

# Federated Composite Optimization

Honglin Yuan<sup>1,2</sup> Manzil Zaheer<sup>3</sup> Sashank Reddi<sup>3</sup>

## Abstract

Federated Learning (FL) is a distributed learning paradigm that scales on-device learning collaboratively and privately. Standard FL algorithms such as FEDAVG are primarily geared towards *smooth unconstrained* settings. In this paper, we study the *Federated Composite Optimization* (FCO) problem, in which the loss function contains a non-smooth regularizer. Such problems arise naturally in FL applications that involve sparsity, low-rank, monotonicity, or more general constraints. We first show that straightforward extensions of primal algorithms such as FEDAVG are not well-suited for FCO since they suffer from the “curse of primal averaging,” resulting in poor convergence. As a solution, we propose a new primal-dual algorithm, *Federated Dual Averaging* (FEDDUALAVG), which by employing a novel server dual averaging procedure circumvents the curse of primal averaging. Our theoretical analysis and empirical experiments demonstrate that FEDDUALAVG outperforms the other baselines.

## 1. Introduction

Federated Learning (FL, Konečný et al. 2015; McMahan et al. 2017) is a novel distributed learning paradigm in which a large number of clients collaboratively train a shared model without disclosing their private local data. The two most distinct features of FL, when compared to classic distributed learning settings, are (1) heterogeneity in data amongst the clients and (2) very high cost to communicate with a client. Due to these aspects, classic distributed optimization algorithms have been rendered ineffective in FL settings (Kairouz et al., 2019). Several algorithms specifically catered towards FL settings have been proposed to address these issues. The most prominent amongst them is Federated Averaging (FEDAVG) algorithm, which by em-

<sup>1</sup>Stanford University <sup>2</sup>Based on work performed at Google Research <sup>3</sup>Google Research. Correspondence to: Honglin Yuan <yuanhl@stanford.edu>.

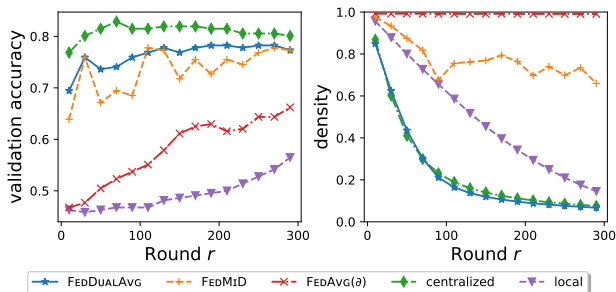


Figure 1. Results on sparse ( $\ell_1$ -regularized) logistic regression for a federated fMRI dataset based on (Haxby, 2001). centralized corresponds to training on the centralized dataset gathered from **all** the training clients. local corresponds to training on the local data from only **one** training client without communication. FEDAVG ( $\partial$ ) corresponds to running FEDAVG algorithms with subgradient in lieu of SGD to handle the non-smooth  $\ell_1$ -regularizer. FEDMID is another straightforward extension of FEDAVG running local proximal gradient method (see Section 3.1 for details). We show that using our proposed algorithm FEDDUALAVG, one can 1) achieve performance comparable to the centralized baseline without the need to gather client data, and 2) significantly outperforms the local baseline on the isolated data and the FEDAVG baseline. See Section 5.3 for details.

ploying local SGD updates, significantly reduces the communication overhead under moderate client heterogeneity. Several follow-up works have focused on improving the FEDAVG in various ways (e.g., Li et al. 2020a; Karimireddy et al. 2020; Reddi et al. 2020; Yuan & Ma 2020).

Existing FL research primarily focuses on the *unconstrained smooth* objectives; however, many FL applications involve non-smooth objectives. Such problems arise naturally in the context of regularization (e.g., sparsity, low-rank, monotonicity, or additional constraints on the model). For instance, consider the problem of cross-silo biomedical FL, where medical organizations collaboratively aim to learn a global model on their patients’ data without sharing. In such applications, sparsity constraints are of paramount importance due to the nature of the problem as it involves only a few data samples (e.g., patients) but with very high dimensions (e.g., fMRI scans). For the purpose of illustration, in Fig. 1, we present results on a federated sparse ( $\ell_1$ -regularized) logistic regression task for an fMRI dataset

(Haxby, 2001). As shown, using a federated approach that can handle non-smooth objectives enables us to find a highly accurate sparse solution without sharing client data.

In this paper, we propose to study the *Federated Composite Optimization* (FCO) problem. As in standard FL, the losses are distributed to  $M$  clients. In addition, we assume all the clients share the same, possibly non-smooth, non-finite regularizer  $\psi$ . Formally, (FCO) is of the following form

$$\min_{w \in \mathbb{R}^d} \Phi(w) := F(w) + \psi(w) := \frac{1}{M} \sum_{m=1}^M F_m(w) + \psi(w), \quad (\text{FCO})$$

where  $F_m(w) := \mathbb{E}_{\xi^m \sim \mathcal{D}_m} f(w; \xi^m)$  is the loss at the  $m$ -th client, assuming  $\mathcal{D}_m$  is its local data distribution. We assume that each client  $m$  can access  $\nabla f(w; \xi^m)$  by drawing independent samples  $\xi^m$  from its local distribution  $\mathcal{D}_m$ . Common examples of  $\psi(w)$  include  $\ell_1$ -regularizer or more broadly  $\ell_p$ -regularizer, nuclear-norm regularizer (for matrix variable), total variation (semi-)norm, etc. The (FCO) reduces to the standard federated optimization problem if  $\psi \equiv 0$ . The (FCO) also covers the constrained federated optimization if one takes  $\psi$  to be the following constraint characteristics  $\chi_{\mathcal{C}}(w) := 0$  if  $w \in \mathcal{C}$  or  $+\infty$  otherwise.

Standard FL algorithms such as FEDAVG (see Algorithm 1) and its variants (e.g., Li et al. 2020a; Karimireddy et al. 2020) are primarily tailored to *smooth unconstrained* settings, and are therefore, not well-suited for FCO. The most straightforward extension of FEDAVG towards (FCO) is to apply local subgradient method (Shor, 1985) in lieu of SGD. This approach is largely ineffective due to the intrinsic slow convergence of subgradient approach (Boyd et al., 2003), which is also demonstrated in Fig. 1 (marked FEDAVG ( $\partial$ )).

A more natural extension of FEDAVG is to replace the local SGD with proximal SGD (Parikh & Boyd 2014, a.k.a. projected SGD for constrained problems), or more generally, mirror descent (Duchi et al., 2010). We refer to this algorithm as *Federated Mirror Descent* (FEDMID, see Algorithm 2). The most noticeable drawback of a primal-averaging method like FEDMID is the ‘‘curse of primal averaging,’’ where the desired regularization of FCO may be rendered completely ineffective due to the server averaging step typically used in FL. For instance, consider a  $\ell_1$ -regularized logistic regression setting. Although each client is able to obtain a sparse solution, simply averaging the client states will inevitably yield a dense solution. See Fig. 2 for an illustrative example.

To overcome this challenge, we propose a novel primal-dual algorithm named *Federated Dual Averaging* (FEDDUALAVG, see Algorithm 3). Unlike FEDMID (or its precursor FEDAVG), the server averaging step of FEDDUALAVG operates in the dual space instead of the primal. Locally, each client runs dual averaging algorithm (Nesterov, 2009) by tracking of a pair of primal and dual states.

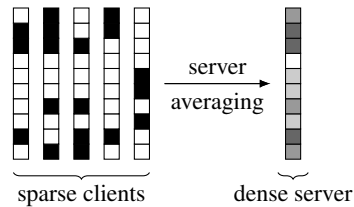


Figure 2. **Illustration of ‘‘curse of primal averaging’’**. While each client of FEDMID can locate a sparse solution, simply averaging them will yield a much denser solution on the server side.

During communication, the dual states are averaged across the clients.

Thus, FEDDUALAVG employs a novel double averaging procedure — averaging of dual states across clients (as in FEDAVG), and the averaging of gradients in dual space (as in the sequential dual averaging). Since both levels of averaging operate in the dual space, we can show that FEDDUALAVG provably overcomes the curse of primal averaging. Specifically, we prove that FEDDUALAVG can attain significantly lower communication complexity when deployed with a large client learning rate.

**Contributions.** In light of the above discussion, let us summarize our key contributions below:

- We propose a generalized federated learning problem, namely *Federated Composite Optimization* (FCO), with non-smooth regularizers and constraints.
- We first propose a natural extension of FEDAVG, namely *Federated Mirror Descent* (FEDMID). We show that FEDMID can attain the mini-batch rate in the small client learning rate regime (Section 4.1). We argue that FEDMID may suffer from the effect of ‘‘curse of primal averaging,’’ which results in poor convergence, especially in the large client learning rate regime (Section 3.2).
- We propose a novel primal-dual algorithm named *Federated Dual Averaging* (FEDDUALAVG), which provably overcomes the curse of primal averaging (Section 3.3). Under certain realistic conditions, we show that by virtue of ‘‘double averaging’’ property, FEDDUALAVG can have significantly lower communication complexity (Section 4.2).
- We demonstrate the empirical performance of FEDMID and FEDDUALAVG on various tasks, including  $\ell_1$ -regularization, nuclear-norm regularization, and various constraints in FL (Section 5).

**Notations.** We use  $[n]$  to denote the set  $\{1, \dots, n\}$ . We use  $\langle \cdot, \cdot \rangle$  to denote the inner product,  $\| \cdot \|$  to denote an arbitrary norm, and  $\| \cdot \|_*$  to denote its dual norm, unless

otherwise specified. We use  $\|\cdot\|_2$  to denote the  $\ell_2$  norm of a vector or the operator norm of a matrix, and  $\|\cdot\|_A$  to denote the vector norm induced by positive definite matrix  $A$ , namely  $\|w\|_A := \sqrt{\langle w, Aw \rangle}$ . For any convex function  $g(w)$ , we use  $g^*(z)$  to denote its convex conjugate  $g^*(z) := \sup_{w \in \mathbb{R}^d} \{\langle z, w \rangle - g(w)\}$ . We use  $w^*$  to denote the optimum of the problem (FCO). We use  $\mathcal{O}, \Theta$  to hide multiplicative absolute constants only and  $x \lesssim y$  to denote  $x = \mathcal{O}(y)$ .

### 1.1. Related Work

In this subsection, we briefly discuss the main related work. We provide a more detailed literature review in Appendix A, including the relation to classic composite optimization and distributed consensus optimization literature.

The first analysis of general FEDAVG was established by Stich (2019) for the homogeneous client dataset. This result was improved by Haddadpour et al. (2019b); Khaled et al. (2020); Woodworth et al. (2020b); Yuan & Ma (2020) via tighter analysis and accelerated algorithms. For heterogeneous clients, numerous recent papers (Haddadpour et al., 2019b; Khaled et al., 2020; Li et al., 2020b; Koloskova et al., 2020; Woodworth et al., 2020a) studied the convergence of FEDAVG under various notions of heterogeneity measure. FEDAVG has also been studied for non-convex objectives (Zhou & Cong, 2018; Haddadpour et al., 2019a; Wang & Joshi, 2018; Yu & Jin, 2019; Yu et al., 2019a;b). Other variants of FEDAVG have been proposed to overcome heterogeneity challenges (e.g., Mohri et al. 2019; Liang et al. 2019; Li et al. 2020a; Wang et al. 2020; Karimireddy et al. 2020; Pathak & Wainwright 2020; Fallah et al. 2020; Hanzely et al. 2020; T. Dinh et al. 2020; Lin et al. 2020; He et al. 2020; Bistriz et al. 2020; Zhang et al. 2020). We refer readers to (Kairouz et al., 2019) for a comprehensive survey of recent advances in FL.

We note that none of the aforementioned works can handle non-smooth problems such as (FCO). Furthermore, the contributions of this work can potentially be integrated with other emerging techniques in FL (e.g., acceleration, adaptivity, variance reduction) to overcome challenges in FL such as communication efficiency and client heterogeneity.

## 2. Preliminaries

In this section, we review the necessary background for composite optimization and federated learning. A detailed technical exposition of these topics is relegated to Appendix C.

### 2.1. Composite Optimization

Composite optimization covers a variety of statistical inference, machine learning, signal processing problems. Standard (non-distributed) composite optimization is defined as

$$\min_{w \in \mathbb{R}^d} \mathbb{E}_{\xi \sim \mathcal{D}} f(w; \xi) + \psi(w), \quad (\text{CO})$$

where  $\psi$  is a non-smooth, possibly non-finite regularizer.

**Proximal Gradient Method.** A natural extension of SGD for (CO) is the following *proximal gradient method* (PGM):

$$\begin{aligned} w_{t+1} &\leftarrow \mathbf{prox}_{\eta\psi}(w_t - \eta\nabla f(w_t; \xi_t)) \\ &= \arg \min_w \left( \eta \langle \nabla f(w_t; \xi_t), w \rangle + \frac{1}{2} \|w - w_t\|_2^2 + \eta\psi(w) \right). \end{aligned} \quad (2.1)$$

The sub-problem Eq. (2.1) can be motivated by optimizing a quadratic upper bound of  $f$  together with the original  $\psi$ . This problem can often be efficiently solved by virtue of the special structure of  $\psi$ . For instance, one can verify that PGM reduces to projected gradient descent if  $\psi$  is a constraint characteristic  $\chi_C$ , soft thresholding if  $\psi(w) = \lambda\|w\|_1$ , or weight decay if  $\psi(w) := \lambda\|w\|_2^2$ .

**Mirror Descent / Bregman-PGM.** PGM can be generalized to the Bregman-PGM if one replaces the Euclidean proximity term by the general Bregman divergence, namely

$$w_{t+1} \leftarrow \arg \min_w (\eta \langle \nabla f(w_t; \xi_t), w \rangle + \eta\psi(w) + D_h(w, w_t)), \quad (2.2)$$

where  $h$  is a strongly convex distance-generating function,  $D_h$  is the Bregman divergence which reduces to Euclidean distance if one takes  $h(w) = \frac{1}{2}\|w\|_2^2$ . We will still refer to this step as a proximal step for ease of reference. This general formulation (2.2) enables an equivalent primal-dual interpretation:

$$w_{t+1} \leftarrow \nabla(h + \eta\psi)^*(\nabla h(w_t) - \nabla f(w_t; \xi_t)). \quad (2.3)$$

A common interpretation of (2.3) is to decompose it into the following three sub-steps (Nemirovski & Yudin, 1983):

- Apply  $\nabla h$  to carry  $w_t$  to a dual state (denoted as  $z_t$ ).
- Update  $z_t$  to  $y_{t+1}$  with the gradient queried at  $w_t$ .
- Map  $y_{t+1}$  back to primal via  $\nabla(h + \eta\psi)^*$ .

This formulation is known as the *composite objective mirror descent* (COMID, Duchi et al. 2010), or simply *mirror descent* in the literature (Flammarion & Bach, 2017).

**Dual Averaging.** An alternative approach for (CO) is the following *dual averaging* algorithm (Nesterov, 2009):

$$z_{t+1} \leftarrow z_t - \eta \nabla f(\nabla(h + \eta t\psi)^*(z_t); \xi_t). \quad (2.4)$$

Similarly, we can decompose (2.4) into two sub-steps:

- Apply  $\nabla(h + \eta t\psi)^*$  to map dual state  $z_t$  to primal  $w_t$ . Note that this sub-step can be reformulated into

$$w_t = \arg \min_w (\langle -z_t, w \rangle + \eta t\psi(w) + h(w)),$$

which allows for efficient computation for many  $\psi$ .

(b) Update  $z_t$  to  $z_{t+1}$  with the gradient queried at  $w_t$ .

Dual averaging is also known as the “lazy” mirror descent algorithm (Bubeck, 2015) since it skips the forward mapping ( $\nabla h$ ) step. Theoretically, mirror descent and dual averaging often share the similar convergence rates for sequential (CO) (e.g., for smooth convex  $f$ , c.f. Flammarion & Bach 2017).

**Remark.** *There are other algorithms that are popular for certain types of (CO) problems. For example, Frank-Wolfe method (Frank & Wolfe, 1956; Jaggi, 2013) solves constrained optimization with a linear optimization oracle. Smoothing method (Nesterov, 2005) can also handle non-smoothness in objectives, but is in general less efficient than specialized CO algorithms such as dual averaging (c.f., Nesterov 2018). In this work, we mostly focus on Mirror Descent and Dual Averaging algorithms since they only employ simple proximal oracles such as projection and soft-thresholding. We refer readers to Appendix A.2 for additional related work in composite optimization.*

## 2.2. Federated Averaging

Federated Averaging (FEDAVG, McMahan et al. 2017) is the *de facto* standard algorithm for Federated Learning with unconstrained smooth objectives (namely  $\psi = 0$  for (FCO)). In this work, we follow the exposition of (Reddi et al., 2020) which splits the client learning rate and server learning rate, offering more flexibility (see Algorithm 1).

FEDAVG involves a series of *rounds* in which each round consists of a client update phase and server update phase. We denote the total number of rounds as  $R$ . At the beginning of each round  $r$ , a subset of clients  $\mathcal{S}_r$  are sampled from the client pools of size  $M$ . The server state is then broadcast to the sampled client as the client initialization. During the client update phase (highlighted in blue shade), each sampled client runs local SGD for  $K$  steps with client learning rate  $\eta_c$  with their own data. We use  $w_{r,k}^m$  to denote the  $m$ -th client state at the  $k$ -th local step of the  $r$ -th round. During the server update phase, the server averages the updates of the sampled clients and treats it as a pseudo-anti-gradient  $\Delta_r$  (Line 9). The server then takes a server update step to update its server state with server learning rate  $\eta_s$  and the pseudo-anti-gradient  $\Delta_r$  (Line 10).

## 3. Proposed Algorithms for FCO

In this section, we explore the possible solutions to approach (FCO). As mentioned earlier, existing FL algorithms such as FEDAVG and its variants do not solve (FCO). Although it is possible to apply FEDAVG to non-smooth settings by using subgradient in place of the gradient, such an approach is usually ineffective owing to the intrinsic slow convergence of subgradient methods (Boyd et al., 2003).

### Algorithm 1 Federated Averaging (FEDAVG)

```

1: procedure FEDAVG ( $w_0, \eta_c, \eta_s$ )
2: for  $r = 0, \dots, R - 1$  do
3:   sample a subset of clients  $\mathcal{S}_r \subseteq [M]$ 
4:   on client for  $m \in \mathcal{S}_r$  in parallel do
5:      $w_{r,0}^m \leftarrow w_r$   $\triangleright$  broadcast client initialization
6:     for  $k = 0, \dots, K - 1$  do
7:        $g_{r,k}^m \leftarrow \nabla f(w_{r,k}^m; \xi_{r,k}^m)$   $\triangleright$  query gradient
8:        $w_{r,k+1}^m \leftarrow w_{r,k}^m - \eta_c \cdot g_{r,k}^m$   $\triangleright$  client update
9:    $\Delta_r = \frac{1}{|\mathcal{S}_r|} \sum_{m \in \mathcal{S}_r} (w_{r,K}^m - w_{r,0}^m)$ 
10:   $w_{r+1} \leftarrow w_r + \eta_s \cdot \Delta_r$   $\triangleright$  server update
    
```

### Algorithm 2 Federated Mirror Descent (FEDMID)

```

1: procedure FEDMID ( $w_0, \eta_c, \eta_s$ )
2: for  $r = 0, \dots, R - 1$  do
3:   sample a subset of clients  $\mathcal{S}_r \subseteq [M]$ 
4:   on client for  $m \in \mathcal{S}_r$  in parallel do
5:      $w_{r,0}^m \leftarrow w_r$   $\triangleright$  broadcast primal initialization
6:     for  $k = 0, \dots, K - 1$  do
7:        $g_{r,k}^m \leftarrow \nabla f(w_{r,k}^m; \xi_{r,k}^m)$   $\triangleright$  query gradient
8:        $w_{r,k+1}^m \leftarrow \nabla (h + \eta_c \psi)^* (\nabla h(w_{r,k}^m) - \eta_c \cdot g_{r,k}^m)$ 
9:    $\Delta_r = \frac{1}{|\mathcal{S}_r|} \sum_{m \in \mathcal{S}_r} (w_{r,K}^m - w_{r,0}^m)$ 
10:   $w_{r+1} \leftarrow \nabla (h + \eta_s \eta_c K \psi)^* (\nabla h(w_r) + \eta_s \cdot \Delta_r)$ 
    
```

### 3.1. Federated Mirror Descent (FEDMID)

A more natural extension of FEDAVG towards (FCO) is to replace the local SGD steps in FEDAVG with local proximal gradient (mirror descent) steps (2.3). The resulting algorithm, which we refer to as *Federated Mirror Descent* (FEDMID)<sup>1</sup>, is outlined in Algorithm 2.

Specifically, we make two changes compared to FEDAVG:

- The client local SGD steps in FEDAVG are replaced with proximal gradient steps (Line 8).
- The server update step is replaced with another proximal step (Line 10).

As a sanity check, for constrained (FCO) with  $\psi = \chi_C$ , if one takes server learning rate  $\eta_s = 1$  and Euclidean distance  $h(w) = \frac{1}{2} \|w\|_2^2$ , FEDMID will simply reduce to the following parallel projected SGD with periodic averaging:

- (a) Each sampled client runs  $K$  steps of projected SGD following  $w_{r,k+1}^m \leftarrow \mathbf{Proj}_C(w_{r,k}^m - \eta_c g_{r,k}^m)$ .

<sup>1</sup>Despite sharing the same term “prox”, FEDMID is fundamentally different from FEDPROX (Li et al., 2020a). The proximal step in FEDPROX was to regularize the client drift caused by heterogeneity, whereas the proximal step in this work is to overcome the non-smoothness of  $\psi$ . The problems approached by the two methods are also different – FEDPROX still solves an unconstrained smooth problem, whereas ours concerns with approaches (FCO).

(b) After  $K$  local steps, the server simply average the client states following  $w_{r+1} \leftarrow \frac{1}{|\mathcal{S}_r|} \sum_{m \in \mathcal{S}_r} w_{r,K}^m$ .

### 3.2. Limitation of FEDMID: Curse of Primal Averaging

Despite its simplicity, FEDMID exhibits a major limitation, which we refer to as ‘‘curse of primal averaging’’: the server averaging step in FEDMID may severely impede the optimization progress. To understand this phenomenon, let us consider the following two illustrative examples:

- **Constrained problem:** Suppose the optimum of the aforementioned constrained problem resides on a non-flat boundary  $\mathcal{C}$ . Even when each client is able to obtain a local solution *on* the boundary, the average of them will almost surely be *off* the boundary (and hence away from the optimum) due to the curvature.
- **Federated  $\ell_1$ -regularized logistic regression problem:** Suppose each client obtains a local *sparse* solution, simply averaging them across clients will invariably yield a non-sparse solution.

As we will see theoretically (Section 4) and empirically (Section 5), the ‘‘curse of primal averaging’’ indeed hampers the performance of FEDMID.

### 3.3. Federated Dual Averaging (FEDDUALAVG)

Before we look into the solution of the curse of primal averaging, let us briefly investigate the cause of this effect. Recall that in standard smooth FL settings, server averaging step is helpful because it implicitly pools the stochastic gradients and thereby reduces the variance (Stich, 2019). In FEDMID, however, the server averaging operates on the post-proximal **primal** states, but the gradient is updated in the **dual** space (recall the primal-dual interpretation of mirror descent in Section 2.1). This primal/dual mismatch creates an obstacle for primal averaging to benefit from the pooling of stochastic gradients in dual space. This thought experiment suggests the importance of aligning the gradient update and server averaging.

Building upon this intuition, we propose a novel primal-dual algorithm, named *Federated Dual Averaging* (FEDDUALAVG, Algorithm 3), which provably addresses the curse of primal averaging. The major novelty of FEDDUALAVG, in comparison with FEDMID or its precursor FEDAVG, is to operate the server averaging in the dual space instead of the primal. This facilitates the server to aggregate the gradient information since the gradients are also accumulated in the dual space.

Formally, each client maintains a pair of primal and dual states  $(w_{r,k}^m, z_{r,k}^m)$ . At the beginning of each client update

#### Algorithm 3 Federated Dual Averaging (FEDDUALAVG)

```

1: procedure FEDDUALAVG ( $w_0, \eta_c, \eta_s$ )
2:  $z_0 \leftarrow \nabla h(w_0)$   $\triangleright$  server dual initialization
3: for  $r = 0, \dots, R - 1$  do
4:   sample a subset of clients  $\mathcal{S}_r \subseteq [M]$ 
5:   on client for  $m \in \mathcal{S}_r$  in parallel do
6:      $z_{r,0}^m \leftarrow z_r$   $\triangleright$  broadcast dual initialization
7:     for  $k = 0, \dots, K - 1$  do
8:        $\tilde{\eta}_{r,k} \leftarrow \eta_s \eta_c r K + \eta_c k$ 
9:        $w_{r,k}^m \leftarrow \nabla (h + \tilde{\eta}_{r,k} \psi)^*(z_{r,k}^m)$   $\triangleright$  retrieve primal
10:       $g_{r,k}^m \leftarrow \nabla f(w_{r,k}^m; \xi_{r,k}^m)$   $\triangleright$  query gradient
11:       $z_{r,k+1}^m \leftarrow z_{r,k}^m - \eta_c g_{r,k}^m$   $\triangleright$  client dual update
12:    $\Delta_r = \frac{1}{|\mathcal{S}_r|} \sum_{m \in \mathcal{S}_r} (z_{r,K}^m - z_{r,0}^m)$ 
13:    $z_{r+1} \leftarrow z_r + \eta_s \Delta_r$   $\triangleright$  server dual update
14:    $w_{r+1} \leftarrow \nabla (h + \eta_s \eta_c (r + 1) K \psi)^*(z_{r+1})$ 
15:    $\triangleright$  (optional) retrieve server primal state
    
```

round, the client dual state is initialized with the server dual state. During the client update stage, each client runs dual averaging steps following (2.4) to update its primal and dual state (highlighted in blue shade). The coefficient of  $\psi$ , namely  $\tilde{\eta}_{r,k}$ , is to balance the contribution from  $F$  and  $\psi$ . At the end of each client update phase, the *dual updates* (instead of primal updates) are returned to the server. The server then averages the dual updates of the sampled clients and updates the server dual state. We observe that the averaging in FEDDUALAVG is two-fold: (1) averaging of gradients in dual space within a client and (2) averaging of dual states across clients at the server. As we shall see shortly in our theoretical analysis, this novel ‘‘double’’ averaging of FEDDUALAVG in the non-smooth case enables lower communication complexity and faster convergence of FEDDUALAVG under realistic assumptions.

## 4. Theoretical Results

In this section, we demonstrate the theoretical results of FEDMID and FEDDUALAVG. We assume the following assumptions throughout the paper. The convex analysis definitions in Assumption 1 are reviewed in Appendix C.

