Federated Learning on Riemannian Manifolds with Differential Privacy

Speaker: Wen Huang

Xiamen University

Nov. 1 2024

Joint work with Zhenwei Huang, Pratik Jawanpuria, and Bamdev Mishra

大湾区大学优化理论与方法研讨会 东莞

Outline

- 1. Federated Learning
- 2. Private Federated Learning
- 3. Convergence Analysis
- 4. Numerical Experiments

5. Summary

Outline

1. Federated Learning

- 2. Private Federated Learning
- 3. Convergence Analysis
- 4. Numerical Experiments

5. Summary

Machine Learning (Empirical risk minimization):

$$\min_{x\in\mathbb{R}^n}\frac{1}{n}\sum_{i=1}^n f_i(x,z_i)$$

- Parameter x;
- Data $D = \{z_1, ..., z_n\};$

Machine Learning (Empirical risk minimization):

$$\min_{x\in\mathbb{R}^n}\frac{1}{n}\sum_{i=1}^n f_i(x,z_i)$$

- Parameter x;
- Data $D = \{z_1, ..., z_n\};$

Federated Learning:

- Privacy
- Computation ability

Applications:

- Healthcare
- Smart Home
- Financial Services
- Transportation



Federated Learning Optimization:

$$\min_{x \in \mathbb{R}^n} f(x) = \sum_{i=1}^N p_i f_i(x) = \sum_{i=1}^N p_i \left(\frac{1}{N_i} \sum_{j=1}^{N_i} f_i(x, z_{i,j}) \right), \text{ with } p_i = \frac{N_i}{n}, \qquad (1.1)$$

- *N* is the number of agents;
- $n = \sum_{i=1}^{N} N_i;$
- *f_i* is the local objective of agent *i*;

•
$$D_i = \{z_{i,1}, \ldots, z_{i,N_i}\}$$
 is the local data held by agent *i* with $|D_i| = N_i$.

Federated Learning Optimization:

$$\min_{x \in \mathcal{M}} f(x) = \sum_{i=1}^{N} p_i f_i(x) = \sum_{i=1}^{N} p_i \left(\frac{1}{N_i} \sum_{j=1}^{N_i} f_i(x, z_{i,j}) \right), \text{ with } p_i = \frac{N_i}{n}, \quad (1.1)$$

- *N* is the number of agents;
- $n = \sum_{i=1}^{N} N_i;$
- *f_i* is the local objective of agent *i*;
- $D_i = \{z_{i,1}, \ldots, z_{i,N_i}\}$ is the local data held by agent *i* with $|D_i| = N_i$.

Riemannian Federated Learning considers (1.1) with x in a manifold \mathcal{M}



Federated Learning Optimization:

$$\min_{x \in \mathcal{M}} f(x) = \sum_{i=1}^{N} p_i f_i(x) = \sum_{i=1}^{N} p_i \left(\frac{1}{N_i} \sum_{j=1}^{N_i} f_i(x, z_{i,j}) \right), \text{ with } p_i = \frac{N_i}{n}, \quad (1.1)$$

- *N* is the number of agents;
- $n = \sum_{i=1}^{N} N_i;$
- *f_i* is the local objective of agent *i*;
- $D_i = \{z_{i,1}, \ldots, z_{i,N_i}\}$ is the local data held by agent *i* with $|D_i| = N_i$.

Riemannian Federated Learning considers (1.1) with x in a manifold \mathcal{M}

Applications:

- Matrix completion
- Principal component analysis
- Online learning
- Taxonomy embedding
- etc



Euclidean version:

Algorithm: A representative federated averaging algorithm [McM+17]

- 1. for $t = 0, 1, \dots, T 1$ do
- The server uniformly selects a subset S_t of S agents at random; 2.
- The server upload global parameter $x^{(t)}$ to all agents in S_t , i.e., $x_i^{(t)} \leftarrow x^{(t)}$; 3.
- for $i \in S_t$ in parallel **do** 4.
- Agent *j* updates a local parameter $x_i^{(t+1)}$ by *K*-step SGD with x_t being 5. initial iterate: $\min_{x_i \in \mathbb{R}^n} \frac{1}{N_i} \sum_{i=1}^{N_i} f_i(x, z_{i,j})$
- Sent $x_i^{(t+1)}$ to the server; 6.
- 7. end for
- Server aggregates the received local parameters $\{x_i^{(t+1)}\}_{i \in S_t}$ by averaging 8.

$$x^{(t+1)} \leftarrow \sum_{j \in \mathcal{S}_t} \frac{N_j}{\sum_{j \in \mathcal{S}_t} N_j} x_j^{(t+1)};$$

9. end for

- Sever: Steps 2, 3, and 8;
- Agents: Steps 5 and 6;

[McM+17] B. McMahan, E. Moore, D. Ramage, B. A. v Arcas, Communication-Efficient Learning of Deep Networks from Decentralized Data, Proceedings of Machine Learning Research, 54, P.1273-1282, 2017.

Fuclidean version:

Algorithm: A representative federated averaging algorithm [McM+17]

- 1. for $t = 0, 1, \dots, T 1$ do
- The server uniformly selects a subset S_t of S agents at random; 2.
- The server upload global parameter $x^{(t)}$ to all agents in S_t , i.e., $x_i^{(t)} \leftarrow x^{(t)}$; 3.
- for $i \in S_t$ in parallel **do** 4.
- Agent *j* updates a local parameter $x_i^{(t+1)}$ by *K*-step SGD with x_t being 5. initial iterate: $\overline{\min_{x_i \in \mathbb{R}^n} \frac{1}{N_i} \sum_{i=1}^{N_i} f_i(x, z_{i,j})}$
- Sent $x_i^{(t+1)}$ to the server; 6.
- 7. end for
- Server aggregates the received local parameters $\{x_i^{(t+1)}\}_{i \in S_t}$ by averaging 8.

$$\mathbf{x}^{(t+1)} \leftarrow \sum_{j \in \mathcal{S}_t} \frac{N_j}{\sum_{j \in \mathcal{S}_t} N_j} \mathbf{x}_j^{(t+1)};$$

9. end for

- Sever: Steps 2, 3, and 8;
- Agents: Steps 5 and 6;

[McM+17] B. McMahan, E. Moore, D. Ramage, B. A. v Arcas, Communication-Efficient Learning of Deep Networks from Decentralized Data, Proceedings of Machine Learning Research, 54, P.1273-1282, 2017.

Euclidean to Riemannian:

Algorithm: A Riemannian federated learning algorithm [LM23]

1. for $t = 0, 1, \ldots, T - 1$ do

- The server uniformly selects a subset S_t of S agents at random; 2.
- The server upload global parameter $x^{(t)}$ to all agents in S_t , i.e., $x_i^{(t)} \leftarrow x^{(t)}$; 3.
- for $j \in S_t$ in parallel do 4.
- Agent *j* updates a local parameter $x_i^{(t+1)}$ by *K*-step Riemannian SGD with $x^{(t)}$ 5. being initial iterate; $\min_{x_i \in \mathcal{M}} \frac{1}{N_i} \sum_{i=1}^{N_i} f_i(x, z_{i,j})$
- Sent $x_i^{(t+1)}$ to the server; 6.
- 7. end for
- Server aggregates the received local parameters $\{x_i^{(t+1)}\}_{i \in S_t}$ by averaging 8.

$$x^{(t+1)} \leftarrow \operatorname{ave}(x_j^{(t+1)} \mid j \in \mathcal{S}_j);$$

- Agents: Riemannian SGD [Bon13]
- Sever: Aggregation

[[]LM23] J. Li and S. Ma. Federated learning on Riemannian manifolds. Applied Set-Valued Analysis and Optimization, 5(2), 2023. (Local training uses RSVRG)

Euclidean to Riemannian:

Algorithm: A Riemannian federated learning algorithm [LM23]

1. for $t = 0, 1, \ldots, T - 1$ do

- The server uniformly selects a subset S_t of S agents at random; 2.
- The server upload global parameter $x^{(t)}$ to all agents in S_t , i.e., $x_i^{(t)} \leftarrow x^{(t)}$; 3.
- for $j \in S_t$ in parallel do 4.
- Agent *j* updates a local parameter $x_i^{(t+1)}$ by *K*-step Riemannian SGD with $x^{(t)}$ 5. being initial iterate; $\min_{x_i \in \mathcal{M}} \frac{1}{N_i} \sum_{i=1}^{N_i} f_i(x, z_{i,j})$
- Sent $x_i^{(t+1)}$ to the server; 6.
- 7. end for
- Server aggregates the received local parameters $\{x_i^{(t+1)}\}_{i \in S_t}$ by averaging 8.

