

# Resilient to Byzantine Attacks

## Finite-Sum Optimization over Networks

Zhaoxian Wu

School of Data and Computer Science, Sun Yat-Sen University

Joint work with  
Qing Ling (SYSU) Tianyi Chen (RPI) Georgios B. Giannakis (UMN)

45th International Conference on Acoustics, Speech, and Signal Processing (ICASSP 2020)

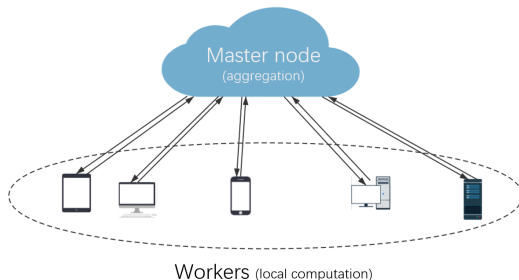
# Outline

- 1 Background
- 2 Problem Statement
- 3 Algorithm Development
- 4 Theoretical Analysis
- 5 Numerical Experiments
- 6 Conclusions

# Background

- Federated learning: promising distributed machine learning framework
  - (Distributed) Major computations are carried at workers locally
  - (Privacy-preserving) Data in workers are kept private
- There may be Byzantine attackers biasing the learning process

How to alleviate the negative effect caused by Byzantine attacks?



# Problem Statement

- Consider a network with one master node and  $W$  workers, among which  $B$  workers are Byzantine attackers
- The goal is to find the solution of the following optimization problem:

$$x^* = \arg \min_x f(x) := \frac{1}{W - B} \sum_{w \notin \mathcal{B}} f_w(x) \quad (1)$$

- Notations

- $f_w(x) := \frac{1}{J} \sum_{j=1}^J f_{w,j}(x)$ : local finite-sum objective
- $\mathcal{B}$ : set of Byzantine workers, with  $|\mathcal{B}| = B$
- $x \in \mathbb{R}^p$ : optimization variable

# Revisiting SGD

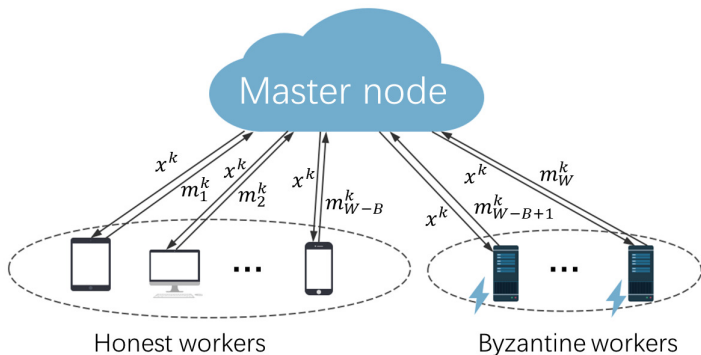
- Stochastic Gradient Descent (SGD): a popular solver of problem (1)
- SGD updates at time slot  $k$ 
  - Master node broadcasts  $x^k$  to workers
  - Worker  $w$  uniformly at random chooses a local data sample (or a mini batch) with index  $i_w^k$  and computes stochastic gradient
  - Worker  $w$  communicates  $f'_{w,i_w^k}(x^k)$  back to the master node
  - Master node updates the model as

$$x^{k+1} = x^k - \gamma^k \frac{1}{W} \sum_{w=1}^W f'_{w,i_w^k}(x^k) \quad (2)$$

where  $\gamma^k$  is the non-negative step size

- SGD is vulnerable to Byzantine attacks

# Illustration of Byzantine Attacks



**Figure 1:** Illustration of SGD in federated learning framework and Byzantine attacks model. In practice, the identities of Byzantine attackers are unknown to the master node

# Byzantine Attacks Model

- Byzantine workers can send arbitrary malicious messages to the master node
- Let  $m_w^k$  denote the message worker  $w$  sends to the master node at slot  $k$ , given by

$$m_w^k = \begin{cases} f'_{w,i_w^k}(x^k), & w \notin \mathcal{B} \\ *, & w \in \mathcal{B} \end{cases} \quad (3)$$

- The update rule of SGD can be written as

$$x^{k+1} = x^k - \gamma^k \frac{1}{W} \sum_{w=1}^W m_w^k \quad (4)$$

- Even only one Byzantine attacker can lead SGD to fail
  - $m_{w_b}^k = -\sum_{w \neq w_b} m_w^k$  yields  $x^{k+1} = x^k$
  - $m_{w_b}^k = +\infty$  yields  $x^{k+1} = +\infty$