**Assumption 1.** Let  $\|\cdot\|$  be a norm and  $\|\cdot\|_*$  be its dual.

- $\psi : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$  is a closed convex function with closed  $\mathbf{dom} \psi$ . Assume that  $\Phi(w) = F(w) + \psi(w)$  attains a finite optimum at  $w^* \in \mathbf{dom} \psi$ .
- $h : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$  is a Legendre function that is 1-strongly-convex w.r.t.  $\|\cdot\|$ . Assume  $\mathbf{dom} h \supset \mathbf{dom} \psi$ .
- $f(\cdot, \xi) : \mathbb{R}^d \rightarrow \mathbb{R}$  is a closed convex function that is differentiable on  $\mathbf{dom} \psi$  for any fixed  $\xi$ . In addition,  $f(\cdot, \xi)$  is  $L$ -smooth w.r.t.  $\|\cdot\|$  on  $\mathbf{dom} \psi$ , namely for

any  $u, w \in \text{dom } \psi$ ,

$$f(u; \xi) \leq f(w; \xi) + \langle \nabla f(w; \xi), u - w \rangle + \frac{1}{2} L \|u - w\|^2.$$

(d)  $\nabla f$  has  $\sigma^2$ -bounded variance over  $\mathcal{D}_m$  under  $\|\cdot\|_*$  within  $\text{dom } \psi$ , namely for any  $w \in \text{dom } \psi$ ,

$$\mathbb{E}_{\xi \sim \mathcal{D}_m} \|\nabla f(w, \xi) - \nabla F_m(w)\|_*^2 \leq \sigma^2, \text{ for any } m \in [M]$$

(e) Assume that all the  $M$  clients participate in the client updates for every round, namely  $\mathcal{S}_r = [M]$ .

Assumption 1(a) & (b) are fairly standard for composite optimization analysis (c.f. [Flammarion & Bach 2017](#)). Assumption 1(c) & (d) are standard assumptions in stochastic federated optimization literature ([Khaled et al., 2020](#); [Woodworth et al., 2020b](#)). (e) is assumed to simplify the exposition of the theoretical results. All results presented can be easily generalized to the partial participation case.

**Remark.** *This work focuses on convex settings because the non-convex composite optimization (either  $F$  or  $\psi$  non-convex) is noticeably challenging and under-developed even for non-distributed settings. This is in sharp contrast to non-convex smooth optimization for which simple algorithms such as SGD can readily work. Existing literature on non-convex CO (e.g., [Attouch et al. 2013](#); [Chouzenoux et al. 2014](#); [Li & Pong 2015](#); [Bredies et al. 2015](#)) typically relies on non-trivial additional assumptions (such as  $K$ - $L$  conditions) and sophisticated algorithms. Hence, it is beyond the scope of this work to study non-convex FCO.*<sup>2</sup>

#### 4.1. FEDMID and FEDDUALAVG: Small Client Learning Rate Regime

We first show that both FEDMID and FEDDUALAVG are (asymptotically) at least as good as stochastic mini-batch algorithms with  $R$  iterations and batch-size  $MK$  when client learning rate  $\eta_c$  is sufficiently small.

**Theorem 4.1** (Simplified from [Theorem F.1](#)). *Assuming Assumption 1, then for sufficiently small client learning rate  $\eta_c$ , and server learning rate  $\eta_s = \Theta(\min\{\frac{1}{\eta_c K L}, \frac{B^{\frac{1}{2}} M^{\frac{1}{2}}}{\eta_c K^{\frac{1}{2}} R^{\frac{1}{2}} \sigma}\})$ , both FEDDUALAVG and FEDMID can output  $\hat{w}$  such that*

$$\mathbb{E}[\Phi(\hat{w})] - \Phi(w^*) \lesssim \frac{LB}{R} + \frac{\sigma B^{\frac{1}{2}}}{\sqrt{MKR}}, \quad (4.1)$$

where  $B := D_h(w^*, w_0)$ .

The intuition is that when  $\eta_c$  is small, the client update will not drift too far away from its initialization of the round. Due to space constraints, the proof is relegated to [Appendix F](#).

<sup>2</sup>However, we conjecture that for simple non-convex settings (e.g., optimize non-convex  $f$  on a convex set, as tested in [Appendix B.5](#)), it is possible to show the convergence and obtain similar advantageous results for FEDDUALAVG.

#### 4.2. FEDDUALAVG with a Larger Client Learning Rate: Usefulness of Local Step

In this subsection, we show that FEDDUALAVG may attain stronger results with a larger client learning rate. In addition to possible faster convergence, [Theorems 4.2](#) and [4.3](#) also indicate that FEDDUALAVG allows for much broader searching scope of efficient learning rates configurations, which is of key importance for practical purpose.

**Bounded Gradient.** We first consider the setting with bounded gradient. Unlike unconstrained, the gradient bound may be particularly useful when the constraint is finite.

**Theorem 4.2** (Simplified from [Theorem D.1](#)). *Assuming Assumption 1 and  $\sup_{w \in \text{dom } \psi} \|\nabla f(w, \xi)\|_* \leq G$ , then for FEDDUALAVG with  $\eta_s = 1$  and  $\eta_c \leq \frac{1}{4L}$ , considering*

$$\hat{w} := \frac{1}{KR} \sum_{r=0}^{R-1} \sum_{k=1}^K \left[ \nabla (h + \tilde{\eta}_{r,k} \psi)^* \left( \frac{1}{M} \sum_{m=1}^M z_{r,k}^m \right) \right], \quad (4.2)$$

the following inequality holds

$$\mathbb{E}[\Phi(\hat{w})] - \Phi(w^*) \lesssim \frac{B}{\eta_c KR} + \frac{\eta_c \sigma^2}{M} + \eta_c^2 L K^2 G^2,$$

where  $B := D_h(w^*, w_0)$ . Moreover, there exists  $\eta_c$  such that

$$\mathbb{E}[\Phi(\hat{w})] - \Phi(w^*) \lesssim \frac{LB}{KR} + \frac{\sigma B^{\frac{1}{2}}}{\sqrt{MKR}} + \frac{L^{\frac{1}{3}} B^{\frac{2}{3}} G^{\frac{2}{3}}}{R^{\frac{2}{3}}}. \quad (4.3)$$

We refer the reader to [Appendix D](#) for complete proof details of [Theorem 4.2](#).

**Remark.** *The result in [Theorem 4.2](#) not only matches the rate by [Stich \(2019\)](#) for smooth, unconstrained FEDAVG but also allows for a general non-smooth composite  $\psi$ , general Bregman divergence induced by  $h$ , and arbitrary norm  $\|\cdot\|$ . Compared with the small learning rate result [Theorem 4.1](#), the first term in [Eq. \(4.3\)](#) is improved from  $\frac{LB}{KR}$  to  $\frac{LB}{KR}$ , whereas the third term incurs an additional loss regarding infrequent communication. One can verify that the bound [Eq. \(4.3\)](#) is better than [Eq. \(4.1\)](#) if  $R \lesssim \frac{L^2 B}{G^2}$ . Therefore, the larger client learning rate may be preferred when the communication is not too infrequent.*

**Bounded Heterogeneity.** Next, we consider the settings with bounded heterogeneity. For simplicity, we focus on the case when the loss  $F$  is quadratic, as shown in [Assumption 2](#). We will discuss other options to relax the quadratic assumption in [Section 4.3](#).

**Assumption 2** (Bounded heterogeneity, quadratic).

(a) *The heterogeneity of  $\nabla F_m$  is bounded, namely*

$$\sup_{w \in \text{dom } \psi} \|\nabla F_m(w) - \nabla F(w)\|_* \leq \zeta^2, \text{ for any } m \in [M]$$

(b)  $F(w) := \frac{1}{2}w^\top Qw + c^\top w$  for some  $Q \succ 0$ .

(c) Assume Assumption 1 is satisfied in which the norm  $\|\cdot\|$  is taken to be the  $\frac{Q}{\|Q\|_2}$ -norm, namely  $\|w\| = \sqrt{\frac{w^\top Qw}{\|Q\|_2}}$ .

**Remark.** Assumption 2(a) is a standard assumption to bound the heterogeneity among clients (e.g., Woodworth et al. 2020a). Note that Assumption 2 only assumes the objective  $F$  to be quadratic. We do not impose any stronger assumptions on either the composite function  $\psi$  or the distance-generating function  $h$ . Therefore, this result still applies to a broad class of problems such as LASSO.

The following results hold under Assumption 2. We sketch the proof in Section 4.3 and defer the details to Appendix E.

**Theorem 4.3** (Simplified from Theorem E.1). *Assuming Assumption 2, then for FEDDUALAVG with  $\eta_s = 1$  and  $\eta_c \leq \frac{1}{4L}$ , the following inequality holds*

$$\mathbb{E}[\Phi(\hat{w})] - \Phi(w^*) \lesssim \frac{B}{\eta_c KR} + \frac{\eta_c \sigma^2}{M} + \eta_c^2 LK\sigma^2 + \eta_c^2 LK^2\zeta^2,$$

where  $\hat{w}$  is the same as defined in Eq. (4.2), and  $B := D_h(w^*, w_0)$ . Moreover, there exists  $\eta_c$  such that

$$\mathbb{E}[\Phi(\hat{w})] - \Phi(w^*) \lesssim \frac{LB}{KR} + \frac{\sigma B^{\frac{1}{2}}}{\sqrt{MKR}} + \frac{L^{\frac{1}{3}} B^{\frac{2}{3}} \sigma^{\frac{2}{3}}}{K^{\frac{1}{3}} R^{\frac{2}{3}}} + \frac{L^{\frac{1}{3}} B^{\frac{2}{3}} \zeta^{\frac{2}{3}}}{R^{\frac{2}{3}}}. \quad (4.4)$$

**Remark.** The result in Theorem 4.3 matches the best-known convergence rate for smooth, unconstrained FEDAVG (Khaled et al., 2020; Woodworth et al., 2020a), while our results allow for general composite  $\psi$  and non-Euclidean distance. Compared with Theorem 4.2, the overhead in Eq. (4.4) involves variance  $\sigma$  and heterogeneity  $\zeta$  but no longer depends on  $G$ . The bound Eq. (4.4) could significantly outperform the previous ones when the variance  $\sigma$  and heterogeneity  $\zeta$  are relatively mild.

### 4.3. Proof Sketch and Discussions

In this subsection, we demonstrate our proof framework by sketching the proof for Theorem 4.3.

**Step 1: Convergence of Dual Shadow Sequence.** We start by characterizing the convergence of the dual shadow sequence  $\overline{z_{r,k}} := \frac{1}{M} \sum_{m=1}^M z_{r,k}^m$ . The key observation for FEDDUALAVG when  $\eta_s = 1$  is the following relation

$$\overline{z_{r,k+1}} = \overline{z_{r,k}} - \eta_c \cdot \frac{1}{M} \sum_{m=1}^M \nabla f(w_{r,k}^m; \xi_{r,k}^m). \quad (4.5)$$

This suggests that the shadow sequence  $\overline{z_{r,k}}$  almost executes a dual averaging update (2.4), but with some perturbed gradient  $\frac{1}{M} \sum_{m=1}^M \nabla f(w_{r,k}^m; \xi_{r,k}^m)$ . To this end, we extend the perturbed iterate analysis framework (Mania et al., 2017) to the dual space. Theoretically we show the following Lemma 4.4, with proof relegated to Appendix D.2.

**Lemma 4.4** (Convergence of dual shadow sequence of FEDDUALAVG, simplified version of Lemma D.2). *Assuming Assumption 1, then for FEDDUALAVG with  $\eta_s = 1$  and  $\eta_c \leq \frac{1}{4L}$ , the following inequality holds*

$$\begin{aligned} & \mathbb{E} \left[ \Phi \left( \frac{1}{KR} \sum_{r=0}^{R-1} \sum_{k=1}^K \nabla (h + \tilde{\eta}_{r,k} \psi)^* (\overline{z_{r,k}}) \right) \right] - \Phi(w^*) \\ & \leq \underbrace{\frac{B}{\eta_c KR} + \frac{\eta_c \sigma^2}{M}}_{\text{Rate if synchronized every iteration}} + \underbrace{\frac{L}{MKR} \left[ \sum_{r=0}^{R-1} \sum_{k=0}^{K-1} \sum_{m=1}^M \mathbb{E} \|\overline{z_{r,k}} - z_{r,k}^m\|_*^2 \right]}_{\text{Discrepancy overhead}}. \end{aligned} \quad (4.6)$$

The first two terms correspond to the rate when FEDDUALAVG is synchronized every step. The last term corresponds to the overhead for not synchronizing every step, which we call “discrepancy overhead”. Lemma 4.4 can serve as a general interface towards the convergence of FEDDUALAVG as it only assumes the blanket Assumption 1.

**Remark.** Note that the relation (4.5) is not satisfied by FEDMID due to the incommutability of the proximal operator and the averaging operator, which thereby breaks Lemma 4.4. Intuitively, this means FEDMID fails to pool the gradients properly (up to a high-order error) in the absence of communication. FEDDUALAVG overcomes the incommutability issue because all the gradients are accumulated and averaged in the dual space, whereas the proximal step only operates at the interface from dual to primal. This key difference explains the “curse of primal averaging” from the theoretical perspective.

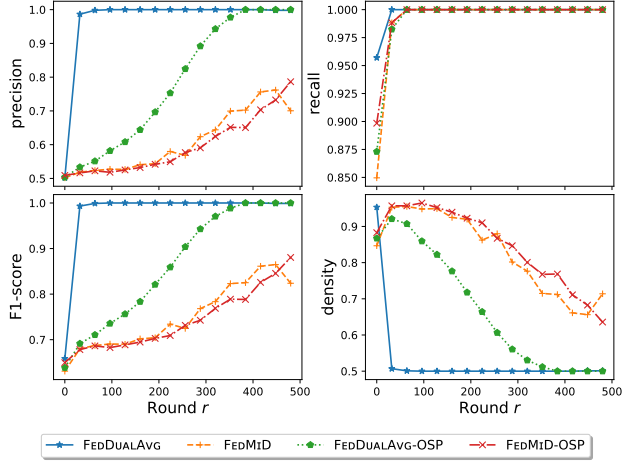
**Step 2: Bounding Discrepancy Overhead via Stability Analysis.** The next step is to bound the discrepancy term introduced in Eq. (4.6). Intuitively, this term characterizes the stability of FEDDUALAVG, in the sense that how far away a single client can deviate from the average (in dual space) if there is no synchronization for  $k$  steps.

However, unlike the smooth convex unconstrained settings in which the stability of SGD is known to be well-behaved (Hardt et al., 2016), the stability analysis for composite optimization is challenging and absent from the literature. We identify that the main challenge originates from the asymmetry of the Bregman divergence. In this work, we provide a set of simple conditions, namely Assumption 2, such that the stability of FEDDUALAVG is well-behaved.

**Lemma 4.5** (Dual stability of FEDDUALAVG under Assumption 2, simplified version of Lemma E.2). *Under the same settings of Theorem 4.3, the following inequality holds*

$$\frac{1}{M} \sum_{m=1}^M \mathbb{E} \left\| \overline{z_{r,k}} - z_{r,k}^m \right\|_*^2 \lesssim \eta_c^2 K\sigma^2 + \eta_c^2 K^2\zeta^2.$$

**Step 3: Deciding  $\eta_c$ .** The final step is to plug in the bound in step 2 back to step 1, and find appropriate  $\eta_c$  to optimize such upper bound. We defer the details to Appendix E.



**Figure 3. Sparsity recovery on a synthetic LASSO problem with 50% sparse ground truth.** Observe that FEDDUALAVG not only identifies most of the sparsity pattern but also is fastest. It is also worth noting that the less-principled FEDDUALAVG-OSP is also very competitive. The poor performance of FEDMiD can be attributed to the “curse of primal averaging”, as the server averaging step “smooths out” the sparsity pattern, which is corroborated empirically by the least sparse solution obtained by FEDMiD.

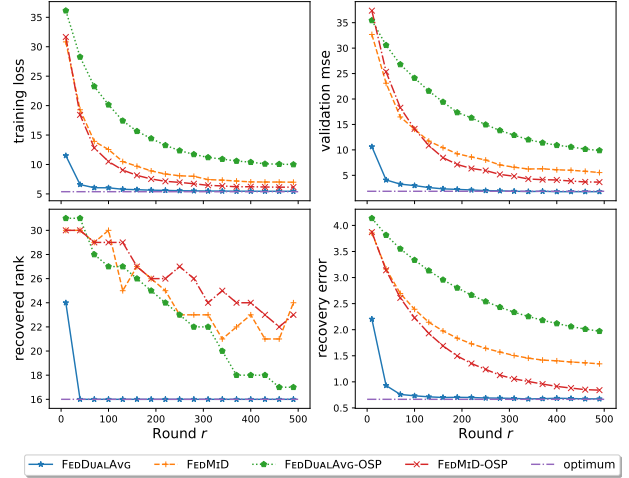
## 5. Numerical Experiments

In this section, we validate our theory and demonstrate the efficiency of the algorithms via numerical experiments. We mostly compare FEDDUALAVG with FEDMiD since the latter serves a natural baseline. We do not present subgradient-FEDAVG in this section due to its consistent ineffectiveness, as demonstrated in Fig. 1 (marked FEDAVG ( $\partial$ )). To examine the necessity of client proximal step, we also test two less-principled versions of FEDMiD and FEDDUALAVG, in which the proximal steps are only performed on the server-side. We refer to these two versions as FEDMiD-OSP and FEDDUALAVG-OSP, where “OSP” stands for “only server proximal,” with pseudo-code provided in Appendix B.1. We provide the complete setup details in Appendix B, including but not limited to hyper-parameter tuning, dataset processing and evaluation metrics. The source code is available at <https://github.com/hongliny/FCO-ICML21>.

### 5.1. Federated LASSO for Sparse Feature Recovery

In this subsection, we consider the LASSO ( $\ell_1$ -regularized least-squares) problem on a synthetic dataset, motivated by models from biomedical and signal processing literature (e.g., Ryalı et al. 2010; Chen et al. 2012). The goal is to recover the sparse signal  $w$  from noisy observations  $(x, y)$ .

$$\min_{w \in \mathbb{R}^d, b \in \mathbb{R}} \frac{1}{M} \sum_{m=1}^M \mathbb{E}_{(x,y) \sim \mathcal{D}_m} (x^\top w + b - y)_2^2 + \lambda \|w\|_1.$$



**Figure 4. Low-rank matrix estimation comparison on a synthetic dataset with the ground truth of rank 16.** We observe that FEDDUALAVG finds the solution with exact rank in less than 100 communication rounds. FEDMiD and FEDMiD-OSP converge slower in loss and rank. The unprincipled FEDDUALAVG-OSP can generate low-rank solutions but is far less accurate.

To generate the synthetic dataset, we first fix a sparse ground truth  $w_{\text{real}} \in \mathbb{R}^d$  and some bias  $b_{\text{real}} \in \mathbb{R}$ , and then sample the dataset  $(x, y)$  following  $y = x^\top w_{\text{real}} + b_{\text{real}} + \varepsilon$  for some noise  $\varepsilon$ . We let the distribution of  $(x, y)$  vary over clients to simulate the heterogeneity. We select  $\lambda$  so that the centralized solver (on gathered data) can successfully recover the sparse pattern. Since the ground truth  $w_{\text{real}}$  is known, we can assess the quality of the sparse features recovered by comparing it with the ground truth.

We evaluate the performance by recording precision, recall, sparsity density, and F1-score. We tune the client learning rate  $\eta_c$  and server learning rate  $\eta_s$  only to attain the best F1-score. The results are presented in Fig. 3. The best learning rates configuration is  $\eta_c = 0.01, \eta_s = 1$  for FEDDUALAVG, and  $\eta_c = 0.001, \eta_s = 0.3$  for other algorithms (including FEDMiD). This matches our theory that FEDDUALAVG can benefit from larger learning rates. We defer the rest of the setup details and further experiments to Appendix B.2.

### 5.2. Federated Low-Rank Matrix Estimation via Nuclear-Norm Regularization

In this subsection, we consider a low-rank matrix estimation problem via the nuclear-norm regularization

$$\min_{W, b} \frac{1}{M} \sum_{m=1}^M \mathbb{E}_{(X,y) \sim \mathcal{D}_m} (\langle X, W \rangle + b - y)^2 + \lambda \|W\|_{\text{nuc}},$$

where  $\|W\|_{\text{nuc}}$  denotes the matrix nuclear norm. The goal is to recover a low-rank matrix  $W$  from noisy observations



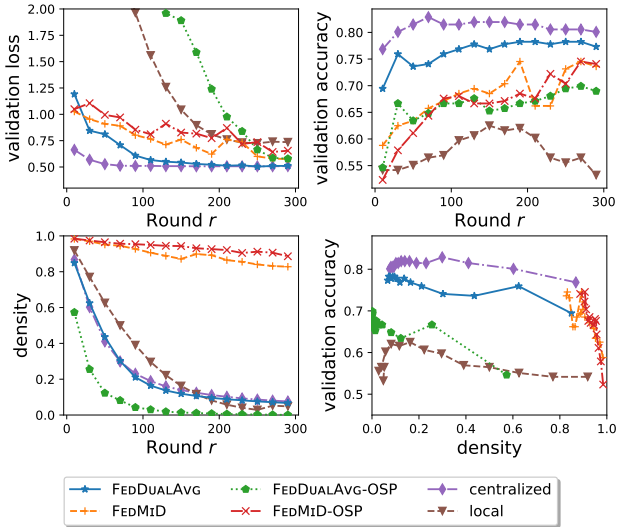


Figure 5. Results on  $\ell_1$ -regularized logistic regression for fMRI data from (Haxby, 2001). We observe that FEDDUALAVG yields sparse and accurate solutions that are comparable with the centralized baseline. FEDMID and FEDMID-OSP provides denser solutions that are relatively less accurate. The unprincipled FEDDUALAVG-OSP can provide sparse solutions but far less accurate.

$(X, y)$ . This formulation captures a variety of problems such as low-rank matrix completion and recommendation systems (Candès & Recht, 2009). Note that the proximal operator with respect to the nuclear-norm regularizer reduces to singular-value thresholding operation (Cai et al., 2010).

We evaluate the algorithms on a synthetic federated dataset with known low-rank ground truth  $W_{\text{real}} \in \mathbb{R}^{d_1 \times d_2}$  and bias  $b_{\text{real}} \in \mathbb{R}$ , similar to the above LASSO experiments. We focus on four metrics for this task: the training (regularized) loss, the validation mean-squared-error, the recovered rank, and the recovery error in Frobenius norm  $\|W_{\text{output}} - W_{\text{real}}\|_F$ . We tune the client learning rate  $\eta_c$  and server learning rate  $\eta_s$  only to attain the best recovery error. We also record the results obtained by the deterministic solver on centralized data, marked as optimum. The results are presented in Fig. 4. We provide the rest of the setup details and more experiments in Appendix B.3.

### 5.3. Sparse Logistic Regression for fMRI Scan

In this subsection, we consider the cross-silo setup of learning a binary classifier on fMRI scans. For this purpose, we use the data collected by Haxby (2001), to understand the pattern of response in the ventral temporal (vt) area of the brain given a visual stimulus. There were six subjects doing image recognition in a block-design experiment over 11 to 12 sessions, with a total of 71 sessions. Each session consists of 18 fMRI scans under the stimuli of a picture

of either a house or a face. We use the `nilearn` package (Abraham et al., 2014) to normalize and transform the four-dimensional raw fMRI scan data into an array with 39,912 volumetric pixels (voxels) using the standard mask. We plan to learn a sparse ( $\ell_1$ -regularized) binary logistic regression on the voxels to classify the stimuli given the voxels input. Enforcing sparsity is crucial for this task as it allows domain experts to understand which part of the brain is differentiating between the stimuli. We select five (out of six) subjects as the training set and the last subject as the held-out validation set. We treat each session as a client, with a total of 59 training clients and 12 validation clients, where each client possesses the voxel data of 18 scans. As in the previous experiment, we tune the client learning rate  $\eta_c$  and server learning rate  $\eta_s$  only. We set the  $\ell_1$ -regularization strength to be  $10^{-3}$ . For each setup, we run the federated algorithms for 300 communication rounds.

We compare the algorithms with two non-federated baselines: 1) centralized corresponds to training on the centralized dataset gathered from **all** the training clients; 2) local corresponds to training on the local data from only **one** training client without communication. The results are shown in Fig. 5. In Appendix B.4.2, we provide another presentation of this experiment to visualize the progress of federated algorithms and understand the robustness to learning rate configurations. The results suggest FEDDUALAVG not only recovers sparse and accurate solutions, but also behaves most robust to learning-rate configurations. We defer the rest of the setup details to Appendix B.4.

In Appendix B.5, we provide another set of experiments on federated constrained optimization for Federated EMNIST dataset (Caldas et al., 2019).

## 6. Conclusion

In this paper, we have shown the shortcomings of primal FL algorithms for FCO and proposed a primal-dual method (FEDDUALAVG) to tackle them. Our theoretical and empirical analysis provide strong evidence to support the superior performance of FEDDUALAVG over natural baselines. Potential future directions include control variates and acceleration based methods for FCO, and applying FCO to personalized settings.

## Acknowledgement

We would like to thank Zachary Charles, Zheng Xu, Andrew Hard, Ehsan Amid, Amr Ahmed, Aranyak Mehta, Qian Li, Junzi Zhang, Tengyu Ma, and TensorFlow Federated team for helpful discussions at various stages of this work. Honglin Yuan would like to thank the support by the TOTAL Innovation Scholars program. We would like to thank the anonymous reviewers for their suggestions and comments.