$$x^{(t+1)} \leftarrow \operatorname{ave}(x_j^{(t+1)} \mid j \in \mathcal{S}_j);$$

- Agents: Riemannian SGD [Bon13]
- Sever: Aggregation

[[]LM23] J. Li and S. Ma. Federated learning on Riemannian manifolds. Applied Set-Valued Analysis and Optimization, 5(2), 2023. (Local training uses RSVRG)

Euclidean to Riemannian:

Algorithm: A Riemannian federated learning algorithm [LM23]

1. for $t = 0, 1, \ldots, T - 1$ do

- 2. The server uniformly selects a subset S_t of S agents at random;
- The server upload global parameter $x^{(t)}$ to all agents in S_t , i.e., $x_i^{(t)} \leftarrow x^{(t)}$; 3.
- for $j \in S_t$ in parallel do 4.
- Agent *j* updates a local parameter $x_i^{(t+1)}$ by *K*-step Riemannian SGD with $x^{(t)}$ 5. being initial iterate; $\min_{x_i \in \mathcal{M}} \frac{1}{N_i} \sum_{i=1}^{N_i} f_i(x, z_{i,j})$
- Sent $x_i^{(t+1)}$ to the server; 6.
- 7. end for
- Server aggregates the received local parameters $\{x_i^{(t+1)}\}_{i \in S_t}$ by averaging 8.

$$x^{(t+1)} \leftarrow \operatorname{ave}(x_j^{(t+1)} \mid j \in \mathcal{S}_j);$$

- Agents: Riemannian SGD [Bon13]
- Sever: Aggregation

[[]LM23] J. Li and S. Ma. Federated learning on Riemannian manifolds. Applied Set-Valued Analysis and Optimization, 5(2), 2023. (Local training uses RSVRG)

Euclidean to Riemannian:

Algorithm: A Riemannian federated learning algorithm [LM23]

1. for $t = 0, 1, \ldots, T - 1$ do

- The server uniformly selects a subset S_t of S agents at random; 2.
- The server upload global parameter $x^{(t)}$ to all agents in S_t , i.e., $x_i^{(t)} \leftarrow x^{(t)}$; 3.
- for $j \in S_t$ in parallel do 4.
- Agent *j* updates a local parameter $x_i^{(t+1)}$ by *K*-step Riemannian SGD with $x^{(t)}$ 5. being initial iterate; $\min_{x_i \in \mathcal{M}} \frac{1}{N_i} \sum_{i=1}^{N_i} f_i(x, z_{i,j})$
- Sent $x_i^{(t+1)}$ to the server; 6.
- 7. end for
- Server aggregates the received local parameters $\{x_i^{(t+1)}\}_{i \in S_t}$ by averaging 8.

$$x^{(t+1)} \leftarrow \operatorname{ave}(x_j^{(t+1)} \mid j \in \mathcal{S}_j);$$

9. end for

- Agents: Riemannian SGD [Bon13]
- Sever: Aggregation

[LM23] J. Li and S. Ma. Federated learning on Riemannian manifolds. Applied Set-Valued Analysis and Optimization, 5(2), 2023. (Local training uses RSVRG)

Euclidean to Riemannian:

Algorithm: A Riemannian federated learning algorithm [LM23]

1. for $t = 0, 1, \ldots, T - 1$ do

- The server uniformly selects a subset S_t of S agents at random; 2.
- The server upload global parameter $x^{(t)}$ to all agents in S_t , i.e., $x_i^{(t)} \leftarrow x^{(t)}$; 3.
- for $j \in S_t$ in parallel do 4.
- Agent *j* updates a local parameter $x_i^{(t+1)}$ by *K*-step Riemannian SGD with $x^{(t)}$ 5. being initial iterate; $\min_{x_i \in \mathcal{M}} \frac{1}{N_i} \sum_{i=1}^{N_i} f_i(x, z_{i,j})$
- Sent $x_i^{(t+1)}$ to the server; 6.
- 7. end for
- Server aggregates the received local parameters $\{x_i^{(t+1)}\}_{i \in S_t}$ by averaging 8.

$$x^{(t+1)} \leftarrow \operatorname{ave}(x_j^{(t+1)} \mid j \in \mathcal{S}_j);$$

9. end for

- Agents: Riemannian SGD [Bon13]
- Sever: Aggregation

How to aggregates
$$\{x_i^{(t+1)}\}_{i \in S_t}$$
 on a manifold?

[LM23] J. Li and S. Ma. Federated learning on Riemannian manifolds. Applied Set-Valued Analysis and Optimization, 5(2), 2023. (Local training uses RSVRG)

Euclidean to Riemannian (Aggregation):

• Naive generalization:

$$x^{(t+1)} \leftarrow \sum_{j \in S_t} rac{N_j}{\sum_{j \in S_t} N_j} x_j^{(t+1)}
eq \mathsf{Riemannian setting}$$

Euclidean to Riemannian (Aggregation):

• Naive generalization:

$$x^{(t+1)} \leftarrow \sum_{j \in S_t} \frac{N_j}{\sum_{j \in S_t} N_j} x_j^{(t+1)}
ightarrow Riemannian setting$$

• An alternative approach:

$$\begin{aligned} \mathbf{x}^{(t+1)} \leftarrow \sum_{j \in \mathcal{S}_t} \frac{N_j}{\sum_{j \in \mathcal{S}_t} N_j} \mathbf{x}_j^{(t+1)} &\iff \mathbf{x}^{(t+1)} = \arg\min_{\mathbf{x}} \sum_{j \in \mathcal{S}_j} \frac{N_j}{\sum_{j \in \mathcal{S}_t} N_j} \|\mathbf{x} - \mathbf{x}_j^{(t+1)}\|_F^2 \\ &\iff \mathbf{x}^{(t+1)} = \arg\min_{\mathbf{x}} \sum_{j \in \mathcal{S}_j} \frac{N_j}{\sum_{j \in \mathcal{S}_t} N_j} \text{dist}^2(\mathbf{x}, \mathbf{x}_j^{(t+1)}) \implies \text{Riemannian setting}; \end{aligned}$$

Euclidean to Riemannian (Aggregation):

• Naive generalization:

$$x^{(t+1)} \leftarrow \sum_{j \in S_t} \frac{N_j}{\sum_{j \in S_t} N_j} x_j^{(t+1)}
ightarrow Riemannian setting$$

• An alternative approach:

$$\begin{aligned} x^{(t+1)} \leftarrow \sum_{j \in \mathcal{S}_t} \frac{N_j}{\sum_{j \in \mathcal{S}_t} N_j} x_j^{(t+1)} \iff x^{(t+1)} &= \arg\min_x \sum_{j \in \mathcal{S}_j} \frac{N_j}{\sum_{j \in \mathcal{S}_t} N_j} \|x - x_j^{(t+1)}\|_F^2 \\ \iff x^{(t+1)} &= \arg\min_x \sum_{j \in \mathcal{S}_j} \frac{N_j}{\sum_{j \in \mathcal{S}_t} N_j} \operatorname{dist}^2(x, x_j^{(t+1)}) \Longrightarrow \text{Riemannian setting}; \\ x^{(t+1)} &= \arg\min_x \sum_{j \in \mathcal{S}_j} \frac{N_j}{\sum_{j \in \mathcal{S}_t} N_j} \operatorname{dist}^2(x, x_j^{(t+1)}): \text{ computationally} \\ &= \operatorname{expensive}; \end{aligned}$$

Euclidean to Riemannian (Aggregation):

• Naive generalization:

$$x^{(t+1)} \leftarrow \sum_{j \in \mathcal{S}_t} \frac{N_j}{\sum_{j \in \mathcal{S}_t} N_j} x_j^{(t+1)}
ightarrow \mathsf{Riemannian setting}$$

