# Sample-based and Feature-based Federated Learning for Unconstrained and Constrained Nonconvex Optimization via Mini-batch SSCA

Ying Cui, Yangchen Li, and Chencheng Ye

*Abstract*—Federated learning (FL) has become a hot research area in enabling the collaborative training of machine learning models among multiple clients that hold sensitive local data. Nevertheless, unconstrained federated optimization has been studied mainly using stochastic gradient descent (SGD), which may converge slowly, and constrained federated optimization, which is more challenging, has not been investigated so far. This paper investigates sample-based and feature-based federated optimization, respectively, and considers both unconstrained and constrained nonconvex problems for each of them. First, we propose FL algorithms using stochastic successive convex approximation (SSCA) and mini-batch techniques. These algorithms can adequately exploit the structures of the objective and constraint functions and incrementally utilize samples. We show that the proposed FL algorithms converge to stationary points and Karush-Kuhn-Tucker (KKT) points of the respective unconstrained and constrained nonconvex problems, respectively. Next, we provide algorithm examples with appealing computational complexity and communication load per communication round. We show that the proposed algorithm examples for unconstrained federated optimization are identical to FL algorithms via momentum SGD and provide an analytical connection between SSCA and momentum SGD. Finally, numerical experiments demonstrate the inherent advantages of the proposed algorithms in convergence speeds, communication and computation costs, and model specifications.

*Index Terms*—Federated learning, nonconvex optimization, stochastic optimization, stochastic successive convex approximation.

## I. INTRODUCTION

Machine learning with distributed databases has been a hot research area [2]. The amount of data at each client can be large, and hence the data uploading to a central server may be constrained by energy and bandwidth limitations. Besides, local data may contain highly sensitive information, e.g., travel records, health information, and web browsing history, and thus a client may be unwilling to share it. Recent years have witnessed the growing interest in federated learning (FL), where data is maintained locally during the collaborative training of the server and clients [3], [4]. FL can protect data privacy for privacy-sensitive applications and improve communication efficiency.

Model aggregation, cryptographic methods, and differential privacy are three main privacy mechanisms in FL. They provide different privacy guarantees. Specifically, model aggregation, including model averaging and gradient averaging, is a basic privacy mechanism that reduces privacy risk by sharing model-related intermediate results computed based on local data [5]–[15]. Note that communicating locally computed results generally reveals much less information than communicating local data. Cryptographic methods, such as homomorphic encryption [10], [13] and secret sharing [16], further enhance privacy protection by encrypting locally computed results before sharing, at the cost of communication and computation efficiency reduction. Finally, differential privacy [11], [15] enhances privacy protection by adding random noise to locally computed results at the cost of model performance decline.

Depending on whether data is distributed over the sample space or feature space, FL can be classified into sample-based (horizontal) FL and feature-based (vertical) FL. Specifically, in sample-based FL [5]–[11], the datasets of different clients have the same feature space but no (or little) intersection on the sample space. On the contrary, in feature-based FL [12]–[15], the datasets of different clients share the same sample space but differ in the feature space. As a client cannot evaluate the impact of the model on the loss for a particular sample relying purely on its local data, feature-based FL is more challenging and hence less studied.

Existing works on FL [5]–[15] investigate only unconstrained optimization problems mainly using mini-batch stochastic gradient descent (SGD). In sample-based FL via mini-batch SGD [5]–[11], the global model is iteratively updated at the server by aggregating and averaging the clients' locally computed models or model-related results. Specifically, at one communication round, each client downloads the latest global model parameters and conducts one (e.g., in FedSGD [5]) or multiple (e.g., in FedAvg[1] [5] and PR-SGD [6]) local SGD updates to refine its local model. Multiple local SGD updates can reduce the communication cost (required number of communication rounds) with possibly increased computation cost. To further reduce communication cost, some recent works carefully design SGD update directions (e.g., momentum term [7]) or the numbers of local SGD updates at all clients [8], [9].

In contrast, the existing feature-based FL algorithms via

[1]In FedAvg, all local samples are utilized during local updates in each communication round.

mini-batch SGD [12]–[15] conduct only one SGD update in each communication round and impose additional restrictions on the structure of the loss function to guarantee privacy risk reduction. Specifically, the feature-based FL algorithms in [12]–[14] are designed only for two clients and some particular loss functions. In contrast, the feature-based FL algorithm in [15] applies to an arbitrary number of clients and a more general loss function. Besides, the feature-based FL algorithms in [12], [14], [15] do not maintain the global model at any node.

SGD has long been used for obtaining stationary points of unconstrained stochastic optimization problems [17] or Karush-Kuhn-Tucker (KKT) points of stochastic optimization problems with deterministic convex constraints [18]. Recently, stochastic successive convex approximation (SSCA) has been proposed to obtain KKT points of stochastic optimization problems with deterministic convex constraints [19] and with general stochastic nonconvex constraints [20], [21]. Apparently, SSCA applies to more types of constraints. Besides, SSCA empirically achieves a higher convergence speed than SGD [19].[2] Notice that [19]–[21] use only one sample at each iteration and may converge slowly when applied to machine learning problems with large datasets. Some recent works [22]–[24] have combined the SSCA algorithm in [19] and mini-batch techniques to solve unconstrained or convex constrained machine learning problems. However, SSCA has never been used for solving machine learning problems with nonconvex constraints or federated optimization problems.

In summary, there are several interesting questions: 1) whether mini-batch SSCA can apply to a broader range of federated optimization problems than mini-batch SGD, 2) whether mini-batch SSCA can converge faster than mini-batch SGD, and 3) whether mini-batch SSCA can reduce privacy risk in FL, like mini-batch SGD. In this paper, we would like to address the above questions. Specifically, we investigate general sample-based and feature-based federated optimization, respectively. For each of them, we consider both unconstrained and constrained nonconvex problems. The main contributions are summarized as follows.

- We propose FL algorithms for solving four federated optimization problems: unconstrained sample-based, constrained sample-based, unconstrained feature-based, and constrained feature-based federated optimization, using mini-batch SSCA. We show that the proposed FL algorithms converge to stationary points and KKT points of the respective unconstrained and constrained problems, respectively. Moreover, the proposed FL algorithms can adequately exploit the structures of the objective and constraint functions and incrementally utilize samples to improve convergence speeds. They can also reduce privacy risk through the model

aggregation mechanism, and their security can be enhanced via additional privacy mechanisms.
- We provide an example for each proposed FL algorithm. The algorithm examples for unconstrained sample-based and feature-based federated optimization have closed-form updates and achieve the same computational complexity (in order) and communication load per communication round as the corresponding SGD-based ones in [5]–[7] and [13], respectively. Besides, the algorithm examples and FL algorithms via momentum SGD with diminishing stepsizes perform identically, which is a rather surprising result.
- We consider two application examples in classification and customize the proposed FL algorithms to them. We show that the updates in the algorithms for the four federated optimization problems all have closed-form expressions. We also characterize the relationship between the two formulations.
- Numerical experiments demonstrate that in general, the proposed mini-batch SSCA-based FL algorithms for unconstrained federated optimization converge faster and achieve better computation and communication tradeoffs than the existing SGD-based ones [5]–[7], [13]. Furthermore, numerical experiments show that the proposed mini-batch SSCA-based FL algorithms for constrained federated optimization can more flexibly specify a training model.

To the best of our knowledge, this is the first work that applies SSCA to solve federated optimization, resolves constrained nonconvex federated optimization, and establishes an analytical connection between SSCA and momentum SGD. The key notation used in this paper is listed in Table I.

| Notation | Description |
|---|---|
| $I$ ($\mathcal{I}$) | number (index set) of clients |
| $N$ ($\mathcal{N}$) | number (index set) of samples |
| $K$ | dimension of the vector for each sample |
| $K_i$ | dimension of the $i$-th subvector for each sample |
| $\mathcal{N}_i$ | index set of samples at client $i$ |
| $B$ | batch size |
| $\mathbf{x}_n$ | vector for the $n$-th sample |
| $\mathbf{x}_{n,i}$ | the $i$-th subvector for the $n$-th sample |
| $\boldsymbol{\omega}$ | model parameters |
| $F_{a,m}(\boldsymbol{\omega})$ | objective or constraint function |
| $f_{a,m}(\boldsymbol{\omega}; \mathbf{x}_n)$ | loss for the $n$-th sample |
| $\bar{F}_{a,m}^{(t)}(\boldsymbol{\omega})$ | convex approximation of $F_{a,m}(\boldsymbol{\omega})$ at iteration $t$ |
| $\bar{f}_{a,m}(\boldsymbol{\omega}; \boldsymbol{\omega}', \mathbf{x}_n)$ | convex approximation of $f_{a,m}(\boldsymbol{\omega}; \mathbf{x}_n)$ around $\boldsymbol{\omega}'$ |

TABLE I. Key notation. $a = s$ and $a = f$ represent sample-based and feature-based, respectively. $m = 0$ and $m = 1, 2 \cdots$ represent the objective and $m$-th constraint, respectively.

## II. SYSTEM SETTING

Consider $N$ data samples, denoted by $\mathbf{x}_n \in \mathbb{R}^K, n \in \mathcal{N} \triangleq \{1, \cdots, N\}$. Consider a central server connected with $I$ local clients, each maintaining a local dataset.[3] Assume that the server and clients are honest-but-curious.[4] The server and $I$

---

[2]SGD utilizes first-order information of a sample estimate of the objective function and usually oscillates across narrow ravines. In contrast, SSCA uses a convex approximation of an incremental sample estimate of the objective/constraint function (reflecting more information) and effectively mitigates oscillations.

[3]The proposed SSCA-based algorithms can be used for solving federated optimization problems over streaming data and have theoretical convergence guarantees if the properties of the data stream do not change over time.

[4]The nodes will follow a predetermined algorithm but will attempt to infer private data using information received throughout the algorithm execution [25].

clients conduct FL, i.e., collaboratively train a model from the local datasets stored on the $I$ clients under the condition that each client cannot expose its local raw data to the server or the other clients. Depending on whether data is distributed over the sample space or feature space, FL can be typically classified into sample-based FL and feature-based FL.

In sample-based FL, the clients have the same feature space but differ in the sample space. Specifically, partition $\mathcal{N}$ into $I$ disjoint subsets, denoted by $\mathcal{N}_i$, $i \in \mathcal{I} \triangleq \{1, \cdots, I\}$, where $N_i \triangleq |\mathcal{N}_i|$ denotes the cardinality of the $i$-th subset and $\sum_{i\in\mathcal{I}} N_i = N$. For all $i \in \mathcal{I}$, the $i$-th client maintains a local dataset containing $N_i$ samples, i.e., $\mathbf{x}_n$, $n \in \mathcal{N}_i$. For example, two companies with similar businesses in different cities may have different user groups (from their respective regions) but the same type of data, e.g., users' occupations, ages, incomes, deposits, etc. The underlying optimization, termed sample-based federated optimization, is to minimize the following loss function:

$$F_{s,0}(\boldsymbol{\omega}) \triangleq \frac{1}{N} \sum_{n\in\mathcal{N}} f_{s,0}(\boldsymbol{\omega}; \mathbf{x}_n) \tag{1}$$

with respect to (w.r.t.) model parameters $\boldsymbol{\omega} \in \mathbb{R}^d$. Here, $f_{s,0}(\boldsymbol{\omega}; \mathbf{x}_n)$ represents the loss function for sample $\mathbf{x}_n$.

In feature-based FL, the clients have the same sample space but differ in the feature space. Specifically, for all $n \in \mathcal{N}$, $\mathbf{x}_n$ can be equivalently expressed by $I$ subvectors of it, denoted by $\mathbf{x}_{n,i} \in \mathbb{R}^{K_i}$, $i \in \mathcal{I}$, where $\sum_{i\in\mathcal{I}} K_i \geq K$.[5] With a slight abuse of notation, we write $\mathbf{x}_n = (\mathbf{x}_{n,i})_{i\in\mathcal{I}}$. For all $i \in \mathcal{I}$, the $i$-th client maintains $\mathbf{x}_{n,i}$, $n \in \mathcal{N}$.[6] For example, two companies in the same city with different businesses may have the same user group but different data types (from different types of businesses), e.g., one stores users' occupations and ages, and the other stores users' incomes and deposits. The underlying optimization, termed feature-based federated optimization, is to minimize the following loss function:

$$F_{f,0}(\boldsymbol{\omega}) \triangleq \frac{1}{N} \sum_{n\in\mathcal{N}} \underbrace{g_0\left(\boldsymbol{\omega}_0, (\mathbf{h}_{0,i}(\boldsymbol{\omega}_i, \mathbf{x}_{n,i}))_{i\in\mathcal{I}}\right)}_{\triangleq f_{f,0}(\boldsymbol{\omega}; \mathbf{x}_n)} \tag{2}$$

w.r.t. model parameters $\boldsymbol{\omega} \triangleq (\boldsymbol{\omega}_i)_{i=0,1,\cdots,I} \in \mathbb{R}^d$, where $\boldsymbol{\omega}_i \in \mathbb{R}^{d_i}$, $i = 0, 1, \cdots, I$ and $\sum_{i=0}^{I} d_i = d$. Here, $f_{f,0}(\boldsymbol{\omega}; \mathbf{x}_n)$ represents the loss function for sample $\mathbf{x}_n$, formed by composing $g_0 : \mathbb{R}^{d_0+H_0 I} \to \mathbb{R}$ with functions $\mathbf{h}_{0,i} : \mathbb{R}^{d_i+K_i} \to \mathbb{R}^{H_0}$, $i \in \mathcal{I}$, for some positive integer $H_0$. That is, we assume that the $i$-th block of model parameters, $\boldsymbol{\omega}_i$, and the $i$-th subvector for the $n$-th sample, $\mathbf{x}_{n,i}$, influence the loss of the $n$-th sample only via $\mathbf{h}_{0,i}(\boldsymbol{\omega}_i, \mathbf{x}_{n,i})$. We impose this additional restriction to enable privacy risk reduction via model aggregation in feature-based FL. It is worth noting that the existing works on feature-based FL impose the same restriction [15] or even stronger restrictions (e.g., $I = 2$ [12]–[14] and the loss function is the mean square error function [12] or cross-entropy function [13],

[14]).

In Section III and Section IV, we investigate sample-based FL and feature-based FL, respectively. To be general, we do not assume $F_{s,0}(\boldsymbol{\omega})$ and $F_{f,0}(\boldsymbol{\omega})$ to be convex in $\boldsymbol{\omega}$. To guarantee the convergence of the proposed FL algorithms, we assume that $f_{s,0}(\boldsymbol{\omega}; \mathbf{x}_n)$ and $f_{f,0}(\boldsymbol{\omega}; \mathbf{x}_n)$ satisfy the following assumption in the rest of the paper.[7]

*Assumption 1 (Assumption on $f(\boldsymbol{\omega}; \mathbf{x})$):* For any $\mathbf{x} \in \mathbb{R}^K$, $f(\boldsymbol{\omega}; \mathbf{x})$ is continuously differentiable, and its gradient is Lipschitz continuous on any compact set.

## III. SAMPLE-BASED FEDERATED LEARNING

In this section, we propose FL algorithms for unconstrained and constrained sample-based federated optimization problems, respectively, using mini-batch SSCA. In sample-based FL, the batch size $B$ satisfies $B \leq N_i, i \in \mathcal{I}$.

### A. Sample-based Federated Learning for Unconstrained Optimization

In this part, we consider the following unconstrained sample-based federated optimization problem:

*Problem 1 (Unconstrained Sample-based Federated Optimization):*

$$\min_{\boldsymbol{\omega}} \quad F_{s,0}(\boldsymbol{\omega})$$

where $F_{s,0}(\boldsymbol{\omega})$ is given by (1).

In [5]–[7], SGD is utilized to obtain a stationary point of Problem 1. SSCA can empirically achieve a higher convergence speed than SGD, as illustrated in Section I. In the following, we propose a sample-based FL algorithm, i.e., Algorithm 1, to obtain a stationary point of Problem 1 using mini-batch SSCA.[8]

*1) Algorithm Description:* The main idea of Algorithm 1 is to solve a sequence of successively refined convex problems, each of which is obtained by approximating $F_{s,0}(\boldsymbol{\omega})$ with a convex function based on its structure and randomly selected samples. Specifically, at iteration $t$, we choose an incremental sample estimate:

$$\bar{F}_{s,0}^{(t)}(\boldsymbol{\omega}) = (1 - \rho^{(t)}) \bar{F}_{s,0}^{(t-1)}(\boldsymbol{\omega}) + \rho^{(t)} \sum_{i\in\mathcal{I}} \frac{N_i}{BN} \sum_{n\in\mathcal{N}_i^{(t)}} \bar{f}_{s,0}(\boldsymbol{\omega}; \boldsymbol{\omega}_s^{(t)}, \mathbf{x}_n) \tag{3}$$

with $\bar{F}_{s,0}^{(0)}(\boldsymbol{\omega}) = 0$ as a convex approximation function of $F_{s,0}(\boldsymbol{\omega})$, where $\rho^{(t)}$ is a stepsize satisfying:

$$0 < \rho^{(t)} \leq 1, \ \lim_{t\to\infty} \rho^{(t)} = 0, \ \sum_{t=1}^{\infty} \rho^{(t)} = \infty, \tag{4}$$

---

[5]For unsupervised learning, $\mathbf{x}_{n,i} \in \mathbb{R}^{K_i}$, $i \in \mathcal{I}$ do not share any common coordinates of $\mathbf{x}_n$. For supervised learning, $\mathbf{x}_{n,i} \in \mathbb{R}^{K_i}$, $i \in \mathcal{I}$ share some common coordinates of $\mathbf{x}_n$, which represent the label of $\mathbf{x}_n$.

[6]The assumption that the $I$ local datasets share the same set of $N$ samples can be easily met using private set intersection techniques [26], [27].

[7]In Assumptions 1 and 2, we omit the subscripts $s, f$ for notation simplicity. Note that Assumptions 1 and 2 are necessary for the convergence of SSCA [19]–[21], and Assumption 1 is necessary for the convergences of SGD [6], [17], [18] and its variants [7].

[8]A machine learning problem involving a huge number of samples is usually transformed to an equivalent stochastic optimization problem and solved using stochastic optimization algorithms.

**Algorithm 1** Mini-batch SSCA for Problem 1
___
1: **initialize**: choose any $\boldsymbol{\omega}_s^1$ at the server.
2: **for** $t = 1, 2, \cdots, T-1$ **do**
3:    the server sends $\boldsymbol{\omega}_s^{(t)}$ to all clients.
4:    for all $i \in \mathcal{I}$, client $i$ randomly selects a mini-batch $\mathcal{N}_i^{(t)} \subseteq \mathcal{N}_i$, computes $\mathbf{q}_{s,0}\left(\boldsymbol{\omega}_s^{(t)}, (\mathbf{x}_n)_{n \in \mathcal{N}_i^{(t)}}\right)$, and sends it to the server.
5:    the server obtains $\bar{\boldsymbol{\omega}}_s^{(t)}$ by solving Problem 2.
6:    the server updates $\boldsymbol{\omega}_s^{(t+1)}$ according to (5).
7: **end for**
8: **Output**: $\boldsymbol{\omega}_s^T$
___

$\mathcal{N}_i^{(t)} \subseteq \mathcal{N}_i$ is a randomly selected mini-batch by client $i$ at iteration $t$, and $\bar{f}_{s,0}(\boldsymbol{\omega}; \boldsymbol{\omega}_s^{(t)}, \mathbf{x}_n)$ is a convex approximation[9] of $f_{s,0}(\boldsymbol{\omega}; \mathbf{x}_n)$ around $\boldsymbol{\omega}_s^{(t)}$ satisfying the following assumptions.

*Assumption 2 (Assumptions on $\bar{f}(\boldsymbol{\omega}; \boldsymbol{\omega}', \mathbf{x})$ for Approximating $f(\boldsymbol{\omega}; \mathbf{x})$ Around $\boldsymbol{\omega}'$):* 1) For any $\boldsymbol{\omega} \in \mathbb{R}^d$ and $\mathbf{x} \in \mathbb{R}^K$, $\nabla \bar{f}(\boldsymbol{\omega}; \boldsymbol{\omega}, \mathbf{x}) = \nabla f(\boldsymbol{\omega}; \mathbf{x})$; 2) For any $\boldsymbol{\omega}' \in \mathbb{R}^d$ and $\mathbf{x} \in \mathbb{R}^K$, $\bar{f}(\boldsymbol{\omega}; \boldsymbol{\omega}', \mathbf{x})$ is strongly convex w.r.t. $\boldsymbol{\omega}$; 3) For any $\mathbf{x} \in \mathbb{R}^K$, $\bar{f}(\boldsymbol{\omega}; \boldsymbol{\omega}', \mathbf{x})$ is Lipschitz continuous on any compact set; 4) For any $\boldsymbol{\omega}' \in \mathbb{R}^d$ and $\mathbf{x} \in \mathbb{R}^K$, $\bar{f}(\boldsymbol{\omega}; \boldsymbol{\omega}', \mathbf{x})$, its derivatives w.r.t. $\boldsymbol{\omega}$, and its second-order derivatives w.r.t. $\boldsymbol{\omega}$ are uniformly bounded on any compact set.

Note that for all $i \in \mathcal{I}$, mini-batch $\mathcal{N}_i' \subseteq \mathcal{N}_i$ with batch size $B$, and $\boldsymbol{\omega}' \in \mathbb{R}^d$, $\sum_{n \in \mathcal{N}_i'} \bar{f}_{s,0}(\boldsymbol{\omega}; \boldsymbol{\omega}', \mathbf{x}_n)$, a function of $\boldsymbol{\omega}$ with parameters jointly determined by $\boldsymbol{\omega}'$ and $\mathbf{x}_n, n \in \mathcal{N}_i'$, can be written naturally as $\sum_{n \in \mathcal{N}_i'} \bar{f}_{s,0}(\boldsymbol{\omega}; \boldsymbol{\omega}', \mathbf{x}_n) = p_{s,0}\left(\boldsymbol{\omega}, \mathbf{q}_{s,0}\left(\boldsymbol{\omega}', (\mathbf{x}_n)_{n \in \mathcal{N}_i'}\right)\right)$ with $p_{s,0} : \mathbb{R}^{d+D_0} \to \mathbb{R}$ and $\mathbf{q}_{s,0} : \mathbb{R}^{d+BK} \to \mathbb{R}^{D_0}$, for some positive integer $D_0$. Here, $\mathbf{q}_{s,0}\left(\boldsymbol{\omega}', (\mathbf{x}_n)_{n \in \mathcal{N}_i'}\right)$ represents the $D_0$ parameters of $\sum_{n \in \mathcal{N}_i'} \bar{f}_{s,0}(\boldsymbol{\omega}; \boldsymbol{\omega}', \mathbf{x}_n)$. Assume that the expressions of $\bar{f}_{s,0}$, $p_{s,0}$, and $\mathbf{q}_{s,0}$ are known to the server and $I$ clients. Each client $i \in \mathcal{I}$ computes $\mathbf{q}_{s,0}\left(\boldsymbol{\omega}_s^{(t)}, (\mathbf{x}_n)_{n \in \mathcal{N}_i^{(t)}}\right)$ and sends it to the server. Then, the server solves the following convex approximate problem to obtain $\bar{\boldsymbol{\omega}}_s^{(t)}$.

