}$, and thus

$$\begin{split} \mathbb{E}\|\boldsymbol{g}_t\|_2^2 &= \mathbb{E}\|\sum_{i=1}^N p_i [\nabla \ell(\boldsymbol{\theta}_i^t; Z_i^{t+1}) - \nabla f_i(\boldsymbol{\theta}_i^t; \boldsymbol{\theta}_i^t) + \nabla f_i(\boldsymbol{\theta}_i^t; \boldsymbol{\theta}_i^t) - \nabla f_i(\boldsymbol{\theta}^{PS}; \boldsymbol{\theta}^{PS})]\|_2^2 \\ &\leq 2\mathbb{E}\|\sum_{i=1}^N \nabla \ell(\boldsymbol{\theta}_i^t; Z_i^{t+1}) - \nabla f_i(\boldsymbol{\theta}_i^t; \boldsymbol{\theta}_i^t)\|_2^2 + 2\mathbb{E}\|\sum_{i=1}^N p_i [\nabla f_i(\boldsymbol{\theta}_i^t; \boldsymbol{\theta}_i^t) - \nabla f_i(\boldsymbol{\theta}^{PS}; \boldsymbol{\theta}^{PS})]\|_2^2 \\ &\leq 2\sum_{i=1}^N p_i \mathbb{E}[\|\nabla \ell(\boldsymbol{\theta}_i^t; Z_i^{t+1}) - \nabla f_i(\boldsymbol{\theta}_i^t; \boldsymbol{\theta}_i^t)\|_2^2] + 2\sum_{i=1}^N p_i \mathbb{E}[\|\nabla f_i(\boldsymbol{\theta}_i^t; \boldsymbol{\theta}_i^t) - \nabla f_i(\boldsymbol{\theta}^{PS}; \boldsymbol{\theta}^{PS})\|_2^2] \\ &\leq 2\sum_{i=1}^N p_i \sigma^2 \left(1 + \mathbb{E}\|\boldsymbol{\theta}_i^t - \boldsymbol{\theta}^{PS}\|_2^2\right) + 2\sum_{i=1}^N p_i L^2 (1 + \epsilon_i)^2 \mathbb{E}\|\boldsymbol{\theta}_i^t - \boldsymbol{\theta}^{PS}\|_2^2 \end{split}$$

where the second inequality is due to the convexity of 2-norm and the last inequality is due to Assumption 2.5 and Lemma 2.4. Since $\|\boldsymbol{\theta}_i^t - \boldsymbol{\theta}^{PS}\|_2^2 \le 2\|\boldsymbol{\theta}_i^t - \overline{\boldsymbol{\theta}}^t\|_2^2 + 2\|\overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}\|_2^2$ and $\epsilon_i \le \epsilon_{max}$, we have

$$\mathbb{E}[\|\boldsymbol{g}_{t}\|_{2}^{2}] \leq 2\sigma^{2} + 4[\sigma^{2} + L^{2}(1 + \epsilon_{max})^{2}]\mathbb{E}\|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + 4[\sigma^{2} + L^{2}(1 + \epsilon_{max})^{2}] \sum_{i=1}^{N} p_{i}\mathbb{E}\|\boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t}\|_{2}^{2}$$

$$= 2\sigma^{2} + c_{2}\mathbb{E}\|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + c_{2}\sum_{i=1}^{N} p_{i}\mathbb{E}\|\boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t}\|_{2}^{2}. \tag{9}$$

Next, we focus on establishing a lower bound for $\mathbb{E}[\langle \overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}, \ \boldsymbol{g}_t \rangle]$. By the law of total expectation and $\sum_{i=1}^N p_i \nabla f_i(\boldsymbol{\theta}^{PS}; \boldsymbol{\theta}^{PS}) = \mathbf{0}$, we have

$$\begin{split} \mathbb{E} \langle \overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}, \; \boldsymbol{g}_t \rangle &= \mathbb{E} \big[\mathbb{E}_t \langle \overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}, \; \boldsymbol{g}_t \rangle \big] \\ &= \mathbb{E} \big[\langle \overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}, \; \overline{\boldsymbol{g}}_t \rangle \big] \\ &= \mathbb{E} \bigg[\underbrace{\sum_{i=1}^N p_i \langle \overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}, \; \nabla f_i(\boldsymbol{\theta}_i^t; \boldsymbol{\theta}_i^t) - \nabla f_i(\overline{\boldsymbol{\theta}}^t; \boldsymbol{\theta}^{PS}) \rangle}_{A} \\ &+ \underbrace{\sum_{i=1}^N p_i \langle \overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}, \; \nabla f_i(\overline{\boldsymbol{\theta}}^t; \boldsymbol{\theta}^{PS}) - \nabla f_i(\boldsymbol{\theta}^{PS}; \boldsymbol{\theta}^{PS}) \rangle}_{B} \bigg]. \end{split}$$

On the one hand, applying Cauchy-Schwarz inequality and Lemma 2.4, we have

$$\begin{split} A &\geq - \|\overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}\|_2 \sum_{i=1}^N p_i \left(L \|\boldsymbol{\theta}_i^t - \overline{\boldsymbol{\theta}}^t\|_2 + L\epsilon_i \|\boldsymbol{\theta}_i^t - \boldsymbol{\theta}^{PS}\|_2 \right) \\ &\geq - \|\overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}\|_2 \sum_{i=1}^N p_i \left(L(1+\epsilon_i) \|\boldsymbol{\theta}_i^t - \overline{\boldsymbol{\theta}}^t\|_2 + L\epsilon_i \|\overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}\|_2 \right) \\ &\geq - L\overline{\epsilon} \|\overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}\|_2^2 - L(1+\epsilon_{max}) \sum_{i=1}^N p_i \|\overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}\|_2 \|\boldsymbol{\theta}_i^t - \overline{\boldsymbol{\theta}}^t\|_2. \end{split}$$

On the other hand, with the strong convexity in Assumption 2.1, we have $B \ge \mu \|\overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}\|_2^2$. Therefore, for any $\alpha > 0$, using the lower bounds on A, B, and the Young's inequality shows that

$$\mathbb{E}[\langle \overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}, \boldsymbol{g}_{t} \rangle] \\
\geq (\mu - L\overline{\epsilon}) \mathbb{E}[\|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2} - L(1 + \epsilon_{max}) \sum_{i=1}^{N} p_{i} \mathbb{E}[\|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2} \|\boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t}\|_{2}] \\
\geq (\mu - L\overline{\epsilon} - \frac{\alpha}{2} L(1 + \epsilon_{max})) \mathbb{E}[\|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2} - \frac{L(1 + \epsilon_{max})}{2\alpha} \sum_{i=1}^{N} p_{i} \mathbb{E}[\|\boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t}\|_{2}^{2} \\
\geq (\mu - (1 + \delta) L\overline{\epsilon}) \mathbb{E}[\|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2} - \frac{L(1 + \epsilon_{max})^{2}}{4\delta\overline{\epsilon}} \sum_{i=1}^{N} p_{i} \mathbb{E}[\|\boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t}\|_{2}^{2} \tag{10}$$

where we have set $\alpha:=\frac{2\delta\overline{\epsilon}}{1+\epsilon_{max}}$ in the last line.

Recall that we denote

$$c_1 := \frac{L(1 + \epsilon_{max})^2}{2\delta \overline{\epsilon}}, \ c_2 := 4[\sigma^2 + L^2(1 + \epsilon_{max})^2], \ \tilde{\mu} := \mu - (1 + \delta)\overline{\epsilon}L.$$

Combining (8), (9), (10), we have

$$\mathbb{E}[\|\overline{\boldsymbol{w}}^{t+1} - \boldsymbol{\theta}^{PS}\|_{2}^{2}] \\
\leq \mathbb{E}\|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2} \\
- 2\eta_{t} \left[(\mu - (1+\delta)L\overline{\epsilon})\mathbb{E}\|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2} - \frac{L(1+\epsilon_{max})}{4\delta\overline{\epsilon}} \sum_{i=1}^{N} p_{i}\mathbb{E}\|\boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t}\|_{2}^{2} \right] \\
+ \eta_{t}^{2} \left[2\sigma^{2} + c_{2}\mathbb{E}\|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + c_{2} \sum_{i=1}^{N} p_{i}\mathbb{E}\|\boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t}\|_{2}^{2} \right] \\
= (1 - 2\tilde{\mu}\eta_{t} + c_{2}\eta_{t}^{2})\mathbb{E}\|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + (c_{1}\eta_{t} + c_{2}\eta_{t}^{2}) \sum_{i=1}^{N} p_{i}\mathbb{E}\|\boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t}\|_{2}^{2} + 2\sigma^{2}\eta_{t}^{2} \\
\leq (1 - \tilde{\mu}\eta_{t})\mathbb{E}\|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + (c_{1}\eta_{t} + c_{2}\eta_{t}^{2}) \sum_{i=1}^{N} p_{i}\mathbb{E}\|\boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t}\|_{2}^{2} + 2\sigma^{2}\eta_{t}^{2}$$

where the last inequality is obtained by observing the condition $\eta_t \leq \tilde{\mu}/c_2$.

Proof of Lemma B.2. In this proof, for convenience, we will discuss with respect to t+1 where we assume $t+1 \notin \mathcal{I}_E$ and transfer back to t in the last. First by the update rule, we have

$$\boldsymbol{\theta}_i^{t+1} - \overline{\boldsymbol{\theta}}^{t+1} = \boldsymbol{\theta}_i^t - \overline{\boldsymbol{\theta}}^t - \eta_t(\nabla \ell(\boldsymbol{\theta}_i^t; Z_i^{t+1}) - g_t).$$

Using Young's inequality, we have

$$\sum_{i=1}^{N} p_{i} \mathbb{E} \|\boldsymbol{\theta}_{i}^{t+1} - \overline{\boldsymbol{\theta}}^{t+1}\|_{2}^{2} = \sum_{i=1}^{N} p_{i} \mathbb{E} \|\boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t} - \eta_{t} (\nabla \ell(\boldsymbol{\theta}_{i}^{t}; Z_{i}^{t+1}) - g_{t})\|_{2}^{2}$$

$$\leq (1 + \alpha_{t}) \sum_{i=1}^{N} p_{i} \mathbb{E} \|\boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t}\|_{2}^{2} + \eta_{t}^{2} (1 + \alpha_{t}^{-1}) \underbrace{\sum_{i=1}^{N} p_{i} \mathbb{E} \|\nabla \ell(\boldsymbol{\theta}_{i}^{t}; Z_{i}^{t+1}) - g_{t}\|_{2}^{2}}_{B} \tag{11}$$

where $\alpha_t > 0$ is a free chosen parameter. Next, we are going to establish an upper bound for B. Notice that

$$\begin{split} B &= \sum_{i=1}^{N} p_{i} \mathbb{E} \|\nabla \ell(\boldsymbol{\theta}_{i}^{t}; Z_{i}^{t+1}) - g_{t}\|_{2}^{2} \\ &= \mathbb{E} \bigg[\sum_{i=1}^{N} p_{i} \|\nabla \ell(\boldsymbol{\theta}_{i}^{t}; Z_{i}^{t+1}) - \sum_{j=1}^{N} p_{j} \nabla \ell(\boldsymbol{\theta}_{j}^{t}; Z_{j}^{t+1}) \|_{2}^{2} \bigg] \\ &= \mathbb{E} \bigg[\sum_{i=1}^{N} p_{i} \|\nabla \ell(\boldsymbol{\theta}_{i}^{t}; Z_{i}^{t+1}) - \nabla f_{i} \left(\boldsymbol{\theta}_{i}^{t}, \boldsymbol{\theta}_{i}^{t}\right) + \nabla f_{i} \left(\boldsymbol{\theta}_{i}^{t}, \boldsymbol{\theta}_{i}^{t}\right) - \sum_{j=1}^{N} p_{j} \nabla f_{j} (\boldsymbol{\theta}_{j}^{t}; \boldsymbol{\theta}_{j}^{t}) \\ &+ \sum_{j=1}^{N} p_{j} \nabla f_{j} (\boldsymbol{\theta}_{i}^{t}; Z_{i}^{t+1}) - \nabla f_{i} \left(\boldsymbol{\theta}_{i}^{t}, \boldsymbol{\theta}_{i}^{t}\right) \|_{2}^{2} \bigg] \\ &\leq 3 \mathbb{E} \bigg[\sum_{i=1}^{N} p_{i} \|\nabla \ell \left(\boldsymbol{\theta}_{i}^{t}; Z_{i}^{t+1}\right) - \nabla f_{i} \left(\boldsymbol{\theta}_{i}^{t}, \boldsymbol{\theta}_{i}^{t}\right) \|_{2}^{2} \bigg] + 3 \mathbb{E} \bigg[\sum_{i=1}^{N} p_{i} \|\nabla f_{i} \left(\boldsymbol{\theta}_{i}^{t}, \boldsymbol{\theta}_{i}^{t}\right) - \sum_{j=1}^{N} p_{j} \nabla \ell(\boldsymbol{\theta}_{j}^{t}; Z_{j}^{t+1}) \|_{2}^{2} \bigg] \\ &\leq 3 \mathbb{E} \bigg[\sum_{i=1}^{N} p_{i} \|\sum_{j=1}^{N} p_{j} \nabla f_{j} (\boldsymbol{\theta}_{j}^{t}; \boldsymbol{\theta}_{j}^{t}) - \sum_{j=1}^{N} p_{j} \nabla \ell(\boldsymbol{\theta}_{j}^{t}; Z_{j}^{t+1}) \|_{2}^{2} \bigg] \\ &\leq 3 \mathbb{E} \bigg[\sum_{i=1}^{N} p_{i} \|\nabla \ell \left(\boldsymbol{\theta}_{i}^{t}; Z_{i}^{t+1}\right) - \nabla f_{i} \left(\boldsymbol{\theta}_{i}^{t}, \boldsymbol{\theta}_{i}^{t}\right) \|_{2}^{2} \bigg] + 3 \mathbb{E} \bigg[\sum_{i=1}^{N} p_{i} \|\nabla f_{i} \left(\boldsymbol{\theta}_{i}^{t}, \boldsymbol{\theta}_{i}^{t}\right) - \sum_{j=1}^{N} p_{j} \nabla f_{j} \left(\boldsymbol{\theta}_{j}^{t}, \boldsymbol{\theta}_{j}^{t}\right) \|_{2}^{2} \bigg] \\ &+ 3 \mathbb{E} \bigg[\sum_{i=1}^{N} p_{i} \|\nabla \ell \left(\boldsymbol{\theta}_{i}^{t}; Z_{i}^{t+1}\right) - \nabla f_{i} \left(\boldsymbol{\theta}_{i}^{t}, \boldsymbol{\theta}_{i}^{t}\right) \|_{2}^{2} \bigg] + 3 \mathbb{E} \bigg[\sum_{i=1}^{N} p_{i} \|\nabla f_{i} \left(\boldsymbol{\theta}_{i}^{t}, \boldsymbol{\theta}_{i}^{t}\right) - \sum_{j=1}^{N} p_{j} \nabla f_{j} \left(\boldsymbol{\theta}_{j}^{t}, \boldsymbol{\theta}_{j}^{t}\right) \|_{2}^{2} \bigg] \\ &+ 3 \mathbb{E} \bigg[\sum_{i=1}^{N} p_{i} \|\nabla \ell \left(\boldsymbol{\theta}_{i}^{t}; Z_{i}^{t+1}\right) - \nabla \ell \left(\boldsymbol{\theta}_{i}^{t}; Z_{j}^{t+1}\right) \|_{2}^{2} \bigg] \\ &\leq 6 \sigma^{2} \bigg(1 + \mathbb{E} \bigg[\sum_{i=1}^{N} p_{i} \|\boldsymbol{\theta}_{i}^{t} - \boldsymbol{\theta}^{PS} \|_{2}^{2} \bigg] \bigg) \\ &+ 3 \mathbb{E} \bigg[\sum_{i=1}^{N} p_{i} \|\nabla f_{i} \left(\boldsymbol{\theta}_{i}^{t}, \boldsymbol{\theta}_{i}^{t}\right) - \sum_{j=1}^{N} p_{j} \nabla f_{j} \left(\boldsymbol{\theta}_{j}^{t}, \boldsymbol{\theta}_{j}^{t}\right) \|_{2}^{2} \bigg] \\ &\leq 6 \sigma^{2} \bigg(1 + \mathbb{E} \bigg[\sum_{i=1}^{N} p_{i} \|\boldsymbol{\theta}_{i}^{t} - \boldsymbol{\theta}^{PS} \|_{2}^{2} \bigg] \bigg) \\ \end{aligned}$$

