

# A Unified Analysis of Federated Learning with Arbitrary Client Participation

Shiqiang Wang<sup>1</sup> and Mingyue Ji<sup>2</sup>

<sup>1</sup> IBM T. J. Watson Research Center, Yorktown Heights, NY, USA

<sup>2</sup> Department of Electrical & Computer Engineering, University of Utah, Salt Lake City, UT, USA

# Federated Learning

- Collaborative model training while not sharing raw data
  - Local objective at client  $n$ :

$$F_n(\mathbf{x}) := \mathbb{E}_{\xi_n \sim \mathcal{D}_n} [\ell_n(\mathbf{x}, \xi_n)]$$

Loss function of model with parameter  $\mathbf{x}$  for data sample  $\xi_n$ :  $\ell_n(\mathbf{x}, \xi_n)$

- Global objective (not directly observable):

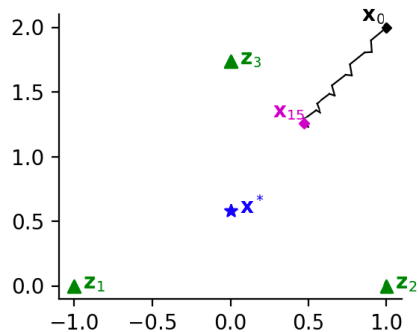
$$f(\mathbf{x}) := \frac{1}{N} \sum_{n=1}^N F_n(\mathbf{x})$$

Find  $\mathbf{x}^*$  to minimize  $f(\mathbf{x})$

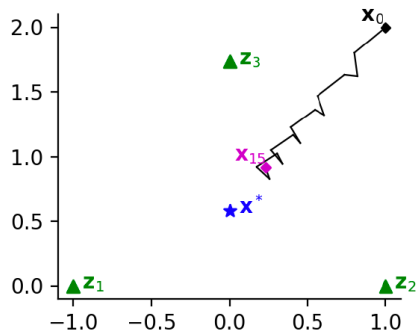
- *Federated averaging (FedAvg)*: local SGD at clients + parameter aggregation via the server
- **Our focus**
  - Clients may be only intermittently available to participate in training
    - For example: mobile devices during charging, edge servers when idle
  - Questions
    - How to effectively train models when clients participate arbitrarily?
    - How do unavailable clients affect the performance of model training?

# Problem with Intermittent Participation

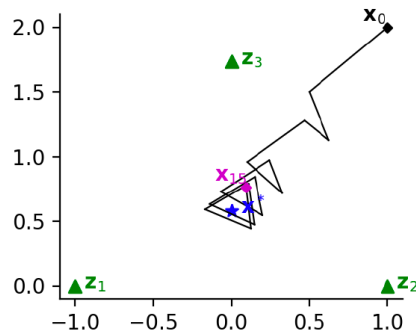
- Motivating example with  $F_n(\mathbf{x}) = \frac{1}{2} \|\mathbf{x} - \mathbf{z}_n\|^2$
- Three clients participating cyclically ( $P = 3$ ), one in each round



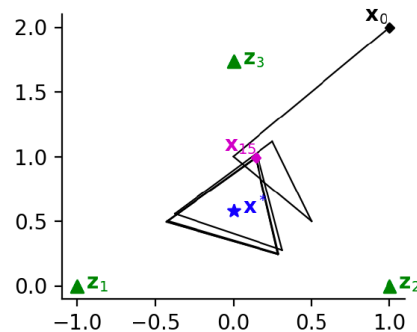
(a)  $\gamma = 0.05, \eta = 1$



(b)  $\gamma = 0.1, \eta = 1$



(c)  $\gamma = 0.25, \eta = 1$



(d)  $\gamma = 0.5, \eta = 1$

$\gamma$ : local learning rate

- Observation
  - Moves slowly to  $\mathbf{x}^*$  when  $\gamma$  is small
  - Circles around  $\mathbf{x}^*$  when  $\gamma$  is large