*Problem 2 (Convex Approximate Problem of Problem 1):*
$$\bar{\boldsymbol{\omega}}_s^{(t)} \triangleq \arg\min_{\boldsymbol{\omega}} \ \bar{F}_{s,0}^{(t)}(\boldsymbol{\omega})$$

Problem 2 is an unconstrained convex problem and can be solved with decent methods such as Newton's method. Given $\bar{\boldsymbol{\omega}}_s^{(t)}$, the server updates $\boldsymbol{\omega}_s^{(t)}$ according to:
$$\boldsymbol{\omega}_s^{(t+1)} = (1 - \gamma^{(t)})\boldsymbol{\omega}_s^{(t)} + \gamma^{(t)}\bar{\boldsymbol{\omega}}_s^{(t)}, \ t = 1, 2, \cdots \quad (5)$$
where $\gamma^{(t)}$ is a stepsize satisfying:
$$0 < \gamma^{(t)} \le 1, \ \lim_{t \to \infty} \gamma^{(t)} = 0, \ \sum_{t=1}^{\infty} \gamma^{(t)} = \infty,$$
$$\sum_{t=1}^{\infty} \left(\gamma^{(t)}\right)^2 < \infty, \ \lim_{t \to \infty} \frac{\gamma^{(t)}}{\rho^{(t)}} = 0. \quad (6)$$

[9]Usually, we preserve all convex terms in $f_{s,0}(\boldsymbol{\omega}; \mathbf{x}_n)$ and properly approximates the remaining nonconvex terms for reducing the approximation error or utilize the first-order approximation of $f_{s,0}(\boldsymbol{\omega}; \mathbf{x}_n)$ (see (7)) for reducing the computational complexity for solving Problem 2.

The detailed procedure is summarized in Algorithm 1.[10] The convergence of Algorithm 1 is summarized below. Algorithm 1 can empirically achieve a high convergence speed (shown in Section VI), as it can adequately exploit the structure of the objective function and incrementally utilize samples.

*Theorem 1 (Convergence of Algorithm 1):* Suppose that $f_{s,0}$ satisfies Assumption 1, $\bar{f}_{s,0}$ satisfies Assumption 2, and the sequence $\{\boldsymbol{\omega}_s^{(t)}\}$ generated by Algorithm 1 is bounded.[11] Then, every limit point of $\{\boldsymbol{\omega}_s^{(t)}\}$ is a stationary point of Problem 1 almost surely.

*Proof:* Please refer to Appendix A. ∎

*2) Security Analysis:* If for all $i \in \mathcal{I}$, mini-batch $\mathcal{N}_i' \subseteq \mathcal{N}_i$, and $\boldsymbol{\omega}' \in \mathbb{R}^d$, the system of equations w.r.t. $\mathbf{z} \in \mathbb{R}^{BK}$, i.e., $\mathbf{q}_{s,0}(\boldsymbol{\omega}', \mathbf{z}) = \mathbf{q}_{s,0}(\boldsymbol{\omega}', (\mathbf{x}_n)_{n \in \mathcal{N}_i'})$, has an infinite (or a sufficiently large) number of solutions, then raw data $\mathbf{x}_n$, $n \in \mathcal{N}_i^{(t)}$ can hardly be extracted by the server from $\mathbf{q}_{s,0}\left(\boldsymbol{\omega}_s^{(t)}, (\mathbf{x}_n)_{n \in \mathcal{N}_i^{(t)}}\right)$ in Step 4 of Algorithm 1, and hence, Algorithm 1 can reduce privacy risk based on model aggregation, like the existing sample-based FL algorithms via SGD [5]–[7]. Otherwise, extra privacy mechanisms can be applied to preserve data privacy. For example, if $\bar{\boldsymbol{\omega}}_s^{(t)}$ is linear in $\mathbf{q}_{s,0}\left(\boldsymbol{\omega}_s^{(t)}, (\mathbf{x}_n)_{n \in \mathcal{N}_i^{(t)}}\right)$, $i \in \mathcal{I}$, then homomorphic encryption [10] can be applied; if $\bar{\boldsymbol{\omega}}_s^{(t)}$ is a polynomial of $\mathbf{x}_n$ and $\boldsymbol{\omega}_s^{(t)}$, then secret sharing [16] can be applied.

*3) Algorithm Example:* We provide an example of $\bar{f}_{s,0}$ which satisfies Assumption 2 and yields an analytical solution of Problem 2:
$$\bar{f}_{s,0}(\boldsymbol{\omega}; \boldsymbol{\omega}_s^{(t)}, \mathbf{x}_n) = \left(\nabla f_{s,0}(\boldsymbol{\omega}_s^{(t)}; \mathbf{x}_n)\right)^T \left(\boldsymbol{\omega} - \boldsymbol{\omega}_s^{(t)}\right)$$
$$+ \tau \left\|\boldsymbol{\omega} - \boldsymbol{\omega}_s^{(t)}\right\|_2^2, \quad (7)$$
where $\tau > 0$ can be any constant, and the term $\tau \left\|\boldsymbol{\omega} - \boldsymbol{\omega}_s^{(t)}\right\|_2^2$ is used to ensure strong convexity. Then, $\sum_{n \in \mathcal{N}_i^{(t)}} \nabla f_{s,0}(\boldsymbol{\omega}_s^{(t)}; \mathbf{x}_n)$ can be viewed as $\mathbf{q}_{s,0}\left(\boldsymbol{\omega}_s^{(t)}, (\mathbf{x}_n)_{n \in \mathcal{N}_i^{(t)}}\right)$ (implying $D_0 = d$). Furthermore, substituting (7) into (3), $\bar{F}_{s,0}^{(t)}(\boldsymbol{\omega})$ can be rewritten as:
$$\bar{F}_{s,0}^{(t)}(\boldsymbol{\omega}) = \left(\hat{\mathbf{f}}_{s,0,1}^{(t)}\right)^T \boldsymbol{\omega} + \tau \|\boldsymbol{\omega}\|_2^2, \quad (8)$$
where $\hat{\mathbf{f}}_{s,0,1}^{(t)} \in \mathbb{R}^d$ is given by:
$$\hat{\mathbf{f}}_{s,0,1}^{(t)} = (1 - \rho^{(t)})\hat{\mathbf{f}}_{s,0,1}^{(t-1)}$$
$$+ \rho^{(t)} \sum_{i \in \mathcal{I}} \frac{N_i}{BN} \sum_{n \in \mathcal{N}_i^{(t)}} \left(\nabla f_{s,0}(\boldsymbol{\omega}_s^{(t)}; \mathbf{x}_n) - 2\tau \boldsymbol{\omega}_s^{(t)}\right) \quad (9)$$

[10]Each iteration of Algorithms 1-4 is implemented in one communication round. The computational complexity and communication load per communication round depend on the specific choices of $\bar{f}_{a,m}(\boldsymbol{\omega}; \boldsymbol{\omega}', \mathbf{x}_n)$, $a = s, f$ and $m = 0, 1, \cdots, M$.

[11]The conclusion of Theorems 1-4 still holds if the boundedness condition of the sequence in the theorem is replaced with the compact set constraint on $\boldsymbol{\omega}$ in the corresponding problem [19]–[21]. Note that the boundedness condition is easily satisfied in numerical experiments, and a simple compact set constraint that is sufficiently large can always be imposed without destroying the optimality [20], [21].

with $\hat{\mathbf{f}}_{s,0,1}^{(0)} = \mathbf{0}$. Apparently, Problem 2 with $\bar{f}_{s,0}$ given by (7) is an unconstrained convex quadratic programming w.r.t. $\boldsymbol{\omega}$. By the first-order optimality condition, it has the following analytical solution:

$$\bar{\boldsymbol{\omega}}_s^{(t)} = -\frac{1}{2\tau}\hat{\mathbf{f}}_{s,0,1}^{(t)}. \tag{10}$$

Therefore, Step 4 and Step 5 of Algorithm 1 with $\bar{f}_{s,0}$ given by (7) (i.e., an example of Algorithm 1) are given below. In Step 4, each client $i \in \mathcal{I}$ computes $\sum_{n \in \mathcal{N}_i^{(t)}} \nabla f_{s,0}(\boldsymbol{\omega}_s^{(t)}; \mathbf{x}_n) \in \mathbb{R}^d$ and sends the $d$-dimensional vector to the server. In Step 5, the server calculates $\bar{\boldsymbol{\omega}}_s^{(t)}$ according to (10). If for all $i \in \mathcal{I}$, $\mathcal{N}_i' \subseteq \mathcal{N}_i$, and $\boldsymbol{\omega}' \in \mathbb{R}^d$, the system of equations w.r.t. $(\mathbf{z}_n)_{n=1\cdots,B}$ with $\mathbf{z}_n \in \mathbb{R}^K$, $n = 1,\cdots,B$, i.e., $\sum_{n=1}^B \nabla f_{s,0}(\boldsymbol{\omega}'; \mathbf{z}_n) = \sum_{n \in \mathcal{N}_i'} \nabla f_{s,0}(\boldsymbol{\omega}'; \mathbf{x}_n)$, has an infinite (or a sufficiently large) number of solutions, then the example of Algorithm 1 can reduce privacy risk. Otherwise, homomorphic encryption [10], [13] can be applied to preserve data privacy, since $\bar{\boldsymbol{\omega}}_s^{(t)}$ is linear in $\nabla f_{s,0}(\boldsymbol{\omega}_s^{(t)}; \mathbf{x}_n)$, $n \in \mathcal{N}$, as shown in (9) and (10).

*Remark 1 (Comparison Between Example of Algorithm 1 and Sample-based FL Algorithms via SGD and Its Variants [5]–[7]):* Algorithm 1 with $\bar{f}_{s,0}$ given by (7) has the same order of computational complexity ($\mathcal{O}(B)$) and communication load per communication round as the sample-based FL algorithms via SGD and its variants [5]–[7], where $B$ samples are utilized by each client per communication round. Besides, it has the same level of privacy protection (due to the same system of equations for inferring private data) as the sample-based algorithm via SGD and its variants with one local SGD update per communication round (e.g., FedSGD [5]).

Finally, by (5), (9), and (10) and by choosing $\rho^{(1)} = 1$, $\{\boldsymbol{\omega}_s^{(t)}\}$ generated by the example of Algorithm 1 satisfies:

$$\boldsymbol{\omega}_s^{(t+1)} = \boldsymbol{\omega}_s^{(t)} - \gamma^{(t)}\mathbf{v}_s^{(t)}, \ t = 1, 2, \cdots \tag{11}$$

$$\mathbf{v}_s^{(t)} = \left(1 - \rho^{(t)}\right)\left(1 - \gamma^{(t-1)}\right)\mathbf{v}_s^{(t-1)}$$
$$+ \frac{\rho^{(t)}}{2\tau}\sum_{i \in \mathcal{I}}\frac{N_i}{BN}\sum_{n \in \mathcal{N}_i^{(t)}}\nabla f_{s,0}(\boldsymbol{\omega}_s^{(t)}; \mathbf{x}_n), \ t = 1, 2, \cdots \tag{12}$$

where $\gamma^{(0)} = 0$, $\mathbf{v}_s^{(0)} = 0$, and $\rho^{(t)}$ and $\gamma^{(t)}$ satisfy (3) and (5), respectively. From (11) and (12), we can make the following remark.

*Remark 2 (Connection Between Example of Algorithm 1 and Sample-based FL Algorithms via Momentum SGD [7]):* Algorithm 1 with $\bar{f}_{s,0}$ given by (7) can also be viewed as sample-based FL algorithm via momentum SGD with the momentum term $\mathbf{v}_s^{(t)}$ and diminishing stepsize $\gamma^{(t)}$ being the update direction and stepsize, respectively. This result also reveals an analytical connection between SSCA and momentum SGD, which is established for the first time. Furthermore, since the existing momentum SGD algorithms [7], [28] with theoretical convergence guarantees all rely on constant stepsizes, this work also enriches the results for momentum SGD.

## B. Sample-based Federated Learning for Constrained Optimization

In this part, we consider the following constrained sample-based federated optimization problem:

*Problem 3 (Constrained Sample-based Federated Optimization):*

$$\min_{\boldsymbol{\omega}} \quad F_{s,0}(\boldsymbol{\omega})$$
$$\text{s.t.} \quad F_{s,m}(\boldsymbol{\omega}) \leq 0, \ m = 1, 2, \cdots, M,$$

where $F_{s,0}(\boldsymbol{\omega})$ is given by (1), and

$$F_{s,m}(\boldsymbol{\omega}) \triangleq \frac{1}{N}\sum_{n \in \mathcal{N}} f_{s,m}(\boldsymbol{\omega}; \mathbf{x}_n), \ m = 1, 2, \cdots, M. \tag{13}$$

To be general, $F_{s,m}(\boldsymbol{\omega})$, $m = 0, \cdots, M$ are not assumed to be convex in $\boldsymbol{\omega}$. Notice that federated optimization with nonconvex constraints has not been investigated so far. In the following, we propose a sample-based FL algorithm, i.e., Algorithm 2, to obtain a KKT point of Problem 3, by combining the exact penalty method for SSCA in our previous work [21] and mini-batch techniques.

*1) Algorithm Description:* Sample convex approximations of Problem 3, obtained by directly approximating $F_{s,m}(\boldsymbol{\omega}), m = 0, 1, \cdots, M$ with the method proposed for $F_{s,0}$ in Section III-A, may not always be feasible, leading to possibly infeasible stochastic iterates [20], [21]. To ensure feasible stochastic iterates, we first transform Problem 3 to the following stochastic optimization problem whose objective function is the weighted sum of the original objective and the penalty for violating the original constraints [21]. We will soon see that its sample convex approximations are always feasible.

*Problem 4 (Transformed Problem of Problem 3):*

$$\min_{\boldsymbol{\omega}, \mathbf{s}} \quad F_{s,0}(\boldsymbol{\omega}) + c\sum_{m=1}^M s_m$$
$$\text{s.t.} \quad F_{s,m}(\boldsymbol{\omega}) \leq s_m, \ m = 1, 2, \cdots, M,$$
$$s_m \geq 0, \ m = 1, 2, \cdots, M,$$

where $\mathbf{s} \triangleq (s_m)_{m=1,\cdots,M}$ are slack variables, and $c > 0$ is a penalty parameter that trades off the original objective function and the slack penalty term.

At iteration $t$, we choose $\bar{F}_{s,0}^{(t)}(\boldsymbol{\omega})$ given in (3) as an approximation function of $F_{s,0}(\boldsymbol{\omega})$ and choose:

$$\bar{F}_{s,m}^{(t)}(\boldsymbol{\omega}) = (1 - \rho^{(t)})\bar{F}_{s,m}^{(t-1)}(\boldsymbol{\omega})$$
$$+ \rho^{(t)}\sum_{i \in \mathcal{I}}\frac{N_i}{BN}\sum_{n \in \mathcal{N}_i^{(t)}}\bar{f}_{s,m}(\boldsymbol{\omega}; \boldsymbol{\omega}_s^{(t)}, \mathbf{x}_n), \ m = 1, \cdots, M \tag{14}$$

with $\bar{F}_{s,m}^{(0)}(\boldsymbol{\omega}) = 0$ as a convex approximation function of $F_{s,m}(\boldsymbol{\omega})$, for all $m = 1, \cdots, M$, where $\rho^{(t)}$ is a stepsize satisfying (4), $\mathcal{N}_i^{(t)}$ is a randomly selected mini-batch by client $i$ at iteration $t$, and $\bar{f}_{s,m}(\boldsymbol{\omega}; \boldsymbol{\omega}_s^{(t)}, \mathbf{x}_n)$ is a convex approximation of $f_{s,m}(\boldsymbol{\omega}; \mathbf{x}_n)$ around $\boldsymbol{\omega}_s^{(t)}$ satisfying $\bar{f}_{s,m}(\boldsymbol{\omega}; \boldsymbol{\omega}, \mathbf{x}) = f_{s,m}(\boldsymbol{\omega}; \mathbf{x})$ and Assumption 2 for all $m = 1, \cdots, M$.

Note that for all $i \in \mathcal{I}$, mini-batch $\mathcal{N}_i' \subseteq \mathcal{N}_i$ with batch size $B$, and $\boldsymbol{\omega}' \in \mathbb{R}^d$, $\sum_{n \in \mathcal{N}_i'} \bar{f}_{s,m}(\boldsymbol{\omega}; \boldsymbol{\omega}', \mathbf{x}_n)$, $m = 0, \cdots, M$ can be written as $\sum_{n \in \mathcal{N}_i'} \bar{f}_{s,m}(\boldsymbol{\omega}; \boldsymbol{\omega}', \mathbf{x}_n) =$

**Algorithm 2** Mini-batch SSCA for Problem 3

---

1: **initialize**: choose any $\boldsymbol{\omega}_s^1$ and $c > 0$ at the server.
2: **for** $t = 1, 2, \cdots, T-1$ **do**
3:      the server sends $\boldsymbol{\omega}_s^{(t)}$ to all clients.
4:      for all $i \in \mathcal{I}$, client $i$ randomly selects a mini-batch $\mathcal{N}_i^{(t)} \subseteq \mathcal{N}_i$, computes $\mathbf{q}_{s,m}\left(\boldsymbol{\omega}_s^{(t)}, (\mathbf{x}_n)_{n \in \mathcal{N}_i^{(t)}}\right)$, $m = 0, 1, \cdots, M$, and sends them to the server.
5:      the server obtains $(\bar{\boldsymbol{\omega}}_s^{(t)}, \mathbf{s}_s^{(t)})$ by solving Problem 5.
6:      the server updates $\boldsymbol{\omega}_s^{(t+1)}$ according to (5).
7: **end for**
8: **Output**: $\boldsymbol{\omega}_s^T$

---

$p_{s,m}\left(\boldsymbol{\omega}, \mathbf{q}_{s,m}\left(\boldsymbol{\omega}', (\mathbf{x}_n)_{n \in \mathcal{N}_i'}\right)\right)$, $m = 0, \cdots, M$ with $p_{s,m} : \mathbb{R}^{D_m + d} \to \mathbb{R}$ and $\mathbf{q}_{s,m} : \mathbb{R}^{BK+d} \to \mathbb{R}^{D_m}$. Assume that the expressions of $\bar{f}_{s,m}$, $p_{s,m}$, $\mathbf{q}_{s,m}$, $m = 0, \cdots, M$ are known to the server and $I$ clients. Each client $i \in \mathcal{I}$ computes $\mathbf{q}_{s,m}\left(\boldsymbol{\omega}_s^{(t)}, (\mathbf{x}_n)_{n \in \mathcal{N}_i^{(t)}}\right)$, $m = 0, \cdots, M$ and sends them to the server. Then, the server solves the following convex approximate problem to obtain $\bar{\boldsymbol{\omega}}_s^{(t)}$.

*Problem 5 (Convex Approximate Problem of Problem 4):*

$$(\bar{\boldsymbol{\omega}}_s^{(t)}, \mathbf{s}_s^{(t)}) \triangleq \arg\min_{\boldsymbol{\omega}, \mathbf{s}} \quad \bar{F}_{s,0}^{(t)}(\boldsymbol{\omega}) + c \sum_{m=1}^{M} s_m$$
$$\text{s.t.} \quad \bar{F}_{s,m}^{(t)}(\boldsymbol{\omega}) \leq s_m, \quad m = 1, 2, \cdots, M,$$
$$s_m \geq 0, \quad m = 1, 2, \cdots, M.$$

Problem 5 is a constrained convex problem that is always feasible and can be readily solved with interior-point methods such as the barrier method.[12] Given $\bar{\boldsymbol{\omega}}_s^{(t)}$, the server updates $\boldsymbol{\omega}_s^{(t)}$ according to (5). The detailed procedure is summarized in Algorithm 2. The convergence of Algorithm 2 is summarized below. Consider a sequence $\{c_j\}$. For all $j$, let $(\boldsymbol{\omega}_{s,j}^\star, \mathbf{s}_{s,j}^\star)$ denote a limit point of $\{(\boldsymbol{\omega}_s^{(t)}, \mathbf{s}_s^{(t)})\}$ generated by Algorithm 2 with $c = c_j$.

*Theorem 2 (Convergence of Algorithm 2):* Suppose that $f_{s,m}$, $m = 0, \cdots, M$ satisfy Assumption 1, $\bar{f}_{s,0}$ satisfies Assumption 2, $\bar{f}_{s,m}$ satisfies $\bar{f}_{s,m}(\boldsymbol{\omega}; \boldsymbol{\omega}, \mathbf{x}) = f_{s,m}(\boldsymbol{\omega}; \mathbf{x})$ and Assumption 2 for all $m = 1, \cdots, M$, the sequence $\{\boldsymbol{\omega}_s^{(t)}\}$ generated by Algorithm 2 with $c = c_j$ is bounded for all $j$, and the sequence $\{c_j\}$ satisfies $0 < c_j < c_{j+1}$ and $\lim_{j \to \infty} c_j = \infty$. Then, the following statements hold. i) For all $j$, if $\mathbf{s}_{s,j}^\star = \mathbf{0}$, then $\boldsymbol{\omega}_{s,j}^\star$ is a KKT point of Problem 3 almost surely; ii) A limit point of $\{(\boldsymbol{\omega}_{s,j}^\star, \mathbf{s}_{s,j}^\star)\}$, denoted by $\{(\boldsymbol{\omega}_{s,\infty}^\star, \mathbf{s}_{s,\infty}^\star)\}$, satisfies that $\mathbf{s}_{s,\infty}^\star = \mathbf{0}$, and $\boldsymbol{\omega}_{s,\infty}^\star$ is a KKT point of Problem 3 almost surely.

*Proof:* Please refer to Appendix B. ∎

*2) Security Analysis:* If for all $i \in \mathcal{I}$, mini-batch $\mathcal{N}_i' \subseteq \mathcal{N}_i$, and $\boldsymbol{\omega}' \in \mathbb{R}^d$, the system of equations w.r.t. $\mathbf{z} \in \mathbb{R}^{BK}$, i.e., $\mathbf{q}_{s,m}(\boldsymbol{\omega}', \mathbf{z}) = \mathbf{q}_{s,m}\left(\boldsymbol{\omega}', (\mathbf{x}_n)_{n \in \mathcal{N}_i'}\right)$, $m = 0, \cdots, M$, has an infinite (or a sufficiently large) number of solutions,

---

[12]Problem 5 can be efficiently solved by the barrier method, regardless of how large $c$ (which influences only the linear terms of the objective function) is. This is because, in each centering step of the barrier method, an unconstrained centering problem is solved by Newton's method, whose convergence rate depends only on the smallest and largest eigenvalues and Lipschitz constant of the Hessian matrix of the objective function.

---

then raw data $\mathbf{x}_n$, $n \in \mathcal{N}_i^{(t)}$ can hardly be extracted from $\mathbf{q}_{s,m}\left(\boldsymbol{\omega}_s^{(t)}, (\mathbf{x}_n)_{n \in \mathcal{N}_i^{(t)}}\right)$, $m = 0, \cdots, M$ in Step 4 of Algorithm 2. Hence, Algorithm 2 can reduce privacy risk. Otherwise, extra privacy mechanisms need to be exploited. Note that FL for constrained optimization has not been studied so far, let alone privacy mechanisms for it.