• An alternative approach:

$$\begin{aligned} x^{(t+1)} \leftarrow \sum_{j \in \mathcal{S}_t} \frac{N_j}{\sum_{j \in \mathcal{S}_t} N_j} x_j^{(t+1)} &\iff x^{(t+1)} = \arg\min_x \sum_{j \in \mathcal{S}_j} \frac{N_j}{\sum_{j \in \mathcal{S}_t} N_j} \|x - x_j^{(t+1)}\|_F^2 \\ &\iff x^{(t+1)} = \arg\min_x \sum_{j \in \mathcal{S}_j} \frac{N_j}{\sum_{j \in \mathcal{S}_t} N_j} \operatorname{dist}^2(x, x_j^{(t+1)}) \implies \text{Riemannian setting}; \\ x^{(t+1)} = \arg\min_x \sum_{j \in \mathcal{S}_j} \frac{N_j}{\sum_{i \in \mathcal{S}_t} N_j} \operatorname{dist}^2(x, x_j^{(t+1)}): \text{ computationally} \end{aligned}$$

- expensive;
- One step of Riemannian gradient descent (called tangent mean) [LM23]:

$$\boldsymbol{x}^{(t+1)} \leftarrow \operatorname{Exp}_{\boldsymbol{x}^{(t)}}\left(\sum_{j \in \mathcal{S}_t} \frac{N_j}{\sum_{i \in \mathcal{S}_t} N_i} \operatorname{Exp}_{\boldsymbol{x}^{(t)}}^{-1}(\boldsymbol{x}_j^{(t+1)})\right);$$

[LM23] Jiaxiang Li and Shiqian Ma. Federated learning on Riemannian manifolds. Applied Set-Valued Analysis and Optimization, 5(2), 2023.

Existing Riemannian Federated Learning:

- The reviewed Riemannian federated learning [LM23]
- Riemannian Federated Learning on Compact Submanifolds with Heterogeneous Data [Zha+24]
 - Use projection onto the manifold
 - Allow multiple agents and multiple local updates
- Riemannian Federated Learning via Averaging Gradient Stream [Hua+24]
 - Send tangent vectors rather than the local parameters
 - Allow abstract manifold
 - Allow multiple agents and multiple local updates

[[]LM23] J. Li and S. Ma. Federated Learning on Riemannian Manifolds.

[[]Hua+24] Z. Huang, W. Huang, P. Jawanpuria, B. Mishra. Riemannian Federated Learning via Averaging Gradient Stream. Applied Set-Valued Analysis and Optimization, 2023.

[[]Zha+24] J. Zhang and J. Hu and A. M.-C. So and M. Johansson. Nonconvex Federated Learning on Compact Smooth Submanifolds With Heterogeneous Data. arxiv:2406.08465, 2024.

Existing Riemannian Federated Learning:

- The reviewed Riemannian federated learning [LM23]
- Riemannian Federated Learning on Compact Submanifolds with Heterogeneous Data [Zha+24]
 - Use projection onto the manifold
 - Allow multiple agents and multiple local updates
- Riemannian Federated Learning via Averaging Gradient Stream [Hua+24]
 - Send tangent vectors rather than the local parameters
 - Allow abstract manifold
 - Allow multiple agents and multiple local updates

Propose Riemannian federated learning with differential privacy.

[[]LM23] J. Li and S. Ma. Federated Learning on Riemannian Manifolds.

[[]Hua+24] Z. Huang, W. Huang, P. Jawanpuria, B. Mishra. Riemannian Federated Learning via Averaging Gradient Stream. Applied Set-Valued Analysis and Optimization, 2023.

[[]Zha+24] J. Zhang and J. Hu and A. M.-C. So and M. Johansson. Nonconvex Federated Learning on Compact Smooth Submanifolds With Heterogeneous Data. arxiv:2406.08465, 2024.

Existing Riemannian Federated Learning:

- The reviewed Riemannian federated learning [LM23]
- Riemannian Federated Learning on Compact Submanifolds with Heterogeneous Data [Zha+24]
 - Use projection onto the manifold
 - Allow multiple agents and multiple local updates
- Riemannian Federated Learning via Averaging Gradient Stream [Hua+24]
 - Send tangent vectors rather than the local parameters
 - Allow abstract manifold
 - Allow multiple agents and multiple local updates

Propose Riemannian federated learning with differential privacy.

The later two works appear after the work in this talk. No differential privacy is considered.

[[]LM23] J. Li and S. Ma. Federated Learning on Riemannian Manifolds.

[[]Hua+24] Z. Huang, W. Huang, P. Jawanpuria, B. Mishra. Riemannian Federated Learning via Averaging Gradient Stream. Applied Set-Valued Analysis and Optimization, 2023.

[[]Zha+24] J. Zhang and J. Hu and A. M.-C. So and M. Johansson. Nonconvex Federated Learning on Compact Smooth Submanifolds With Heterogeneous Data. arxiv:2406.08465, 2024.

- 1. Federated Learning
- 2. Private Federated Learning
- 3. Convergence Analysis
- 4. Numerical Experiments
- 5. Summary

- (Advantage) Data is stored in each agent.
- (Downside) It is possibly attacked by inference and thus agents' information is leaked.

Encryption method:

- k-anonymity
- Secure multiparty computing
- Homomorphic encryption
- Differential privacy

- (Advantage) Data is stored in each agent.
- (Downside) It is possibly attacked by inference and thus agents' information is leaked.

Encryption method:

- k-anonymity
- Secure multiparty computing
- Homomorphic encryption
- Differential privacy

Differential Privacy (DP): Differential Privacy offers a rigorous framework for addressing data privacy by precisely quantifying the deviation in the output distribution of a mechanism when a small number of input datasets are modified.

Differential Privacy (DP): Differential Privacy offers a rigorous framework for addressing data privacy by precisely quantifying the deviation in the output distribution of a mechanism when a small number of input datasets are modified.

Definition 1

Let \mathcal{A} be a manifold-valued randomized mechanism, namely, $\mathcal{A} : \mathcal{Z}^n \to \mathcal{M}$. We say that \mathcal{A} is (ϵ, δ) -differential privacy if for any two adjacent datasets D, D', which differ in at most one data point, and for any measurable sets $\mathcal{S} \subseteq \mathcal{M}$, it holds that

 $\mathrm{P}\{\mathcal{A}(D)\subseteq \mathcal{S}\}\leq \exp(\epsilon)\mathrm{P}\{\mathcal{A}(D')\subseteq \mathcal{S}\}+\delta.$

Private Federated Learning

Algorithm: A Riemannian federated learning algorithm [LM23]

1. for $t = 0, 1, \dots, T - 1$ do

- 2. The server uniformly selects a subset S_t of S agents at random;
- The server upload global parameter $x^{(t)}$ to all agents in S_t , i.e., $x_i^{(t)} \leftarrow x^{(t)}$; 3.
- 4. for $i \in S_t$ in parallel do
- Agent *j* updates a local parameter $x_i^{(t+1)}$ by *K*-step Riemannian SGD with $x^{(t)}$ 5. being initial iterate; $\min_{x_i \in \mathcal{M}} \frac{1}{N_i} \sum_{i=1}^{N_i} f_i(x, z_{i,j})$
- Sent $x_i^{(t+1)}$ to the server; 6.

Server aggregates the received local parameters $\{x_i^{(t+1)}\}_{i \in S_t}$ by averaging 8.

$$\boldsymbol{x}^{(t+1)} \leftarrow \operatorname{Exp}_{\boldsymbol{x}^{(t)}}\left(\sum_{j \in \mathcal{S}_t} \frac{N_j}{\sum_{i \in \mathcal{S}_t} N_i} \operatorname{Exp}_{\boldsymbol{x}^{(t)}}^{-1}(\boldsymbol{x}_j^{(t)})\right);$$

- For *t*-th outer iteration on *j*-th agent: $x_i^{(t+1)}$ is an output of a mechanism on $D_i = \{z_{i,1}, \ldots, z_{i,N_i}\};$
- For outer iterations from 0 to t: x_{t+1} is an output of a mechanism on D = $\cup_{i=1}$ $_{N}D_{i} = \{z_{1}, z_{2}, \ldots, z_{n}\};$
- Differential privacy: unlikely to infer any data from $x_i^{(t+1)}$ or $x^{(t)}$;

Algorithm: Riemannian federated learning with differential privacy (PriRFed)