where the last inequality is by Assumption 2.5. On the other hand, we have

$$3\mathbb{E}\left[\sum_{i=1}^{N} p_{i} \|\nabla f_{i}(\boldsymbol{\theta}_{i}^{t}, \boldsymbol{\theta}_{i}^{t}) - \sum_{j=1}^{N} p_{j} \nabla f_{j}(\boldsymbol{\theta}_{j}^{t}, \boldsymbol{\theta}_{j}^{t}) \|_{2}^{2}\right]$$

$$= 3\sum_{i=1}^{N} p_{i} \mathbb{E}\left\|\nabla f_{i}(\boldsymbol{\theta}_{i}^{t}, \boldsymbol{\theta}_{i}^{t}) - \nabla f_{i}(\overline{\boldsymbol{\theta}}^{t}, \overline{\boldsymbol{\theta}}^{t}) + \nabla f_{i}(\overline{\boldsymbol{\theta}}^{t}, \overline{\boldsymbol{\theta}}^{t}) - \sum_{j=1}^{N} p_{j} \nabla f_{j}(\overline{\boldsymbol{\theta}}^{t}, \overline{\boldsymbol{\theta}}^{t}) - \sum_{j=1}^{N} p_{j} \left(\nabla f_{j}(\boldsymbol{\theta}_{j}^{t}, \boldsymbol{\theta}_{j}^{t}) - \nabla f_{j}(\overline{\boldsymbol{\theta}}^{t}, \overline{\boldsymbol{\theta}}^{t})\right) \|_{2}^{2}$$

$$\leq 9\sum_{i=1}^{N} p_{i} \mathbb{E}\|\nabla f_{i}(\boldsymbol{\theta}_{i}^{t}, \boldsymbol{\theta}_{i}^{t}) - \nabla f_{i}(\overline{\boldsymbol{\theta}}^{t}, \overline{\boldsymbol{\theta}}^{t}) \|_{2}^{2} + 9\sum_{i=1}^{N} p_{i} \mathbb{E}\|\nabla f_{i}(\overline{\boldsymbol{\theta}}^{t}, \overline{\boldsymbol{\theta}}^{t}) - \sum_{j=1}^{N} p_{j} \nabla f_{j}(\overline{\boldsymbol{\theta}}^{t}, \overline{\boldsymbol{\theta}}^{t}) \|_{2}^{2}$$

$$+ 9\sum_{i=1}^{N} p_{i} \mathbb{E}\|\sum_{j=1}^{N} p_{j} \left(\nabla f_{j}(\boldsymbol{\theta}_{j}^{t}, \boldsymbol{\theta}_{j}^{t}) - \nabla f_{j}(\overline{\boldsymbol{\theta}}^{t}, \overline{\boldsymbol{\theta}}^{t})\right) \|_{2}^{2}.$$

Using Lemma 2.4, we have

$$B_{1} \leq 9 \sum_{i=1}^{N} p_{i} L^{2} (1 + \epsilon_{i})^{2} \mathbb{E} \|\boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t}\|_{2}^{2} \leq 9 \sum_{i=1}^{N} p_{i} L^{2} (1 + \epsilon_{max})^{2} \mathbb{E} \|\boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t}\|_{2}^{2}.$$

Using Assumption 2.6, we have

$$B_{2} = 9 \sum_{i=1}^{N} p_{i} \mathbb{E} \| \nabla f_{i}(\overline{\boldsymbol{\theta}}^{t}, \overline{\boldsymbol{\theta}}^{t}) - \nabla f(\overline{\boldsymbol{\theta}}^{t}, \overline{\boldsymbol{\theta}}^{t}) \|_{2}^{2}$$

$$\leq 9 \sum_{i=1}^{N} p_{i} \varsigma^{2} (1 + \mathbb{E} \| \overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS} \|_{2}^{2})$$

$$= 9 \varsigma^{2} + 9 \varsigma^{2} \mathbb{E} \| \overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS} \|_{2}^{2}.$$

Using Lemma 2.4, we have

$$B_3 \leq 9 \sum_{i=1}^N p_i \sum_{j=1}^N p_j \mathbb{E} \|\nabla f_j(\boldsymbol{\theta}_j^t, \boldsymbol{\theta}_j^t) - \nabla f_j(\overline{\boldsymbol{\theta}}^t, \overline{\boldsymbol{\theta}}^t)\|_2^2 \leq 9 \sum_{i=1}^N p_i L^2 (1 + \epsilon_{max})^2 \mathbb{E} \|\boldsymbol{\theta}_i^t - \overline{\boldsymbol{\theta}}^t\|_2^2.$$

Therefore,

$$B_1 + B_2 + B_3 \le 18 \sum_{i=1}^{N} p_i L^2 (1 + \epsilon_{max})^2 \mathbb{E} \|\boldsymbol{\theta}_i^t - \overline{\boldsymbol{\theta}}^t\|_2^2 + 9\varsigma^2 + 9\varsigma^2 \sum_{i=1}^{N} p_i \mathbb{E} \|\overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}\|_2^2,$$

which results in that

$$B \leq 6\sigma^{2} \left(1 + \mathbb{E} \left[\sum_{i=1}^{N} p_{i} \| \boldsymbol{\theta}_{i}^{t} - \boldsymbol{\theta}^{PS} \|_{2}^{2} \right] \right) + 18 \sum_{i=1}^{N} p_{i} L^{2} (1 + \epsilon_{max})^{2} \mathbb{E} \| \boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t} \|_{2}^{2} + 9\varsigma^{2} \mathbb{E} \| \overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS} \|_{2}^{2}$$

$$\leq 6\sigma^{2} \left(1 + 2 \sum_{i=1}^{N} p_{i} \mathbb{E} \| \boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t} \|_{2}^{2} + 2 \mathbb{E} \| \overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS} \|_{2}^{2} + 18 \sum_{i=1}^{N} p_{i} L^{2} (1 + \epsilon_{max})^{2} \mathbb{E} \| \boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t} \|_{2}^{2} + 9\varsigma^{2} \mathbb{E} \| \overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS} \|_{2}^{2} \right)$$

$$= 6\sigma^{2} + 9\varsigma^{2} + \left(12\sigma^{2} + 18L^{2} (1 + \epsilon_{max})^{2} \right) \sum_{i=1}^{N} p_{i} \mathbb{E} \| \boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t} \|_{2}^{2} + (12\sigma^{2} + 9\varsigma^{2}) \mathbb{E} \| \overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS} \|_{2}^{2}.$$

Inserting this formula into (11), we obtain

$$\begin{split} &\sum_{i=1}^{N} p_{i} \mathbb{E} \| \boldsymbol{\theta}_{i}^{t+1} - \overline{\boldsymbol{\theta}}^{t+1} \|_{2}^{2} \\ &\leq (1 + \alpha_{t}) \sum_{i=1}^{N} p_{i} \mathbb{E} \| \boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t} \|_{2}^{2} \\ &\quad + \eta_{t}^{2} (1 + \alpha_{t}^{-1}) \bigg(6\sigma^{2} + 9\varsigma^{2} + \left(12\sigma^{2} + 18L^{2}(1 + \epsilon_{\max})^{2} \right) \sum_{i=1}^{N} p_{i} \mathbb{E} \| \boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t} \|_{2}^{2} + \left(12\sigma^{2} + 9\varsigma^{2} \right) \mathbb{E} \| \overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS} \|_{2}^{2} \bigg) \\ &= \left(1 + \alpha_{t} + \eta_{t}^{2} (1 + \alpha_{t}^{-1}) \left(12\sigma^{2} + 18L^{2}(1 + \epsilon_{\max})^{2} \right) \right) \sum_{i=1}^{N} p_{i} \mathbb{E} \| \boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t} \|_{2}^{2} + \eta_{t}^{2} (1 + \alpha_{t}^{-1}) (12\sigma^{2} + 9\varsigma^{2}) \mathbb{E} \| \overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS} \|_{2}^{2} \\ &\quad + \eta_{t}^{2} (1 + \alpha_{t}^{-1}) (6\sigma^{2} + 9\varsigma^{2}) \end{split}$$

where $\alpha_t > 0$ is a free chosen parameter. Let $t_0 := \max\{s \mid s < t+1, s \in \mathcal{I}_E\}$ and $c_3 := 12\sigma^2 + 18L^2(1+\epsilon_{\max})^2$. Then we choose $\alpha_t = \frac{1}{2(t+1-t_0)}$, if we have

$$\eta_t^2 (1 + \alpha_t^{-1}) \left(12\sigma^2 + 18L^2 (1 + \epsilon_{\text{max}})^2 \right) = \eta_t^2 (1 + \alpha_t^{-1}) c_3 \le \frac{1}{2(t + 1 - t_0)}$$

$$\iff \eta_t^2 \le \frac{1}{2c_3(t + 1 - t_0) \left(1 + 2(t + 1 - t_0) \right)},$$
(12)

then note that $1 + \alpha_t^{-1} = 1 + 2(t+1-t_0) \le 2E - 1$,

$$\begin{split} & \sum_{i=1}^{N} p_{i} \mathbb{E} \|\boldsymbol{\theta}_{i}^{t+1} - \overline{\boldsymbol{\theta}}^{t+1}\|_{2}^{2} \\ & \leq \frac{t+2-t_{0}}{t+1-t_{0}} \sum_{i=1}^{N} p_{i} \mathbb{E} \|\boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t}\|_{2}^{2} + \eta_{t}^{2} (2E-1)(12\sigma^{2}+9\varsigma^{2}) \mathbb{E} \|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + \eta_{t}^{2} (2E-1)(6\sigma^{2}+9\varsigma^{2}). \end{split}$$

Continuing the above expansion until t_0 and leveraging $\eta_s \leq \eta_{t_0} \leq 2\eta_{t_0+E} \leq 2\eta_t$ gives us

$$\sum_{i=1}^{N} p_{i} \mathbb{E} \|\boldsymbol{\theta}_{i}^{t+1} - \overline{\boldsymbol{\theta}}^{t+1}\|_{2}^{2} \\
\leq \frac{t+2-t_{0}}{t_{0}+1-t_{0}} \sum_{i=1}^{N} p_{i} \mathbb{E} \|\boldsymbol{\theta}_{i}^{t_{0}} - \overline{\boldsymbol{\theta}}^{t_{0}}\|_{2}^{2} + \sum_{s=t_{0}}^{t} \frac{t+2-t_{0}}{s+2-t_{0}} \eta_{s}^{2} (2E-1)(12\sigma^{2}+9\varsigma^{2}) \mathbb{E} \|\overline{\boldsymbol{\theta}}^{s} - \boldsymbol{\theta}^{PS}\|_{2}^{2} \\
+ \sum_{s=t_{0}}^{t} \frac{t+2-t_{0}}{s+2-t_{0}} \eta_{s}^{2} (2E-1)(6\sigma^{2}+9\varsigma^{2}) \\
= \sum_{s=0}^{t-t_{0}} \frac{t+2-t_{0}}{s+2} \eta_{s}^{2} (2E-1)(12\sigma^{2}+9\varsigma^{2}) \mathbb{E} \|\overline{\boldsymbol{\theta}}^{s} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + \sum_{s=0}^{t-t_{0}} \frac{t+2-t_{0}}{s+2} \eta_{s}^{2} (2E-1)(6\sigma^{2}+9\varsigma^{2}) \\
\leq \sum_{s=0}^{t-t_{0}} \frac{t+2-t_{0}}{s+2} \eta_{t}^{2} (2E-1)(48\sigma^{2}+36\varsigma^{2}) \mathbb{E} \|\overline{\boldsymbol{\theta}}^{s} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + \sum_{s=0}^{t-t_{0}} \frac{t+2-t_{0}}{s+2} \eta_{t}^{2} (2E-1)(24\sigma^{2}+36\varsigma^{2}). \tag{13}$$

With the above formula and Lemma B.1, we now prove that if η_0 is sufficiently small, the for any t, we have $\mathbb{E}\|\overline{\boldsymbol{\theta}}^t - {\boldsymbol{\theta}}^{PS}\|_2^2 \leq \mathbb{E}\|\overline{\boldsymbol{\theta}}^0 - {\boldsymbol{\theta}}^{PS}\|_2^2$. We first derive the following inequality, which we will use later. Note that for any t where $t_0 := \max\{s \mid s < t+1, s \in \mathcal{I}_E\}$, we have

$$\sum_{s=0}^{t-t_0} \frac{t+2-t_0}{s+2} = (t+2-t_0)(\frac{1}{2}+\ldots+\frac{1}{t-t_0+2}) \le (t+2-t_0)\log(t+2-t_0) \le E\log E.$$
 (14)

We prove $\mathbb{E}\|\overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}\|_2^2 \leq \mathbb{E}\|\overline{\boldsymbol{\theta}}^0 - \boldsymbol{\theta}^{PS}\|_2^2$ by induction.