Apparent gap between  $\mathbf{x}_{15}$  and  $\mathbf{x}^*$

# Generalized FedAvg

```
1 Input:  $\gamma, \eta, \mathbf{x}_0, I, P, T$ ; Output:  $\{\mathbf{x}_t : \forall t\}$ 
2 Initialize  $t_0 \leftarrow 0, \mathbf{u} \leftarrow \mathbf{0}$ ;
3 for  $t \leftarrow 0, \dots, T - 1$  do
4   for  $n \leftarrow 1, \dots, N$  in parallel do
5      $\mathbf{y}_{t,0}^n \leftarrow \mathbf{x}_t$ ;
6     for  $i \leftarrow 0, \dots, I - 1$  do
7        $\mathbf{y}_{t,i+1}^n \leftarrow \mathbf{y}_{t,i}^n - \gamma \mathbf{g}_n(\mathbf{y}_{t,i}^n)$ 
8      $\Delta_t^n \leftarrow \mathbf{y}_{t,I}^n - \mathbf{x}_t$ ;
9      $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t + \sum_{n=1}^N q_t^n \Delta_t^n$ ; //update
10     $\mathbf{u} \leftarrow \mathbf{u} + \sum_{n=1}^N q_t^n \Delta_t^n$ ; //accumulate
11    if  $t + 1 - t_0 = P$  then
12       $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_{t+1} + (\eta - 1)\mathbf{u}$ ; //amplify
13       $t_0 \leftarrow t + 1$ ;
14       $\mathbf{u} \leftarrow \mathbf{0}$ ;
```

Participation weight

Amplification interval

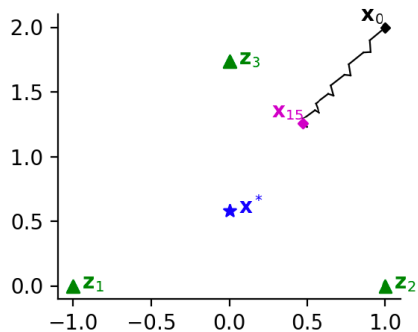
Amplification factor

Same as standard FedAvg

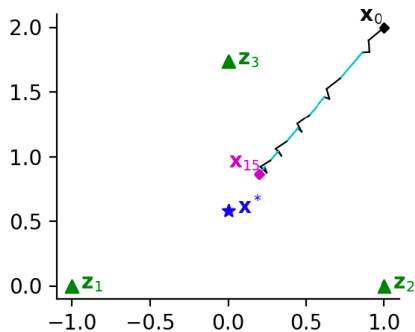
Accumulate and amplify updates every  $P$  rounds  
(no additional communication, minimal additional computation)

# Amplification Helps!

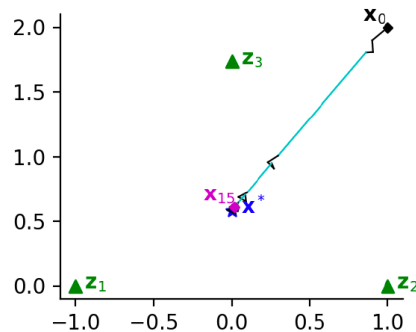
- Motivating example with  $F_n(\mathbf{x}) = \frac{1}{2} \|\mathbf{x} - \mathbf{z}_n\|^2$
- Three clients participating cyclically ( $P = 3$ ), one in each round



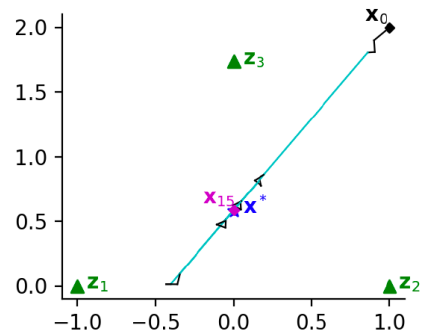
(a)  $\gamma = 0.05, \eta = 1$



(b)  $\gamma = 0.05, \eta = 2$



(c)  $\gamma = 0.05, \eta = 5$



(d)  $\gamma = 0.05, \eta = 10$

Change in  $x$  due to amplification shown in cyan color

- Observation
  - By choosing a smaller  $\gamma$  and a larger  $\eta$ , we can get very close to  $\mathbf{x}^*$  within only a few rounds