01. for $t = 0, 1, \ldots, T - 1$ do

- 02. Sample S_t from [N] without replacement of size $|S_t| = s_t$;
- 03. Broadcast the global parameter $x^{(t)}$ to selected agents;
- 04. while $i \in S_t$ in parallel do
- 05. $\tilde{x}_i^{(t+1)} \leftarrow \text{PrivateLocalTraining}(x^{(t)}, D_i, K, \sigma_i, \alpha_t);$
- 06. Send $\tilde{x}_i^{(t+1)}$ to the server;
- 07. end while
- Aggregate the received local parameters to produce a global parameter x^(t+1) using tangent mean;
- 09. end for
- 10. output option 1: $\tilde{x} = x^{(T)}$;
- 11. output option 2: \tilde{x} is uniformly sampled from $\{x^{(t)}\}_{t=1}^{T}$ at random;

• Assumption: Local training satisfies (ϵ, δ) -DP;

Algorithm: Riemannian federated learning with differential privacy (PriRFed)

01. for $t = 0, 1, \ldots, T - 1$ do

- 02. Sample S_t from [N] without replacement of size $|S_t| = s_t$;
- 03. Broadcast the global parameter $x^{(t)}$ to selected agents;
- 04. while $i \in S_t$ in parallel do
- 05. $\tilde{x}_i^{(t+1)} \leftarrow \text{PrivateLocalTraining}(x^{(t)}, D_i, K, \sigma_i, \alpha_t);$
- 06. Send $\tilde{x}_i^{(t+1)}$ to the server;
- 07. end while
- Aggregate the received local parameters to produce a global parameter x^(t+1) using tangent mean;
- 09. end for
- 10. output option 1: $\tilde{x} = x^{(T)}$;
- 11. output option 2: \tilde{x} is uniformly sampled from $\{x^{(t)}\}_{t=1}^{T}$ at random;
 - Assumption: Local training satisfies (ϵ, δ)-DP;
 - Then the output \tilde{x} satisfies $(\tilde{\epsilon}, \tilde{\delta})$ -DP;

Algorithm: Riemannian federated learning with differential privacy (PriRFed)

01. for $t = 0, 1, \ldots, T - 1$ do

- 02. Sample S_t from [N] without replacement of size $|S_t| = s_t$;
- 03. Broadcast the global parameter $x^{(t)}$ to selected agents;
- 04. while $i \in S_t$ in parallel do
- 05. $\tilde{x}_i^{(t+1)} \leftarrow \text{PrivateLocalTraining}(x^{(t)}, D_i, K, \sigma_i, \alpha_t);$
- 06. Send $\tilde{x}_i^{(t+1)}$ to the server;
- 07. end while
- Aggregate the received local parameters to produce a global parameter x^(t+1) using tangent mean;
- 09. end for
- 10. output option 1: $\tilde{x} = x^{(T)}$;
- 11. **output option 2:** \tilde{x} is uniformly sampled from $\{x^{(t)}\}_{t=1}^{T}$ at random;
 - Assumption: Local training satisfies (ε, δ)-DP;
 - Then the output *x̃* satisfies (*ϵ̃*, *δ̃*)-DP;

The derivation relies on the existing properties of differential privacy

Properties of differential privacy

Theorem 2 (Post-processing [DR+14, Proposition 2.1])

Suppose that $\mathcal{A} : \mathcal{Z}^n \to \mathcal{M}$ is a randomized algorithm that is (ϵ, δ) -DP. Let $P : \mathcal{M} \to \mathcal{M}'$ be an arbitrary mapping. Then $P \circ \mathcal{A} : \mathcal{Z}^n \to \mathcal{M}'$ is (ϵ, δ) -DP.

Theorem 3 (Sequential composition theorem [DR+14, Theorem 3.16])

Suppose that $\mathcal{A}_i : \mathcal{Z}^n \to \mathcal{M}_i$ is an (ϵ_i, δ_i) -DP algorithm for $i \in [k]$. Then if $\mathcal{A}_{[k]} : \mathcal{Z}^n \to \prod_{i=1}^k \mathcal{M}_i$ is defined to be $\mathcal{A}_{[k]} := (\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_k)$, then $\mathcal{A}_{[k]}$ is $(\sum_{i=1}^k \epsilon_i, \sum_{i=1}^k \delta_i)$ -DP.

Theorem 4 (Advanced composition theorem [DRV10; KOV17; Whi+23])

Suppose that A_i is an (ϵ, δ) -DP algorithm for $i \in [k]$. Then the adaptive composition $A_{[k]}$ of A_1, \ldots, A_k , i.e., $A_{[k]} = A_k \circ A_{k-1} \circ \cdots \circ A_1$, is $(\epsilon', \delta' + k\delta)$ -DP with $\epsilon' = \sqrt{2k \log(1/\delta')}\epsilon + k\epsilon(\exp(\epsilon) - 1)$ and any $\delta' > 0$.

Definition 5 (Subsample)

Given a dataset *D* of *n* points, namely $D = \{D_1, D_2, ..., D_n\}$ (where D_i may be a dataset), an operator S_{ρ}^{wor} , called subsample, randomly chooses a sample from the uniform distribution over all subsets of *D* of size *m* without replacement. The ratio $\rho := m/n$ is called the sampling rate of the operator S_{ρ}^{wor} . Formally, letting $\{i_1, i_2, ..., i_m\}$ be the sampled indices, S_{ρ}^{wor} is defined by

$$\begin{split} \mathcal{S}_{\rho}^{\text{wor}} : & \mathcal{Z}^{N_1} \times \mathcal{Z}^{N_2} \times \ldots \times \mathcal{Z}^{N_n} \to \mathcal{Z}^{N_{i_1}} \times \mathcal{Z}^{N_{i_2}} \times \ldots \times \mathcal{Z}^{N_{i_m}} : \\ & (D_1, D_2, \ldots, D_n) \mapsto (D_{i_1}, D_{i_2}, \ldots, D_{i_m}). \end{split}$$

Lemma 6 (Subsampling lemma)

Let $S_{\rho}^{\text{wor}} : \mathbb{Z}^{N_1} \times \mathbb{Z}^{N_2} \times \ldots \times \mathbb{Z}^{N_n} \to \mathbb{Z}^{N_{i_1}} \times \mathbb{Z}^{N_{i_2}} \times \cdots \times \mathbb{Z}^{N_{i_m}}$ be a subsample operator (defined by Definition 5) with sampling rate $\rho = m/n$ and sampled indices $\{i_1, i_2, \ldots, i_m\}$. Suppose that $\mathcal{A} : \mathbb{Z}^{N_{i_1}} \times \mathbb{Z}^{N_{i_2}} \times \cdots \times \mathbb{Z}^{N_{i_m}} \to \mathcal{M}$ is a mechanism obeying (ϵ, δ) -DP. Then $\mathcal{A}' := \mathcal{A} \circ S_{\rho}^{\text{wor}} : \mathbb{Z}^{N_1} \times \mathbb{Z}^{N_2} \times \ldots \times \mathbb{Z}^{N_n} \to \mathcal{M}$ provides (ϵ', δ') -DP, where $\epsilon' = \log(1 + \rho(\exp(\epsilon) - 1))$ and $\delta' \leq \rho \delta$.

Definition 7 (c-stable transformation [McS09, Definition 2])

Let $\overline{\mathcal{Z}} := \mathcal{Z}^{N_1} \times \mathcal{Z}^{N_2} \times \ldots \times \mathcal{Z}^{N_n}$. A transformation $T : \overline{\mathcal{Z}} \to \overline{\mathcal{Z}}$ is said *c*-stable if for any two data sets $D_1, D_2 \in \overline{\mathcal{Z}}$, the following holds

 $|T(D_1)\oplus T(D_2)|\leq c|D_1\oplus D_2|.$

where c is a constant.

Theorem 8 (Pre-processing)

Let $\overline{Z} := Z^{N_1} \times Z^{N_2} \times \ldots \times Z^{N_n}$. Suppose that $\mathcal{A} : \overline{Z} \to \mathcal{M}$ be (ϵ, δ) -DP, and $T : \overline{Z} \to \overline{Z}$ be an arbitrary c-stable transformation. Then the composition $\mathcal{A} \circ T$ provides $(c\epsilon, \operatorname{cexp}((c-1)\epsilon)\delta)$ -DP.