First, this inequality clearly holds for t=0. Suppose it holds for $0 \le s \le t$ where $t \le E-1$. Then by Lemma B.1 and (13), we have

$$\begin{split} & \mathbb{E}[\|\overline{\boldsymbol{\theta}}^{t+1} - \boldsymbol{\theta}^{PS}\|_{2}^{2}] \\ &= \mathbb{E}[\|\overline{\boldsymbol{w}}^{t+1} - \boldsymbol{\theta}^{PS}\|_{2}^{2}] \\ &\leq (1 - \tilde{\mu}\eta_{t})\mathbb{E}\|\overline{\boldsymbol{\theta}}^{0} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + 2\sigma^{2}\eta_{t}^{2} + (c_{1}\eta_{t} + c_{2}\eta_{t}^{2}) \bigg(\sum_{s=0}^{t-t_{0}-1} \frac{t+1-t_{0}}{s+2} \eta_{t-1}^{2} (2E-1)(48\sigma^{2} + 36\varsigma^{2}) \mathbb{E}\|\overline{\boldsymbol{\theta}}^{0} - \boldsymbol{\theta}^{PS}\|_{2}^{2} \\ &+ \sum_{s=0}^{t-t_{0}-1} \frac{t+1-t_{0}}{s+2} \eta_{t-1}^{2} (2E-1)(24\sigma^{2} + 36\varsigma^{2}) \bigg) \\ &= (1 - \tilde{\mu}\eta_{t}) \mathbb{E}\|\overline{\boldsymbol{\theta}}^{0} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + 2\sigma^{2}\eta_{t}^{2} \\ &+ (c_{1}\eta_{t} + c_{2}\eta_{t}^{2}) \eta_{t-1}^{2} (2E^{2} - E) \log E \bigg((48\sigma^{2} + 36\varsigma^{2}) \mathbb{E}\|\overline{\boldsymbol{\theta}}^{0} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + (24\sigma^{2} + 36\varsigma^{2}) \bigg). \end{split}$$

By (12), we need $\eta_0 \leq \frac{1}{6c_3}$, together with $\eta_{t-1} \leq 2\eta_t$ implies

$$(c_1\eta_t + c_2\eta_t^2)\eta_{t-1}^2 \le \eta_t^2(c_1 + c_2(6c_3)^{-1})\frac{c_3}{3}.$$

Therefore, we have

$$\mathbb{E}[\|\overline{\boldsymbol{\theta}}^{t+1} - \boldsymbol{\theta}^{PS}\|_{2}^{2}] \\
\leq (1 - \tilde{\mu}\eta_{t})\mathbb{E}\|\overline{\boldsymbol{\theta}}^{0} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + 2\sigma^{2}\eta_{t}^{2} \\
+ \eta_{t}^{2}(c_{1}c_{3} + c_{2}/6)(2E^{2} - E)\log E\left((16\sigma^{2} + 12\varsigma^{2})\mathbb{E}\|\overline{\boldsymbol{\theta}}^{0} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + (8\sigma^{2} + 12\varsigma^{2})\right)$$

whose right-hand side is no larger than $\mathbb{E}\|\overline{\boldsymbol{\theta}}^0 - \boldsymbol{\theta}^{PS}\|_2^2$ if

$$\eta_{t} \leq \eta_{0}$$

$$\leq \frac{\tilde{\mu}\mathbb{E}\|\overline{\boldsymbol{\theta}}^{0} - \boldsymbol{\theta}^{PS}\|_{2}^{2}}{2\sigma^{2} + (c_{1}c_{3} + c_{2}/6)(2E^{2} - E)\log E((16\sigma^{2} + 12\varsigma^{2})\mathbb{E}\|\overline{\boldsymbol{\theta}}^{0} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + (8\sigma^{2} + 12\varsigma^{2}))} = \hat{\eta}_{0}.$$
(15)

Thus we have proved that for any $0 \le t \le E$, if (15) holds, then $\mathbb{E}\|\overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}\|_2^2 \le \mathbb{E}\|\overline{\boldsymbol{\theta}}^0 - \boldsymbol{\theta}^{PS}\|_2^2$. The same proof technique can be extended to any $nE \le t \le (n+1)E$ where $n \in \mathbb{N}_+$ and thus for any t, if (15) holds, then $\mathbb{E}\|\overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}\|_2^2 \le \mathbb{E}\|\overline{\boldsymbol{\theta}}^0 - \boldsymbol{\theta}^{PS}\|_2^2$.

Therefore, under (12), (15) and $\eta_{t-1} \leq 2\eta_t$, by (11) and (13), if $t \notin \mathcal{I}_E$, we have

$$\sum_{i=1}^{N} p_{i} \mathbb{E} \|\boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t}\|_{2}^{2} \leq \eta_{t-1}^{2} (2E^{2} - E) \log E(48\sigma^{2} + 36\varsigma^{2}) \mathbb{E} \|\overline{\boldsymbol{\theta}}^{0} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + \eta_{t-1}^{2} (2E^{2} - E) \log E(24\sigma^{2} + 36\varsigma^{2}) \\
\leq 4\eta_{t}^{2} (2E^{2} - E) \log E(48\sigma^{2} + 36\varsigma^{2}) \mathbb{E} \|\overline{\boldsymbol{\theta}}^{0} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + 4\eta_{t}^{2} (2E^{2} - E) \log E(24\sigma^{2} + 36\varsigma^{2}).$$

Proof of Lemma B.3. We discuss in two cases, $t \in \mathcal{I}_E$ and $t \notin \mathcal{I}_E$. If $t \in \mathcal{I}_E$, then we have $\theta_i^t = \overline{\theta}^t$, and by Lemma B.1,

$$\begin{split} & \mathbb{E}[\|\overline{\boldsymbol{\theta}}^{t+1} - \boldsymbol{\theta}^{PS}\|_{2}^{2}] \\ & \leq (1 - \tilde{\mu}\eta_{t})\mathbb{E}\|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + 2\sigma^{2}\eta_{t}^{2} + (c_{1}\eta_{t} + c_{2}\eta_{t}^{2}) \sum_{i=1}^{N} p_{i}\mathbb{E}\|\boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t}\|_{2}^{2} \\ & = (1 - \tilde{\mu}\eta_{t})\mathbb{E}\|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + 2\sigma^{2}\eta_{t}^{2} \\ & \leq (1 - \tilde{\mu}\eta_{t})\mathbb{E}\|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + B\eta_{t}^{2}. \end{split}$$

L

If $t \notin \mathcal{I}_E$, combining Lemma B.1 and Lemma B.2, we have

$$\begin{split} & \mathbb{E}[\|\overline{\boldsymbol{\theta}}^{t+1} - \boldsymbol{\theta}^{PS}\|_{2}^{2}] \\ & \leq (1 - \tilde{\mu}\eta_{t})\mathbb{E}\|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + 2\sigma^{2}\eta_{t}^{2} + (c_{1}\eta_{t} + c_{2}\eta_{t}^{2}) \sum_{i=1}^{N} p_{i}\mathbb{E}\|\boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t}\|_{2}^{2} \\ & \leq (1 - \tilde{\mu}\eta_{t})\mathbb{E}\|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + 2\sigma^{2}\eta_{t}^{2} \\ & \quad + (4c_{1}\eta_{t} + 4c_{2}\eta_{t}^{2}) \left(\eta_{t}^{2}(2E^{2} - E) \log E(48\sigma^{2} + 36\varsigma^{2})\mathbb{E}\|\overline{\boldsymbol{\theta}}^{0} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + \eta_{t}^{2}(2E^{2} - E) \log E(24\sigma^{2} + 36\varsigma^{2})\right) \\ & \leq (1 - \tilde{\mu}\eta_{t})\mathbb{E}\|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + 2B\eta_{t}^{2}. \end{split}$$

C Proof of Theorem 3.2 and Theorem 3.3

C.1 Additional Notations

Similar to Appendix B, in our analysis, for the sake of convenience, we will define two additional sequences as $\overline{\boldsymbol{w}}^t := \sum_{i=1}^N p_i \boldsymbol{w}_i^t$ and $\overline{\boldsymbol{\theta}}^t := \sum_{i=1}^N p_i \boldsymbol{\theta}_i^t$, following that of Li et al. [2020b]. We note that $\overline{\boldsymbol{w}}^t$ results from a single step of SGD from $\overline{\boldsymbol{\theta}}^t$. When $t+1 \notin \mathcal{I}_E$, both $\overline{\boldsymbol{w}}^t$ and $\overline{\boldsymbol{\theta}}^t$ are unaccessible. When $t+1 \in \mathcal{I}_E$, we can obtain $\overline{\boldsymbol{\theta}}^t$. In addition, we also define $\overline{\boldsymbol{g}}_t := \sum_{i=1}^N p_i \nabla f_i(\boldsymbol{\theta}_i^t; \boldsymbol{\theta}_i^t), \quad \boldsymbol{g}_t := \sum_{i=1}^N p_i \nabla \ell(\boldsymbol{\theta}_i^t; \mathcal{Z}_i^{t+1})$ where $Z_i^{t+1} \sim \mathcal{D}_i(\boldsymbol{\theta}_i^t)$. It is clear that in full participation, $\overline{\boldsymbol{w}}^{t+1} = \overline{\boldsymbol{w}}^t - \eta_t \boldsymbol{g}_t$ and $\mathbb{E} \boldsymbol{g}_t = \overline{\boldsymbol{g}}_t$. Notice now we do not have $\overline{\boldsymbol{\theta}}^t = \overline{\boldsymbol{w}}^t$ for any t. But we will show later that they are equal with expectation to the choice of \mathcal{S}_t .

In particular, in our analysis, there would be two types of randomness, one from the stochastic gradients and one from the random sampling of the devices. All analysis in Appendix B only involves the former. To make a distinguishment, we use \mathbb{E}_{S_t} to denote the latter.

C.2 Key Lemmas

We show in this subsection the key lemmas for proving the convergence and defer proofs to later parts. We first show that the sampling schemes I & II are unbiased.

Lemma C.1. Li et al. [2020b] (Unbiased sampling scheme). If $t + 1 \in \mathcal{I}_E$, for Scheme I and Scheme II, we have

$$\mathbb{E}_{\mathcal{S}_t} \big[\overline{m{ heta}}^{t+1} \big] = \overline{m{w}}^{t+1}.$$

Similar to Lemma B.1, we are going to establish an upper bound for $\mathbb{E}[\|\overline{\boldsymbol{\theta}}^{t+1} - \boldsymbol{\theta}^{PS}\|_2^2]$. When $t+1 \notin \mathcal{I}_E$, we have $\overline{\boldsymbol{\theta}}^{t+1} = \overline{\boldsymbol{w}}^{t+1}$ for both schemes, and therefore this is equivalent to establishing an upper bound for $\mathbb{E}[\|\overline{\boldsymbol{w}}^{t+1} - \boldsymbol{\theta}^{PS}\|_2^2]$. However, when $t+1 \in \mathcal{I}_E$, we only have $\mathbb{E}_{\mathcal{S}_t}[\overline{\boldsymbol{\theta}}^{t+1}] = \overline{\boldsymbol{w}}^{t+1}$ and we need other upper-bounding strategies.

Lemma C.2. Under Assumptions 2.1, 2.2, 2.3, 2.5, for scheme I & II:

1. if
$$t+1 \notin \mathcal{I}_E$$
,

$$\mathbb{E}[\|\overline{\boldsymbol{\theta}}^{t+1} - \boldsymbol{\theta}^{PS}\|_2^2] = \mathbb{E}[\|\overline{\boldsymbol{w}}^{t+1} - \boldsymbol{\theta}^{PS}\|_2^2] \leq (1 - \tilde{\mu}\eta_t)\mathbb{E}\|\overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}\|_2^2 + 2\sigma^2\eta_t^2 + (c_1\eta_t + c_2\eta_t^2)\sum_{i=1}^N p_i\mathbb{E}\|\boldsymbol{\theta}_i^t - \overline{\boldsymbol{\theta}}^t\|_2^2.$$

2. if $t + 1 \in \mathcal{I}_E$: for scheme I,

$$\begin{split} & \mathbb{E}[\|\overline{\boldsymbol{\theta}}^{t+1} - \boldsymbol{\theta}^{PS}\|_2^2] \\ & \leq \frac{1}{K} \sum_{k=1}^N p_k \mathbb{E}\|\boldsymbol{w}_k^{t+1} - \overline{\boldsymbol{w}}^{t+1}\|_2^2 + (1 - \tilde{\mu}\eta_t) \mathbb{E}\|\overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}\|_2^2 + 2\sigma^2 \eta_t^2 + (c_1\eta_t + c_2\eta_t^2) \sum_{i=1}^N p_i \mathbb{E}\|\boldsymbol{\theta}_i^t - \overline{\boldsymbol{\theta}}^t\|_2^2, \end{split}$$

while for scheme II,

$$\mathbb{E}[\|\overline{\boldsymbol{\theta}}^{t+1} - \boldsymbol{\theta}^{PS}\|_{2}^{2}]$$

$$\leq \frac{N}{K(N-1)} \left(1 - \frac{K}{N}\right) \sum_{i=1}^{N} p_{k} \mathbb{E}\|\boldsymbol{w}_{k}^{t+1} - \overline{\boldsymbol{w}}^{t+1}\|_{2}^{2} + (1 - \tilde{\mu}\eta_{t}) \mathbb{E}\|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + 2\sigma^{2}\eta_{t}^{2} + (c_{1}\eta_{t} + c_{2}\eta_{t}^{2}) \sum_{i=1}^{N} p_{i} \mathbb{E}\|\boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t}\|_{2}^{2}.$$

where
$$\epsilon_{max} := \max_i \epsilon_i, \overline{\epsilon} := \sum_{i=1}^N p_i \epsilon_i, c_1 := \frac{L(1+\epsilon_{max})^2}{2\delta \overline{\epsilon}}, c_2 := 4[\sigma^2 + L^2(1+\epsilon_{max})^2], \tilde{\mu} := \mu - (1+\delta)\overline{\epsilon}L.$$

To really give a descent lemma as in SGD analysis, we have to bound $\sum_{i=1}^{N} p_i \mathbb{E} \|\boldsymbol{\theta}_i^t - \overline{\boldsymbol{\theta}}^t\|_2^2$ for $t \notin \mathcal{I}_E$ and $\sum_{i=1}^{N} p_i \mathbb{E} \|\boldsymbol{w}_i^t - \overline{\boldsymbol{w}}^t\|_2^2$ for $t \in \mathcal{I}_E$, given by the following lemma.