## References

- Abraham, A., Pedregosa, F., Eickenberg, M., Gervais, P., Mueller, A., Kossaifi, J., Gramfort, A., Thirion, B., and Varoquaux, G. Machine learning for neuroimaging with scikit-learn. *Frontiers in Neuroinformatics*, 8, 2014.
- Al-Shedivat, M., Gillenwater, J., Xing, E., and Ros-tamizadeh, A. Federated Learning via Posterior Averaging: A New Perspective and Practical Algorithms. In *International Conference on Learning Representations*, 2021.
- Attouch, H., Bolte, J., and Svaiter, B. F. Convergence of descent methods for semi-algebraic and tame problems: Proximal algorithms, forward–backward splitting, and regularized Gauss–Seidel methods. *Mathematical Programming*, 137(1-2), 2013.
- Bauschke, H. H., Borwein, J. M., et al. Legendre functions and the method of random Bregman projections. *Journal of convex analysis*, 4(1), 1997.
- Beck, A. and Teboulle, M. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Research Letters*, 31(3), 2003.
- Bistriz, I., Mann, A., and Bambos, N. Distributed Distillation for On-Device Learning. In *Advances in Neural Information Processing Systems 33*, volume 33, 2020.
- Boyd, S., Xiao, L., and Mutapic, A. Subgradient methods. *lecture notes of EE392o, Stanford University, Autumn Quarter*, 2004, 2003.
- Bredies, K., Lorenz, D. A., and Reiterer, S. Minimization of Non-smooth, Non-convex Functionals by Iterative Thresholding. *Journal of Optimization Theory and Applications*, 165(1), 2015.
- Bregman, L. The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming. *USSR Computational Mathematics and Mathematical Physics*, 7(3), 1967.
- Bubeck, S. *Convex Optimization: Algorithms and Complexity*. 2015.
- Cai, J.-F., Candès, E. J., and Shen, Z. A Singular Value Thresholding Algorithm for Matrix Completion. *SIAM Journal on Optimization*, 20(4), 2010.
- Caldas, S., Duddu, S. M. K., Wu, P., Li, T., Konečný, J., McMahan, H. B., Smith, V., and Talwalkar, A. LEAF: A Benchmark for Federated Settings. In *NeurIPS 2019 Workshop on Federated Learning for Data Privacy and Confidentiality*, 2019.
- Candès, E. J. and Recht, B. Exact Matrix Completion via Convex Optimization. *Foundations of Computational Mathematics*, 9(6), 2009.
- Chen, F., Luo, M., Dong, Z., Li, Z., and He, X. Federated Meta-Learning with Fast Convergence and Efficient Communication. *arXiv:1802.07876 [cs]*, 2019.
- Chen, X., Lin, Q., and Pena, J. Optimal regularized dual averaging methods for stochastic optimization. In *Advances in Neural Information Processing Systems 25*. Curran Associates, Inc., 2012.
- Chouzenoux, E., Pesquet, J.-C., and Repetti, A. Variable Metric Forward–Backward Algorithm for Minimizing the Sum of a Differentiable Function and a Convex Function. *Journal of Optimization Theory and Applications*, 162(1), 2014.
- Daubechies, I., Defrise, M., and De Mol, C. An iterative thresholding algorithm for linear inverse problems with a sparsity constraint. *Communications on Pure and Applied Mathematics*, 57(11), 2004.
- Dekel, O., Gilad-Bachrach, R., Shamir, O., and Xiao, L. Optimal distributed online prediction using mini-batches. *Journal of Machine Learning Research*, 13(6), 2012.
- Deng, Y., Kamani, M. M., and Mahdavi, M. Adaptive Personalized Federated Learning. *arXiv:2003.13461 [cs, stat]*, 2020.
- Diakonikolas, J. and Orecchia, L. The Approximate Duality Gap Technique: A Unified Theory of First-Order Methods. *SIAM Journal on Optimization*, 29(1), 2019.
- Duchi, J. and Singer, Y. Efficient online and batch learning using forward backward splitting. *Journal of Machine Learning Research*, 10(99), 2009.
- Duchi, J. C., Shalev-shwartz, S., Singer, Y., and Tewari, A. Composite objective mirror descent. In *COLT 2010*, 2010.
- Duchi, J. C., Agarwal, A., and Wainwright, M. J. Dual Averaging for Distributed Optimization: Convergence Analysis and Network Scaling. *IEEE Transactions on Automatic Control*, 57(3), 2012.
- Fallah, A., Mokhtari, A., and Ozdaglar, A. E. Personalized federated learning with theoretical guarantees: A model-agnostic meta-learning approach. In *Advances in Neural Information Processing Systems 33*, 2020.
- Flammarion, N. and Bach, F. Stochastic composite least-squares regression with convergence rate  $O(1/n)$ . In *Proceedings of the 2017 Conference on Learning Theory*, volume 65. PMLR, 2017.

- Frank, M. and Wolfe, P. An algorithm for quadratic programming. *Naval Research Logistics Quarterly*, 3(1-2), 1956.
- Godichon-Baggioni, A. and Saadane, S. On the rates of convergence of parallelized averaged stochastic gradient algorithms. *Statistics*, 54(3), 2020.
- Haddadpour, F., Kamani, M. M., Mahdavi, M., and Cadambe, V. Trading redundancy for communication: Speeding up distributed SGD for non-convex optimization. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97. PMLR, 2019a.
- Haddadpour, F., Kamani, M. M., Mahdavi, M., and Cadambe, V. Local SGD with periodic averaging: Tighter analysis and adaptive synchronization. In *Advances in Neural Information Processing Systems 32*. Curran Associates, Inc., 2019b.
- Hallac, D., Leskovec, J., and Boyd, S. Network Lasso: Clustering and Optimization in Large Graphs. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 2015.
- Hanzely, F., Hanzely, S., Horváth, S., and Richtárik, P. Lower bounds and optimal algorithms for personalized federated learning. In *Advances in Neural Information Processing Systems 33*, 2020.
- Hard, A., Rao, K., Mathews, R., Ramaswamy, S., Beaufays, F., Augenstein, S., Eichner, H., Kiddon, C., and Ramage, D. Federated Learning for Mobile Keyboard Prediction. *arXiv:1811.03604 [cs]*, 2018.
- Hard, A., Partridge, K., Nguyen, C., Subrahmanya, N., Shah, A., Zhu, P., Lopez-Moreno, I., and Mathews, R. Training keyword spotting models on non-iid data with federated learning. In *Interspeech 2020, 21st Annual Conference of the International Speech Communication Association, Virtual Event, Shanghai, China, 25-29 October 2020*. ISCA, 2020.
- Hardt, M., Recht, B., and Singer, Y. Train faster, generalize better: Stability of stochastic gradient descent. In *Proceedings of the 33rd International Conference on Machine Learning*, volume 48. PMLR, 2016.
- Hartmann, F., Suh, S., Komarzewski, A., Smith, T. D., and Segall, I. Federated Learning for Ranking Browser History Suggestions. *arXiv:1911.11807 [cs, stat]*, 2019.
- Haxby, J. V. Distributed and Overlapping Representations of Faces and Objects in Ventral Temporal Cortex. *Science*, 293(5539), 2001.
- He, C., Annavaram, M., and Avestimehr, S. Group Knowledge Transfer: Federated Learning of Large CNNs at the Edge. In *Advances in Neural Information Processing Systems 33*, volume 33, 2020.
- Hiriart-Urruty, J.-B. and Lemaréchal, C. *Fundamentals of Convex Analysis*. Springer Berlin Heidelberg, 2001.
- Ingerman, A. and Ostrowski, K. Introducing TensorFlow Federated, 2019.
- Jaggi, M. Revisiting Frank-Wolfe: Projection-free sparse convex optimization. In *Proceedings of the 30th International Conference on Machine Learning*, volume 28. PMLR, 2013.
- Jain, P., Kakade, S. M., Kidambi, R., Netrapalli, P., and Sidford, A. Accelerating stochastic gradient descent for least squares regression. In *Proceedings of the 31st Conference on Learning Theory*, volume 75. PMLR, 2018.
- Jiang, Y., Konečný, J., Rush, K., and Kannan, S. Improving Federated Learning Personalization via Model Agnostic Meta Learning. *arXiv:1909.12488 [cs, stat]*, 2019.
- Jung, A. On the Complexity of Sparse Label Propagation. *Frontiers in Applied Mathematics and Statistics*, 4, 2018.
- Kairouz, P., McMahan, H. B., Avent, B., Bellet, A., Bennis, M., Bhagoji, A. N., Bonawitz, K., Charles, Z., Cormode, G., Cummings, R., D'Oliveira, R. G. L., Rouayheb, S. E., Evans, D., Gardner, J., Garrett, Z., Gascón, A., Ghazi, B., Gibbons, P. B., Gruteser, M., Harchaoui, Z., He, C., He, L., Huo, Z., Hutchinson, B., Hsu, J., Jaggi, M., Javidi, T., Joshi, G., Khodak, M., Konečný, J., Korolova, A., Koushanfar, F., Koyejo, S., Lepoint, T., Liu, Y., Mittal, P., Mohri, M., Nock, R., Özgür, A., Pagh, R., Raykova, M., Qi, H., Ramage, D., Raskar, R., Song, D., Song, W., Stich, S. U., Sun, Z., Suresh, A. T., Tramèr, F., Vepakomma, P., Wang, J., Xiong, L., Xu, Z., Yang, Q., Yu, F. X., Yu, H., and Zhao, S. Advances and Open Problems in Federated Learning. *arXiv:1912.04977 [cs, stat]*, 2019.
- Karimireddy, S. P., Kale, S., Mohri, M., Reddi, S. J., Stich, S. U., and Suresh, A. T. SCAFFOLD: Stochastic Controlled Averaging for Federated Learning. In *Proceedings of the International Conference on Machine Learning 1 Pre-Proceedings (ICML 2020)*, 2020.
- Khaled, A., Mishchenko, K., and Richtárik, P. Tighter Theory for Local SGD on Identical and Heterogeneous Data. In *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108. PMLR, 2020.
- Koloskova, A., Loizou, N., Boreiri, S., Jaggi, M., and Stich, S. U. A Unified Theory of Decentralized SGD with Changing Topology and Local Updates. In *Proceedings of the International Conference on Machine Learning 1 Pre-Proceedings (ICML 2020)*, 2020.

- Konečný, J., McMahan, B., and Ramage, D. Federated optimization: Distributed optimization beyond the datacenter. In *8th NIPS Workshop on Optimization for Machine Learning*, 2015.
- Langford, J., Li, L., and Zhang, T. Sparse online learning via truncated gradient. *Journal of Machine Learning Research*, 10(28), 2009.
- Li, G. and Pong, T. K. Global Convergence of Splitting Methods for Nonconvex Composite Optimization. *SIAM Journal on Optimization*, 25(4), 2015.
- Li, T., Sahu, A. K., Zaheer, M., Sanjabi, M., Talwalkar, A., and Smith, V. Federated optimization in heterogeneous networks. In *Proceedings of Machine Learning and Systems 2020*, 2020a.
- Li, X., Huang, K., Yang, W., Wang, S., and Zhang, Z. On the convergence of FedAvg on non-iid data. In *International Conference on Learning Representations*, 2020b.
- Liang, X., Shen, S., Liu, J., Pan, Z., Chen, E., and Cheng, Y. Variance Reduced Local SGD with Lower Communication Complexity. *arXiv:1912.12844 [cs, math, stat]*, 2019.
- Lin, T., Kong, L., Stich, S. U., and Jaggi, M. Ensemble Distillation for Robust Model Fusion in Federated Learning. In *Advances in Neural Information Processing Systems 33*, 2020.
- Liu, S., Chen, P.-Y., and Hero, A. O. Accelerated Distributed Dual Averaging Over Evolving Networks of Growing Connectivity. *IEEE Transactions on Signal Processing*, 66(7), 2018.
- Lu, H., Freund, R. M., and Nesterov, Y. Relatively Smooth Convex Optimization by First-Order Methods, and Applications. *SIAM Journal on Optimization*, 28(1), 2018.
- Mahdavi, M., Yang, T., Jin, R., Zhu, S., and Yi, J. Stochastic gradient descent with only one projection. In *Advances in Neural Information Processing Systems 25*. Curran Associates, Inc., 2012.
- Mania, H., Pan, X., Papailiopoulos, D., Recht, B., Ramchandran, K., and Jordan, M. I. Perturbed Iterate Analysis for Asynchronous Stochastic Optimization. *SIAM Journal on Optimization*, 27(4), 2017.
- McDonald, R., Mohri, M., Silberman, N., Walker, D., and Mann, G. S. Efficient large-scale distributed training of conditional maximum entropy models. In *Advances in Neural Information Processing Systems 22*. Curran Associates, Inc., 2009.
- McMahan, B. Follow-the-regularized-leader and mirror descent: Equivalence theorems and L1 regularization. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15. PMLR, 2011.
- McMahan, B., Moore, E., Ramage, D., Hampson, S., and y Arcas, B. A. Communication-efficient learning of deep networks from decentralized data. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, volume 54. PMLR, 2017.
- Mohri, M., Sivek, G., and Suresh, A. T. Agnostic federated learning. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97. PMLR, 2019.
- Mokhtari, A. and Ribeiro, A. DSA: Decentralized double stochastic averaging gradient algorithm. *Journal of Machine Learning Research*, 17(61), 2016.
- Nedic, A. and Ozdaglar, A. Distributed Subgradient Methods for Multi-Agent Optimization. *IEEE Transactions on Automatic Control*, 54(1), 2009.
- Nedich, A. *Convergence Rate of Distributed Averaging Dynamics and Optimization in Networks*, volume 2. 2015.
- Nemirovski, A. and Yudin, D. B. *Problem complexity and method efficiency in optimization*. Wiley, 1983.
- Nesterov, Y. Smooth minimization of non-smooth functions. *Mathematical Programming*, 103(1), 2005.
- Nesterov, Y. Primal-dual subgradient methods for convex problems. *Mathematical Programming*, 120(1), 2009.
- Nesterov, Y. Gradient methods for minimizing composite functions. *Mathematical Programming*, 140(1), 2013.
- Nesterov, Y. *Lectures on Convex Optimization*. 2018.
- Parikh, N. and Boyd, S. P. *Proximal Algorithms*, volume 1. Now Publishers Inc., 2014.
- Pathak, R. and Wainwright, M. J. FedSplit: An algorithmic framework for fast federated optimization. In *Advances in Neural Information Processing Systems 33*, 2020.
- Rabbat, M. Multi-agent mirror descent for decentralized stochastic optimization. In *2015 IEEE 6th International Workshop on Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP)*. IEEE, 2015.
- Reddi, S., Charles, Z., Zaheer, M., Garrett, Z., Rush, K., Konečný, J., Kumar, S., and McMahan, H. B. Adaptive Federated Optimization. *arXiv:2003.00295 [cs, math, stat]*, 2020.

- Rockafellar, R. T. *Convex Analysis*. Number 28. Princeton University Press, 1970.
- Rosenblatt, J. D. and Nadler, B. On the optimality of averaging in distributed statistical learning. *Information and Inference*, 5(4), 2016.
- Ryali, S., Supekar, K., Abrams, D. A., and Menon, V. Sparse logistic regression for whole-brain classification of fMRI data. *NeuroImage*, 51(2), 2010.
- Shamir, O. and Srebro, N. Distributed stochastic optimization and learning. In *2014 52nd Annual Allerton Conference on Communication, Control, and Computing (Allerton)*. IEEE, 2014.
- Shi, W., Ling, Q., Wu, G., and Yin, W. EXTRA: An Exact First-Order Algorithm for Decentralized Consensus Optimization. *SIAM Journal on Optimization*, 25(2), 2015a.
- Shi, W., Ling, Q., Wu, G., and Yin, W. A Proximal Gradient Algorithm for Decentralized Composite Optimization. *IEEE Transactions on Signal Processing*, 63(22), 2015b.
- Shor, N. Z. *Minimization Methods for Non-Differentiable Functions*, volume 3. Springer Berlin Heidelberg, 1985.
- Smith, V., Chiang, C.-K., Sanjabi, M., and Talwalkar, A. S. Federated multi-task learning. In *Advances in Neural Information Processing Systems 30*. Curran Associates, Inc., 2017.
- Stich, S. U. Local SGD converges fast and communicates little. In *International Conference on Learning Representations*, 2019.
- Sundhar Ram, S., Nedić, A., and Veeravalli, V. V. Distributed Stochastic Subgradient Projection Algorithms for Convex Optimization. *Journal of Optimization Theory and Applications*, 147(3), 2010.
- T. Dinh, C., Tran, N., and Nguyen, T. D. Personalized Federated Learning with Moreau Envelopes. In *Advances in Neural Information Processing Systems 33*, 2020.
- Tong, Q., Liang, G., Zhu, T., and Bi, J. Federated Nonconvex Sparse Learning. *arXiv:2101.00052 [cs]*, 2020.
- Tsianos, K. I. and Rabbat, M. G. Distributed dual averaging for convex optimization under communication delays. In *2012 American Control Conference (ACC)*. IEEE, 2012.
- Tsianos, K. I., Lawlor, S., and Rabbat, M. G. Push-Sum Distributed Dual Averaging for convex optimization. In *2012 IEEE 51st IEEE Conference on Decision and Control (CDC)*. IEEE, 2012.
- Wang, J. and Joshi, G. Cooperative SGD: A unified Framework for the Design and Analysis of Communication-Efficient SGD Algorithms. *arXiv:1808.07576 [cs, stat]*, 2018.
- Wang, J., Tantia, V., Ballas, N., and Rabbat, M. SlowMo: Improving communication-efficient distributed SGD with slow momentum. In *International Conference on Learning Representations*, 2020.
- Woodworth, B., Patel, K. K., and Srebro, N. Minibatch vs Local SGD for Heterogeneous Distributed Learning. In *Advances in Neural Information Processing Systems 33*, 2020a.
- Woodworth, B., Patel, K. K., Stich, S. U., Dai, Z., Bullins, B., McMahan, H. B., Shamir, O., and Srebro, N. Is Local SGD Better than Minibatch SGD? In *Proceedings of the International Conference on Machine Learning 1 Pre-Proceedings (ICML 2020)*, 2020b.
- Xiao, L. Dual averaging methods for regularized stochastic learning and online optimization. *Journal of Machine Learning Research*, 11(88), 2010.
- Yang, T., Lin, Q., and Zhang, L. A richer theory of convex constrained optimization with reduced projections and improved rates. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70. PMLR, 2017.
- Yu, H. and Jin, R. On the computation and communication complexity of parallel SGD with dynamic batch sizes for stochastic non-convex optimization. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97. PMLR, 2019.
- Yu, H., Jin, R., and Yang, S. On the linear speedup analysis of communication efficient momentum SGD for distributed non-convex optimization. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97. PMLR, 2019a.
- Yu, H., Yang, S., and Zhu, S. Parallel Restarted SGD with Faster Convergence and Less Communication: Demystifying Why Model Averaging Works for Deep Learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, 2019b.
- Yuan, D., Hong, Y., Ho, D. W., and Jiang, G. Optimal distributed stochastic mirror descent for strongly convex optimization. *Automatica*, 90, 2018.
- Yuan, D., Hong, Y., Ho, D. W. C., and Xu, S. Distributed Mirror Descent for Online Composite Optimization. *IEEE Transactions on Automatic Control*, 2020.

- Yuan, H. and Ma, T. Federated Accelerated Stochastic Gradient Descent. In *Advances in Neural Information Processing Systems 33*, 2020.
- Yuan, K., Ling, Q., and Yin, W. On the Convergence of Decentralized Gradient Descent. *SIAM Journal on Optimization*, 26(3), 2016.
- Zhang, X., Hong, M., Dhople, S., Yin, W., and Liu, Y. FedPD: A Federated Learning Framework with Optimal Rates and Adaptivity to Non-IID Data. *arXiv:2005.11418 [cs, stat]*, 2020.
- Zhou, F. and Cong, G. On the convergence properties of a k-step averaging stochastic gradient descent algorithm for nonconvex optimization. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence*, 2018.
- Zinkevich, M. Online convex programming and generalized infinitesimal gradient ascent. In *Machine Learning, Proceedings of the Twentieth International Conference (ICML 2003), August 21-24, 2003, Washington, DC, USA*. AAAI Press, 2003.
- Zinkevich, M., Weimer, M., Li, L., and Smola, A. J. Parallelized stochastic gradient descent. In *Advances in Neural Information Processing Systems 23*. Curran Associates, Inc., 2010.

# Appendices

The appendices are structured as follows. In Appendix A, we discuss the additional related literature of this work. In Appendix B, we include additional experiments and detailed setups. In Appendix C, we provide the necessary backgrounds for our theoretical results. We prove two of our main results, namely Theorems 4.2 and 4.3, in Appendices D and E, respectively. The proof of Theorem 4.1 is sketched in Appendix F.

## List of Appendices

<b>A</b>	<b>Additional Related Work</b>	<b>15</b>
A.1	Federated Learning . . . . .	15
A.2	Composite Optimization, Dual Averaging, and Mirror Descent . . . . .	16
A.3	Classic Decentralized Consensus Optimization . . . . .	16
<b>B</b>	<b>Additional Experiments and Setup Details</b>	<b>17</b>
B.1	General Setup . . . . .	17
B.2	Federated LASSO for Sparse Feature Selection . . . . .	18
B.3	Nuclear-Norm-Regularization for Low-Rank Matrix Estimation . . . . .	20
B.4	Sparse Logistic Regression for fMRI . . . . .	22
B.5	Constrained Federated Optimization for Federated EMNIST . . . . .	23
<b>C</b>	<b>Theoretical Background and Technicalities</b>	<b>26</b>
<b>D</b>	<b>Proof of Theorem 4.2: Convergence of FEDDUALAVG Under Bounded Gradient Assumption</b>	<b>28</b>
D.1	Main Theorem and Lemmas . . . . .	28
D.2	Perturbed Iterate Analysis of FEDDUALAVG: Proof of Lemma D.2 . . . . .	29
D.3	Stability of FEDDUALAVG Under Bounded Gradient Assumptions: Proof of Lemma D.3 . . . . .	33
<b>E</b>	<b>Proof of Theorem 4.3: Convergence of FEDDUALAVG Under Bounded Heterogeneity and Quadratic Assumptions</b>	<b>33</b>
E.1	Main Theorem and Lemmas . . . . .	33
E.2	Stability of FEDDUALAVG Under Quadratic Assumptions: Proof of Lemma E.2 . . . . .	34
<b>F</b>	<b>Proof of Theorem 4.1</b>	<b>39</b>
F.1	Main Theorem and Lemmas . . . . .	39
F.2	Deferred Proof of Lemma F.2 . . . . .	40
F.3	Deferred Proof of Lemma F.3 . . . . .	41

## A. Additional Related Work

### A.1. Federated Learning

Recent years have witnessed a growing interest in various aspects of Federated Learning. The early analysis of FEDAVG preceded the inception of Federated Learning, which was studied under the names of parallel SGD and local SGD (Zinkevich et al., 2010; Zhou & Cong, 2018). Early results on FEDAVG mostly focused on the “one-shot” averaging case, in which the clients are only synchronized once at the end of the procedure (e.g., McDonald et al. 2009; Shamir & Srebro 2014; Rosenblatt & Nadler 2016; Jain et al. 2018; Godichon-Baggioni & Saadane 2020). The first analysis of general FEDAVG was established by (Stich, 2019) for the homogeneous client dataset. This result was improved by (Haddadpour et al., 2019b);

Khaled et al., 2020; Woodworth et al., 2020b; Yuan & Ma, 2020) via tighter analysis and accelerated algorithms. FEDAVG has also been studied for non-convex objectives (Zhou & Cong, 2018; Haddadpour et al., 2019a; Wang & Joshi, 2018; Yu & Jin, 2019; Yu et al., 2019a;b). For heterogeneous clients, numerous recent papers (Haddadpour et al., 2019b; Khaled et al., 2020; Li et al., 2020b; Koloskova et al., 2020; Woodworth et al., 2020a) studied the convergence of FEDAVG under various notions of heterogeneity measure. Other variants of FEDAVG have been proposed to overcome heterogeneity (e.g., Mohri et al. 2019; Zhang et al. 2020; Li et al. 2020a; Wang et al. 2020; Karimireddy et al. 2020; Reddi et al. 2020; Pathak & Wainwright 2020; Al-Shedivat et al. 2021). A recent line of work has studied the behavior of Federated algorithms for personalized multi-task objectives (Smith et al., 2017; Hanzely et al., 2020; T. Dinh et al., 2020; Deng et al., 2020) and meta-learning objectives (Fallah et al., 2020; Chen et al., 2019; Jiang et al., 2019). Federated Learning techniques have been successfully applied in a broad range of practical applications (Hard et al., 2018; Hartmann et al., 2019; Hard et al., 2020). We refer readers to (Kairouz et al., 2019) for a comprehensive survey of the recent advances in Federated Learning. However, none of the aforementioned work allows for non-smooth or constrained problems such as (FCO). To the best of our knowledge, the present work is the first work that studies non-smooth or constrained problems in Federated settings.

Shortly after the initial preprint release of the present work, Tong et al. (2020) proposed a related federated  $\ell_0$ -constrained problem (which does not belong to FCO due to the non-convexity of  $\ell_0$ ), and two algorithms to solve (similar to FEDMID (-OSP) but with hard-thresholding instead). As in most hard-thresholding work, the convergence is weaker since it depends on the sparsity level  $\tau$  (worsens as  $\tau$  gets tighter).

### A.2. Composite Optimization, Dual Averaging, and Mirror Descent

Composite optimization has been a classic problem in convex optimization, which covers a variety of statistical inference, machine learning, signal processing problems. Mirror Descent (MD, a generalization of proximal gradient method) and Dual Averaging (DA, a.k.a. lazy mirror descent) are two representative algorithms for convex composite optimization. The *Mirror Descent* (MD) method was originally introduced by Nemirovski & Yudin (1983) for the constrained case and reinterpreted by Beck & Teboulle (2003). MD was generalized to the composite case by Duchi et al. (2010) under the name of COMID, though numerous preceding work had studied the special case of COMID under a variety of names such as gradient mapping (Nesterov, 2013), forward-backward splitting method (FOBOS, Duchi & Singer 2009), iterative shrinkage and thresholding (ISTA, Daubechies et al. 2004), and truncated gradient (Langford et al., 2009). The *Dual Averaging* (DA) method was introduced by Nesterov (2009) for the constrained case, which is also known as *Lazy Mirror Descent* in the literature (Bubeck, 2015). The DA method was generalized to the composite (regularized) case by (Xiao, 2010; Dekel et al., 2012) under the name of Regularized Dual Averaging, and extended by recent works (Flammarion & Bach, 2017; Lu et al., 2018) to account for non-Euclidean geometry induced by an arbitrary distance-generating function  $h$ . DA also has its roots in online learning (Zinkevich, 2003), and is related to the follow-the-regularized-leader (FTRL) algorithms (McMahan, 2011). Other variants of MD or DA (such as delayed / skipped proximal step) have been investigated to mitigate the expensive proximal oracles (Mahdavi et al., 2012; Yang et al., 2017). Composite optimization has also been studied for non-convex objective (Attouch et al., 2013; Chouzenoux et al., 2014; Bredies et al., 2015; Li & Pong, 2015). These works are typically limited to special cases due to the hardness of non-convex composite optimization, which is in sharp contrast to smooth non-convex settings. In addition to MD and DA, there are other algorithms that are popular for certain types of composite optimization problems. For example, Frank-Wolfe method (Frank & Wolfe, 1956; Jaggi, 2013) solves constrained optimization with a linear optimization oracle, which is different from the proximal oracle applied by MD and DA. We refer readers to (Flammarion & Bach, 2017; Diakonikolas & Orecchia, 2019) for more detailed discussions on the recent advances of MD and DA.

### A.3. Classic Decentralized Consensus Optimization

A related distributed setting is the *decentralized consensus optimization*, also known as *multi-agent optimization* or *optimization over networks* in the literature (Nedich, 2015). Unlike the federated settings, in decentralized consensus optimization, each client can communicate every iteration, but the communication is limited to its graphic neighborhood. Standard algorithms for unconstrained consensus optimization include decentralized (sub)gradient methods (Nedic & Ozdaglar, 2009; Yuan et al., 2016) and EXTRA (Shi et al., 2015a; Mokhtari & Ribeiro, 2016). For constrained or composite consensus problems, people have studied both mirror-descent type methods (with primal consensus), e.g., (Sundhar Ram et al., 2010; Shi et al., 2015b; Rabbat, 2015; Yuan et al., 2018; 2020); and dual-averaging type methods (with dual consensus), e.g., (Duchi et al., 2012; Tsianos et al., 2012; Tsianos & Rabbat, 2012; Liu et al., 2018). In particular, the distributed dual averaging (Duchi et al., 2012) has gained great popularity since its dual consensus scheme elegantly handles the constraints,



and overcomes the technical difficulties of primal consensus, as noted by the original paper. We identify that while the federated settings share certain backgrounds with the decentralized consensus optimization, the motivations, techniques, challenges, and results are quite dissimilar due to the fundamental difference of communication protocol, as noted by (Kairouz et al., 2019). We refer readers to (Nedich, 2015) for a more detailed introduction to the classic decentralized consensus optimization.