Theorem 9 (Privacy guarantee of PriRFed)

If the privately local training satisfies (ϵ, δ) -DP for agent $i \in [N]$, then PriRFed obeys (ϵ', δ') -DP with $\epsilon' = \min(T\tilde{\epsilon}, \sqrt{2T \ln(1/\delta)}\tilde{\epsilon} + T\tilde{\epsilon}(\exp(\tilde{\epsilon}) - 1))$ and $\delta' = \hat{\delta} + T\tilde{\delta}$ for any $\hat{\delta} \in (0, 1)$, where $\tilde{\epsilon} = \log(1 + \rho(\exp(s\epsilon) - 1))$, $\tilde{\delta} = \rho s \delta$, $S_t \subseteq [N]$ with $s_t = s \le N$, and $\rho = s/N$ is the sampling rate.

Theorem 9 (Privacy guarantee of PriRFed)

If the privately local training satisfies (ϵ, δ) -DP for agent $i \in [N]$, then PriRFed obeys (ϵ', δ') -DP with $\epsilon' = \min(T\tilde{\epsilon}, \sqrt{2T \ln(1/\delta)}\tilde{\epsilon} + T\tilde{\epsilon}(\exp(\tilde{\epsilon}) - 1))$ and $\delta' = \hat{\delta} + T\tilde{\delta}$ for any $\hat{\delta} \in (0, 1)$, where $\tilde{\epsilon} = \log(1 + \rho(\exp(s\epsilon) - 1))$, $\tilde{\delta} = \rho s \delta$, $S_t \subseteq [N]$ with $s_t = s \le N$, and $\rho = s/N$ is the sampling rate.

Differences to existing Riemannian algorithms with DP in [Han+24; Upt+23]:

- Not federated learning framework/no aggregation;
- Guarantee DP in any communication between agents and the server;

[HMJG2024] A. Han, B. Mishra, P. Jawanpuria, and J. Gao, Differentially private Riemannian optimization, Machine Learning, 113, pp. 1133–1161, 2024

[Upt+23] S. Uptake, A. Han, P. Jawanpuria, and B. Mishra, Improved differentially private Riemannian optimization: Fast sampling and variance reduction, Transactions on Machine Learning Research, 2023
How to make the local training algorithm satisfy (ϵ, δ) -DP?

How to make the local training algorithm satisfy (ϵ, δ) -DP? Add noise to the local training algorithm. How to make the local training algorithm satisfy (ϵ, δ) -DP? Add noise to the local training algorithm.

Gaussian Mechanism:

- In Euclidean setting:
 - Given a function H : Zⁿ → ℝ^m, the Gaussian mechanism associated with H is given by H(D) = H(D) + ε, where ε is a Gaussian noise.
 - Letting $\Delta H = \sup_{D \sim D'} ||H(D) H(D')||_2$ be the sensitivity of H, if $\sigma^2 \geq 2 \log(1.25/\delta) \Delta_H^2/\epsilon^2$, then \mathcal{H} obeys (ϵ, δ) -DP.

How to make the local training algorithm satisfy (ϵ, δ) -DP? Add noise to the local training algorithm.

Gaussian Mechanism:

- In Euclidean setting:
 - Given a function H : Zⁿ → ℝ^m, the Gaussian mechanism associated with H is given by H(D) = H(D) + ε, where ε is a Gaussian noise.
 - Letting $\Delta H = \sup_{D \sim D'} ||H(D) H(D')||_2$ be the sensitivity of H, if $\sigma^2 \geq 2 \log(1.25/\delta) \Delta_H^2/\epsilon^2$, then \mathcal{H} obeys (ϵ, δ) -DP.
- In Riemannian setting:
 - Given any *x* ∈ *M*, for a function *H* : Zⁿ → T_x*M*, the Gaussian mechanism associated with *f* is given by *H*(*D*) = *H*(*D*) + ε, where ε is a tangent Gaussian noise;
 - Letting $\Delta H = \sup_{D \sim D'} ||H(D) H(D')||_x$ be the sensitivity of H, if $\sigma^2 \geq 2 \log(1.25/\delta) \Delta_H^2/\epsilon^2$, then \mathcal{H} obeys (ϵ, δ) -DP.

Algorithm: A Riemannian SGD Algorithm with DP (DP-RSGD) [Han+24]

• Add tangent Gaussian noise $\varepsilon_{i,k}^{(t)}$ to the stochastic gradient direction;

• If
$$\sigma_i^2 = o_i \frac{K \log(1/\delta) \tau^2}{N_i^2 \epsilon^2}$$
, then DP-RSGD satisfies (ϵ, δ) -DP.

[[]Han+24] A. Han, B. Mishra, P. Jawanpuria, J, Gao. Differentially private Riemannian optimization. Machine Learning, 113, P.1133–1161, 2024.

Private Federated Learning

Algorithm: A Riemannian SVRG Algorithm with DP (DP-RSVRG) [Upt+23]

01. Set
$$\tilde{x}_{i,0}^{(t)} \leftarrow x^{(t)}$$

02. for $k = 0, 1, ..., K - 1$ do
03. $x_{k+1,0}^{(t)} \leftarrow \tilde{x}_{i,k}^{(t)}$ and $g_{k+1}^{(t)} \leftarrow \frac{1}{N_i} \sum_{j=1}^{N_i} \operatorname{clip}_{\tau}(\operatorname{grad} f_{i,j}(\tilde{x}_{i,k}^{(t)}));$
04. for $j = 0, 1, ..., m - 1$ do
05. Randomly pick $l_j \in [N_i]$ and $\varepsilon_{k+1,j}^{(t)} \sim \mathcal{N}_{x_{k+1,j}^{(t)}}(0, \sigma_i^2);$
06. $\eta_{k+1,j}^{(t)} \leftarrow \operatorname{clip}_{\tau}(\operatorname{grad} f_{i,l_j}(x_{k+1,j}^{(t)})) - \Gamma_{\tilde{x}_{i,k}^{(t)}}^{\tilde{x}_{k+1,j}^{(t)}}(\operatorname{clip}_{\tau}(\operatorname{grad} f_{i,l_j}(\tilde{x}_{i,k}^{(t)}) - g_{k+1}^{(t)})) + \varepsilon_{k+1,j}^{(t)};$
07. $x_{k+1,j+1}^{(t)} \leftarrow \operatorname{Exp}_{x_{k+1,j}^{(t)}}(-\alpha_t \eta_{k+1,j}^{(t)});$
08. end for
09. $\tilde{x}_{i,k+1}^{(t)} \leftarrow x_{k+1,m}^{(t)};$
10. end for
11. output option 1: $\tilde{x}_i^{(t+1)} \leftarrow x_{K,m}^{(t)};$
12. output option 2: $\tilde{x}_i^{(t+1)}$ is uniformly randomly selected from $\{\{x_{k+1,j}^{(t)}\}_{j=0}^{m-1}\}_{k=0}^{K-1};$

• Add tangent Gaussian noise $\varepsilon_{i,k}^{(t)}$ to the stochastic gradient direction;

• If
$$\sigma_i^2 = \tilde{o}_i \frac{mK \log(1/\delta)\tau^2}{N^2 \epsilon^2}$$
, then DP-RSVRG satisfies (ϵ, δ) -DP.

[Upt+23] A. Utpala, A. Han, P. Jawanpuria, B. Mishra. Improved Differentially Private Riemannian Optimization: Fast Sampling and Variance Reduction. Transactions on Machine Learning Research, ISSN:2835–8856, 2023.

- 1. Federated Learning
- 2. Private Federated Learning
- Convergence Analysis
 PriRFed with DP-RSGD
 PriRFed with DP-RSVRG
- 4. Numerical Experiments
- 5. Summary

Convergence Analysis

PriRFed with DP-RSGD

Assumptions:

- The function *f_{ij}* is geodesic *L_f*-Lipschitz continuous for any *i*, *j*;
- The function *f_{ij}* is geodesically *L_g*-smooth, i.e., Lipschitz continuous Riemannian gradient for any *i*, *j*;

Convergence Analysis

PriRFed with DP-RSGD

Assumptions:

- The function *f_{ij}* is geodesic *L_f*-Lipschitz continuous for any *i*, *j*;
- The function *f_{ij}* is geodesically *L_g*-smooth, i.e., Lipschitz continuous Riemannian gradient for any *i*, *j*;

Definition 10 (Geodesic Lipschitz continuity)

A function $f : \mathcal{M} \to \mathbb{R}$ is called geodesically L_f -Lipschitz continuous if for any $x, y \in \mathcal{M}$, there exists $L_f \ge 0$ such that $|f(x) - f(y)| \le L_f \operatorname{dist}(x, y)$.

