FedADMM: A federated primal-dual algorithm allowing partial participation

Han Wang, Siddartha Marella, James Anderson

Abstract—Federated learning is a framework for distributed optimization that places emphasis on communication efficiency. In particular, it follows a client-server broadcast model and is appealing because of its ability to accommodate heterogeneity in client compute and storage resources, non-i.i.d. data assumptions, and data privacy. Our contribution is to offer a new federated learning algorithm, FedADMM, for solving non-convex composite optimization problems with non-smooth regularizers. We prove the convergence of FedADMM for the case when not all clients are able to participate in a given communication round under a very general sampling model.

I. Introduction

Federated learning (FL) [14, 19], a novel distributed learning paradigm has attracted significant attention over the past few years. Federated algorithms take a client/server computation model and provide scope to train large-scale machine learning models over an edge-based distributed computing architecture. Models are trained collaboratively under the coordination of a central server while storing data locally on the edge/clients. Typically, clients (devices and entities ranging from mobile phones to hospitals, to an internet of things [21, 11]) are assumed to be heterogeneous; each client is subject to its own constraints on available computational and storage resources. By allowing data to be stored client-side, the FL paradigm has many favorable privacy properties.

In contrast to "traditional" distributed optimization, the FL framework has its own unique challenges and characteristics. First, communication becomes problematic when the number of edge devices/clients is large, or the connection between the central server and a device is slow, e.g., limited bandwidth cellular networks. Second, datasets stored in each client may be highly heterogeneous in that they are sampled from different population distributions, or the amount of data belonging to each client is unbalanced. Third, device/client heterogeneity can severely hinder algorithm performance; differences in hardware, software, and power (connectivity) lead to varying computation speeds among clients, leading to global performance being dominated by the slowest participant. This is known as the "straggler" effect. Additionally, the server may lose control over the clients when they power down or lose connectivity. It is thus common for

This work is supported by awards from the NSF; award ID CAREER-2144634, and DOE; award ID DE-SC0022234. The authors are with the Department of Electrical Engineering, Columbia University in the City of New York, New York, NY, 10027, USA (e-mail {hw2786, sm4940, james.anderson}@columbia.edu).

only a fraction of clients to participate in each round of the training (optimization) process, and federated optimization algorithms must accommodate this *partial participation*.

A wealth of algorithms have been developed to address the aforementioned challenges. Notably, work in [19] proposed the now popular FedAvg algorithm, where each client performs multiple stochastic gradient descent (SGD) steps before sending the model to the server for aggregation. Subsequent efforts [14, 30, 33, 16, 25] provided theoretical analysis and further empirical performance evaluations. Since the proposal of FedAvg, there has been a rich body of work concentrating on developing federated optimization algorithms, such as; FedProx [27], FedSplit [23], Scaffold [12], FedLin [20], FedDyn [1], FedDR [32] and FedPD [38].

We consider a general unconstrained, composite optimization models formulated as

$$\frac{1}{n}\sum_{i=1}^{n}f_{i}(x) + g(x). \tag{1}$$

No convexity assumptions on f_i are made and g can be non-smooth. Of the previously mentioned federated algorithms, we restrict our attention to FedDR and FedPD. These algorithms are designed to alleviate the unrealistic assumptions required by FedAvg in order to realize desirable theoretical convergence rates. As described in [32], FedDR combines the nonconvex Douglas-Rachford splitting (DRS) algorithm [5] with a randomized block-coordinate strategy. FedDR provably converges when only a subset of clients participate in any given communication round. In contrast, FedPD is a primal-dual algorithm which requires either full participation or no participation by all clients at every round. Unlike FedDR, FedPD cannot handle optimization problems of the form of (1) for $q \neq 0$.

The key observation of this paper is to note that the updating rules of FedPD share a similar form to those of the alternating direction method of multipliers (ADMM) [9], but specifies how the local models are updated to satisfy the flexibility need of FL. Motivated by the fact that ADMM is the dual formulation of DRS [35, 7], we provide a new algorithm called FedADMM. Specifically, our contributions are:

1) By applying FedDR to the dual formulation of problem (1), we propose a new algorithm called FedADMM, which allows partial participation and solves the federated composite optimization problems as in [36].

- 2) When $g \equiv 0$ in problem (1), we find that FedADMM reduces to FedPD but requires only partial participation.
- 3) We prove equivalence between FedDR and FedADMM and provide a one-to-one and onto mapping between the iterates of both algorithms.
- 4) We provide convergence guarantees for FedADMM using the equivalence established in point 3.

Since FedADMM is the dual formulation of FedDR, it inherits all the desirable properties from FedDR. First, it can handle both *statistical* and *system* heterogeneity. Second, it allows inexact evaluation of users' proximal operators as in FedProx and FedPD. Third, by considering $g \not\equiv 0$ in (1), more general applications and problems with constraints can be considered [36].

A. Related Work

ADMM and DRS: DRS was first proposed in [5] in the context of providing numerical solutions to heat conduction partial differential equations. Subsequently, it found applications in the solution of convex optimization problems [18, 28] and later non-convex problems [15, 31]. ADMM [10, 3] is a very popular iterative algorithm for solving composite optimization problems. The equivalence between DRS and ADMM has been subject of a lot of work [8, 6, 35, 39]. It was first established for convex problems where ADMM is equivalent to applying DRS to the dual problem [8, 6]. Recently, these ideas were extended in [31] to show equivalence in the non-convex regime. Inspired by the fact that FedDR can be viewed as a variant of nonconvex DRS applied to the FL framework, we propose a new algorithm, FedADMM, and further extend the equivalence of these two algorithms to the FL paradigm.