Another related (but different) non-smooth distribution optimization setting is the network Lasso (NLASSO, Hallac et al. 2015) problem. NLASSO considers  $\min \sum_i F_i(w_i) + \sum_{(j,k) \in \mathcal{E}} A_{jk} \|w_j - w_k\|_2$ , where each client  $i$  only optimizes a disjoint subset of parameters  $w_i$ , with non-smooth TV-regularization to control the difference between clients. A similar inter-client regularized objective has also been studied by Smith et al. (2017) for multi-task learning purpose. Due to its decentralized nature, NLASSO can often be efficiently solved by decomposition methods (e.g., ADMM by Hallac et al. (2015), preconditioned Chambolle-Pock method by Jung (2018)). These decomposition methods often require (almost) exact solution of a subproblem involving  $f_i$ , which is typically not required in FL algorithms and analyses. *In contrast*, the present work considers (FCO) where all clients jointly optimize the shared  $w$ . The non-smooth  $\psi$  is to regularize the desired property of the shared  $w$  (e.g., sparsity). Algorithmically, decomposition methods do *not* apply to FCO since  $w$  is shared. As in standard FL practice (Kairouz et al., 2019), our algorithms follow the local optimization + periodic synchronization framework to save communication.

## B. Additional Experiments and Setup Details

### B.1. General Setup

**Algorithms.** In this paper we mainly test four Federated algorithms, namely Federated Mirror Descent (FEDMID, see Algorithm 2), Federated Dual Averaging (FEDDUALAVG, see Algorithm 3), as well as two less-principled algorithms which skip the client-side proximal operations. We refer to these two algorithms as FEDMID-OSP and FEDDUALAVG-OSP, where ‘‘OSP’’ stands for ‘‘only server proximal’’. We formally state these two OSP algorithms in Algorithms 4 and 5. We study these two OSP algorithms mainly for ablation study purpose, those they might be of special interest if the proximal step is computationally intensive. For instance, in FEDMID-OSP, the client proximal step is replaced by  $w_{r,k+1}^m \leftarrow \nabla h^*(\nabla h(w_{r,k}^m) - \eta_c g_{r,k}^m)$  with no  $\psi$  involved (see line 8 of Algorithm 4). This step reduces to the ordinary SGD  $w_{r,k+1}^m \leftarrow w_{r,k}^m - \eta_c g_{r,k}^m$  if  $h(w) = \frac{1}{2} \|w\|_2^2$  in which case both  $\nabla h$  and  $\nabla h^*$  are identity mapping. Theoretically, it is not hard to establish similar rates of Theorem 4.1 for FEDMID-OSP with finite  $\psi$ . For infinite  $\psi$ , we need extension of  $f$  outside  $\text{dom}\psi$  to fix regularity. To keep this paper focused, we will not establish these results formally. There is no theoretical guarantee on the convergence of FEDDUALAVG-OSP.

---

#### Algorithm 4 Federated Mirror Descent Only Server Proximal (FEDMID-OSP)

---

```

1: procedure FEDMID-OSP ( $w_0, \eta_c, \eta_s$ )
2: for  $r = 0, \dots, R - 1$  do
3:   sample a subset of clients  $\mathcal{S}_r \subseteq [M]$ 
4:   on client for  $m \in \mathcal{S}_r$  in parallel do
5:     client initialization  $w_{r,0}^m \leftarrow w_r$  ▷ Broadcast primal initialization for round  $r$ 
6:     for  $k = 0, \dots, K - 1$  do
7:        $g_{r,k}^m \leftarrow \nabla f(w_{r,k}^m; \xi_{r,k}^m)$  ▷ Query gradient
8:        $w_{r,k+1}^m \leftarrow \nabla h^*(\nabla h(w_{r,k}^m) - \eta_c g_{r,k}^m)$  ▷ Client (primal) update – proximal operation skipped
9:        $\Delta_r = \frac{1}{|\mathcal{S}_r|} \sum_{m \in \mathcal{S}_r} (w_{r,K}^m - w_{r,0}^m)$  ▷ Compute pseudo-anti-gradient
10:     $w_{r+1} \leftarrow \nabla (h + \eta_s \eta_c K \psi)^*(\nabla h(w_r) + \eta_s \Delta_r)$  ▷ Server (primal) update

```

---

**Environments.** We simulate the algorithms in the TensorFlow Federated (TFF) framework (Ingerman & Ostrowski, 2019). The implementation is based on the Federated Research repository.<sup>3</sup>

**Tasks.** We experiment the following four tasks in this work.

1. Federated Lasso ( $\ell_1$ -regularized least squares) for sparse feature selection, see Appendix B.2.

<sup>3</sup><https://github.com/google-research/federated>

**Algorithm 5** Federated Dual Averaging Only Server Proximal (FEDDUALAVG-OSP)

---

```

1: procedure FEDDUALAVG-OSP( $w_0, \eta_c, \eta_s$ )
2: server initialization  $z_0 \leftarrow \nabla h(w_0)$ 
3: for  $r = 0, \dots, R - 1$  do
4:   sample a subset of clients  $\mathcal{S}_r \subseteq [M]$ 
5:   on client for  $m \in \mathcal{S}_r$  in parallel do
6:     client initialization  $z_{r,0}^m \leftarrow z_r$  ▷ Broadcast dual initialization for round  $r$ 
7:     for  $k = 0, \dots, K - 1$  do
8:        $w_{r,k}^m \leftarrow \nabla h^*(z_{r,k}^m)$  ▷ Compute primal point  $w_{r,k}^m$  – proximal operation skipped
9:        $g_{r,k}^m \leftarrow \nabla f(w_{r,k}^m; \xi_{r,k}^m)$  ▷ Query gradient
10:       $z_{r,k+1}^m \leftarrow z_{r,k}^m - \eta_c g_{r,k}^m$  ▷ Client (dual) update
11:       $\Delta_r = \frac{1}{|\mathcal{S}_r|} \sum_{m \in \mathcal{S}_r} (z_{r,K}^m - z_{r,0}^m)$  ▷ Compute pseudo-anti-gradient
12:       $z_{r+1} \leftarrow z_r + \eta_s \Delta_r$  ▷ Server (dual) update
13:       $w_{r+1} \leftarrow \nabla (h + \eta_s \eta_c (r + 1) K \psi)^*(z_{r+1})$  ▷ (Optional) Compute server primal state

```

---

2. Federated low-rank matrix recovery via nuclear-norm regularization, see Appendix B.3.
3. Federated sparse ( $\ell_1$ -regularized) logistic regression for fMRI dataset (Haxby, 2001), see Appendix B.4.
4. Federated constrained optimization for Federated EMNIST dataset (Caldas et al., 2019), see Appendix B.5.

We take the distance-generating function  $h$  to be  $h(w) := \frac{1}{2} \|w\|_2^2$  for all the four tasks. The detailed setups of each experiment are stated in the corresponding subsections.

## B.2. Federated LASSO for Sparse Feature Selection

### B.2.1. SETUP DETAILS

In this experiment, we consider the federated LASSO ( $\ell_1$ -regularized least squares) on a synthetic dataset inspired by models from biomedical and signal processing literature (e.g., Ryali et al. 2010; Chen et al. 2012)

$$\min_{w,b} \frac{1}{M} \sum_{m=1}^M \mathbb{E}_{(x,y) \sim \mathcal{D}_m} (x^\top w + b - y)_2^2 + \lambda \|w\|_1.$$

The goal is to retrieve sparse features of  $w$  from noisy observations.

**Synthetic Dataset Descriptions.** We first generate the ground truth  $w_{\text{real}}$  with  $d_1$  ones and  $d_0$  zeros for some  $d_1 + d_0 = d$ , namely

$$w_{\text{real}} = \begin{bmatrix} \mathbf{1}_{d_1} \\ \mathbf{0}_{d_0} \end{bmatrix} \in \mathbb{R}^d,$$

and ground truth  $b_{\text{real}} \sim \mathcal{N}(0, 1)$ .

The observations  $(x, y)$  are generated as follows to simulate the heterogeneity among clients. Let  $(x_m^{(i)}, y_m^{(i)})$  denotes the  $i$ -th observation of the  $m$ -th client. For each client  $m$ , we first generate and fix the mean  $\mu_m \sim \mathcal{N}(0, I_{d \times d})$ . Then we sample  $n_m$  pairs of observations following

$$\begin{aligned} x_m^{(i)} &= \mu_m + \delta_m^{(i)}, & \text{where } \delta_m^{(i)} &\sim \mathcal{N}(\mathbf{0}_d, I_{d \times d}) \text{ are i.i.d., for } i = 1, \dots, n_m; \\ y_m^{(i)} &= w_{\text{real}}^\top x_m^{(i)} + b_{\text{real}} + \varepsilon_m^{(i)}, & \text{where } \varepsilon_m^{(i)} &\sim \mathcal{N}(0, 1) \text{ are i.i.d., for } i = 1, \dots, n_m. \end{aligned}$$

We test four configurations of the above synthetic dataset.

- (I) The ground truth  $w_{\text{real}}$  has  $d_1 = 512$  ones and  $d_0 = 512$  zeros. We generate  $M = 64$  training clients where each client possesses 128 pairs of samples. There are 8,192 training samples in total.

- (II) (sparse ground truth) The ground truth  $w_{\text{real}}$  has  $d_1 = 64$  ones and  $d_0 = 960$  zeros. The rest of the configurations are the same as dataset (I).
- (III) (sparser ground truth) The ground truth  $w_{\text{real}}$  has  $d_1 = 8$  ones and  $d_0 = 1016$  zeros. The rest of the configurations are the same as dataset (I).
- (IV) (more distributed data) The ground truth is the same as (I). We generate  $M = 256$  training clients where each client possesses 32 pairs of samples. The total number of training examples are the same.

**Evaluation Metrics.** Since the ground truth of the synthetic dataset is known, we can evaluate the quality of the sparse features retrieved by comparing it with the ground truth. To numerically evaluate the sparsity, we treat all the features in  $w$  with absolute values smaller than  $10^{-2}$  as zero elements, and non-zero otherwise. We evaluate the performance by recording precision, recall, F1-score, and sparse density.

**Hyperparameters.** For all algorithms, we tune the client learning rate  $\eta_c$  and server learning rate  $\eta_s$  only. We test 49 different combinations of  $\eta_c$  and  $\eta_s$ .  $\eta_c$  is selected from  $\{0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1\}$ , and  $\eta_s$  is selected from  $\{0.01, 0.03, 0.1, 0.3, 1, 3, 10\}$ . All methods are tuned to achieve the best averaged recovery error over the last 100 communication rounds. We claim that the best learning rate combination falls in this range for all the algorithms tested. We draw 10 clients uniformly at random at each communication round and let the selected clients run local algorithms with batch size 10 for one epoch (of its local dataset) for this round. We run 500 rounds in total, though FEDDUALAVG usually converges to almost perfect solutions in much fewer rounds.

The Fig. 4 presented in the main paper (Section 5.1) is for the synthetic dataset (I). Now we test the performance on the other three datasets.

B.2.2. RESULTS ON SYNTHETIC DATASET (II) AND (III) WITH SPARSER GROUND TRUTH

We repeat the experiments on the dataset (II) and (III) with  $1/2^4$  and  $1/2^7$  ground truth density, respectively. The results are shown in Figs. 6 and 7. We observe that FEDDUALAVG converges to the perfect F1-score in less than 100 rounds, which outperforms the other baselines by a margin. The F1-score of FEDDUALAVG-OSP converges faster on these sparser datasets than (I), which makes it comparably more competitive. The convergence of FEDMID and FEDMID-OSP remains slow.

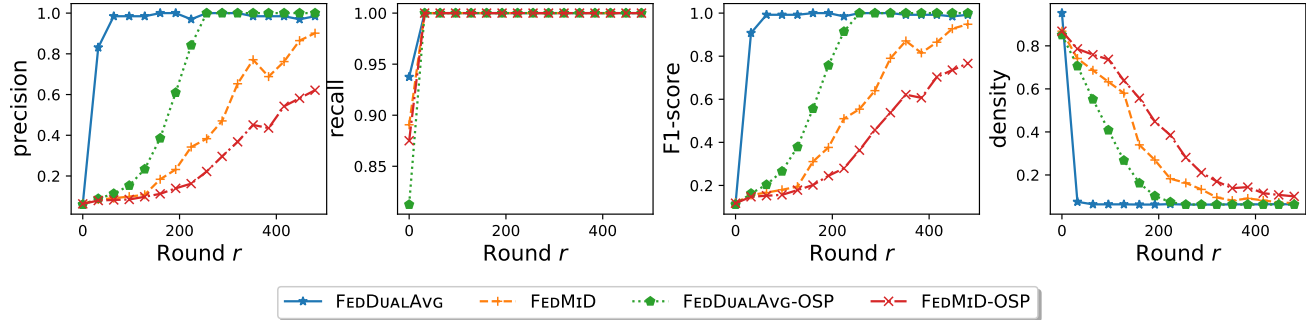


Figure 6. Results on Dataset (II):  $1/2^4$  Ground Truth Density. See Appendix B.2.2 for discussions.

B.2.3. RESULTS ON SYNTHETIC DATASET (IV): MORE DISTRIBUTED DATA (256 CLIENTS)

We repeat the experiments on the dataset (IV) with more distributed data (256 clients). The results are shown in Fig. 8. We observe that all the four algorithms take more rounds to converge in that each client has fewer data than the previous configurations. FEDDUALAVG manages to find perfect F1-score in less than 200 rounds, which outperforms the other algorithms significantly. FEDDUALAVG-OSP can recover an almost perfect F1-score after 500 rounds, but is much slower than on the less distributed dataset (I). FEDMID and FEDMID-OSP have very limited progress within 500 rounds. This is because the server averaging step in FEDMID and FEDMID-OSP fails to aggregate the sparsity patterns properly. Since each client is subject to larger noise due to the limited amount of local data, simply averaging the primal updates will “smooth out” the sparsity pattern.

### B.3. Nuclear-Norm-Regularization for Low-Rank Matrix Estimation

#### B.3.1. SETUP DETAILS

In this subsection, we consider a low-rank matrix estimation problem via the nuclear-norm regularization

$$\min_{W \in \mathbb{R}^{d_1 \times d_2}, b \in \mathbb{R}} \frac{1}{M} \sum_{m=1}^M \mathbb{E}_{(X,y) \sim \mathcal{D}_m} (\langle X, W \rangle + b - y)^2 + \lambda \|W\|_{\text{nuc}},$$

where  $\|W\|_{\text{nuc}} := \sum_i \sigma_i(W)$  denotes the nuclear norm (a.k.a. trace norm) defined by the summation of all the singular values. The goal is to recover a low-rank matrix  $W$  from noisy observations  $(X, y)$ . This formulation captures a variety of problems, such as low-rank matrix completion and recommendation systems (c.f. [Candès & Recht 2009](#)). Note that the proximal operator with respect to the nuclear-norm regularizer  $\|\cdot\|_{\text{nuc}}$  reduces to the well-known singular-value thresholding operation ([Cai et al., 2010](#)).

**Synthetic Dataset Descriptions.** We first generate the following ground truth  $W_{\text{real}} \in \mathbb{R}^{d \times d}$  of rank  $r$

$$W_{\text{real}} = \begin{bmatrix} I_{r \times r} & \mathbf{0}_{r \times (d-r)} \\ \mathbf{0}_{(d-r) \times r} & \mathbf{0}_{(d-r) \times (d-r)} \end{bmatrix},$$

and ground truth  $b_{\text{real}} \sim \mathcal{N}(0, 1)$ .

The observations  $(X, y)$  are generated as follows to simulate the heterogeneity among clients. Let  $(X_m^{(i)}, y_m^{(i)})$  denotes the  $i$ -th observation of the  $m$ -th client. For each client  $m$ , we first generate and fix the mean  $\mu_m \in \mathbb{R}^{d \times d}$  where all coordinates are i.i.d. standard Gaussian  $\mathcal{N}(0, 1)$ . Then we sample  $n_m$  pairs of observations following

$$\begin{aligned} X_m^{(i)} &= \mu_m + \delta_m^{(i)}, \text{ where } \delta_m^{(i)} \in \mathbb{R}^{d \times d} \text{ is a matrix with all coordinates from standard Gaussian;} \\ y_m^{(i)} &= \langle w_{\text{real}}, X_m^{(i)} \rangle + b_{\text{real}} + \varepsilon_m^{(i)}, \text{ where } \varepsilon_m^{(i)} \sim \mathcal{N}(0, 1) \text{ are i.i.d.} \end{aligned}$$

We tested four configurations of the above synthetic dataset.

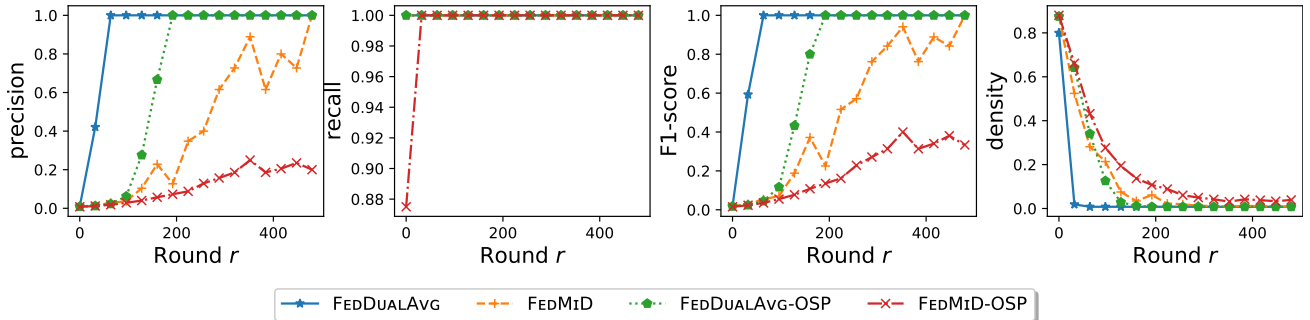


Figure 7. Results on Dataset (III):  $1/2^7$  Ground Truth Density. See Appendix B.2.2 for discussions.

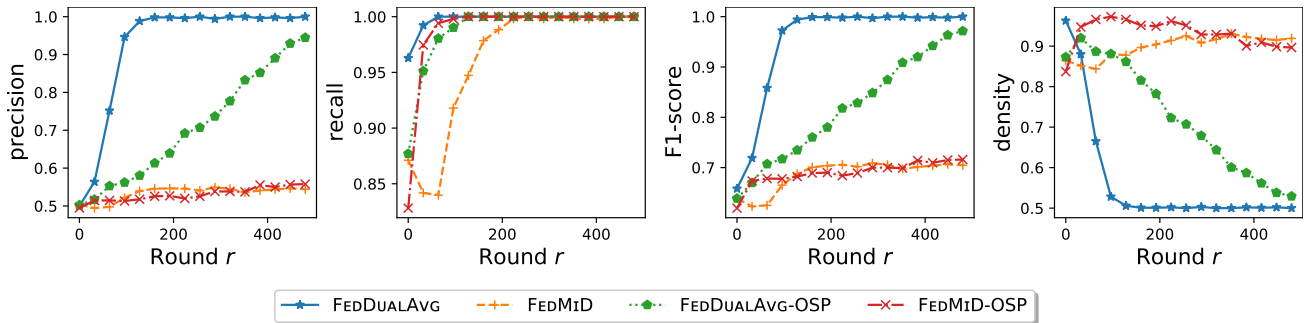


Figure 8. Results on Dataset (IV): More Distributed Data. See Appendix B.2.3 for discussions.

- (I) The ground truth  $W_{\text{real}}$  is a matrix of dimension  $32 \times 32$  with rank  $r = 16$ . We generate  $M = 64$  training clients where each client possesses 128 pairs of samples. There are 8,192 training samples in total.
- (II) (rank-4 ground truth) The ground truth  $W_{\text{real}}$  has rank  $r = 4$ . The other configurations are the same as the dataset (I).
- (III) (rank-1 ground truth) The ground truth  $W_{\text{real}}$  has rank  $r = 1$ . The other configurations are the same as the dataset (I).
- (IV) (more distributed data) The ground truth is the same as (I). We generate  $M = 256$  training clients where each client possesses 32 samples. The total number of training examples remains the same.

**Evaluation Metrics.** We focus on four metrics for this task: the training (regularized) loss, the validation mean-squared-error, the recovered rank, and the recovery error in Frobenius norm  $\|W_{\text{output}} - W_{\text{real}}\|_F$ . To numerically evaluate the rank, we count the number of singular values that are greater than  $10^{-2}$ .

**Hyperparameters.** For all algorithms, we tune the client learning rate  $\eta_c$  and server learning rate  $\eta_s$  only. We test 49 different combinations of  $\eta_c$  and  $\eta_s$ .  $\eta_c$  is selected from  $\{0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1\}$ , and  $\eta_s$  is selected from  $\{0.01, 0.03, 0.1, 0.3, 1, 3, 10\}$ . All methods are tuned to achieve the best averaged recovery error on the last 100 communication rounds. We claim that the best learning rate combination falls in this range for all algorithms tested. We draw 10 clients uniformly at random at each communication round and let the selected clients run local algorithms with batch size 10 for one epoch (of its local dataset) for this round. We run 500 rounds in total, though FEDDUALAVG usually converges to perfect F1-score in much fewer rounds.

The Fig. 4 presented in the main paper (Section 5.2) is for the synthetic dataset (I). Now we test the performance of the algorithms on the other three datasets.

### B.3.2. RESULTS ON SYNTHETIC DATASET (II) AND (III) WITH GROUND TRUTH OF LOWER RANK

We repeat the experiments on the dataset (II) and (III) with 4 and 1 ground truth rank, respectively. The results are shown in Figs. 9 and 10. The results are qualitatively reminiscent of the previous experiments on the dataset (I). FEDDUALAVG can recover the exact rank in less than 100 rounds, which outperforms the other baselines by a margin. FEDDUALAVG-OSP can recover a low-rank solution but is less accurate. The convergence of FEDMID and FEDMID-OSP remains slow.

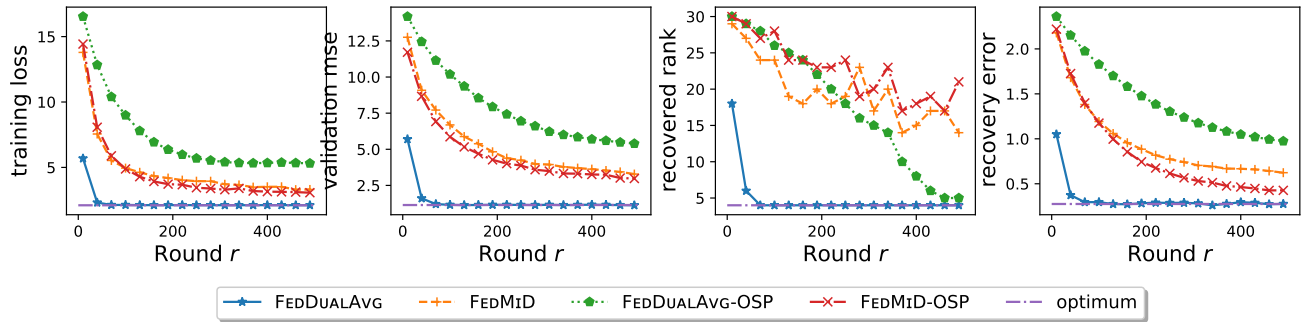


Figure 9. Results on Dataset (II): Ground Truth Rank 4. See Appendix B.3.2 for discussions.

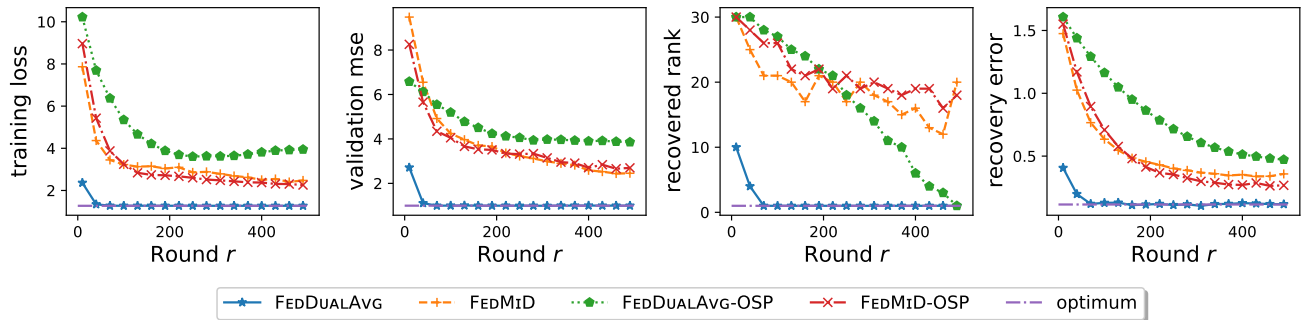


Figure 10. Results on Dataset (III): Ground Truth Rank 1. See Appendix B.3.2 for discussions.

B.3.3. RESULTS ON SYNTHETIC DATASET (IV): MORE DISTRIBUTED DATA (256 CLIENTS)

We repeat the experiments on the dataset (IV) with more distributed data. The results are shown in Fig. 11. We observe that all four algorithms take more rounds to converge in that each client has fewer data than the previous configurations. The other messages are qualitatively similar to the previous experiments – FEDDUALAVG manages to find exact rank in less than 200 rounds, which outperforms the other algorithms significantly.

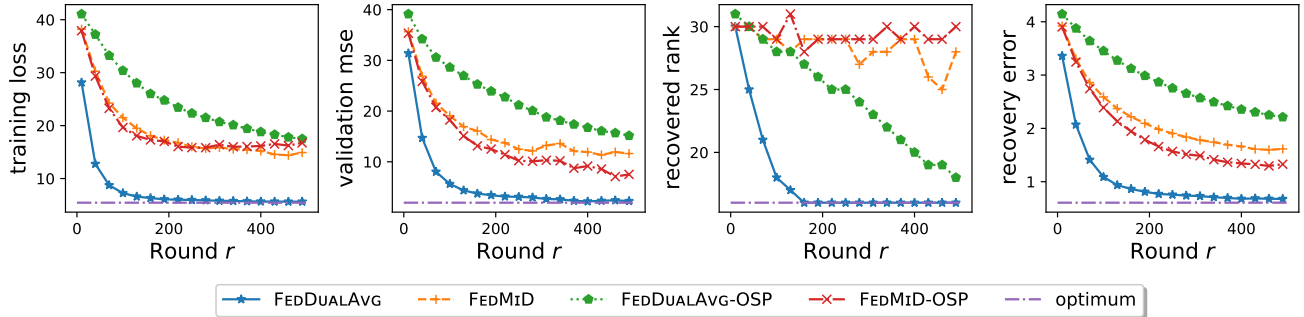


Figure 11. Results on Dataset (IV): More Distributed Data. See Appendix B.3.3 for discussions.

B.4. Sparse Logistic Regression for fMRI

B.4.1. SETUP DETAILS

In this subsection, we provide the additional setup details for the fMRI experiment presented in Fig. 5. The goal is to understand the pattern of response in ventral temporal area of the brain given a visual stimulus. Enforcing sparsity is important as it allows domain experts to understand which part of the brain is differentiating between the stimuli. We apply  $\ell_1$ -regularized logistic regression on the voxels to classify the visual stimuli.

**Dataset Descriptions and Preprocessing.** We use data collected by Haxby (2001). There were 6 subjects doing binary image recognition (from a horse and a face) in a block-design experiment over 11-12 sessions per subject, in which each session consists of 18 scans. We use nilearn package (Abraham et al., 2014) to normalize and transform the 4-dimensional raw fMRI scan data into an array with 39,912 volumetric pixels (voxels) using the standard mask. We choose the first 5 subjects as training set and the last subject as validation set. To simulate the cross-silo federated setup, we treat each session as a client. There are 59 training clients and 12 test clients, where each client possesses the voxel data of 18 scans.