*3) Algorithm Example:* We provide an example of $\bar{f}_{s,m}$, $m = 0, \cdots, M$ with $\bar{f}_{s,0}$ satisfying Assumption 2 and $\bar{f}_{s,m}$ satisfying $\bar{f}_{s,m}(\boldsymbol{\omega}; \boldsymbol{\omega}, \mathbf{x}) = f_{s,m}(\boldsymbol{\omega}; \mathbf{x})$ and Assumption 2 for all $m = 1, \cdots, M$. Specifically, we can choose $\bar{f}_{s,0}$ given by (7) and choose $\bar{f}_{s,m}$, $m = 1, \cdots, M$ as follows:

$$\bar{f}_{s,m}(\boldsymbol{\omega}, \boldsymbol{\omega}_s^{(t)}; \mathbf{x}_n) = f_{s,m}(\boldsymbol{\omega}_s^{(t)}; \mathbf{x}_n)$$
$$+ \left(\nabla f_{s,m}(\boldsymbol{\omega}_s^{(t)}; \mathbf{x}_n)\right)^T \left(\boldsymbol{\omega} - \boldsymbol{\omega}_s^{(t)}\right) + \tau \left\|\boldsymbol{\omega} - \boldsymbol{\omega}_s^{(t)}\right\|_2^2, \quad (15)$$

where $\tau > 0$ can be any constant. Then, $\sum_{n \in \mathcal{N}_i^{(t)}} \nabla f_{s,0}(\boldsymbol{\omega}_s^{(t)}; \mathbf{x}_n)$ can be viewed as $\mathbf{q}_{s,0}\left(\boldsymbol{\omega}_s^{(t)}, (\mathbf{x}_n)_{n \in \mathcal{N}_i^{(t)}}\right)$ (implying $D_0 = d$), and $\left(\sum_{n \in \mathcal{N}_i^{(t)}} f_{s,m}(\boldsymbol{\omega}_s^{(t)}; \mathbf{x}_n), \sum_{n \in \mathcal{N}_i^{(t)}} \nabla f_{s,m}(\boldsymbol{\omega}_s^{(t)}; \mathbf{x}_n)\right)$ can be viewed as $\mathbf{q}_{s,m}\left(\boldsymbol{\omega}_s^{(t)}, (\mathbf{x}_n)_{n \in \mathcal{N}_i^{(t)}}\right)$ (implying $D_m = 1+d$), for all $m = 1, \cdots, M$. Recall that with $f_{s,0}$ given in (7), $\bar{F}_{s,0}^{(t)}(\boldsymbol{\omega})$ is given in (8). In addition, for all $m = 1, \cdots, M$, substituting (15) into (14), $\bar{F}_{s,m}^{(t)}(\boldsymbol{\omega})$ can be rewritten as:

$$\bar{F}_{s,m}^{(t)}(\boldsymbol{\omega}) = \hat{f}_{s,m,0}^{(t)} + \left(\hat{\mathbf{f}}_{s,m,1}^{(t)}\right)^T \boldsymbol{\omega} + \tau \|\boldsymbol{\omega}\|_2^2, m = 1, \cdots, M,$$

where $\hat{f}_{s,m,0}^{(t)}$ and $\hat{\mathbf{f}}_{s,m,1}^{(t)} \in \mathbb{R}^d$ are given by:

$$\hat{f}_{s,m,0}^{(t)} = (1 - \rho^{(t)})\hat{f}_{s,m,0}^{(t-1)} + \rho^{(t)} \sum_{i \in \mathcal{I}} \frac{N_i}{BN} \sum_{n \in \mathcal{N}_i^{(t)}} \left(f_{s,m}(\boldsymbol{\omega}_s^{(t)}; \mathbf{x}_n)\right.$$
$$\left. - \left(\nabla f_{s,m}(\boldsymbol{\omega}_s^{(t)}; \mathbf{x}_n)\right)^T \boldsymbol{\omega}_s^{(t)} + \tau \left\|\boldsymbol{\omega}_s^{(t)}\right\|_2^2\right), m = 1, \cdots, M,$$
$$\hat{\mathbf{f}}_{s,m,1}^{(t)} = (1 - \rho^{(t)})\hat{\mathbf{f}}_{s,m,1}^{(t-1)}$$
$$+ \rho^{(t)} \sum_{i \in \mathcal{I}} \frac{N_i}{BN} \sum_{n \in \mathcal{N}_i^{(t)}} \left(\nabla f_{s,m}(\boldsymbol{\omega}_s^{(t)}; \mathbf{x}_n) - 2\tau \boldsymbol{\omega}_s^{(t)}\right), m = 1, \cdots, M,$$

with $\hat{f}_{s,m,0}^{(0)} = 0$ and $\hat{\mathbf{f}}_{s,m,1}^{(0)} = \mathbf{0}$. Apparently, Problem 5 with $\bar{f}_{s,0}$ given by (7) and $\bar{f}_{s,m}$, $m = 1, \cdots, M$ given by (15) is a convex quadratically constrained quadratic programming and can be solved using an interior-point method.

Therefore, Step 4 and Step 5 of Algorithm 2 with $\bar{f}_{s,0}$ given by (7) and $\bar{f}_{s,m}$, $m = 1, \cdots, M$ given by (15) (i.e., an example of Algorithm 2) are given below. In Step 4, each client $i \in \mathcal{I}$ computes $\sum_{n \in \mathcal{N}_i^{(t)}} \nabla f_{s,0}(\boldsymbol{\omega}_s^{(t)}; \mathbf{x}_n) \in \mathbb{R}^d$ and $\left(\sum_{n \in \mathcal{N}_i^{(t)}} f_{s,m}(\boldsymbol{\omega}_s^{(t)}; \mathbf{x}_n), \sum_{n \in \mathcal{N}_i^{(t)}} \nabla f_{s,m}(\boldsymbol{\omega}_s^{(t)}; \mathbf{x}_n)\right) \in \mathbb{R}^{1+d}, m = 1, \cdots, M$ and sends the $d$-dimensional vector and $M$ $(1+d)$-dimensional vectors to the server. In Step 5, the server calculates $(\bar{\boldsymbol{\omega}}_s^{(t)}, \mathbf{s}_s^{(t)})$ using an interior-point method. If for all $i \in \mathcal{I}$, $\mathcal{N}_i' \subseteq \mathcal{N}_i$, and $\boldsymbol{\omega}' \in \mathbb{R}^d$, the system of equations w.r.t. $(\mathbf{z}_n)_{n=1\cdots,B}$ with $\mathbf{z}_n \in \mathbb{R}^K$, i.e., $\sum_{n=1}^B f_{s,m}(\boldsymbol{\omega}'; \mathbf{z}_n) = \sum_{n \in \mathcal{N}_i'} f_{s,m}(\boldsymbol{\omega}'; \mathbf{x}_n)$, $m = 1, \cdots, M$ and $\sum_{n=1}^B \nabla f_{s,m}(\boldsymbol{\omega}'; \mathbf{z}_n) = \sum_{n \in \mathcal{N}_i'} \nabla f_{s,m}(\boldsymbol{\omega}'; \mathbf{x}_n)$, $m = 0, \cdots, M$, has an infinite (or a sufficiently large) number of

**Algorithm 3** Mini-batch SSCA for Problem 6

1: **initialize**: choose any $\boldsymbol{\omega}_f^1$ at the server.
2: **for** $t = 1, 2, \cdots, T - 1$ **do**
3:    the server randomly selects a mini-batch with the index set denoted by $\mathcal{N}^{(t)} \subset \mathcal{N}$ and sends $\mathcal{N}^{(t)}$ and $(\boldsymbol{\omega}_0^{(t)}, \boldsymbol{\omega}_i^{(t)})$ to client $i$ for all $i \in \mathcal{I}$.
4:    for all $i \in \mathcal{I}$, client $i$ computes $\mathbf{h}_{0,i}(\boldsymbol{\omega}_i^{(t)}, \mathbf{x}_{n,i})$, $n \in \mathcal{N}^{(t)}$ and sends them to the other clients.
5:    the client with the highest computation speed (or any client) computes $\mathbf{q}_{f,0,0}\left(\boldsymbol{\omega}_0^{(t)}, \left(\mathbf{h}_{0,i}(\boldsymbol{\omega}_i^{(t)}, \mathbf{x}_{n,i})\right)_{n \in \mathcal{N}^{(t)}, i \in \mathcal{I}}\right)$ and sends it to the server.
6:    for all $i \in \mathcal{I}$, client $i$ computes $\mathbf{q}_{f,0,i}\left(\boldsymbol{\omega}_0^{(t)}, \boldsymbol{\omega}_i^{(t)}, (\mathbf{x}_{n,i})_{n \in \mathcal{N}^{(t)}}, \left(\mathbf{h}_{0,j}(\boldsymbol{\omega}_j^{(t)}, \mathbf{x}_{n,j})\right)_{n \in \mathcal{N}^{(t)}, j \in \mathcal{I}}\right)$ and sends it to the server.
7:    the server obtains $\bar{\boldsymbol{\omega}}_f^{(t)}$ by solving Problem 7.
8:    the server updates $\boldsymbol{\omega}_f^{(t+1)}$ according to (18).
9: **end for**
10: **Output**: $\boldsymbol{\omega}_f^T$

---

solutions, then the example of Algorithm 2 can reduce privacy risk.

## IV. FEATURE-BASED FEDERATED LEARNING

In this section, we propose FL algorithms for unconstrained and constrained feature-based federated optimization problems, respectively, using mini-batch SSCA. In feature-based FL, the batch size $B$ satisfies $B \leq N$.

### A. Feature-based Federated Learning for Unconstrained Optimization

In this part, we consider the following unconstrained feature-based federated optimization problem:

*Problem 6 (Unconstrained Feature-based Federated Optimization):*

$$\min_{\boldsymbol{\omega}} \quad F_{f,0}(\boldsymbol{\omega})$$

where $F_{f,0}(\boldsymbol{\omega})$ is given by (2).

In [13], SGD is utilized to obtain a stationary point of Problem 6 only with $I = 2$ and $F_{f,0}(\boldsymbol{\omega})$ being the cross-entropy function. In the following, we propose a feature-based FL algorithm, i.e., Algorithm 3, to obtain a stationary point of Problem 6 using mini-batch SSCA, which empirically achieves a higher convergence speed than SGD.

*1) Algorithm Description:* At iteration $t$, we choose:

$$\bar{F}_{f,0}^{(t)}(\boldsymbol{\omega}) = (1 - \rho^{(t)})\bar{F}_{f,0}^{(t-1)}(\boldsymbol{\omega}) + \rho^{(t)}\frac{1}{B}\sum_{n \in \mathcal{N}^{(t)}} \bar{f}_{f,0}(\boldsymbol{\omega}; \boldsymbol{\omega}_f^{(t)}, \mathbf{x}_n) \tag{16}$$

with $\bar{F}_{f,0}^{(0)}(\boldsymbol{\omega}) = 0$ as a convex approximation function of $F_{f,0}(\boldsymbol{\omega})$, where $\rho^{(t)}$ is a stepsize satisfying (4), $\mathcal{N}^{(t)} \in \mathcal{N}$ is a randomly selected mini-batch by the server at iteration $t$, and $\bar{f}_{f,0}(\boldsymbol{\omega}; \boldsymbol{\omega}_f^{(t)}, \mathbf{x}_n)$ is a convex approximation of $f_{f,0}(\boldsymbol{\omega}; \mathbf{x}_n)$ around $\boldsymbol{\omega}_f^{(t)}$ satisfying Assumption 2.

Suppose that for any mini-batch $\mathcal{N}' \subseteq \mathcal{N}$ with batch size $B$, $\sum_{n \in \mathcal{N}'} \bar{f}_{f,0}(\boldsymbol{\omega}; \boldsymbol{\omega}', \mathbf{x}_n)$, a function of $\boldsymbol{\omega}$ with parameters jointly determined by $\boldsymbol{\omega}'$ and $\mathbf{x}_n, n \in \mathcal{N}_i'$, can be

written as (17), as shown at the top of the next page, with $p_{f,0} : \mathbb{R}^{d + \sum_{i=1}^I E_{0,i}} \to \mathbb{R}$, $\mathbf{q}_{f,0,0} : \mathbb{R}^{d_0 + H_0 BI} \to \mathbb{R}^{E_{0,0}}$, and $\mathbf{q}_{f,0,i} : \mathbb{R}^{d_0 + d_i + K_i B + H_0 BI} \to \mathbb{R}^{E_{0,i}}$, $i \in \mathcal{I}$, for some positive integers $E_i, i = 0, 1, \cdots, I$.[13] Assume that the expressions of $\bar{f}_{f,0}$, $p_{f,0}$, $\mathbf{q}_{f,0,0}$, $\mathbf{q}_{f,0,i}, i \in \mathcal{I}$, and $\mathbf{h}_{0,i}, i \in \mathcal{I}$ are known to the server and $I$ clients. Each client $i \in \mathcal{I}$ computes $\mathbf{h}_{0,i}(\boldsymbol{\omega}_i^{(t)}, \mathbf{x}_{n,i})$, $n \in \mathcal{N}^{(t)}$ and sends them to the other clients. The client with the highest computation speed (or any client) computes $\mathbf{q}_{f,0,0}\left(\boldsymbol{\omega}_0^{(t)}, \left(\mathbf{h}_{0,i}(\boldsymbol{\omega}_i^{(t)}, \mathbf{x}_{n,i})\right)_{n \in \mathcal{N}^{(t)}, i \in \mathcal{I}}\right)$ based on $\mathbf{h}_{0,i}(\boldsymbol{\omega}_i^{(t)}, \mathbf{x}_{n,i})$, $i \in \mathcal{I}$, $n \in \mathcal{N}^{(t)}$ and sends it to the server. Moreover, each client $i \in \mathcal{I}$ computes $\mathbf{q}_{f,0,i}\left(\boldsymbol{\omega}_0^{(t)}, \boldsymbol{\omega}_i^{(t)}, (\mathbf{x}_{n,i})_{n \in \mathcal{N}^{(t)}}, \left(\mathbf{h}_{0,j}(\boldsymbol{\omega}_j^{(t)}, \mathbf{x}_{n,j})\right)_{n \in \mathcal{N}^{(t)}, j \in \mathcal{I}}\right)$ and sends it to the server.[14] Then, the server solves the following convex approximate problem to obtain $\bar{\boldsymbol{\omega}}_f^{(t)}$.

*Problem 7 (Convex Approximate Problem of Problem 6):*

$$\bar{\boldsymbol{\omega}}_f^{(t)} \triangleq \arg\min_{\boldsymbol{\omega}} \bar{F}_{f,0}^{(t)}(\boldsymbol{\omega})$$

Like Problem 2, Problem 7 is an unconstrained convex problem and can be readily solved. Given $\bar{\boldsymbol{\omega}}_f^{(t)}$, the server updates $\boldsymbol{\omega}_f^{(t)}$ according to:

$$\boldsymbol{\omega}_f^{(t+1)} = (1 - \gamma^{(t)})\boldsymbol{\omega}_f^{(t)} + \gamma^{(t)}\bar{\boldsymbol{\omega}}_f^{(t)}, \ t = 1, 2, \cdots \tag{18}$$

where $\gamma^{(t)}$ is a stepsize satisfying (6). The detailed procedure is summarized in Algorithm 3. The convergence of Algorithm 3 is summarized below.

*Theorem 3 (Convergence of Algorithm 3):* Suppose that $f_{f,0}$ satisfies Assumption 1, $\bar{f}_{f,0}$ satisfies Assumption 2, and the sequence $\{\boldsymbol{\omega}_f^{(t)}\}$ generated by Algorithm 3 is bounded almost surely. Then, every limit point of $\{\boldsymbol{\omega}_f^{(t)}\}$ is a stationary point of Problem 6 almost surely.

*Proof:* Please refer to Appendix A. ∎

*2) Security Analysis:* Suppose 1) for all $i \in \mathcal{I}$, mini-batch $\mathcal{N}' \subseteq \mathcal{N}$, and $\boldsymbol{\omega}_i' \in \mathbb{R}^{d_i}$, the system of equations w.r.t. $(\boldsymbol{\theta}, (\mathbf{z}_n)_{n=1,\cdots,B}) \in \mathbb{R}^{d_i + BK_i}$ with $\boldsymbol{\theta} \in \mathbb{R}^{d_i}$ and $\mathbf{z}_n \in \mathbb{R}^{K_i}$, i.e., $\mathbf{h}_{0,i}(\boldsymbol{\theta}, \mathbf{z}_n) = \mathbf{h}_{0,i}(\boldsymbol{\omega}_i', \mathbf{x}_{n,i})$, $n \in \mathcal{N}'$, has an infinite (or a sufficiently large) number of solutions; 2) for any mini-batch $\mathcal{N}' \subseteq \mathcal{N}$ and $\boldsymbol{\omega}' \in \mathbb{R}^d$, the system of equations w.r.t. $(\mathbf{z}_{n,i})_{n=1,\cdots,B,i \in \mathcal{I}} \in \mathbb{R}^{BK}$ with $\mathbf{z}_{n,i} \in \mathbb{R}^{K_i}$, i.e., $\mathbf{q}_{f,0,0}\left(\boldsymbol{\omega}_0', (\mathbf{h}_{0,i}(\boldsymbol{\omega}_i', \mathbf{z}_{n,i}))_{n=1,\cdots,B,i \in \mathcal{I}}\right) = \mathbf{q}_{f,0,0}\left(\boldsymbol{\omega}_0', (\mathbf{h}_{0,i}(\boldsymbol{\omega}_i', \mathbf{x}_{n,i}))_{n \in \mathcal{N}', i \in \mathcal{I}}\right)$, $\mathbf{q}_{f,0,i}\left(\boldsymbol{\omega}_0', \boldsymbol{\omega}_i', (\mathbf{z}_{n,i})_{n \in \mathcal{N}'}, (\mathbf{h}_{0,j}(\boldsymbol{\omega}_j', \mathbf{z}_{n,j}))_{n=1,\cdots,B,j \in \mathcal{I}}\right) = \mathbf{q}_{f,0,i}\left(\boldsymbol{\omega}_0', \boldsymbol{\omega}_i', (\mathbf{x}_{n,i})_{n \in \mathcal{N}'}, (\mathbf{h}_{0,j}(\boldsymbol{\omega}_j', \mathbf{x}_{n,j}))_{n \in \mathcal{N}', j \in \mathcal{I}}\right)$, $i \in \mathcal{I}$, has an infinite (or a sufficiently large) number of solutions. In that case, raw data $\mathbf{x}_n$, $n \in \mathcal{N}^{(t)}$ can hardly be extracted by any client from $\mathbf{h}_{0,i}(\boldsymbol{\omega}_i^{(t)}, \mathbf{x}_{n,i})$, $n \in \mathcal{N}^{(t)}$, $i \in \mathcal{I}$ or by the server from $\mathbf{q}_{f,0,0}\left(\boldsymbol{\omega}_0^{(t)}, \left(\mathbf{h}_{0,i}(\boldsymbol{\omega}_i^{(t)}, \mathbf{x}_{n,i})\right)_{n \in \mathcal{N}^{(t)}, i \in \mathcal{I}}\right)$,

---

[13]This assumption is met by commonly used loss functions such as those in [12]–[14] and Section V.

[14]The information collection mechanism in Algorithm 3 can be viewed as an extension of that in the feature-based FL algorithm via SGD [13].

$$\sum_{n \in \mathcal{N}'} \bar{f}_{f,0}(\boldsymbol{\omega}; \boldsymbol{\omega}', \mathbf{x}_n) = p_{f,0}\bigg(\boldsymbol{\omega}, \mathbf{q}_{f,0,0}\Big(\boldsymbol{\omega}_0', (\mathbf{h}_{0,i}(\boldsymbol{\omega}_i', \mathbf{x}_{n,i}))_{n \in \mathcal{N}', i \in \mathcal{I}}\Big), \Big(\mathbf{q}_{f,0,i}\big(\boldsymbol{\omega}_0', \boldsymbol{\omega}_i', (\mathbf{x}_{n,i})_{n \in \mathcal{N}'}, (\mathbf{h}_{0,j}(\boldsymbol{\omega}_j', \mathbf{x}_{n,j}))_{n \in \mathcal{N}', j \in \mathcal{I}}\big)\Big)_{i \in \mathcal{I}}\bigg) \tag{17}$$

$\mathbf{q}_{f,0,i}\Big(\boldsymbol{\omega}_0^{(t)}, \boldsymbol{\omega}_i^{(t)}, (\mathbf{x}_{n,i})_{n \in \mathcal{N}^{(t)}}, \big(\mathbf{h}_{0,j}(\boldsymbol{\omega}_j^{(t)}, \mathbf{x}_{n,j})\big)_{n \in \mathcal{N}^{(t)}, j \in \mathcal{I}}\Big)$, $i \in \mathcal{I}$ in Steps 4-6 of Algorithm 3, and hence Algorithm 3 can reduce privacy risk. However, if the two assumptions mentioned above are not satisfied, extra privacy mechanisms are required. For instance, if $\bar{\boldsymbol{\omega}}_f^{(t)}$ is linear in $q_{f,0,0}$ and $q_{f,0,i}$, $i \in \mathcal{I}$, then homomorphic encryption [13] can be applied.

