

Problem-Parameter-Free Federated Learning

Wenjing Yan
The Chinese University of Hong Kong

(Joint work with Kai Zhang, Xiaolu Wang, Xuanyu Cao)

Outline

1 Background
and
Motivation



2. Algorithm
Design



3. Simulation
Results

Outline

1 Background
and
Motivation



2. Algorithm
Design



3. Simulation
Results

Problem Formulation of Federated Learning (FL)

• Federated Learning:

$$\min_{\theta \in \mathbb{R}^d} f(\theta) := \frac{1}{N} \sum_{i=1}^N f_i(\theta)$$

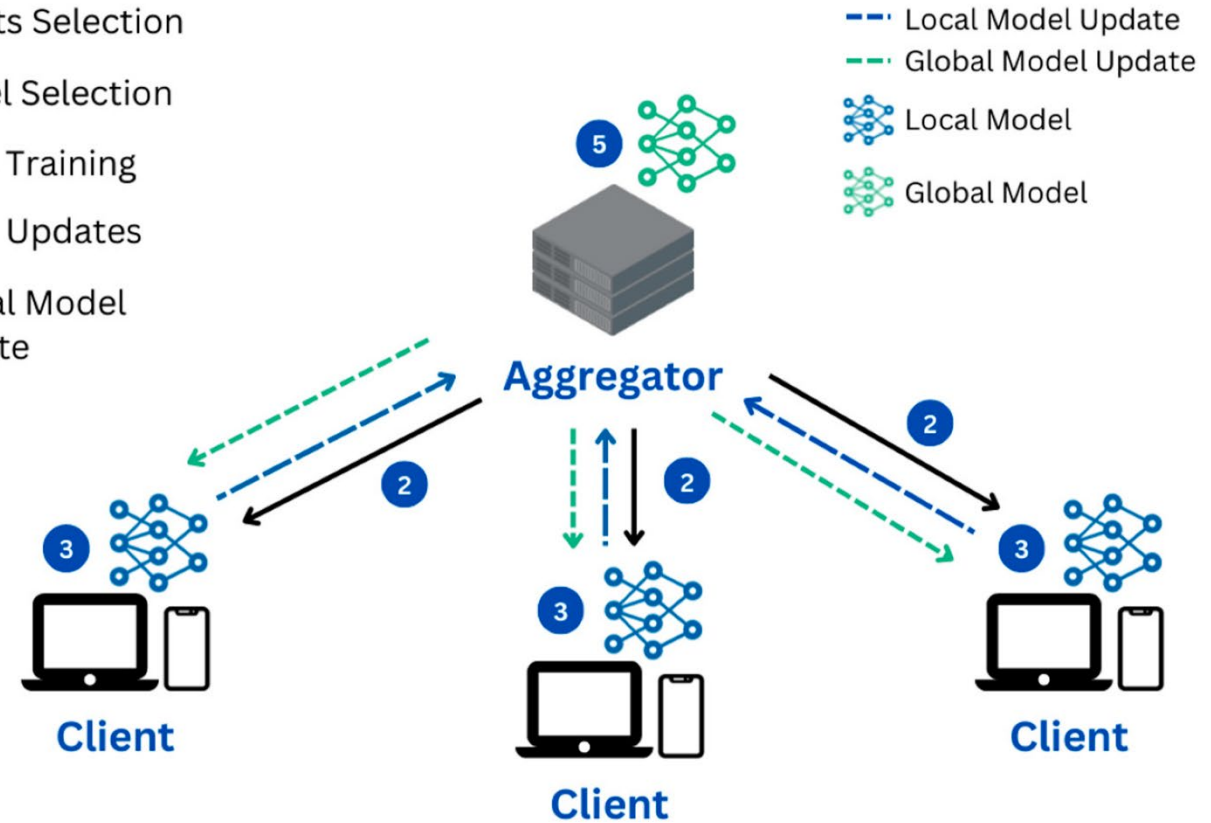
where $f_i(\theta) := \mathbb{E}_{\xi_i \sim \mathcal{D}_i} [F(\theta; \xi_i)]$

- Multiple clients collaboratively train a machine learning model with the help of a central server.
- Each client performs multiple local update based on private data
- Server aggregates the global model

Advantages:

- Ensures privacy by avoiding raw data sharing
- Offers scalability and communication efficiency

- 1 Clients Selection
- 2 Model Selection
- 3 Local Training
- 4 Local Updates
- 5 Global Model Update



Problem Formulation of Federated Learning (FL)

• Federated Learning:

$$\min_{\theta \in \mathbb{R}^d} f(\theta) := \frac{1}{N} \sum_{i=1}^N f_i(\theta)$$

where $f_i(\theta) := \mathbb{E}_{\xi_i \sim \mathcal{D}_i} [F(\theta; \xi_i)]$

FedAvg Key Steps:

- K steps local update at client i :

$$g_i^{t,k} = \nabla F(\theta_i^{t,k}; \xi_i^{t,k})$$

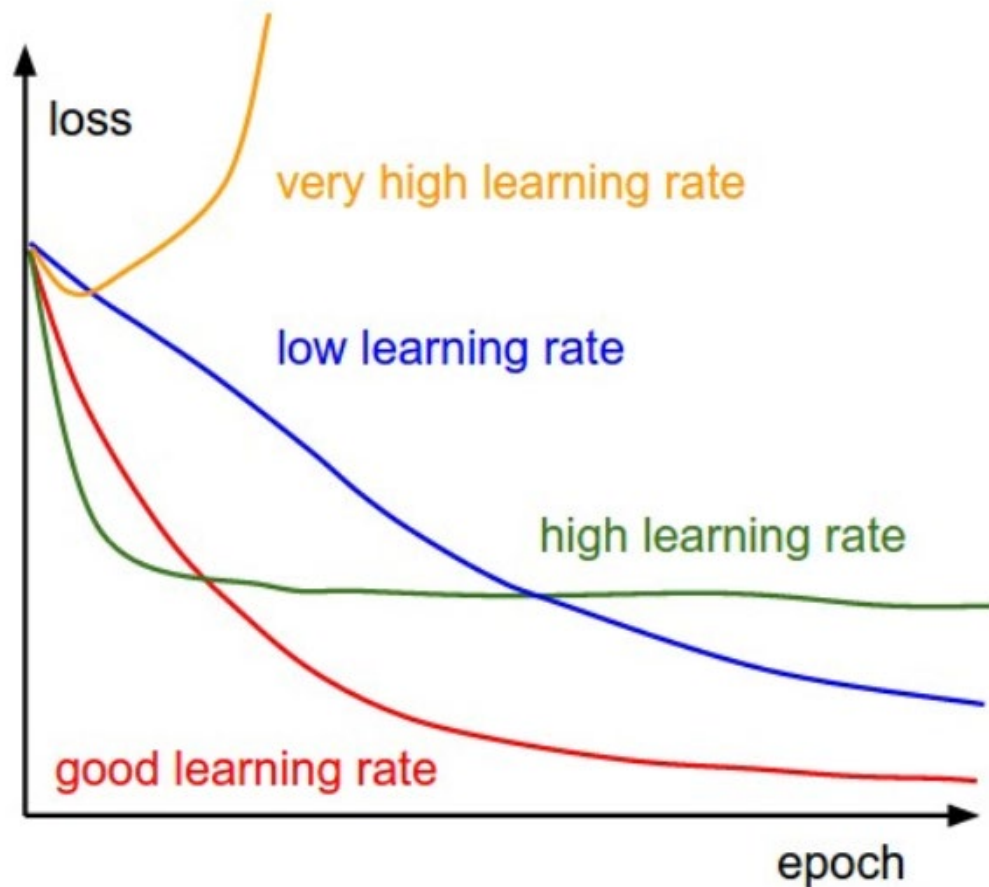
$$\theta_i^{t,k+1} = \theta_i^{t,k} - \eta_t g_i^{t,k}$$