Federated Learning: FedAvg was first proposed in [19]. However, it works well only with a homogeneous set of clients. It is difficult to analyze the convergence of FedAvg for the heterogeneous setting unless additional assumptions are made [16, 17, 13, 34]. The main reason for this is that the algorithm suffers from client-drift [41] under objective heterogeneity. To address the data and system heterogeneity, FedProx [27] was proposed by adding an extra proximal term [22] to the objective. However, this extra term might degrade the training performance so that FedProx doesn't converge to the global or local stationery points unless the step-size is carefully tuned. Another method called Scaffold [12] uses control variates (or variance reduction) to reduce client-drift at the cost of increased communication incurred by sending extra variables to the server. Fed-Split [23] applied the operator splitting schemes to remedy the objective heterogeneity issues, while it only considered the convex problems and required the full participation of clients. As mentioned earlier, FedDR [32] was inspired from DRS, and allowed partial participation. From the primaldual optimization perspective, FedPD [38] proposed a new concept of participation, which restricted its potential application to real problems. It is also worthwhile to mention that FedDyn [1] is equivalent to FedPD [38] from [37] under the full participation setting, but it allows partial participation. Unlike [36], FedPD and FedDyn can't solve non-smooth or constrained problems. Finally, we refer readers to [11] for a comprehensive understanding of the recent advances in FL.

II. PRELIMINARIES AND PROBLEM FORMULATION

We consider the canonical Federated learning optimization problem defined as

$$\min_{x \in \mathbb{R}^d} \left\{ F(x) = f(x) + g(x) \equiv \frac{1}{n} \sum_{i=1}^n f_i(x) + g(x) \right\}$$
 (2)

where n is the number of clients, f_i denotes the loss function associated to the i-th client. Each f_i is nonconvex and Lipschitz differentiable (see Assumptions 2.1 and 2.2 below), and g is a proper, closed, and convex function and is not necessarily smooth. For example, g could be any ℓ_p norm or an indicator function.

Assumption 1. F(x) is bounded below, i.e.,

$$\inf_{x \in \mathbb{R}^d} F(x) > -\infty \text{ and } dom(F) \neq \emptyset.$$

Assumption 2. (*Lipschitz differentiability*) Each $f_i(\cdot)$ in (2) has L-Lipschitz gradient, i.e.,

$$\|\nabla f_i(x) - \nabla f_i(y)\| \le L\|x - y\|$$

for all $i \in [n]$ and $x, y \in \mathbb{R}^d$.

The notation [n] above defines the set $\{1, 2, ..., n\}$. All norms are ℓ_2 - norms. We will frequently make use of the proximal operator [22]. Although typically defined for convex functions, we make no such assumptions.

Definition 1. (*Proximal operator*) Given an L-Lipschitz (possibly nonconvex and nonsmooth) function f, then the proximal mapping $\mathbb{R}^d \to (-\infty, \infty]$ is defined as

$$\operatorname{prox}_{\eta f}(x) = \arg\min_{y} \left\{ f(y) + \frac{1}{2\eta} ||x - y||^{2} \right\}.$$
 (3)

with parameter $\eta > 0$.

If f is nonconvex but L-Lipschitz, $\operatorname{prox}_{\eta f}(x)$ is still well-defined with $0 < \eta < 1/L$.

Definition 2. (Conjugate function) Let $f : \mathbb{R}^d \to \mathbb{R}$. The function $f^* : \mathbb{R}^d \to \mathbb{R}$ defined as

$$f^*(y) \triangleq \sup_{x \in \text{dom } f} (y^T x - f(x))$$

is called the conjugate function of f.

The conjugate function is is the piecewise supremum of a set of affine functions, hence it is closed and convex even when f is not.

Definition 3. (ε -stationarity) A vector x is said to be an ε -stationery solution to (2) if

$$\mathbb{E}\left[\left\|\nabla F(x)\right\|^{2}\right] \leq \varepsilon^{2},$$

where expectation is taken with respect to all random variables in the respective algorithm.

III. DOUGLAS-RACHFORD ALGORITHM

A. Douglas-Rachford Splitting

Douglas-Rachford Splitting (DRS) [5] is an iterative splitting algorithm for solving the optimization problems that can be written as

$$minimize_{x \in \mathbb{R}^d} \quad f(x) + g(x). \tag{4}$$

Although originally used for solving convex problems, it has been shown to work well on certain non-convex problems with additional structure. DRS solves problem (4) by producing a series of iterates (y_k, z_k, x_k) for $k = 1, 2, \ldots$ given by

$$\begin{cases} y_k &= \operatorname{prox}_{\eta f}(x_k) \\ z_k &= \operatorname{prox}_{\eta g}(2y_k - x_k) \\ x_{k+1} &= x_k + \alpha(z_k - y_k) \end{cases}$$
 (5)

where α is a relaxation parameter. When $\alpha = 1$, (5) is the classical Douglas-Rachford splitting and when $\alpha = 2$, (5) is a related splitting algorithm called Peaceman-Rachford splitting [24].

If f in problem (4) can be decomposed as $f(x) = \frac{1}{n} \sum_{i=1}^{n} f_i(x)$, then (5) can be modified so as to run in parallel if we include a global averaging step. The resulting algorithm is given below:

$$\begin{cases} y_{i}^{k+1} &= y_{i}^{k} + \alpha \left(\bar{x}^{k} - x_{i}^{k} \right), \quad \forall i \in [n] \\ x_{i}^{k+1} &= \operatorname{prox}_{\eta f_{i}} \left(y_{i}^{k+1} \right), \quad \forall i \in [n] \\ \hat{x}_{i}^{k+1} &= 2x_{i}^{k+1} - y_{i}^{k+1}, \quad \forall i \in [n] \\ \tilde{x}^{k+1} &= \frac{1}{n} \sum_{i=1}^{n} \hat{x}_{i}^{k+1}, \\ \bar{x}^{k+1} &= \operatorname{prox}_{\eta g} \left(\tilde{x}^{k+1} \right). \end{cases}$$
(6)

A full derivation is provided in [32]. Equation (6) is called full parallel Douglas-Rachford splitting (DRS).

B. FedDR

Implicit in the full parallel DRS (6), is the fact that all users are required to participate in every iteration. Instead of requiring all users $i \in [n]$ to participate as in (6), work in [32] proposed an inexact randomized block-coordinate DRS algorithm, called FedDR. Here, a subset \mathcal{S}_k of clients is sampled from a "proper" sampling scheme $\hat{\mathcal{S}}$ (See Definition 4 below for details) at each iteration. Each client, $i \in \mathcal{S}_k$ performs a local update (i.e., executes the first three steps in (6)), then sends its local model to server for aggregation. Each client $i \notin \mathcal{S}_k$ does noting. The complete FedDR algorithm is shown in Algorithm 1.

