

# Fast Composite Optimization and Statistical Recovery in Federated Learning

Yajie Bao <sup>1</sup>   Michael Crawshaw <sup>2</sup>   Shan Luo <sup>1</sup>   Mingrui Liu <sup>2</sup>

<sup>1</sup>School of Mathematical Sciences, Shanghai Jiao Tong University

<sup>2</sup>Department of Computer Science, George Mason University

June 26, 2022



SHANGHAI JIAO TONG  
UNIVERSITY



# Outline

- 1. Introduction**
2. Fast Composite Optimization in FL
3. Fast Statistical Recovery in FL
4. Experiment Results

# Federated Learning Environment

Consider a federated learning problem with  $K$  clients,

- Local loss:  $\mathcal{L}_k(\mathbf{w})$ ;
- Local weight:  $\pi_k$ ;
- Non-smooth regularizer  $h(\cdot)$ .

We will consider two instances of the FL problem:

$$\arg \min_{\mathbf{w} \in \mathcal{W}} \left\{ \sum_{k=1}^K \pi_k \mathcal{L}_k(\mathbf{w}) + h(\mathbf{w}) \right\}. \quad (1.1)$$

We denote the **global loss** as  $\mathcal{L} = \sum_{k=1}^K \pi_k \mathcal{L}_k$  and the **global composite objective** as  $\phi = \mathcal{L} + h$ .

Examples: federated sparse linear regression, federated low-rank matrix estimation.

# Federated Composite Optimization

For general **composite optimization** in FL,

- Local population distribution:  $\mathcal{P}_k$ ;
- Local population loss:  $\mathcal{L}_k(\mathbf{w}) = \mathbb{E}_{\xi \sim \mathcal{P}_k} [f(\mathbf{w}; \xi)]$ ;
- Non-smooth convex regularizer:  $h(\cdot)$ .

With corresponding objective

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathcal{W}} \left\{ \sum_{k=1}^K \pi_k \mathcal{L}_k(\mathbf{w}) + h(\mathbf{w}) \right\}. \quad (1.2)$$

Matches previous analyses which ignore the statistical estimation problem.

# Federated Statistical Recovery

For **statistical recovery** problem in FL,

- Local empirical distribution:  $\mathcal{D}_k$ ;
- Local empirical loss:  $\mathcal{L}_k(\mathbf{w}) = \mathbb{E}_{\xi \sim \mathcal{D}_k} [f(\mathbf{w}; \xi)]$ ;
- Non-smooth norm regularizer:  $\mathcal{R}(\cdot)$ , which is decomposable [Negahban et al., 2012];
- The ground-truth parameter:  $\mathbf{w}^* = \arg \min_{\mathbf{w} \in \mathcal{W}} \sum_{k=1}^K \pi_k \mathbb{E}_{\xi \sim \mathcal{P}_k} [f(\mathbf{w}; \xi)]$ .

We want to obtain the ground-truth parameter  $\mathbf{w}^*$  through solving:

$$\hat{\mathbf{w}}_{\text{opt}} = \arg \min_{\mathbf{w} \in \mathcal{W}} \left\{ \sum_{k=1}^K \pi_k \mathcal{L}_k(\mathbf{w}) + \lambda_{\text{opt}} \mathcal{R}(\mathbf{w}) \right\}, \quad (1.3)$$

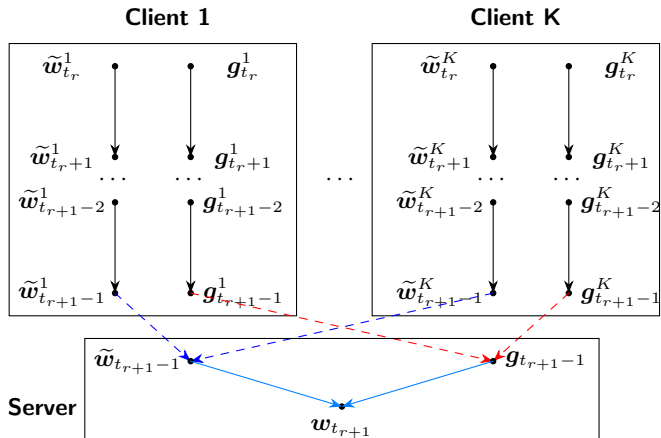
where  $\lambda_{\text{opt}}$  is the *optimal regularization parameter*.

Our goal is to design an algorithm to recover  $\mathbf{w}^*$  with the **optimal statistical precision**  
 $\|\hat{\mathbf{w}}_{\text{opt}} - \mathbf{w}^*\| \lesssim \epsilon_{\text{stat}} := \Psi(\bar{\mathcal{M}}) \lambda_{\text{opt}} / \mu$  [Negahban et al., 2012].