- Global model aggregation at server:

$$g^t = \frac{1}{NK} \sum_{i=1}^N (\theta^t - \theta_i^{t,K})$$

$$\theta^{t+1} = \theta^t - \gamma g^t$$

Stepsize setting is crucial



Challenge 1: Problem-Specific Hyperparameter Tuning

- Federated Learning:

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} f(\boldsymbol{\theta}) := \frac{1}{N} \sum_{i=1}^N f_i(\boldsymbol{\theta})$$

where $f_i(\boldsymbol{\theta}) := \mathbb{E}_{\boldsymbol{\xi}_i \sim \mathcal{D}_i} [F(\boldsymbol{\theta}; \boldsymbol{\xi}_i)]$



- Problem-Specific Constants:

- L : Smoothness constant
- σ^2 : Stochastic gradient variance
- L, σ_h^2 : coefficients on gradient dissimilarity bound
- $\Delta := f(\boldsymbol{\theta}_0) - f^*$: Initial suboptimality gap

Assumptions

 L -Smoothness

$$\|\nabla F(\boldsymbol{\theta}; \boldsymbol{\xi}_i) - \nabla F(\boldsymbol{\delta}; \boldsymbol{\xi}_i)\| \leq L \|\boldsymbol{\theta} - \boldsymbol{\delta}\|$$

Stochastic Gradient Variance

$$\mathbb{E}_{\boldsymbol{\xi}_i} \|\nabla F(\boldsymbol{\theta}; \boldsymbol{\xi}_i) - \nabla f_i(\boldsymbol{\theta})\|^2 \leq \sigma^2$$

Bounded Gradient Dissimilarity

$$\frac{1}{N} \sum_{i=1}^N \|\nabla f_i(\boldsymbol{\theta})\|^2 \leq B \|\nabla f(\boldsymbol{\theta})\|^2 + \sigma_h^2$$

Challenge 1: Problem-Specific Hyperparameter Tuning

Federated Learning:

Assumptions

- Algorithms require careful tuning of learning rates, momentum, and other coefficients
- Hyperparameter Tuning depends on problem-specific factors (e.g., smoothness constants L , gradient variances σ^2 , etc.)
- Estimating those constants is difficult in FL due to privacy and data constraints.
- Problem-specific tuning limits applicability in dynamic environments (e.g., IoT, edge devices).



- L : Smoothness
- σ^2 : Stochastic gradient variance
- L, σ_h^2 : coefficients on gradient dissimilarity bound
- $\Delta := f(\theta_0) - f^*$: Initial suboptimality gap

$$\frac{1}{N} \sum_{i=1}^N \|\nabla f_i(\theta)\|^2 \leq B \|\nabla f(\theta)\|^2 + \sigma_h^2$$

Challenge 2: Requirement on Data Heterogeneity Bounds

- Federated Learning:

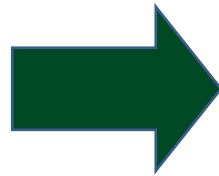
$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} f(\boldsymbol{\theta}) := \frac{1}{N} \sum_{i=1}^N f_i(\boldsymbol{\theta})$$

where $f_i(\boldsymbol{\theta}) := \mathbb{E}_{\boldsymbol{\xi}_i \sim \mathcal{D}_i} [F(\boldsymbol{\theta}; \boldsymbol{\xi}_i)]$



- Impact of Data Heterogeneity:

- Clients often have non-IID
- Leads to inconsistent local updates and “client drift”



Assumptions

L-Smoothness

$$\|\nabla F(\boldsymbol{\theta}; \boldsymbol{\xi}_i) - \nabla F(\boldsymbol{\delta}; \boldsymbol{\xi}_i)\| \leq L \|\boldsymbol{\theta} - \boldsymbol{\delta}\|$$

Stochastic Gradient Variance

$$\mathbb{E}_{\boldsymbol{\xi}_i} \|\nabla F(\boldsymbol{\theta}; \boldsymbol{\xi}_i) - \nabla f_i(\boldsymbol{\theta})\|^2 \leq \sigma^2$$

Bounded Gradient Dissimilarity

$$\frac{1}{N} \sum_{i=1}^N \|\nabla f_i(\boldsymbol{\theta})\|^2 \leq B \|\nabla f(\boldsymbol{\theta})\|^2 + \sigma_h^2$$

Requirement on Data Heterogeneity Bounds

- Federated Learning:

$$\min_{\theta \in \mathbb{R}^d} f(\theta) := \frac{1}{N} \sum_{i=1}^N f_i(\theta)$$

where $f_i(\theta) := \mathbb{E}_{\xi_i \sim \mathcal{D}_i} [F(\theta; \xi_i)]$

Assumptions

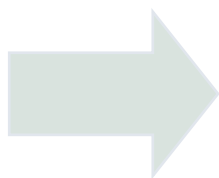
L-Smoothness

$$\|\nabla F(\theta; \xi) - \nabla F(\delta; \xi)\| \leq L \|\theta - \delta\|$$

- Quantifying gradient dissimilarity bound is difficult in FL due to privacy and data constraints.
- Data Heterogeneity Bounds limit the applicability of FL dynamic environments with varying data distributions.