**Evaluation Metrics.** We focus on three metrics for this task: validation (regularized) loss, validation accuracy, and (sparsity) density. To numerically evaluate the density, we treat all weights with absolute values smaller than  $10^{-4}$  as zero elements. The density is computed as non-zero parameters divided by the total number of parameters.

**Hyperparameters.** For all algorithms, we adjust only client learning rate  $\eta_c$  and server learning rate  $\eta_s$ . For each federated setup, we tested 49 different combinations of  $\eta_c$  and  $\eta_s$ .  $\eta_c$  is selected from  $\{0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1\}$ , and  $\eta_s$  is selected from  $\{0.01, 0.03, 0.1, 0.3, 1, 3, 10\}$ . We let each client run its local algorithm with batch-size one for one epoch per round. At the beginning of each round, we draw 20 clients uniformly at random. We run each configuration for 300 rounds and present the configuration with the lowest validation (regularized) loss at the last round.

We also tested two non-federated baselines for comparison, marked as centralized and local. centralized corresponds to training on the centralized dataset gathered from all the 59 training clients. local corresponds to training on the local data from only one training client without communication. We run proximal gradient descent for these two baselines for 300 epochs. The learning rate is tuned from  $\{0.0001, 0.0003, 0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1\}$  to attain the best validation loss at the last epoch. The results are presented in Fig. 5.

B.4.2. PROGRESS VISUALIZATION ACROSS VARIOUS LEARNING RATE CONFIGURATIONS

In this subsection, we present an alternative viewpoint to visualize the progress of federated algorithms and understand the robustness to hyper-parameters. To this end, we run four algorithms for various learning rate configurations (we present all the combinations of learning rates mentioned above such that  $\eta_c \eta_s \in [0.003, 0.3]$ ) and record the validation accuracy and (sparsity) density after 10th, 30th, 100th, and 300th round. The results are presented in Fig. 12. Each dot stands for a

learning rate configuration (client and server). We can observe that most FEDDUALAVG configurations reach the upper-left region of the box, which indicates sparse and accurate solutions. FEDDUALAVG-OSP reaches to the mid-left region of the box, which indicates sparse but less accurate solutions. The majority of FEDMiD and FEDMiD-OSP lands on the right side region box, which reflects the hardness for FEDMiD and FEDMiD-OSP to find the sparse solutions.

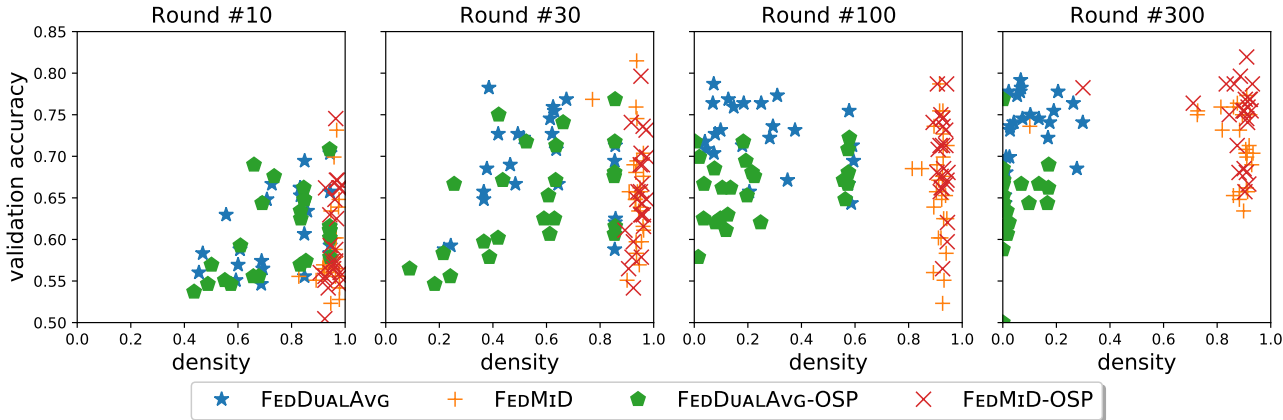


Figure 12. Progress of Federated Algorithms Under Various Learning Rate Configurations for fMRI. Each dot stands for a learning rate configuration (client and server). FEDDUALAVG recovers sparse and accurate solutions, and is robust to learning-rate configurations.

### B.5. Constrained Federated Optimization for Federated EMNIST

#### B.5.1. SETUP DETAILS

In this task we test the performance of the algorithms when the composite term  $\psi$  is taken to be convex characteristics

$$\chi_{\mathcal{C}}(w) := \begin{cases} 0 & \text{if } w \in \mathcal{C}, \\ +\infty & \text{if } w \notin \mathcal{C}. \end{cases} \text{ which encodes a hard constraint.}$$

**Dataset Descriptions and Models.** We tested on the Federated EMNIST (FEMNIST) dataset provided by TensorFlow Federated, which was derived from the Leaf repository (Caldas et al., 2019). EMNIST is an image classification dataset that extends MNIST dataset by incorporating alphabetical classes. The Federated EMNIST dataset groups the examples from EMNIST by writers.

We tested two versions of FEMNIST in this work:

- (I) FEMNIST-10: digits-only version of FEMNIST which contains 10 label classes. We experiment the logistic regression models with  $\ell_1$ -ball-constraint or  $\ell_2$ -ball-constraint on this dataset. Note that for this task we only trained on 10% of the examples in the original FEMNIST-10 dataset because the original FEMNIST-10 has an unnecessarily large number (340k) of examples for the logistic regression model.
- (II) FEMNIST-62: full version of FEMNIST which contains 62 label classes (including 52 alphabetical classes and 10 digital classes). We test a two-hidden-layer fully connected neural network model where all fully connected layers are simultaneously subject to  $\ell_1$ -ball-constraint. Note that there is no theoretical guarantee for either of the four algorithms on non-convex objectives. We directly implement the algorithms as if the objectives were convex. We defer the study of FEDMiD and FEDDUALAVG for non-convex objectives to the future work.

**Evaluation Metrics.** We focused on three metrics for this task: training error, training accuracy, and test accuracy. Note that the constraints are always satisfied because all the trajectories of all the four algorithms are always in the feasible region.

**Hyperparameters.** For all algorithms, we tune only the client learning rate  $\eta_c$  and server learning rate  $\eta_s$ . For each setup, we tested 25 different combinations of  $\eta_c$  and  $\eta_s$ .  $\eta_c$  is selected from  $\{0.001, 0.003, 0.01, 0.03, 0.1\}$ , and  $\eta_s$  is selected from  $\{0.01, 0.03, 0.1, 0.3, 1\}$ . We draw 10 clients uniformly at random at each communication round and let the selected clients run local algorithms with batch size 10 for 10 epochs (of its local dataset) for this round. We run 5,000 communication rounds in total and evaluate the training loss every 100 rounds. All methods are tuned to achieve the best averaged training loss on the last 10 checkpoints.

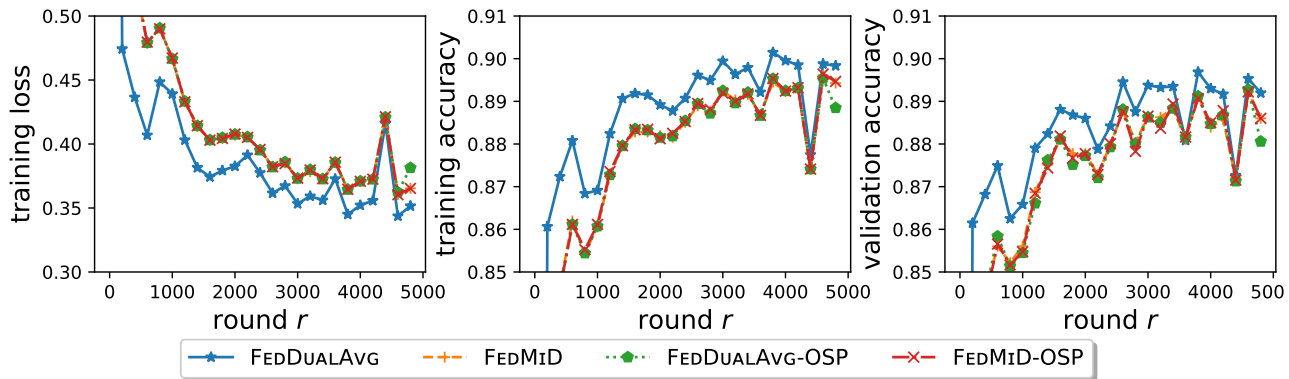


Figure 13.  $\ell_1$ -Constrained logistic regression. Dataset: FEMNIST-10. Constraint:  $\|w\|_1 \leq 1000$ .

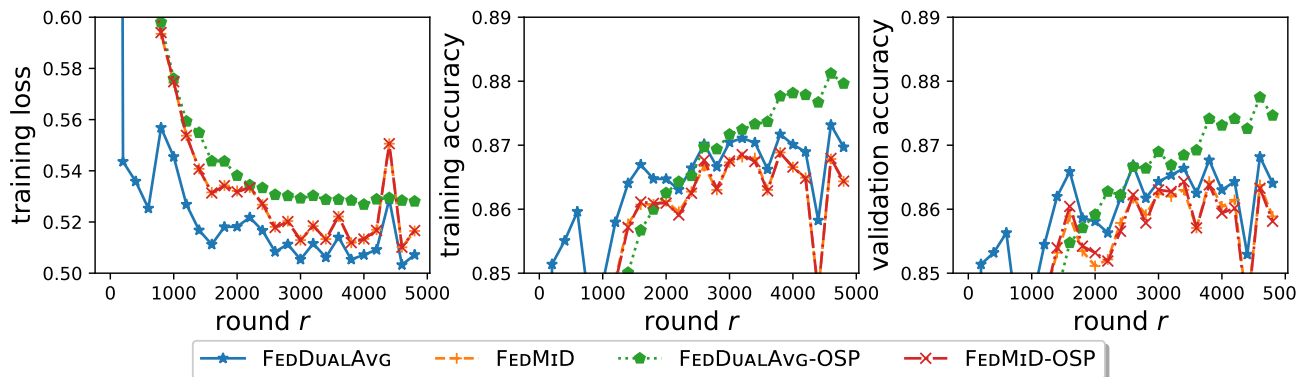


Figure 14.  $\ell_2$ -constrained logistic regression. Dataset: FEMNIST-10. Constraint:  $\|w\|_2 \leq 10$ .

### B.5.2. EXPERIMENTAL RESULTS

**$\ell_1$ -Constrained Logistic Regression** We first test the  $\ell_1$ -regularized logistic regression. The results are shown in Fig. 13. We observe that FEDDUALAVG outperforms the other three algorithms by a margin. Somewhat surprisingly, we observe that the other three algorithms behave very closely in terms of the three metrics tested. This seems to suggest that the client proximal step (in this case projection step) might be saved in FEDMiD.

**$\ell_2$ -Constrained Logistic Regression** Next, we test the  $\ell_2$ -regularized logistic regression. The results are shown in Fig. 14. We observe that FEDDUALAVG outperforms the FEDMiD and FEDMiD-OSP in all three metrics (note again that FEDMiD and FEDMiD-OSP share very similar trajectories). Interestingly, the FEDDUALAVG-OSP behaves much worse in training loss than the other three algorithms, but the training accuracy and validation accuracy are better. We conjecture that this effect might be attributed to the homogeneous property of  $\ell_2$ -constrained logistic regression which FEDDUALAVG-OSP can benefit from.

**$\ell_1$ -Constrained Two-Hidden-Layer Neural Network** Finally, we test on the two-hidden-layer neural network with  $\ell_1$ -constraints. The results are shown in Fig. 15. We observe that FEDDUALAVG outperforms FEDMiD and FEDMiD-OSP in all three metrics (once again, note that FEDMiD and FEDMiD-OSP share similar trajectories). On the other hand, FEDDUALAVG-OSP behaves much worse (which is out of the plotting ranges). This is not quite surprising because FEDDUALAVG-OSP does not have any theoretical guarantees.



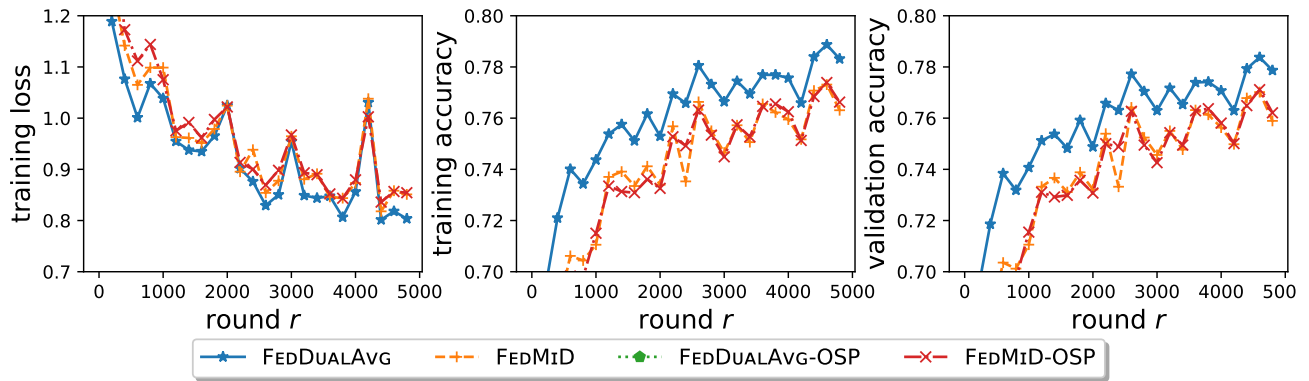


Figure 15.  $\ell_1$ -Constrained Two-Hidden-Layer Neural Network. Dataset: FEMNIST-62. Constraint: all three dense kernels  $w^{[l]}$  simultaneously satisfy  $\|w^{[l]}\|_1 \leq 1000$ .

## C. Theoretical Background and Technicalities

In this section, we introduce some definitions and propositions that are necessary for the proof of our theoretical results. Most of the definitions and results are standard and can be found in the classic convex analysis literature (e.g., [Rockafellar 1970](#); [Hiriart-Urruty & Lemaréchal 2001](#)), unless otherwise noted.

The following definition of the *effective domain* extends the notion of *domain* (of a finite-valued function) to an extended-valued convex function  $\mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ .

**Definition C.1** (Effective domain). *Let  $g : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$  be an extended-valued convex function. The **effective domain** of  $g$ , denoted by  $\mathbf{dom} g$ , is defined by*

$$\mathbf{dom} g := \{w \in \mathbb{R}^d : g(w) < +\infty\}.$$

In this work we assume all extended-valued convex functions discussed are **proper**, namely the effective domain is nonempty.

Next, we formally define the concept of *strict* and *strong convexity*. Note that the strong convexity is parametrized by some parameter  $\mu > 0$  and therefore implies strict convexity.

**Definition C.2** (Strict and Strong convexity ([Hiriart-Urruty & Lemaréchal, 2001](#), Definition B.1.1.1)). *A convex function  $g : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$  is **strictly convex** if for any  $w_1, w_2 \in \mathbf{dom} g$ , for any  $\alpha \in (0, 1)$ , it is the case that*

$$g(\alpha w_1 + (1 - \alpha)w_2) < \alpha g(w_1) + (1 - \alpha)g(w_2).$$

*Moreover,  $g$  is  $\mu$ -strongly convex with respect to  $\|\cdot\|$  norm if for any  $w_1, w_2 \in \mathbf{dom} g$ , for any  $\alpha \in (0, 1)$ , it is the case that*

$$g(\alpha w_1 + (1 - \alpha)w_2) \leq \alpha g(w_1) + (1 - \alpha)g(w_2) - \frac{1}{2}\mu\alpha(1 - \alpha)\|w_2 - w_1\|^2.$$

The notion of *convex conjugate* (a.k.a. *Legendre-Fenchel transformation*) is defined as follows. The outcome of convex conjugate is always convex and closed.

**Definition C.3** (Convex conjugate). *Let  $g : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$  be a convex function. The convex conjugate is defined as*

$$g^*(z) := \sup_{w \in \mathbb{R}^d} \{\langle z, w \rangle - g(w)\}.$$

The following result shows that the differentiability of the conjugate function and the strict convexity of the original function is linked.

**Proposition C.4** (Differentiability of the conjugate of strictly convex function ([Hiriart-Urruty & Lemaréchal, 2001](#), Theorem E.4.1.1)). *Let  $g : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$  be a closed, strictly convex function. Then we have  $\mathbf{int} \mathbf{dom} g^* \neq \emptyset$  and  $g^*$  is continuously differentiable on  $\mathbf{int} \mathbf{dom} g^*$  (where  $\mathbf{int}$  stands for interior).*

*Moreover, for  $z \in \mathbf{int} \mathbf{dom} g^*$ , it is the case that*

$$\nabla g^*(z) = \arg \min_w \{-\langle z, w \rangle + g(w)\}.$$

The differentiability in Proposition C.4 can be strengthened to smoothness if we further assume the strong convexity of the original function  $g$ .

**Proposition C.5** (Smoothness of the conjugate of strongly convex function ([Hiriart-Urruty & Lemaréchal, 2001](#), Theorem E.4.2.1)). *Let  $g : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$  be a closed,  $\mu$ -strongly convex function. Then  $g^*$  is continuously differentiable on  $\mathbb{R}^d$ , and  $g^*$  is  $\frac{1}{\mu}$ -smooth on  $\mathbb{R}^d$ , namely  $\|\nabla g^*(z) - \nabla g^*(y)\|^* \leq \frac{1}{\mu}\|z - y\|$ .*

Next we define the *Legendre function class*.

**Definition C.6** (Legendre function ([Rockafellar, 1970](#), §26)). *A proper, convex, closed function  $h : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$  is **of Legendre type** if*

- (a)  $h$  is *strictly convex*.
- (b)  $h$  is *essentially smooth*, namely  $h$  is differentiable on  $\text{int dom } h$ , and  $\|\nabla h(w_k)\| \rightarrow \infty$  for every sequence  $\{w_k\}_{k=0}^{\infty} \subset \text{int dom } h$  converging to a boundary point of  $\text{dom } h$  as  $k \rightarrow +\infty$ .

An important property of the Legendre function is the following proposition (Bauschke et al., 1997).

**Proposition C.7** (Rockafellar (1970, Theorem 26.5)). *A convex function  $g$  is of Legendre type if and only if its conjugate  $g^*$  is. In this case, the gradient mapping  $\nabla g$  is a topological isomorphism with inverse mapping, namely  $(\nabla g)^{-1} = \nabla g^*$ .*

Next, recall the definition of Bregman divergence:

**Definition C.8** (Bregman divergence (Bregman, 1967)). *Let  $g : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$  be a closed, strictly convex function that is differentiable in  $\text{int dom } g$ . The **Bregman divergence**  $D_g(w, u)$  for  $w \in \text{dom } g$ ,  $u \in \text{int dom } g$  is defined by*

$$D_g(w, u) = g(w) - g(u) - \langle \nabla g(u), w - u \rangle.$$

Note the definition of Bregman divergence requires the differentiability of the base function  $g$ . To extend the concept of Bregman divergence to non-differentiable function  $g$ , we consider the following generalized Bregman divergence (slightly modified from (Flammarion & Bach, 2017)). The generalized Bregman divergence plays an important role in the analysis of FEDDUALAVG.

**Definition C.9** (Generalized Bregman divergence (slightly modified from Flammarion & Bach, 2017, Section B.2)). *Let  $g : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$  be a closed strictly convex function (which may not be differentiable). The **Generalized Bregman divergence**  $\tilde{D}_g(w, z)$  for  $w \in \text{dom } g$ ,  $z \in \text{int dom } g^*$  is defined by*

$$\tilde{D}_g(w, z) = g(w) - g(\nabla g^*(z)) - \langle z, w - \nabla g^*(z) \rangle.$$

Note that  $\nabla g^*$  is well-defined because  $g^*$  is differentiable in  $\text{int dom } g^*$  according to Proposition C.4.

The generalized Bregman divergence is lower bounded by the ordinary Bregman divergence in the following sense.

**Proposition C.10** ((Flammarion & Bach, 2017, Lemma 6)). *Let  $h : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$  be a Legendre function. Let  $\psi : \mathbb{R}^d \rightarrow \mathbb{R}$  be a convex function (which may not be differentiable). Then for any  $w \in \text{dom } h$ , for any  $z \in \text{int dom } (h+\psi)^*$ , the following inequality holds*

$$\tilde{D}_{h+\psi}(w, z) \geq D_h(w, \nabla(h+\psi)^*(z)).$$

*Proof of Proposition C.10.* The proof is very similar to Lemma 6 of (Flammarion & Bach, 2017), and we include for completeness. By definition of the generalized Bregman divergence (Definition C.9),

$$\tilde{D}_{h+\psi}(w, z) = (h+\psi)(w) - (h+\psi)(\nabla(h+\psi)^*(z)) - \langle z, w - \nabla(h+\psi)^*(z) \rangle.$$

By definition of the (ordinary) Bregman divergence (Definition C.8),

$$D_h(w, \nabla(h+\psi)^*(z)) = h(w) - h(\nabla(h+\psi)^*(z)) - \langle \nabla h(\nabla(h+\psi)^*(z)), w - \nabla(h+\psi)^*(z) \rangle.$$

Taking difference,

$$\tilde{D}_{h+\psi}(w, z) - D_h(w, \nabla(h+\psi)^*(z)) = \psi(w) - \psi(\nabla(h+\psi)^*(z)) - \langle z - \nabla h(\nabla(h+\psi)^*(z)), w - \nabla(h+\psi)^*(z) \rangle. \quad (\text{C.1})$$

By Proposition C.4, one has  $z \in \partial(h+\psi)(\nabla(h+\psi)^*(z))$ . Since  $h$  is differentiable in  $\text{int dom } h$ , we have (by subgradient calculus)

$$z - \nabla h(\nabla(h+\psi)^*(z)) \in \partial\psi(\nabla(h+\psi)^*(z)).$$

Therefore by the property of subgradient as the supporting hyperplane,

$$\psi(w) \geq \psi(\nabla(h+\psi)^*(z)) + \langle z - \nabla h(\nabla(h+\psi)^*(z)), w - \nabla(h+\psi)^*(z) \rangle \quad (\text{C.2})$$

Combining Eq. (C.1) and Eq. (C.2) yields

$$\tilde{D}_{h+\psi}(w, z) - D_h(w, \nabla(h+\psi)^*(z)) \geq 0,$$

completing the proof.  $\square$

## D. Proof of Theorem 4.2: Convergence of FEDDUALAVG Under Bounded Gradient Assumption

In this section, we provide a complete, non-asymptotic version of Theorem 4.2 with detailed proof.

We now formally state the assumptions of Theorem 4.2 for ease of reference.

**Assumption 3** (Bounded gradient). *In addition to Assumption 1, assume that the gradient is  $G$ -uniformly-bounded, namely*

$$\sup_{w \in \text{dom } \psi} \|\nabla f(w, \xi)\|_* \leq G$$

This is a standard assumption in analyzing classic distributed composite optimization (Duchi et al., 2012).

Before we start, we introduce a few more notations to simplify the exposition and analysis throughout this section. Let  $h_{r,k}(w) = h(w) + (rK + k)\eta_c\psi(w)$ . Let  $\bar{z}_{r,k} := \frac{1}{M} \sum_{m=1}^M z_{r,k}^m$  denote the average over clients, and  $\widehat{w}_{r,k} := \nabla h_{r,k}^*(\bar{z}_{r,k})$  denote the primal image of  $\bar{z}_{r,k}$ . Formally, we use  $\mathcal{F}_{r,k}$  to denote the  $\sigma$ -algebra generated by  $\{z_{\rho,\kappa}^m : \rho < r \text{ or } (\rho = r \text{ and } \kappa \leq k), m \in [M]\}$ .

### D.1. Main Theorem and Lemmas

Now we introduce the full version of Theorem 4.2 regarding the convergence of FEDDUALAVG with unit server learning rate  $\eta_s = 1$  under bounded gradient assumptions.

**Theorem D.1** (Detailed version of Theorem 4.2). *Assume Assumption 3, then for any initialization  $w_0 \in \text{dom } \psi$ , for unit server learning rate  $\eta_s = 1$  and any client learning rate  $\eta_c \leq \frac{1}{4L}$ , FEDDUALAVG yields*

$$\mathbb{E} \left[ \Phi \left( \frac{1}{KR} \sum_{r=0}^{R-1} \sum_{k=1}^K \widehat{w}_{r,k} \right) - \Phi(w^*) \right] \leq \frac{B}{\eta_c KR} + \frac{\eta_c \sigma^2}{M} + 4\eta_c^2 L(K-1)^2 G^2, \quad (\text{D.1})$$

where  $B := D_h(w^*, w_0)$  is the Bregman divergence between the optimal  $w^*$  and the initial  $w_0$ .

Particularly for

$$\eta_c = \min \left\{ \frac{1}{4L}, \frac{M^{\frac{1}{2}} B^{\frac{1}{2}}}{\sigma K^{\frac{1}{2}} R^{\frac{1}{2}}}, \frac{B^{\frac{1}{3}}}{L^{\frac{1}{3}} K R^{\frac{1}{3}} G^{\frac{2}{3}}} \right\},$$

one has

$$\mathbb{E} \left[ \Phi \left( \frac{1}{KR} \sum_{r=0}^{R-1} \sum_{k=1}^K \widehat{w}_{r,k} \right) - \Phi(w^*) \right] \leq \frac{4LB}{KR} + \frac{2\sigma B^{\frac{1}{2}}}{M^{\frac{1}{2}} K^{\frac{1}{2}} R^{\frac{1}{2}}} + \frac{5L^{\frac{1}{3}} B^{\frac{2}{3}} G^{\frac{2}{3}}}{R^{\frac{2}{3}}}.$$

The proof of Theorem D.1 is based on the following two lemmas regarding perturbed convergence and stability respectively.

**Lemma D.2** (Perturbed iterate analysis of FEDDUALAVG). *Assume Assumption 1, then for any initialization  $w_0 \in \text{dom } \psi$ , for any reference point  $w \in \text{dom } \psi$ , for  $\eta_s = 1$ , for any  $\eta_c \leq \frac{1}{4L}$ , FEDDUALAVG yields*

$$\mathbb{E} \left[ \Phi \left( \frac{1}{KR} \sum_{r=0}^{R-1} \sum_{k=1}^K \widehat{w}_{r,k} \right) - \Phi(w) \right] \leq \frac{1}{\eta_c KR} D_h(w, w_0) + \frac{\eta_c \sigma^2}{M} + \frac{L}{MKR} \left[ \sum_{r=0}^{R-1} \sum_{k=0}^{K-1} \sum_{m=1}^M \mathbb{E} \|\bar{z}_{r,k} - z_{r,k}^m\|_*^2 \right].$$

Lemma D.2 decomposes the convergence of FEDDUALAVG into two parts: the first part  $\frac{1}{\eta_c KR} D_h(w, w_0) + \frac{\eta_c \sigma^2}{2M} + \frac{L}{MKR}$  corresponds to the convergence rate if all clients were synchronized every iteration. The second part  $\frac{L}{MKR} \sum_{r=0}^{R-1} \sum_{k=0}^{K-1} \sum_{m=1}^M \mathbb{E} \|z_{r,k}^m - \bar{z}_{r,k}\|_*^2$  characterizes the stability of the algorithm. Note that Lemma D.2 only assumes the blanket Assumption 1. We defer the proof of Lemma D.2 to Appendix D.2.