Convergence to an ϵ -stationary point of FedDR is guaranteed when the sampling scheme \hat{S} is proper and Assumption 1 and 2 hold [32].

Definition 4. Let $p = (p_1, p_2, \dots, p_n)$, where $p_i = \mathbb{P}(i \in \hat{S})$. If $p_i > 0$ for all $i \in [n]$, we call the sampling scheme \hat{S} proper, i.e., every client has a nonzero probability to be selected.

Assumption 3. All partial participation algorithms in this paper use a proper sampling scheme.

Algorithm 1 FL with Randomized DR (FedDR) [32]

1: **Initialize** $x^0, \eta, \alpha > 0, K$, and tolerances $\epsilon_{i,0} \geq 0$. 2: **Initialize** the server with $\bar{x}^0 = x^0$ and $\tilde{x}^0 = x^0$ **Initialize** each client $i \in [n]$ with $y_i^0 = x^0, x_i^0 \approx \operatorname{prox}_{nf_i}\left(y_i^0\right)$, and $\hat{x}_i^0 = 2x_i^0 - y_i^0$. 3: **for** k = 0, ..., K **do** Randomly sample $S_k \subseteq [n]$ with size S. 5: for each user $i \in \mathcal{S}_k$ do receive \bar{x}^k from the server. choose $\epsilon_{i,k+1} \geq 0$ and update choose $\epsilon_{i,k+1} \geq 0$ and update $y_i^{k+1} = y_i^k + \alpha \left(\bar{x}^k - x_i^k \right),$ $x_i^{k+1} \approx \operatorname{prox}_{\eta f_i} \left(y_i^{k+1} \right),$ $\hat{x}_i^{k+1} = 2x_i^{k+1} - y_i^{k+1}.$ send $\Delta \hat{x}_i^k = \hat{x}_i^{k+1} - \hat{x}_i^k$ back to the server . 10: 11: 12: 14: aggregation $\tilde{x}^{k+1} = \tilde{x}^k + \frac{1}{n} \sum_{i \in \mathcal{S}_k} \Delta \hat{x}_i^k$ update $\bar{x}^{k+1} = \text{prox}_{\eta g} \left(\tilde{x}^{k+1} \right)$

From the analysis in [26], this assumption includes a lot of sampling schemes such as non-overlapping uniform and doubly uniform sampling as special cases. The intuition behind proper sampling is to ensure that on average every client has a chance to be selected at every iteration.

In FedDR there are three variables that get updated: \bar{x}^k, x_i^k and y_i^k . The variable \bar{x}^k denotes the consensus/average variable to minimize the global model F, x_i^k denotes the local variable associated to f_i , while y_i^k measures the distance between the global variable \bar{x}^k and local model x_i^k . To account for the limitations on computation resources for local users, FedDR allows the inexact calculation of the proximal step, i.e.,

$$x_i^{k+1} \approx \operatorname{prox}_{\eta f_i} \left(y_i^{k+1} \right) \iff \|x_i^{k+1} - \operatorname{prox}_{\eta f_i} \left(y_i^{k+1} \right)\| \le \epsilon_{i,k+1}.$$

Thus \approx defines an ϵ -close solution. After local clients $i \in \mathcal{S}_k$ update their model and send them back to the server, the server aggregates the updates to update the global model by executing steps 16 and 17 in Algorithm 1.

IV. FROM FEDDR TO FEDADMM

Our first contribution is to derive the FedADMM algorithm from FedDR.

A. An equivalent formulation

17: end for

We begin by rewriting problem (2) as the equivalent constrained problem:

$$\min_{x \in \mathbb{R}^{nd}, \bar{x}} \left\{ F(x) = \frac{1}{n} \sum_{i=1}^{n} f_i(x_i) + g(\bar{x}) \right\}$$
s.t. $\mathbb{I}_{nd} x = \mathbb{1}\bar{x}$ (7)

where $x = \begin{bmatrix} x_1^T, x_2^T, \cdots, x_n^T \end{bmatrix}^T \in \mathbb{R}^{nd}$, \mathbb{I}_d is the $d \times d$ identity matrix, and $\mathbb{1} = \begin{bmatrix} \mathbb{I}_d & \cdots & \mathbb{I}_d \end{bmatrix}^T$. Here \bar{x} should be interpreted as the global consensus variable.

Forming the Lagrangian of (7) and using the definition of the conjugate function, the dual formulation of (7) is

$$\max_{z \in \mathbb{R}^{nd}} \left\{ F^*(z) = -f^*(-\mathbb{I}_{nd}z) - g^*(\mathbb{1}^T z) \right\}$$
 (8)

where $z = \begin{bmatrix} z_1^T, z_2^T, \cdots, z_n^T \end{bmatrix}^T \in \mathbb{R}^{nd}$ is the vector of dual variables. Problem (8) is clearly equivalent to

$$\min_{z_1, z_2, \dots, z_n} \left\{ \frac{1}{n} \sum_{i=1}^n f_i^* \left(-z_i \right) + g^* \left(\sum_{i=1}^n z_i \right) \right\}. \tag{9}$$

Before proceeding to develop an algorithm for solving (9), we first rewrite the full parallel DRS algorithm 6. Changing the execution order of (6) and choosing $\alpha = 1$ give

$$\begin{cases} \hat{x}_{i}^{k} &= 2x_{i}^{k} - y_{i}^{k}, \quad \forall i \in [n] \\ \tilde{x}^{k} &= \frac{1}{n} \sum_{i=1}^{n} \hat{x}_{i}^{k}, \quad \forall i \in [n] \\ \bar{x}^{k} &= \operatorname{prox}_{\eta g} \left(\tilde{x}^{k} \right), \\ x_{i}^{k+1} &= \operatorname{prox}_{\eta f_{i}} \left(y_{i}^{k} + \bar{x}^{k} - x_{i}^{k} \right), \quad \forall i \in [n] \\ y_{i}^{k+1} &= y_{i}^{k} + \bar{x}^{k} - x_{i}^{k}, \quad \forall i \in [n]. \end{cases}$$
(10)