*3) Algorithm Example:* We provide an example of $\bar{f}_{f,0}$ which satisfies Assumption 2 and yields an analytical solution of Problem 7:

$$\bar{f}_{f,0}(\boldsymbol{\omega}; \boldsymbol{\omega}_f^{(t)}, \mathbf{x}_n) = \Big(\nabla_{\boldsymbol{\omega}} f_{f,0}(\boldsymbol{\omega}_f^{(t)}; \mathbf{x}_n)\Big)^T \Big(\boldsymbol{\omega} - \boldsymbol{\omega}_f^{(t)}\Big)$$
$$+ \tau \left\| \boldsymbol{\omega} - \boldsymbol{\omega}_f^{(t)} \right\|_2^2, \ n \in \mathcal{N}^{(t)}, \tag{19}$$

where $\tau > 0$ can be any constant. By the chain rule, we have:

$$\nabla_{\boldsymbol{\omega}_0} f_{f,0}(\boldsymbol{\omega}_f^{(t)}; \mathbf{x}_n) = \nabla_{\boldsymbol{\omega}_0} g_0\Big(\boldsymbol{\omega}_0^{(t)}, \big(\mathbf{h}_{0,i}(\boldsymbol{\omega}_i^{(t)}, \mathbf{x}_{n,i})\big)_{i \in \mathcal{I}}\Big),$$
$$n \in \mathcal{N}^{(t)}, \tag{20}$$

$$\nabla_{\boldsymbol{\omega}_i} f_{f,0}(\boldsymbol{\omega}_f^{(t)}; \mathbf{x}_n)$$
$$= \nabla_{\mathbf{h}_{0,i}} g_0\Big(\boldsymbol{\omega}_0^{(t)}, \big(\mathbf{h}_{0,i}(\boldsymbol{\omega}_i^{(t)}, \mathbf{x}_{n,i})\big)_{i \in \mathcal{I}}\Big)^T \frac{\partial \mathbf{h}_{0,i}(\boldsymbol{\omega}_i^{(t)}, \mathbf{x}_{n,i})}{\partial \boldsymbol{\omega}_i},$$
$$n \in \mathcal{N}^{(t)}, i \in \mathcal{I}. \tag{21}$$

Substituting (20) and (21) into (19), we know that $\sum_{n \in \mathcal{N}^{(t)}} \nabla_{\boldsymbol{\omega}_0} f_{f,0}(\boldsymbol{\omega}_f^{(t)}; \mathbf{x}_n)$ and $\sum_{n \in \mathcal{N}^{(t)}} \nabla_{\boldsymbol{\omega}_i} f_{f,0}(\boldsymbol{\omega}_f^{(t)}; \mathbf{x}_n), i \in \mathcal{I}$ can be viewed as $\mathbf{q}_{f,0,0}\Big(\boldsymbol{\omega}_0^{(t)}, \big(\mathbf{h}_{0,i}(\boldsymbol{\omega}_i^{(t)}, \mathbf{x}_{n,i})\big)_{n \in \mathcal{N}^{(t)}, i \in \mathcal{I}}\Big)$ (implying $E_{0,0} = d_0$) and $\mathbf{q}_{f,0,i}\Big(\boldsymbol{\omega}_0^{(t)}, \boldsymbol{\omega}_i^{(t)}, (\mathbf{x}_{n,i})_{n \in \mathcal{N}^{(t)}}, \big(\mathbf{h}_{0,j}(\boldsymbol{\omega}_j^{(t)}, \mathbf{x}_{n,j})\big)_{n \in \mathcal{N}^{(t)}, j \in \mathcal{I}}\Big)$ (implying $E_{0,i} = d_i$), $i \in \mathcal{I}$, respectively. Besides, substituting (19) into (16), $\bar{F}_{f,0}^{(t)}(\boldsymbol{\omega})$ can be rewritten as:

$$\bar{F}_{f,0}^{(t)}(\boldsymbol{\omega}) = \Big(\hat{\mathbf{f}}_{f,0,1}^{(t)}\Big)^T \boldsymbol{\omega} + \tau \|\boldsymbol{\omega}\|_2^2, \tag{22}$$

where $\hat{\mathbf{f}}_{f,0,1}^{(t)} \in \mathbb{R}^d$ is given by:

$$\hat{\mathbf{f}}_{f,0,1}^{(t)} = (1 - \rho^{(t)}) \hat{\mathbf{f}}_{f,0,1}^{(t-1)} + \frac{\rho^{(t)}}{B} \sum_{n \in \mathcal{N}^{(t)}} \Big(\nabla f_{f,0}(\boldsymbol{\omega}_f^{(t)}; \mathbf{x}_n) - 2\tau \boldsymbol{\omega}_f^{(t)}\Big) \tag{23}$$

with $\hat{\mathbf{f}}_{f,0,1}^{(0)} = \mathbf{0}$. Similar to Problem 2 with $\bar{f}_{s,0}$ given by (7), Problem 7 with $\bar{f}_{f,0}$ given by (19) is an unconstrained convex quadratic programming w.r.t. $\boldsymbol{\omega}$ and hence has the following analytical solution:

$$\bar{\boldsymbol{\omega}}_f^{(t)} = -\frac{1}{2\tau} \hat{\mathbf{f}}_{f,0,1}^{(t)}. \tag{24}$$

Therefore, Steps 5-7 of Algorithm 3 with $\bar{f}_{f,0}$ given by (19) (i.e., an example of Algorithm 3) are given below. In Step 5, the client with the highest computation speed (or any client) computes $\sum_{n \in \mathcal{N}^{(t)}} \nabla_{\boldsymbol{\omega}_0} f_{f,0}(\boldsymbol{\omega}_f^{(t)}; \mathbf{x}_n) \in \mathbb{R}^{d_0}$ and sends the

$d_0$-dimensional vector to the server. In Step 6, each client $i \in \mathcal{I}$ computes $\sum_{n \in \mathcal{N}^{(t)}} \nabla_{\boldsymbol{\omega}_i} f_{f,0}(\boldsymbol{\omega}_f^{(t)}; \mathbf{x}_n) \in \mathbb{R}^{d_i}$ and sends the $d_i$-dimensional vector to the server. In Step 7, the server calculates $\bar{\boldsymbol{\omega}}_f^{(t)}$ according to (24). Suppose that for all $i \in \mathcal{I}$, mini-batch $\mathcal{N}' \subseteq \mathcal{N}$, and $\boldsymbol{\omega}_i' \in \mathbb{R}^{d_i}$, the system of equations w.r.t. $(\boldsymbol{\theta}, (\mathbf{z}_n)_{n=1,\cdots,B}) \in \mathbb{R}^{d_i + BK_i}$ with $\boldsymbol{\theta} \in \mathbb{R}^{d_i}$ and $\mathbf{z}_n \in \mathbb{R}^{K_i}$, i.e., $\mathbf{h}_{0,i}(\boldsymbol{\theta}, \mathbf{z}_n) = \mathbf{h}_{0,i}(\boldsymbol{\omega}_i', \mathbf{x}_{n,i})$, $n \in \mathcal{N}'$, has an infinite (or a sufficiently large) number of solutions, and for any $\mathcal{N}' \subseteq \mathcal{N}$ and $\boldsymbol{\omega}' \in \mathbb{R}^d$, the system of equations w.r.t. $(\mathbf{z}_n)_{n=1\cdots,B}$ with $\mathbf{z}_n \in \mathbb{R}^K$, i.e., $\sum_{n=1}^B \nabla f_{f,0}(\boldsymbol{\omega}'; \mathbf{z}_n) = \sum_{n \in \mathcal{N}'} \nabla f_{f,0}(\boldsymbol{\omega}'; \mathbf{x}_n)$, has an infinite (or a sufficiently large) number of solutions. In that case, the example of Algorithm 3 can reduce privacy risk. If the two assumptions are not satisfied, homomorphic encryption [10], [13] can be applied to preserve data privacy, since $\bar{\boldsymbol{\omega}}_f^{(t)}$ is linear in $\nabla f_{f,0}(\boldsymbol{\omega}_f^{(t)}; \mathbf{x}_n)$, $n \in \mathcal{N}$, as shown in (23) and (24). Similarly, the example of Algorithm 3 can be viewed as a feature-based FL algorithm via momentum SGD with diminishing stepsize $\gamma^{(t)}$.

*Remark 3 (Comparison Between Example of Algorithm 3 and Feature-based FL Algorithm via SGD [13]):* Algorithm 3 with $\bar{f}_{f,0}$ given by (19) and the extension of the feature-based FL algorithm via SGD [13] (without extra privacy mechanisms) to the general case with $I > 2$ and $F_{f,0}(\boldsymbol{\omega})$ given in (2) have the same order of computational complexity ($\mathcal{O}(B)$) and communication load per communication round and the same level of privacy protection (due to the same system of equations for inferring private data).

*Remark 4 (Information Collection for Example of Algorithm 3):* When choosing $\bar{f}_{f,0}$ given by (19), another option for collecting information is to let each client $i \in \mathcal{I}$ directly send $\mathbf{h}_{0,i}(\boldsymbol{\omega}_i^{(t)}, \mathbf{x}_{n,i})$, $\nabla_{\boldsymbol{\omega}_i} \mathbf{h}_{0,i}(\boldsymbol{\omega}_i^{(t)}, \mathbf{x}_{n,i})$, $n \in \mathcal{N}^{(t)}$ to the sever.[15] In general, it has a lower communication load but higher privacy risk than the information collection mechanism in Steps 4-6 of the example of Algorithm 3, without using additional privacy mechanisms.[16]

*B. Feature-based Federated Learning for Constrained Optimization*

In this part, we consider the following constrained feature-based federated optimization problem:

*Problem 8 (Constrained Feature-based Federated Optimization):*

$$\min_{\boldsymbol{\omega}} \quad F_{f,0}(\boldsymbol{\omega})$$
$$\text{s.t.} \quad F_{f,m}(\boldsymbol{\omega}) \leq 0, \ m = 1, 2, \cdots, M,$$

---

[15]This one-step information collection mechanism can be viewed as an extension of that in the feature-based FL algorithm via SGD [15] to the case where the server maintains the global model.

[16]For the loss function given in (28), the one-step information collection mechanism exposes raw data (as $\nabla_{\boldsymbol{\omega}_i} \mathbf{h}_{0,i}(\boldsymbol{\omega}_i^{(t)}, \mathbf{x}_{n,i}) = \mathbf{x}_{n,i}$), whereas the one adopted in Algorithm 3 does not, as shown in Section V.

**Algorithm 4** Mini-batch SSCA for Problem 8

1: **initialize**: choose any $\boldsymbol{\omega}_f^1$ and $c > 0$ at the server.
2: **for** $t = 1, 2, \cdots, T-1$ **do**
3:    the server randomly selects a mini-batch with the index set denoted by $\mathcal{N}^{(t)} \subset \mathcal{N}$ and sends $\mathcal{N}^{(t)}$ and $(\boldsymbol{\omega}_0^{(t)}, \boldsymbol{\omega}_i^{(t)})$ to client $i$ for all $i \in \mathcal{I}$.
4:    for all $i \in \mathcal{I}$, client $i$ computes $\mathbf{h}_{m,i}(\boldsymbol{\omega}_i^{(t)}, \mathbf{x}_{n,i})$, $n \in \mathcal{N}^{(t)}$, $m = 0, 1, \cdots, M$ and sends them to the other clients.
5:    the client with the highest computation speed (or any client) computes $\quad \mathbf{q}_{f,m,0}\left(\boldsymbol{\omega}_0^{(t)}, \left(\mathbf{h}_{m,i}(\boldsymbol{\omega}_i^{(t)}, \mathbf{x}_{n,i})\right)_{n \in \mathcal{N}^{(t)}, i \in \mathcal{I}}\right)$, $m = 0, 1, \cdots, M$ and sends them to the server.
6:    for all $i \in \mathcal{I}$, client $i$ computes $\mathbf{q}_{f,m,i}\left(\boldsymbol{\omega}_0^{(t)}, \boldsymbol{\omega}_i^{(t)}, (\mathbf{x}_{n,i})_{n \in \mathcal{N}^{(t)}}, \left(\mathbf{h}_{m,j}(\boldsymbol{\omega}_j^{(t)}, \mathbf{x}_{n,j})\right)_{n \in \mathcal{N}^{(t)}, j \in \mathcal{I}}\right)$, $m = 0, 1, \cdots, M$ and sends them to the server.
7:    the server obtains $(\bar{\boldsymbol{\omega}}_f^{(t)}, \mathbf{s}_f^{(t)})$ by solving Problem 10.
8:    the server updates $\boldsymbol{\omega}_f^{(t+1)}$ according to (18).
9: **end for**
10: **Output**: $\boldsymbol{\omega}_f^T$

---

where $F_{f,0}(\boldsymbol{\omega})$ is given by (2), and

$$F_{f,m}(\boldsymbol{\omega}) \triangleq \frac{1}{N} \sum_{n \in \mathcal{N}} \underbrace{g_m\left(\boldsymbol{\omega}_0, (\mathbf{h}_{m,i}(\boldsymbol{\omega}_i, \mathbf{x}_{n,i}))_{i \in \mathcal{I}}\right)}_{\triangleq f_{f,m}(\boldsymbol{\omega}; \mathbf{x}_n)}.$$

Here, $f_{f,m}(\boldsymbol{\omega}; \mathbf{x}_n)$ is formed by composing $g_m : \mathbb{R}^{d_0 + H_m I} \to \mathbb{R}$ with functions $\mathbf{h}_{m,i} : \mathbb{R}^{d_i + K_i} \to \mathbb{R}^{H_m}, i \in \mathcal{I}$, for some positive integer $H_m$.

To be general, $F_{f,m}(\boldsymbol{\omega})$, $m = 0, \cdots, M$ are not assumed to be convex in $\boldsymbol{\omega}$. Analogously to Algorithm 2, we propose a feature-based FL algorithm, i.e., Algorithm 4, to obtain a KKT point of Problem 8, by combining the exact penalty method for SSCA in our previous work [21] and mini-batch techniques.

*1) Algorithm Description:* Similarly, to ensure feasible stochastic iterates, we first transform Problem 8 to the following stochastic optimization problem with a slack penalty term.

*Problem 9 (Transformed Problem of Problem 8):*

$$\min_{\boldsymbol{\omega}, \mathbf{s}} \quad F_{f,0}(\boldsymbol{\omega}) + c \sum_{m=1}^{M} s_m$$
$$\text{s.t.} \quad F_{f,m}(\boldsymbol{\omega}) \le s_m, \ m = 1, 2, \cdots, M,$$
$$s_m \ge 0, \ m = 1, 2, \cdots, M.$$

At iteration $t$, we choose $\bar{F}_{f,0}^{(t)}(\boldsymbol{\omega})$ given in (16) as an approximation function of $F_{f,0}(\boldsymbol{\omega})$ and choose:

$$\bar{F}_{f,m}^{(t)}(\boldsymbol{\omega}) = (1 - \rho^{(t)})\bar{F}_{f,m}^{(t-1)}(\boldsymbol{\omega}) + \rho^{(t)}\frac{1}{B}\sum_{n \in \mathcal{N}^{(t)}} \bar{f}_{f,m}(\boldsymbol{\omega}; \boldsymbol{\omega}_f^{(t)}, \mathbf{x}_n),$$
$$m = 1, \cdots, M \qquad (25)$$

with $\bar{F}_{f,m}^{(0)}(\boldsymbol{\omega}) = 0$ as a convex approximation function of $F_{f,m}(\boldsymbol{\omega})$, for all $m = 1, \cdots, M$, where $\rho^{(t)}$ is a stepsize satisfying (4), $\mathcal{N}^{(t)}$ is a randomly selected mini-batch by the server at iteration $t$, and $\bar{f}_{f,m}(\boldsymbol{\omega}; \boldsymbol{\omega}_f^{(t)}, \mathbf{x}_n)$ is a convex approximation of $f_{f,m}(\boldsymbol{\omega}; \mathbf{x}_n)$ around $\boldsymbol{\omega}_f^{(t)}$ satisfying $\bar{f}_{s,m}(\boldsymbol{\omega}; \boldsymbol{\omega}, \mathbf{x}) = f_{s,m}(\boldsymbol{\omega}; \mathbf{x})$ and Assumption 2 for all $m = 1 \cdots, M$.

Note that for any mini-batch $\mathcal{N}' \subseteq \mathcal{N}$ with

batch size $B$, $\sum_{n \in \mathcal{N}'} \bar{f}_{f,m}(\boldsymbol{\omega}; \boldsymbol{\omega}', \mathbf{x}_n)$ can be written as (26), as shown at the top of the next page, with $p_{f,m} : \mathbb{R}^{d + \sum_{i=0}^{I} E_{m,i}} \to \mathbb{R}$, $\mathbf{q}_{f,m,0} : \mathbb{R}^{d_0 + H_m BI} \to \mathbb{R}^{E_{m,0}}$, and $\mathbf{q}_{f,m,i} : \mathbb{R}^{d_0 + d_i + K_i B + H_m BI} \to \mathbb{R}^{E_{m,i}}$, $i \in \mathcal{I}$, for some positive integers $E_{m,i}, i = 0, 1, \cdots, I$. Assume that the expressions of $\bar{f}_{f,m}$, $p_{f,m}$, $\mathbf{q}_{f,m,0}$, $\mathbf{q}_{f,m,i}, i \in \mathcal{I}$, and $\mathbf{h}_{m,i}$, $m = 0, \cdots, M$, $i \in \mathcal{I}$ are known to the server and $I$ clients. Each client $i \in \mathcal{I}$ computes $\mathbf{h}_{m,i}(\boldsymbol{\omega}_i^{(t)}, \mathbf{x}_{n,i})$, $n \in \mathcal{N}^{(t)}$, $m = 0, \cdots, M$ and sends them to the other clients. Based on $\mathbf{h}_{m,i}(\boldsymbol{\omega}_i^{(t)}, \mathbf{x}_{n,i})$, $m = 0, \cdots, M$, $n \in \mathcal{N}^{(t)}$, $i \in \mathcal{I}$, the client with the highest computation speed (or any client) computes $\quad \mathbf{q}_{f,m,0}\left(\boldsymbol{\omega}_0^{(t)}, \left(\mathbf{h}_{m,i}(\boldsymbol{\omega}_i^{(t)}, \mathbf{x}_{n,i})\right)_{n \in \mathcal{N}^{(t)}, i \in \mathcal{I}}\right)$, $m = 0, \cdots, M$ and sends them to the server. Moreover, each client $i \in \mathcal{I}$ computes $\mathbf{q}_{f,m,i}\left(\boldsymbol{\omega}_0^{(t)}, \boldsymbol{\omega}_i^{(t)}, (\mathbf{x}_{n,i})_{n \in \mathcal{N}^{(t)}}, \left(\mathbf{h}_{m,j}(\boldsymbol{\omega}_j^{(t)}, \mathbf{x}_{n,j})\right)_{n \in \mathcal{N}^{(t)}, j \in \mathcal{I}}\right)$, $m = 0, \cdots, M$ and sends them to the server. Then, the server solves the following convex approximate problem to obtain $\bar{\boldsymbol{\omega}}_f^{(t)}$.

*Problem 10 (Convex Approximate Problem of Problem 9):*

$$(\bar{\boldsymbol{\omega}}_f^{(t)}, \mathbf{s}_f^{(t)}) \triangleq \arg\min_{\boldsymbol{\omega}, \mathbf{s}} \bar{F}_{f,0}^{(t)}(\boldsymbol{\omega}) + c \sum_{m=1}^{M} s_m$$
$$\text{s.t.} \quad \bar{F}_{f,m}^{(t)}(\boldsymbol{\omega}) \le s_m, \ m = 1, 2, \cdots, M,$$
$$s_m \ge 0, \ m = 1, 2, \cdots, M.$$

Like Problem 5, Problem 10 is a constrained convex problem that is always feasible and can be readily solved. Given $\bar{\boldsymbol{\omega}}_f^{(t)}$, the server updates $\boldsymbol{\omega}_f^{(t)}$ according to (18). The detailed procedure is summarized in Algorithm 4. The convergence of Algorithm 4 is summarized below. Consider a sequence $\{c_j\}$. For all $j$, let $(\boldsymbol{\omega}_{f,j}^{\star}, \mathbf{s}_{f,j}^{\star})$ denote a limit point of $\{(\boldsymbol{\omega}_f^{(t)}, \mathbf{s}_f^{(t)})\}$ generated by Algorithm 4 with $c = c_j$.

*Theorem 4 (Convergence of Algorithm 4):* Suppose that $f_{f,m}$ satisfies Assumption 1 for all $m = 0, \cdots, M$, $\bar{f}_{f,0}$ satisfies Assumption 2, $\bar{f}_{f,m}$ satisfies $\bar{f}_{f,m}(\boldsymbol{\omega}; \boldsymbol{\omega}, \mathbf{x}) = f_{f,m}(\boldsymbol{\omega}; \mathbf{x})$ and Assumption 2 for all $m = 1, \cdots, M$, the sequence $\{\boldsymbol{\omega}_f^{(t)}\}$ generated by Algorithm 4 with $c = c_j$ is bounded for all $j$, and the sequence $\{c_j\}$ satisfies $0 < c_j < c_{j+1}$ and $\lim_{j \to \infty} c_j = \infty$. Then, the following statements hold. i) For all $j$, if $\mathbf{s}_{f,j}^{\star} = \mathbf{0}$, then $\boldsymbol{\omega}_{f,j}^{\star}$ is a KKT point of Problem 8 almost surely; ii) A limit point of $\{(\boldsymbol{\omega}_{f,j}^{\star}, \mathbf{s}_{f,j}^{\star})\}$, denoted by $\{(\boldsymbol{\omega}_{f,\infty}^{\star}, \mathbf{s}_{f,\infty}^{\star})\}$, satisfies that $\mathbf{s}_{f,\infty}^{\star} = \mathbf{0}$, and $\boldsymbol{\omega}_{f,\infty}^{\star}$ is a KKT point of Problem 8 almost surely.

*Proof:* Please refer to Appendix B. ∎

*2) Security Analysis:* Suppose 1) for all $i \in \mathcal{I}$, mini-batch $\mathcal{N}' \subseteq \mathcal{N}$, and $\boldsymbol{\omega}_i' \in \mathbb{R}^{d_i}$, the system of equations w.r.t. $(\boldsymbol{\theta}, (\mathbf{z}_n)_{n=1, \cdots, B}) \in \mathbb{R}^{d_i + BK_i}$ with $\boldsymbol{\theta} \in \mathbb{R}^{d_i}$ and $\mathbf{z}_n \in \mathbb{R}^{K_i}$, i.e., $\mathbf{h}_{m,i}(\boldsymbol{\theta}, \mathbf{z}_n) = \mathbf{h}_{m,i}(\boldsymbol{\omega}_i', \mathbf{x}_{n,i})$, $n \in \mathcal{N}'$, $m = 0, \cdots, M$, has an infinite (or a sufficiently large) number of solutions; 2) for any mini-batch $\mathcal{N}' \subseteq \mathcal{N}$ and $\boldsymbol{\omega}' \in \mathbb{R}^d$, the system of equations w.r.t. $(\mathbf{z}_{n,i})_{n=1, \cdots, B, i \in \mathcal{I}} \in \mathbb{R}^{BK}$ with $\mathbf{z}_{n,i} \in \mathbb{R}^{K_i}$, i.e., $\mathbf{q}_{f,m,0}\left(\boldsymbol{\omega}_0', (\mathbf{h}_{m,i}(\boldsymbol{\omega}_i', \mathbf{z}_{n,i}))_{n=1, \cdots, B, i \in \mathcal{I}}\right) = \mathbf{q}_{f,m,0}\left(\boldsymbol{\omega}_0', (\mathbf{h}_{m,i}(\boldsymbol{\omega}_i', \mathbf{x}_{n,i}))_{n \in \mathcal{N}', i \in \mathcal{I}}\right)$, $m = 0, \cdots, M$ and

$$\sum_{n\in\mathcal{N}'}\bar{f}_{f,m}(\boldsymbol{\omega};\boldsymbol{\omega}',\mathbf{x}_n)=p_{f,m}\Big(\boldsymbol{\omega},\mathbf{q}_{f,m,0}(\boldsymbol{\omega}_0',(\mathbf{h}_{m,i}(\boldsymbol{\omega}_i',\mathbf{x}_{n,i}))_{n\in\mathcal{N}',i\in\mathcal{I}}),\Big(\mathbf{q}_{f,m,i}\Big(\boldsymbol{\omega}_0',\boldsymbol{\omega}_i',(\mathbf{x}_{n,i})_{n\in\mathcal{N}'},\big(\mathbf{h}_{m,j}(\boldsymbol{\omega}_j',\mathbf{x}_{n,j})\big)_{n\in\mathcal{N}',j\in\mathcal{I}}\Big)\Big)_{i\in\mathcal{I}}\Big),$$
$$m=0,\cdots,M \qquad\qquad (26)$$

$\mathbf{q}_{f,m,i}\left(\boldsymbol{\omega}_0',\boldsymbol{\omega}_i',(\mathbf{z}_{n,i})_{n\in\mathcal{N}'},\big(\mathbf{h}_{m,j}(\boldsymbol{\omega}_j',\mathbf{z}_{n,j})\big)_{n=1,\cdots,B,j\in\mathcal{I}}\right)=$ $\mathbf{q}_{f,m,i}\left(\boldsymbol{\omega}_0',\boldsymbol{\omega}_i',(\mathbf{x}_{n,i})_{n\in\mathcal{N}'},\big(\mathbf{h}_{m,j}(\boldsymbol{\omega}_j',\mathbf{x}_{n,j})\big)_{n\in\mathcal{N}',j\in\mathcal{I}}\right)$, $m=0,\cdots,M$, $i\in\mathcal{I}$, has an infinite (or a sufficiently large) number of solutions. In that case, raw data $\mathbf{x}_n$, $n\in\mathcal{N}^{(t)}$ can hardly be extracted by any client from $\mathbf{h}_{m,i}(\boldsymbol{\omega}_i^{(t)},\mathbf{x}_{n,i})$, $m=0,\cdots,M$, $n\in\mathcal{N}^{(t)}$, $i\in\mathcal{I}$ or by the server from $\mathbf{q}_{f,m,0}\left(\boldsymbol{\omega}_0^{(t)},\big(\mathbf{h}_{m,i}(\boldsymbol{\omega}_i^{(t)},\mathbf{x}_{n,i})\big)_{n\in\mathcal{N}^{(t)},i\in\mathcal{I}}\right)$, $\mathbf{q}_{f,m,i}\left(\boldsymbol{\omega}_0^{(t)},\boldsymbol{\omega}_i^{(t)},(\mathbf{x}_{n,i})_{n\in\mathcal{N}^{(t)}},\big(\mathbf{h}_{m,j}(\boldsymbol{\omega}_j^{(t)},\mathbf{x}_{n,j})\big)_{n\in\mathcal{N}^{(t)},j\in\mathcal{I}}\right)$, $m=0,\cdots,M$, $i\in\mathcal{I}$ in Steps 4-6 of Algorithm 4. Hence, Algorithm 4 can reduce privacy risk. However, extra privacy mechanisms need to be investigated if the two assumptions mentioned above are not satisfied.