Definition 11 (Geodesic smoothness [ZS16])

A differentiable function $f : \mathcal{M} \to \mathbb{R}$ is call geodesically L_g -smooth if its gradient is L_g -Lipschitz continuous, that is, $\|\operatorname{grad} f(x) - \Gamma_y^x(\operatorname{grad} f(y))\|_x \le L_g \operatorname{dist}(x, y)$.

Assumptions:

- The function *f_{ij}* is geodesic *L_f*-Lipschitz continuous for any *i*, *j*;
- The function *f_{ij}* is geodesically *L_g*-smooth, i.e., Lipschitz continuous Riemannian gradient for any *i*, *j*;

Theorem 10 (Nonconvex, K = 1, fixed step size)

Consider PriRFed-DP-RSGD with **Option 1**. Set $s_t = N$, K = 1, $b_i = N_i$, and $\alpha_t = \alpha \le 1/L_g$. It holds that $\min_{0 \le t \le T} \mathbb{E}[\|\operatorname{grad} f(x^{(t)})\|^2] \le \frac{2}{T\alpha}(f(x^{(0)}) - f(x^*)) + \frac{dL_g\alpha}{(\sum_{i=1}^N N_i)^2} \sum_{i=1}^N N_i^2 \sigma_i^2$, where the expectation is taken over $\varepsilon_{i,0}^{(t)}$, and $x^* = \arg\min_{x \in M} f(x)$.

Theorem 11 (Nonconvex, K = 1, decaying step sizes)

Same parameters except that the step sizes $\{\alpha_t\}_{t=0}^{T-1}$ satisfies $\lim_{T\to\infty} \sum_{t=0}^{T-1} \alpha_t = \infty, \lim_{T\to\infty} \sum_{t=0}^{T-1} \alpha_t^2 < \infty, \text{ and } L_g \alpha_t \le 2\delta,$ for a constant $\delta \in (0, 1)$. It holds that $\liminf_{T\to\infty} \mathbb{E}\left[\| \operatorname{grad} f(x^{(t)}) \|^2 \right] = 0.$

Assumptions:

- The function *f_{ij}* is geodesic *L_f*-Lipschitz continuous for any *i*, *j*;
- The function *f_{ij}* is geodesically *L_g*-smooth, i.e., Lipschitz continuous Riemannian gradient for any *i*, *j*;

Theorem 10 (Nonconvex, K = 1, fixed step size)

Consider PriRFed-DP-RSGD with **Option 1**. Set $s_t = N$, K = 1, $b_i = N_i$, and $\alpha_t = \alpha \le 1/L_g$. It holds that $\min_{0 \le t \le T} \mathbb{E}[\|\operatorname{grad} f(x^{(t)})\|^2] \le \frac{2}{T\alpha}(f(x^{(0)}) - f(x^*)) + \frac{dL_g\alpha}{(\sum_{i=1}^N N_i)^2} \sum_{i=1}^N N_i^2 \sigma_i^2$, where the expectation is taken over $\varepsilon_{i,0}^{(t)}$, and $x^* = \arg \min_{x \in \mathcal{M}} f(x)$.

Theorem 11 (Nonconvex, K = 1, decaying step sizes)

Same parameters except that the step sizes $\{\alpha_t\}_{t=0}^{T-1}$ satisfies $\lim_{T\to\infty} \sum_{t=0}^{T-1} \alpha_t = \infty, \lim_{T\to\infty} \sum_{t=0}^{T-1} \alpha_t^2 < \infty, \text{ and } L_g \alpha_t \le 2\delta,$ for a constant $\delta \in (0, 1)$. It holds that $\liminf_{T\to\infty} \mathbb{E} \left[\| \operatorname{grad} f(x^{(t)}) \|^2 \right] = 0.$

The noise does not influence the convergences for decaying step sizes.

Assumptions:

- The function *f_{ij}* is geodesic *L_f*-Lipschitz continuous for any *i*, *j*;
- The function *f_{ij}* is geodesically *L_g*-smooth, i.e., Lipschitz continuous Riemannian gradient for any *i*, *j*;

Theorem 10 (Nonconvex, K = 1, fixed step size)

Consider PriRFed-DP-RSGD with **Option 1**. Set $s_t = N$, K = 1, $b_i = N_i$, and $\alpha_t = \alpha \le 1/L_g$. It holds that $\min_{0 \le t \le T} \mathbb{E}[\|\operatorname{grad} f(x^{(t)})\|^2] \le \frac{2}{T\alpha}(f(x^{(0)}) - f(x^*)) + \frac{dL_g\alpha}{(\sum_{i=1}^N N_i)^2} \sum_{i=1}^N N_i^2 \sigma_i^2$, where the expectation is taken over $\varepsilon_{i,n}^{(t)}$, and $x^* = \arg\min_{x \in M} f(x)$.

Theorem 11 (Nonconvex, K = 1, decaying step sizes)

Same parameters except that the step sizes $\{\alpha_t\}_{t=0}^{T-1}$ satisfies $\lim_{T\to\infty} \sum_{t=0}^{T-1} \alpha_t = \infty, \lim_{T\to\infty} \sum_{t=0}^{T-1} \alpha_t^2 < \infty, \text{ and } L_g \alpha_t \le 2\delta,$ for a constant $\delta \in (0, 1)$. It holds that $\liminf_{T\to\infty} \mathbb{E}\left[\| \operatorname{grad} f(x^{(t)}) \|^2 \right] = 0.$

Decaying step sizes: not a generalization of any Riemannian FL result.

Assumptions:

- The function *f_{ij}* is geodesic *L_f*-Lipschitz continuous for any *i*, *j*;
- The function *f_{ij}* is geodesically *L_g*-smooth;
- Manifold is complete;
- All iterates stay in a compact subset $W \subset M$ and the sectional curvature in W is bounded in $[\kappa_{\min}, \kappa_{\max}]$;
- The function *f_i* is geodesically convex in *W* for any *i*;

Assumptions:

- The function *f_{ij}* is geodesic *L_f*-Lipschitz continuous for any *i*, *j*;
- The function *f_{ij}* is geodesically *L_g*-smooth;
- Manifold is complete;
- All iterates stay in a compact subset $W \subset M$ and the sectional curvature in W is bounded in $[\kappa_{\min}, \kappa_{\max}]$;
- The function *f_i* is geodesically convex in *W* for any *i*;

Definition 12 (Geodesic convexity [ZS16])

A function $\mathcal{M} \to \mathbb{R}$ is called geodesically convex on \mathcal{W} if for any $x, y \in \mathcal{W}$, a geodesic $\gamma : [0, 1] \to \mathcal{M}$ satisfying $\gamma(0) = x, \gamma(1) = y$, and $\gamma(t) \in \mathcal{W}, \forall t \in [0, 1]$, it holds that $f(\gamma(t)) \leq (1 - t)f(x) + tf(y)$.

Assumptions:

- The function *f_{ij}* is geodesic *L_f*-Lipschitz continuous for any *i*, *j*;
- The function *f_{ij}* is geodesically *L_g*-smooth;
- Manifold is complete;
- All iterates stay in a compact subset $W \subset M$ and the sectional curvature in W is bounded in $[\kappa_{\min}, \kappa_{\max}]$;
- The function *f_i* is geodesically convex in *W* for any *i*;

Theorem 12 (Convex, K = 1)

Consider PriRFed-DP-RSGD with **Option 1**. Set $s_t = N$, K = 1, $b_i = N_i$, and $\alpha_t = \alpha \le 1/(2L_g)$. It holds that $\mathbb{E}[f(x^{(T)}) - f(x^*)] \le \frac{\zeta \operatorname{dist}^2(x^{(0)}, x^*)}{2\alpha(\zeta + T - 1)} + \frac{T(6\zeta + T - 1)d}{16L_g(\zeta + T - 1)(\sum_{i=1}^N N_i)^2} \sum_{i=1}^N N_i^2 \sigma_i^2,$ where the expectation is taken over $\varepsilon_{i,0}^{(t)}$, $x^* = \arg \min_{x \in \mathcal{M}} f(x)$, and ζ depends on κ_{\min} .

- An optimal value of *T*;
- Consistent with [LM23, Theorem 9] if non-DP (different algorithm though);

Assumptions:

- The function *f_{ij}* is geodesic *L_f*-Lipschitz continuous for any *i*, *j*;
- The function *f_{ij}* is geodesically *L_g*-smooth;
- The function *f_i* is geodesically convex in *W* for any *i*;
- f satisfies the Riemannian Polyak-Lojasiewicz (RPL) condition with a positive constant μ;
- For any *i* and $x \in \mathcal{M}$, $\|\operatorname{grad} f_i(x) \operatorname{grad} f(x)\|^2 \le v_i$ and $\mathbb{E}[v_i] = v$.