Introducing the change of variables $w_i^k = x_i^k - y_i^k$, we have the following parallel DR algorithm

$$\begin{cases} \hat{x}_{i}^{k} &= x_{i}^{k} + w_{i}^{k}, \quad \forall i \in [n] \\ \tilde{x}^{k} &= \frac{1}{n} \sum_{i=1}^{n} \hat{x}_{i}^{k}, \quad \forall i \in [n] \\ \bar{x}^{k} &= \operatorname{prox}_{\eta g} \left(\tilde{x}^{k} \right), \\ x_{i}^{k+1} &= \operatorname{prox}_{\eta f_{i}} \left(\bar{x}^{k} - w_{i}^{k} \right), \quad \forall i \in [n]. \\ w_{i}^{k+1} &= w_{i}^{k} + x_{i}^{k+1} - \bar{x}^{k}, \quad \forall i \in [n] \end{cases}$$
(11)

Remark 1. Note that (6),(10) and (11) are essentially the same parallel algorithm under a change of execution order and variables.

B. FedDR-II

From section III, we observe that the only difference between full parallel DRS and FedDR is that FedDR only requires a subset of clients to update their variables, while full parallel DRS requires full participation. Similarly, by only considering partial participation in (11), we introduce the intermediate FedDR-II algorithm. We now describe each step of a single epoch of FedDR-II:

- 1) **Initialization:** Given an initial vector $x^0 \in \text{dom}(F)$ and tolerances $\epsilon_{i,0} \geq 0$. Initialize the server with $\bar{x}^0 =$ x^0 . Initialize all users $i \in [n]$ with $w_i^0 = 0$ and $x_i^0 = x^0$.
- 2) The k-th iteration: $(k \ge 0)$ Sample a proper subset $\mathcal{S}_k \subseteq [n]$ so that \mathcal{S}_k represents the subset of active clients.
- 3) Client update (Local): For each client $i \in S_k$, update $\hat{x}_i^k = x_i^k + w_i^k$. Clients $i \notin \mathcal{S}_k$ do nothing, i.e.

$$\begin{cases} \hat{x}_{i}^{k} = \hat{x}_{i}^{k-1} \\ x_{i}^{k} = x_{i}^{k-1} \\ w_{i}^{k} = w_{i}^{k-1} \end{cases}$$

- 4) Communication: Each user $i \in S_k$ sends only \hat{x}_i^k to
- 5) **Server update:** The server aggregates \tilde{x}^k $\frac{1}{n}\sum_{i=1}^{n}\hat{x}_{i}^{k}$, and then compute $\bar{x}^{k}=\operatorname{prox}_{\eta g}\left(\tilde{x}^{k}\right)$. Communication (Broadcast): Each user $i\in\mathcal{S}_{k}$ re-
- ceives \bar{x}^k from the server.
- 7) Client update (Local): For each user $i \in \mathcal{S}_k$, given $\epsilon_{i,k+1} \geq 0$, it updates

$$\begin{cases} x_i^{k+1} \approx \operatorname{prox}_{\eta f_i} \left(\bar{x}^k - w_i^k \right) \\ w_i^{k+1} = w_i^k + x_i^{k+1} - \bar{x}^k. \end{cases}$$

Each user $i \notin S_k$ does nothing, i.e.

$$\left\{\begin{array}{lcl} w_i^{k+1} & = & w_i^k \\ x_i^{k+1} & = & x_i^k \end{array}\right.$$

Remark 2. FedDR and FedDR-II are equivalent because they are partial participation versions of (6) and (11) respectively.

C. Solving the dual problem using FedDR-II

In this subsection, we use FedDR-II to solve the dual problem (9), introducing a new algorithm called FedADMM. We call this algorithm FedADMM because it is derived from applying FedDR-II to the dual problem (9). Let us define the augmented Lagrangian functions associated to (7) as

$$\mathcal{L}_{i}(x_{i}, \bar{x}^{k}, z_{i}) = f_{i}(x_{i}) + g(\bar{x}^{k}) + \left\langle z_{i}^{k}, x_{i} - \bar{x}^{k} \right\rangle + \frac{\eta}{2} \left\| x_{i} - \bar{x}^{k} \right\|^{2}.$$

$$(12)$$

where η denotes penalty parameter. Finally, we define $\Delta \hat{x}_i^k = \hat{x}_i^{k+1} - \hat{x}_i^k$. With everything defined, FedADMM is presented in Algorithm 2 below.

Algorithm 2 Federated ADMM Algorithm (FedADMM)

- 1: **Initialize** $x^0, \eta > 0, K$, and tolerances $\epsilon_{i,0} (i \in [n])$.
- 2: **Initialize** the server with $\bar{x}^0 = x^0$
- 3: **Initialize** all clients with $z_i^0 = 0$ and $x_i^0 = \hat{x}_i^0 = x^0$.
- 4: **for** k = 0, ..., K **do**
- Randomly sample $S_k \subseteq [n]$ with size S.
- for each client $i \in \mathcal{S}_k$ do
- receive \bar{x}^k from the server.
- $x_i^{k+1} \approx \arg\min \mathcal{L}_i \left(x_i, \bar{x}^k, z_i^k \right)$
- $\begin{aligned} z_i^{k+1} &= z_i^k + \eta \left(x_i^{k+1} \bar{x}^k \right) \\ z_i^{k+1} &= z_i^k + \eta \left(x_i^{k+1} \bar{x}^k \right) \quad \diamondsuit \text{Dual updates} \\ \hat{x}_i^{k+1} &= x_i^{k+1} + \frac{1}{\eta} z_i^{k+1} \\ \text{send } \Delta \hat{x}_i^k &= \hat{x}_i^{k+1} \hat{x}_i^k \text{ back to the server} \end{aligned}$
- 11:
- end for 13:
- 14:
- aggregation $\tilde{x}^{k+1} = \tilde{x}^k + \frac{1}{n} \sum_{i \in \mathcal{S}_k} \Delta \hat{x}_i^k$ update $\bar{x}^{k+1} = \text{prox}_{g/n} \left(\tilde{x}^{k+1} \right)$
- 16:
- 17: end for