The following Lemma D.3 bounds the stability term under the additional bounded gradient assumptions.

**Lemma D.3** (Stability of FEDDUALAVG under bounded gradient assumption). *In the same settings of Theorem D.1, it is the case that*

$$\frac{1}{M} \sum_{m=1}^M \mathbb{E} \|z_{r,k}^m - \bar{z}_{r,k}\|_*^2 \leq 4\eta_c^2 (K-1)^2 G^2.$$

We defer the proof of Lemma D.3 to Appendix D.3. With Lemmas D.2 and D.3 at hands the proof of Theorem D.1 is immediate.

*Proof of Theorem D.1.* Eq. (D.1) follows immediately from Lemmas D.2 and D.3 by putting  $w = w^*$  in Lemma D.2.

Now put

$$\eta_c = \min \left\{ \frac{1}{4L}, \frac{M^{\frac{1}{2}} B^{\frac{1}{2}}}{\sigma K^{\frac{1}{2}} R^{\frac{1}{2}}}, \frac{B^{\frac{1}{3}}}{L^{\frac{1}{3}} K R^{\frac{1}{3}} G^{\frac{2}{3}}} \right\},$$

which yields

$$\frac{B}{\eta_c K R} = \max \left\{ \frac{4LB}{KR}, \frac{\sigma B^{\frac{1}{2}}}{M^{\frac{1}{2}} K^{\frac{1}{2}} R^{\frac{1}{2}}}, \frac{L^{\frac{1}{3}} B^{\frac{2}{3}} G^{\frac{2}{3}}}{R^{\frac{2}{3}}} \right\} \leq \frac{4LB}{KR} + \frac{\sigma B^{\frac{1}{2}}}{M^{\frac{1}{2}} K^{\frac{1}{2}} R^{\frac{1}{2}}} + \frac{L^{\frac{1}{3}} B^{\frac{2}{3}} G^{\frac{2}{3}}}{R^{\frac{2}{3}}},$$

and

$$\frac{\eta_c \sigma^2}{2M} \leq \frac{M^{\frac{1}{2}} B^{\frac{1}{2}}}{\sigma T^{\frac{1}{2}}} \cdot \frac{\sigma^2}{2M} = \frac{\sigma B^{\frac{1}{2}}}{2M^{\frac{1}{2}} K^{\frac{1}{2}} R^{\frac{1}{2}}}, \quad 4\eta_c^2 L K^2 G^2 \leq 4 \left( \frac{B^{\frac{1}{3}}}{L^{\frac{1}{3}} K R^{\frac{1}{3}} G^{\frac{2}{3}}} \right)^2 L K^2 G^2 = \frac{4L^{\frac{1}{3}} B^{\frac{2}{3}} G^{\frac{2}{3}}}{R^{\frac{2}{3}}}.$$

Summarizing the above three inequalities completes the proof of Theorem D.1.  $\square$

## D.2. Perturbed Iterate Analysis of FEDDUALAVG: Proof of Lemma D.2

In this subsection, we prove Lemma D.2. We start by showing the following Proposition D.4 regarding the one step improvement of the shadow sequence  $\overline{z}_{r,k}$ .

**Proposition D.4** (One step analysis of FEDDUALAVG). *Under the same assumptions of Lemma D.2, the following inequality holds*

$$\begin{aligned} \mathbb{E} \left[ \tilde{D}_{h_{r,k+1}}(w, \overline{z}_{r,k+1}) \middle| \mathcal{F}_{r,k} \right] &\leq \tilde{D}_{h_{r,k}}(w, \overline{z}_{r,k}) - \eta_c \mathbb{E} \left[ \Phi(\widehat{w}_{r,k+1}) - \Phi(w) \middle| \mathcal{F}_{r,k} \right] \\ &\quad + \eta_c L \cdot \mathbb{E} \left[ \frac{1}{M} \sum_{m=1}^M \|\overline{z}_{r,k} - z_{r,k}^m\|_*^2 \middle| \mathcal{F}_{r,k} \right] + \frac{\eta_c^2 \sigma^2}{M}, \end{aligned}$$

where  $\tilde{D}$  is the generalized Bregman divergence defined in Definition C.9.

The proof of Proposition D.4 relies on the following two claims regarding the deterministic analysis of FEDDUALAVG. We defer the proof of Claims D.5 and D.6 to Appendices D.2.1 and D.2.2, respectively.

**Claim D.5.** *Under the same assumptions of Lemma D.2, the following inequality holds*

$$\begin{aligned} &\tilde{D}_{h_{r,k+1}}(w, \overline{z}_{r,k+1}) \\ &= \tilde{D}_{h_{r,k}}(w, \overline{z}_{r,k}) - \tilde{D}_{h_{r,k}}(\widehat{w}_{r,k+1}, \overline{z}_{r,k}) - \eta_c (\psi(\widehat{w}_{r,k+1}) - \psi(w)) - \eta_c \left\langle \frac{1}{M} \sum_{m=1}^M \nabla f(w_{r,k}^m; \xi_{r,k}^m), \widehat{w}_{r,k+1} - w \right\rangle. \end{aligned} \quad (\text{D.2})$$

**Claim D.6.** *Under the same assumptions of Lemma D.2, it is the case that*

$$\begin{aligned} F(\widehat{w}_{r,k+1}) - F(w) &\leq \left\langle \frac{1}{M} \sum_{m=1}^M \nabla f(w_{r,k}^m; \xi_{r,k}^m), \widehat{w}_{r,k+1} - w \right\rangle \\ &\quad + \left\langle \frac{1}{M} \sum_{m=1}^M (\nabla F_m(w_{r,k}^m) - \nabla f(w_{r,k}^m; \xi_{r,k}^m)), \widehat{w}_{r,k+1} - w \right\rangle + L \|\widehat{w}_{r,k+1} - \widehat{w}_{r,k}\|^2 + \frac{L}{M} \sum_{m=1}^M \|\overline{z}_{r,k} - z_{r,k}^m\|_*^2. \end{aligned} \quad (\text{D.3})$$

With Claims D.5 and D.6 at hand we are ready to prove the one step analysis Proposition D.4.

*Proof of Proposition D.4.* Applying Claims D.5 and D.6 yields (summing Eq. (D.2) with  $\eta_c$  times of Eq. (D.3)),

$$\begin{aligned} \tilde{D}_{h_{r,k+1}}(w, \overline{z_{r,k+1}}) &\leq \tilde{D}_{h_{r,k}}(w, \overline{z_{r,k}}) - \tilde{D}_{h_{r,k}}(\widehat{w_{r,k+1}}, \overline{z_{r,k}}) + \eta_c L \|\widehat{w_{r,k+1}} - \widehat{w_{r,k}}\|^2 - \eta_c (\Phi(\widehat{w_{r,k+1}}) - \Phi(w)) \\ &\quad + \eta_c \left\langle \frac{1}{M} \sum_{m=1}^M (\nabla F_m(w_{r,k}^m) - \nabla f(w_{r,k}^m; \xi_{r,k}^m)), \widehat{w_{r,k+1}} - w \right\rangle \\ &\quad + \eta_c L \cdot \frac{1}{M} \sum_{m=1}^M \|\overline{z_{r,k}} - z_{r,k}^m\|_*^2. \end{aligned} \quad (\text{D.4})$$

Note that

$$\tilde{D}_{h_{r,k}}(\widehat{w_{r,k+1}}, \overline{z_{r,k}}) \geq D_h(\widehat{w_{r,k+1}}, \nabla h_{r,k}^*(\overline{z_{r,k}})) = D_h(\widehat{w_{r,k+1}}, \widehat{w_{r,k}}) \geq \frac{1}{2} \|\widehat{w_{r,k+1}} - \widehat{w_{r,k}}\|^2,$$

and

$$\eta_c L \|\widehat{w_{r,k+1}} - \widehat{w_{r,k}}\|^2 \leq \frac{1}{4} \|\widehat{w_{r,k+1}} - \widehat{w_{r,k}}\|^2,$$

since  $\eta_c \leq \frac{1}{4L}$  by assumption. Therefore,

$$-\tilde{D}_{h_{r,k}}(\widehat{w_{r,k+1}}, \overline{z_{r,k}}) + \eta_c L \|\widehat{w_{r,k+1}} - \widehat{w_{r,k}}\|^2 \leq -\frac{1}{4} \|\widehat{w_{r,k+1}} - \widehat{w_{r,k}}\|^2. \quad (\text{D.5})$$

Plugging Eq. (D.5) to Eq. (D.4) gives

$$\begin{aligned} \tilde{D}_{h_{r,k+1}}(w, \overline{z_{r,k+1}}) &\leq \tilde{D}_{h_{r,k}}(w, \overline{z_{r,k}}) - \frac{1}{4} \|\widehat{w_{r,k+1}} - \widehat{w_{r,k}}\|^2 - \eta_c (\Phi(\widehat{w_{r,k+1}}) - \Phi(w)) \\ &\quad + \eta_c \left\langle \frac{1}{M} \sum_{m=1}^M (\nabla F_m(w_{r,k}^m) - \nabla f(w_{r,k}^m; \xi_{r,k}^m)), \widehat{w_{r,k+1}} - w \right\rangle \\ &\quad + \eta_c L \cdot \frac{1}{M} \sum_{m=1}^M \|\overline{z_{r,k}} - z_{r,k}^m\|_*^2. \end{aligned} \quad (\text{D.6})$$

Now we take conditional expectation. Note that

$$\begin{aligned} &\mathbb{E} \left[ \left\langle \frac{1}{M} \sum_{m=1}^M \nabla F_m(w_{r,k}^m) - \nabla f(w_{r,k}^m; \xi_{r,k}^m), \widehat{w_{r,k+1}} - w \right\rangle \middle| \mathcal{F}_{r,k} \right] \\ &= \mathbb{E} \left[ \left\langle \frac{1}{M} \sum_{m=1}^M \nabla F_m(w_{r,k}^m) - \nabla f(w_{r,k}^m; \xi_{r,k}^m), \widehat{w_{r,k+1}} - \widehat{w_{r,k}} \right\rangle \middle| \mathcal{F}_{r,k} \right] \\ &\quad \text{(since } \mathbb{E}_{\xi_{r,k}^m \sim \mathcal{D}_m} [\nabla f(w_{r,k}^m; \xi_{r,k}^m)] = \nabla F_m(w_{r,k}^m)\text{)} \\ &\leq \mathbb{E} \left[ \left\| \frac{1}{M} \sum_{m=1}^M \nabla F_m(w_{r,k}^m) - \nabla f(w_{r,k}^m; \xi_{r,k}^m) \right\|_* \middle| \mathcal{F}_{r,k} \right] \cdot \mathbb{E} [\|\widehat{w_{r,k+1}} - \widehat{w_{r,k}}\| \middle| \mathcal{F}_{r,k}] \quad \text{(by definition of dual norm } \|\cdot\|_*\text{)} \\ &\leq \frac{\sigma}{\sqrt{M}} \mathbb{E} [\|\widehat{w_{r,k+1}} - \widehat{w_{r,k}}\| \middle| \mathcal{F}_{r,k}]. \quad \text{(by bounded variance assumption and independence)} \end{aligned}$$

Plugging the above inequality to Eq. (D.6) gives

$$\begin{aligned} &\mathbb{E} \left[ \tilde{D}_{h_{r,k+1}}(w, \overline{z_{r,k+1}}) \middle| \mathcal{F}_{r,k} \right] \\ &\leq \tilde{D}_{h_{r,k}}(w, \overline{z_{r,k}}) - \eta_c \mathbb{E} [\Phi(\widehat{w_{r,k+1}}) - \Phi(w) \middle| \mathcal{F}_{r,k}] + \eta_c L \cdot \mathbb{E} \left[ \frac{1}{M} \sum_{m=1}^M \|\overline{z_{r,k}} - z_{r,k}^m\|_*^2 \middle| \mathcal{F}_{r,k} \right] \\ &\quad + \frac{\eta_c \sigma}{\sqrt{M}} \mathbb{E} [\|\widehat{w_{r,k+1}} - \widehat{w_{r,k}}\| \middle| \mathcal{F}_{r,k}] - \frac{1}{4} \mathbb{E} [\|\widehat{w_{r,k+1}} - \widehat{w_{r,k}}\|^2 \middle| \mathcal{F}_{r,k}] \\ &\leq \tilde{D}_{h_{r,k}}(w, \overline{z_{r,k}}) - \eta_c \mathbb{E} [\Phi(\widehat{w_{r,k+1}}) - \Phi(w) \middle| \mathcal{F}_{r,k}] + \eta_c L \cdot \mathbb{E} \left[ \frac{1}{M} \sum_{m=1}^M \|\overline{z_{r,k}} - z_{r,k}^m\|_*^2 \middle| \mathcal{F}_{r,k} \right] + \frac{\eta_c^2 \sigma^2}{M}, \\ &\quad \text{(by quadratic maximum)} \end{aligned}$$

completing the proof of Proposition D.4.  $\square$

The Lemma D.2 then follows by telescoping the one step analysis Proposition D.4.

*Proof of Lemma D.2.* Let us first telescope Proposition D.4 within the same round  $r$ , from  $k = 0$  to  $K$ , which gives

$$\begin{aligned} \mathbb{E} \left[ \tilde{D}_{h_{r,K}}(w, \overline{z_{r,K}}) \middle| \mathcal{F}_{r,0} \right] &\leq \tilde{D}_{h_{r,0}}(w, \overline{z_{r,0}}) - \eta_c \sum_{k=1}^K \mathbb{E} [\Phi(\widehat{w_{r,k}}) - \Phi(w) | \mathcal{F}_{r,0}] \\ &\quad + \eta_c L \cdot \mathbb{E} \left[ \frac{1}{M} \sum_{k=0}^{K-1} \sum_{m=1}^M \|\overline{z_{r,k}} - z_{r,k}^m\|_*^2 \middle| \mathcal{F}_{r,0} \right] + \frac{\eta_c^2 K \sigma^2}{M}. \end{aligned}$$

Since server learning rate  $\eta_s = 1$  we have  $\overline{z_{r,K}} = \overline{z_{r+1,0}}$ . Therefore we can telescope the round from  $r = 0$  to  $R$ , which gives

$$\begin{aligned} \mathbb{E} \left[ \tilde{D}_{h_{R,0}}(w, \overline{z_{R,0}}) \right] &\leq \tilde{D}_{h_{0,0}}(w, \overline{z_{0,0}}) - \eta_c \sum_{r=0}^{R-1} \sum_{k=1}^K \mathbb{E} [\Phi(\widehat{w_{r,k}}) - \Phi(w)] \\ &\quad + \eta_c L \cdot \mathbb{E} \left[ \frac{1}{M} \sum_{r=0}^{R-1} \sum_{k=0}^{K-1} \sum_{m=1}^M \|\overline{z_{r,k}} - z_{r,k}^m\|_*^2 \right] + \frac{\eta_c^2 K R \sigma^2}{M}. \end{aligned}$$

Dividing both sides by  $\eta_c \cdot KR$  and rearranging

$$\begin{aligned} \frac{1}{KR} \sum_{r=0}^{R-1} \sum_{k=1}^K \mathbb{E} [\Phi(\widehat{w_{r,k}}) - \Phi(w)] &\leq \frac{1}{\eta_c KR} \left( \tilde{D}_{h_{0,0}}(w, \overline{z_{0,0}}) - \mathbb{E} \left[ \tilde{D}_{h_{R,0}}(w, \overline{z_{R,0}}) \right] \right) \\ &\quad + L \cdot \mathbb{E} \left[ \frac{1}{MKR} \sum_{r=0}^{R-1} \sum_{k=0}^{K-1} \sum_{m=1}^M \|\overline{z_{r,k}} - z_{r,k}^m\|_*^2 \right] + \frac{\eta_c \sigma^2}{M}. \end{aligned}$$

Applying Jensen's inequality on the LHS and dropping the negative term on the RHS yield

$$\mathbb{E} \left[ \Phi \left( \frac{1}{KR} \sum_{r=0}^{R-1} \sum_{k=1}^K \widehat{w_{r,k}} \right) - \Phi(w) \right] \leq \frac{1}{\eta_c KR} \tilde{D}_{h_{0,0}}(w, \overline{z_{0,0}}) + \frac{L}{MKR} \left[ \sum_{r=0}^{R-1} \sum_{k=0}^{K-1} \sum_{m=1}^M \mathbb{E} \|\overline{z_{r,k}} - z_{r,k}^m\|_*^2 \right] + \frac{\eta_c \sigma^2}{M}. \quad (\text{D.7})$$

Since  $\overline{z_{0,0}} = \nabla h(w_0)$  and  $w_0 \in \text{dom } \psi$ , we have  $\nabla h_{0,0}^*(\nabla h(w_0)) = w_0$  by Proposition C.7 since  $h$  is assumed to be of Legendre type. Consequently

$$\begin{aligned} \tilde{D}_{h_{0,0}}(w, \overline{z_{0,0}}) &= h(w) - h(\nabla h_{0,0}^*(\nabla h(w_0))) - \langle z_0, w - \nabla h_{0,0}^*(\nabla h(w_0)) \rangle \\ &= h(w) - h(w_0) - \langle \nabla h(w_0), w - w_0 \rangle = D_h(w, w_0). \end{aligned} \quad (\text{D.8})$$

Plugging Eq. (D.8) back to Eq. (D.7) completes the proof of Lemma D.2.  $\square$

### D.2.1. DEFERRED PROOF OF CLAIM D.5

*Proof of Claim D.5.* By definition of FEDDUALAVG procedure, for all  $m \in [M]$ ,  $k \in \{0, 1, \dots, K-1\}$ , we have

$$z_{r,k+1}^m = z_{r,k}^m - \eta_c \nabla f(w_{r,k}^m; \xi_{r,k}^m).$$

Taking average over  $m \in [M]$  gives (recall the overline denotes the average over clients)

$$\overline{z_{r,k+1}} = \overline{z_{r,k}} - \eta_c \cdot \frac{1}{M} \sum_{m=1}^M \nabla f(w_{r,k}^m; \xi_{r,k}^m). \quad (\text{D.9})$$

Now we study generalized Bregman divergence  $\tilde{D}_{h,k+1}(w, \overline{z_{r,k+1}})$  for any arbitrary pre-fixed  $w \in \text{dom } h_{r,k}$

$$\begin{aligned}
 & \tilde{D}_{h_{r,k+1}}(w, \overline{z_{r,k+1}}) \\
 &= h_{r,k+1}(w) - h_{r,k+1}(\nabla h_{r,k+1}^*(\overline{z_{r,k+1}})) - \langle \overline{z_{r,k+1}}, w - \nabla h_{r,k+1}^*(\overline{z_{r,k+1}}) \rangle && \text{(By definition of } \tilde{D}) \\
 &= h_{r,k+1}(w) - h_{r,k+1}(\widehat{w_{r,k+1}}) - \langle \overline{z_{r,k+1}}, w - \widehat{w_{r,k+1}} \rangle && \text{(By definition of } \widehat{w_{r,k+1}}) \\
 &= h_{r,k+1}(w) - h_{r,k+1}(\widehat{w_{r,k+1}}) - \left\langle \overline{z_{r,k}} - \eta_c \cdot \frac{1}{M} \sum_{m=1}^M \nabla f(w_{r,k}^m; \xi_{r,k}^m), w - \widehat{w_{r,k+1}} \right\rangle && \text{(By Eq. (D.9))} \\
 &= (h_{r,k}(w) + \eta_c \psi(w)) - (h_{r,k}(\widehat{w_{r,k+1}}) + \eta_c \psi(\widehat{w_{r,k+1}})) - \left\langle \overline{z_{r,k}} - \eta_c \cdot \frac{1}{M} \sum_{m=1}^M \nabla f(w_{r,k}^m; \xi_{r,k}^m), w - \widehat{w_{r,k+1}} \right\rangle \\
 & \hspace{15em} \text{(Since } h_{r,k+1} = h_{r,k} + \eta_c \psi \text{ by definition of } h_{r,k+1}) \\
 &= [h_{r,k}(w) - h_{r,k}(\widehat{w_{r,k}}) - \langle \overline{z_{r,k}}, w - \widehat{w_{r,k}} \rangle] - [h_{r,k}(\widehat{w_{r,k+1}}) - h_{r,k}(\widehat{w_{r,k}}) - \langle \overline{z_{r,k}}, \widehat{w_{r,k+1}} - \widehat{w_{r,k}} \rangle] \\
 & \quad - \eta_c (\psi(\widehat{w_{r,k+1}}) - \psi(w)) - \eta_c \left\langle \frac{1}{M} \sum_{m=1}^M \nabla f(w_{r,k}^m; \xi_{r,k}^m), \widehat{w_{r,k+1}} - w \right\rangle && \text{(Rearranging)} \\
 &= \tilde{D}_{h_{r,k}}(w, \overline{z_{r,k}}) - \tilde{D}_{h_{r,k}}(\widehat{w_{r,k+1}}, \overline{z_{r,k}}) - \eta_c (\psi(\widehat{w_{r,k+1}}) - \psi(w)) - \eta_c \left\langle \frac{1}{M} \sum_{m=1}^M \nabla f(w_{r,k}^m; \xi_{r,k}^m), \widehat{w_{r,k+1}} - w \right\rangle,
 \end{aligned}$$

where the last equality is by definition of  $\tilde{D}$ .  $\square$

#### D.2.2. DEFERRED PROOF OF CLAIM D.6

*Proof of Claim D.6.* By smoothness and convexity of  $F_m$ , we know

$$\begin{aligned}
 F_m(\widehat{w_{r,k+1}}) &\leq F_m(w_{r,k}^m) + \langle \nabla F_m(w_{r,k}^m), \widehat{w_{r,k+1}} - w_{r,k}^m \rangle + \frac{L}{2} \|\widehat{w_{r,k+1}} - w_{r,k}^m\|^2 && \text{(smoothness)} \\
 &\leq F_m(w) + \langle \nabla F_m(w_{r,k}^m), \widehat{w_{r,k+1}} - w \rangle + \frac{L}{2} \|\widehat{w_{r,k+1}} - w_{r,k}^m\|^2. && \text{(convexity)}
 \end{aligned}$$

Taking summation over  $m$  gives

$$\begin{aligned}
 F(\widehat{w_{r,k+1}}) - F(w) &= \frac{1}{M} \sum_{m=1}^M (F_m(\widehat{w_{r,k+1}}) - F_m(w)) \\
 &\leq \left\langle \frac{1}{M} \sum_{m=1}^M \nabla F_m(w_{r,k}^m), \widehat{w_{r,k+1}} - w \right\rangle + \frac{L}{2M} \sum_{m=1}^M \|\widehat{w_{r,k+1}} - w_{r,k}^m\|^2 \\
 &= \left\langle \frac{1}{M} \sum_{m=1}^M \nabla f(w_{r,k}^m; \xi_{r,k}^m), \widehat{w_{r,k+1}} - w \right\rangle + \left\langle \frac{1}{M} \sum_{m=1}^M (\nabla F_m(w_{r,k}^m) - \nabla f(w_{r,k}^m; \xi_{r,k}^m)), \widehat{w_{r,k+1}} - w \right\rangle \\
 & \quad + \frac{L}{2M} \sum_{m=1}^M \|\widehat{w_{r,k+1}} - w_{r,k}^m\|^2 \\
 &\leq \left\langle \frac{1}{M} \sum_{m=1}^M \nabla f(w_{r,k}^m; \xi_{r,k}^m), \widehat{w_{r,k+1}} - w \right\rangle + \left\langle \frac{1}{M} \sum_{m=1}^M (\nabla F_m(w_{r,k}^m) - \nabla f(w_{r,k}^m; \xi_{r,k}^m)), \widehat{w_{r,k+1}} - w \right\rangle \\
 & \quad + L \|\widehat{w_{r,k+1}} - \widehat{w_{r,k}}\|^2 + \frac{L}{M} \sum_{m=1}^M \|\widehat{w_{r,k}} - w_{r,k}^m\|^2, && \text{(D.10)}
 \end{aligned}$$

where in the last inequality we applied the triangle inequality (for an arbitrary norm  $\|\cdot\|$ ):

$$\|\widehat{w_{r,k+1}} - w_{r,k}^m\|^2 \leq (\|\widehat{w_{r,k+1}} - \widehat{w_{r,k}}\| + \|\widehat{w_{r,k}} - w_{r,k}^m\|)^2 \leq 2\|\widehat{w_{r,k+1}} - \widehat{w_{r,k}}\|^2 + 2\|\widehat{w_{r,k}} - w_{r,k}^m\|^2.$$



Since  $\psi$  is convex and  $h$  is 1-strongly-convex according to Assumption 1, we know that  $h_{r,k} = h + \eta_c(rK + k)\psi$  is also 1-strongly-convex. Therefore  $h_{r,k}^*$  is 1-smooth by Proposition C.5. Consequently,

$$\|w_{r,k}^m - \widehat{w}_{r,k}\|^2 = \|\nabla h_{r,k}^*(z_{r,k}^m) - \nabla h_{r,k}^*(\overline{z}_{r,k})\|^2 \leq \|z_{r,k}^m - \overline{z}_{r,k}\|_*^2, \quad (\text{D.11})$$

where the first equality is by definition of  $w_{r,k}^m$  and  $\widehat{w}_{r,k}$  and the second inequality is by 1-smoothness. Plugging Eq. (D.11) back to Eq. (D.10) completes the proof of Claim D.6.  $\square$

### D.3. Stability of FEDDUALAVG Under Bounded Gradient Assumptions: Proof of Lemma D.3

The proof of Lemma D.3 is straightforward given the assumption of bounded gradient and the fact that  $z_{r,0}^{m_1} = z_{r,0}^{m_2}$  for all  $m_1, m_2 \in [M]$ .

*Proof of Lemma D.3.* Let  $m_1, m_2 \in [M]$  be two arbitrary clients, then

$$\begin{aligned} \mathbb{E} \left[ \|z_{r,k}^{m_1} - z_{r,k}^{m_2}\|_*^2 \middle| \mathcal{F}_{r,0} \right] &= \eta_c^2 \mathbb{E} \left[ \left\| \sum_{\kappa=0}^{k-1} (\nabla f(w_{r,\kappa}^{m_1}; \xi_{r,\kappa}^{m_1}) - \nabla f(w_{r,\kappa}^{m_2}; \xi_{r,\kappa}^{m_2})) \right\|_*^2 \middle| \mathcal{F}_{r,0} \right] \quad (\text{since } z_{r,0}^{m_1} = z_{r,0}^{m_2}) \\ &\leq \eta_c^2 \mathbb{E} \left[ \left( \sum_{\kappa=0}^{k-1} \|\nabla f(w_{r,\kappa}^{m_1}; \xi_{r,\kappa}^{m_1})\|_* + \sum_{\kappa=0}^{k-1} \|\nabla f(w_{r,\kappa}^{m_2}; \xi_{r,\kappa}^{m_2})\|_* \right)^2 \middle| \mathcal{F}_{r,0} \right] \\ &\hspace{15em} (\text{triangle inequality of } \|\cdot\|_*) \\ &\leq \eta_c^2 (2(k-1)G)^2 = 4\eta_c^2 (K-1)^2 G^2. \end{aligned}$$

By convexity of  $\|\cdot\|_*$ ,

$$\frac{1}{M} \sum_{m=1}^M \mathbb{E} \|z_{r,k}^m - \overline{z}_{r,k}\|_*^2 \leq \mathbb{E} \|z_{r,k}^{m_1} - z_{r,k}^{m_2}\|_*^2 \leq 4\eta_c^2 (K-1)^2 G^2,$$

completing the proof of Lemma D.3.  $\square$

## E. Proof of Theorem 4.3: Convergence of FEDDUALAVG Under Bounded Heterogeneity and Quadratic Assumptions

In this section, we study the convergence of FEDDUALAVG under Assumption 2 (quadraticness) with unit server learning rate  $\eta_s = 1$ . We provide a complete, non-asymptotic version of Theorem 4.3 with detailed proof, which expands the proof sketch in Section 4.3. We will reuse the notations ( $\overline{z}_{r,k}$ ,  $\widehat{w}_{r,k}$ , etc.) introduced at the beginning of Appendix D.