- Impact of Data Heterogeneity:

- Clients often have non-IID
- Leads to inconsistent local updates and “client drift”



Bounded Gradient Dissimilarity

$$\frac{1}{N} \sum_{i=1}^N \|\nabla f_i(\theta)\|^2 \leq B \|\nabla f(\theta)\|^2 + \sigma_h^2$$

Our Method: Problem-Parameter-Free Federated Learning

- Federated Learning:

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} f(\boldsymbol{\theta}) := \frac{1}{N} \sum_{i=1}^N f_i(\boldsymbol{\theta})$$

where $f_i(\boldsymbol{\theta}) := \mathbb{E}_{\boldsymbol{\xi}_i \sim \mathcal{D}_i} [F(\boldsymbol{\theta}; \boldsymbol{\xi}_i)]$



Independent of all problem-specific parameters, enabling **tuning-free**

Eliminating the requirement on data heterogeneity bounds

Assumptions

 L -Smoothness

$$\|\nabla F(\boldsymbol{\theta}; \boldsymbol{\xi}_i) - \nabla F(\boldsymbol{\delta}; \boldsymbol{\xi}_i)\| \leq L \|\boldsymbol{\theta} - \boldsymbol{\delta}\|$$

Stochastic Gradient Variance

$$\mathbb{E}_{\boldsymbol{\xi}_i} \|\nabla F(\boldsymbol{\theta}; \boldsymbol{\xi}_i) - \nabla f_i(\boldsymbol{\theta})\|^2 \leq \sigma^2$$

Bounded Gradient Dissimilarity

$$\frac{1}{N} \sum_{i=1}^N \|\nabla f_i(\boldsymbol{\theta})\|^2 \leq B \|\nabla f(\boldsymbol{\theta})\|^2 + \sigma_h^2$$

Outline

1 Background
and
Motivation



2. Algorithm
Design



3. Simulation
Results

Key Techniques: Normalized Gradient Descent

- Traditional gradient descent:

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \nabla f(\mathbf{x}_t)$$

- The stepsize constraint is:

$$\text{e.g. } \eta \leq 1/L$$

The step size η must be small enough to account for the gradient's magnitude, which is governed by L . Large gradients can lead to overshooting and instability.

- Normalized gradient descent:

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla f(\mathbf{x}_t)}{\|\nabla f(\mathbf{x}_t)\|}$$

The gradient is normalized, so the step size η no longer depends on how large the gradient is (irrespective of L). Consistent step sizes allows NGD to navigate flat regions, steep regions, and saddle points more effectively.



Key Techniques: Normalized SGD

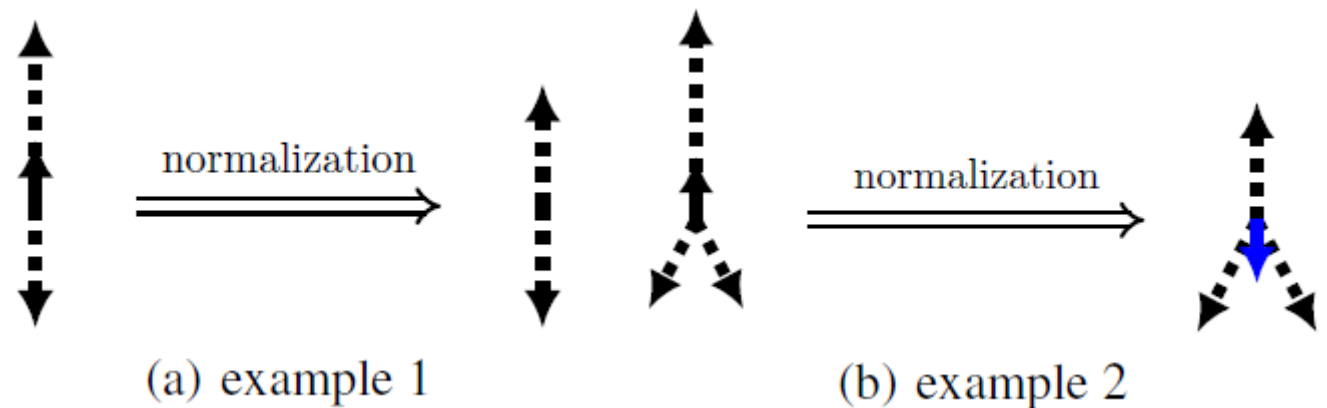
Normalized SGD: The stepsize is normalized by the current gradient.

- K steps local update at client i :

$$\mathbf{g}_i^{t,k} = \nabla F \left(\boldsymbol{\theta}_i^{t,k}; \boldsymbol{\xi}_i^{t,k} \right)$$

$$\boldsymbol{\theta}_i^{t,k+1} = \boldsymbol{\theta}_i^{t,k} - \eta_t \frac{\mathbf{g}_i^{t,k}}{\|\mathbf{g}_i^{t,k}\|}$$

Problems for Stochastic Gradients *



Normalization loses magnitude information, inefficient for stochastic gradients

*Yang J, Li X, Fatkhullin I, et al. Two sides of one coin: the limits of untuned SGD and the power of adaptive methods [J]. Advances in Neural Information Processing Systems, 2024, 36.

Key Techniques: Momentum

How to maintain the descent direction for normalized SGD: **Momentum**

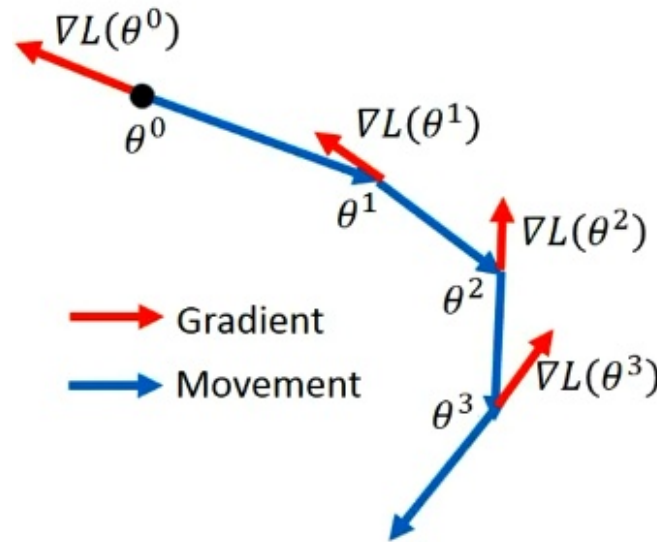
- Local update at client i :

$$\mathbf{g}_i^{t,k} = \beta \nabla F(\theta_i^t; \xi_i^t) + (1 - \beta) \mathbf{g}^t$$