Assumptions:

- The function *f_{ij}* is geodesic *L_f*-Lipschitz continuous for any *i*, *j*;
- The function *f_{ij}* is geodesically *L_g*-smooth;
- The function *f_i* is geodesically convex in *W* for any *i*;
- f satisfies the Riemannian Polyak-Lojasiewicz (RPL) condition with a positive constant μ;
- For any *i* and $x \in \mathcal{M}$, $\|\operatorname{grad} f_i(x) \operatorname{grad} f(x)\|^2 \le v_i$ and $\mathbb{E}[v_i] = v$.

Definition 13 (Riemannian Polyak-Lojasiewicz condition)

A function $\mathcal{M} \to \mathbb{R}$ satisfies the Riemannian Polyak-Lojasiewicz (RPL) condition with a positive constant μ if $f(x) - f(x^*) \le \frac{1}{2\mu} \| \operatorname{grad} f(x) \|^2$, where x^* is a global minimizer of f.

Assumptions:

- The function *f_{ij}* is geodesic *L_f*-Lipschitz continuous for any *i*, *j*;
- The function *f_{ij}* is geodesically *L_g*-smooth;
- The function *f_i* is geodesically convex in *W* for any *i*;
- f satisfies the Riemannian Polyak-Lojasiewicz (RPL) condition with a positive constant μ;
- For any *i* and $x \in \mathcal{M}$, $\|\operatorname{grad} f_i(x) \operatorname{grad} f(x)\|^2 \le v_i$ and $\mathbb{E}[v_i] = v$.

Theorem 13 (Convex and RPL, K = 1)

Consider PriRFed-DP-RSGD with **Option 1**. Set $N_i = N_j$ ($\forall i, j \in [N]$), $s_t = S \leq N, K = 1$, and $b_i = N_i$. Then it holds that $\mathbb{E}[f(x^{(T)})] - f(x^*) \leq P^T(f(x^{(0)}) - f(x^*)) + (1 - P^T) \left(\kappa_0 \frac{\sum_{i \in S} \sigma_i^2}{S^2} + \kappa_1 \frac{N-S}{S(N-1)}\right),$ where $P = (1 - 2\mu\alpha + \mu\alpha^2 L_g), \kappa_0 = \frac{L_g \alpha d}{2\mu(2 - L_g \alpha)},$ and $\kappa_1 = v\kappa_0 = \frac{L_g \alpha dv}{2\mu(2 - L_g \alpha)}.$

- An optimal value of T;
- The second term \neq 0 even if $\sigma_i = 0$ (non-DP);
- Consistent with the Euclidean FL in [Wei+22];

Assumptions:

- The function *f_{ij}* is geodesic *L_f*-Lipschitz continuous for any *i*, *j*;
- The function *f_{ij}* is geodesically *L_g*-smooth;
- Manifold is complete;
- All iterates stay in a compact subset $W \subset M$ and the sectional curvature in W is bounded in $[\kappa_{\min}, \kappa_{\max}]$;

Theorem 14 (Nonconvex, K > 1)

Consider PriRFed-DP-RSGD and **Option 2**. Set $s_t = 1, K > 1$, $\alpha_t = \alpha < (\beta^{\frac{1}{K-1}} - 1)/\beta$ with $\beta > 1$ being a free constant. It holds that

$$\mathbb{E}[\|\operatorname{grad} f(\tilde{x})\|^2] \leq \frac{f(x^{(0)}) - f(x^*)}{KT\delta_{\min}} + q_0 q,$$

where $\delta_{\min} = \alpha(1 - \frac{(1+\alpha\beta)^{K-1}}{\beta})$, $q_0 = \frac{\alpha\beta(L_g+2(1+\alpha\beta)^{K-1}\zeta)}{2(\beta-(1+\alpha\beta)^{K-1})}$, $q = L_f^2 + d\sigma^2$, and $\sigma = \max_{i=1,2,...,N} \sigma_i$.

Allow multiple inner iteration, i.e., K > 1, but only one agent is used, i.e., $s_t = 1$;

Assumptions:

- The function *f_{ij}* is geodesic *L_f*-Lipschitz continuous for any *i*, *j*;
- The function *f_{ij}* is geodesically *L_g*-smooth;
- Manifold is complete;
- All iterates stay in a compact subset $W \subset M$ and the sectional curvature in W is bounded in $[\kappa_{\min}, \kappa_{\max}]$;

Theorem 15 (Nonconvex, K > 1, m > 1)

Consider PriRFed-DP-RSVRG with **Option 2**. Set $s_t = 1, K > 1,$ $m = \lfloor 10N/(3\zeta^{1/2}) \rfloor > 1$, and $\alpha_t = \alpha \le 1/(10L_g N^{2/3}\zeta^{1/2})$. It holds that $\mathbb{E}[\|\operatorname{grad} f(\tilde{x})\|^2] \le c \frac{L_g \zeta^{1/2}}{N^{1/3}KT} [f(x^{(0)}) - f(x^*)] + \frac{d\sigma^2}{100\zeta^{1/2}} \left(\frac{1}{N^{2/3}} + \frac{1}{N}\right),$ where c > 10N/(3m) is constant, $\sigma = \max_{i=1,2,...,N} \sigma_i$, and ζ depends on κ_{\min} . Difficulty in convergence analysis: Control the noise propagation on manifolds

- Existing result: FL-DP on linear spaces, FL on manifolds, or DP on manifolds;
- Consider the distortion caused by linearity of manifolds, involving curvature;

Difficulty in convergence analysis: Control the noise propagation on manifolds

- Existing result: FL-DP on linear spaces, FL on manifolds, or DP on manifolds;
- Consider the distortion caused by linearity of manifolds, involving curvature;

A summary

- Consider "nonconvex", "convex", "RPL" scenarios;
- Some consistent with existing Riemannian FL in [LM23];
- Some consistent with Euclidean FL-DP in [Wei+22];
- Not guarantee convergence for all parameter settings;
- Importantly, $s_t > 1$ and K > 1 not considered;

- 1. Federated Learning
- 2. Private Federated Learning
- 3. Convergence Analysis
- 4. Numerical Experiments
- 4.1 Principal eigenvector computation
- 4.2 Hyperbolic structured prediction

5. Summary

Principal Eigenvector over Sphere Manifold

$$\min_{x \in S^{d}} f(x) := -\sum_{i=1}^{N} p_{i} \left(\frac{1}{N_{i}} \sum_{j=1}^{N_{i}} x^{T}(z_{i,j} z_{i,j}^{T}) x \right),$$
(4.1)
where $S^{d} = \{ x \in \mathbb{R}^{d+1} : \|x\|_{2} = 1 \}, p_{i} = \frac{N_{i}}{\sum_{i=1}^{N} N_{i}}, \text{ and } D_{i} = \{ z_{i,1}, \dots, z_{i,N_{i}} \}$
with $z_{i,j} \in \mathbb{R}^{d+1}$.

Principal Eigenvector over Sphere Manifold

$$\min_{x\in S^d} f(x) := -\sum_{i=1}^N p_i\left(\frac{1}{N_i}\sum_{j=1}^{N_i} x^T(z_{i,j}z_{i,j}^T)x\right)$$

Synthetic data:

- construct a $(d + 1) \times (d + 1)$ diagonal matrix $\Sigma_i = \text{diag}\{1, 1 - 1.1\nu, \dots, 1 - 1.4\nu, |y_1|/(d + 1), |y_2|/(d + 1), \dots\}$ where ν is the eigengap and y_1, y_2, \dots are drawn from the standard Gaussian distribution;
- construct $Z_i = U_i \Sigma_i V_i$ where U_i , V_i of size $N_i \times (d + 1)$ and $(d + 1) \times (d + 1)$ are random column orthonormal matrices

MNIST:

• 60,000 hand-written images of size 28×28 (in the case, $d + 1 = 28^2 = 784$);

Synthetic data: An average result of 10 random runs (d = 24)



- The smaller ε is, the larger the σ_i is;
- Left column: $s_t = N$, K = 1, $b_i = N_i$ and $\alpha_t = 1/(2L_g)$;
- Middle column: $s_t = 1$, K = 5, $b_i = N_i/2$, and $\alpha_t = 1.0$;
- Right column: $s_t = 1$, K = 2, $m = \lfloor 10N/3 \rfloor$, and $\alpha_t = 1/(10N^{2/3}L_g)$;

Federated Learning on Riemannian Manifolds with Differential Privacy

Synthetic data: An average result of 10 random runs (d = 24)



Large noise slows down the convergence speed and increases the noise floor.