$\gamma$ : local learning rate  
 $\eta$ : amplification factor

# Main Building Block of Unified Analysis

**Assumption 1** (Lipschitz gradient).

$$\|\nabla F_n(\mathbf{x}) - \nabla F_n(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|, \forall \mathbf{x}, \mathbf{y}, n.$$

**Assumption 2** (Unbiased stochastic gradient with bounded variance).

$$\mathbb{E}[\mathbf{g}_n(\mathbf{x})] = \nabla F_n(\mathbf{x}) \text{ and } \mathbb{E}[\|\mathbf{g}_n(\mathbf{x}) - \nabla F_n(\mathbf{x})\|^2] \leq \sigma^2, \forall \mathbf{x}, n.$$

**Assumption 3** (Bounded gradient divergence).

$$\|\nabla F_n(\mathbf{x}) - \nabla f(\mathbf{x})\|^2 \leq d^2, \forall \mathbf{x}, n.$$

Decomposition of gradient divergence:

$$\left\| \sum_{n=1}^N q_t^n [\nabla F_n(\mathbf{x}) - \nabla f(\mathbf{x})] \right\|^2 \leq \tilde{\beta}^2, \forall \mathbf{x}, t,$$

$$\left\| \sum_{n=1}^N q_t^n \left\| \nabla F_n(\mathbf{x}) - \sum_{n'=1}^N q_t^{n'} \nabla F_{n'}(\mathbf{x}) \right\|^2 \leq \tilde{\nu}^2, \forall \mathbf{x}, t,$$

$$\left\| \frac{1}{P} \sum_{t=t_0}^{t_0+P-1} \sum_{n=1}^N q_t^n (\nabla F_n(\mathbf{x}) - \nabla f(\mathbf{x})) \right\|^2 \leq \tilde{\delta}^2(P), \forall \mathbf{x}, t_0.$$

Effect of partial participation

Choose depending on whether  $P$  scales in  $T$



**Corollary 3.2.** Choosing  $\gamma = \frac{1}{12LIP\sqrt{T}}$  and  $\eta = \min \left\{ \frac{12P\sqrt{LIF}}{\sigma\rho}; 12\sqrt{T} \right\}$ , for  $P \leq \frac{T}{2}$ , we have

$$\min_t \mathbb{E} \left[ \|\nabla f(\mathbf{x}_t)\|^2 \mid \mathcal{Q} \right] \leq \mathcal{O} \left( \frac{\sigma\rho\sqrt{L\mathcal{F}}}{\sqrt{IT}} + \frac{LP\mathcal{F}}{T} + \frac{\tilde{\nu}^2}{P^2T} + \frac{\tilde{\beta}^2}{T} + \tilde{\delta}^2(P) + \frac{\sigma^2}{IPT} \right).$$

**Corollary 3.3.** If  $\frac{\sqrt{\mathcal{F}}}{\rho\sqrt{LIT}} \leq \frac{1}{LIP}$ , choosing  $\gamma = \frac{1}{12LIP\sqrt{T}}$  and  $\eta = \frac{12P\sqrt{LIF}}{\rho}$ , for  $P \leq \frac{T}{2}$ , we have

$$\min_t \mathbb{E} \left[ \|\nabla f(\mathbf{x}_t)\|^2 \mid \mathcal{Q} \right] \leq \mathcal{O} \left( \frac{(1+\sigma^2)\rho\sqrt{L\mathcal{F}}}{\sqrt{IT}} + \frac{\tilde{\nu}^2}{P^2T} + \frac{\tilde{\beta}^2}{T} + \tilde{\delta}^2(P) + \frac{\sigma^2}{IPT} \right).$$

$$\begin{aligned} \sum_{n=1}^N q_t^n &= 1 \\ \sum_{n=1}^N (q_t^n)^2 &\leq \rho^2 \end{aligned}$$