### E.1. Main Theorem and Lemmas

Now we state the full version of Theorem 4.3 on FEDDUALAVG with unit server learning rate  $\eta_s = 1$  under quadratic assumptions.

**Theorem E.1** (Detailed version of Theorem 4.3). *Assuming Assumption 2, then for any initialization  $w_0 \in \text{dom } \psi$ , for unit server learning rate  $\eta_s = 1$  and any client learning rate  $\eta_c \leq \frac{1}{4L}$ , FEDDUALAVG yields*

$$\mathbb{E} \left[ \Phi \left( \frac{1}{KR} \sum_{r=0}^{R-1} \sum_{k=1}^K \widehat{w}_{r,k} \right) - \Phi(w^*) \right] \leq \frac{B}{\eta_c KR} + \frac{\eta_c \sigma^2}{M} + 7\eta_c^2 LK \sigma^2 + 14\eta_c^2 LK^2 \zeta^2,$$

where  $B := D_h(w^*, w_0)$  is the Bregman divergence between the optimal  $w^*$  and the initialization  $w_0$ .

Particularly for

$$\eta_c = \min \left\{ \frac{1}{4L}, \frac{M^{\frac{1}{2}} B^{\frac{1}{2}}}{\sigma K^{\frac{1}{2}} R^{\frac{1}{2}}}, \frac{B^{\frac{1}{3}}}{L^{\frac{1}{3}} K^{\frac{2}{3}} R^{\frac{1}{3}} \sigma^{\frac{2}{3}}}, \frac{B^{\frac{1}{3}}}{L^{\frac{1}{3}} K R^{\frac{1}{3}} \zeta^{\frac{2}{3}}} \right\},$$

we have

$$\mathbb{E} \left[ \Phi \left( \frac{1}{KR} \sum_{r=0}^{R-1} \sum_{k=1}^K \widehat{w}_{r,k} \right) - \Phi(w^*) \right] \leq \frac{4LB}{KR} + \frac{2\sigma B^{\frac{1}{2}}}{M^{\frac{1}{2}} K^{\frac{1}{2}} R^{\frac{1}{2}}} + \frac{8L^{\frac{1}{3}} B^{\frac{2}{3}} \sigma^{\frac{2}{3}}}{K^{\frac{1}{3}} R^{\frac{2}{3}}} + \frac{15L^{\frac{1}{3}} B^{\frac{2}{3}} \zeta^{\frac{2}{3}}}{R^{\frac{2}{3}}}.$$

The proof of Theorem E.1 relies on the perturbed iterate analysis Lemma D.2 of FEDDUALAVG and a stability bound for quadratic objectives, as stated below in Lemma E.2. Note that Lemma D.2 only assumes Assumption 1 and therefore applicable to Theorem E.1.

**Lemma E.2.** *In the same settings of Theorem E.1, the following inequality holds for any  $k \in \{0, 1, \dots, K\}$  and  $r \in \{0, 1, \dots, R\}$ ,*

$$\frac{1}{M} \sum_{m=1}^M \mathbb{E} \left[ \|\overline{z}_{r,k} - z_{r,k}^m\|_*^2 \right] \leq 7\eta_c^2 K \sigma^2 + 14\eta_c^2 K^2 \zeta^2.$$

The proof of Lemma E.2 is deferred to Appendix E.2. With Lemma E.2 at hand we are ready to prove Theorem E.1.

*Proof of Theorem E.1.* Applying Lemmas D.2 and E.2 one has

$$\begin{aligned} & \mathbb{E} \left[ \Phi \left( \frac{1}{KR} \sum_{r=0}^{R-1} \sum_{k=1}^K \widehat{w}_{r,k} \right) - \Phi(w^*) \right] \\ & \leq \frac{1}{\eta_c KR} D_h(w^*, w_0) + \frac{\eta_c \sigma^2}{M} + \frac{L}{MKR} \left[ \sum_{r=0}^{R-1} \sum_{k=0}^{K-1} \sum_{m=1}^M \mathbb{E} \|\overline{z}_{r,k} - z_{r,k}^m\|_*^2 \right] \quad (\text{by Lemma D.2}) \\ & \leq \frac{1}{\eta_c KR} D_h(w^*, w_0) + \frac{\eta_c \sigma^2}{M} + L \cdot (7\eta_c^2 K \sigma^2 + 14\eta_c^2 K^2 \zeta^2) \quad (\text{by Lemma E.2}) \\ & = \frac{B}{\eta_c KR} + \frac{\eta_c \sigma^2}{M} + 7\eta_c^2 LK \sigma^2 + 14\eta_c^2 LK^2 \zeta^2, \end{aligned}$$

which gives the first inequality in Theorem E.1.

Now set

$$\eta_c = \min \left\{ \frac{1}{4L}, \frac{M^{\frac{1}{2}} B^{\frac{1}{2}}}{\sigma K^{\frac{1}{2}} R^{\frac{1}{2}}}, \frac{B^{\frac{1}{3}}}{L^{\frac{1}{3}} K^{\frac{2}{3}} R^{\frac{1}{3}} \sigma^{\frac{2}{3}}}, \frac{B^{\frac{1}{3}}}{L^{\frac{1}{3}} K R^{\frac{1}{3}} \zeta^{\frac{2}{3}}} \right\}.$$

We have

$$\frac{B}{\eta_c KR} \leq \max \left\{ \frac{4LB}{KR}, \frac{\sigma B^{\frac{1}{2}}}{M^{\frac{1}{2}} K^{\frac{1}{2}} R^{\frac{1}{2}}}, \frac{L^{\frac{1}{3}} B^{\frac{2}{3}} \sigma^{\frac{2}{3}}}{K^{\frac{1}{3}} R^{\frac{2}{3}}}, \frac{L^{\frac{1}{3}} B^{\frac{2}{3}} \zeta^{\frac{2}{3}}}{R^{\frac{2}{3}}} \right\},$$

and

$$\frac{\eta_c \sigma^2}{M} \leq \frac{\sigma B^{\frac{1}{2}}}{M^{\frac{1}{2}} K^{\frac{1}{2}} R^{\frac{1}{2}}}, \quad 7\eta_c^2 LK \sigma^2 \leq \frac{7L^{\frac{1}{3}} B^{\frac{2}{3}} \sigma^{\frac{2}{3}}}{K^{\frac{1}{3}} R^{\frac{2}{3}}}, \quad 14\eta_c^2 LK^2 \zeta^2 \leq \frac{14L^{\frac{1}{3}} B^{\frac{2}{3}} \zeta^{\frac{2}{3}}}{R^{\frac{2}{3}}}.$$

Consequently

$$\frac{B}{\eta_c KR} + \frac{\eta_c \sigma^2}{M} + 7\eta_c^2 LK \sigma^2 + 14\eta_c^2 LK^2 \zeta^2 \leq \frac{4LB}{KR} + \frac{2\sigma B^{\frac{1}{2}}}{M^{\frac{1}{2}} K^{\frac{1}{2}} R^{\frac{1}{2}}} + \frac{8L^{\frac{1}{3}} B^{\frac{2}{3}} \sigma^{\frac{2}{3}}}{K^{\frac{1}{3}} R^{\frac{2}{3}}} + \frac{15L^{\frac{1}{3}} B^{\frac{2}{3}} \zeta^{\frac{2}{3}}}{R^{\frac{2}{3}}},$$

completing the proof of Theorem E.1.  $\square$

## E.2. Stability of FEDDUALAVG Under Quadratic Assumptions: Proof of Lemma E.2

In this subsection, we prove Lemma E.2 on the stability of FEDDUALAVG for quadratic  $F$ . We first state and prove the following Proposition E.3 on the one-step analysis of stability.

**Proposition E.3.** *In the same settings of Theorem E.1, let  $m_1, m_2 \in [M]$  be two arbitrary clients. Then the following inequality holds*

$$\mathbb{E} \left[ \left\| z_{r,k+1}^{m_1} - z_{r,k+1}^{m_2} \right\|_{Q^{-1}}^2 \middle| \mathcal{F}_{r,k} \right] \leq \left( 1 + \frac{1}{K} \right) \left\| z_{r,k}^{m_1} - z_{r,k}^{m_2} \right\|_{Q^{-1}}^2 + 2 \left( 1 + \frac{1}{K} \right) \eta_c^2 \sigma^2 \|Q\|_2^{-1} + 4(1+K)\eta_c^2 \zeta^2 \|Q\|_2^{-1}.$$

The proof of Proposition E.3 relies on the following three claims. To simplify the exposition we introduce two more notations for this subsection. For any  $r, k, m$ , let

$$\varepsilon_{r,k}^m := \nabla f(w_{r,k}^m; \zeta_{r,k}^m) - \nabla F_m(w_{r,k}^m), \quad \delta_{r,k}^m := \nabla F_m(w_{r,k}^m) - \nabla F(w_{r,k}^m).$$

The following claim upper bounds the growth of  $\|z_{r,k+1}^{m_1} - z_{r,k+1}^{m_2}\|_{Q^{-1}}^2$ . The proof of Claim E.4 is deferred to Appendix E.2.1.

**Claim E.4.** *In the same settings of Proposition E.3, the following inequality holds*

$$\begin{aligned} \left\| z_{r,k+1}^{m_1} - z_{r,k+1}^{m_2} \right\|_{Q^{-1}}^2 &\leq \left( 1 + \frac{1}{K} \right) \left\| z_{r,k}^{m_1} - z_{r,k}^{m_2} - \eta_c \cdot Q \left( w_{r,k}^{m_1} - w_{r,k}^{m_2} \right) - \eta_c \left( \varepsilon_{r,k}^{m_1} - \varepsilon_{r,k}^{m_2} \right) \right\|_{Q^{-1}}^2 \\ &\quad + (1 + K) \eta_c^2 \left\| \delta_{r,k}^{m_1} - \delta_{r,k}^{m_2} \right\|_{Q^{-1}}^2. \end{aligned} \quad (\text{E.1})$$

The next claim upper bounds the growth of the first term in Eq. (E.1) in conditional expectation. We extend the stability technique in (Flammarion & Bach, 2017) to bound this term. The proof of Claim E.5 is deferred to Appendix E.2.2.

**Claim E.5.** *In the same settings of Proposition E.3, the following inequality holds*

$$\mathbb{E} \left[ \left\| z_{r,k}^{m_1} - z_{r,k}^{m_2} - \eta_c Q \left( w_{r,k}^{m_1} - w_{r,k}^{m_2} \right) - \eta_c \left( \varepsilon_{r,k}^{m_1} - \varepsilon_{r,k}^{m_2} \right) \right\|_{Q^{-1}}^2 \middle| \mathcal{F}_{r,k} \right] \leq \left\| z_{r,k}^{m_1} - z_{r,k}^{m_2} \right\|_{Q^{-1}}^2 + 2\eta_c^2 \sigma^2 \|Q\|_2^{-1}.$$

The third claim upper bounds the growth of the second term in Eq. (E.1) under conditional expectation. This is a result of the bounded heterogeneity assumption (Assumption 2(c)). The proof of Claim E.6 is deferred to Appendix E.2.3.

**Claim E.6.** *In the same settings of Proposition E.3, the following inequality holds*

$$\mathbb{E} \left[ \left\| \delta_{r,k}^{m_1} - \delta_{r,k}^{m_2} \right\|_{Q^{-1}}^2 \middle| \mathcal{F}_{r,k} \right] \leq 4\|Q\|_2^{-1} \zeta^2.$$

The proof of the above claims as well as the main lemma require the following helper claim which we also state here. The proof is also deferred to Appendix E.2.3.

**Claim E.7.** *In the same settings of Proposition E.3, the dual norm  $\|\cdot\|_*$  corresponds to the  $\|Q\|_2 \cdot Q^{-1}$ -norm, namely  $\|z\|_* = \sqrt{\|Q\|_2 \cdot z^\top Q^{-1} z}$ .*

The proof of Proposition E.3 is immediate once we have Claims E.4, E.5 and E.6.

*Proof of Proposition E.3.* By Claims E.4, E.5 and E.6,

$$\begin{aligned} &\mathbb{E} \left[ \left\| z_{r,k+1}^{m_1} - z_{r,k+1}^{m_2} \right\|_{Q^{-1}}^2 \middle| \mathcal{F}_{r,k} \right] \\ &\leq \left( 1 + \frac{1}{K} \right) \mathbb{E} \left[ \left\| z_{r,k}^{m_1} - z_{r,k}^{m_2} - \eta_c Q \left( w_{r,k}^{m_1} - w_{r,k}^{m_2} \right) - \eta_c \left( \varepsilon_{r,k}^{m_1} - \varepsilon_{r,k}^{m_2} \right) \right\|_{Q^{-1}}^2 \middle| \mathcal{F}_{r,k} \right] \\ &\quad + (1 + K) \eta_c^2 \mathbb{E} \left[ \left\| \delta_{r,k}^{m_1} - \delta_{r,k}^{m_2} \right\|_{Q^{-1}}^2 \middle| \mathcal{F}_{r,k} \right] \quad (\text{by Claim E.4}) \\ &\leq \left( 1 + \frac{1}{K} \right) \left\| z_{r,k}^{m_1} - z_{r,k}^{m_2} \right\|_{Q^{-1}}^2 + 2 \left( 1 + \frac{1}{K} \right) \eta^2 \sigma^2 \|Q\|_2^{-1} + 4(1 + K) \eta_c^2 \zeta^2 \|Q\|_2^{-1}, \quad (\text{by Claims E.5 and E.6}) \end{aligned}$$

completing the proof of Proposition E.3.  $\square$

The main Lemma E.2 then follows by telescoping Proposition E.3.

*Proof of Lemma E.2.* Let  $m_1, m_2$  be two arbitrary clients. Telescoping Proposition E.3 from  $\mathcal{F}_{r,0}$  to  $\mathcal{F}_{r,k}$  gives

$$\begin{aligned}
 & \mathbb{E} \left[ \left\| z_{r,k}^{m_1} - z_{r,k}^{m_2} \right\|_{Q^{-1}}^2 \right] \\
 & \leq \frac{\left(1 + \frac{1}{K}\right)^k - 1}{\frac{1}{K}} \left( 2 \left(1 + \frac{1}{K}\right) \eta_c^2 \sigma^2 \|Q\|_2^{-1} + 4(1+K)\eta_c^2 \zeta^2 \|Q\|_2^{-1} \right) \quad (\text{telescoping of Proposition E.3}) \\
 & \leq (e-1)K \left( 2 \left(1 + \frac{1}{K}\right) \eta_c^2 \sigma^2 \|Q\|_2^{-1} + 4(1+K)\eta_c^2 \zeta^2 \|Q\|_2^{-1} \right) \quad (\text{since } \left(1 + \frac{1}{K}\right)^k \leq \left(1 + \frac{1}{K}\right)^K < e) \\
 & \leq (e-1)K \left( 4\eta_c^2 \sigma^2 \|Q\|_2^{-1} + 8K\eta_c^2 \zeta^2 \|Q\|_2^{-1} \right) \quad (\text{since } 1 + \frac{1}{K} \leq 2 \text{ and } 1 + K \leq 2K) \\
 & \leq 7\eta_c^2 K \sigma^2 \|Q\|_2^{-1} + 14\eta_c^2 K^2 \zeta^2 \|Q\|_2^{-1} \quad (\text{since } 4(e-1) < 7 \text{ and } 8(e-1) < 14)
 \end{aligned}$$

By convexity of  $\|\cdot\|_{Q^{-1}}^2$  and Proposition E.3 one has

$$\frac{1}{M} \sum_{m=1}^M \mathbb{E} \left[ \left\| \bar{z}_{r,k} - z_{r,k}^m \right\|_{Q^{-1}}^2 \right] \leq \mathbb{E} \left[ \left\| z_{r,k}^{m_1} - z_{r,k}^{m_2} \right\|_{Q^{-1}}^2 \right] \leq 7\eta_c^2 K \sigma^2 \|Q\|_2^{-1} + 14\eta_c^2 K^2 \zeta^2 \|Q\|_2^{-1}.$$

Finally, we switch back to  $\|\cdot\|_*$  norm following Claim E.7

$$\frac{1}{M} \sum_{m=1}^M \mathbb{E} \left[ \left\| \bar{z}_{r,k} - z_{r,k}^m \right\|_*^2 \right] \leq 7\eta_c^2 K \sigma^2 + 14\eta_c^2 K^2 \zeta^2,$$

completing the proof of Lemma E.2.  $\square$

### E.2.1. DEFERRED PROOF OF CLAIM E.4

*Proof of Claim E.4.* By definition of FEDDUALAVG procedure one has

$$\begin{aligned}
 z_{r,k+1}^m &= z_{r,k}^m - \eta_c \nabla f(w_{r,k}^m; \xi_{r,k}^m) \\
 &= z_{r,k}^m - \eta_c \nabla F(w_{r,k}^m) + \eta_c \left( \nabla F_m(w_{r,k}^m) - \nabla F(w_{r,k}^m) \right) + \eta_c \left( \nabla f(w_{r,k}^m; \xi_{r,k}^m) - \nabla F_m(w_{r,k}^m) \right) \\
 &= z_{r,k}^m - \eta_c \nabla F(w_{r,k}^m) - \eta_c \varepsilon_{r,k}^m - \eta_c \delta_{r,k}^m,
 \end{aligned} \tag{E.2}$$

where the last equality is by definition of  $\varepsilon_{r,k}^m$  and  $\delta_{r,k}^m$ . Therefore

$$\begin{aligned}
 & \left\| z_{r,k+1}^{m_1} - z_{r,k+1}^{m_2} \right\|_{Q^{-1}}^2 \\
 &= \left\| z_{r,k}^{m_1} - z_{r,k}^{m_2} - \eta_c Q \left( w_{r,k}^{m_1} - w_{r,k}^{m_2} \right) - \eta_c \left( \varepsilon_{r,k}^{m_1} - \varepsilon_{r,k}^{m_2} \right) - \eta_c \left( \delta_{r,k}^{m_1} - \delta_{r,k}^{m_2} \right) \right\|_{Q^{-1}}^2 \quad (\text{by Eq. (E.2)}) \\
 &= \left\| z_{r,k}^{m_1} - z_{r,k}^{m_2} - \eta_c Q \left( w_{r,k}^{m_1} - w_{r,k}^{m_2} \right) - \eta_c \left( \varepsilon_{r,k}^{m_1} - \varepsilon_{r,k}^{m_2} \right) \right\|_{Q^{-1}}^2 + \eta_c^2 \left\| \delta_{r,k}^{m_1} - \delta_{r,k}^{m_2} \right\|_{Q^{-1}}^2 \\
 & \quad + 2 \left\langle z_{r,k}^{m_1} - z_{r,k}^{m_2} - \eta_c Q \left( w_{r,k}^{m_1} - w_{r,k}^{m_2} \right) - \eta_c \left( \varepsilon_{r,k}^{m_1} - \varepsilon_{r,k}^{m_2} \right), \eta_c Q^{-1} \left( \delta_{r,k}^{m_1} - \delta_{r,k}^{m_2} \right) \right\rangle. \tag{E.3}
 \end{aligned}$$

By Cauchy-Schwartz inequality and AM-GM inequality one has (for any  $\gamma > 0$ )

$$\begin{aligned}
 & \left\langle z_{r,k}^{m_1} - z_{r,k}^{m_2} - \eta_c Q \left( w_{r,k}^{m_1} - w_{r,k}^{m_2} \right) - \eta_c \left( \varepsilon_{r,k}^{m_1} - \varepsilon_{r,k}^{m_2} \right), \eta_c Q^{-1} \left( \delta_{r,k}^{m_1} - \delta_{r,k}^{m_2} \right) \right\rangle \\
 & \leq \left\| z_{r,k}^{m_1} - z_{r,k}^{m_2} - \eta_c Q \left( w_{r,k}^{m_1} - w_{r,k}^{m_2} \right) - \eta_c \left( \varepsilon_{r,k}^{m_1} - \varepsilon_{r,k}^{m_2} \right) \right\|_{Q^{-1}} \left\| \eta_c \left( \delta_{r,k}^{m_1} - \delta_{r,k}^{m_2} \right) \right\|_{Q^{-1}} \quad (\text{Cauchy-Schwarz inequality}) \\
 & \leq \frac{1}{2\gamma} \left\| z_{r,k}^{m_1} - z_{r,k}^{m_2} - \eta_c Q \left( w_{r,k}^{m_1} - w_{r,k}^{m_2} \right) - \eta_c \left( \varepsilon_{r,k}^{m_1} - \varepsilon_{r,k}^{m_2} \right) \right\|_{Q^{-1}}^2 + \frac{1}{2} \gamma \left\| \eta_c \left( \delta_{r,k}^{m_1} - \delta_{r,k}^{m_2} \right) \right\|_{Q^{-1}}^2. \quad (\text{AM-GM inequality}) \\
 & \leq \frac{1}{2\gamma} \left\| z_{r,k}^{m_1} - z_{r,k}^{m_2} - \eta_c Q \left( w_{r,k}^{m_1} - w_{r,k}^{m_2} \right) - \eta_c \left( \varepsilon_{r,k}^{m_1} - \varepsilon_{r,k}^{m_2} \right) \right\|_{Q^{-1}}^2 + \frac{1}{2} \gamma \eta_c^2 \left\| \left( \delta_{r,k}^{m_1} - \delta_{r,k}^{m_2} \right) \right\|_{Q^{-1}}^2. \tag{E.4}
 \end{aligned}$$

Plugging Eq. (E.4) to Eq. (E.3) with  $\gamma = K$  gives

$$\begin{aligned} & \left\| z_{r,k+1}^{m_1} - z_{r,k+1}^{m_2} \right\|_{Q^{-1}}^2 \\ & \leq \left(1 + \frac{1}{K}\right) \left\| z_{r,k}^{m_1} - z_{r,k}^{m_2} - \eta_c Q \left( w_{r,k}^{m_1} - w_{r,k}^{m_2} \right) - \eta_c \left( \varepsilon_{r,k}^{m_1} - \varepsilon_{r,k}^{m_2} \right) \right\|_{Q^{-1}}^2 + (1+K) \eta_c^2 \left\| \delta_{r,k}^{m_1} - \delta_{r,k}^{m_2} \right\|_{Q^{-1}}^2, \end{aligned}$$

completing the proof of Claim E.4.  $\square$

### E.2.2. DEFERRED PROOF OF CLAIM E.5

The proof technique of this claim is similar to (Flammarion & Bach, 2017, Lemma 8) which we adapt to fit into our settings.

*Proof of Claim E.5.* Let us first expand the  $\|\cdot\|_{Q^{-1}}^2$ :

$$\begin{aligned} & \left\| z_{r,k}^{m_1} - z_{r,k}^{m_2} - \eta_c Q \left( w_{r,k}^{m_1} - w_{r,k}^{m_2} \right) - \eta_c \left( \varepsilon_{r,k}^{m_1} - \varepsilon_{r,k}^{m_2} \right) \right\|_{Q^{-1}}^2 \\ & = \left\| z_{r,k}^{m_1} - z_{r,k}^{m_2} \right\|_{Q^{-1}}^2 + \left\| \eta_c Q \left( w_{r,k}^{m_1} - w_{r,k}^{m_2} \right) \right\|_{Q^{-1}}^2 + \left\| \eta_c \left( \varepsilon_{r,k}^{m_1} - \varepsilon_{r,k}^{m_2} \right) \right\|_{Q^{-1}}^2 + 2 \left\langle \eta_c \left( w_{r,k}^{m_1} - w_{r,k}^{m_2} \right), \eta_c \left( \varepsilon_{r,k}^{m_1} - \varepsilon_{r,k}^{m_2} \right) \right\rangle \\ & \quad + 2 \left\langle z_{r,k}^{m_1} - z_{r,k}^{m_2}, -\eta_c \left( w_{r,k}^{m_1} - w_{r,k}^{m_2} \right) \right\rangle + 2 \left\langle z_{r,k}^{m_1} - z_{r,k}^{m_2}, -\eta_c Q^{-1} \left( \varepsilon_{r,k}^{m_1} - \varepsilon_{r,k}^{m_2} \right) \right\rangle. \end{aligned}$$

Now we take conditional expectation. Note that by bounded variance assumption one has

$$\mathbb{E} \left[ \left\| \eta_c \left( \varepsilon_{r,k}^{m_1} - \varepsilon_{r,k}^{m_2} \right) \right\|_{Q^{-1}}^2 \middle| \mathcal{F}_{r,k} \right] = \|Q\|_2^{-1} \cdot \mathbb{E} \left[ \left\| \eta_c \left( \varepsilon_{r,k}^{m_1} - \varepsilon_{r,k}^{m_2} \right) \right\|_*^2 \middle| \mathcal{F}_{r,k} \right] \leq 2\eta_c^2 \sigma^2 \|Q\|_2^{-1},$$

where in the first equality we applied Claim E.7.