When $g \equiv 0$, the server-side steps 15-16 of FedADMM reduce to the single step:

$$\bar{x}^{k+1} = \tilde{x}^{k+1} = \tilde{x}^k + \frac{1}{n} \sum_{i \in \mathcal{S}_k} \Delta \hat{x}_i^k = \frac{1}{n} \sum_{i=1}^n \hat{x}_i^{k+1}.$$

In this case, the updating rules of FedADMM are essentially the same as FedPD in [38]. Both compute the local model x_i^{k+1} by first minimizing (12), followed by updating the dual variable λ_i^{k+1} , and then aggregating \hat{x}_i^{k+1} to achieve the global model \bar{x}^{k+1} . However, FedADMM allows for partial participation (only chooses a subset of clients to update) while FedPD requires all clients to update at each communication round, making it less practical and applicable in real world scenarios.

Note that FedADMM can handle the case where $g\not\equiv 0$ whereas FedPD didn't consider this more general formulation. Just like step 11 (approximately evaluating $\operatorname{prox}_{\eta f_i}$) in FedDR, FedADMM obtains the new local model x_i^{k+1} by inexactly solving (12). Note that we do not specify how to (approximately) solve the proximal steps or Langrangian minimization step in (either) algorithm. Various oracles are specified in [38].

V. ANALYSIS

We now present the main theoretical results of the paper. Namely, an equivalence between FedDR and FedADMM. We leverage the FedDR convergence results [32] to show that FedADMM converges under partial participation.

We say that two iterative optimization algorithms are "equivalent" if they produce sequences $(x^k)_{k\geq 0}$ and $(y^k)_{k\geq 0}$ such that there exists a unique linear mapping between the two sequences. More general equivalence classes are defined and studied in [40].

Theorem 1. (Equivalence between FedDR and FedADMM) Let $(x_i^k, z_i^k, \bar{x}^k)_{k\geq 0}$ be a sequence generated by FedADMM with penalty parameter η , and $(s_i^k, u_i^k, \hat{u}_i^k, \bar{v}^k)$ a sequence generated by FedDR with parameter $\frac{1}{\eta}$. Then FedADMM and FedDR are equivalent.

Proof. For each triplet $(x_i^k, z_i^k, \bar{x}^k)$ at the k-th iteration of FedADMM with stepsize η , define

$$\begin{cases} s_i^k &= x_i^k - z_i^k/\eta \\ u_i^k &= x_i^k \\ \hat{u}_i^k &= x_i^k + z_i^k/\eta \end{cases} \quad \text{and} \quad \begin{cases} s_i^{k+1} &= x_i^{k+1} - z_i^{k+1}/\eta \\ u_i^{k+1} &= x_i^{k+1} \\ \hat{u}_i^{k+1} &= x_i^{k+1} + z_i^{k+1}/\eta \\ \bar{v}^{k+1} &= \bar{x}^{k+1} \end{cases}$$

Then $(s_i^k,u_i^k,\hat{u}_i^k,\bar{v}^k)$ and $(s_i^{k+1},u_i^{k+1},\hat{u}_i^{k+1},\bar{v}^{k+1})$ satisfy the updating rule of FedDR

$$\begin{cases} s_i^{k+1} &= s_i^k + (\bar{v}^k - u_i^k), \quad \forall i \in \mathcal{S}_k, \\ u_i^{k+1} &= \operatorname{prox}_{rf_i} \left(s_i^{k+1} \right), \quad \forall i \in \mathcal{S}_k, \\ \hat{u}_i^{k+1} &= 2u_i^{k+1} - s_i^{k+1}, \quad \forall i \in \mathcal{S}_k, \\ \bar{v}^{k+1} &= \operatorname{prox}_{rg} \left(\frac{1}{n} \sum_{i=1}^n \hat{u}_i^{k+1} \right), \end{cases}$$

where $r = 1/\eta$ and when $i \notin \mathcal{S}_k$

$$\begin{cases} s_i^{k+1} &= s_i^k, \\ u_i^{k+1} &= u_i^k, \\ \hat{u}_i^{k+1} &= \hat{u}_i^k \end{cases}$$

where the same sampling realizations S_k are used at each iteration for both algorithms.

We have

$$\begin{split} s_i^k + (\bar{v}^k - u_i^k) &= x_i^k - z_i^k / \eta + (\bar{x}^k - x_i^k) \\ &= x_i^{k+1} - z_i^k / \eta + \bar{x}^k - x_i^{k+1} \\ &\stackrel{(a)}{=} x_i^{k+1} - z_i^{k+1} / \eta = s_i^{k+1}. \end{split}$$

where (a) is due to the dual updates (line 10) in FedADMM algorithm. Moreover,

$$\begin{aligned} u_i^{k+1} &= x_i^{k+1} = \underset{x_i}{\operatorname{arg\,min}} \ \mathcal{L}_i \left(x_i, \bar{x}^k, z_i^k \right) \\ &= \underset{rf_i}{\operatorname{prox}}_{rf_i} (\bar{x}^k - z_i^k / \eta) \overset{(b)}{=} \underset{rf_i}{\operatorname{prox}}_{rf_i} (s_i^{k+1}) \end{aligned}$$

where (b) uses the fact that

$$\bar{x}^k - z_i^k / \eta = s_i^k + (\bar{v}^k - u_i^k) = s_i^{k+1}.$$

Finally, note that

$$\hat{u}_i^{k+1} = 2u_i^{k+1} - s_i^{k+1} = x_i^{k+1} + z_i^{k+1} / \eta,$$

which gives

$$\bar{v}^{k+1} = \bar{x}^{k+1} \stackrel{(c)}{=} \operatorname{prox}_{rg} \left(\sum_{i=1}^{n} \left(x_i^{k+1} + \frac{1}{\eta} z_i^{k+1} \right) \right)$$
$$= \operatorname{prox}_{rg} \left(\frac{1}{n} \sum_{i=1}^{n} \hat{u}_i^{k+1} \right)$$

where (c) comes from the FedADMM updating rule (line 11-16 in Alg 2).