Principal Eigenvector over Sphere Manifold

MNIST: An average result of 10 random runs (d = 783)



ε = 0.15, δ = 10⁻⁴, s_t = 1, K = 3, b_i = N_i/2, and α_t = 0.1;
DP-SVRG: m = |10N/3|;

$$h(x) := \underset{x \in \mathcal{H}^{d}}{\arg\min} f(x) = \sum_{i=1}^{N} p_{i} \sum_{j=1}^{N_{i}} \frac{1}{N_{i}} \alpha_{i,j}(w) \text{dist}^{2}(x, y_{i,j}), \quad (4.2)$$

where features $w = \{w_{i,j}\}$ with $w_{i,j} \in \mathbb{R}^r$ and hyperbolic embeddings $y_{i,j} \in \mathcal{H}^d$ are given,

- $\alpha_i(w) = (\alpha_{i,1}(w), \ldots, \alpha_{i,N_i}(w))^T \in \mathbb{R}^{N_i};$
- $\alpha_i(w) = (K_i + \gamma I)^{-1} K_{i,w}$ with a given $\gamma > 0$;

•
$$K_i \in \mathbb{R}^{N_i \times N_i}$$
 and $(K_i)_{i,j} = k(w_{i,i}, w_{i,j});$

•
$$K_{i,w} \in \mathbb{R}^{N_i}$$
 and $(K_{i,w})_j = k(w_{i,j}, w);$

• The kernel function $k(w, w') = \exp(-\|w - w'\|_2^2/(2\bar{v})^2)$.

$$h(x) := \underset{x \in \mathcal{H}^{d}}{\arg\min} f(x) = \sum_{i=1}^{N} p_{i} \sum_{j=1}^{N_{i}} \frac{1}{N_{i}} \alpha_{i,j}(w) \text{dist}^{2}(x, y_{i,j}), \quad (4.2)$$

where features $w = \{w_{i,j}\}$ with $w_{i,j} \in \mathbb{R}^r$ and hyperbolic embeddings $y_{i,j} \in \mathcal{H}^d$ are given,

- $\alpha_i(w) = (\alpha_{i,1}(w), \ldots, \alpha_{i,N_i}(w))^T \in \mathbb{R}^{N_i};$
- $\alpha_i(w) = (K_i + \gamma I)^{-1} K_{i,w}$ with a given $\gamma > 0$;

•
$$K_i \in \mathbb{R}^{N_i \times N_i}$$
 and $(K_i)_{i,j} = k(w_{i,i}, w_{i,j});$

•
$$K_{i,w} \in \mathbb{R}^{N_i}$$
 and $(K_{i,w})_j = k(w_{i,j}, w);$

- The kernel function $k(w, w') = \exp(-\|w w'\|_2^2/(2\bar{v})^2)$.
- Hyperbolic manifold $\mathcal{H}^d := \{x \in \mathbb{R}^{d+1} : \langle x, x \rangle_{\mathcal{L}} = -1\}$ with $\langle x, y \rangle_{\mathcal{L}} = x^T y 2x_1 y_1$ (Riemannian metric);
- The exponential map: $\operatorname{Exp}_{X}(v) = \operatorname{cosh}(\|v\|_{\mathcal{L}})X + v \frac{\sinh(\|v\|_{\mathcal{L}})}{\|v\|_{\mathcal{L}}};$
- The logarithm map $\operatorname{Exp}_{X}^{-1}(y) = \frac{\cosh^{-1}(-\langle x, y \rangle_{\mathcal{L}})}{\sinh(\cosh^{-1}(-\langle x, y \rangle_{\mathcal{L}}))}(y + \langle x, y \rangle_{\mathcal{L}} x);$
- The distance function $dist(x, y) = \cosh^{-1}(-\langle x, y \rangle_{\mathcal{L}}).$

Mammals Substree of WordNet

- Mammals Substree are embedded on \mathcal{H}^2 ;
- Transitive closure containing n = 1180 nodes and 6540 edges;
- The feature are stemmed from Laplacian eigenmap to dimension r = 3;
- The word "primate" as the test sample;



Plot in Poincaré ball model



- Introduced the federated learning;
- Introduced the differential privacy;
- Proposed a Riemannian federated learning framework with differential privacy guaranteed (PriRFed);
- Convergence analysis for two instances of PriRFed, i.e., with DP-RSGD and DP-RSVRG;
- Numerical experiments verify the performance;

For more details, see

Zhenwei Huang, Wen Huang, Pratik Jawanpuria, and Bamdev Mishra. Federated learning on Riemannian manifolds with differential privacy. arXiv:2404.10029, 2024.

Thank you for your attention!

References I

S. Bonnabel. "Stochastic Gradient Descent on Riemannian Manifolds". In: *IEEE Transactions on Automatic Control* 58.9 (2013), pp. 2217–2229. DOI: 10.1109/TAC.2013.2254619.

Cynthia Dwork, Aaron Roth, et al. "The algorithmic foundations of differential privacy". In: *Foundations and Trends*® *in Theoretical Computer Science* 9.3–4 (2014), pp. 211–407.

Cynthia Dwork, Guy N Rothblum, and Salil Vadhan. "Boosting and differential privacy". In: 2010 IEEE 51st Annual Symposium on Foundations of Computer Science. IEEE. 2010, pp. 51–60.

Andi Han et al. "Differentially private Riemannian optimization". In: *Machine Learning* 113.3 (2024), pp. 1133–1161.

Zhenwei Huang et al. *Riemannian Federated Learning via Averaging Gradient Stream*. 2024. arXiv: 2409.07223 [cs.LG]. URL: https://arxiv.org/abs/2409.07223.

Peter Kairouz, Sewoong Oh, and Pramod Viswanath. "The Composition Theorem for Differential Privacy". In: *IEEE Transactions on Information Theory* 63.6 (2017), pp. 4037–4049.

References II

Jiaxiang Li and Shiqian Ma. "Federated Learning on Riemannian Manifolds". In: *Applied Set-Valued Analysis and Optimization* 5.2 (2023).

Brendan McMahan et al. "Communication-Efficient Learning of Deep Networks from Decentralized Data". In: *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*. Ed. by Aarti Singh and Jerry Zhu. Vol. 54. Proceedings of Machine Learning Research. PMLR, 20–22 Apr 2017, pp. 1273–1282. URL: https: //proceedings.mlr.press/v54/mcmahan17a.html.

//proceedings.mlr.press/v54/mcmahanl/a.htm

Frank D McSherry. "Privacy integrated queries: an extensible platform for privacy-preserving data analysis". In: *Proceedings of the 2009 ACM SIGMOD International Conference on Management of data.* 2009, pp. 19–30.

Saiteja Uptake et al. "Improved Differentially Private Riemannian Optimization: Fast Sampling and Variance Reduction". In: *Transactions on Machine Learning Research* (2023). ISSN: 2835-8856.
References III

Kang Wei et al. "User-Level Privacy-Preserving Federated Learning: Analysis and Performance Optimization". In: *IEEE Transactions on Mobile Computing* 21.9 (2022), pp. 3388–3401.

Justin Whitehouse et al. "Fully-Adaptive Composition in Differential Privacy". In: *Proceedings of the 40th International Conference on Machine Learning*. Ed. by Andreas Krause et al. Vol. 202. Proceedings of Machine Learning Research. PMLR, 23–29 Jul 2023, pp. 36990–37007.

Jiaojiao Zhang et al. Nonconvex Federated Learning on Compact Smooth Submanifolds With Heterogeneous Data. 2024. arXiv: 2406.08465 [cs.LG]. URL: https://arxiv.org/abs/2406.08465.

Hongyi Zhang and Suvrit Sra. "First-order Methods for Geodesically Convex Optimization". In: 29th Annual Conference on Learning Theory. Vol. 49. Proceedings of Machine Learning Research. Columbia University, New York, New York, USA: PMLR, 23–26 Jun 2016, pp. 1617–1638.