Lemma C.3. Under Assumption 2.1, 2.2, 2.3, 2.5, 2.6, if $\{\eta_t\}$ is non-increasing, $\eta_t \leq 2\eta_{t+E}$, $t \notin \mathcal{I}_E$, $\eta_t^2 \leq 1/(2c_3(t+1-t_0)(1+2(t+1-t_0)))$, and

$$\eta_0 \leq \tilde{\eta}_0 := \frac{\tilde{\mu} \mathbb{E} \|\overline{\boldsymbol{\theta}}^0 - {\boldsymbol{\theta}}^{PS}\|_2^2}{2\sigma^2 + (c_1c_3 + c_2/6)(2E^2 + 3E + 1)\log(E + 1)((16\sigma^2 + 12\varsigma^2)\mathbb{E} \|\overline{\boldsymbol{\theta}}^0 - {\boldsymbol{\theta}}^{PS}\|_2^2 + (8\sigma^2 + 12\varsigma^2))},$$

then

1. for scheme I,

$$\mathbb{E}[\|\overline{\boldsymbol{\theta}}^{t+1} - \boldsymbol{\theta}^{PS}\|_{2}^{2}] \\
\leq (1 - \tilde{\mu}\eta_{t})\mathbb{E}\|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + 2\sigma^{2}\eta_{t}^{2} \\
+ (4c_{1}\eta_{t} + 4c_{2}\eta_{t}^{2} + K^{-1})(2E^{2} + 3E + 1)\log(E + 1)\eta_{t}^{2} \Big((48\sigma^{2} + 36\varsigma^{2})\mathbb{E}\|\overline{\boldsymbol{\theta}}^{0} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + (24\sigma^{2} + 36\varsigma^{2}) \Big),$$

2. for scheme II,

$$\mathbb{E}[\|\overline{\boldsymbol{\theta}}^{t+1} - \boldsymbol{\theta}^{PS}\|_{2}^{2}] \\
\leq (1 - \tilde{\mu}\eta_{t})\mathbb{E}\|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + 2\sigma^{2}\eta_{t}^{2} \\
+ \left(4c_{1}\eta_{t} + 4c_{2}\eta_{t}^{2} + \frac{N - K}{K(N - 1)}\right)(2E^{2} + 3E + 1)\log(E + 1)\eta_{t}^{2}\left((48\sigma^{2} + 36\varsigma^{2})\mathbb{E}\|\overline{\boldsymbol{\theta}}^{0} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + (24\sigma^{2} + 36\varsigma^{2})\right).$$

where for any t, where $\epsilon_{max} := \max_i \epsilon_i, \bar{\epsilon} := \sum_{i=1}^N p_i \epsilon_i, c_1 := \frac{L(1+\epsilon_{max})^2}{2\delta \bar{\epsilon}}, c_2 := 4[\sigma^2 + L^2(1+\epsilon_{max})^2], \tilde{\mu} := \mu - (1+\delta)\bar{\epsilon}L, c_3 := 12\sigma^2 + 18L^2(1+\epsilon_{max})^2.$

(One should note that $4c_1\eta_t + 4c_2\eta_t^2$, $(48\sigma^2 + 36\varsigma^2)$, and $(24\sigma^2 + 36\varsigma^2)$ comes from several times of applying $\eta_{t-1} \leq 2\eta_t$ and the real constants could be much smaller by choosing stepsizes carefully.)

The following lemma gives us a standard descent lemma in SGD analysis under technical conditions for establishing the $\mathcal{O}(\frac{1}{T})$ convergence in Theorem 3.2 and 3.3.

Lemma C.4. Under Assumption 2.1, 2.2, 2.3, 2.5, 2.6, if $\{\eta_t\}$ is non-increasing, $\eta_t \leq 2\eta_{t+E}$, $t \notin \mathcal{I}_E$, $\eta_t^2 \leq 1/(2c_3(t+1-t_0)(1+2(t+1-t_0)))$, and

$$\eta_0 \leq \tilde{\eta}_0 := \frac{\tilde{\mu} \mathbb{E} \|\overline{\boldsymbol{\theta}}^0 - {\boldsymbol{\theta}}^{PS}\|_2^2}{2\sigma^2 + (c_1c_3 + c_2/6)(2E^2 + 3E + 1)\log(E + 1)\left((16\sigma^2 + 12\varsigma^2)\mathbb{E} \|\overline{\boldsymbol{\theta}}^0 - {\boldsymbol{\theta}}^{PS}\|_2^2 + (8\sigma^2 + 12\varsigma^2)\right)},$$

then

1. for scheme I,

$$\mathbb{E}[\|\overline{\boldsymbol{\theta}}^{t+1} - \boldsymbol{\theta}^{PS}\|_{2}^{2}] \leq (1 - \tilde{\mu}\eta_{t})\mathbb{E}\|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + B_{1}\eta_{t}^{2},$$

with

$$B_1 := 2\sigma^2 + (4c_1\eta_t + 4c_2\eta_t^2 + K^{-1})(2E^2 + 3E + 1)\log(E + 1)\left((48\sigma^2 + 36\varsigma^2)\mathbb{E}\|\overline{\boldsymbol{\theta}}^0 - {\boldsymbol{\theta}}^{PS}\|_2^2 + (24\sigma^2 + 36\varsigma^2)\right),$$

2. for scheme II,

$$\mathbb{E}[\|\overline{\boldsymbol{\theta}}^{t+1} - \boldsymbol{\theta}^{PS}\|_2^2] \le (1 - \tilde{\mu}\eta_t)\mathbb{E}\|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_2^2 + B_2\eta_t^2$$

with

$$B_2 := 2\sigma^2 + \left(4c_1\eta_t + 4c_2\eta_t^2 + \frac{N - K}{K(N - 1)}\right)(2E^2 + 3E + 1)\log(E + 1)\left((48\sigma^2 + 36\varsigma^2)\mathbb{E}\|\overline{\boldsymbol{\theta}}^0 - \boldsymbol{\theta}^{PS}\|_2^2 + (24\sigma^2 + 36\varsigma^2)\right)$$

where for any t, where $\epsilon_{max} := \max_{i} \epsilon_{i}, \bar{\epsilon} := \sum_{i=1}^{N} p_{i} \epsilon_{i}, c_{1} := \frac{L(1+\epsilon_{max})^{2}}{2\delta \bar{\epsilon}}, c_{2} := 4[\sigma^{2} + L^{2}(1+\epsilon_{max})^{2}], \tilde{\mu} := \mu - (1+\delta)\bar{\epsilon}L, c_{3} := 12\sigma^{2} + 18L^{2}(1+\epsilon_{max})^{2}.$

C.3 Completing the proof of Theorem 3.2 and 3.3

We restate the definitions of all the constants here:

Constants independent of system design.

$$\epsilon_{max} := \max_i \epsilon_i,$$

$$\overline{\epsilon} := \sum_{i=1}^{N} p_i \epsilon_i,$$

$$\tilde{\mu} := \mu - (1 + \delta)\overline{\epsilon}L$$

$$c_1 := (L(1 + \epsilon_{max})^2)/(2\delta\overline{\epsilon}),$$

$$c_2 := 4 \left[\sigma^2 + L^2 (1 + \epsilon_{max})^2 \right],$$

$$c_3 := 6[2\sigma^2 + 3L^2(1 + \epsilon_{\text{max}})^2],$$

$$c_4 := 16\sigma^2 + 12\varsigma^2 + (8\sigma^2 + 12\varsigma^2)/\mathbb{E}\|\overline{\boldsymbol{\theta}}^0 - \boldsymbol{\theta}^{PS}\|_{2}^2,$$

$$c_5 := (48\sigma^2 + 36\varsigma^2) \mathbb{E} \|\overline{\boldsymbol{\theta}}^0 - {\boldsymbol{\theta}}^{PS}\|_2^2 + (24\sigma^2 + 36\varsigma^2).$$

Constants related to system design (e.g., E, K).

$$\hat{\eta}_0 := \tilde{\mu}/(2\sigma^2 + (c_1c_3 + c_2/6)c_4(2E^2 - E)\log E),$$

$$B := 2\sigma^2 + (4c_1\hat{\eta}_0 + 4c_2\hat{\eta}_0^2)c_5(2E^2 - E)\log E$$

$$c_6 := (2E^2 + 3E + 1)\log(E + 1),$$

$$\tilde{\eta}_0 := \tilde{\mu}/(2\sigma^2 + (c_1c_3 + c_2/6)c_4c_6),$$

$$B_1 := 2\sigma^2 + (4c_1\tilde{\eta}_0 + 4c_2\tilde{\eta}_0^2 + 1/K)c_5c_6,$$

$$B_2 := 2\sigma^2 + \left(4c_1\tilde{\eta}_0 + 4c_2\tilde{\eta}_0^2 + \frac{N-K}{K(N-1)}\right)c_5c_6.$$

Instead of proving Theorem 3.2 and 3.3 directly, we prove a more general version of convergence results suppose that some conditions about the stepsize are satisfied. Then we will show that the stepsizes given in Theorem 3.2 and 3.3 satisfy the conditions. The proof of lemmas will be deferred

Theorem C.5. Under Assumption 2.1, 2.2, 2.3, 2.5, 2.6, for a diminishing stepsize $\eta_t = \frac{\beta}{t+\gamma}$ where $\beta > \frac{1}{\tilde{\mu}}$, $\gamma > 0$ such that $\eta_0 \leq \tilde{\eta}_0$, $\eta_t \leq 2\eta_{t+E}$, and $\eta_t^2 \leq 1/(2c_3(t+1-t_0)(1+2(t+1-t_0)))$, then

1. for scheme I.

$$\mathbb{E}[\|\overline{\boldsymbol{\theta}}^{t+1} - \boldsymbol{\theta}^{PS}\|_2^2] \le \frac{\upsilon}{\gamma + t},$$

where
$$v = \max\left\{\frac{4B_1}{\tilde{\mu}^2}, \gamma \mathbb{E}[\|\overline{\boldsymbol{\theta}}^0 - \boldsymbol{\theta}^{PS}\|_2^2]\right\};$$

2. for scheme II,

$$\mathbb{E}[\|\overline{\boldsymbol{\theta}}^{t+1} - \boldsymbol{\theta}^{PS}\|_{2}^{2}] \le \frac{\upsilon}{\gamma + t}$$

where
$$\upsilon = \max\Big\{\frac{4B_2}{\tilde{\mu}^2}, \gamma \mathbb{E}[\|\overline{\boldsymbol{\theta}}^0 - \boldsymbol{\theta}^{PS}\|_2^2]\Big\}.$$

Proof. We give a proof for scheme I and the proof for scheme II follows exactly the same way.

Let $\Delta_t := \mathbb{E}[\|\overline{m{ heta}}^t - m{ heta}^{PS}\|_2^2]$, then from Lemma C.4, we have

$$\Delta_{t+1} \le (1 - \tilde{\mu}\eta_t) \mathbb{E} \|\overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}\|_2^2 + B_1 \eta_t^2.$$

For a diminishing stepsize $\eta_t = \frac{\beta}{t+\gamma}$ where $\beta > \frac{1}{\tilde{\mu}}$, $\gamma > 0$ such that $\eta_t^2 \le 1/\left(2c_3(t+1-t_0)(1+2(t+1-t_0))\right)$, $\eta_0 \le \hat{\eta}_0$, and $\eta_t \le 2\eta_{t+E}$, we will prove that $\Delta_t \le \frac{\upsilon}{\gamma+t}$ where $\upsilon = \max\left\{\frac{\beta^2 B_1}{\beta\tilde{\mu}-1}, \gamma\Delta_0\right\} = \max\left\{\frac{4B_1}{\tilde{\mu}^2}, \gamma\Delta_0\right\}$ by induction.