*3) Algorithm Example:* We provide an example of $\bar{f}_{f,m}$, $m=0,\cdots,M$ with $\bar{f}_{f,0}$ satisfying Assumption 2 and $\bar{f}_{f,m}$ satisfying $\bar{f}_{f,m}(\boldsymbol{\omega};\boldsymbol{\omega},\mathbf{x})=f_{f,m}(\boldsymbol{\omega};\mathbf{x})$ and Assumption 2 for all $m=1,\cdots,M$. Specifically, we can choose $\bar{f}_{f,0}$ given by (19) and choose $\bar{f}_{f,m}$, $m=1,\cdots,M$ as follows:

$$\bar{f}_{f,m}(\boldsymbol{\omega};\boldsymbol{\omega}_f^{(t)},\mathbf{x}_n)=f_{f,m}(\boldsymbol{\omega}_f^{(t)};\mathbf{x}_n)+\Big(\nabla f_{f,m}(\boldsymbol{\omega}_f^{(t)};\mathbf{x}_n)\Big)^T\Big(\boldsymbol{\omega}-\boldsymbol{\omega}_f^{(t)}\Big)$$
$$+\tau\left\|\boldsymbol{\omega}-\boldsymbol{\omega}_f^{(t)}\right\|_2^2,\ m=1,\cdots,M, \qquad (27)$$

where $\tau>0$ can be any constant. Note that $\nabla_{\boldsymbol{\omega}}f_{f,m}(\boldsymbol{\omega}_f^{(t)};\mathbf{x}_n)$ can be computed according to the chain rule, similarly to (20) and (21). Thus, $\sum_{n\in\mathcal{N}^{(t)}}\nabla_{\boldsymbol{\omega}_0}f_{f,0}(\boldsymbol{\omega}_f^{(t)};\mathbf{x}_n)$ can be viewed as $\mathbf{q}_{f,0,0}\left(\boldsymbol{\omega}_0^{(t)},\big(\mathbf{h}_{0,i}(\boldsymbol{\omega}_i^{(t)},\mathbf{x}_{n,i})\big)_{n\in\mathcal{N}^{(t)},i\in\mathcal{I}}\right)$ (implying $E_{0,0}=d_0$); $\left(\sum_{n\in\mathcal{N}^{(t)}}f_{f,m}(\boldsymbol{\omega}_f^{(t)};\mathbf{x}_n),\sum_{n\in\mathcal{N}^{(t)}}\nabla_{\boldsymbol{\omega}_0}f_{f,m}(\boldsymbol{\omega}_f^{(t)};\mathbf{x}_n)\right)$ can be viewed as $\mathbf{q}_{f,m,0}\left(\boldsymbol{\omega}_0^{(t)},\big(\mathbf{h}_{m,i}(\boldsymbol{\omega}_i^{(t)},\mathbf{x}_{n,i})\big)_{n\in\mathcal{N}^{(t)},i\in\mathcal{I}}\right)$ (implying $E_{m,0}=1+d_0$), for all $m=1,\cdots,M$; and $\sum_{n\in\mathcal{N}^{(t)}}\nabla_{\boldsymbol{\omega}_i}f_{f,m}(\boldsymbol{\omega}_f^{(t)};\mathbf{x}_n)$ can be viewed as $\mathbf{q}_{f,m,i}\left(\boldsymbol{\omega}_0^{(t)},\boldsymbol{\omega}_i^{(t)},(\mathbf{x}_{n,i})_{n\in\mathcal{N}^{(t)}},\big(\mathbf{h}_{m,j}(\boldsymbol{\omega}_j^{(t)},\mathbf{x}_{n,j})\big)_{n\in\mathcal{N}^{(t)},j\in\mathcal{I}}\right)$ (implying $E_{m,i}=d_i$), for all $m=0,\cdots,M$, $i\in\mathcal{I}$. Recall that $\bar{F}_{f,0}^{(t)}(\boldsymbol{\omega})$ is given in (22) with $f_{f,0}$ given in (19). In addition, for all $m=1,\cdots,M$, substituting (27) into (25), $\bar{F}_{f,m}^{(t)}(\boldsymbol{\omega})$ can be rewritten as:

$$\bar{F}_{f,m}^{(t)}(\boldsymbol{\omega})=\hat{f}_{f,m,0}^{(t)}+\Big(\hat{\mathbf{f}}_{f,m,1}^{(t)}\Big)^T\boldsymbol{\omega}+\tau\|\boldsymbol{\omega}\|_2^2,\ m=1,\cdots,M,$$

where $\hat{f}_{f,m,0}^{(t)}$ and $\hat{\mathbf{f}}_{f,m,1}^{(t)}\in\mathbb{R}^d$ are given by:

$$\hat{f}_{f,m,0}^{(t)}=(1-\rho^{(t)})\hat{f}_{f,m,0}^{(t-1)}+\rho^{(t)}\frac{1}{B}\sum_{n\in\mathcal{N}^{(t)}}\Big(f_{f,m}(\boldsymbol{\omega}_f^{(t)};\mathbf{x}_n)$$
$$-\Big(\nabla f_{f,m}(\boldsymbol{\omega}_f^{(t)};\mathbf{x}_n)\Big)^T\boldsymbol{\omega}_f^{(t)}+\tau\left\|\boldsymbol{\omega}_f^{(t)}\right\|_2^2\Big),\ m=1,\cdots,M,$$

$$\hat{\mathbf{f}}_{f,m,1}^{(t)}=(1-\rho^{(t)})\hat{\mathbf{f}}_{f,m,1}^{(t-1)}+\rho^{(t)}\frac{1}{B}$$
$$\times\sum_{n\in\mathcal{N}^{(t)}}\Big(\nabla f_{f,m}(\boldsymbol{\omega}_f^{(t)};\mathbf{x}_n)-2\tau\boldsymbol{\omega}_f^{(t)}\Big),\ m=1,\cdots,M$$

with $\hat{f}_{f,m,0}^{(0)}=0$ and $\hat{\mathbf{f}}_{f,m,1}^{(0)}=\mathbf{0}$. Problem 10 with $\bar{f}_{f,0}$ given by (19) and $\bar{f}_{f,m}$, $m=1,\cdots,M$ given by (27) is a convex quadratically constrained quadratic programming and can be solved using an interior-point method.

Therefore, Steps 5-7 of Algorithm 4 with $\bar{f}_{f,0}$ given by (19) and $\bar{f}_{f,m},m=1,\cdots,M$ given by (27) (i.e., an example of Algorithm 4) are given below. In Step 5, the client with the highest computation speed (or any client) computes $\sum_{n\in\mathcal{N}^{(t)}}\nabla_{\boldsymbol{\omega}_0}f_{f,0}(\boldsymbol{\omega}_f^{(t)};\mathbf{x}_n)\in\mathbb{R}^{d_0}$ and $\left(\sum_{n\in\mathcal{N}^{(t)}}f_{f,m}(\boldsymbol{\omega}_f^{(t)};\mathbf{x}_n),\sum_{n\in\mathcal{N}^{(t)}}\nabla_{\boldsymbol{\omega}_0}f_{f,m}(\boldsymbol{\omega}_f^{(t)};\mathbf{x}_n)\right)\in\mathbb{R}^{1+d_0}$, $m=1,\cdots,M$ and sends the $d_0$-dimensional vector and $M$ $(1+d_0)$-dimensional vectors to the server. In Step 6, each client $i\in\mathcal{I}$ computes $\sum_{n\in\mathcal{N}^{(t)}}\nabla_{\boldsymbol{\omega}_i}f_{f,m}(\boldsymbol{\omega}_f^{(t)};\mathbf{x}_n)$, $m=0,\cdots,M$ and sends the $(M+1)$ $d_i$-dimensional vectors to the server. In Step 7, the server calculates $(\bar{\boldsymbol{\omega}}_f^{(t)},\mathbf{s}_f^{(t)})$ using an interior-point method. Suppose that for all $i\in\mathcal{I}$, mini-batch $\mathcal{N}'\subseteq\mathcal{N}$, and $\boldsymbol{\omega}_i'\in\mathbb{R}^{d_i}$, the system of equations w.r.t $(\boldsymbol{\theta},(\mathbf{z}_n)_{n=1,\cdots,B})\in\mathbb{R}^{d_i+BK_i}$ with $\boldsymbol{\theta}\in\mathbb{R}^{d_i}$ and $\mathbf{z}_n\in\mathbb{R}^{K_i}$, i.e., $\mathbf{h}_{m,i}(\boldsymbol{\theta},\mathbf{z}_n)=\mathbf{h}_{m,i}(\boldsymbol{\omega}_i',\mathbf{x}_{n,i})$, $n\in\mathcal{N}'$, $m=0,\cdots,M$, has an infinite (or a sufficiently large) number of solutions; and for all $\mathcal{N}'\subseteq\mathcal{N}$ and $\boldsymbol{\omega}'\in\mathbb{R}^d$, the system of equations w.r.t. $(\mathbf{z}_n)_{n=1,\cdots,B}$ with $\mathbf{z}_n\in\mathbb{R}^K$, i.e., $\sum_{n=1}^B f_{f,m}(\boldsymbol{\omega}';\mathbf{z}_n)=\sum_{n\in\mathcal{N}'}f_{f,m}(\boldsymbol{\omega}';\mathbf{x}_n)$, $m=1,\cdots,M$ and $\sum_{n=1}^B\nabla f_{f,m}(\boldsymbol{\omega}';\mathbf{z}_n)=\sum_{n\in\mathcal{N}'}\nabla f_{f,m}(\boldsymbol{\omega}';\mathbf{x}_n)$, $m=0,\cdots,M$, has an infinite (or a sufficiently large) number of solutions. In that case, the example of Algorithm 4 can reduce privacy risk.

## V. APPLICATION EXAMPLES

In this section, we customize the proposed algorithmic frameworks to some applications and provide detailed solutions for the specific problems. The server and $I$ clients collaboratively solve an $L$-class classification problem with a dataset of $N$ samples using FL. Denote $\mathcal{P}\triangleq\{1,\cdots,P\}$ and $\mathcal{L}\triangleq\{1,\cdots,L\}$. The $n$-th sample is represented by $\mathbf{x}_n\triangleq(\mathbf{z}_n,\mathbf{y}_n)\in\mathbb{R}^K$, where $K=P+L$, and $\mathbf{z}_n\triangleq(z_{n,p})_{p\in\mathcal{P}}\in\mathbb{R}^P$ and $\mathbf{y}_n\triangleq(y_{n,l})_{l\in\mathcal{L}}\in\{0,1\}^L$ represent the $P$ features and label of the $n$-th sample, respectively. In feature-based FL, $\mathcal{P}$ is partitioned into $I$ subsets, denoted by $\mathcal{P}_i,i\in\mathcal{I}$, and for each sample $n\in\mathcal{N}$, client $i$ maintains the $P_i$ features $\mathbf{z}_{n,i}\triangleq(z_{n,p})_{p\in\mathcal{P}_i}\in\mathbb{R}^{P_i}$ and the label $\mathbf{y}_n$. Note that $P=\sum_{i\in\mathcal{I}}P_i$. Thus, the $i$-th subvector for the $n$-th sample is given by $\mathbf{x}_{n,i}\triangleq(\mathbf{z}_{n,i},\mathbf{y}_n)$.

Consider a two-layer neural network, including an input layer composed of $P$ cells, a hidden layer composed of $J$ cells, and an output layer composed of $L$ cells. Denote $\mathcal{J} \triangleq \{1, \cdots, J\}$. The model parameters are represented by $\boldsymbol{\omega} \triangleq ((\omega_{0,l,j})_{l \in \mathcal{L}, j \in \mathcal{J}}, (\omega_{1,j,p})_{j \in \mathcal{J}, p \in \mathcal{P}}) \in \mathbb{R}^d$, where $d = J(P + L)$. For feature-based FL, $\boldsymbol{\omega}$ is also expressed as $\boldsymbol{\omega} = (\boldsymbol{\omega}_0, (\boldsymbol{\omega}_i)_{i \in \mathcal{I}})$, where $\boldsymbol{\omega}_0 \triangleq (\omega_{0,l,j})_{l \in \mathcal{L}, j \in \mathcal{J}}$ and $\boldsymbol{\omega}_i \triangleq (\omega_{1,j,p})_{j \in \mathcal{J}, p \in \mathcal{P}_i}$, $i \in \mathcal{I}$. We use the swish activation function $S(z) = z/(1 + \exp(-z))$ [29] for the hidden layer and the softmax activation function for the output layer. Note that $S'(z) = \frac{1}{1+\exp(-z)}\left(1 + \frac{z \exp(-z)}{1+\exp(-z)}\right)$. We consider the cross-entropy loss function. Thus, the resulting loss function for sample-based and feature-based FL is given by:

$$F(\boldsymbol{\omega}) \triangleq -\frac{1}{N} \sum_{n \in \mathcal{N}} \sum_{l \in \mathcal{L}} y_{n,l} \log\left(Q_l(\boldsymbol{\omega}; \mathbf{x}_n)\right), \qquad (28)$$

where

$$Q_l(\boldsymbol{\omega}; \mathbf{x}_n) \triangleq \frac{\exp(\sum_{j \in \mathcal{J}} \omega_{0,l,j} S(\sum_{p \in \mathcal{P}} \omega_{1,j,p} z_{n,p}))}{\sum_{h=1}^{L} \exp(\sum_{j \in \mathcal{J}} \omega_{0,h,j} S(\sum_{p \in \mathcal{P}} \omega_{1,j,p} z_{n,p}))}.$$

For ease of exposition, in the rest of this section, we denote:

$$\bar{A}_{a,l,j}^{(t)} \triangleq \begin{cases} \sum_{i \in \mathcal{I}} \frac{N_i}{BN} \sum_{n \in \mathcal{N}_i^{(t)}} \bar{a}_{a,n,l,j}, & a = s \\ \frac{1}{B} \sum_{n \in \mathcal{N}^{(t)}} \bar{a}_{a,n,l,j}, & a = f \end{cases}, \qquad (29)$$

$$\bar{B}_{a,j,p}^{(t)} \triangleq \begin{cases} \sum_{i \in \mathcal{I}} \frac{N_i}{BN} \sum_{n \in \mathcal{N}_i^{(t)}} \bar{b}_{a,n,j,p}, & a = s \\ \frac{1}{B} \sum_{n \in \mathcal{N}^{(t)}} \bar{b}_{a,n,j,p}, & a = f \end{cases}, \qquad (30)$$

$$\bar{C}_a^{(t)} \triangleq \begin{cases} \sum_{i \in \mathcal{I}} \frac{N_i}{BN} \sum_{n \in \mathcal{N}_i^{(t)}} \bar{c}_{a,n} + \tau \left\|\boldsymbol{\omega}_a^{(t)}\right\|_2^2, & a = s \\ \frac{1}{B} \sum_{n \in \mathcal{N}^{(t)}} \bar{c}_{a,n} + \tau \left\|\boldsymbol{\omega}_a^{(t)}\right\|_2^2, & a = f \end{cases}, \qquad (31)$$

where

$$\bar{a}_{a,n,l,j} \triangleq (Q_l(\boldsymbol{\omega}_a^{(t)}; \mathbf{x}_n) - y_{n,l}) S\left(\sum_{p'=1}^{P} \omega_{a,1,j,p'}^{(t)} x_{n,p'}\right),$$

$$\bar{b}_{a,n,j,p} \triangleq \sum_{l \in \mathcal{L}} (Q_l(\boldsymbol{\omega}_a^{(t)}; \mathbf{x}_n) - y_{n,l}) S'\left(\sum_{p'=1}^{P} \omega_{a,1,j,p'}^{(t)} x_{n,p'}\right) \omega_{a,0,l,j}^{(t)} x_{n,p},$$

$$\bar{c}_{a,n} \triangleq \sum_{l \in \mathcal{L}} y_{n,l} \log(Q_l(\boldsymbol{\omega}_a^{(t)}; \mathbf{x}_n)).$$

### A. Unconstrained Federated Optimization

For $a = s, f$, one unconstrained federated optimization formulation for the $L$-class classification problem is to minimize the weighted sum of the loss function $F(\boldsymbol{\omega})$ in (28) and the $\ell_2$-norm regularization term $\|\boldsymbol{\omega}\|_2^2$:

$$\min_{\boldsymbol{\omega}} \quad F_{a,0}(\boldsymbol{\omega}) \triangleq F(\boldsymbol{\omega}) + \lambda \|\boldsymbol{\omega}\|_2^2 \qquad (32)$$

where $\lambda > 0$ is the regularization parameter that trades off the cost and model sparsity. Obviously, $F(\boldsymbol{\omega}) + \lambda \|\boldsymbol{\omega}\|_2^2$ satisfies the additional restrictions on the structure of $F_{f,0}(\boldsymbol{\omega})$. We can view $-\sum_{l \in \mathcal{L}} y_{n,l} \log(Q_l(\boldsymbol{\omega}; \mathbf{x}_n))$ as $f_{a,0}(\boldsymbol{\omega}; \mathbf{x}_n)$, apply Algorithm 1 with $\bar{f}_{s,0}(\boldsymbol{\omega}; \boldsymbol{\omega}_s^{(t)}, \mathbf{x}_n)$ given by (7) to solve the problem in (32) for $a = s$, and apply Algorithm 3 with $\bar{f}_{f,0}(\boldsymbol{\omega}; \boldsymbol{\omega}_f^{(t)}, \mathbf{x}_n)$ given by (19) to solve the problem in (32) for $a = f$.

First, we present the details of Step 4 in Algorithm 1 and the details of Steps 4-6 in Algorithm 3. In Step 4 of Algorithm 1, each client $i$ computes $((\sum_{n \in \mathcal{N}_i^{(t)}} \bar{a}_{s,n,l,j})_{l \in \mathcal{L}, j \in \mathcal{J}}, (\sum_{n \in \mathcal{N}_i^{(t)}} \bar{b}_{s,n,j,p})_{j \in \mathcal{J}, p \in \mathcal{P}})$ and sends it to the server. In Steps 4-6 of Algorithm 3, each client $i$ computes $(\omega_{f,1,j,p}^{(t)} x_{n,p})_{j \in \mathcal{J}, p \in \mathcal{P}_i}$ $n \in \mathcal{N}^{(t)}$ and sends them to the other clients; based on $(\omega_{f,1,j,p}^{(t)} x_{n,p})_{j \in \mathcal{J}, p \in \mathcal{P}_i}$, $n \in \mathcal{N}^{(t)}$, the client with the highest computation speed (or any client) computes $(\sum_{n \in \mathcal{N}^{(t)}} \bar{a}_{f,n,l,j})_{j \in \mathcal{J}, l \in \mathcal{L}}$ and sends it to the server; each client $i$ computes $(\sum_{n \in \mathcal{N}^{(t)}} \bar{b}_{f,n,j,p})_{j \in \mathcal{J}, p \in \mathcal{P}_i}$ and sends it to the server.

Next, we present the details of Step 5 in Algorithm 1 and the details of Step 7 in Algorithm 3. For $a = s, f$, the convex approximate problem is given by:

$$\min_{\boldsymbol{\omega}} \quad \bar{F}_{a,0}^{(t)}(\boldsymbol{\omega}) = \bar{F}_a^{(t)}(\boldsymbol{\omega}) + 2\lambda(\boldsymbol{\beta}^{(t)})^T \boldsymbol{\omega} \qquad (33)$$

where $\bar{F}_a^{(t)}(\boldsymbol{\omega})$ is given by

$$\bar{F}_a^{(t)}(\boldsymbol{\omega}) = \sum_{l \in \mathcal{L}} \sum_{j \in \mathcal{J}} A_{a,l,j}^{(t)} \omega_{0,l,j} + \sum_{j \in \mathcal{J}} \sum_{p \in \mathcal{P}} B_{a,j,p}^{(t)} \omega_{1,j,p} + \tau \|\boldsymbol{\omega}\|_2^2, \qquad (34)$$

and $\boldsymbol{\beta}^{(t)} \in \mathbb{R}^d$, $A_{a,l,j}^{(t)} \in \mathbb{R}$, and $B_{a,j,p}^{(t)} \in \mathbb{R}$ are updated according to:

$$\boldsymbol{\beta}^{(t)} = (1 - \rho^{(t)})\boldsymbol{\beta}^{(t-1)} + \rho^{(t)}\boldsymbol{\omega}_a^{(t)}, \qquad (35)$$

$$A_{a,l,j}^{(t)} = (1 - \rho^{(t)})A_{a,l,j}^{(t-1)} + \rho^{(t)}\left(\bar{A}_{a,l,j}^{(t)} - 2\tau\omega_{a,0,l,j}^{(t)}\right), \qquad (36)$$

$$B_{a,j,p}^{(t)} = (1 - \rho^{(t)})B_{a,j,p}^{(t)} + \rho^{(t)}\left(\bar{B}_{a,j,p}^{(t)} - 2\tau\omega_{a,1,j,p}^{(t)}\right), \qquad (37)$$

respectively, with $\boldsymbol{\beta}^{(0)} = \mathbf{0}$ and $A_{a,l,j}^{(0)} = B_{a,j,p}^{(0)} = 0$. Here, $\bar{A}_{a,l,j}^{(t)}$ and $\bar{B}_{a,j,p}^{(t)}$ are given by (29) and (30) respectively. By (10) for $a = s$ and (24) for $a = f$, the closed-form solutions of the problem in (33) for $a = s, f$ are given by:

$$\bar{\omega}_{a,0,l,j}^{(t)} = -\frac{1}{2\tau}\left(A_{a,l,j}^{(t)} + 2\lambda\beta_{2,l,j}^{(t)}\right), \ l \in \mathcal{L}, \ j \in \mathcal{J}, \quad (38)$$

$$\bar{\omega}_{a,1,j,p}^{(t)} = -\frac{1}{2\tau}\left(B_{a,j,p}^{(t)} + 2\lambda\beta_{1,j,p}^{(t)}\right), \ j \in \mathcal{J}, \ p \in \mathcal{P}. \quad (39)$$

Thus, in Step 5 in Algorithm 1 and Step 7 in Algorithm 3, the server only needs to compute $\bar{\omega}_a^{(t)}$ according to (38) and (39).