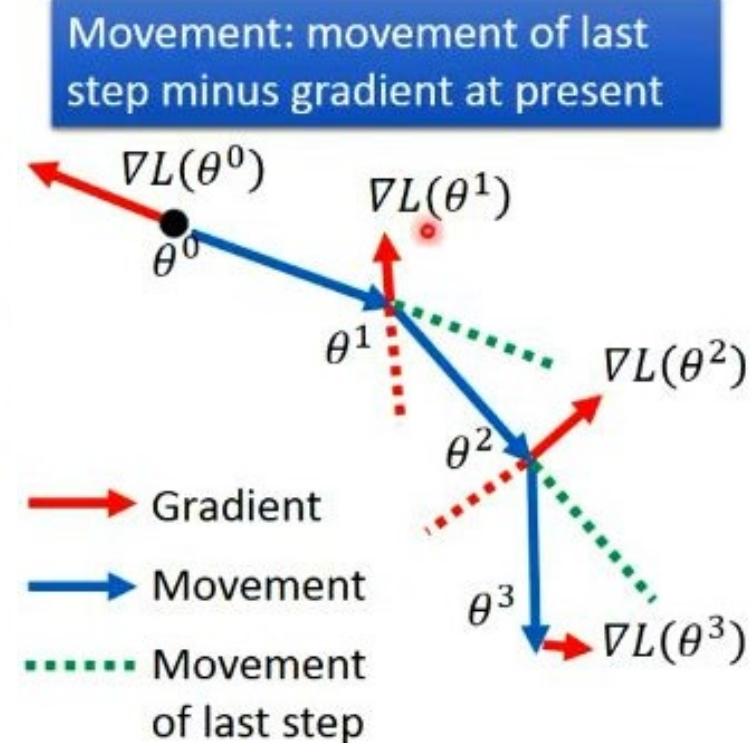
$$\theta_i^{t,k+1} = \theta_i^{t,k} - \eta_t \frac{\mathbf{g}_i^{t,k}}{\|\mathbf{g}_i^{t,k}\|}$$

- Global aggregation at server:

$$\mathbf{g}^t = \frac{1}{NK} \sum_{i,k} \mathbf{g}_{i,k}^t$$



→ Gradient
→ Movement



→ Gradient
→ Movement
→ Movement of last step

Momentum: Accumulating past gradients across iterations and clients

Key Techniques: Momentum

- Data Heterogeneity in FL:

$$f_i(\theta) := \mathbb{E}_{\xi_i \sim \mathcal{D}_i} [F(\theta; \xi_i)]$$

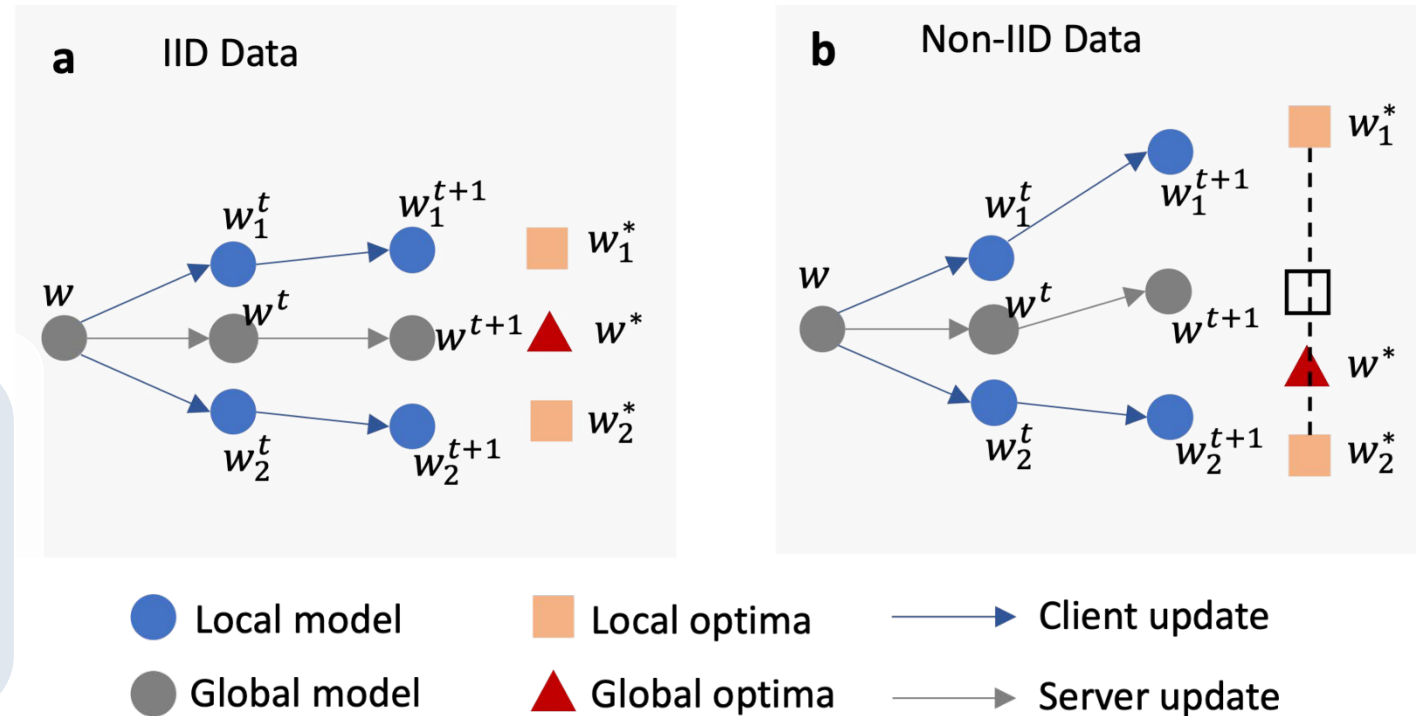
$$\mathcal{D}_i \neq \mathcal{D}_j \text{ for any } i \neq j$$

Bounded Gradient Dissimilarity

$$\frac{1}{N} \sum_{i=1}^N \|\nabla f_i(\theta) - \nabla f(\theta)\|^2 \leq \sigma_h^2$$



“Client drift”: Local updates from individual clients diverge significantly from one another and from the global objective