# Results for Different Participation Patterns

- Assuming  $S$  clients participate in each round with equal weight, we have

$$q_t^n = 1/S \quad \rho = \left[ \sum_{n=1}^N (q_t^n)^2 \right]^{1/2} = 1/\sqrt{S}$$

“Linear speedup”

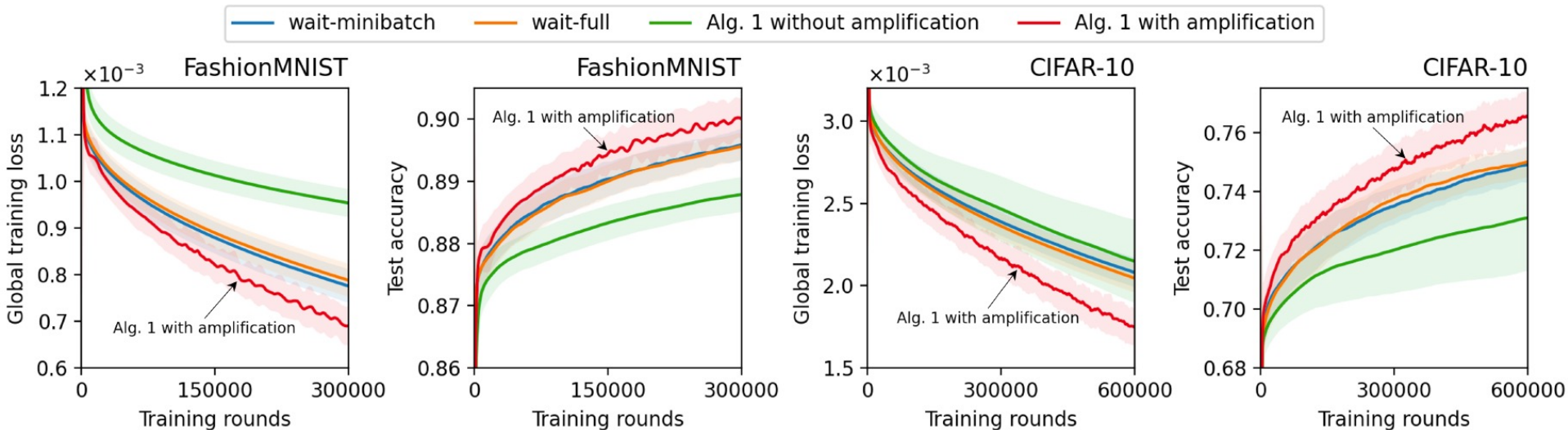
	Participation	Convergence error upper bound	Remark
From Corollary 3.2	Regularized	$\mathcal{O}\left(\frac{\sigma}{\sqrt{SIT}}\right)$	Matching centralized SGD lower bound
	Ergodic	Approaches zero as $T \rightarrow \infty$	
From Corollary 3.3	Mixing	$\mathcal{O}\left(\frac{1+\sigma^2}{\sqrt{SIT}}\right)$ i.E.; $\mathcal{O}\left(\frac{1+\sigma^2}{c\sqrt{SIT}}\right)$ w.p. $1-c$	Matching known bound with idealized participation
	Independent	$\mathcal{O}\left(\frac{1+\sigma^2}{\sqrt{SIT}}\right)$ i.E.; $\mathcal{O}\left(\frac{(1+\sigma^2)\ln(2/c)}{\sqrt{SIT}}\right)$ w.p. $1-c$	Matching known bound with idealized participation

The dominant term does not depend on  $P$

i.E. = in expectation  
w.p. = with probability

# Experiments

- Cyclic participation of clients with heterogeneous data, where each full cycle includes 500 rounds
- Optimized learning rates from grid search for each method
- $\eta = 10$  and  $P = 500$  for the generalized FedAvg algorithm with amplification





# Recap

- Generalized FedAvg with amplification
- A unified framework for convergence analysis with arbitrary participation
- Theoretical convergence bounds for different participation patterns
- Experiment confirming improvement compared to baselines

# Thank You!

Email: [wangshiq@us.ibm.com](mailto:wangshiq@us.ibm.com)  
Homepage: <https://shiqiang.wang/>