# Existing Algorithm: SGD with Geometric Median

- To replace the mean with geometric median, we can robustify distributed SGD
- With  $\{z, z \in \mathcal{Z}\}$  denoting a subset in a normed space, the geometric median of  $\{z, z \in \mathcal{Z}\}$  is

$$\text{geomed}\{z\} := \arg \min_y \sum_{z \in \mathcal{Z}} \|y - z\| \quad (5)$$

- With geometric median, the distributed SGD in (2) can be modified to its Byzantine attack resilient form as

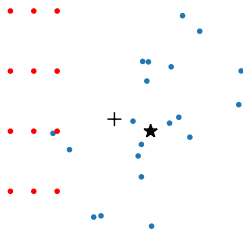
$$x^{k+1} = x^k - \gamma^k \cdot \text{geomed}\{m_w^k\} \quad (6)$$

- Geometric median has majority-voting property, which makes it insensitive to a few outliers

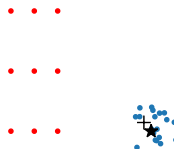


# Impact of Gradient Noise of SGD

gradients with large variance



gradients with small variance



- While geometric median can resist Byzantine attacks to some extent, its performance strongly depends on the amplitude of gradient noise
- Smaller noise may make the same Byzantine attacks less effective

# Revisiting Variance Reduction Techniques

- Our key idea is to reduce the variance of stochastic gradients in order to enhance robustness to Byzantine attacks
- An effective approach to alleviating stochastic gradient noise in SGD is through variance reduction
- Existing variance reduction techniques in stochastic optimization include mini-batch, SAG, SVRG, [SAGA](#), SDCA, SARAH, Katyusha, to list a few

# Revisiting Distributed SAGA with Mean

- In distributed SAGA, worker  $w$  store a gradient table  $f'_{w,j}(\phi_{w,j}^k)$  locally

$$\phi_{w,j}^0 = x^0, \quad \phi_{w,j}^{k+1} = \begin{cases} \phi_{w,j}^k, & j \neq i_w^k \\ x^k, & j = i_w^k \end{cases} \quad (7)$$

- $f'_{w,j}(\phi_{w,j}^k)$  refers to the previously stored stochastic gradient of the  $j$ -th data sample prior to slot  $k$  on worker  $w$
- Worker  $w$  will send the corrected stochastic gradient  $g_w^k$  to center

$$g_w^k := f'_{w,i_w^k}(x^k) - f'_{w,i_w^k}(\phi_{w,i_w^k}^k) + \frac{1}{J} \sum_{j=1}^J f'_{w,j}(\phi_{w,j}^k)$$

- The model update of SAGA is hence

$$x^{k+1} = x^k - \gamma \cdot \frac{1}{W} \sum_{w=1}^W g_w^k \quad (8)$$

where  $\gamma > 0$  is the constant step size

# Byrd-SAGA: Distributed SAGA with Geometric Median

- Let  $m_w^k$  denote the message from worker  $w$  at time slot  $k$

$$m_w^k = \begin{cases} g_w^k, & w \notin \mathcal{B} \\ *, & w \in \mathcal{B} \end{cases} \quad (9)$$

- Similar to distributed SGD, we can equip distributed SAGA with geometric median to defend against Byzantine attacks

$$x^{k+1} = x^k - \gamma \cdot \underset{w \in \mathcal{W}}{\text{geomed}}\{m_w^k\} \quad (10)$$

- This leads to the proposed **Byzantine-attack resilient distributed form of SAGA (Byrd-SAGA)**

# Convergence Analysis: Assumptions

## Assumption 1. *(Strong convexity and Lipschitz continuous gradients)*

*The function  $f$  is  $\mu$ -strongly convex and has  $L$ -Lipschitz continuous gradients.*

## Assumption 2. *(Bounded outer variation)*

*For any  $x \in \mathbb{R}^p$ , variation of the aggregated gradients at the honest workers with respect to the overall gradient is upper-bounded by*  
$$E_{w \notin \mathcal{B}} \|f'_w(x) - f'(x)\|^2 \leq \delta^2.$$