Since we have proved the equivalence of FedDR and FedADMM for arbitrary (nonconvex) problems, FedADMM will directly inherit the convergence properties of FedDR, specifically at rate $\mathcal{O}(\frac{1}{k})$. The explicit convergence rate of FedADMM is characterized in the following theorem which is a direct application of Theorem 3.1 in [32].

Theorem 2. Suppose that Assumptions 1, 2, and 3 hold and $\gamma_1, \gamma_2, \gamma_3, \gamma_4 > 0$ are constants. Let $(x_i^k, z_i^k, \hat{x}_i^k, \bar{x}^k)_{k \geq 0}$ be generated by Alg 2 (FedADMM) using penalty parameter η that satisfies

$$\eta>\frac{4L\left(1+2\gamma_{4}\right)}{\sqrt{9-16\gamma_{4}\left(1+4\gamma_{4}\right)}-1}.$$

Then when $g \equiv 0$, the following holds

$$\frac{1}{K+1} \sum_{k=0}^{K} \mathbb{E}\left[\left\|\nabla f\left(\bar{x}^{k}\right)\right\|^{2}\right] \leq \frac{C_{1}\left[F\left(x^{0}\right) - F^{\star}\right]}{K+1} + \frac{1}{n(K+1)} \sum_{k=0}^{K} \sum_{i=1}^{n} \left(C_{2} \epsilon_{i,k}^{2} + C_{3} \epsilon_{i,k+1}^{2}\right)$$

where $\hat{\eta} = 1/\eta$, β , ρ_1 , and ρ_2 are defined as

$$\begin{cases} \beta &= \frac{\hat{\mathbf{p}}\left[2 - (L\hat{\eta} + 1) - 2L^2\hat{\eta}^2 - 4\gamma_4\left(1 + L^2\hat{\eta}^2\right)\right]}{2\hat{\eta}(1 + \gamma_1)(1 + L^2\hat{\eta}^2)} > 0 \\ \rho_2 &= \frac{2(1 + \hat{\eta}L)^2}{\gamma_4\hat{\eta}} + \frac{\left(1 + \hat{\eta}^2L^2\right)}{\hat{\eta}} \\ &+ \frac{\left[2 - (L\hat{\eta} + 1) - 2L^2\hat{\eta}^2 - 4\gamma_4\left(1 + L^2\hat{\eta}^2\right)\right]}{2\hat{\eta}(1 + L^2\hat{\eta}^2)\gamma_1} \\ \rho_1 &= \rho_2 + \frac{\left(1 + \hat{\eta}^2L^2\right)}{\hat{\eta}} \end{cases}$$

and the constants are

$$C_{1} = \frac{2(1+\hat{\eta}L)^{2}(1+\gamma_{2})}{\hat{\eta}^{2}\beta}, \quad C_{2} = \rho_{1}C_{1},$$

$$C_{3} = \rho_{2}C_{1} + \frac{(1+\hat{\eta}L)^{2}(1+\gamma_{2})}{\hat{\eta}^{2}\gamma_{2}}.$$

and $\hat{p} = \min \{p_i : i \in [n]\} > 0$ in Assumption 3.

Corollary 1. If the accuracy sequence $\epsilon_{i,k}$ (for all $i \in [n]$ and k > 0) at Step 8 in Alg 2 satisfies $\frac{1}{n} \sum_{i=1}^{n} \sum_{k=0}^{K+1} \epsilon_{i,k}^2 \leq D$ for a given constant D > 0 and all $K \geq 0$. Then, FedADMM needs

$$K = \left\lfloor \frac{C_1 \left[F\left(x^0\right) - F^{\star} \right] + \left(C_2 + C_3\right) D}{\varepsilon^2} \right\rfloor \equiv \mathcal{O}\left(\varepsilon^{-2}\right)$$

iterations to achieve $\frac{1}{K+1}\sum_{k=0}^K \mathbb{E}\left[\left\|\nabla f\left(\tilde{x}^k\right)\right\|^2\right] \leq \varepsilon^2$, where \tilde{x}^K is randomly selected from $\{\bar{x}^0, \bar{x}^1, \cdots, \bar{x}^K\}$. In other words, after $K = \mathcal{O}(\varepsilon^{-2})$ iterations, \tilde{x}^K is an ε -stationary solution of problem (2) when $g \equiv 0$.

Remark 3. Our convergence analysis can be easily extended to $g \not\equiv 0$, as long as we change the suboptimal condition into the gradient mapping as in [32]. To make $\frac{1}{n} \sum_{i=1}^{n} \sum_{k=0}^{K+1} \epsilon_{i,k}^2 \leq D$ hold, interested readers could refer to Remark 3.1 in [32].

Remark 4. Although FedADMM is a partial participation version of FedPD when $g \equiv 0$, its communication complexity is still $\mathcal{O}(\varepsilon^{-2})$, which matches the lower bound (up to constant factors) in [38].

VI. NUMERICAL SIMULATIONS

To demonstrate the equivalence of FedDR and FedADMM, we conduct simulations on both synthetic and real datasets. It is worthwhile to mention that our goal is to show the equivalence of the algorithms and the effect of partial participation, not to compare their performance with other algorithms. Performance profiling of FedPD and FedDR can be found in [32, 38]. We have not attempted to optimize any hyper-parameters. All the experiments are run on a Lambda machine with CPU(AMD Threadripper Pro 3975wx) and GPU(RTX A6000). The code is available at https://github.com/jd-anderson/Federated. **Datasets:** We first generate synthetic non-iid datasets by following the same setup as in [29] and denote them as synthetic- (α, β) . Here α controls how much local models differ from each other and β controls how much the local data at each device differs from that of other devices. We run the experiments by using the unbalanced datasets: synthetic-(0, 0), synthetic-(0.5, 0.5) and synthetic-(1, 1). We then compare FedADMM with FedDR on the FEMNIST data set [4]. FEMNIST is a more complex 62-class Federated Extended MNIST dataset. It consists of handwritten characters including: numbers 1-10, 26 upper-and lower-case letters A-Z and a-z from different writers and is also separated by the writers, therefore the dataset is non-iid.