Theorem 1 and Theorem 3 guarantee the convergences of Algorithm 1 and Algorithm 3, respectively, as Assumption 1 and Assumption 2 are satisfied.

### B. Constrained Federated Optimization

For $a = s, f$, one constrained federated optimization formulation for the $L$-class classification problem is to minimize the $\ell_2$-norm of the network parameters $\|\boldsymbol{\omega}\|_2^2$ under a constraint on the loss function $F(\boldsymbol{\omega})$ in (28):

$$\min_{\boldsymbol{\omega}} \quad F_{a,0}(\boldsymbol{\omega}) \triangleq \|\boldsymbol{\omega}\|_2^2 \qquad (40)$$

$$\text{s.t.} \quad F_{a,1}(\boldsymbol{\omega}) \triangleq F(\boldsymbol{\omega}) - U \le 0,$$

where $U$ represents the limit on the cost. We can view $0$ and $-\sum_{l \in \mathcal{L}} y_{n,l} \log(Q_l(\boldsymbol{\omega}; \mathbf{x}_n))$ as $f_{a,0}(\boldsymbol{\omega}; \mathbf{x}_n)$ and $f_{a,1}(\boldsymbol{\omega}; \mathbf{x}_n)$, respectively. Then, we can apply Algorithm 2 with $\bar{f}_{s,0}(\boldsymbol{\omega}; \boldsymbol{\omega}_s^{(t)}, \mathbf{x}_n)$ given by (7) and $\bar{f}_{s,1}(\boldsymbol{\omega}; \boldsymbol{\omega}_s^{(t)}, \mathbf{x}_n)$ given by (15) to solve the problem in (40) for $a = s$ and ap-

ply Algorithm 4 with $\bar{f}_{f,0}(\boldsymbol{\omega}; \boldsymbol{\omega}_f^{(t)}, \mathbf{x}_n)$ given by (19) and $\bar{f}_{f,1}(\boldsymbol{\omega}; \boldsymbol{\omega}_f^{(t)}, \mathbf{x}_n)$ given by (27) to solve the problem in (40) for $a = f$.

First, we present the details of Step 4 in Algorithm 2 and the details of Steps 4-6 in Algorithm 4. In Step 4 of Algorithm 2, each client $i$ computes $((\sum_{n \in \mathcal{N}_i^{(t)}} \bar{a}_{s,n,l,j})_{l \in \mathcal{L}, j \in \mathcal{J}}, (\sum_{n \in \mathcal{N}_i^{(t)}} \bar{b}_{s,n,j,p})_{j \in \mathcal{J}, p \in \mathcal{P}})$ and $\sum_{n \in \mathcal{N}_i^{(t)}} \bar{c}_{s,n}$ and sends them to the server. In Steps 4-6 of Algorithm 4, each client $i$ computes $(\omega_{f,1,j,p}^{(t)} x_{n,p})_{j \in \mathcal{J}, p \in \mathcal{P}_i}$, $n \in \mathcal{N}^{(t)}$ and sends them to the other clients; based on $(\omega_{f,1,j,p}^{(t)} x_{n,p})_{j \in \mathcal{J}, p \in \mathcal{P}_i}$, $n \in \mathcal{N}^{(t)}$, the client with the highest computation speed (or any client) computes $(\sum_{n \in \mathcal{N}^{(t)}} \bar{a}_{f,n,l,j})_{l \in \mathcal{L}, j \in \mathcal{J}}$ and $\sum_{n \in \mathcal{N}^{(t)}} \bar{c}_{f,n}$ and sends them to the server; each client $i$ computes $(\sum_{n \in \mathcal{N}^{(t)}} \bar{b}_{f,n,j,p})_{j \in \mathcal{J}, p \in \mathcal{P}_i}$ and sends it to the server.

Next, we present the details of Step 5 in Algorithm 2 and the details of Step 7 in Algorithm 4. For $a = s, f$, the convex approximate problem is given by:

$$\min_{\boldsymbol{\omega}, s} \quad \|\boldsymbol{\omega}\|_2^2 + cs \tag{41}$$
$$\text{s.t.} \quad \bar{F}_a^{(t)}(\boldsymbol{\omega}) + C_a^{(t)} - U \leq s,$$
$$s \geq 0,$$

where $\bar{F}_a^{(t)}(\boldsymbol{\omega})$ is given by (34) with $A_{a,l,j}^{(t)}$, $B_{a,j,p}^{(t)}$, and $C_a^{(t)}$ updated according to (36), (37), and

$$C_a^{(t)} = (1 - \rho^{(t)}) C_a^{(t-1)} +$$
$$\rho^{(t)} \left( \bar{C}_a^{(t)} - \sum_{l \in \mathcal{L}} \sum_{j \in \mathcal{J}} \bar{A}_{a,l,j}^{(t)} \omega_{a,0,l,j}^{(t)} - \sum_{j \in \mathcal{J}} \sum_{p \in \mathcal{P}} \bar{B}_{a,j,p}^{(t)} \omega_{a,1,j,p}^{(t)} \right), \tag{42}$$

respectively, with $C_a^{(0)} = 0$ and $\bar{C}_a^{(t)}$ given by (31). By the KKT conditions, the closed-form solutions of the problem in (41) for $a = s, f$ are given as follows.

*Lemma 1 (Optimal Solution of Problem in (41)):*

$$\bar{\omega}_{a,0,l,j}^{(t)} = -\frac{\nu A_{a,l,j}^{(t)}}{2(1 + \nu\tau)}, \quad l \in \mathcal{L}, \ j \in \mathcal{J}, \tag{43}$$

$$\bar{\omega}_{a,1,j,p}^{(t)} = -\frac{\nu B_{a,j,p}^{(t)}}{2(1 + \nu\tau)}, \quad j \in \mathcal{J}, \ p \in \mathcal{P}, \tag{44}$$

where

$$\nu = \begin{cases} \left[ \frac{1}{\tau} \left( \sqrt{\frac{b}{b + 4\tau(U - C_a^{(t)})}} - 1 \right) \right]_0^c, & b + 4\tau(U - C_a^{(t)}) > 0 \\ c, & b + 4\tau(U - C_a^{(t)}) \leq 0 \end{cases},$$

$$b = \sum_{l \in \mathcal{L}} \sum_{j \in \mathcal{J}} (A_{a,l,j}^{(t)})^2 + \sum_{j \in \mathcal{J}} \sum_{p \in \mathcal{P}} (B_{a,j,p}^{(t)})^2. \tag{45}$$

Here, $[x]_0^c \triangleq \min \{\max\{x, 0\}, c\}$.

*Proof:* Please refer to Appendix C. ∎

Thus, in Step 5 of Algorithm 2 and Step 7 of Algorithm 4, the server only needs to compute $\bar{\boldsymbol{\omega}}_a^{(t)}$ according to (43) and (44).

The convergences of Algorithm 2 and Algorithm 4 are guaranteed by Theorem 2 and Theorem 4, respectively, as Assumption 1 and Assumption 2 are satisfied.

## C. Comparisons of Two Formulations

Both the unconstrained federated optimization formulation in (32) and constrained federated optimization formulation in (40) allow tradeoffs between the cost and model sparsity [30]. The equivalence between the two formulations is summarized in the following theorem.

*Theorem 5 (Equivalence between Problems in (32) and (40)):* i) If $\boldsymbol{\omega}^*$ is a locally optimal solution of the problem in (32) with $\lambda > 0$, then there exists $U \geq 0$ such that $\boldsymbol{\omega}^*$ is a locally optimal solution of the problem in (40). ii) If $\boldsymbol{\omega}^\dagger$ is a locally optimal solution of the problem in (40) with $U > 0$, which is regular and satisfies the KKT conditions together with a corresponding Lagrange multiplier $\xi > 0$, then there exists $\lambda > 0$ such that $\boldsymbol{\omega}^\dagger$ is a stationary point of the problem in (32). If, in addition, $\lambda$ and $\boldsymbol{\omega}^\dagger$ satisfy $\nabla^2 F(\boldsymbol{\omega}^\dagger) + \lambda I \succeq 0$, then $\boldsymbol{\omega}^\dagger$ is a locally optimal solution of the problem in (32).

*Proof:* Please refer to Appendix D. ∎

By the above theorem, we know that the problem in (32) and the problem in (40) have the same locally optimal solution for certain $\lambda$ and $U$ under some conditions. Besides, we can tradeoff between the training accuracy and model sparsity of each formulation. It is evident that with the constrained federated optimization formulation in (40), one can set an explicit constraint on the training cost to control the test accuracy effectively.
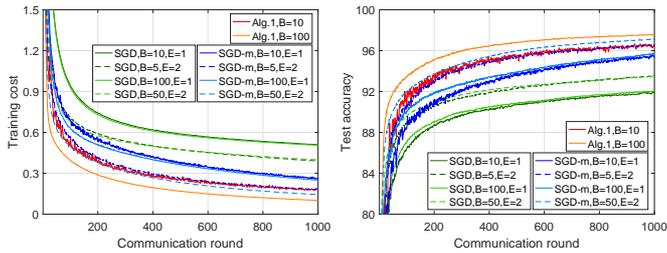
## VI. NUMERICAL RESULTS

In this section, we numerically evaluate the proposed examples of Algorithms 1-4 using the application examples in Section V.[17] For unconstrained federated optimization, we adopt the existing SGD-based [5], [6], [13] and momentum SGD-based [7] FL algorithms, called SGD and SGD-m, respectively, as the baseline algorithms for the proposed examples of Algorithm 1 and Algorithm 3. Let $E$ denote the number of local SGD (momentum SGD) updates for sample-based SGD (SGD-m). Note that sample-based SGD with $B \times E = N$ becomes FedAvg [5]. Feature-based SGD and SGD-m adopt the information collection mechanism used in Algorithms 3 (i.e., the extension of the one in [13]). In each communication round, each proposed algorithm executes one iteration, each sample-based SGD (SGD-m) executes one global iteration and $E$ local SGD (momentum SGD) updates, and each feature-based SGD (SGD-m) executes one global iteration. Algorithm 1 (Algorithm 3) and its baseline algorithms have the same communication load per communication round. Besides, if the value of $B$ for Algorithm 1 (Algorithm 3) and the value of $B \times E$ for each sample-based (B for each feature-based) baseline algorithm are equal, the two algorithms have the same order of computational complexity per communication round.[18]

We set $\lambda = 10^{-5}$ and $U = 0.13$ for the unconstrained and constrained federated optimization problems in (32) and
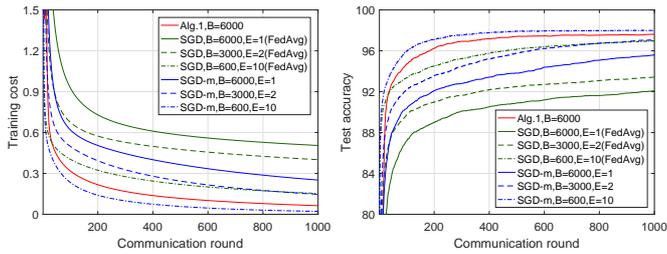
---

[17] Source code for the experiments is available at [31].

[18] The example of Algorithm 1 (Algorithm 3) has the same level of privacy protection as its baseline algorithms, as illustrated in Section III-A (Section IV-A).
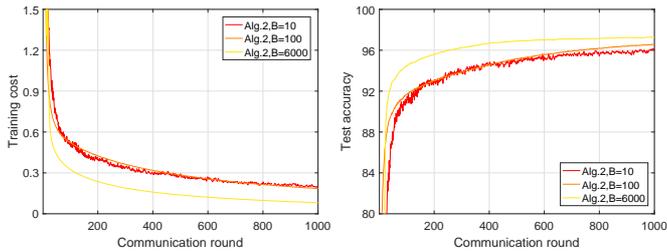
(40), respectively, unless otherwise specified. We carry our experiments on Mnist dataset. For the training model, we set $N = 60000$, $I = 10$, $K = 784$, $J = 128$, and $L = 10$. For the proposed algorithms, we choose $T = 1000$, $c = 10^5$, $\rho^{(t)} = a_1/t^\alpha$ and $\gamma^{(t)} = a_2/t^\alpha$ with $a_1 = 0.9, 0.3, 0.2$, $a_2 = 0.5, 0.3, 0.3$, $\alpha = 0.1, 0.1, 0.1$, and $\tau = 0.2, 0.05, 0.03$ for batch sizes $B = 10, 100, 6000$ in sample-based FL and $a_1 = 0.9, 0.9, 0.3$, $a_2 = 0.3, 0.5, 0.3$, $\alpha = 0.3, 0.1, 0.1$, and $\tau = 0.1, 0.2, 0.05$ for batch sizes $B = 10, 100, 1000$ in feature-based FL. For SGD, the learning rate is set as $r = \bar{a}/t^{\bar{\alpha}}$ with $\bar{a} = 0.3$ and $\bar{\alpha} = 0.3$. For SGD-m, the learning rate is set as $r = \bar{a}$ with $\bar{a} = 0.3$ and the momentum parameter is set as $\bar{\beta} = 0.1$. Note that all the algorithm parameters are selected using a grid search method, and all the results are given by averaging over ten runs.



(a) Training cost for unconstrained federated optimization at $B, B \times E = 10, 100$.

(b) Test accuracy for unconstrained federated optimization at $B, B \times E = 10, 100$.
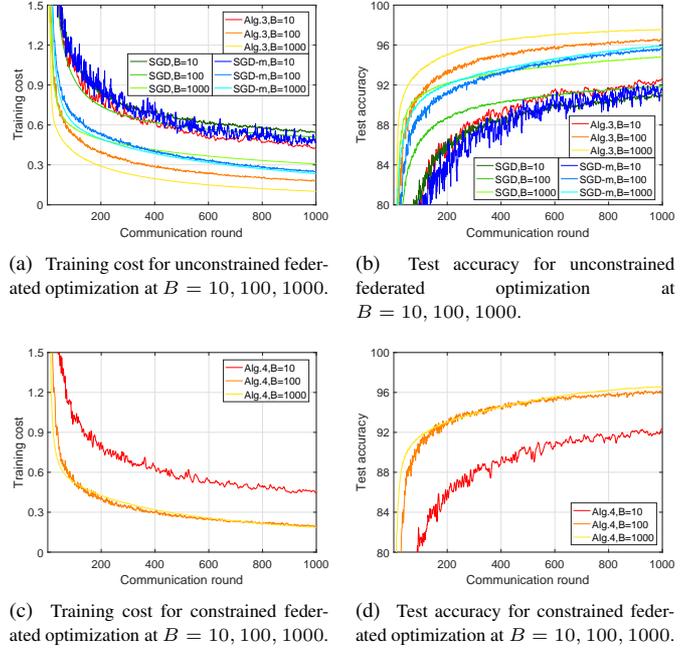
(c) Training cost for unconstrained federated optimization at $B, B \times E = 6000$.

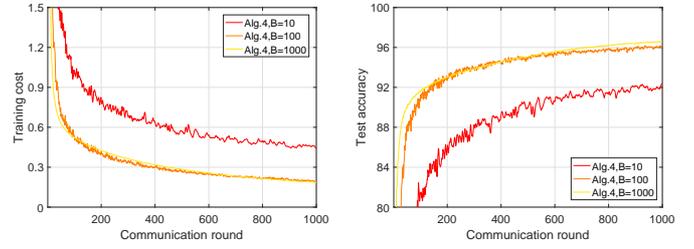(d) Test accuracy for unconstrained federated optimization at $B, B \times E = 6000$.

(e) Training cost for constrained federated optimization at $B = 10, 100, 6000$.

(f) Test accuracy for constrained federated optimization at $B = 10, 100, 6000$.

Fig. 1. Training cost $F(\boldsymbol{\omega}_s^{(t)})$ and test accuracy at $\boldsymbol{\omega}_s^{(t)}$ versus communication round index $t$ for sample-based FL.

Fig. 1 and Fig. 2 illustrate the training cost and test accuracy versus the communication round index in sample-based FL and feature-based FL, respectively. From Fig. 1 (a), (c), (e) and Fig. 2 (a), (c), we can see that each proposed algorithm with larger $B$, sample-based SGD (SGD-m) with larger $B \times$
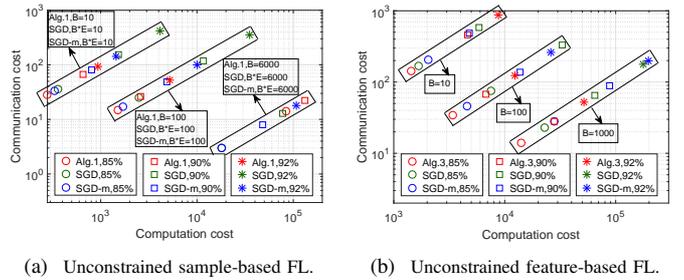


(a) Training cost for unconstrained federated optimization at $B = 10, 100, 1000$.

(b) Test accuracy for unconstrained federated optimization at $B = 10, 100, 1000$.

(c) Training cost for constrained federated optimization at $B = 10, 100, 1000$.

(d) Test accuracy for constrained federated optimization at $B = 10, 100, 1000$.

Fig. 2. Training cost $F(\boldsymbol{\omega}_f^{(t)})$ and test accuracy at $\boldsymbol{\omega}_f^{(t)}$ versus communication round index $t$ for feature-based FL.



(a) Unconstrained sample-based FL.

(b) Unconstrained feature-based FL.

Fig. 3. Tradeoff between communication cost and computation cost for solving unconstrained federated optimization with a specific test accuracy.

$E$, and feature-based SGD (SGD-m) with larger $B$ converge faster at higher computation costs per communication round. We can also observe that Algorithm 1 (Algorithm 3) converges faster than all the baseline algorithms with the same order of computational complexity per communication round in most (all) cases. The only exception for Algorithm 1 is that in Fig. 1 (c), Algorithm 1 with $B = 6000$ converges slightly slower than sample-based SGD-m with $B = 600$ and $E = 10$.

Fig. 3 shows the tradeoff between the communication and computation costs for solving unconstrained federated optimization. Here, the communication cost of each algorithm is measured by the number of communication rounds, the computation costs of Algorithm 1, Algorithm 3, and feature-based SGD (SGD-m) are measured by $B$, and the communication cost of sample-based SGD (SGD-m) is measured by $B \times E$. From Fig. 3(a), we see that the proposed algorithms achieve the best tradeoff between the communication cost

(a) $\ell_2$-norm $\|\boldsymbol{\omega}_s^{(T)}\|_2^2$ vs. training cost obtained by Algorithm 1.

(b) $\ell_2$-norm $\|\boldsymbol{\omega}_s^{(T)}\|_2^2$ vs. training cost obtained by Algorithm 2.
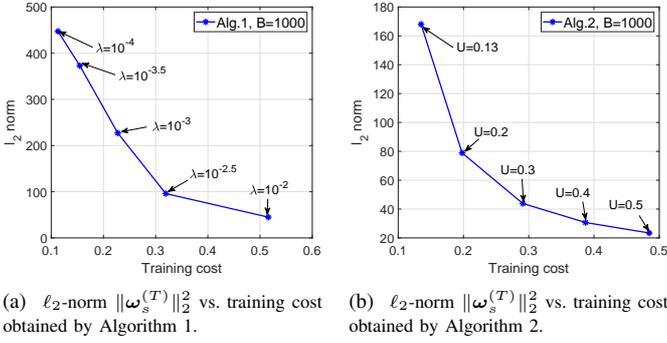
Fig. 4. Tradeoff between model sparsity and training cost for sample-based FL at $T = 1000$.

and computation cost in all cases except the case where all local samples are utilized per communication round for solving for sample-based FL. Thus, Fig. 3(a) indicates that Algorithm 1 (Algorithm 3) achieves the lowest communication and computation costs for reaching a specific convergence performance in most (all) cases.

Fig. 4 shows the tradeoff curve between the model sparsity and training cost of each proposed algorithm for sample-based FL. From Fig. 4(b), we see that with constrained sample-based federated optimization, one can set an explicit constraint on the training cost to control the test accuracy effectively.

## VII. CONCLUSIONS

In this paper, we investigated sample-based and feature-based federated optimization, respectively, and considered both the unconstrained problem and constrained problem for each of them. We proposed FL algorithms that converge to stationary points or KKT points using SSCA and mini-batch techniques. We also provided algorithm examples that have appealing computational complexities and communication loads per communication round and connect to FL algorithms via momentum SGD. Numerical experiments demonstrated that the proposed mini-batch SSCA-based FL algorithms for unconstrained sample-based and feature-based federated optimization generally converge faster than existing FL algorithms, and the proposed mini-batch SSCA-based FL algorithms for constrained sample-based and feature-based federated optimization problems obtain models that strictly satisfy nonconvex constraints. To the best of our knowledge, this is the first work that provides an SSCA framework for federated optimization, highlights the value of constrained federated optimization, and establishes an analytical connection between SSCA and momentum SGD. This paper opens up several directions for future research. An important direction is to design advanced SSCA-based FL algorithms that allow multiple local updates to reduce communication costs further. Another interesting direction is to design more privacy mechanisms for SSCA-based FL algorithms.

## APPENDIX A: PROOFS OF THEOREM 1 AND THEOREM 3

The proofs of Theorem 1 and Theorem 3 are identical. In the following proof, we omit the subscripts $s, f$ for notation simplicity. First, we introduce the following preliminary results.

*Lemma 2:* Let $\{\boldsymbol{\omega}^{(t)}\}$ be the sequence generated by Algorithm 1 (Algorithm 3). Then, we have:

$$\lim_{t \to \infty} \left\| \nabla \bar{F}_0^{(t)}(\boldsymbol{\omega}^{(t)}) - \nabla F_0(\boldsymbol{\omega}^{(t)}) \right\|_2 = 0,$$

$$\lim_{t \to \infty} \left| \bar{F}_0^{(t)}(\boldsymbol{\omega}) - G_0(\boldsymbol{\omega}; \boldsymbol{\omega}^{(t)}) \right| = 0, \quad \boldsymbol{\omega} \in \mathbb{R}^d,$$

almost surely, where $G_0(\boldsymbol{\omega}; \boldsymbol{\omega}^{(t)}) \triangleq \frac{1}{N} \sum_{n \in \mathcal{N}} \bar{f}_0(\boldsymbol{\omega}; \boldsymbol{\omega}^{(t)}, \mathbf{x}_n)$.

*Proof:* Lemma 2 is a consequence of [32, Lemma 1]. We only need to verify that all the technical conditions therein are satisfied. Specifically, Condition (a) of [32, Lemma 1] is satisfied because $\{\boldsymbol{\omega}^{(t)}\}$ is assumed to be bounded. Condition (b) of [32, Lemma 1] comes from Assumption 2.4. Conditions (c)-(d) of [32, Lemma 1] come from the stepsize rules in (4) and (6). Condition (e) of [32, Lemma 1] comes from the Lipschitz property of $F_0(\boldsymbol{\omega})$ from Assumption 1 and the stepsize rule in (6). ∎

*Lemma 3:* Let $\{\boldsymbol{\omega}^{(t)}\}$ be the sequence generated by Algorithm 1 (Algorithm 3). Then, there exists a constant $\bar{L}$ such that

$$\left\| \bar{\boldsymbol{\omega}}^{(t_1)} - \bar{\boldsymbol{\omega}}^{(t_2)} \right\|_2 \leq \bar{L} \left\| \boldsymbol{\omega}^{(t_1)} - \boldsymbol{\omega}^{(t_2)} \right\|_2 + e(t_1, t_2), \quad (46)$$

and $\lim_{t_1, t_2 \to \infty} e(t_1, t_2) = 0$ almost surely.