## Assumption 3. *(Bounded inner variation)*

*For every honest worker  $w$  and any  $x \in \mathbb{R}^p$ , the variation of its stochastic gradients with respect to its aggregated gradient is upper-bounded by*  
$$E_{i_w^k} \|f'_{w,i_w^k}(x) - f'_w(x)\|^2 \leq \sigma^2, \forall w \notin \mathcal{B}.$$

# Convergence Analysis: SAGA vs SGD

## Theorem 1: SAGA with geometric median (Byrd-SAGA)

Under Assumptions 1 and 2, if the number of Byzantine attackers satisfies  $B < \frac{W}{2}$  and the step size satisfies  $\gamma \leq \frac{\mu}{8J^2 C_\alpha L^2}$ , then for Byrd-SAGA with geometric median aggregation, it holds that

$$E\|x^k - x^*\|^2 \leq O((1 - \frac{\gamma\mu}{2})^k) + O(\delta^2). \quad (11)$$

## Theorem 2: SGD with geometric median

Under Assumptions 1, 2 and 3, if the number of Byzantine attackers is  $B < \frac{W}{2}$  and the step size satisfies  $\gamma < \frac{L}{2\mu^2}$ , then for Byzantine attack resilient SGD with geometric median aggregation, it holds that

$$E\|x^k - x^*\|^2 \leq O((1 - \gamma\mu)^k) + O(\sigma^2 + \delta^2). \quad (12)$$

Asymptotic error: from  $O(\sigma^2 + \delta^2)$  to  $O(\delta^2)$

# Numerical Experiments: $\ell_2$ -regularized Logistic Regression

- $\ell_2$ -regularized logistic regression on IJCNN1 dataset
- $W - B = 50$  honest workers and  $B = 20$  Byzantine workers (if any)
- Benchmark algorithms: [SGD](#), mini-Batch SGD with batch size 50 ([BSGD](#)), [SAGA](#) equipped with mean or geometric median
- Attacks
  - [Gaussian attacks](#): draws its  $m_w^k$  from a Gaussian distribution with mean  $\frac{1}{W-B} \sum_{w' \notin \mathcal{B}} m_{w'}^k$  and variance 30
  - [Sign-flipping attacks](#):  $m_w^k = u \cdot \frac{1}{W-B} \sum_{w' \notin \mathcal{B}} m_{w'}^k$ , where the magnitude  $u = -3$
  - [Zero gradient attacks](#):  $m_w^k = -\frac{1}{B} \sum_{w' \notin \mathcal{B}} m_{w'}^k$

# Numerical Experiments: $\ell_2$ -regularized Logistic Regression

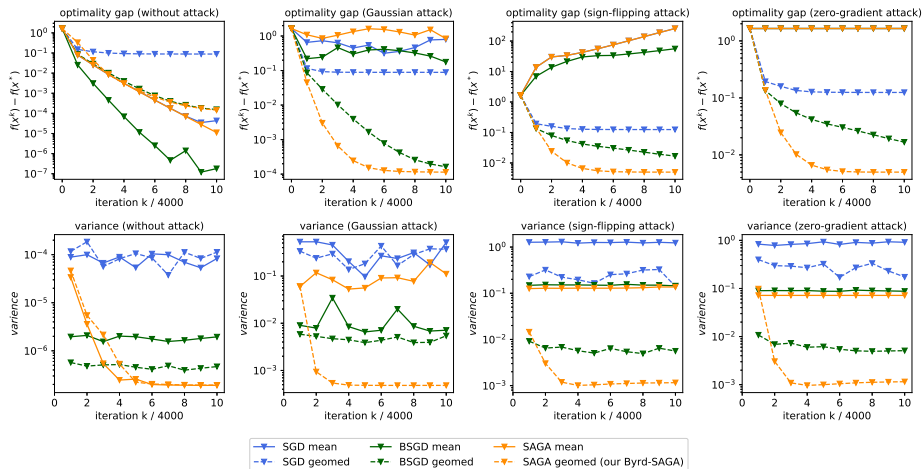


Figure 2: The first row shows the optimal gap  $f(x^k) - f(x^*)$  and the second row shows the variance of honest workers



- For the first time, we reveal the relation between gradient noise and resilience to Byzantine attacks
- We design a novel algorithm, Byrd-SAGA, to defend against Byzantine attacks in federated learning
- Byrd-SAGA reduces the gradient noise in SGD and achieves better resilience to Byzantine attacks

## Thanks