Key Techniques: Momentum

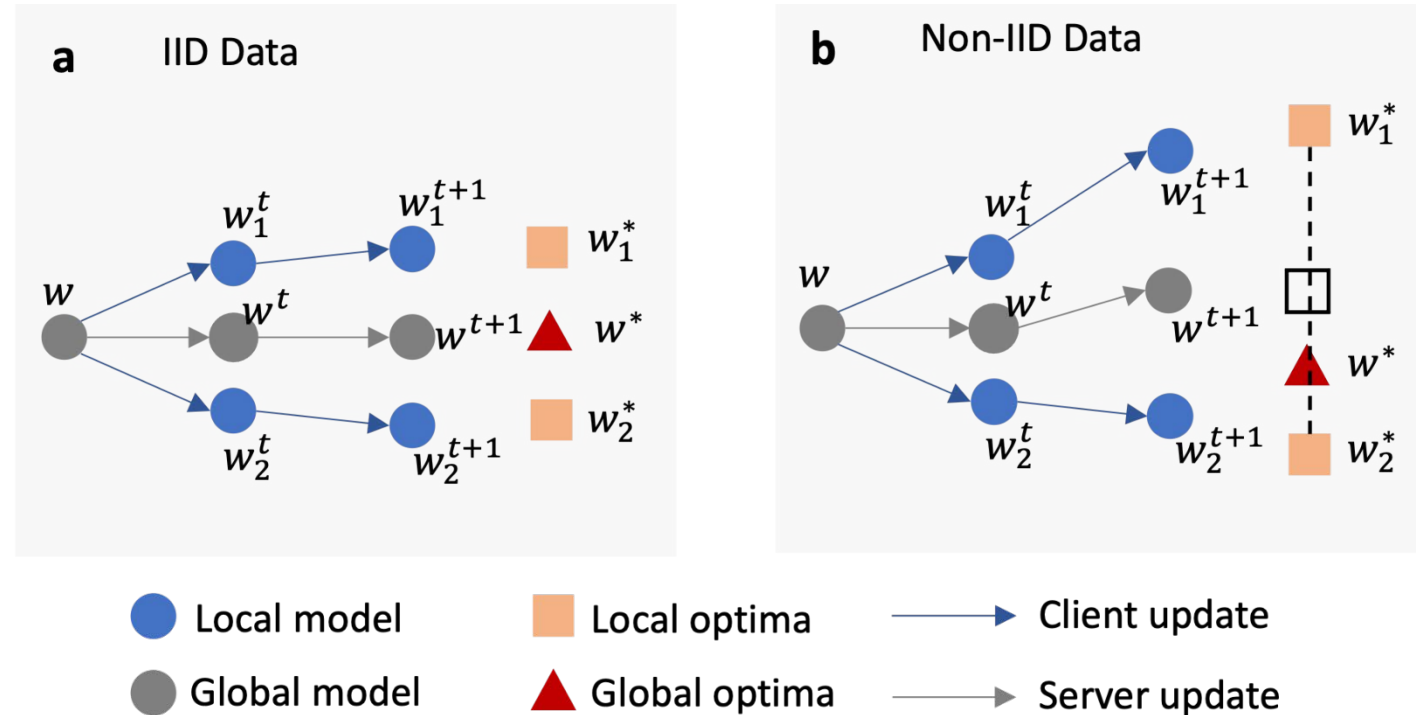
Momentum helps eliminating data heterogeneity bounds *

$$g_i^{t,k} = \beta \nabla F(\theta_i^t; \xi_i^t) + (1 - \beta) g^t$$

$$\theta_i^{t,k+1} = \theta_i^{t,k} - \eta_t g_i^{t,k}$$

$$g^t = \frac{1}{NK} \sum_{i,k} g_{i,k}^t$$

“Client drift”: Local updates from individual clients diverge significantly from one another and from the global objective



Momentum: Accumulating past gradients across iterations and clients

- * Cheng, Z., Huang, X., Wu, P., & Yuan, K. Momentum Benefits Non-iid Federated Learning Simply and Provably. In The Twelfth International Conference on Learning Representations (ICLR 2024).
- Figure from <https://arxiv.org/abs/2103.00710>

Key Techniques: Momentum

However...

$$\mathbf{g}_i^{t,k} = \beta \nabla F(\boldsymbol{\theta}_i^t; \boldsymbol{\xi}_i^t) + (1 - \beta) \mathbf{g}^t$$

$$\boldsymbol{\theta}_i^{t,k+1} = \boldsymbol{\theta}_i^{t,k} - \eta_t \frac{\mathbf{g}_i^{t,k}}{\|\mathbf{g}_i^{t,k}\|}$$

$$\mathbf{g}^t = \frac{1}{NK} \sum_{i,k} \mathbf{g}_{i,k}^t$$

To achieve parameter-free, our derivative needs to handle the following unavoidable term:

$$\begin{aligned} & \left\| \mathbf{g}_i^{t,k} - \mathbf{g}^t \right\|, \forall i, k \\ &= \beta \mathbb{E} \left\| \nabla F(\boldsymbol{\theta}_i^{t,k}; \boldsymbol{\xi}_i^{t,k}) - \frac{1}{NK} \sum_{i,k} \nabla F(\boldsymbol{\theta}_i^{t,k}; \boldsymbol{\xi}_i^{t,k}) \right\| \end{aligned}$$

Direct reflection of gradient dissimilarity



Momentum is **insufficient** to eliminate gradient dissimilarity bounds for normalized SGD case to achieve parameter-parameter-free

Key Techniques: Control Variates

Add control variates* to control data heterogeneity

- Local update at client i :

$$\mathbf{g}_i^{t,k} = \beta \left(\nabla F \left(\boldsymbol{\theta}_i^{t,k}; \boldsymbol{\xi}_i^{t,k} \right) - \mathbf{c}_i^{t-1} + \mathbf{c}^{t-1} \right) + (1 - \beta) \mathbf{g}^{t-1}$$

$$\mathbf{c}_i^t = \frac{1}{K} \sum_{k=0}^{K-1} \nabla F \left(\boldsymbol{\theta}_i^{t,k}; \boldsymbol{\xi}_i^{t,k} \right)$$

- Global aggregation at server:

➤ Full participation: $\mathbf{c}^t = \frac{1}{N} \sum_{i=1}^N \mathbf{c}_i^t$

➤ Partial participation: $\mathbf{c}^t = \mathbf{c}^{t-1} + \frac{1}{N} \sum_{i \in \mathcal{S}_t} (\mathbf{c}_i^t - \mathbf{c}_i^{t-1})$

* Karimireddy, Sai Praneeth, et al. "Scaffold: Stochastic controlled averaging for federated learning." International conference on machine learning. PMLR, 2020.