Firstly, $\Delta_0 \leq \frac{v}{\gamma}$ by the definition of v. Assume that for some $0 \leq t$, $\Delta_t \leq \frac{v}{\gamma + t}$, then

$$\Delta_{t+1} \leq (1 - \eta_t \tilde{\mu}) \Delta_t + \eta_t^2 B_1$$

$$\leq \left(1 - \frac{\beta \tilde{\mu}}{t + \gamma}\right) \frac{v}{t + \gamma} + \frac{\beta^2 B_1}{(t + \gamma)^2}$$

$$= \frac{t + \gamma - 1}{(t + \gamma)^2} v + \left[\frac{\beta^2 B_1}{(t + \gamma)^2} - \frac{\beta \tilde{\mu} - 1}{(t + \gamma)^2} v\right]$$

$$\leq \frac{v}{t + \gamma + 1}.$$

Specifically, if we choose $\beta=\frac{2}{\tilde{\mu}}, \gamma=\max\{\frac{2}{\tilde{\mu}\tilde{\eta}_0}, E, \frac{2}{\tilde{\mu}}\sqrt{(4E^2+10E+6)(12\sigma^2+18L^2(1+\epsilon_{\max})^2)}\}$, then we have

$$\eta_0 = \frac{\beta}{\gamma} \le \frac{2}{\tilde{\mu} \frac{2}{\tilde{\mu} \tilde{\eta}_0}} = \hat{\eta}_0$$

and

$$\eta_t - 2\eta_{t+E} = \frac{\beta}{\gamma + t} - \frac{2\beta}{\gamma + t + E} = \frac{\beta(E - \gamma - t)}{(\gamma + t)(\gamma + t + E)} \le \frac{\beta(E - \gamma)}{(\gamma + t)(\gamma + t + E)} \le 0.$$

To prove that $\eta_t^2 \leq 1/\big(2c_3(t+1-t_0)(1+2(t+1-t_0))\big)$ for any t, it suffices to prove that for $0 \leq t \leq E$ because $\{\eta_t\}$, i.e., $t_0=0$, is non-increasing and $t+1-t_0$ is periodic with period E. When $t_0=0$, we need to prove $\eta_t^2 \leq 1/\big(2c_3(t+1)(1+2(t+1))\big)$ for $0 \leq t \leq E$, which is satisfied if

$$\begin{split} \max_{0 \le t \le E} \eta_t & \le \min_{0 \le t \le E} \sqrt{1/\big(2c_3(t+1)(1+2(t+1))\big)} \\ \iff & \eta_0 \le \sqrt{\frac{1}{(4E^2+10E+6)c_3}} \\ \iff & \gamma \ge \beta \sqrt{(4E^2+10E+6)c_3} = \frac{2}{\tilde{\mu}} \sqrt{(4E^2+10E+6)c_3} = \frac{2}{\tilde{\mu}} \sqrt{(4E^2+10E+6)(12\sigma^2+18L^2(1+\epsilon_{\max})^2)}. \end{split}$$

C.4 Deferred Proofs of Key Lemmas

Proof of C.1. Let $\{x_i\}_{i=1}^N$ denote any fixed deterministic sequence. We sample a multiset \mathcal{S}_t with $|\mathcal{S}_t| = K$ by the procedure where each sampling time, we sample x_k with probability q_k for each time. Note that two samples are not necessarily independent. We only require each sampling distribution is identical. Let $\mathcal{S}_t = \{i_1, \dots, i_K\} \subset [N]$ (some i_k 's may have the same value if sampling with replacement). Then

$$\mathbb{E}_{\mathcal{S}_t} \sum_{k \in \mathcal{S}_t} x_k = \mathbb{E}_{\mathcal{S}_t} \sum_{k=1}^K x_{i_k} = K \mathbb{E}_{\mathcal{S}_t} x_{i_1} = K \sum_{k=1}^K q_k x_k.$$

For Scheme I, $q_k = p_k$ and for Scheme II, $q_k = \frac{1}{N}$, replacing the values into the above proves the lemma.

Proof of Lemma C.2. When $t+1 \notin \mathcal{I}_E$, because $\overline{\boldsymbol{\theta}}^{t+1} = \overline{\boldsymbol{w}}^{t+1}$ for both schemes, by Lemma B.1, we got the conclusion. When $t+1 \in \mathcal{I}_E$, we have

$$\mathbb{E}\|\overline{\boldsymbol{\theta}}^{t+1} - \boldsymbol{\theta}^{PS}\|_2^2 = \mathbb{E}\|\overline{\boldsymbol{\theta}}^{t+1} - \overline{\boldsymbol{w}}^{t+1}\|_2^2 + \mathbb{E}\|\overline{\boldsymbol{w}}^{t+1} - \boldsymbol{\theta}^{PS}\|_2^2 + 2\mathbb{E}\langle\overline{\boldsymbol{\theta}}^{t+1} - \overline{\boldsymbol{w}}^{t+1}, \overline{\boldsymbol{w}}^{t+1} - \boldsymbol{\theta}^{PS}\rangle.$$

By Lemma C.1 and the law of total expectation, we have

$$\mathbb{E}\langle \overline{\boldsymbol{\theta}}^{t+1} - \overline{\boldsymbol{w}}^{t+1}, \overline{\boldsymbol{w}}^{t+1} - \boldsymbol{\theta}^{PS} \rangle = \mathbb{E}\big[\mathbb{E}_{\mathcal{S}_{t+1}}\langle \overline{\boldsymbol{\theta}}^{t+1} - \overline{\boldsymbol{w}}^{t+1}, \overline{\boldsymbol{w}}^{t+1} - \boldsymbol{\theta}^{PS} \rangle\big] = 0.$$

Next we focus on upper bounding $\mathbb{E}[\|\overline{m{ heta}}^{t+1}-\overline{m{w}}^{t+1}\|_2^2]$ under two sampling schemes.

Denote $S_{t+1} = \{i_1, \dots, i_K\}$, then for scheme $I, \overline{\boldsymbol{\theta}}^{t+1} = \frac{1}{K} \sum_{l=1}^K \boldsymbol{w}_{i_l}^{t+1}$. Thus by the law of total expectation, we have

$$\begin{split} \mathbb{E} \| \overline{\boldsymbol{\theta}}^{t+1} - \overline{\boldsymbol{w}}^{t+1} \|_{2}^{2} &= \mathbb{E} \big[\mathbb{E}_{\mathcal{S}_{t+1}} \| \overline{\boldsymbol{\theta}}^{t+1} - \overline{\boldsymbol{w}}^{t+1} \|_{2}^{2} \big] \\ &= \mathbb{E} \big[\mathbb{E}_{\mathcal{S}_{t+1}} \| \frac{1}{K} \sum_{l=1}^{K} \boldsymbol{w}_{i_{l}}^{t+1} - \overline{\boldsymbol{w}}^{t+1} \|_{2}^{2} \big] \\ &\leq \mathbb{E} \big[\mathbb{E}_{\mathcal{S}_{t+1}} \frac{1}{K^{2}} \sum_{l=1}^{K} \| \boldsymbol{w}_{i_{l}}^{t+1} - \overline{\boldsymbol{w}}^{t+1} \|_{2}^{2} \big] \\ &= \frac{1}{K} \sum_{k=1}^{N} p_{k} \mathbb{E} \| \boldsymbol{w}_{k}^{t+1} - \overline{\boldsymbol{w}}^{t+1} \|_{2}^{2}. \end{split}$$

Again with $\overline{\theta}^{t+1} = \frac{1}{K} \sum_{l=1}^{K} w_{i_l}^{t+1}$, for scheme II, by the law of total expectation, we have

$$\begin{split} \mathbb{E}\|\overline{\theta}^{t+1} - \overline{w}^{t+1}\|_{2}^{2} &= \mathbb{E}\left[\mathbb{E}_{\mathcal{S}_{t+1}}\|\overline{\theta}^{t+1} - \overline{w}^{t+1}\|_{2}^{2}\right] \\ &= \mathbb{E}\left[\mathbb{E}_{\mathcal{S}_{t+1}}\left[\|\frac{1}{K}\sum_{l=1}^{K} \boldsymbol{w}_{i_{l}}^{t+1} - \overline{w}^{t+1}\|_{2}^{2}\right]\right] \\ &= \mathbb{E}\left[\frac{1}{K^{2}}\mathbb{E}_{\mathcal{S}_{t+1}}\left[\|\sum_{i=1}^{N} \mathbf{1}\{i \in \mathcal{S}_{t+1}\}(\boldsymbol{w}_{i}^{t+1} - \overline{\boldsymbol{w}}^{t+1})\|_{2}^{2}\right]\right] \\ &\leq \frac{1}{K^{2}}\mathbb{E}\left[\sum_{i=1}^{N} \mathbb{P}(i \in \mathcal{S}_{t+1})\|\boldsymbol{w}_{i}^{t+1} - \overline{\boldsymbol{w}}^{t+1}\|_{2}^{2} + \sum_{j \neq i} \mathbb{P}(i, j \in \mathcal{S}_{t+1})\left\langle\boldsymbol{w}_{i}^{t+1} - \overline{\boldsymbol{w}}^{t+1}, \; \boldsymbol{w}_{j}^{t+1} - \overline{\boldsymbol{w}}^{t+1}\right\rangle\right] \\ &= \frac{1}{KN}\sum_{i=1}^{N} \mathbb{E}\|\boldsymbol{w}_{i}^{t+1} - \overline{\boldsymbol{w}}^{t+1}\|_{2}^{2} + \frac{K-1}{KN(N-1)}\sum_{i \neq j} \mathbb{E}\left\langle\boldsymbol{w}_{i}^{t+1} - \overline{\boldsymbol{w}}^{t+1}, \; \boldsymbol{w}_{j}^{t+1} - \overline{\boldsymbol{w}}^{t+1}\right\rangle \\ &= \frac{N}{K(N-1)}\left(1 - \frac{K}{N}\right)\sum_{i=1}^{N} p_{i}\mathbb{E}\|\boldsymbol{w}_{i}^{t+1} - \overline{\boldsymbol{w}}^{t+1}\|_{2}^{2} \end{split}$$

where we use the following equalities: (1) $\mathbb{P}(i \in S_{t+1}) = \frac{K}{N}$ and $\mathbb{P}(i, j \in S_{t+1}) = \frac{K(K-1)}{N(N-1)}$ for all $i \neq j$ and (2) $\sum_{i=1}^{N} \|\boldsymbol{w}_{i}^{t} - \overline{\boldsymbol{w}}^{t}\|^{2} + \sum_{i \neq j} \langle \boldsymbol{w}_{i}^{t+1} - \overline{\boldsymbol{w}}^{t+1}, \boldsymbol{w}_{j}^{t+1} - \overline{\boldsymbol{w}}^{t+1} \rangle = 0$.

The conclusion follows from the above discussion.

Proof of Lemma C.4. In this proof, for convenience, we will discuss with respect to t+1 where we assume $t+1 \notin \mathcal{I}_E$ and transfer back to t in the last. First by the update rule, we have when $t+1 \notin \mathcal{I}_E$

$$\boldsymbol{\theta}_i^{t+1} - \overline{\boldsymbol{\theta}}^{t+1} = \boldsymbol{\theta}_i^t - \overline{\boldsymbol{\theta}}^t - \eta_t(\nabla \ell(\boldsymbol{\theta}_i^t; Z_i^{t+1}) - g_t)$$

and when $t+1 \in \mathcal{I}_E$,

$$\boldsymbol{w}_{i}^{t+1} - \overline{\boldsymbol{w}}^{t+1} = \boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t} - \eta_{t}(\nabla \ell(\boldsymbol{\theta}_{i}^{t}; Z_{i}^{t+1}) - g_{t}).$$

Then with the same method in Lemma B.2, let $t_0 := \max\{s \mid s < t+1, s \in \mathcal{I}_E\}$ and $c_3 := 12\sigma^2 + 18L^2(1+\epsilon_{\max})^2$, if $\eta_t^2 \le \frac{1}{2c_3(t+1-t_0)\left(1+2(t+1-t_0)\right)}$, we will have: if $t+1 \notin \mathcal{I}_E$,

$$\sum_{i=1}^{N} p_{i} \mathbb{E} \|\boldsymbol{\theta}_{i}^{t+1} - \overline{\boldsymbol{\theta}}^{t+1}\|_{2}^{2}$$

$$\leq \sum_{s=0}^{t-t_{0}} \frac{t+2-t_{0}}{s+2} \eta_{t}^{2} (2E+1) (48\sigma^{2}+36\varsigma^{2}) \mathbb{E} \|\overline{\boldsymbol{\theta}}^{s} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + \sum_{s=0}^{t-t_{0}} \frac{t+2-t_{0}}{s+2} \eta_{t}^{2} (2E+1) (24\sigma^{2}+36\varsigma^{2})$$

and if $t+1 \in \mathcal{I}_E$,

$$\sum_{i=1}^{N} p_{i} \mathbb{E} \| \boldsymbol{w}_{i}^{t+1} - \overline{\boldsymbol{w}}^{t+1} \|_{2}^{2} \\
\leq \sum_{s=0}^{t-t_{0}} \frac{t+2-t_{0}}{s+2} \eta_{t}^{2} (2E+1) (48\sigma^{2}+36\varsigma^{2}) \mathbb{E} \| \overline{\boldsymbol{\theta}}^{s} - \boldsymbol{\theta}^{PS} \|_{2}^{2} + \sum_{s=0}^{t-t_{0}} \frac{t+2-t_{0}}{s+2} \eta_{t}^{2} (2E+1) (24\sigma^{2}+36\varsigma^{2}).$$

With the above formula and Lemma C.2, we now prove that if η_0 is sufficiently small, then for any t, we have $\mathbb{E}\|\overline{\boldsymbol{\theta}}^t - {\boldsymbol{\theta}}^{PS}\|_2^2 \leq \mathbb{E}\|\overline{\boldsymbol{\theta}}^0 - {\boldsymbol{\theta}}^{PS}\|_2^2$. We first derive the following inequality, which we will use later. Note that for any t where $t_0 := \max\{s \mid s < t+1, s \in \mathcal{I}_E\}$, we have

$$\sum_{s=0}^{t-t_0} \frac{t+2-t_0}{s+2} = (t+2-t_0)(\frac{1}{2}+\ldots+\frac{1}{t-t_0+2}) \le (t+2-t_0)\log(t+2-t_0) \le (E+1)\log(E+1).$$

Then again by the same induction method in Lemma B.2, we have if

$$\eta_t \leq \eta_0 \\ \leq \frac{\tilde{\mu} \mathbb{E} \|\overline{\boldsymbol{\theta}}^0 - \boldsymbol{\theta}^{PS}\|_2^2}{2\sigma^2 + (c_1c_3 + c_2/6)(2E^2 + 3E + 1)\log(E + 1)\left((16\sigma^2 + 12\varsigma^2)\mathbb{E} \|\overline{\boldsymbol{\theta}}^0 - \boldsymbol{\theta}^{PS}\|_2^2 + (8\sigma^2 + 12\varsigma^2)\right)} = \tilde{\eta}_0,$$

then for any t, we have $\mathbb{E}\|\overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}\|_2^2 \leq \mathbb{E}\|\overline{\boldsymbol{\theta}}^0 - \boldsymbol{\theta}^{PS}\|_2^2$.