Models and Hyper-parameters: For all the synthetic datasets, we use the model described in [32]: a neural network with a single hidden layer. The network architecture is $60 \times 32 \times 10$ corresponding to input layer \times hidden layer × output layer size. For FEMNIST data, we use the same model as [4], which consists of 2 convolutional layers and two fully connected layers, with 62 neurons in the output layer matching the number of classes in the FEMNIST dataset. For all the experiments, we use $\eta = 1$ (in FedADMM) and $\alpha = 1$ (in FedDR). As in [38], we choose stochastic gradient descent as a local solver with 300 local iterations to solve step 11 in FedDR and step 9 in FedADMM. The mini-batch size in calculating the stochastic gradient is 2 and the learning rate is 0.01. We stress that we do not attempt to optimize these parameters. Implementation: We use the uniform sampling scheme to select the clients in each round. We set the number of active clients in each round as 10 out of 30 in Fig 2 and test the algorithms with 10 out of 20 users in Fig 1. We also compare the performance of FedADMM using different participation numbers in Fig 3. To provide a fair comparison, we use the same random seeds across all algorithms.

After running multiple experiments on different datasets and models, from figure 1 and 2 we observe that the training accuracy and loss of FedDR and FedADMM coincide at each iteration, which verifies our analysis in section V.

VII. CONCLUSION

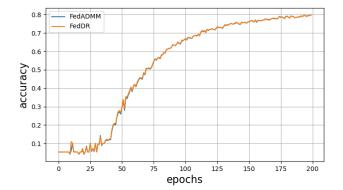
We have developed a new federated learning algorithm, FedADMM, for finding stationary points in non-convex composite optimization problems. Current work is focused on incorporating convex constraints into the algorithm and applying it to non-localizable model predictive control problems where communication efficiency is necessary [2].

ACKNOWLEDGEMENTS

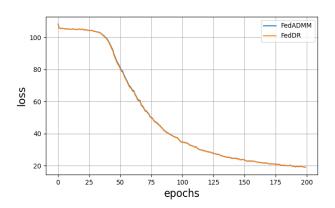
HW is generously funded by a Wei family fellowship and the Columbia Data Science Institute. We also acknowledge funding from the DoE under grant DE-SC0022234 and NSF career award 2144634.

REFERENCES

- D. A. E. Acar, Y. Zhao, R. M. Navarro, M. Mattina, P. N. Whatmough, and V. Saligrama, "Federated learning based on dynamic regularization," arXiv preprint arXiv:2111.04263, 2021.
- [2] C. A. Alonso, J. Shuang, J. Anderson, and N. Matni, "Distributed and localized model predictive control. part i: Synthesis and implementation," arXiv preprint arXiv:2110.07010, 2021.
- [3] S. Boyd, N. Parikh, E. Chu, B. Peleato, J. Eckstein *et al.*, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Foundations and Trends*® *in Machine learning*, vol. 3, no. 1, pp. 1–122, 2011.
- [4] S. Caldas, S. M. K. Duddu, P. Wu, T. Li, J. Konečný, H. B. McMahan, V. Smith, and A. Talwalkar, "Leaf: A benchmark for federated settings," arXiv preprint arXiv:1812.01097, 2018.
- [5] J. Douglas and H. H. Rachford, "On the numerical solution of heat conduction problems in two and three space variables," *Transactions of the American mathematical Society*, vol. 82, no. 2, pp. 421–439, 1956.



(a) FEMNIST accuracy



(b) FEMNIST loss

Figure 1: Performance of FL algorithms with partial participation(10 of 20 users) in terms of training accuracy and cross-entropy training loss of FEMNIST dataset

- [6] J. Eckstein, "Splitting methods for monotone operators with applications to parallel optimization," Ph.D. dissertation, Massachusetts Institute of Technology, 1989.
- [7] M. Fukushima, "Application of the alternating direction method of multipliers to separable convex programming problems," *Computational Optimization and Applications*, vol. 1, no. 1, pp. 93–111, 1992.
- [8] D. Gabay, "Chapter ix applications of the method of multipliers to variational inequalities," in *Studies in mathematics* and its applications. Elsevier, 1983, vol. 15, pp. 299–331.
- [9] D. Gabay and B. Mercier, "A dual algorithm for the solution of nonlinear variational problems via finite element approximation," *Computers & mathematics with applications*, vol. 2, no. 1, pp. 17–40, 1976.
- [10] P. Giselsson and S. Boyd, "Linear convergence and metric selection for Douglas-Rachford splitting and ADMM," *IEEE Transactions on Automatic Control*, vol. 62, no. 2, pp. 532–544, 2016.
- [11] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings et al., "Advances and open problems in federated learning," Foundations and Trends® in Machine Learning, vol. 14, no. 1–2, pp. 1–210, 2021.
- [12] S. P. Karimireddy, S. Kale, M. Mohri, S. Reddi, S. Stich, and