# Outline

1. Introduction
- 2. Fast Composite Optimization in FL**
3. Fast Statistical Recovery in FL
4. Experiment Results

## Fast-FedDA



# Convergence Rate of Fast-FedDA

We impose the following assumptions:

- 1  $\mathcal{L}_k$  for  $k \in [K]$  are  $\mu$ -strongly convex and  $L$ -smooth.
- 2  $h : \mathcal{W} \rightarrow \mathbb{R}$  is a closed convex function.
- 3 Bounded heterogeneity:  $\|\nabla \mathcal{L}_k(\mathbf{w}) - \nabla \mathcal{L}(\mathbf{w})\| \leq H$  for any  $\mathbf{w} \in \mathcal{W}$ .
- 4 Bounded variance of stochastic gradient:  $\mathbb{E}_{\xi \sim \mathcal{P}_k} [\|\nabla f(\mathbf{w}; \xi) - \nabla \mathcal{L}_k(\mathbf{w})\|^2] \leq \sigma^2$ .

## Theorem

Under Assumptions 1-4, with  $\alpha_t = (t + a)^2$  and  $\gamma = \mu a^3$  with  $a \geq 4L/\mu$  in Algorithm Fast-FedDA,

$$\mathbb{E}_{\mathcal{P}} [\phi(\hat{\mathbf{w}}_{\text{Fast-FedDA}}) - \phi(\hat{\mathbf{w}})] \leq O\left(\frac{\sigma^2}{K\mu T}\right) \quad (2.1)$$

for equal-weighted case ( $\pi_1 = \dots = \pi_K = 1/K$ ).



# Outline

1. Introduction
2. Fast Composite Optimization in FL
- 3. Fast Statistical Recovery in FL**
4. Experiment Results

## RSC and RSM

Widely used in statistical recovery literature [Agarwal et al., 2012, Wang et al., 2014, Loh and Wainwright, 2015, Cai et al., 2020]. Denote

$$\mathcal{T}_k(\mathbf{w}, \mathbf{w}') = \mathcal{L}_k(\mathbf{w}) - \mathcal{L}_k(\mathbf{w}') - \langle \nabla \mathcal{L}_k(\mathbf{w}'), \mathbf{w} - \mathbf{w}' \rangle.$$

### Assumption (RSC)

*The local loss functions  $\mathcal{L}_k$  for  $k = 1, \dots, K$  are convex and there exist  $\mu > 0$  and  $\tau_k \geq 0$  such that for any  $\mathbf{w}, \mathbf{w}' \in \mathcal{W}$ :*

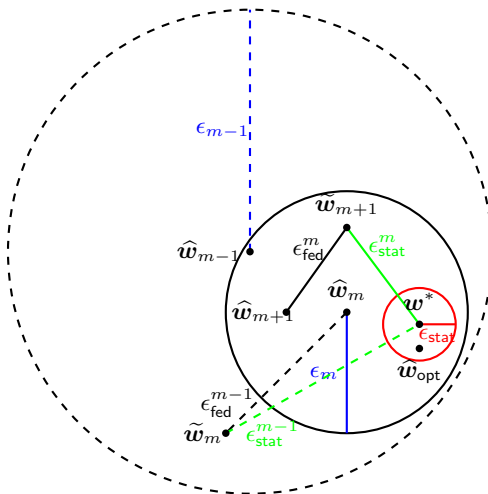
$$\mathcal{T}_k(\mathbf{w}, \mathbf{w}') \geq \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}'\|^2 - \tau_k \mathcal{R}^2(\mathbf{w} - \mathbf{w}').$$

### Assumption (RSM)

*For the local loss functions  $\mathcal{L}_k$  for  $k = 1, \dots, K$ , there exist  $L > 0$  and  $\nu_k \geq 0$  such that for any  $\mathbf{w}, \mathbf{w}' \in \mathcal{W}$ :*

$$\mathcal{T}_k(\mathbf{w}, \mathbf{w}') \leq \frac{L}{2} \|\mathbf{w} - \mathbf{w}'\|^2 + \nu_k \mathcal{R}^2(\mathbf{w} - \mathbf{w}').$$

## Multi-stage Constrained FedDA



**Figure:** Optimization landscape in  $\mathcal{R}$  norm at the  $m$ -th stage of MC-FedDA.

## Complexity of MC-FedDA

With high probability, we can guarantee bounds on optimization error and estimation error:

- Optimization:  $\phi(\hat{\mathbf{w}}_{\text{MC-FedDA}}) - \phi(\hat{\mathbf{w}}_{\text{opt}}) \leq \frac{\Psi^2(\bar{\mathcal{M}})\lambda_{\text{opt}}^2}{\mu}.$
- Estimation:  $\|\hat{\mathbf{w}}_{\text{MC-FedDA}} - \mathbf{w}^*\| \leq \frac{4\Psi(\bar{\mathcal{M}})\lambda_{\text{opt}}}{\mu} = 4\epsilon_{\text{stat}}.$

When  $\pi_1 = \dots = \pi_K = 1/K$ :

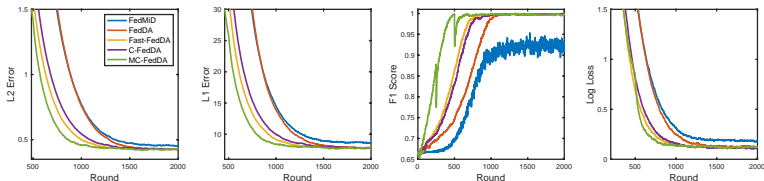
- $\epsilon_{\text{stat}} = \Psi(\bar{\mathcal{M}})\lambda_{\text{opt}}/\mu$  converges to 0 as the total sample size  $N \rightarrow \infty$ .
- Let  $\epsilon = \mu\epsilon_{\text{stat}}^2$ , if the total number of iterations satisfies  $T = \tilde{\mathcal{O}}(\Psi(\bar{\mathcal{M}})^2\sigma^2/(K\mu\epsilon))$ , then we are guaranteed that  $\phi(\hat{\mathbf{w}}_{\text{MC-FedDA}}) - \phi(\hat{\mathbf{w}}_{\text{opt}}) \leq \epsilon$ .
- The total communication complexity is bounded by  $\tilde{\mathcal{O}}(T^{1/2}K^{1/2})$ , matching the best known result of FedAvg for unconstrained problem [Woodworth et al., 2020, Stich, 2019].

# Outline

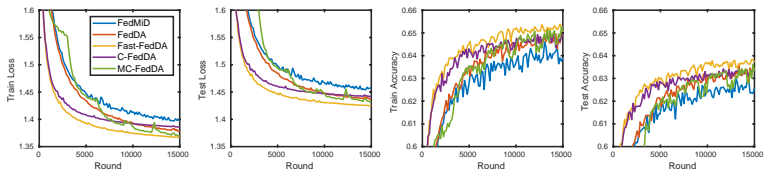
1. Introduction
2. Fast Composite Optimization in FL
3. Fast Statistical Recovery in FL
- 4. Experiment Results**

## Experiment Results

We compare our proposed algorithms with Federated Mirror Descent (FedMiD) and Federated Dual Averaging (FedDA) algorithms introduced in [Yuan et al. \[2021\]](#).



**Figure:** Recovery results for federated sparse linear regression problem.



**Figure:** Results for federated sparse logistic regression on EMNIST-62 dataset.

## Reference I

- A. Agarwal, S. Negahban, and M. J. Wainwright. Fast global convergence of gradient methods for high-dimensional statistical recovery. *The Annals of Statistics*, pages 2452–2482, 2012.
- T. T. Cai, Y. Wang, and L. Zhang. The cost of privacy in generalized linear models: Algorithms and minimax lower bounds. *arXiv preprint arXiv:2011.03900*, 2020.
- P.-L. Loh and M. J. Wainwright. Regularized M-estimators with nonconvexity: Statistical and algorithmic theory for local optima. *Journal of Machine Learning Research*, 16(1):559–616, 2015.
- S. N. Negahban, P. Ravikumar, M. J. Wainwright, and B. Yu. A unified framework for high-dimensional analysis of  $M$ -estimators with decomposable regularizers. *Statistical Science*, 27(4):538–557, 2012.
- S. U. Stich. Local SGD converges fast and communicates little. In *International Conference on Learning Representations*, 2019.
- Z. Wang, H. Liu, and T. Zhang. Optimal computational and statistical rates of convergence for sparse nonconvex learning problems. *The Annals of statistics*, 42(6):2164, 2014.
- B. Woodworth, K. K. Patel, and N. Srebro. Minibatch vs local SGD for heterogeneous distributed learning. *arXiv preprint arXiv:2006.04735*, 2020.
- H. Yuan, M. Zaheer, and S. Reddi. Federated composite optimization. In *International Conference on Machine Learning*, pages 12253–12266. PMLR, 2021.