Under all these conditions, if $t \notin \mathcal{I}_E$, we have

$$\begin{split} \sum_{i=1}^{N} p_{i} \mathbb{E} \|\boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t}\|_{2}^{2} &\leq \eta_{t-1}^{2} (2E^{2} + 3E + 1) \log(E + 1) (48\sigma^{2} + 36\varsigma^{2}) \mathbb{E} \|\overline{\boldsymbol{\theta}}^{0} - \boldsymbol{\theta}^{PS}\|_{2}^{2} \\ &+ \eta_{t-1}^{2} (2E^{2} + 3E + 1) \log(E + 1) (24\sigma^{2} + 36\varsigma^{2}) \\ &\leq 4\eta_{t}^{2} (2E^{2} + 3E + 1) \log(E + 1) (48\sigma^{2} + 36\varsigma^{2}) \mathbb{E} \|\overline{\boldsymbol{\theta}}^{0} - \boldsymbol{\theta}^{PS}\|_{2}^{2} \\ &+ 4\eta_{t}^{2} (2E^{2} + 3E + 1) \log(E + 1) (24\sigma^{2} + 36\varsigma^{2}), \end{split}$$

and if $t + 1 \in \mathcal{I}_E$, we have

$$\sum_{i=1}^{N} p_{i} \mathbb{E} \| \boldsymbol{w}_{i}^{t+1} - \overline{\boldsymbol{w}}^{t+1} \|_{2}^{2} \leq \eta_{t-1}^{2} (2E^{2} + 3E + 1) \log(E+1) (48\sigma^{2} + 36\varsigma^{2}) \mathbb{E} \| \overline{\boldsymbol{\theta}}^{0} - \boldsymbol{\theta}^{PS} \|_{2}^{2} + \eta_{t-1}^{2} (2E^{2} + 3E + 1) \log(E+1) (24\sigma^{2} + 36\varsigma^{2}).$$

Note that in Lemma C.2, the inequality for $t \notin \mathcal{I}_E$ is looser than the inequality for $t \in \mathcal{I}_E$. Therefore, we can apply the inequality for $t \in \mathcal{I}_E$ for all t. Combining this inequality with the above formula gives us that:

1. for scheme I,

$$\mathbb{E}[\|\overline{\boldsymbol{\theta}}^{t+1} - \boldsymbol{\theta}^{PS}\|_{2}^{2}]$$

$$\leq \frac{1}{K} \sum_{k=1}^{N} p_{k} \mathbb{E} \|\boldsymbol{w}_{k}^{t+1} - \overline{\boldsymbol{w}}^{t+1}\|_{2}^{2} + (1 - \tilde{\mu}\eta_{t}) \mathbb{E} \|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + 2\sigma^{2}\eta_{t}^{2} + (c_{1}\eta_{t} + c_{2}\eta_{t}^{2}) \sum_{i=1}^{N} p_{i} \mathbb{E} \|\boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t}\|_{2}^{2}$$

$$\leq (1 - \tilde{\mu}\eta_{t}) \mathbb{E} \|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + 2\sigma^{2}\eta_{t}^{2}$$

$$+ (4c_{1}\eta_{t} + 4c_{2}\eta_{t}^{2} + K^{-1})(2E^{2} + 3E + 1) \log(E + 1)\eta_{t}^{2} \left((48\sigma^{2} + 36\varsigma^{2}) \mathbb{E} \|\overline{\boldsymbol{\theta}}^{0} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + (24\sigma^{2} + 36\varsigma^{2}) \right),$$

2. for scheme II,

$$\mathbb{E}[\|\overline{\boldsymbol{\theta}}^{t+1} - \boldsymbol{\theta}^{PS}\|_{2}^{2}] \\
\leq \frac{N - K}{KN(N-1)} \sum_{k=1}^{N} p_{k} \mathbb{E}\|\boldsymbol{w}_{k}^{t+1} - \overline{\boldsymbol{w}}^{t+1}\|_{2}^{2} + (1 - \tilde{\mu}\eta_{t}) \mathbb{E}\|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + 2\sigma^{2}\eta_{t}^{2} + (c_{1}\eta_{t} + c_{2}\eta_{t}^{2}) \sum_{i=1}^{N} p_{i} \mathbb{E}\|\boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t}\|_{2}^{2} \\
\leq (1 - \tilde{\mu}\eta_{t}) \mathbb{E}\|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + 2\sigma^{2}\eta_{t}^{2} \\
+ \left(4c_{1}\eta_{t} + 4c_{2}\eta_{t}^{2} + \frac{N - K}{K(N-1)}\right) (2E^{2} + 3E + 1) \log(E + 1)\eta_{t}^{2} \left((48\sigma^{2} + 36\varsigma^{2})\mathbb{E}\|\overline{\boldsymbol{\theta}}^{0} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + (24\sigma^{2} + 36\varsigma^{2})\right).$$

Proof of Lemma C.4. The conclusion follows directly from Lemma C.3.

D Proof of convergence under the alternative assumption in (3)

Assumption D.1. Suppose the following hold

$$\mathbb{E}_{Z_i \sim \mathcal{D}_i(\boldsymbol{\theta})}[\|\nabla \ell(\boldsymbol{\theta}; Z_i)\|_2^2] \le G^2. \tag{16}$$

Lemma D.2. (Bound on the divergence of parameters, i.e., consensus error bound)

When E > 1, under Assumption 2.1, 2.2, 2.3, 2.5, and if η_t is non-increasing and $\eta_t \leq 2\eta_{t+E}$ holds for all $t \geq 0$, we have

$$\mathbb{E}\left[\sum_{i=1}^{N} p_{i} \|\boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t}\|_{2}^{2}\right] \leq 4(E-1)^{2} \eta_{t}^{2} G^{2}$$
(17)

Proof. FedAvg requires a communication every E steps, so for any $t \geq 0$, there exists a $t_0 \leq t, t_0 \in \mathcal{I}_E$, such that $t - t_0 \leq E - 1$ and $\boldsymbol{\theta}_i^{t_0} = \overline{\boldsymbol{\theta}}^{t_0}$, $\forall i$. Also, we use the fact that $\eta_{t_0} \leq 2\eta_t$ for all $t - t_0 \leq E - 1$, then

$$\mathbb{E}\left[\sum_{i=1}^{N} p_i \|\boldsymbol{\theta}_i^t - \overline{\boldsymbol{\theta}}^t\|_2^2\right] = \mathbb{E}\left[\sum_{i=1}^{N} p_i \|(\boldsymbol{\theta}_i^t - \overline{\boldsymbol{\theta}}^{t_0}) - (\overline{\boldsymbol{\theta}}^t - \overline{\boldsymbol{\theta}}^{t_0})\|_2^2\right] \le \mathbb{E}\left[\sum_{i=1}^{N} p_i \|\boldsymbol{\theta}_i^t - \overline{\boldsymbol{\theta}}^{t_0}\|_2^2\right],\tag{18}$$

since $\mathbb{E}\|X - \mathbb{E}X\|_2^2 \leq \mathbb{E}\|X\|_2^2$ where $X = \boldsymbol{\theta}_i^t - \overline{\boldsymbol{\theta}}^{t_0}$. Using Jensen's inequality, we further have

$$\|\boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t_{0}}\|_{2}^{2} = \left\| \sum_{s=t_{0}}^{t-1} \eta_{s} \nabla \ell(\boldsymbol{\theta}_{i}^{s}; Z_{i}^{s+1}) \right\|_{2}^{2} \le (t - t_{0}) \sum_{s=t_{0}}^{t-1} \eta_{s}^{2} \|\nabla \ell(\boldsymbol{\theta}_{i}^{s}; Z_{i}^{s+1})\|_{2}^{2}, \tag{19}$$

$$\mathbb{E}\left[\sum_{i=1}^{N} p_{i} \|\boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t_{0}}\|_{2}^{2}\right] = \mathbb{E}\left[\left\|\sum_{i=1}^{N} p_{i} \sum_{s=t_{0}}^{t-1} \eta_{s}^{2} \nabla \ell(\boldsymbol{\theta}_{i}^{s}; Z_{i}^{s+1})\right\|_{2}^{2}\right] \leq (t - t_{0}) \sum_{s=t_{0}}^{t-1} \eta_{s}^{2} \sum_{i=1}^{N} p_{i} \mathbb{E}\left[\left\|\nabla \ell(\boldsymbol{\theta}_{i}^{s}; Z_{i}^{s+1})\right\|_{2}^{2}\right], \tag{20}$$

]

$$\mathbb{E}\left[\sum_{i=1}^{N} p_{i} \|\boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t_{0}}\|_{2}\right] = \mathbb{E}\left[\left\|\sum_{i=1}^{N} p_{i} \sum_{s=t_{0}}^{t-1} \eta_{s} \nabla \ell(\boldsymbol{\theta}_{i}^{s}; Z_{i}^{s+1})\right\|_{2}\right] \leq \sum_{s=t_{0}}^{t-1} \eta_{s} \sum_{i=1}^{N} p_{i} \mathbb{E}\left[\left\|\nabla \ell(\boldsymbol{\theta}_{i}^{s}; Z_{i}^{s+1})\right\|_{2}\right], \quad (21)$$

where we used $\eta_s \leq \eta_{t_0}$. Therefore, based on **A5**, we have

$$\mathbb{E}\left[\sum_{i=1}^{N} p_{i} \|\boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t}\|_{2}^{2}\right] \leq \sum_{i=1}^{N} p_{i} \mathbb{E}\left[\sum_{s=t_{0}}^{t-1} (E-1)\eta_{s}^{2} \|\nabla \ell(\boldsymbol{\theta}_{i}^{s}; Z_{i}^{s+1})\|_{2}^{2}\right] \\
\leq \sum_{i=1}^{N} p_{i} \left[\sum_{s=t_{0}}^{t-1} (E-1)\eta_{s}^{2} G^{2}\right] \\
\leq \sum_{i=1}^{N} p_{i} (E-1)^{2} \eta_{t_{0}}^{2} G^{2} \\
\leq 4(E-1)^{2} \eta_{t}^{2} G^{2} \tag{22}$$

since $\eta_s \leq \eta_{t_0} \leq 2\eta_{t_0+E} \leq 2\eta_t$ in the last two inequalities.

Lemma D.3. Li et al. [2022] Consider a sequence of non-negative, non-increasing step sizes $\{\eta_t\}_{t\geq 1}$. Let $a>0, p\in \mathbb{Z}_+$ and $\eta_1<2/a$. If $\eta_t^p/\eta_{t+1}^p\leq 1+(a/2)\eta_{t+1}^p$ for any $t\geq 1$, then

$$\sum_{j=1}^{t} \eta_j^{p+1} \prod_{\ell=j+1}^{t} (1 - \eta_{\ell} a) \le \frac{2}{a} \eta_t^p, \quad \forall t \ge 1$$
 (23)

Lemma D.4. Under Assumptions 2.1, 2.2, 2.3, D.1 and the condition that $\eta_t \leq \tilde{\mu}/c_2$, , $\eta_t \leq \eta_{t_0} \leq 2\eta_t$ where $t_0 = \max_s \{s \in \mathbb{N} | Es \leq t\}$, $\eta_{t+1} < \eta_t$ for any $t \geq 0$, $\eta_1 < \frac{2}{\tilde{\mu}}$ and $\eta_t^q/\eta_{t+1}^q \leq 1 + (\tilde{\mu}/2)\eta_{t+1}^q$ for any $t \geq 0$ and q = 1, 2, 3.

$$\mathbb{E}[\|\overline{\boldsymbol{w}}^{t+1} - \boldsymbol{\theta}^{PS}\|_{2}^{2}] \leq \prod_{i=0}^{t} (1 - \tilde{\mu}\eta_{i}) \|\overline{\boldsymbol{\theta}}^{0} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + \frac{2c_{2}c_{7}}{\tilde{\mu}}\eta_{t}^{3} + \frac{2c_{1}c_{7}}{\tilde{\mu}}\eta_{t}^{2} + \frac{4\sigma^{2}}{\tilde{\mu}}\eta_{t}, \tag{24}$$

where $c_7 := 4(E-1)G^2$.

Proof. From Lemma B.1, we have

$$\mathbb{E}\left[\|\overline{\boldsymbol{w}}^{t+1} - \boldsymbol{\theta}^{PS}\|_{2}^{2}\right] \leq (1 - \tilde{\mu}\eta_{t})\mathbb{E}\left[\|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2}\right] + (c_{1}\eta_{t} + c_{2}\eta_{t}^{2})\mathbb{E}\left[\sum_{i=1}^{N} p_{i}\|\boldsymbol{\theta}_{i}^{t} - \overline{\boldsymbol{\theta}}^{t}\|_{2}^{2}\right] + 2\sigma^{2}\eta_{t}^{2} \\
\leq (1 - \tilde{\mu}\eta_{t})\mathbb{E}\left[\|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\|_{2}^{2}\right] + c_{2}c_{7}\eta_{t}^{4} + c_{1}c_{7}\eta_{t}^{3} + 2\sigma^{2}\eta_{t}^{2} \\
= \prod_{i=0}^{t} (1 - \tilde{\mu}\eta_{i})\|\overline{\boldsymbol{\theta}}^{0} - \boldsymbol{\theta}^{PS}\|_{2}^{2} + \sum_{s=1}^{t} \prod_{i=s+1}^{t} (1 - \tilde{\mu}\eta_{i})\left(c_{2}c_{7}\eta_{s}^{4} + c_{1}c_{7}\eta_{s}^{3} + 2\sigma^{2}\eta_{s}^{2}\right). \tag{25}$$

The second inequality holds because of Lemma D.2. Using Lemma D.3,

$$\sum_{s=1}^{t} \prod_{i=s+1}^{t} (1 - \tilde{\mu}\eta_i) \left(c_2 c_7 \eta_s^4 + c_1 c_7 \eta_s^3 + 2\sigma^2 \eta_s^2 \right) \le \frac{2c_2 c_7}{\tilde{\mu}} \eta_t^3 + \frac{2c_1 c_7}{\tilde{\mu}} \eta_t^2 + \frac{4\sigma^2}{\tilde{\mu}} \eta_t. \tag{26}$$

Theorem D.5. (Full participation convergence theorem, alternative assumption)

Under Assumption 2.1, 2.2, 2.3, 2.5, the full participation scheme has convergence rate $\mathcal{O}(\frac{1}{T})$, i.e., denote $\Delta_t := \mathbb{E}[\|\overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}\|_2^2]$, then for some $\gamma > 0$,

$$\Delta_t \le \frac{v}{\gamma + t},\tag{27}$$

where $v := \max\{\frac{c_2 c_7 \beta^4 \gamma^{-2} + c_1 c_7 \beta^3 \gamma^{-1} + 2 \sigma^2 \beta^2}{\beta \mu - 1}, (\gamma + 1) \Delta_1\}.$

Proof. We will show it on the partial participation algorithm, and the proof for the full participation is similar.