By unbiased and independence assumptions

$$\mathbb{E} \left[ \varepsilon_{r,k}^{m_1} - \varepsilon_{r,k}^{m_2} \middle| \mathcal{F}_{r,k} \right] = 0.$$

Thus

$$\begin{aligned} & \mathbb{E} \left[ \left\| z_{r,k}^{m_1} - z_{r,k}^{m_2} - \eta_c Q \left( w_{r,k}^{m_1} - w_{r,k}^{m_2} \right) - \eta_c \left( \varepsilon_{r,k}^{m_1} - \varepsilon_{r,k}^{m_2} \right) \right\|_{Q^{-1}}^2 \middle| \mathcal{F}_{r,k} \right] \\ & \leq \underbrace{\left\| z_{r,k}^{m_1} - z_{r,k}^{m_2} \right\|_{Q^{-1}}^2}_{\text{(I)}} + \underbrace{\eta_c^2 \left\| Q \left( w_{r,k}^{m_1} - w_{r,k}^{m_2} \right) \right\|_{Q^{-1}}^2}_{\text{(II)}} - 2\eta_c \left\langle z_{r,k}^{m_1} - z_{r,k}^{m_2}, w_{r,k}^{m_1} - w_{r,k}^{m_2} \right\rangle + 2\eta_c^2 \sigma^2 \|Q\|_2^{-1}. \quad (\text{E.5}) \end{aligned}$$

Now we analyze (I), (II) in Eq. (E.5). First note that

$$\begin{aligned} \text{(I)} & = \eta_c^2 \left\| Q \left( w_{r,k}^{m_1} - w_{r,k}^{m_2} \right) \right\|_{Q^{-1}}^2 \\ & = \eta_c^2 \left\langle w_{r,k}^{m_1} - w_{r,k}^{m_2}, Q \left( w_{r,k}^{m_1} - w_{r,k}^{m_2} \right) \right\rangle \quad (\text{by definition of } \|\cdot\|_{Q^{-1}}) \\ & = \eta_c \left\langle w_{r,k}^{m_1} - w_{r,k}^{m_2}, \eta_c \left( \nabla F(w_{r,k}^{m_1}) - \nabla F(w_{r,k}^{m_2}) \right) \right\rangle \quad (\text{since } F \text{ is quadratic}) \\ & = \eta_c \left\langle w_{r,k}^{m_1} - w_{r,k}^{m_2}, \nabla(\eta_c F - 2h)(w_{r,k}^{m_1}) - \nabla(\eta_c F - 2h)(w_{r,k}^{m_2}) \right\rangle + 2\eta_c \left\langle w_{r,k}^{m_1} - w_{r,k}^{m_2}, \nabla h(w_{r,k}^{m_1}) - \nabla h(w_{r,k}^{m_2}) \right\rangle \end{aligned}$$

By  $L$ -smoothness of  $F_m$  (Assumption 1(c)) we know that  $F := \frac{1}{M} \sum_{m=1}^M F_m$  is also  $L$ -smooth. Thus  $\eta_c F$  is  $\frac{1}{4}$ -smooth since  $\eta_c \leq \frac{1}{4L}$ . Thus  $\eta_c F - 2h$  is concave since  $h$  is 1-strongly convex, which implies

$$\left\langle w_{r,k}^{m_1} - w_{r,k}^{m_2}, \nabla(\eta_c F - 2h)(w_{r,k}^{m_1}) - \nabla(\eta_c F - 2h)(w_{r,k}^{m_2}) \right\rangle \leq 0.$$

We obtain

$$(I) \leq 2\eta_c \left\langle w_{r,k}^{m_1} - w_{r,k}^{m_2}, \nabla h(w_{r,k}^{m_1}) - \nabla h(w_{r,k}^{m_2}) \right\rangle. \quad (E.6)$$

Now we study (I)+(II) in Eq. (E.5):

$$\begin{aligned} (I) + (II) &= \eta_c^2 \left\| Q \left( w_{r,k}^{m_1} - w_{r,k}^{m_2} \right) \right\|_{Q^{-1}}^2 - 2\eta_c \left\langle z_{r,k}^{m_1} - z_{r,k}^{m_2}, w_{r,k}^{m_1} - w_{r,k}^{m_2} \right\rangle \\ &\leq 2\eta_c \left\langle w_{r,k}^{m_1} - w_{r,k}^{m_2}, \nabla h(w_{r,k}^{m_1}) - \nabla h(w_{r,k}^{m_2}) \right\rangle - 2\eta_c \left\langle w_{r,k}^{m_1} - w_{r,k}^{m_2}, z_{r,k}^{m_1} - z_{r,k}^{m_2} \right\rangle \quad (\text{by inequality Eq. (E.6)}) \\ &= -2\eta_c \left\langle w_{r,k}^{m_1} - w_{r,k}^{m_2}, \left( z_{r,k}^{m_1} - \nabla h(w_{r,k}^{m_1}) \right) - \left( z_{r,k}^{m_2} - \nabla h(w_{r,k}^{m_2}) \right) \right\rangle \end{aligned} \quad (E.7)$$

On the other hand, by definition of  $w_{r,k}^m$  we have

$$w_{r,k}^m = \nabla (h + (rK + k)\eta_c\psi)^*(z_{r,k}^m) = \arg \min_w \left\{ \left\langle -z_{r,k}^m, w \right\rangle + (rK + k)\eta_c\psi(w) + h(w) \right\}.$$

By subdifferential calculus one has

$$z_{r,k}^m - \nabla h(w_{r,k}^m) \in \partial [\eta_c(rK + k)\psi(w_{r,k}^m)].$$

By monotonicity of subgradients one has

$$\left\langle w_{r,k}^{m_1} - w_{r,k}^{m_2}, \left( z_{r,k}^{m_1} - \nabla h(w_{r,k}^{m_1}) \right) - \left( z_{r,k}^{m_2} - \nabla h(w_{r,k}^{m_2}) \right) \right\rangle \geq 0. \quad (E.8)$$

Combining Eqs. (E.7) and (E.8) gives

$$(I) + (II) \leq 0. \quad (E.9)$$

Combining Eqs. (E.5) and (E.9) completes the proof as

$$\mathbb{E} \left[ \left\| z_{r,k}^{m_1} - z_{r,k}^{m_2} - \eta_c Q \left( w_{r,k}^{m_1} - w_{r,k}^{m_2} \right) - \eta_c \left( \varepsilon_{r,k}^{m_1} - \varepsilon_{r,k}^{m_2} \right) \right\|_{Q^{-1}}^2 \middle| \mathcal{F}_{r,k} \right] \leq \left\| z_{r,k}^{m_1} - z_{r,k}^{m_2} \right\|_{Q^{-1}}^2 + 2\eta_c^2 \sigma^2 \|Q\|_2^{-1}.$$

□

### E.2.3. DEFERRED PROOF OF CLAIMS E.6 AND E.7

*Proof of Claim E.6.* By triangle inequality and AM-GM inequality,

$$\begin{aligned} &\mathbb{E} \left[ \left\| \delta_{r,k}^{m_1} - \delta_{r,k}^{m_2} \right\|_{Q^{-1}}^2 \middle| \mathcal{F}_{r,k} \right] \\ &\leq \mathbb{E} \left[ \left( \left\| \delta_{r,k}^{m_1} \right\|_{Q^{-1}} + \left\| \delta_{r,k}^{m_2} \right\|_{Q^{-1}} \right)^2 \middle| \mathcal{F}_{r,k} \right] \quad (\text{triangle inequality}) \\ &\leq 2 \mathbb{E} \left[ \left\| \delta_{r,k}^{m_1} \right\|_{Q^{-1}}^2 + \left\| \delta_{r,k}^{m_2} \right\|_{Q^{-1}}^2 \middle| \mathcal{F}_{r,k} \right]. \quad (\text{AM-GM inequality}) \end{aligned}$$

By Claim E.7,

$$\mathbb{E} \left[ \left\| \delta_{r,k}^{m_1} - \delta_{r,k}^{m_2} \right\|_{Q^{-1}}^2 \middle| \mathcal{F}_{r,k} \right] \leq 2 \|Q\|_2^{-1} \mathbb{E} \left[ \left\| \delta_{r,k}^{m_1} \right\|_*^2 + \left\| \delta_{r,k}^{m_2} \right\|_*^2 \middle| \mathcal{F}_{r,k} \right] \leq 4 \|Q\|_2^{-1} \zeta^2,$$

where the last inequality is due to bounded heterogeneity Assumption 2(c). This completes the proof of Claim E.6. □

*Proof of Claim E.7.* Since the primal norm  $\| \cdot \|$  is  $(\|Q\|_2^{-1} \cdot Q)$ -norm by Assumption 2(b), the dual norm  $\| \cdot \|_*$  is  $(\|Q\|_2^{-1} \cdot Q)^{-1} = \|Q\|_2 \cdot Q^{-1}$ -norm. □

## F. Proof of Theorem 4.1

In this section, we state and prove Theorem 4.1 on the convergence of FEDDUALAVG for small client learning rate  $\eta_c$ . The intuition is that for sufficiently small client learning rate, FEDDUALAVG is almost as good as stochastic mini-batch with  $R$  iterations and batch-size  $MK$ . The proof technique is very similar to the above sections and (Karimireddy et al., 2020) so we skip a substantial amount of the proof details. We present the proof for FEDDUALAVG only since the analysis of FEDMiD is very similar.

To facilitate the analysis we re-parameterize the hyperparameters by letting  $\eta := \eta_s \eta_c$ , and we treat  $(\eta, \eta_c)$  as independent hyperparameters (rather than  $(\eta_c, \eta_s)$ ). We use the notation  $h_{r,k} := h + \tilde{\eta}_{r,k} \cdot \psi = h + (\eta r K + \eta_c k) \psi$ ,  $\bar{z}_{r,k} := \frac{1}{M} \sum_{m=1}^M z_{r,k}^m$ , and  $\widehat{w}_{r,k} := \nabla h_{r,k}^*(\bar{z}_{r,k})$ . Note that  $\widehat{w}_{r,0} = w_{r,0}^m$  for all  $m \in [M]$  by definition.

### F.1. Main Theorem and Lemmas

Now we state the full version of Theorem 4.1 on FEDDUALAVG with small client learning rate  $\eta_c$ .

**Theorem F.1** (Detailed version of Theorem 4.1). *Assuming Assumption 1, then for any  $\eta \in (0, \frac{1}{4KL}]$ , for any initialization  $w_0 \in \text{dom } \psi$ , there exists an  $\eta_c^{\max} > 0$  (which may depend on  $\eta$  and  $w_0$ ) such that for any  $\eta_c \in (0, \eta_c^{\max}]$ , FEDDUALAVG yields*

$$\mathbb{E} \left[ \Phi \left( \frac{1}{R} \sum_{r=1}^R \widehat{w}_{r,0} \right) - \Phi(w^*) \right] \leq \frac{B}{\eta KR} + \frac{3\eta\sigma^2}{M},$$

where  $B := D_h(w^*, w_0)$  is the Bregman divergence between the optimal  $w^*$  and the initialization  $w_0$ .

In particular for

$$\eta = \min \left\{ \frac{1}{4KL}, \frac{B^{\frac{1}{2}} M^{\frac{1}{2}}}{K^{\frac{1}{2}} R^{\frac{1}{2}} \sigma} \right\},$$

one has

$$\mathbb{E} \left[ \Phi \left( \frac{1}{R} \sum_{r=1}^R \widehat{w}_{r,0} \right) - \Phi(w^*) \right] \leq \frac{4LB}{R} + \frac{4\sigma B^{\frac{1}{2}}}{M^{\frac{1}{2}} K^{\frac{1}{2}} R^{\frac{1}{2}}}.$$

The proof of Theorem F.1 relies on the following lemmas.

The first Lemma F.2 analyzes  $\tilde{D}_{h_{r+1,0}}(w, \bar{z}_{r+1,0})$ . The proof of Lemma F.2 is deferred to Appendix F.2.

**Lemma F.2.** *Under the same settings of Theorem F.1, the following inequality holds.*

$$\begin{aligned} & \tilde{D}_{h_{r+1,0}}(w, \bar{z}_{r+1,0}) - \tilde{D}_{h_{r,0}}(w, \bar{z}_{r,0}) \\ & \leq -\tilde{D}_{h_{r,0}}(\widehat{w}_{r+1,0}, \bar{z}_{r,0}) - \eta K (\Phi(\widehat{w}_{r+1,0}) - \Phi(w)) + \frac{L}{2} \eta K \|\widehat{w}_{r+1,0} - \widehat{w}_{r,0}\|^2 \\ & \quad + \eta K \left\langle \nabla F(\widehat{w}_{r,0}) - \frac{1}{MK} \sum_{m=1}^M \sum_{k=0}^{K-1} \nabla f(w_{r,k}^m; \xi_{r,k}^m), \widehat{w}_{r+1,0} - w \right\rangle \end{aligned}$$

The second lemma analyzes  $\tilde{D}_{h_{r+1,0}}(w, \bar{z}_{r+1,0})$  under conditional expectation. The proof of Lemma F.3 is deferred to Appendix F.3.

**Lemma F.3.** *Under the same settings of Theorem F.1, there exists an  $\eta_c^{\max} > 0$  (which may depend on  $\eta$  and  $w_0$ ) such that for any  $\eta_c \in (0, \eta_c^{\max}]$ , FEDDUALAVG yields*

$$\mathbb{E} \left[ \tilde{D}_{h_{r+1,0}}(w, \bar{z}_{r+1,0}) \middle| \mathcal{F}_{r,0} \right] - \tilde{D}_{h_{r,0}}(w, \bar{z}_{r,0}) \leq -\eta K \mathbb{E} [(\Phi(\widehat{w}_{r+1,0}) - \Phi(w)) \middle| \mathcal{F}_{r,0}] + \frac{3\eta^2 K \sigma^2}{M}.$$

With Lemmas F.2 and F.3 at hand we are ready to prove Theorem F.1.

*Proof of Theorem F.1.* Telescoping Lemma F.3 and dropping the negative terms gives

$$\frac{1}{R} \sum_{r=1}^R \mathbb{E} [\Phi(\widehat{w}_{r,0}) - \Phi(w)] \leq \frac{1}{\eta KR} \tilde{D}_{h_{r,0}}(w, \bar{z}_{r,0}) + \frac{3\eta\sigma^2}{M} = \frac{B}{\eta KR} + \frac{3\eta\sigma^2}{M}.$$

The second inequality of Theorem F.1 follows immediately once we plug in the specified  $\eta$ .  $\square$

## F.2. Deferred Proof of Lemma F.2

*Proof of Lemma F.2.* The proof of this lemma is very similar to Claims D.5 and D.6 so we skip most of the details.

We start by analyzing  $\tilde{D}_{h_{r+1,0}}(w, \overline{z_{r+1,0}})$ .

$$\begin{aligned}
 & \tilde{D}_{h_{r+1,0}}(w, \overline{z_{r+1,0}}) \\
 &= h_{r+1,0}(w) - h_{r+1,0}(\nabla h_{r+1,0}^*(\overline{z_{r+1,0}})) - \langle \overline{z_{r+1,0}}, w - \nabla h_{r+1,0}^*(\overline{z_{r+1,0}}) \rangle \\
 & \hspace{15em} \text{(By definition of generalized Bregman divergence } \tilde{D}) \\
 &= h_{r+1,0}(w) - h_{r+1,0}(\widehat{w_{r+1,0}}) - \langle \overline{z_{r+1,0}}, w - \widehat{w_{r+1,0}} \rangle \\
 & \hspace{15em} \text{(By definition of } \widehat{w_{r+1,0}}) \\
 &= h_{r+1,0}(w) - h_{r+1,0}(\widehat{w_{r+1,0}}) - \left\langle \overline{z_{r,0}} - \eta K \cdot \frac{1}{MK} \sum_{m=1}^M \sum_{k=0}^{K-1} \nabla f(w_{r,k}^m; \xi_{r,k}^m), w - \widehat{w_{r+1,0}} \right\rangle \\
 & \hspace{15em} \text{(By FEDDUALAVG procedure)} \\
 &= (h_{r,0}(w) + \eta K \psi(w)) - (h_{r,0}(\widehat{w_{r+1,0}}) + \eta K \psi(\widehat{w_{r+1,0}})) - \left\langle \overline{z_{r,0}} - \eta K \cdot \frac{1}{MK} \sum_{m=1}^M \sum_{k=0}^{K-1} \nabla f(w_{r,k}^m; \xi_{r,k}^m), w - \widehat{w_{r+1,0}} \right\rangle \\
 & \hspace{15em} \text{(By definition of } h_{r+1,0}) \\
 &= (h_{r,0}(w) - h_{r,0}(\widehat{w_{r,0}}) - \langle \overline{z_{r,0}}, w - \widehat{w_{r,0}} \rangle) - (h_{r,0}(\widehat{w_{r+1,0}}) - h_{r,0}(\widehat{w_{r,0}}) - \langle \overline{z_{r,0}}, \widehat{w_{r+1,0}} - \widehat{w_{r,0}} \rangle) \\
 & \quad - \eta K (\psi(\widehat{w_{r+1,0}}) - \psi(w)) - \eta K \left\langle \frac{1}{MK} \sum_{m=1}^M \sum_{k=0}^{K-1} \nabla f(w_{r,k}^m; \xi_{r,k}^m), \widehat{w_{r+1,0}} - w \right\rangle \\
 & \hspace{15em} \text{(Rearranging)} \\
 &= \tilde{D}_{h_{r,0}}(w, \overline{z_{r,0}}) - \tilde{D}_{h_{r,0}}(\widehat{w_{r+1,0}}, \overline{z_{r,0}}) - \eta K (\psi(\widehat{w_{r+1,0}}) - \psi(w)) - \eta K \left\langle \frac{1}{MK} \sum_{m=1}^M \sum_{k=0}^{K-1} \nabla f(w_{r,k}^m; \xi_{r,k}^m), \widehat{w_{r+1,0}} - w \right\rangle \\
 & \hspace{15em} \text{(By definition of } \tilde{D})
 \end{aligned}$$

By smoothness and convexity of  $F$  we have

$$\begin{aligned}
 F(\widehat{w_{r+1,0}}) &\leq F(\widehat{w_{r,0}}) + \langle \nabla F(\widehat{w_{r,0}}), \widehat{w_{r+1,0}} - \widehat{w_{r,0}} \rangle + \frac{L}{2} \|\widehat{w_{r+1,0}} - \widehat{w_{r,0}}\|^2 && \text{(by } L\text{-smoothness of } F) \\
 &\leq F(w) + \langle \nabla F(\widehat{w_{r,0}}), \widehat{w_{r+1,0}} - w \rangle + \frac{L}{2} \|\widehat{w_{r+1,0}} - \widehat{w_{r,0}}\|^2 && \text{(by convexity of } F)
 \end{aligned}$$

Combining the above two (in)equalities gives

$$\begin{aligned}
 \tilde{D}_{h_{r+1,0}}(w, \overline{z_{r+1,0}}) - \tilde{D}_{h_{r,0}}(w, \overline{z_{r,0}}) &\leq -\tilde{D}_{h_{r,0}}(\widehat{w_{r+1,0}}, \overline{z_{r,0}}) - \eta K (\Phi(\widehat{w_{r+1,0}}) - \Phi(w)) + \frac{L}{2} \eta K \|\widehat{w_{r+1,0}} - \widehat{w_{r,0}}\|^2 \\
 &\quad + \eta K \left\langle \nabla F(\widehat{w_{r,0}}) - \frac{1}{MK} \sum_{m=1}^M \sum_{k=0}^{K-1} \nabla f(w_{r,k}^m; \xi_{r,k}^m), \widehat{w_{r+1,0}} - w \right\rangle.
 \end{aligned}$$

$\square$



**F.3. Deferred Proof of Lemma F.3**

*Proof of Lemma F.3.* We start by splitting the inner product term in the inequality of Lemma F.2:

$$\begin{aligned}
 & \left\langle \nabla F(\widehat{w}_{r,0}) - \frac{1}{MK} \sum_{m=1}^M \sum_{k=0}^{K-1} \nabla f(w_{r,k}^m; \xi_{r,k}^m), \widehat{w}_{r+1,0} - w \right\rangle \\
 = & \underbrace{\left\langle \nabla F(\widehat{w}_{r,0}) - \frac{1}{MK} \sum_{m=1}^M \sum_{k=0}^{K-1} \nabla f(\widehat{w}_{r,0}; \xi_{r,k}^m), \widehat{w}_{r,0} - w \right\rangle}_{\text{(I)}} \\
 & + \underbrace{\left\langle \nabla F(\widehat{w}_{r,0}) - \frac{1}{MK} \sum_{m=1}^M \sum_{k=0}^{K-1} \nabla f(\widehat{w}_{r,0}; \xi_{r,k}^m), \widehat{w}_{r+1,0} - \widehat{w}_{r,0} \right\rangle}_{\text{(II)}} \\
 & + \underbrace{\frac{1}{MK} \sum_{m=1}^M \sum_{k=0}^{K-1} \langle \nabla f(\widehat{w}_{r,0}; \xi_{r,k}^m) - \nabla f(w_{r,k}^m; \xi_{r,k}^m), \widehat{w}_{r+1,0} - w \rangle}_{\text{(III)}}.
 \end{aligned}$$

Now we investigate the terms (I)-(III). By conditional independence we know  $\mathbb{E}[\text{(I)}|\mathcal{F}_{r,0}] = 0$ . For (II), we know that

$$\begin{aligned}
 \mathbb{E}[\text{(II)}|\mathcal{F}_{r,0}] & \leq \mathbb{E} \left[ \left\| \nabla F(\widehat{w}_{r,0}) - \frac{1}{MK} \sum_{m=1}^M \sum_{k=0}^{K-1} \nabla f(\widehat{w}_{r,0}; \xi_{r,k}^m) \right\|_* \middle| \mathcal{F}_{r,0} \right] \mathbb{E} [\|\widehat{w}_{r+1,0} - \widehat{w}_{r,0}\| | \mathcal{F}_{r,0}] \\
 & \leq \frac{\sigma}{\sqrt{MK}} \cdot \mathbb{E} [\|\widehat{w}_{r+1,0} - \widehat{w}_{r,0}\| | \mathcal{F}_{r,0}]
 \end{aligned}$$

For (III) we observe that (by smoothness assumption)

$$\begin{aligned}
 \text{(III)} & \leq \frac{1}{MK} \sum_{m=1}^M \sum_{k=0}^{K-1} \|\nabla f(\widehat{w}_{r,0}; \xi_{r,k}^m) - \nabla f(w_{r,k}^m; \xi_{r,k}^m)\|_* \|\widehat{w}_{r+1,0} - w\| \\
 & \leq \frac{L}{MK} \sum_{m=1}^M \sum_{k=0}^{K-1} \|\widehat{w}_{r,0} - w_{r,k}^m\| \|\widehat{w}_{r+1,0} - w\|.
 \end{aligned}$$

Taking conditional expectation,

$$\begin{aligned}
 \mathbb{E}[\text{(III)}|\mathcal{F}_{r,0}] & \leq \frac{1}{MK} \sum_{m=1}^M \sum_{k=0}^{K-1} \mathbb{E} [\|\nabla f(\widehat{w}_{r,0}; \xi_{r,k}^m) - \nabla f(w_{r,k}^m; \xi_{r,k}^m)\|_* | \mathcal{F}_{r,0}] \mathbb{E} [\|\widehat{w}_{r+1,0} - w\| | \mathcal{F}_{r,0}] \\
 & \leq \frac{L}{MK} \left( \sum_{m=1}^M \sum_{k=0}^{K-1} \mathbb{E} [\|\widehat{w}_{r,0} - w_{r,k}^m\| | \mathcal{F}_{r,0}] \right) \mathbb{E} [\|\widehat{w}_{r+1,0} - w\| | \mathcal{F}_{r,0}]
 \end{aligned}$$

Combining the above inequalities with Lemma F.2 gives

$$\begin{aligned}
 & \mathbb{E} \left[ \tilde{D}_{h_{r+1,0}}(w, \overline{z_{r+1,0}}) \middle| \mathcal{F}_t \right] - \tilde{D}_{h_{r,0}}(w, \overline{z_{r,0}}) \\
 & \leq -\eta K \mathbb{E} [(\Phi(\widehat{w}_{r+1,0}) - \Phi(w)) | \mathcal{F}_{r,0}] - \left( \frac{1}{2} - \frac{L}{2} \eta K \right) \mathbb{E} [\|\widehat{w}_{r+1,0} - \widehat{w}_{r,0}\|^2 | \mathcal{F}_{r,0}] \\
 & \quad + \frac{\eta \sigma \sqrt{K}}{\sqrt{M}} \cdot \mathbb{E} [\|\widehat{w}_{r+1,0} - \widehat{w}_{r,0}\| | \mathcal{F}_{r,0}] + \frac{\eta L}{M} \left( \sum_{m=1}^M \sum_{k=0}^{K-1} \mathbb{E} [\|\widehat{w}_{r,0} - w_{r,k}^m\| | \mathcal{F}_{r,0}] \right) \mathbb{E} [\|\widehat{w}_{r+1,0} - w\| | \mathcal{F}_{r,0}]
 \end{aligned}$$

Note that

$$\begin{aligned}
 & - \left( \frac{1}{2} - \frac{L}{2} \eta K \right) \mathbb{E} [\|\widehat{w}_{r+1,0} - \widehat{w}_{r,0}\|^2 | \mathcal{F}_{r,0}] + \frac{\eta \sigma \sqrt{K}}{\sqrt{M}} \cdot \mathbb{E} [\|\widehat{w}_{r+1,0} - \widehat{w}_{r,0}\| | \mathcal{F}_{r,0}] \\
 & \leq - \frac{3}{8} \mathbb{E} [\|\widehat{w}_{r+1,0} - \widehat{w}_{r,0}\|^2 | \mathcal{F}_{r,0}] + \frac{\eta \sigma \sqrt{K}}{\sqrt{M}} \cdot \mathbb{E} [\|\widehat{w}_{r+1,0} - \widehat{w}_{r,0}\| | \mathcal{F}_{r,0}] \quad (\text{since } \eta \leq \frac{1}{4KL}) \\
 & \leq - \frac{1}{4} \mathbb{E} [\|\widehat{w}_{r+1,0} - \widehat{w}_{r,0}\|^2 | \mathcal{F}_{r,0}] + \frac{2\eta^2 K \sigma^2}{M}. \quad (\text{by quadratic optimum})
 \end{aligned}$$

Therefore

$$\begin{aligned}
 & \mathbb{E} \left[ \tilde{D}_{h_{r+1,0}}(w, \overline{z_{r+1,0}}) \Big| \mathcal{F}_t \right] - \tilde{D}_{h_{r,0}}(w, \overline{z_{r,0}}) \\
 & \leq - \eta K \mathbb{E} [(\Phi(\widehat{w}_{r+1,0}) - \Phi(w)) | \mathcal{F}_{r,0}] - \frac{1}{4} \mathbb{E} [\|\widehat{w}_{r+1,0} - \widehat{w}_{r,0}\|^2 | \mathcal{F}_{r,0}] + \frac{2\eta^2 K \sigma^2}{M} \\
 & \quad + \frac{\eta L}{M} \left( \sum_{m=1}^M \sum_{k=0}^{K-1} \mathbb{E} [\|\widehat{w}_{r,0} - w_{r,k}^m\| | \mathcal{F}_{r,0}] \right) \mathbb{E} [\|\widehat{w}_{r+1,0} - w\| | \mathcal{F}_{r,0}].
 \end{aligned}$$

Since  $w_{r,k}^m$  is generated by running local dual averaging with learning rate  $\eta_c$ , one has

$$\lim_{\eta_c \downarrow 0} \left[ \left( \sum_{m=1}^M \sum_{k=0}^{K-1} \mathbb{E} [\|\widehat{w}_{r,0} - w_{r,k}^m\| | \mathcal{F}_{r,0}] \right) \mathbb{E} [\|\widehat{w}_{r+1,0} - w\| | \mathcal{F}_{r,0}] \right] = 0.$$

There exists an upper bound  $\eta_c^{\max}$  such that for any  $\eta_c \in (0, \eta_c^{\max}]$ , it is the case that

$$\left( \sum_{m=1}^M \sum_{k=0}^{K-1} \mathbb{E} [\|\widehat{w}_{r,0} - w_{r,k}^m\| | \mathcal{F}_{r,0}] \right) \mathbb{E} [\|\widehat{w}_{r+1,0} - w\| | \mathcal{F}_{r,0}] \leq \frac{\eta K \sigma^2}{L}.$$

Therefore, for any  $\eta_c \in (0, \eta_c^{\max}]$ ,

$$\mathbb{E} \left[ \tilde{D}_{h_{r+1,0}}(w, \overline{z_{r+1,0}}) \Big| \mathcal{F}_t \right] - \tilde{D}_{h_{r,0}}(w, \overline{z_{r,0}}) \leq - \eta K \mathbb{E} [\Phi(\widehat{w}_{r+1,0}) - \Phi(w) | \mathcal{F}_{r,0}] + \frac{3\eta^2 K \sigma^2}{M}.$$

□