*Proof:* It follows from Lemma 2 that

$$\bar{F}_0^{(t)}(\boldsymbol{\omega}) = G_0(\boldsymbol{\omega}; \boldsymbol{\omega}^{(t)}) + \bar{e}(t), \quad (47)$$

where $\bar{e}(t)$ satisfies $\lim_{t \to \infty} \bar{e}(t) = 0$. From Assumption 2.3, $G_0(\boldsymbol{\omega}; \boldsymbol{\omega}^{(t)})$ is Lipschitz continuous in $\boldsymbol{\omega}^{(t)}$ and thus

$$\left| G_0(\boldsymbol{\omega}; \boldsymbol{\omega}^{(t_1)}) - G_0(\boldsymbol{\omega}; \boldsymbol{\omega}^{(t_2)}) \right| \leq \tilde{L} \left\| \boldsymbol{\omega}^{(t_1)} - \boldsymbol{\omega}^{(t_2)} \right\|_2, \boldsymbol{\omega} \in \mathbb{R}^d, \quad (48)$$

for some constant $\tilde{L} > 0$. Combining (47) and (48), we have:

$$\left| \bar{F}_0^{(t_1)}(\boldsymbol{\omega}) - \bar{F}_0^{(t_2)}(\boldsymbol{\omega}) \right| \leq \tilde{L} \left\| \boldsymbol{\omega}^{(t_1)} - \boldsymbol{\omega}^{(t_2)} \right\|_2 + \tilde{e}(t_1, t_2), \boldsymbol{\omega} \in \mathbb{R}^d, \quad (49)$$

where $\tilde{e}(t_1, t_2)$ satisfies $\lim_{t_1, t_2 \to \infty} \tilde{e}(t_1, t_2) = 0$. From Assumption 2.3, there exists constant $\mu > 0$ such that for all $t = 1, 2, \cdots, \infty$, $\bar{F}_0^{(t)}(\boldsymbol{\omega})$ is strongly convex with $\mu$. Due to the strong convexity of $\bar{F}_0^{(t_1)}(\boldsymbol{\omega})$ and the optimality of $\bar{\boldsymbol{\omega}}^{(t_1)}$, we have:

$$\bar{F}_0^{(t_1)}(\boldsymbol{\omega}) - \bar{F}_0^{(t_1)}(\bar{\boldsymbol{\omega}}^{(t_1)}) \geq \frac{\mu}{2} \left\| \boldsymbol{\omega} - \bar{\boldsymbol{\omega}}^{(t_1)} \right\|_2, \quad \boldsymbol{\omega} \in \mathbb{R}^d. \quad (50)$$

Setting $\boldsymbol{\omega} = \bar{\boldsymbol{\omega}}^{(t_2)}$ in (50), we have:

$$\bar{F}_0^{(t_1)}(\bar{\boldsymbol{\omega}}^{(t_2)}) - \bar{F}_0^{(t_1)}(\bar{\boldsymbol{\omega}}^{(t_1)}) \geq \frac{\mu}{2} \left\| \bar{\boldsymbol{\omega}}^{(t_2)} - \bar{\boldsymbol{\omega}}^{(t_1)} \right\|_2. \quad (51)$$

Similarly, by the strong convexity of $\bar{F}_0^{(t_2)}(\boldsymbol{\omega})$ and the optimality of $\bar{\boldsymbol{\omega}}^{(t_2)}$, we have:

$$\bar{F}_0^{(t_2)}(\bar{\boldsymbol{\omega}}^{(t_1)}) - \bar{F}_0^{(t_2)}(\bar{\boldsymbol{\omega}}^{(t_2)}) \geq \frac{\mu}{2} \left\| \bar{\boldsymbol{\omega}}^{(t_1)} - \bar{\boldsymbol{\omega}}^{(t_2)} \right\|_2. \quad (52)$$

Thus, we have:

$$\left\| \bar{\boldsymbol{\omega}}^{(t_1)} - \bar{\boldsymbol{\omega}}^{(t_2)} \right\|_2$$

$$\overset{(a)}{\leq}\frac{1}{\mu}\Big(\Big|\bar{F}_0^{(t_1)}(\bar{\boldsymbol{\omega}}^{(t_1)})-\bar{F}_0^{(t_2)}(\bar{\boldsymbol{\omega}}^{(t_1)})\Big|+\Big|\bar{F}_0^{(t_1)}(\bar{\boldsymbol{\omega}}^{(t_2)})-\bar{F}_0^{(t_2)}(\bar{\boldsymbol{\omega}}^{(t_2)})\Big|\Big)$$

$$\overset{(b)}{\leq}\frac{2\tilde{L}}{\mu}\left\|\boldsymbol{\omega}^{(t_1)}-\boldsymbol{\omega}^{(t_2)}\right\|_2+\frac{2}{\mu}\tilde{e}(t_1,t_2),\tag{53}$$

where $(a)$ follows from (51) and (52), and $(b)$ follows from (49). Finally, (46) follows from (53) immediately. ∎

*Lemma 4:* Let $\{\boldsymbol{\omega}^{(t)}\}$ be the sequence generated by Algorithm 1 (Algorithm 3). Then, we have:

$$F_0(\boldsymbol{\omega}^{(t+1)})-F_0(\boldsymbol{\omega}^{(t)})$$
$$\leq\gamma^{(t)}\left\|\bar{\boldsymbol{\omega}}^{(t)}-\boldsymbol{\omega}^{(t)}\right\|_2\left\|\nabla F_0(\boldsymbol{\omega}^{(t)})-\nabla\bar{F}_0^{(t)}(\boldsymbol{\omega}^{(t)})\right\|_2$$
$$-\gamma^{(t)}\left(\mu-\frac{\hat{L}}{2}\gamma^{(t)}\right)\left\|\bar{\boldsymbol{\omega}}^{(t)}-\boldsymbol{\omega}^{(t)}\right\|_2^2.\tag{54}$$

*Proof:* From Assumption 2.2, $\bar{F}_0^{(t)}(\boldsymbol{\omega})$ is uniformly strongly convex, and thus:

$$(\bar{\boldsymbol{\omega}}^{(t)}-\boldsymbol{\omega}^{(t)})^T\nabla\bar{F}_0^{(t)}(\boldsymbol{\omega}^{(t)})$$
$$\leq-\mu\left\|\bar{\boldsymbol{\omega}}^{(t)}-\boldsymbol{\omega}^{(t)}\right\|_2+\bar{F}_0^{(t)}(\bar{\boldsymbol{\omega}}^{(t)})-\bar{F}_0^{(t)}(\boldsymbol{\omega}^{(t)})$$
$$\leq-\mu\left\|\bar{\boldsymbol{\omega}}^{(t)}-\boldsymbol{\omega}^{(t)}\right\|_2,\tag{55}$$

where the last inequality follows from the optimality of $\bar{\boldsymbol{\omega}}^{(t)}$. Suppose $\nabla F_0(\boldsymbol{\omega})$ is Lipschitz continuous with constant $\hat{L}>0$, we have:

$$F_0(\boldsymbol{\omega}^{(t+1)})-F_0(\boldsymbol{\omega}^{(t)})$$
$$\leq(\boldsymbol{\omega}^{(t+1)}-\boldsymbol{\omega}^{(t)})^T\nabla F_0(\boldsymbol{\omega}^{(t)})+\frac{\hat{L}}{2}\left\|\boldsymbol{\omega}^{(t+1)}-\boldsymbol{\omega}^{(t)}\right\|_2^2$$
$$\leq\gamma^{(t)}(\bar{\boldsymbol{\omega}}^{(t)}-\boldsymbol{\omega}^{(t)})^T\nabla F_0(\boldsymbol{\omega}^{(t)})+\frac{\hat{L}}{2}(\gamma^{(t)})^2\left\|\bar{\boldsymbol{\omega}}^{(t)}-\boldsymbol{\omega}^{(t)}\right\|_2^2$$
$$\leq\gamma^{(t)}(\bar{\boldsymbol{\omega}}^{(t)}-\boldsymbol{\omega}^{(t)})^T\Big(\nabla F_0(\boldsymbol{\omega}^{(t)})-\nabla\bar{F}_0^{(t)}(\boldsymbol{\omega}^{(t)})+\nabla\bar{F}_0^{(t)}(\boldsymbol{\omega}^{(t)})\Big)$$
$$+\frac{\hat{L}}{2}(\gamma^{(t)})^2\left\|\bar{\boldsymbol{\omega}}^{(t)}-\boldsymbol{\omega}^{(t)}\right\|_2^2$$
$$\leq\gamma^{(t)}\left\|\bar{\boldsymbol{\omega}}^{(t)}-\boldsymbol{\omega}^{(t)}\right\|_2\left\|\nabla F_0(\boldsymbol{\omega}^{(t)})-\nabla\bar{F}_0^{(t)}(\boldsymbol{\omega}^{(t)})\right\|_2$$
$$-\gamma^{(t)}\left(\mu-\frac{\hat{L}}{2}\gamma^{(t)}\right)\left\|\bar{\boldsymbol{\omega}}^{(t)}-\boldsymbol{\omega}^{(t)}\right\|_2^2,\tag{56}$$

where the last inequality follows form (55). ∎

Then, we show by contradiction that $\liminf_{t\to\infty}\left\|\bar{\boldsymbol{\omega}}^{(t)}-\boldsymbol{\omega}^{(t)}\right\|_2=0$ almost surely. Suppose $\liminf_{t\to\infty}\left\|\bar{\boldsymbol{\omega}}^{(t)}-\boldsymbol{\omega}^{(t)}\right\|_2\geq\chi>0$ with a positive probability. Then we can find a realization such that $\left\|\bar{\boldsymbol{\omega}}^{(t)}-\boldsymbol{\omega}^{(t)}\right\|_2\geq\chi>0$ for all $t$. We focus next on such a realization. By $\left\|\bar{\boldsymbol{\omega}}^{(t)}-\boldsymbol{\omega}^{(t)}\right\|_2\geq\chi>0$ and Lemma 4, we have:

$$F_0(\boldsymbol{\omega}^{(t+1)})-F_0(\boldsymbol{\omega}^{(t)})$$
$$\leq-\gamma^{(t)}\left(\mu-\frac{\hat{L}}{2}\gamma^{(t)}-\frac{1}{\chi}\left\|\nabla F_0(\boldsymbol{\omega}^{(t)})-\nabla\bar{F}_0^{(t)}(\boldsymbol{\omega}^{(t)})\right\|_2\right)$$
$$\times\left\|\bar{\boldsymbol{\omega}}^{(t)}-\boldsymbol{\omega}^{(t)}\right\|_2^2.\tag{57}$$

Since $\lim_{t\to\infty}\left\|\nabla\bar{F}_0^{(t)}(\boldsymbol{\omega}^{(t)})-\nabla F_0(\boldsymbol{\omega}^{(t)})\right\|_2=0$, $\lim_{t\to\infty}\gamma^{(t)}$ and $\mu>0$, there exists a $t_0$ sufficiently large such that

$$\mu-\frac{\hat{L}}{2}\gamma^{(t)}-\frac{1}{\chi}\left\|\nabla F_0(\boldsymbol{\omega}^{(t)})-\nabla\bar{F}_0^{(t)}(\boldsymbol{\omega}^{(t)})\right\|_2\geq\bar{\mu},\ \forall t\geq t_0,\tag{58}$$

for some $\bar{\mu}\in(0,\mu)$. Therefore, it follows from (57), (58) and $\left\|\bar{\boldsymbol{\omega}}^{(t)}-\boldsymbol{\omega}^{(t)}\right\|_2\geq\chi$ for all $t$ that

$$F_0(\boldsymbol{\omega}^{(t)})-F_0(\boldsymbol{\omega}^{(t_0)})\leq-\bar{\mu}\chi^2\sum_{n=t_0}^{(t)}\gamma^{(t)},\tag{59}$$

which, in view of $\sum_{n=t_0}^{\infty}\gamma^{(t)}=\infty$, contradicts the boundedness of $\{F_0(\boldsymbol{\omega}^{(t)})\}$. Therefore, it must be $\liminf_{t\to\infty}\left\|\bar{\boldsymbol{\omega}}^{(t)}-\boldsymbol{\omega}^{(t)}\right\|_2=0$ almost surely.

Next, we show by contradiction that $\limsup_{t\to\infty}\left\|\bar{\boldsymbol{\omega}}^{(t)}-\boldsymbol{\omega}^{(t)}\right\|_2=0$ almost surely. Suppose $\limsup_{t\to\infty}\left\|\bar{\boldsymbol{\omega}}^{(t)}-\boldsymbol{\omega}^{(t)}\right\|_2>0$ with a positive probability. We focus next on a realization along with $\limsup_{t\to\infty}\left\|\bar{\boldsymbol{\omega}}^{(t)}-\boldsymbol{\omega}^{(t)}\right\|_2>0$, $\liminf_{t\to\infty}\left\|\bar{\boldsymbol{\omega}}^{(t)}-\boldsymbol{\omega}^{(t)}\right\|_2=0$, and $\lim_{t_1,t_2\to\infty}e(t_1,t_2)=0$, where $e(t_1,t_2)$ is defined in Lemma 3. It follows from $\limsup_{t\to\infty}\left\|\bar{\boldsymbol{\omega}}^{(t)}-\boldsymbol{\omega}^{(t)}\right\|_2>0$ and $\liminf_{t\to\infty}\left\|\bar{\boldsymbol{\omega}}^{(t)}-\boldsymbol{\omega}^{(t)}\right\|_2=0$ that there exists a $\delta>0$ such that $\left\|\Delta\boldsymbol{\omega}^{(t)}\right\|_2\geq2\delta$ (with $\Delta\boldsymbol{\omega}^{(t)}\triangleq\bar{\boldsymbol{\omega}}^{(t)}-\boldsymbol{\omega}$) for infinitely many $t$ and also $\left\|\Delta\boldsymbol{\omega}^{(t)}\right\|_2\leq\delta$ for infinitely many $t$. Therefore, one can always find an infinite set of indices, say $\mathcal{T}$, having the following properties: for any $t\in\mathcal{T}$, we have:

$$\left\|\Delta\boldsymbol{\omega}^{(t)}\right\|_2\leq\delta,\tag{60}$$

and there exists an integer $i_t>t$ such that

$$\left\|\Delta\boldsymbol{\omega}^{(i_t)}\right\|_2\geq2\delta,\ \delta\leq\left\|\Delta\boldsymbol{\omega}^{(n)}\right\|_2\leq2\delta,\ t<n<i_t.\tag{61}$$

Thus, for all $t\in\mathcal{T}$, we have:

$$\delta\leq\left\|\Delta\boldsymbol{\omega}^{(i_t)}\right\|_2-\left\|\Delta\boldsymbol{\omega}^{(t)}\right\|_2\leq\left\|\Delta\boldsymbol{\omega}^{(i_t)}-\Delta\boldsymbol{\omega}^{(t)}\right\|_2$$
$$=\left\|(\bar{\boldsymbol{\omega}}^{(i_t)}-\boldsymbol{\omega}^{(i_t)})-(\bar{\boldsymbol{\omega}}^{(t)}-\boldsymbol{\omega}^{(t)})\right\|_2$$
$$\leq\left\|\bar{\boldsymbol{\omega}}^{(i_t)}-\bar{\boldsymbol{\omega}}^{(t)}\right\|_2+\left\|\boldsymbol{\omega}^{(i_t)}-\boldsymbol{\omega}^{(t)}\right\|_2$$
$$\overset{(a)}{\leq}(1+\bar{L})\left\|\boldsymbol{\omega}^{(i_t)}-\boldsymbol{\omega}^{(t)}\right\|_2+e(i_t,t)$$
$$\overset{(b)}{\leq}(1+\bar{L})\sum_{n=t}^{i_t-1}\gamma^{(n)}\left\|\Delta\boldsymbol{\omega}^{(n)}\right\|_2+e(i_t,t)$$
$$\leq2\delta(1+\bar{L})\sum_{n=t}^{i_t-1}\gamma^{(n)}+e(i_t,t),\tag{62}$$

where $(a)$ is due to Lemma 3, and $(b)$ is due to (60) and (61). By (62) and $\lim_{t\to\infty}e(i_t,t)=0$, we have:

$$\liminf_{\mathcal{T}\ni t\to\infty}\sum_{n=t}^{i_t-1}\gamma^{(t)}\geq\delta_1\triangleq\frac{1}{2(1+\bar{L})}>0.\tag{63}$$

Proceeding as in (62), for all $t\in\mathcal{T}$, we also have:

$$\left\|\Delta\boldsymbol{\omega}^{(t+1)}\right\|-\left\|\Delta\boldsymbol{\omega}^{(t)}\right\|_2\leq\left\|\Delta\boldsymbol{\omega}^{(t+1)}-\Delta\boldsymbol{\omega}^{(t)}\right\|_2$$
$$\leq(1+\bar{L})\gamma^{(t)}\left\|\Delta\boldsymbol{\omega}^{(t)}\right\|_2+e(t,t+1),\tag{64}$$

which leads to

$$(1+(1+\bar{L})\gamma^{(t)})\left\|\Delta\boldsymbol{\omega}^{(t)}\right\|_2 + e(t,t+1) \geq \left\|\Delta\boldsymbol{\omega}^{(t+1)}\right\|_2 \geq \delta, \tag{65}$$

where the second inequality follows from (61). It follows from (65) and $\lim_{t\to\infty} e(t,t+1)=0$ that there exists a $\delta_2 > 0$ such that for a sufficiently large $t \in \mathcal{T}$,

$$\left\|\Delta\boldsymbol{\omega}^{(t)}\right\|_2 \geq \frac{\delta - e(t,t+1)}{1 + (1+\bar{L})\gamma^{(t)}} \geq \delta_2 > 0. \tag{66}$$

Here after we assume w.l.o.g. that (66) holds for all $t \in \mathcal{T}$ (in fact one can always restrict $\{\boldsymbol{\omega}^{(t)}\}_{t\in\mathcal{T}}$ to a proper subsequence). We show now that (63) is in contradiction with the convergence of $\{F_0(\boldsymbol{\omega}^{(t)})\}$. By Lemma 4, for all $t \in \mathcal{T}$, we have:

$$F_0(\boldsymbol{\omega}^{(t+1)}) - F_0(\boldsymbol{\omega}^{(t)}) \leq -\gamma^{(t)}\left(\mu - \frac{\hat{L}}{2}\gamma^{(t)}\right)\left\|\bar{\boldsymbol{\omega}}^{(t)} - \boldsymbol{\omega}^{(t)}\right\|_2^2$$
$$+ \gamma^{(t)}\delta\left\|\nabla F_0(\boldsymbol{\omega}^{(t)}) - \nabla\bar{F}_0^{(t)}(\boldsymbol{\omega}^{(t)})\right\|_2, \tag{67}$$

and for $t < n < i_t$,

$$F_0(\boldsymbol{\omega}^{(n+1)}) - F_0(\boldsymbol{\omega}^{(n)})$$
$$\leq -\gamma^{(n)}\left(\mu - \frac{\hat{L}}{2}\gamma^{(n)} - \frac{\left\|\nabla F_0(\boldsymbol{\omega}^{(n)}) - \nabla\bar{F}_0^{(n)}(\boldsymbol{\omega}^{(n)})\right\|_2}{\left\|\bar{\boldsymbol{\omega}}^{(n)} - \boldsymbol{\omega}^{(n)}\right\|_2}\right)\left\|\bar{\boldsymbol{\omega}}^{(n)} - \boldsymbol{\omega}^{(n)}\right\|_2^2$$
$$\leq -\gamma^{(n)}\left(\mu - \frac{\hat{L}}{2}\gamma^{(n)} - \frac{\left\|\nabla F_0(\boldsymbol{\omega}^{(n)}) - \nabla\bar{F}_0^{(n)}(\boldsymbol{\omega}^{(n)})\right\|_2}{\delta}\right)\left\|\bar{\boldsymbol{\omega}}^{(n)} - \boldsymbol{\omega}^{(n)}\right\|_2^2, \tag{68}$$

where the second inequality follows from (61). Adding (67) and (68) over $n = t+1, \cdots, i_t - 1$ and, for $t \in \mathcal{T}$ sufficiently large (so that $\mu - \frac{\hat{L}}{2}\gamma^{(t)} - \delta^{-1}\left\|\nabla F_0(\boldsymbol{\omega}^{(t)}) - \nabla\bar{F}_0^{(t)}(\boldsymbol{\omega}^{(t)})\right\|_2 \geq \hat{\mu} > 0$ and $\left\|\nabla F_0(\boldsymbol{\omega}^{(t)}) - \nabla\bar{F}_0^{(t)}(\boldsymbol{\omega}^{(t)})\right\|_2 < \hat{\mu}\delta_2^2\delta^{-1}$), we have:

$$F_0(\boldsymbol{\omega}^{(i_t)}) - F_0(\boldsymbol{\omega}^{(t)})$$
$$\overset{(a)}{\leq} -\hat{\mu}\sum_{n=t}^{i_t-1}\gamma^{(n)}\left\|\bar{\boldsymbol{\omega}}^{(n)} - \boldsymbol{\omega}^{(n)}\right\|_2^2 + \gamma^{(t)}\delta\left\|\nabla F_0(\boldsymbol{\omega}^{(t)}) - \nabla\bar{F}_0^{(t)}(\boldsymbol{\omega}^{(t)})\right\|_2$$
$$\overset{(b)}{\leq} -\hat{\mu}\delta_2^2\sum_{n=t+1}^{i_t-1}\gamma^{(n)} - \gamma^{(t)}\left(\hat{\mu}\delta_2^2 - \delta\left\|\nabla F_0(\boldsymbol{\omega}^{(t)}) - \nabla\bar{F}_0^{(t)}(\boldsymbol{\omega}^{(t)})\right\|_2\right)$$
$$\overset{(c)}{\leq} -\hat{\mu}\delta_2^2\sum_{n=t+1}^{i_t-1}\gamma^{(n)}, \tag{69}$$

where $(a)$ follows from $\mu - \frac{\hat{L}}{2}\gamma^{(t)} - \delta^{-1}\left\|\nabla F_0(\boldsymbol{\omega}^{(t)}) - \nabla\bar{F}_0^{(t)}(\boldsymbol{\omega}^{(t)})\right\|_2 \geq \hat{\mu} > 0$; $(b)$ follows from (66); and $(c)$ follows from $\left\|\nabla F_0(\boldsymbol{\omega}^{(t)}) - \nabla\bar{F}_0^{(t)}(\boldsymbol{\omega}^{(t)})\right\|_2 < \hat{\mu}\delta_2^2\delta^{-1}$. Since $\{F_0(\boldsymbol{\omega}^{(t)})\}$ converges, it must be $\liminf_{\mathcal{T}\ni t\to\infty}\sum_{n=t+1}^{i_t-1}\gamma^{(t)} = 0$, which contradicts (63). Therefore, it must be $\limsup_{t\to\infty}\left\|\bar{\boldsymbol{\omega}}^{(t)} - \boldsymbol{\omega}^{(t)}\right\|_2 = 0$ almost surely.