For a diminishing step size $\eta_t = \frac{\beta}{t+\gamma}$ for some $\beta > \frac{1}{\tilde{\mu}}$ and $\gamma > 0$ such that $\eta_1 \leq \min\{\frac{1}{\tilde{\mu}}, \frac{1}{4L}\} = \frac{1}{4L}$ and $\eta_t \leq 2\eta_{t+E}$. We will prove $\Delta_t := \mathbb{E}[\|\overline{\boldsymbol{\theta}}^t - {\boldsymbol{\theta}}^{PS}\|_2^2] \leq \frac{v}{\gamma+t}$, where $v := \max\{\frac{c_2c_7\beta^4\gamma^{-2} + c_1c_7\beta^3\gamma^{-1} + 2\sigma^2\beta^2}{\beta\mu - 1}, (\gamma+1)\Delta_1\}$. We prove this by induction. Firstly, the definition if v ensures it holds for t=1. Assume it holds for some t, i.e., $\eta_t = \frac{\beta}{t+\gamma}$, then it follows from Lemma D.4 that

$$\Delta_{t+1} \leq (1 - \eta_t \tilde{\mu}) \Delta_t + c_2 c_7 \eta_t^4 + c_1 c_7 \eta_t^3 + 2\sigma^2 \eta_t^2
\leq (1 - \frac{\beta \tilde{\mu}}{t + \gamma}) \frac{v}{t + \gamma} + \frac{c_2 c_7 \beta^4}{(t + \gamma)^4} + \frac{c_1 c_7 \beta^3}{(t + \gamma)^3} + \frac{2\sigma^2 \beta^2}{(t + \gamma)^2}
= \frac{t + \gamma - 1}{(t + \gamma)^2} v + \left[\frac{c_2 c_7 \beta^4}{(t + \gamma)^4} + \frac{c_1 c_7 \beta^3}{(t + \gamma)^3} + \frac{2\sigma^2 \beta^2}{(t + \gamma)^2} - \frac{\beta \mu - 1}{(t + \gamma)^2} v \right]
\leq \frac{t + \gamma - 1}{(t + \gamma)^2} v + \left[\frac{c_2 c_7 \beta^4}{(t + \gamma)^2 \gamma^2} + \frac{c_1 c_7 \beta^3}{(t + \gamma)^2 \gamma} + \frac{2\sigma^2 \beta^2}{(t + \gamma)^2} - \frac{\beta \mu - 1}{(t + \gamma)^2} v \right]
\leq \frac{v}{t + \gamma + 1}$$
(28)

where $\tilde{\mu}$, c_1 , c_2 , c_3 , c_7 are defined the same as in earlier proofs, and thus the $\mathcal{O}(1/T)$ convergence rate is shown. \Box

Lemma D.6. (Bounding the difference $\overline{\boldsymbol{w}}^{t+1} - \overline{\boldsymbol{\theta}}^{t+1}$ in partial participation) Suppose Assumption 2.1, 2.2, 2.3, and 2.5 hold. For $t+1 \in \mathcal{I}_E$, assume that η_t is non-increasing and $\eta_t \leq 2\eta_{t+E}$ for all t, then we have the following results

1. For Scheme I, the expected difference $\overline{m{w}}^{t+1} - \overline{m{ heta}}^{t+1}$ is bounded by

$$\mathbb{E}_{\mathcal{S}_t} \|\overline{\boldsymbol{\theta}}^{t+1} - \overline{\boldsymbol{w}}^{t+1}\|_2^2 \le \frac{4}{K} \eta_t^2 E^2 G^2.$$
 (29)

2. For Scheme II, assuming $p_1=p_2=\cdots=p_N=\frac{1}{N}$, the expected difference $\overline{{\pmb w}}^{t+1}-\overline{{\pmb \theta}}^{t+1}$ is bounded by

$$\mathbb{E}_{\mathcal{S}_t} \|\overline{\boldsymbol{\theta}}^{t+1} - \overline{\boldsymbol{w}}^{t+1}\|_2^2 \le \frac{4(N-K)}{K(N-1)} \eta_t^2 E^2 G^2$$
(30)

Proof. We prove the bound for Scheme I as follows. Since $\overline{\theta}^{t+1} = \frac{1}{K} \sum_{l=1}^{K} w_{i_l}^t$, taking expectation over \mathcal{S}_{t+1} , we have

$$\mathbb{E}_{\mathcal{S}_{t}} \|\overline{\boldsymbol{\theta}}^{t+1} - \overline{\boldsymbol{w}}^{t+1}\|_{2}^{2} = \mathbb{E}_{\mathcal{S}_{t}} \frac{1}{K^{2}} \sum_{l=1}^{K} \|\boldsymbol{w}_{i_{l}}^{t+1} - \overline{\boldsymbol{w}}^{t+1}\|_{2}^{2} = \frac{1}{K} \sum_{k=1}^{N} p_{k} \|\boldsymbol{w}_{k}^{t+1} - \overline{\boldsymbol{w}}^{t+1}\|_{2}^{2}$$
(31)

We note that since $t+1 \in \mathcal{I}_E$, we know that the time $t_0 = t - E + 1 \in \mathcal{I}_E$ is the communication time, which implies $\{\theta_{t_0}^k\}$ is identical. Then

$$\sum_{k=1}^{N} p_{k} \|\boldsymbol{w}_{k}^{t+1} - \overline{\boldsymbol{w}}^{t+1}\|_{2}^{2} = \sum_{i=1}^{N} p_{k} \|(\boldsymbol{w}_{k}^{t+1} - \overline{\boldsymbol{\theta}}^{t_{0}}) - (\overline{\boldsymbol{w}}^{t+1} - \overline{\boldsymbol{\theta}}^{t_{0}})\|_{2}^{2} \leq \sum_{i=1}^{N} p_{k} \|\boldsymbol{w}_{k}^{t+1} - \overline{\boldsymbol{\theta}}^{t_{0}}\|_{2}^{2}$$
(32)

Similar to Lemma D.2, the last inequality is due to $\mathbb{E}\|X - \mathbb{E}X\|_2^2 \leq \mathbb{E}\|X\|_2^2$ where $X = \boldsymbol{w}_k^{t+1} - \overline{\boldsymbol{\theta}}^{t_0}$, and $\sum_{k=1}^N p_k(\boldsymbol{w}_k^{t+1} - \overline{\boldsymbol{\theta}}^{t_0}) = \overline{\boldsymbol{w}}^{t+1} - \overline{\boldsymbol{\theta}}^{t_0}$. Similarly, we have

$$\mathbb{E}_{\mathcal{S}_{t}} \left[\| \overline{\boldsymbol{\theta}}^{t+1} - \overline{\boldsymbol{w}}^{t+1} \|_{2}^{2} \right] \leq \frac{1}{K} \sum_{k=1}^{N} p_{k} \mathbb{E} \left[\| \boldsymbol{w}_{k}^{t+1} - \overline{\boldsymbol{\theta}}^{t_{0}} \|_{2}^{2} \right] \\
= \frac{1}{K} \sum_{k=1}^{N} p_{k} \mathbb{E} \left[\| \boldsymbol{w}_{k}^{t+1} - \boldsymbol{\theta}_{k}^{t_{0}} \|_{2}^{2} \right] \\
= \frac{1}{K} \sum_{k=1}^{N} p_{k} \mathbb{E} \left[\| \sum_{s=t_{0}}^{t} \eta_{s} \nabla l(\boldsymbol{\theta}_{k}^{s}; Z_{k}^{s+1}) \|_{2}^{2} \right] \\
\leq \frac{1}{K} \sum_{k=1}^{N} p_{k} E \sum_{s=t_{0}}^{t} \mathbb{E} \left[\| \eta_{s} \nabla l(\boldsymbol{\theta}_{k}^{s}; Z_{k}^{s+1}) \|_{2}^{2} \right] \\
\leq \frac{1}{K} E^{2} \eta_{t_{0}}^{2} G^{2} \leq \frac{4}{K} \eta_{t}^{2} E^{2} G^{2}. \tag{33}$$

Then we prove the bound for Scheme II. Since $\overline{m{ heta}}^{t+1} = rac{1}{K} \sum_{l=1}^K m{w}_{i_l}^{t+1}$, we have

$$\mathbb{E}_{\mathcal{S}_{t}} \left[\| \overline{\boldsymbol{\theta}}^{t+1} - \overline{\boldsymbol{w}}^{t+1} \|_{2}^{2} \right] \\
= \mathbb{E}_{\mathcal{S}_{t}} \left[\| \frac{1}{K} \sum_{l=1}^{K} \boldsymbol{w}_{i_{l}}^{t+1} - \overline{\boldsymbol{w}}^{t+1} \|_{2}^{2} \right] \\
= \frac{1}{K^{2}} \mathbb{E}_{\mathcal{S}_{t}} \left[\| \sum_{i=1}^{N} \mathbf{1} \{ i \in \mathcal{S}_{t} \} (\boldsymbol{w}_{i}^{t+1} - \overline{\boldsymbol{w}}^{t+1}) \|_{2}^{2} \right] \\
= \frac{1}{K^{2}} \left[\sum_{i=1}^{N} \mathbb{P}(i \in \mathcal{S}_{t}) \| \boldsymbol{w}_{i}^{t+1} - \overline{\boldsymbol{w}}^{t+1} \|_{2}^{2} + \sum_{j \neq i} \mathbb{P}(i, j \in \mathcal{S}_{t}) \left\langle \boldsymbol{w}_{i}^{t+1} - \overline{\boldsymbol{w}}^{t+1}, \ \boldsymbol{w}_{j}^{t+1} - \overline{\boldsymbol{w}}^{t+1} \right\rangle \right] \\
= \frac{1}{KN} \sum_{i=1}^{N} \| \boldsymbol{w}_{i}^{t+1} - \overline{\boldsymbol{w}}^{t+1} \|_{2}^{2} + \frac{K-1}{KN(N-1)} \sum_{i \neq j} \left\langle \boldsymbol{w}_{i}^{t+1} - \overline{\boldsymbol{w}}^{t+1}, \ \boldsymbol{w}_{j}^{t+1} - \overline{\boldsymbol{w}}^{t+1} \right\rangle \\
= \frac{1}{K(N-1)} (1 - \frac{K}{N}) \sum_{i=1}^{N} \| \boldsymbol{w}_{i}^{t+1} - \overline{\boldsymbol{w}}^{t+1} \|_{2}^{2}. \tag{34}$$

Note that the second last equality holds because $\mathbb{P}(i \in \mathcal{S}_t) = \frac{K}{N}$ and $\mathbb{P}(i, j \in \mathcal{S}_t) = \frac{K(K-1)}{N(N-1)}$; and the last equality holds because

$$\begin{split} &\sum_{i=1}^{N} \|\boldsymbol{w}_{i}^{t+1} - \overline{\boldsymbol{w}}^{t+1}\|_{2}^{2} + \sum_{i \neq j} \left\langle \boldsymbol{w}_{i}^{t+1} - \overline{\boldsymbol{w}}^{t+1}, \; \boldsymbol{w}_{j}^{t+1} - \overline{\boldsymbol{w}}^{t+1} \right\rangle \\ &= \sum_{i=1}^{N} \left\langle \boldsymbol{w}_{i}^{t+1} - \overline{\boldsymbol{w}}^{t+1}, \; \left(\sum_{i=1}^{N} \boldsymbol{w}_{j}^{t+1}\right) - N \overline{\boldsymbol{w}}^{t+1} \right\rangle = 0. \end{split}$$

Recall that

$$\sum_{k=1}^{N} p_k \| \boldsymbol{w}_k^{t+1} - \overline{\boldsymbol{w}}^{t+1} \|_2^2 \le \sum_{i=1}^{N} p_k \| \boldsymbol{w}_k^{t+1} - \overline{\boldsymbol{\theta}}^{t_0} \|_2^2,$$

we get

$$\mathbb{E}\left[\|\overline{\boldsymbol{\theta}}^{t+1} - \overline{\boldsymbol{w}}^{t+1}\|_{2}^{2}\right] = \frac{1}{K(N-1)}(1 - \frac{K}{N})\mathbb{E}\left[\sum_{i=1}^{N} \|\boldsymbol{w}_{i}^{t+1} - \overline{\boldsymbol{w}}^{t+1}\|_{2}^{2}\right] \\
\leq \frac{N}{K(N-1)}(1 - \frac{K}{N})\mathbb{E}\left[\sum_{i=1}^{N} \frac{1}{N} \|\boldsymbol{w}_{i}^{t+1} - \overline{\boldsymbol{\theta}}^{t_{0}}\|_{2}^{2}\right] \\
\leq \frac{N}{K(N-1)}(1 - \frac{K}{N})4\eta_{t}^{2}E^{2}G^{2} = \frac{4(N-K)}{K(N-1)}\eta_{t}^{2}E^{2}G^{2} \tag{35}$$

where the last inequality can be found in (22) in the proof of Lemma D.2.