Algorithm 1 PAdaMFed: A Problem-Parameter-Agnostic Algorithm for Nonconvex FL

-
- 1: **Require:** initial model θ^0 , control variates $c_i^{-1} = \frac{1}{K} \sum_{k=0}^{K-1} \nabla F(\theta^0; \xi_i^{-1,k})$ for any i , $c^{-1} = \frac{1}{N} \sum_i c_i^{-1}$, momentum $g^{-1} = c^{-1}$, global learning rate γ , local learning rate η , and momentum parameter β
 - 2: **for** $t = 0, \dots, T - 1$ **do**
 - 3: **Central Server:** Uniformly sample clients $\mathcal{S}_t \subseteq \{1, \dots, N\}$ with $|\mathcal{S}_t| = S$
 - 4: **for** each client $i \in \mathcal{S}_t$ in parallel **do**
 - 5: Initialize local model $\theta_i^{t,0} = \theta^t$
 - 6: **for** $k = 0, \dots, K - 1$ **do**
 - 7: Compute $g_i^{t,k} = \beta \left(\nabla F(\theta_i^{t,k}; \xi_i^{t,k}) - c_i^{t-1} + c^{t-1} \right) + (1 - \beta)g^{t-1}$
 - 8: Update local model $\theta_i^{t,k+1} = \theta_i^{t,k} - \eta \frac{g_i^{t,k}}{\|g_i^{t,k}\|}$
 - 9: **end for**
 - 10: Update control variate $c_i^t = \frac{1}{K} \sum_{k=0}^{K-1} \nabla F(\theta_i^{t,k}; \xi_i^{t,k})$ (set $c_i^t = c_i^{t-1}$ for $i \notin \mathcal{S}_t$)
 - 11: Upload $\theta_i^{t,K}$ and c_i^t to central server
 - 12: **end for**
 - 13: **Central server:**
 - 14: Aggregate local updates $\bar{g}^t = \frac{1}{\eta SK} \sum_{i \in \mathcal{S}_t} (\theta^t - \theta_i^{t,K})$
 - 15: Update global model $\theta^{t+1} = \theta^t - \gamma \bar{g}^t$
 - 16: Aggregate control variate $c^t = c^{t-1} + \frac{1}{N} \sum_{i \in \mathcal{S}_t} (c_i^t - c_i^{t-1})$
 - 17: Aggregate momentum $g^t = \beta \left(\frac{1}{S} \sum_{i \in \mathcal{S}_t} (c_i^t - c_i^{t-1}) + c^{t-1} \right) + (1 - \beta)g^{t-1}$
 - 18: Download θ^{t+1} , $\beta c^t + (1 - \beta)g^t$ to all clients
 - 19: **end for**
-

Convergence of Algorithm 1

Theorem 1. Suppose that Assumptions 1 and 2 hold. Let the local and global learning rates of PAdaMFed be $\eta = \frac{1}{K\sqrt{T}}$ and $\gamma = \frac{(SK)^{1/4}}{T^{3/4}}$, respectively, the momentum parameter be $\beta = \sqrt{\frac{SK}{T}}$, and $\{\theta^t\}_{t \geq 0}$ be the iterates generated by Algorithm 1. Then, it holds for all $T \geq 1$ that

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(\theta^t)\| \leq \mathcal{O} \left(\frac{\Delta + L + \sigma + \sqrt{L}\sigma}{(SKT)^{\frac{1}{4}}} + \frac{\sqrt{SK}\sigma + L}{\sqrt{T}} \right),$$

where $\Delta := f(\theta^0) - \min_{\theta} f(\theta)$.

Remark 1

All hyperparameters (η : local learning rate, γ : global learning rate, and β : momentum parameter) in Algorithm 1 are **explicated determined by system-predefined constants**: S (the number of participation clients), K (local update times), T (iteration times)

Enhanced convergence by STORM* variance reduction

$$g_i^{t,k} = \nabla F \left(\theta_i^{t,k}; \xi_i^{t,k} \right) + \beta \left(c^{t-1} - c_i^{t-1} \right) + (1 - \beta) \left(g^{t-1} - \nabla F \left(\theta^{t-1}; \xi_i^{t,k} \right) \right)$$

Convergence of Algorithm 1 with variance reduction

Theorem 2. *Let the local and global learning rates of PAdaMFed-VR be $\eta = \frac{1}{KT}$ and $\gamma = \frac{(SK)^{1/3}}{T^{2/3}}$, respectively, the momentum parameter be $\beta = \frac{(SK)^{1/3}}{T^{2/3}}$, and $\{\theta^t\}_{t \geq 0}$ be the iterates generated by Algorithm 2. Then, it holds for all $T \geq 1$ that*

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(\theta^t)\| \leq \mathcal{O} \left(\frac{\Delta + L + \sigma}{(SKT)^{\frac{1}{3}}} + \frac{(L + \sigma)(SK)^{\frac{1}{3}}}{T^{\frac{2}{3}}} \right).$$

* Cutkosky, Ashok, and Francesco Orabona. "Momentum-based variance reduction in non-convex sgd." Advances in neural information processing systems 32 (2019).

Algorithms	Add. Assump.	Stepsize Restrictions	Stepsize-Related Problem-Parameters	Communication Complexity
SCAFFOLD (Karimireddy et al., 2020b)	–	$\gamma = \sqrt{S}, \eta \leq \frac{1}{24\gamma KL} \left(\frac{S}{N}\right)^{\frac{2}{3}}$	L	$\mathcal{O}\left(\left(\frac{N}{S}\right)^{\frac{1}{3}} \frac{L}{K\epsilon^4}\right)$
Mime (Karimireddy et al., 2020a)	BDH, BHD	$\eta = \sqrt{\frac{\Delta S}{L\tilde{G}TK^2}}, \tilde{G} = \sigma_h^2 + \frac{\sigma^2}{K}$	$L, \Delta, \sigma^2, \sigma_h^2$	$\mathcal{O}\left(\frac{1}{SK\epsilon^4}\right)$
FedSPS Sohom Mukherjee (2024)	BDH	$\eta_i^{t,k} = \min \left\{ \frac{F(\theta_i^{t,k}; \xi_i^{t,k}) - \ell_i^*}{c \ \nabla F(\theta_i^{t,k}; \xi_i^{t,k})\ ^2}, \eta_b \right\}_1$ $\eta_b \leq \min \left\{ \frac{1}{2cL}, \frac{1}{25LK} \right\}$	$L, \ell_i^*, \forall i$	$\mathcal{O}\left(\frac{1}{NK\epsilon^4}\right)$
SCAFFOLD-M (Cheng et al., 2024)	–	$\beta = \min \left\{ 1, \frac{S}{N^{\frac{2}{3}}}, \sqrt{\frac{L\Delta SK}{\sigma^2 T}}, \sqrt{\frac{L\Delta S^2}{G_0 N}} \right\}_2$ $\gamma = \frac{\beta}{L}, \eta KL \lesssim \min \left\{ \frac{1}{\sqrt{S}}, \frac{1}{\beta K^{\frac{1}{4}}}, \frac{\sqrt{S}}{N} \right\}$	L, Δ, σ^2, G_0	$\mathcal{O}\left(\frac{1}{SK\epsilon^4}\right)$
PAdaMFed (This paper)	–	$\beta = \sqrt{\frac{SK}{T}}, \gamma = \frac{(SK)^{\frac{1}{4}}}{T^{\frac{3}{4}}}, \eta = \frac{1}{K\sqrt{T}}$	–	$\mathcal{O}\left(\frac{1}{SK\epsilon^4}\right)$
Variance Reduction				
FAFED (Wu et al., 2023)	BDH, BG	$\eta_t \propto \frac{N^{\frac{2}{3}}}{Lt^{\frac{1}{3}}}, \beta_t \propto \eta_t^2$	L	$\mathcal{O}\left(\frac{1}{SK\epsilon^3}\right)$
SCAFFOLD-M-VR (Cheng et al., 2024)	–	$\beta = \min \left\{ \frac{S}{N}, \left(\frac{KL\Delta}{\sigma^2 T}\right)^{\frac{2}{3}}, S^{\frac{1}{3}} \right\}$ $\gamma L = \min \left\{ 1, \sqrt{\beta S} \right\}$ $\eta KL \lesssim \min \left\{ \sqrt{\frac{\beta}{S}}, \left(\frac{\beta}{SK}\right)^{\frac{1}{4}} \right\}$	L, Δ, σ^2	$\mathcal{O}\left(\frac{1}{S\sqrt{K}\epsilon^4}\right)$
PAdaMFed-VR (This paper)	–	$\beta = \frac{(SK)^{\frac{1}{3}}}{T^{\frac{2}{3}}}, \gamma = \frac{(SK)^{\frac{1}{3}}}{T^{\frac{2}{3}}}, \eta = \frac{1}{KT}$	–	$\mathcal{O}\left(\frac{1}{SK\epsilon^3}\right)$