- A. T. Suresh, "Scaffold: Stochastic controlled averaging for federated learning," in *International Conference on Machine Learning*. PMLR, 2020, pp. 5132–5143.
- [13] A. Khaled, K. Mishchenko, and P. Richtárik, "First analysis of local gd on heterogeneous data," arXiv preprint arXiv:1909.04715, 2019.
- [14] J. Konečný, H. B. McMahan, F. X. Yu, P. Richtárik, A. T. Suresh, and D. Bacon, "Federated learning: Strategies for improving communication efficiency," arXiv preprint arXiv:1610.05492, 2016.
- [15] G. Li and T. K. Pong, "Global convergence of splitting methods for nonconvex composite optimization," SIAM Journal on Optimization, vol. 25, no. 4, pp. 2434–2460, 2015.
- [16] X. Li, K. Huang, W. Yang, S. Wang, and Z. Zhang, "On the convergence of fedavg on non-iid data," arXiv preprint arXiv:1907.02189, 2019.
- [17] X. Li and F. Orabona, "On the convergence of stochastic gradient descent with adaptive stepsizes," in *The 22nd Inter*national Conference on Artificial Intelligence and Statistics. PMLR, 2019, pp. 983–992.
- [18] P.-L. Lions and B. Mercier, "Splitting algorithms for the sum of two nonlinear operators," SIAM Journal on Numerical Analysis, vol. 16, no. 6, pp. 964–979, 1979.
- [19] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Artificial intelligence* and statistics. PMLR, 2017, pp. 1273–1282.
- [20] A. Mitra, R. Jaafar, G. Pappas, and H. Hassani, "Linear Convergence in Federated Learning: Tackling Client Heterogeneity and Sparse Gradients," *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [21] M. Mohri, G. Sivek, and A. T. Suresh, "Agnostic federated learning," in *International Conference on Machine Learning*. PMLR, 2019, pp. 4615–4625.
- [22] N. Parikh and S. Boyd, "Proximal algorithms," Foundations and Trends in optimization, vol. 1, no. 3, pp. 127–239, 2014.
- [23] R. Pathak and M. J. Wainwright, "FedSplit: An algorithmic framework for fast federated optimization," Advances in Neural Information Processing Systems, vol. 33, pp. 7057–7066, 2020
- [24] D. W. Peaceman and H. H. Rachford, Jr, "The numerical solution of parabolic and elliptic differential equations," *Journal of the Society for industrial and Applied Mathematics*, vol. 3, no. 1, pp. 28–41, 1955.
- [25] A. Reisizadeh, A. Mokhtari, H. Hassani, A. Jadbabaie, and R. Pedarsani, "FedPAQ: A communication-efficient federated learning method with periodic averaging and quantization," in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2020, pp. 2021–2031.
- [26] P. Richtárik and M. Takáč, "Parallel coordinate descent methods for big data optimization," *Mathematical Programming*, vol. 156, no. 1, pp. 433–484, 2016.
- [27] A. K. Sahu, T. Li, M. Sanjabi, M. Zaheer, A. Talwalkar, and V. Smith, "On the convergence of federated optimization in heterogeneous networks," arXiv preprint arXiv:1812.06127, vol. 3, p. 3, 2018.
- [28] J. H. Seidman, M. Fazlyab, V. M. Preciado, and G. J. Pappas, "A control-theoretic approach to analysis and parameter selection of Douglas–Rachford splitting," *IEEE Control Systems Letters*, vol. 4, no. 1, pp. 199–204, 2019.
- [29] O. Shamir, N. Srebro, and T. Zhang, "Communication-efficient distributed optimization using an approximate newton-type method," in *International conference on machine learning*. PMLR, 2014, pp. 1000–1008.
- [30] S. U. Stich, "Local SGD converges fast and communicates little," arXiv preprint arXiv:1805.09767, 2018.
- [31] A. Themelis and P. Patrinos, "Douglas-Rachford splitting

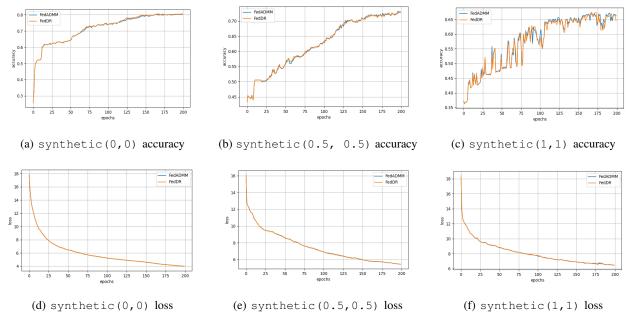


Figure 2: Performance of FL algorithms with partial participation(10 of 30 users) in terms of training accuracy and cross-entropy training loss of synthetic datasets

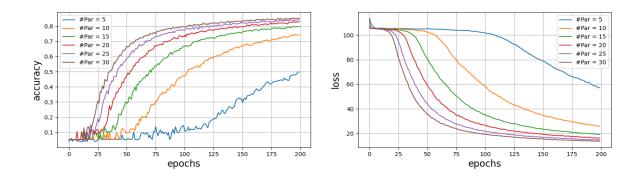


Figure 3: FedADMM performance comparison on FEMNIST dataset with different number of partial users.

- and ADMM for nonconvex optimization: Tight convergence results," *SIAM Journal on Optimization*, vol. 30, no. 1, pp. 149–181, 2020.
- [32] Q. Tran Dinh, N. Pham, D. Phan, and L. Nguyen, "FedDR-randomized Douglas-Rachford splitting algorithms for non-convex federated composite optimization," *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [33] J. Wang and G. Joshi, "Cooperative SGD: A unified framework for the design and analysis of communication-efficient SGD algorithms," arXiv preprint arXiv:1808.07576, 2018.
- [34] S. Wang, T. Tuor, T. Salonidis, K. K. Leung, C. Makaya, T. He, and K. Chan, "Adaptive federated learning in resource constrained edge computing systems," *IEEE Journal on Selected Areas in Communications*, vol. 37, no. 6, pp. 1205– 1221, 2019.
- [35] M. Yan and W. Yin, "Self equivalence of the alternating direction method of multipliers," in Splitting Methods in Communication, Imaging, Science, and Engineering. Springer,

- 2016, pp. 165-194.
- [36] H. Yuan, M. Zaheer, and S. Reddi, "Federated composite optimization," in *International Conference on Machine Learning*. PMLR, 2021, pp. 12253–12266.
- [37] X. Zhang and M. Hong, "On the Connection Between Fed-Dyn and FedPD," 2021.
- [38] X. Zhang, M. Hong, S. Dhople, W. Yin, and Y. Liu, "FedPD: A federated learning framework with adaptivity to non-iid data," *IEEE Transactions on Signal Processing*, vol. 69, pp. 6055–6070, 2021.
- [39] S. Zhao, L. Lessard, and M. Udell, "An automatic system to detect equivalence between iterative algorithms," *arXiv* preprint arXiv:2105.04684, 2021.
- [40] ——, "An automatic system to detect equivalence between iterative algorithms," arXiv preprint arXiv:2105.04684, 2021.
- [41] Y. Zhao, M. Li, L. Lai, N. Suda, D. Civin, and V. Chandra, "Federated learning with non-iid data," arXiv preprint arXiv:1806.00582, 2018.