Lemma D.7. Under Under Assumption 2.1, 2.2, 2.3, 2.5, and the condition that $\eta_t \leq \tilde{\mu}/c_2$, $\eta_t \leq \eta_{t_0} \leq 2\eta_t$ where $t_0 = \max_s \{s \in \mathbb{N} | Es \leq t\}$, $\eta_{t+1} < \eta_t$ for any $t \geq 0$, $\eta_1 < \frac{2}{\tilde{\mu}}$ and $\eta_t^q/\eta_{t+1}^q \leq 1 + (\tilde{\mu}/2)\eta_{t+1}^q$ for any $t \geq 0$ and q = 1, 2, 3, we have

$$\mathbb{E}_{\mathcal{S}_{t}}[||\overline{\boldsymbol{\theta}}^{t+1} - \boldsymbol{\theta}^{PS}||_{2}^{2}] \leq \prod_{i=0}^{t} (1 - \tilde{\mu}\eta_{i}) ||\overline{\boldsymbol{\theta}}^{0} - \boldsymbol{\theta}^{PS}||_{2}^{2} + \frac{2c_{2}c_{3}}{\tilde{\mu}}\eta_{t}^{3} + \frac{2c_{1}c_{3}}{\tilde{\mu}}\eta_{t}^{2} + \frac{2c_{8}}{\tilde{\mu}}\eta_{t}, \tag{36}$$

(c₈ for Scheme I, replace c₈ with c₉ in Scheme II) where we define c₈ := $2\sigma^2 + \frac{4}{K}E^2G^2$ in Scheme I, and c₉ := $2\sigma^2 + \frac{4(N-K)}{K(N-1)}E^2G^2$ in Scheme II.

Proof. Note that

$$||\overline{\boldsymbol{\theta}}^{t+1} - \boldsymbol{\theta}^{PS}||_{2}^{2}$$

$$= ||\overline{\boldsymbol{\theta}}^{t+1} - \overline{\boldsymbol{w}}^{t+1} + \overline{\boldsymbol{w}}^{t+1} - \boldsymbol{\theta}^{PS}||_{2}^{2}$$

$$= \underbrace{||\overline{\boldsymbol{\theta}}^{t+1} - \overline{\boldsymbol{w}}^{t+1}||_{2}^{2}}_{T_{1}} + \underbrace{||\overline{\boldsymbol{w}}^{t+1} - \boldsymbol{\theta}^{PS}||_{2}^{2}}_{T_{2}} + \underbrace{2\langle \overline{\boldsymbol{w}}^{t+1} - \overline{\boldsymbol{\theta}}^{t+1}, \overline{\boldsymbol{\theta}}^{t+1} - \boldsymbol{\theta}^{PS}\rangle}_{T_{2}}$$
(37)

When expectation is taken over S_{t+1} , the last term T_3 vanishes due to Lemma C.1.

If $t+1 \notin \mathcal{I}_E$, T_1 vanishes since $\overline{\boldsymbol{\theta}}^{t+1} = \overline{\boldsymbol{w}}^{t+1}$ by definition when t+1 is not a communication step. For term T_2 , it's not hard to see that we can use Lemma B.1 to derive one step bounds for it (and use (25) in Lemma D.4), and thus we have

$$\mathbb{E}[||\overline{\boldsymbol{\theta}}^{t+1} - \boldsymbol{\theta}^{PS}||_{2}^{2}] = \mathbb{E}[||\overline{\boldsymbol{w}}^{t+1} - \boldsymbol{\theta}^{PS}||_{2}^{2}] \\
\leq (1 - \tilde{\mu}\eta_{t})\mathbb{E}||\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}||_{2}^{2} + c_{2}c_{7}\eta_{t}^{4} + c_{1}c_{7}\eta_{t}^{3} + 2\sigma^{2}\eta_{t}^{2}, \tag{38}$$

and we recall that $c_1 := \frac{L^2(1+\epsilon_{max})^2}{2\delta \overline{\epsilon}}, c_2 := 4[\sigma^2 + L^2(1+\epsilon_{max})^2], c_3 := 4(E-1)^2G^2, \tilde{\mu} := \mu - (1+\delta)\overline{\epsilon}L.$

If $t + 1 \in \mathcal{I}_E$, then we have the following result from Lemma D.6,

$$\mathbb{E}[||\overline{\boldsymbol{\theta}}^{t+1} - \boldsymbol{\theta}^{PS}||_{2}^{2}]
= \mathbb{E}[||\overline{\boldsymbol{\theta}}^{t+1} - \overline{\boldsymbol{w}}^{t+1}||_{2}^{2}] + \mathbb{E}[||\overline{\boldsymbol{w}}^{t+1} - \boldsymbol{\theta}^{PS}||_{2}^{2}]
\leq (1 - \tilde{\mu}\eta_{t})\mathbb{E}||\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}||_{2}^{2} + c_{2}c_{7}\eta_{t}^{4} + c_{1}c_{7}\eta_{t}^{3} + c_{8}\eta_{t}^{2},$$
(39)

where we recall $c_8:=2\sigma^2+\frac{4}{K}E^2G^2$ in Scheme I, and $c_9:=2\sigma^2+\frac{4(N-K)}{K(N-1)}E^2G^2$ in Scheme II.

The only difference between (26) and (39) is in $(c_8 - 2\sigma^2)\eta_t^2$. Therefore, we can use similar techniques to show the convergence,

$$\begin{split} \mathbb{E}[||\overline{\boldsymbol{\theta}}^{t+1} - \boldsymbol{\theta}^{PS}||_{2}^{2}] &\leq (1 - \tilde{\mu}\eta_{t})\mathbb{E}\big|\big|\overline{\boldsymbol{\theta}}^{t} - \boldsymbol{\theta}^{PS}\big|\big|_{2}^{2} + c_{2}c_{7}\eta_{t}^{4} + c_{1}c_{7}\eta_{t}^{3} + c_{8}\eta_{t}^{2} \\ &= \prod_{i=0}^{t} (1 - \tilde{\mu}\eta_{i})\big|\big|\overline{\boldsymbol{\theta}}^{0} - \boldsymbol{\theta}^{PS}\big|\big|_{2}^{2} \\ &+ \sum_{s=1}^{t} \prod_{i=s+1}^{t} (1 - \tilde{\mu}\eta_{i})\Big(c_{2}c_{7}\eta_{s}^{4} + c_{1}c_{7}\eta_{s}^{3} + c_{8}\eta_{s}^{2}\Big) \\ &\leq \prod_{i=0}^{t} (1 - \tilde{\mu}\eta_{i})\big|\big|\overline{\boldsymbol{\theta}}^{0} - \boldsymbol{\theta}^{PS}\big|\big|_{2}^{2} + \frac{2c_{2}c_{7}}{\tilde{\mu}}\eta_{t}^{3} + \frac{2c_{1}c_{7}}{\tilde{\mu}}\eta_{t}^{2} + \frac{2c_{8}}{\tilde{\mu}}\eta_{t} \end{split}$$

(c_8 for Scheme I, replace c_8 with c_9 in Scheme II).

Theorem D.8. (Full participation convergence theorem, alternative assumption)

Under Assumption 2.1, 2.2, 2.3, 2.5, the full participation scheme has convergence rate $\mathcal{O}(\frac{1}{T})$, i.e., denote $\Delta_t := \mathbb{E}[\|\overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}\|_2^2]$, then for some $\gamma > 0$,

$$\Delta_t \le \frac{\upsilon}{\gamma + t},\tag{40}$$

 $\textit{where } \upsilon := \max\{ \frac{c_2c_7\beta^4\gamma^{-2} + c_1c_7\beta^3\gamma^{-1} + c_8\beta^2}{\beta\mu - 1}, (\gamma+1)\Delta_1 \} \textit{ (c_8 for Scheme II, replace c_8 with c_9 in Scheme II)}.$

Proof. We will show it on the partial participation algorithm, and the proof for the full participation is similar.

For a diminishing step size $\eta_t = \frac{\beta}{t+\gamma}$ for some $\beta > \frac{1}{\tilde{\mu}}$ and $\gamma > 0$ such that $\eta_1 \leq \min\{\frac{1}{\tilde{\mu}}, \frac{1}{4L}\} = \frac{1}{4L}$ and $\eta_t \leq 2\eta_{t+E}$. We will prove $\triangle_t := \mathbb{E}[||\overline{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^{PS}||_2^2] \leq \frac{v}{\gamma+t}$, where $v := \max\{\frac{c_2c_7\beta^4\gamma^{-2} + c_1c_7\beta^3\gamma^{-1} + c_8\beta^2}{\beta\mu - 1}, (\gamma+1)\triangle_1\}$. We prove this by induction. Firstly, the definition if v ensures it holds for t=1. Assume it holds for some t, i.e., $\eta_t = \frac{\beta}{t+\gamma}$, then it follows that

$$\Delta_{t+1} \leq (1 - \eta_t \tilde{\mu}) \Delta_t + c_2 c_7 \eta_t^4 + c_1 c_7 \eta_t^3 + c_8 \eta_t^2
\leq (1 - \frac{\beta \tilde{\mu}}{t + \gamma}) \frac{v}{t + \gamma} + \frac{c_2 c_7 \beta^4}{(t + \gamma)^4} + \frac{c_1 c_7 \beta^3}{(t + \gamma)^3} + \frac{c_8 \beta^2}{(t + \gamma)^2}
= \frac{t + \gamma - 1}{(t + \gamma)^2} v + \left[\frac{c_2 c_7 \beta^4}{(t + \gamma)^4} + \frac{c_1 c_7 \beta^3}{(t + \gamma)^3} + \frac{c_8 \beta^2}{(t + \gamma)^2} - \frac{\beta \mu - 1}{(t + \gamma)^2} v \right]
\leq \frac{t + \gamma - 1}{(t + \gamma)^2} v + \left[\frac{c_2 c_7 \beta^4}{(t + \gamma)^2 \gamma^2} + \frac{c_1 c_7 \beta^3}{(t + \gamma)^2 \gamma} + \frac{c_8 \beta^2}{(t + \gamma)^2} - \frac{\beta \mu - 1}{(t + \gamma)^2} v \right]
\leq \frac{v}{t + \gamma + 1}$$
(41)

(c_8 for Scheme I, replace c_8 with c_9 in Scheme II), where $\tilde{\mu}$, c_1 to c_9 are defined the same as in earlier proofs, and thus the $\mathcal{O}(1/T)$ convergence rate is shown.

E Experiments

E.1 Numerical simulation

E.1.1 Learning rate decay (Figure 6)

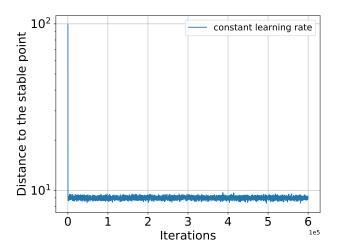


Figure 6: Constant learning rate on ExpGaussian. Full participation with E=10. The learning rate is set to 0.02.

E.1.2 Scheme II with lower learning rate (Figure 7)

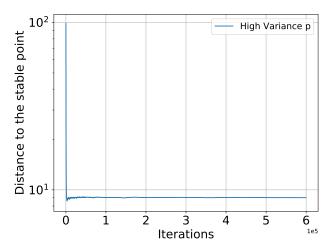


Figure 7: Sampling without replacement on ExpGaussian. The variance of p_i is set to 0.01. The learning rate is set to $\frac{5}{t+10000}$.

E.2 Credit score strategic classification

Figure 8 shows the losses and distances to the PS solution for the full participation, Scheme I and Scheme II using batch size 1 in client gradient descent. Figure 9 and 10 show the same figure with batch size 4 and 16, respectively. The scales of y axes are set equal for convenience of comparison. Using a larger batch size improves the convergence speed for all three schemes, especially for the two schemes of partial participation, both converging as fast as the full participation with batch size 16.

To study how batch batch size affect the convergence, we initialize P-FedAvg with θ^{PS} , the solution that minimizes the performative objective function. Due to the randomness of minibatch stochastic descent, we expect the parameter to

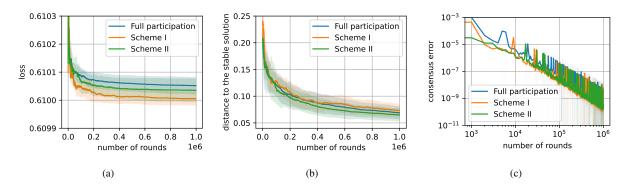


Figure 8: The losses (a) and distances (b) to the PS solution for the full participation, Scheme I and Scheme II using batch size 1 in client gradient descent.

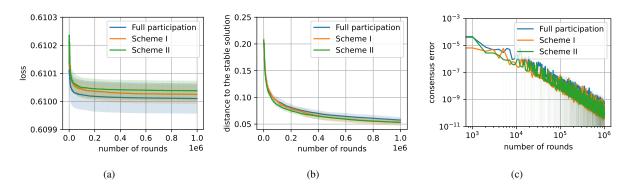


Figure 9: Same as Figure 8, but using batch size 4.

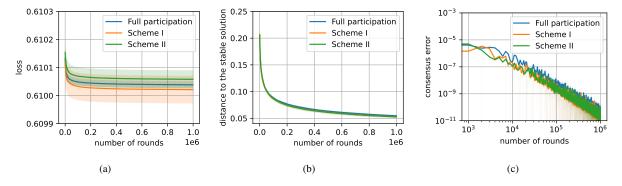


Figure 10: Same as Figure 8, but using batch size 16.

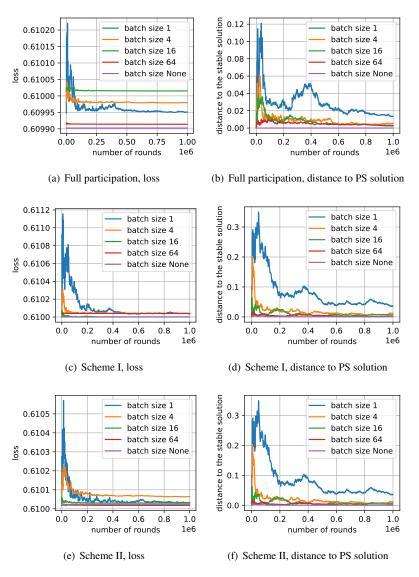


Figure 11: The loss functions and the distance to θ^{PS} of P-FedAvg initialized with θ^{PS} .

deviate from θ^{PS} and gradually stabilize back to θ^{PS} as the algorithm proceeds with decaying step sizes. It can be seen from Figure 11 (b), (c) and (f) that it is indeed the case for batch sizes larger than 1. This motivates our choice of a batch size larger than 1.