Finally, we show that a limit point of the sequence $\{\boldsymbol{\omega}^{(t)}\}$ generated by Algorithm 1 (Algorithm 3), i.e., $\boldsymbol{\omega}^\star$, is a stationary point of Problem 1 (Problem 6). It follows from first-order optimality condition for $\bar{\boldsymbol{\omega}}^{(t)}$ that

$$(\boldsymbol{\omega} - \bar{\boldsymbol{\omega}}^{(t)})^T\nabla\bar{F}_0^{(t)}(\boldsymbol{\omega}^{(t)}) \geq 0, \quad \boldsymbol{\omega} \in \mathbb{R}^d. \tag{70}$$

Taking the limit of (70) over the index set $\mathcal{T}$, we have:

$$\lim_{\mathcal{T}\ni t\to\infty}(\boldsymbol{\omega} - \bar{\boldsymbol{\omega}}^{(t)})^T\nabla\bar{F}_0^{(t)}(\bar{\boldsymbol{\omega}}^{(t)}) = (\boldsymbol{\omega} - \boldsymbol{\omega}^\star)^T\nabla F_0(\boldsymbol{\omega}^\star) \geq 0, \quad \boldsymbol{\omega} \in \mathbb{R}^d,$$

where the equality follows from $\lim_{t\to\infty}\left\|\bar{\boldsymbol{\omega}}^{(t)} - \boldsymbol{\omega}^{(t)}\right\|_2 = 0$ (which is due to $\liminf t \to \infty \left\|\bar{\boldsymbol{\omega}}^{(t)} - \boldsymbol{\omega}^{(t)}\right\|_2 = 0$ and $\limsup_{t\to\infty}\left\|\bar{\boldsymbol{\omega}}^{(t)} - \boldsymbol{\omega}^{(t)}\right\|_2 = 0$) and $\lim_{t\to\infty}\left\|\nabla F_0(\boldsymbol{\omega}^{(t)}) - \nabla\bar{F}_0^{(t)}(\boldsymbol{\omega}^{(t)})\right\|_2 = 0$. This is the desired first-order optimality condition and $\boldsymbol{\omega}^\star$ is a stationary point of Problem 1 (Problem 6).

## APPENDIX B: PROOFS OF THEOREM 2 AND THEOREM 4

The proofs of Theorem 2 and Theorem 4 are identical. In the following proof, we omit the subscripts $s, f$ for notation simplicity. We first introduce the following preliminary results.

*Lemma 5:* Let $(\boldsymbol{\omega}_j^\star, \mathbf{s}_j^\star)$ denote a KKT point of Problem 4 (Problem 9) with $c = c_j$ and let $(\boldsymbol{\omega}_\infty^\star, \mathbf{s}_\infty^\star)$ denote a limit point of $\{(\boldsymbol{\omega}_j^\star, \mathbf{s}_j^\star)\}$. Then, the following statements hold. i) For all $j$, if $\mathbf{s}_{s,j}^\star = \mathbf{0}$, then $\boldsymbol{\omega}_{s,j}^\star$ is a KKT point of Problem 3 (Problem 8); ii) $\mathbf{s}_\infty^\star = \mathbf{0}$, and $\boldsymbol{\omega}_\infty^\star$ is a KKT point of Problem 3 (Problem 8).

*Proof:* i) The KKT conditions of Problem 4 (Problem 9) with $c = c_j$ are given by:

$$F_m(\boldsymbol{\omega}_j^\star) \leq s_{m,j}^\star, \quad s_{m,j}^\star \geq 0, \quad m = 1, 2, \cdots, M, \tag{71a}$$
$$\lambda_m(F_m(\boldsymbol{\omega}_j^\star) - s_{m,j}^\star) = 0, \mu_m s_{m,j}^\star = 0, \quad m = 1, \cdots, M, \tag{71b}$$
$$\nabla_{\boldsymbol{\omega}}F_0(\boldsymbol{\omega}_j^\star) + \sum_{m=1}^M\lambda_m\nabla_{\boldsymbol{\omega}}F_m(\boldsymbol{\omega}_j^\star) = 0, \tag{71c}$$
$$c_j - \lambda_m - \mu_m = 0, \quad m = 1, \cdots, M. \tag{71d}$$

On the other hand, the KKT conditions of Problem 3 (Problem 8) are given by:

$$F_m(\boldsymbol{\omega}_j^\star) \leq 0, \quad m = 1, 2, \cdots, M, \tag{72a}$$
$$\lambda_m(F_m(\boldsymbol{\omega}_j^\star) - 0) = 0, \quad m = 1, \cdots, M, \tag{72b}$$
$$\nabla_{\boldsymbol{\omega}}F_0(\boldsymbol{\omega}_j^\star) + \sum_{m=1}^M\lambda_m\nabla_{\boldsymbol{\omega}}F_m(\boldsymbol{\omega}_j^\star) = 0. \tag{72c}$$

As (71) with $\mathbf{s}_{s,j}^\star = \mathbf{0}$ implies (72), we can show the first statement. ii) Construct a convex approximation of Problem 4 (Problem 9) with $c = c_j$ around $(\boldsymbol{\omega}_j^\star, \mathbf{s}_j^\star)$, which satisfies the assumptions in Theorem 2 and Theorem 4. It is clear that $(\boldsymbol{\omega}_j^\star, \mathbf{s}_j^\star)$ is an optimal solution of the approximate problem for $c = c_j$. Following the proof of [37, Theorem 1], we can show the second statement. ∎

*Lemma 6:* Let $\{\boldsymbol{\omega}^{(t)}\}$ be the sequence generated by Algorithm 2 (Algorithm 4). Then, we have:

$$\lim_{t\to\infty}\left|\bar{F}_m^{(t)}(\boldsymbol{\omega}^{(t)}) - F_m(\boldsymbol{\omega}^{(t)})\right| = 0, \quad m = 1, \cdots, M,$$
$$\lim_{t\to\infty}\left\|\nabla\bar{F}_m^{(t)}(\boldsymbol{\omega}^{(t)}) - \nabla F_m(\boldsymbol{\omega}^{(t)})\right\|_2 = 0, \quad m = 0, \cdots, M,$$
$$\lim_{t\to\infty}\left|\bar{F}_m^{(t)}(\boldsymbol{\omega}) - G_m(\boldsymbol{\omega}; \boldsymbol{\omega}^{(t)})\right| = 0, \quad \boldsymbol{\omega} \in \mathbb{R}^d, \quad m = 0, \cdots, M$$

almost surely, where $G_m(\boldsymbol{\omega}; \boldsymbol{\omega}^{(t)}) \triangleq \frac{1}{N}\sum_{n\in\mathcal{N}}g_m(\boldsymbol{\omega}; \boldsymbol{\omega}^{(t)}, \mathbf{x}_n)$.

*Proof:* Lemma 6 is a consequence of [32, Lemma 1]. We just need to verify that all the technical conditions therein are satisfied. Specifically, Condition (a) of [32, Lemma 1] is

satisfied because $\{\boldsymbol{\omega}^{(t)}\}$ is assumed to be bounded. Condition (b) of [32, Lemma 1] comes from Assumption 2.4. Conditions (c)-(d) of [32, Lemma 1] come from the stepsize rules in (4) and (6). Condition (e) of [32, Lemma 1] comes from the Lipschitz property of $F(\boldsymbol{\omega})$ from Assumption 1.2 and the stepsize rule in (6). ∎

*Lemma 7:* Consider a subsequence $\{\boldsymbol{\omega}^{(t_l)}\}_{l=1}^{\infty}$ generated by Algorithm 2 (Algorithm 4) with $c = c_j$ converging to a limit point $\boldsymbol{\omega}_j^{\star}$. There exist uniformly continuous functions $\tilde{F}_m(\boldsymbol{\omega})$, $m = 0, \cdots, M$ such that

$$\lim_{l \to \infty} \bar{F}_m^{(t_l)}(\boldsymbol{\omega}) = \tilde{F}_m(\boldsymbol{\omega}), \ \boldsymbol{\omega} \in \mathbb{R}^d, \ m = 0, \cdots, M \quad (73)$$

almost surely. Moreover, we have:

$$\tilde{F}_m(\boldsymbol{\omega}_j^{\star}) = F_m(\boldsymbol{\omega}_j^{\star}), \ m = 1, \cdots, M, \quad (74)$$

$$\nabla \tilde{F}_m(\boldsymbol{\omega}_j^{\star}) = \nabla F_m(\boldsymbol{\omega}_j^{\star}), \ m = 0, \cdots, M. \quad (75)$$

*Proof:* It readily follows from Assumption 2 that the families of functions $\{\bar{F}_m^{(t_l)}(\boldsymbol{\omega})\}$ are equicontinuous. Moreover, they are bounded and defined over a compact set. Hence, the Arzela–Ascoli theorem [36] implies that, by restricting to a subsequence, there exists uniformly continuous functions $\tilde{F}_m(\boldsymbol{\omega})$ such that (73) is satisfied. Finally, (74) and (75) follow immediately from (73) and Lemma 6. ∎

By Lemma 5, it remains to show that a limit point of $\{(\boldsymbol{\omega}^{(t)}, \mathbf{s}^{(t)})\}$ generated by Algorithm 2 (Algorithm 4) with $c = c_j$, $(\boldsymbol{\omega}_j^{\star}, \mathbf{s}_j^{\star})$, is a KKT point of Problem 4 (Problem 9). By Assumption 1, Assumption 2, and Lemma 6, we can show $\lim_{t \to \infty} \|\bar{\boldsymbol{\omega}}^{(t)} - \boldsymbol{\omega}^{(t)}\| = 0$. As the proof is similar to that in Appendix A, the details are omitted for conciseness. Consider the subsequence $\{\boldsymbol{\omega}^{(t_l)}\}_{l=1}^{\infty}$ converging to $\boldsymbol{\omega}_j^{\star}$. By $\lim_{t \to \infty} \|\bar{\boldsymbol{\omega}}^{(t)} - \boldsymbol{\omega}^{(t)}\|_2 = 0$ and $\lim_{l \to \infty} \boldsymbol{\omega}^{(t_l)} = \boldsymbol{\omega}_j^{\star}$, we have $\lim_{l \to \infty} \bar{\boldsymbol{\omega}}^{(t_l)} = \boldsymbol{\omega}_j^{\star}$. Then, by $\lim_{l \to \infty} \bar{\boldsymbol{\omega}}^{(t_l)} = \boldsymbol{\omega}_j^{\star}$, (73), and Problem 4 (Problem 9) with $c = c_j$, we have:

$$(\boldsymbol{\omega}_j^{\star}, \mathbf{s}_j^{\star}) \triangleq \arg\min_{\boldsymbol{\omega}, \mathbf{s}} \ \tilde{F}_0(\boldsymbol{\omega}) + c_j \sum_{m=1}^{M} s_m \quad (76)$$

$$\text{s.t.} \quad \tilde{F}_m(\boldsymbol{\omega}) \leq s_m, \ m = 1, 2, \cdots, M.$$

As $(\boldsymbol{\omega}_j^{\star}, \mathbf{s}_j^{\star})$ satisfies the KKT conditions of the problem in (76), and (74) and (75) in Lemma 7 hold, $\{(\boldsymbol{\omega}_j^{\star}, \mathbf{s}_j^{\star})\}$ also satisfies the KKT conditions of Problem 4 (Problem 9) with $c = c_j$, i.e., (71), implying that it is a KKT point of Problem 4 (Problem 9) with $c = c_j$. Therefore, we complete the proof.

## APPENDIX C: PROOF OF LEMMA 1

As the problem in (41) is convex and the Slater's condition holds, we can solve the problem in (41) by solving its dual problem. The Lagrangian function of the problem in (41) is:

$$\mathcal{L}(\boldsymbol{\omega}, s, \nu, \mu) = \|\boldsymbol{\omega}\|_2^2 + cs + \nu\left(\bar{F}^{(t)}(\boldsymbol{\omega}) + C_a^t - U - s\right) + \mu(-s)$$

$$= \|\boldsymbol{\omega}\|_2^2 + \nu\left(\bar{F}^{(t)}(\boldsymbol{\omega}) + C_a^t - U\right) + (c - \nu - \mu)s,$$

where $\nu$ and $\mu$ are the Lagrange multipliers. Thus, the Lagrange dual function is given by:

$$g(\nu, \mu) = \inf_{\boldsymbol{\omega}, s \geq 0} \mathcal{L}(\boldsymbol{\omega}, s, \nu, \mu)$$

$$= \begin{cases} \inf_{\boldsymbol{\omega}} \left( \|\boldsymbol{\omega}\|_2^2 + \nu\left(\bar{F}^{(t)}(\boldsymbol{\omega}) + C_a^t - U\right)\right), & c - \nu - \mu \geq 0, \\ -\infty, & c - \nu - \mu < 0. \end{cases}$$

As $\mathcal{L}(\boldsymbol{\omega}, s, \nu, \mu)$ is convex w.r.t. $\boldsymbol{\omega}$, by taking its derivative and setting it to zero, we can obtain the optimal solution:

$$\bar{\omega}_{a,0,l,j}^{(t)} = \frac{-\nu A_{a,l,j}^{(t)}}{2(1 + \nu\tau)}, \quad \bar{\omega}_{a,1,j,p}^{(t)} = \frac{-\nu B_{a,j,p}^{(t)}}{2(1 + \nu\tau)},$$

and the optimal value $h(\nu) = \nu\left(C_a^{(t)} - U - \frac{b\nu}{4(1+\tau\nu)}\right)$, where $b$ is given in (45). Therefore, the dual problem of the problem in (41) is given by:

$$\max_{\nu, \mu} \quad h(\nu)$$

$$\text{s.t.} \quad c - \nu - \mu \geq 0, \ \nu \geq 0, \ \mu \geq 0,$$

which is equivalent to the following problem:

$$\nu^* \triangleq \arg\max_{\nu} \quad h(\nu) \quad (77)$$

$$\text{s.t.} \quad 0 \leq \nu \leq c.$$

As $h(\nu)$ is convex in $\nu$, and $h'(\nu) = \frac{b - \left(b + 4\tau(U - C_a^{(t)})\right)(1 + \nu\tau)^2}{4\tau(1 + \nu\tau)^2}$, by the optimality conditions of problem in (77), we have:

$$\nu^* = \begin{cases} \left[\frac{1}{\tau}\left(\sqrt{\frac{b}{b + 4\tau(U - C_a^{(t)})}} - 1\right)\right]_0^c, & b + 4\tau(U - C_a^{(t)}) > 0 \\ c, & b + 4\tau(C_a^{(t)}) \leq 0 \end{cases},$$

which completes the proof.

## APPENDIX D: PROOF OF THEOREM 5

As $\boldsymbol{\omega}^*$ is a locally optimal solution of the problem in (32), there exists $\varepsilon > 0$ such that for all $\boldsymbol{\omega}$ with $\|\boldsymbol{\omega} - \boldsymbol{\omega}^*\|_2 < \varepsilon$, we have:

$$F(\boldsymbol{\omega}^*) + \lambda \|\boldsymbol{\omega}^*\|_2^2 \leq F(\boldsymbol{\omega}) + \lambda \|\boldsymbol{\omega}\|_2^2. \quad (78)$$

Set $U = F(\boldsymbol{\omega}^*)$. Then, for all $\boldsymbol{\omega}$ with $\|\boldsymbol{\omega} - \boldsymbol{\omega}^*\|_2 < \varepsilon$ and $F(\boldsymbol{\omega}) \leq U$, $\|\boldsymbol{\omega}^*\|_2^2 \overset{(a)}{\leq} \frac{1}{\lambda}(F(\boldsymbol{\omega}) - F(\boldsymbol{\omega}^*)) + \|\boldsymbol{\omega}\|_2^2 \overset{(b)}{\leq} \|\boldsymbol{\omega}\|_2^2$, where $(a)$ is due to (78) and $(b)$ is due to $F(\boldsymbol{\omega}) \leq U = F(\boldsymbol{\omega}^*)$. Therefore, $\boldsymbol{\omega}^*$ is a locally optimal solution of the problem in (40). The first statement holds.

As $\boldsymbol{\omega}^{\dagger}$ is a locally optimal solution of the problem in (40), the necessary KKT condition $\nabla \|\boldsymbol{\omega}^{\dagger}\|_2^2 + \xi \nabla F(\boldsymbol{\omega}^{\dagger}) = 0$ holds. Set $\lambda = \frac{1}{\xi}$. Then, we have $\nabla F(\boldsymbol{\omega}^{\dagger}) + \lambda \nabla \|\boldsymbol{\omega}^{\dagger}\|_2^2 = 0$. Therefore, $\boldsymbol{\omega}^{\dagger}$ is a stationary point of the problem in (32). If, in addtion, $\lambda$ and $\boldsymbol{\omega}^{\dagger}$ satisfy $\nabla^2 F(\boldsymbol{\omega}^{\dagger}) + \lambda I \succeq 0$, i.e., the Hessian Matrix is semi-definite, then $\boldsymbol{\omega}^{\dagger}$ is a locally optimal solution of the problem in (32). The second statement holds.

## REFERENCES

[1] C. Ye and Y. Cui, "Sample-based federated learning via mini-batch SSCA," in *Proc. IEEE ICC*, 2021, pp. 1–6.

[2] M. Li, D. G. Andersen, A. J. Smola, and K. Yu, "Communication efficient distributed machine learning with the parameter server," in *Adv. Neural Inf. Proces. Syst.*, 2014, pp. 19–27.

[3] Q. Li, Z. Wen, Z. Wu, S. Hu, N. Wang, Y. Li, X. Liu, and B. He, "A survey on federated learning systems: Vision, hype and reality for data privacy and protection," *IEEE Trans. Knowl. Data Eng.*, pp. 1–1, 2021.

[4] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, "Federated learning: Challenges, methods, and future directions," *IEEE Signal Process. Mag.*, vol. 37, no. 3, pp. 50–60, 2020.

[5] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Artificial Intelligence and Statistics*, 2017, pp. 1273–1282.

[6] H. Yu, S. Yang, and S. Zhu, "Parallel restarted SGD with faster convergence and less communication: Demystifying why model averaging works for deep learning," in *Proc. AAAI*, vol. 33, 2019, pp. 5693–5700.

[7] W. Liu, L. Chen, Y. Chen, and W. Zhang, "Accelerating federated learning via momentum gradient descent," *IEEE Trans. Parallel Distrib. Syst.*, vol. 31, no. 8, pp. 1754–1766, 2020.

[8] C. Xie, S. Koyejo, and I. Gupta, "Asynchronous federated optimization," *arXiv preprint arXiv:1903.03934*, 2020. [Online]. Available: https://arxiv.org/abs/1903.03934

[9] Y. Li, Y. Cui, and V. Lau, "Optimization-based GenQSGD for federated edge learning," in *IEEE GLOBECOM*, 2021, pp. 1–6.

[10] L. T. Phong, Y. Aono, T. Hayashi, L. Wang, and S. Moriai, "Privacy-preserving deep learning via additively homomorphic encryption," *IEEE Trans. Inf. Forensics Security*, vol. 13, no. 5, pp. 1333–1345, 2018.

[11] S. Song, K. Chaudhuri, and A. D. Sarwate, "Stochastic gradient descent with differentially private updates," in *Proc. IEEE GlobalSIP*, 2013, pp. 245–248.

[12] Q. Yang, Y. Liu, T. Chen, and Y. Tong, "Federated machine learning: Concept and applications," *ACM Trans. Intell. Syst. Technol.*, vol. 10, no. 2, pp. 1–19, 2019.

[13] S. Hardy, W. Henecka, H. Ivey-Law, R. Nock, G. Patrini, G. Smith, and B. Thorne, "Private federated learning on vertically partitioned data via entity resolution and additively homomorphic encryption," *arXiv preprint arXiv:1711.10677*, 2017. [Online]. Available: https://arxiv.org/abs/1711.10677

[14] S. Yang, B. Ren, X. Zhou, and L. Liu, "Parallel distributed logistic regression for vertical federated learning without third-party coordinator," *arXiv preprint arXiv:1911.09824*, 2019. [Online]. Available: https://arxiv.org/abs/1911.09824

[15] T. Chen, X. Jin, Y. Sun, and W. Yin, "VAFL: a method of vertical asynchronous federated learning," *arXiv preprint arXiv:2007.06081*, 2020. [Online]. Available: https://arxiv.org/abs/2007.06081

[16] P. Mohassel and Y. Zhang, "SecureML: A system for scalable privacy-preserving machine learning," in *Securit and Privacy*, 2017, pp. 19–38.

[17] H. Robbins and S. Monro, "A stochastic approximation method," *The Annals of Mathematical Statistics*, vol. 22, no. 3, pp. 400–407, 1951.

[18] L. Bottou, F. E. Curtis, and J. Nocedal, "Optimization methods for large-scale machine learning," *SIAM Review*, vol. 60, no. 2, pp. 223–311, 2018.

[19] Y. Yang, G. Scutari, D. P. Palomar, and M. Pesavento, "A parallel decomposition method for nonconvex stochastic multi-agent optimization problems," *IEEE Trans. Signal Process.*, vol. 64, no. 11, pp. 2949–2964, 2016.

[20] A. Liu, V. K. Lau, and B. Kananian, "Stochastic successive convex approximation for non-convex constrained stochastic optimization," *IEEE Trans. Signal Process.*, vol. 67, no. 16, pp. 4189–4203, 2019.

[21] C. Ye and Y. Cui, "Stochastic successive convex approximation for general stochastic optimization problems," *IEEE Wireless Commun. Lett.*, vol. 9, no. 6, pp. 755–759, 2019.

[22] P. Di Lorenzo and S. Scardapane, "Parallel and distributed training of neural networks via successive convex approximation," in *Proc. IEEE MLSP Workshop*. IEEE, 2016, pp. 1–6.

[23] S. Scardapane and P. Di Lorenzo, "Stochastic training of neural networks via successive convex approximations," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 29, no. 10, pp. 4947–4956, 2018.

[24] A. Koppel, A. Mokhtari, and A. Ribeiro, "Parallel stochastic successive convex approximation method for large-scale dictionary learning," in *ICASSP 2018*, 2018, pp. 2771–2775.

[25] S. Truex, N. Baracaldo, A. Anwar, T. Steinke, H. Ludwig, R. Zhang, and Y. Zhou, "A hybrid approach to privacy-preserving federated learning," in *Proc. 12th ACM Workshop on AISec*, 2019, pp. 1–11.

[26] H. Chen, K. Laine, and P. Rindal, "Fast private set intersection from homomorphic encryption," in *Proc. ACM CCS*, 2017, pp. 1243–1255.

[27] B. Pinkas, T. Schneider, and M. Zohner, "Scalable private set intersection based on OT extension," *ACM Trans. Priv. Secur.*, vol. 21, no. 2, pp. 1–35, 2018.

[28] N. Qian, "On the momentum term in gradient descent learning algorithms," *Neural Networks*, vol. 12, no. 1, pp. 145–151, 1999.

[29] P. Ramachandran, B. Zoph, and Q. V. Le, "Searching for activation functions," *arXiv preprint arXiv:1710.05941*, 2017. [Online]. Available: https://arxiv.org/abs/1710.05941

[30] S. Foucart and H. Rauhut, "A mathematical introduction to compressive sensing," *Bull. Am. Math*, vol. 54, pp. 151–165, 2017.

[31] Y. Cui, Y. Li, and C. Ye, *GitHub repository*. [Online]. Available: https://github.com/CuiYing123456/SB-and-FB-FL-for-Unconstrained-and-Constrained-Nonconvex-Optimization-via-Mini-batch-SSCA

[32] A. Ruszczyński, "Feasible direction methods for stochastic programming problems," *Mathematical Programming*, vol. 19, pp. 220–229, 1980.