Shorthand notation:

BDH = Bounded data heterogeneity

BG = Bounded gradient that

$$\|\nabla f_i(\theta)\| \leq G, \forall i, \theta$$

BHD = Bounded hessian dissimilarity that

$$\|\nabla^2 f_i(\theta) - \nabla^2 f(\theta)\|^2 \leq \delta, \forall i, \theta$$

 ${}^1\ell_i^* \leq \inf_{\xi_i \in \mathcal{D}_i, \theta} F(\theta; \xi_i)$ for any i , and c is a constant to balance adaptivity and accuracy.

$${}^2G_0 := \frac{1}{N} \sum_{i=1}^N \|\nabla f_i(\theta^0)\|^2.$$

A point θ is said to be ϵ -stationary if

$$\|\nabla f(\theta)\| \leq \epsilon.$$

For any ϵ -stationarypoint defined using $\|\nabla f(\theta)\|^2$, wehave guarantee for $\|\nabla f(\theta)\|$ by:

$$\frac{1}{T} \sum_{t=1}^{T-1} \mathbb{E} \|\nabla f(\theta^t)\|$$

$$= \frac{1}{T} \sum_{t=1}^{T-1} \mathbb{E} \sqrt{\|\nabla f(\theta^t)\|^2}$$

$$\leq \frac{1}{T} \sum_{t=1}^{T-1} \sqrt{\mathbb{E} \|\nabla f(\theta^t)\|^2}$$

$$\leq \sqrt{\frac{1}{T} \sum_{t=1}^{T-1} \mathbb{E} \|\nabla f(\theta^t)\|^2}$$

