

# Byzantine-resilient distributed learning under constraints

**Dongsheng Ding**

a joint work with

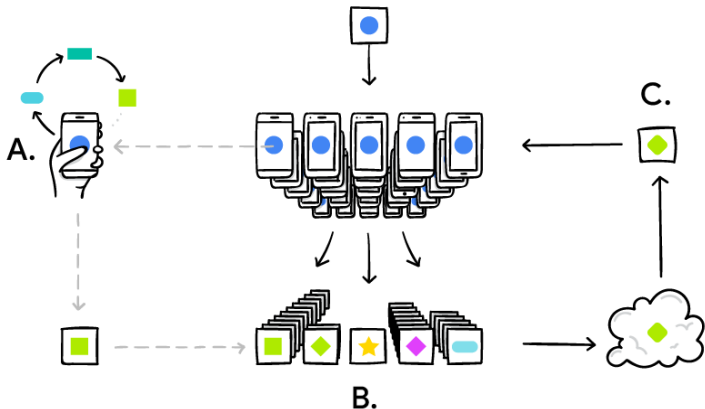
Xiaohan Wei, Hao Yu, Mihailo Jovanović



**2021 American Control Conference**

# Motivating application

- FEDERATED LEARNING



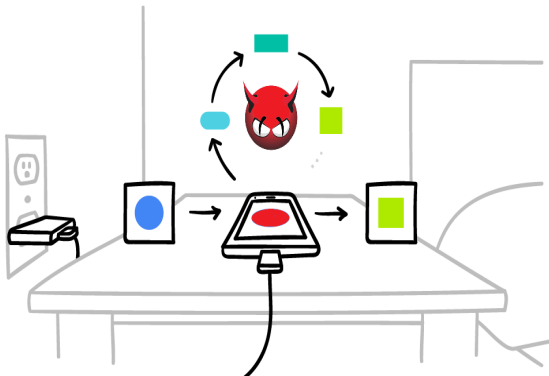
A. Worker machine

B. Master machine

C. Shared model

Google AI, Blog '17

# Byzantine fault



- **FAULT SOURCES**

- ★ Machine failures
- ★ Communication errors
- ★ Malicious users

# Byzantine failure model

- STOCHASTIC LEARNING PROBLEM WITH CONSTRAINTS

$$\begin{aligned} & \underset{w \in \mathcal{W}}{\text{minimize}} && F(w) := \mathbb{E}_{z \sim \mathcal{D}}(f(w; z)) \\ & \text{subject to} && g_j(w) \leq 0, \quad j = 1, \dots, k. \end{aligned}$$

# Byzantine failure model

- STOCHASTIC LEARNING PROBLEM WITH CONSTRAINTS

$$\begin{aligned} & \underset{w \in \mathcal{W}}{\text{minimize}} && F(w) := \mathbb{E}_{z \sim \mathcal{D}}(f(w; z)) \\ & \text{subject to} && g_j(w) \leq 0, \quad j = 1, \dots, k. \end{aligned}$$

## 1 master and $m$ workers

- ★  $g_j$  – deterministic constraints
- ★  $z_t^i \sim \mathcal{D}, i \in \{1, \dots, m\}$

# Byzantine failure model

- STOCHASTIC LEARNING PROBLEM WITH CONSTRAINTS

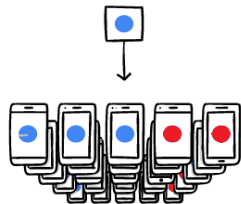
$$\begin{aligned} & \underset{w \in \mathcal{W}}{\text{minimize}} && F(w) := \mathbb{E}_{z \sim \mathcal{D}} (f(w; z)) \\ & \text{subject to} && g_j(w) \leq 0, \quad j = 1, \dots, k. \end{aligned}$$

## 1 master and $m$ workers

- \*  $g_j$  – deterministic constraints
- \*  $z_t^i \sim \mathcal{D}, i \in \{1, \dots, m\}$

$m$  gradients at time  $t$

$$\nabla_t^i := \begin{cases} \nabla f(w_t; z_t^i) & \text{normal machine} \\ \text{arbitrary} & \text{Byzantine machine} \end{cases}$$



$\alpha$  fraction

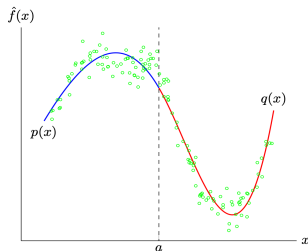
# Example

- CONSTRAINED LINEAR REGRESSION

$$f(w; z) = \frac{1}{2} (x^T w - y)^2$$

$$g_j(w) = A_j w - b_j, \quad j = 1, \dots, k.$$

- ★  $z := (x, y)$  – data from distribution  $\mathcal{D}$
- ★  $g_j(w) \leq 0$  – constraints, e.g., smooth spline fitting



# **Adding Byzantine resilience to primal-dual methods**



# Identification of "good" workers

- MEDIAN AGGREGATION

- ★ robust to outliers

sequence	mean	median
1, 2, 3, 6, 7, 8, 10	5.3	6
$10^{-3}$ , 2, 3, 6, 7, 8, $10^3$	145.6	6

