

# Nonconvex Federated Learning on Compact Smooth Submanifolds With Heterogeneous Data

Jiaojiao Zhang<sup>1</sup>, **Jiang Hu**<sup>2</sup>, Anthony Man-cho So<sup>3</sup>, Mikael Johansson<sup>1</sup>

<sup>1</sup>KTH Royal Institute of Technology,

<sup>2</sup>University of California, Berkeley, <sup>3</sup>The Chinese University of Hong Kong





# Introduction

- ▶ Consider federated learning (FL) on manifolds

$$\underset{x \in \mathcal{M}}{\text{minimize}} \quad f(x) := \frac{1}{n} \sum_{i=1}^n f_i(x), \quad f_i(x) = \frac{1}{m_i} \sum_{l=1}^{m_i} f_{il}(x; \mathcal{D}_{il}) \quad (1)$$

- Each client has local loss  $f_i$  that is smooth but nonconvex
- Local datasets  $\mathcal{D}_i$  across clients  $i$  are heterogeneous
- $\mathcal{M}$  is a **compact smooth submanifold** embedded in  $\mathbb{R}^{d \times k}$ , with Euclidean metric serving as its Riemannian metric. E.g., Stiefel manifold:  $\text{St}(d, k) = \{x \in \mathbb{R}^{d \times k} : x^T x = I_k\}$

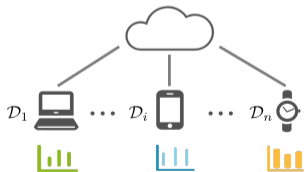


Figure: Federated learning

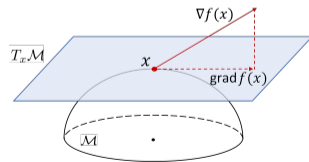


Figure: Optimization on manifolds



# Challenges and Contributions

## Challenges

- ▶ Single-machine optimization on  $\mathcal{M}$  cannot be directly adapted to FL
  - Even if each local model lies on  $\mathcal{M}$ , their average typically does not
- ▶ Extending FL algorithms to manifold optimization is not straightforward
  - $\mathcal{M}$  is nonconvex
- ▶ FL algorithms with local updates need substantial modifications to accommodate  $\mathcal{M}$ 
  - Client drift issue due to local updates and heterogeneous data persists

## Contributions

- ▶ Propose a computation- and communication-efficient algorithm for solving (1)
- ▶ Establish sub-linear convergence to a neighborhood of a first-order optimal solution
- ▶ Demonstrate superior performance over alternative methods



# Algorithm Intuition and Innovations

- ▶ The equivalent and compact form of our algorithm is

$$\begin{cases} \hat{\mathbf{z}}_{t+1}^r = \hat{\mathbf{z}}_t^r - \eta \left( \underbrace{\text{gradf}(\mathbf{z}_t^r; \mathcal{B}_t^r)}_{\text{new}} + \underbrace{\frac{1}{\tau} \sum_{t=0}^{\tau-1} \overline{\text{gradf}(\mathbf{z}_t^{r-1}; \mathcal{B}_t^{r-1})}}_{\text{average}} - \underbrace{\frac{1}{\tau} \sum_{t=0}^{\tau-1} \text{gradf}(\mathbf{z}_t^{r-1}; \mathcal{B}_t^{r-1})}_{\text{old}} \right) \\ \mathbf{z}_{t+1}^r = \mathcal{P}_{\mathcal{M}}(\hat{\mathbf{z}}_{t+1}^r) \\ \mathbf{x}^{r+1} = \mathcal{P}_{\mathcal{M}}(\mathbf{x}^r) - \eta_g \eta \sum_{t=0}^{\tau-1} \overline{\text{gradf}(\mathbf{z}_t^r; \mathcal{B}_t^r)} \end{cases}$$

- ▶ Mimic centralized projected Riemannian gradient descent
  - When  $\tau = 1$  and  $b = m_i$ , we recover  $\tilde{\mathbf{x}}^{r+1} := \mathcal{P}_{\mathcal{M}}\left(\mathcal{P}_{\mathcal{M}}(\bar{\mathbf{x}}^r) - \tilde{\eta} \cdot \text{grad}f(\mathcal{P}_{\mathcal{M}}(\bar{\mathbf{x}}^r))\right)$
- ▶ Feasibility of all iterates at a low computational cost by using  $\mathcal{P}_{\mathcal{M}}$ 
  - Avoid exponential mapping, inverse exponential mapping, and parallel transport
- ▶ Overcome client drift
  - Correction employs “variance reduction” and does not incur extra communication