Outline

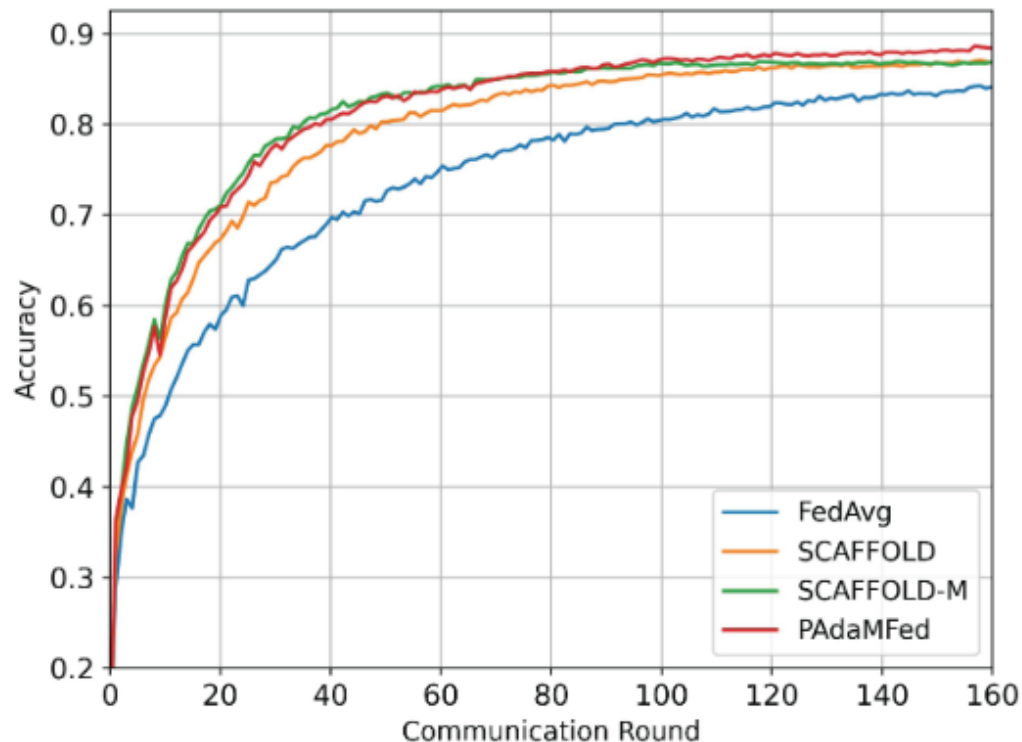
1 Background
and
Motivation



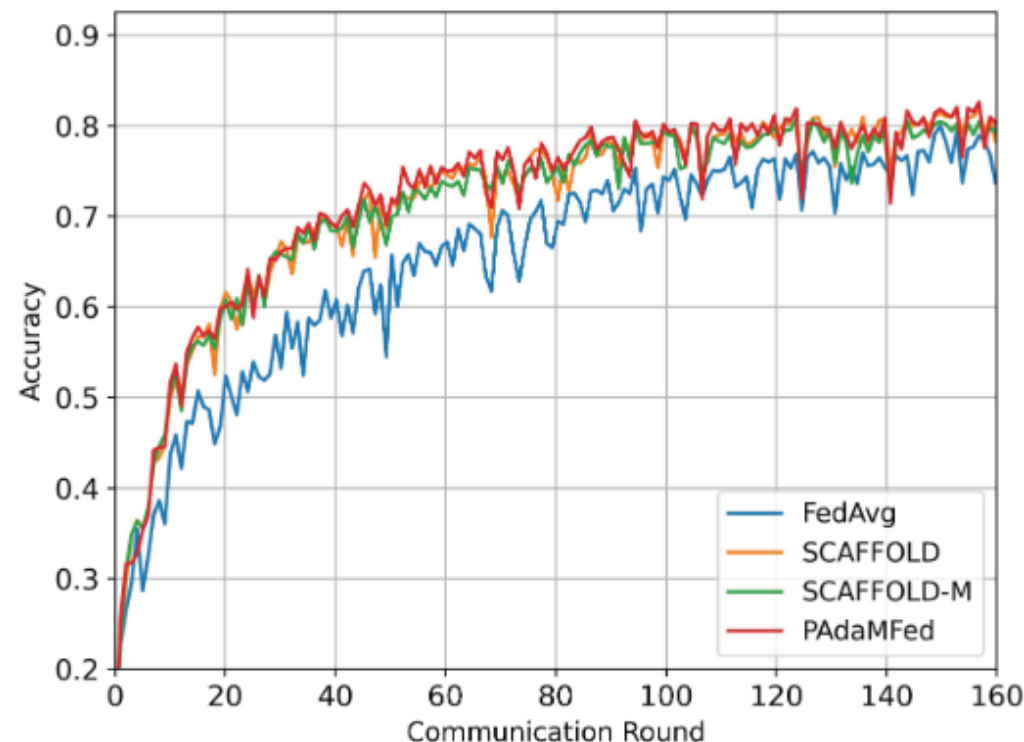
2. Algorithm
Design



3. Simulation
Results



(a) i.i.d

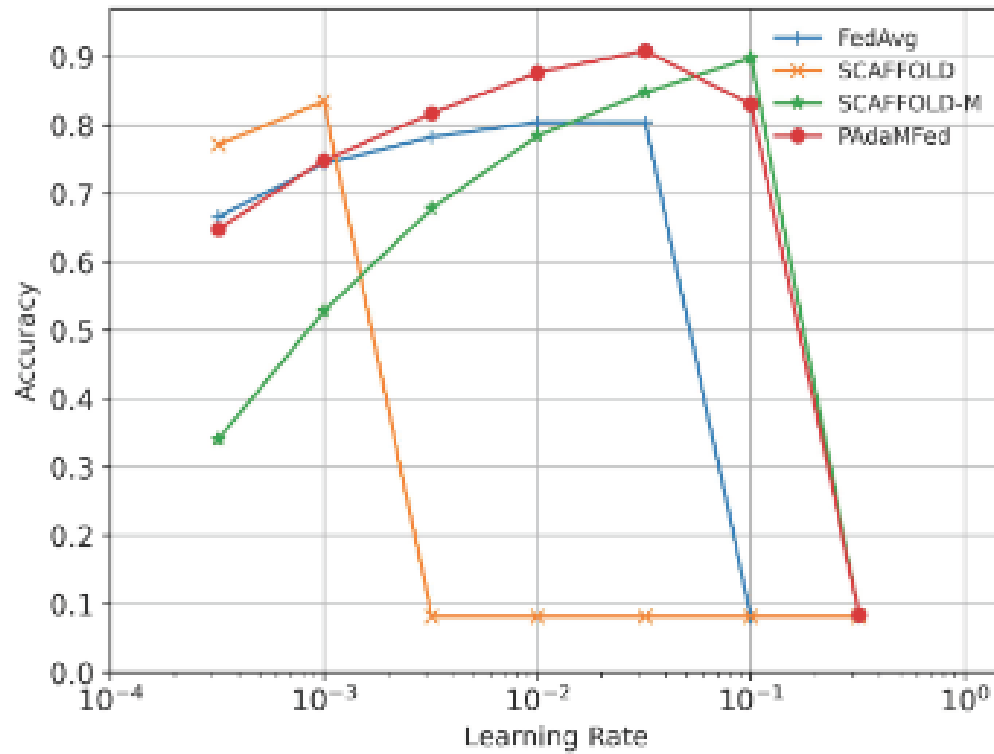


(b) non-i.i.d

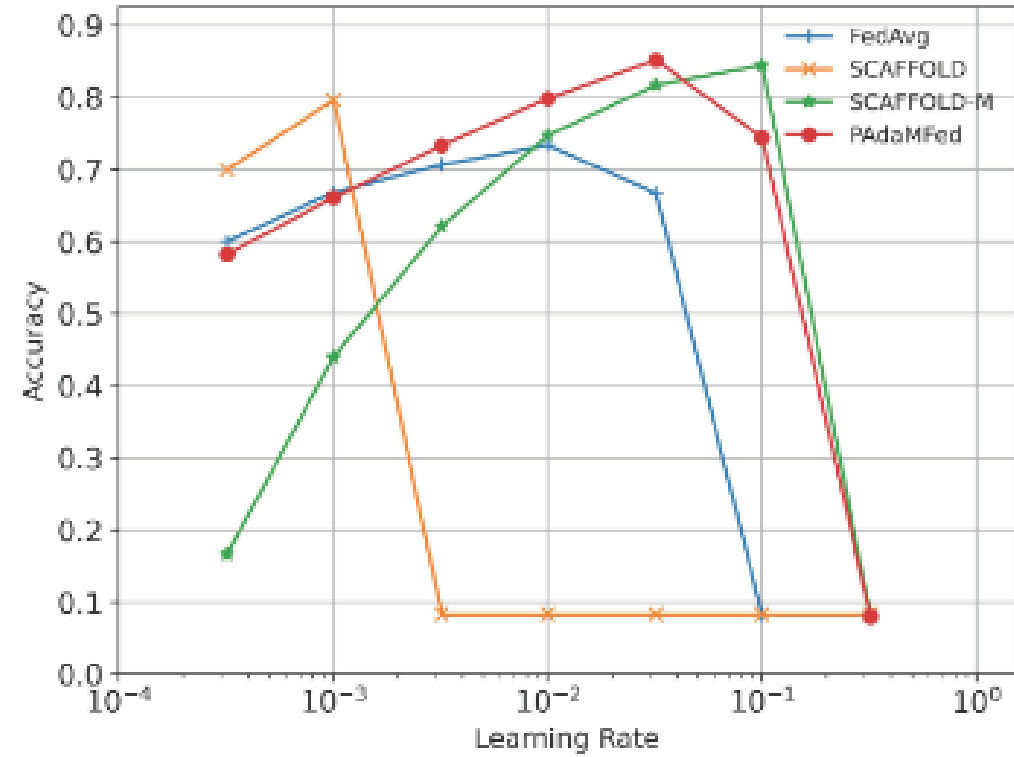
Figure 1. Test accuracy versus the number of communication rounds on the CIFAR-10 dataset (ResNet 18, Dir(0.5)).

Our approach achieves start-of-the-art performance while eliminating the tedious stepsize tuning process

The stepsizes of our algorithm are set directly based on the guidance of Theorem 1. The stepsize of all baselines are perfectly tuned by grid search.



(a) i.i.d



(b) non-i.i.d

Figure 1. Test accuracy versus learning rate on EMNIST dataset (RNN, 3 layer, Dir(1)).

Our algorithm demonstrates superior robustness to stepsize selection, maintaining stable performance across a significantly wider range of learning rates

THANKS

Thanks!