# Algorithm Implementation

- 1: **Input:**  $R, \tau, \eta, \eta_g, \tilde{\eta} = \eta\eta_g\tau, \bar{x}^1$ , and  $c_i^1 = 0$  for all  $i \in [n]$
- 2: **for**  $r = 1, 2, \dots, R$  **do**
- 3:   **Client**  $i$
- 4:   Set  $\hat{z}_{i,0}^r = \mathcal{P}_{\mathcal{M}}(\bar{x}^r)$  and  $z_{i,0}^r = \mathcal{P}_{\mathcal{M}}(\bar{x}^r)$
- 5:   **for**  $t = 0, 1, \dots, \tau - 1$  **do**
- 6:     Sample a mini-batch dataset  $\mathcal{B}_{i,t}^r \subseteq \mathcal{D}_i$  with  $|\mathcal{B}_{i,t}^r| = b$
- 7:     Update  $\text{grad}f_i(z_{i,t}^r; \mathcal{B}_{i,t}^r) = \frac{1}{b} \sum_{\mathcal{D}_{il} \in \mathcal{B}_{i,t}^r} \text{grad}f_{il}(z_{i,t}^r; \mathcal{D}_{il})$
- 8:     Update  $\hat{z}_{i,t+1}^r = \hat{z}_{i,t}^r - \eta (\text{grad}f_i(z_{i,t}^r; \mathcal{B}_{i,t}^r) + c_i^r)$
- 9:     Update  $z_{i,t+1}^r = \mathcal{P}_{\mathcal{M}}(\hat{z}_{i,t+1}^r)$
- 10:   **end for**
- 11:   Send  $\hat{z}_{i,\tau}^r$  to the server
- 12:   **Server**
- 13:   Update  $\bar{x}^{r+1} = \mathcal{P}_{\mathcal{M}}(\bar{x}^r) + \eta_g \left( \frac{1}{n} \sum_{i=1}^n \hat{z}_{i,\tau}^r - \mathcal{P}_{\mathcal{M}}(\bar{x}^r) \right)$
- 14:   Broadcast  $\bar{x}^{r+1}$  to all the clients
- 15:   **Client**  $i$
- 16:   Update  $c_i^{r+1} = \frac{1}{\eta_g\eta\tau} (\mathcal{P}_{\mathcal{M}}(\bar{x}^r) - \bar{x}^{r+1}) - \frac{1}{\tau} \sum_{t=0}^{\tau-1} \text{grad}f_i(z_{i,t}^r; \mathcal{B}_{i,t}^r)$
- 17:   **end for**
- 18: **Output:**  $\mathcal{P}_{\mathcal{M}}(\bar{x}^{R+1})$



## Definitions

- ▶ (Riemannian gradient): The Riemannian gradient  $\text{grad}f(x)$  of a function  $f$  at the point  $x \in \mathcal{M}$  is the unique tangent vector that satisfies

$$\langle \text{grad}f(x), \xi \rangle_x = df(x)[\xi], \quad \forall \xi \in T_x \mathcal{M}$$

- For a submanifold  $\mathcal{M}$ ,  $\text{grad}f(x)$  can be computed as  $\text{grad}f(x) = \mathcal{P}_{T_x \mathcal{M}}(\nabla f(x))$
- ▶ ( $\hat{\gamma}$ -proximal smoothness of  $\mathcal{M}$ ): The  $\hat{\gamma}$ -tube around  $\mathcal{M}$  is  $U_{\mathcal{M}}(\hat{\gamma}) := \{x : \text{dist}(x, \mathcal{M}) < \hat{\gamma}\}$ . We say that  $\mathcal{M}$  is  $\hat{\gamma}$ -proximally smooth if the projection operator  $\mathcal{P}_{\mathcal{M}}(x)$  is a **singleton** whenever  $x \in U_{\mathcal{M}}(\hat{\gamma})$ 
  - Any compact smooth submanifold  $\mathcal{M}$  embedded in  $\mathbb{R}^{d \times k}$  is a proximally smooth set
  - Ensure not only the uniqueness of the projection but also the Lipschitz continuity of  $\mathcal{P}_{\mathcal{M}}$

$$\|\mathcal{P}_{\mathcal{M}}(x) - \mathcal{P}_{\mathcal{M}}(y)\| \leq 2\|x - y\|, \quad \forall x, y \in \bar{U}_{\mathcal{M}}(\hat{\gamma}/2)$$



# Convergence Analysis

## Assumptions

- ▶ The proximal smoothness constant of  $\mathcal{M}$  is  $2\gamma$
- ▶  $L$ -smoothness:  $\|\text{grad}f_{il}(x; \mathcal{D}_{il}) - \text{grad}f_{il}(y; \mathcal{D}_{il})\| \leq L\|x - y\|$
- ▶ Unbiasedness and bounded variance:  $\mathbb{E}[\|\text{grad}f_i(z_{i,t}^r; \mathcal{B}_{i,t}^r) - \text{grad}f_i(z_{i,t}^r)\|^2 \mid \mathcal{F}_t^r] \leq \sigma^2/b$

**(Theorem)** Under some assumptions and conditions on  $\tilde{\eta} := \eta\tau\eta_g$ , we have

$$\frac{1}{R} \sum_{r=1}^R \mathbb{E} \|\mathcal{G}_{\tilde{\eta}}(\mathcal{P}_{\mathcal{M}}(\bar{x}^r))\|^2 \leq \mathcal{O} \left( \frac{1}{\sqrt{n}\tau R\eta} + \frac{\sigma^2}{n\tau b} \right)$$

where  $\mathcal{G}_{\tilde{\eta}}(\mathcal{P}_{\mathcal{M}}(x^r)) := (\mathcal{P}_{\mathcal{M}}(x^r) - \tilde{x}^{r+1})/\tilde{\eta}$



# Numerical Experiments

► kPCA problem with Mnist dataset

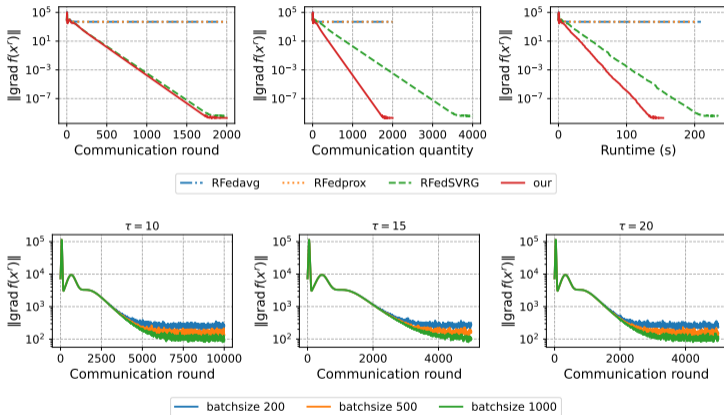


Figure: Comparison with alternative methods (1st row) and impacts of batch size (2nd row)