- ★ breakdown point 50%

# Bregman divergence

$$D(x, y) := \phi(x) - \phi(y) - \nabla\phi(y)^T(x - y)$$

★  $\phi$  – differentiable, 1-strongly convex w.r.t.  $\|\cdot\|$

## • EXAMPLES

★  $\phi(x) = \frac{1}{2} \|x\|_2^2$  strongly convex w.r.t.  $\|\cdot\|_2$

$$D(x, y) = \frac{1}{2} \|x - y\|_2^2$$

★  $\phi(x) = \sum_i x(i) \log x(i)$  strongly convex w.r.t.  $\|\cdot\|_1$

$$D(x, y) = \sum_i x(i) \log \frac{x(i)}{y(i)} \quad \text{KL divergence}$$

# Byzantine primal-dual method

- BYZANTINE PRIMAL MIRROR DESCENT

$$w_{t+1} := \operatorname{argmin}_{w \in \mathcal{W}} \left\langle \xi_t + \sum_{j=1}^k (q_{j,t+1} + g_j(w_t)) \nabla g_j(w_t), w - w_t \right\rangle + \eta_t D(w, w_t)$$

# Byzantine primal-dual method

- BYZANTINE PRIMAL MIRROR DESCENT

$$w_{t+1} := \operatorname{argmin}_{w \in \mathcal{W}} \left\langle \xi_t + \sum_{j=1}^k (q_{j,t+1} + g_j(w_t)) \nabla g_j(w_t), w - w_t \right\rangle + \eta_t D(w, w_t)$$

★  $\xi_t$  – stochastic estimate of the gradient  $\nabla F(w_t)$

$$\xi_t = \frac{1}{m} \sum_{i \in \Omega_t} \nabla_t^i, \quad \Omega_t \text{ – set of "good" workers}$$

# Byzantine primal-dual method

- BYZANTINE PRIMAL MIRROR DESCENT

$$w_{t+1} := \operatorname{argmin}_{w \in \mathcal{W}} \left\langle \xi_t + \sum_{j=1}^k (q_{j,t+1} + g_j(w_t)) \nabla g_j(w_t), w - w_t \right\rangle + \eta_t D(w, w_t)$$

★  $\xi_t$  – stochastic estimate of the gradient  $\nabla F(w_t)$

$$\xi_t = \frac{1}{m} \sum_{i \in \Omega_t} \nabla_t^i, \quad \Omega_t \text{ – set of "good" workers}$$

- DUAL UPDATE

$$q_{j,t+2} = \max(-g_j(w_{t+1}), q_{j,t+1} + g_j(w_{t+1}))$$

## ● IMPORTANT QUANTITIES

★  $\nabla_t^1, \dots, \nabla_t^m$  – gradients (normal or Byzantine)

★  $A_t^1, \dots, A_t^m$  – gradient related values

$$A_t^i := \sum_{\tau=1}^t \langle \nabla_{\tau}^i, w_{\tau} - w_1 \rangle$$

★  $B_t^1, \dots, B_t^m$  – accumulated gradients

$$B_t^i := \sum_{\tau=1}^t \nabla_{\tau}^i$$

used to update the set of "good" workers

Ding, Wei, Jovanović. CDC '19  
Alistarh, Allen-Zhu, Li, NeurIPS '18

- UPDATE THE SET OF 'GOOD' WORKERS

$\Omega_t \leftarrow i \in \Omega_{t-1}$  satisfies

$$\left\{ \begin{array}{l} |A_t^i - A_{\text{med}}| \leq I_A \\ \|B_t^i - B_{\text{med}}\|_* \leq I_B \\ \|\nabla_t^i - \nabla_{\text{med}}\|_* \leq 4C \end{array} \right.$$

$$\Omega_0 = \{1, \dots, m\}$$

$$\|\nabla_t^i - \nabla_t\|_* \leq C$$

$$I_A = 4WC\Delta\sqrt{2T}$$

$$I_B = 4C\Delta\sqrt{2T}$$

$$\Delta = \Theta\left(\sqrt{\log \frac{mT}{\delta}}\right)$$

# Convergence result

- OPTIMALITY GAP

$$F(\bar{w}) - F(w^*) \leq \tilde{O}\left(\frac{1}{T} + \frac{1}{\sqrt{mT}} + \frac{\alpha}{\sqrt{T}}\right) \text{ w.h.p.}$$

- CONSTRAINT VIOLATION

$$g_j(\bar{w}) \leq \tilde{O}\left(\frac{1}{T} + \frac{1}{\sqrt{mT}} + \frac{\alpha}{\sqrt{T}}\right) \text{ w.h.p.}$$

$$\alpha \in [0, 0.5)$$

$$\bar{w} = \frac{1}{T} \sum_{t=0}^{T-1} w_t$$

$F, g$  – convex and smooth



# Convergence result

- OPTIMALITY GAP

$$F(\bar{w}) - F(w^*) \leq \tilde{O}\left(\frac{1}{T} + \frac{1}{\sqrt{mT}} + \frac{\alpha}{\sqrt{T}}\right) \text{ w.h.p.}$$

- CONSTRAINT VIOLATION

$$g_j(\bar{w}) \leq \tilde{O}\left(\frac{1}{T} + \frac{1}{\sqrt{mT}} + \frac{\alpha}{\sqrt{T}}\right) \text{ w.h.p.}$$

$$\alpha \in [0, 0.5)$$

$$\bar{w} = \frac{1}{T} \sum_{t=0}^{T-1} w_t$$

$F, g$  – convex and smooth

matches the rate of  $\left\{ \begin{array}{l} \text{Byzantine MD \& SGD, } \alpha \neq 0 \\ \text{batch SGD, } \alpha = 0 \end{array} \right.$

# Summary

- RESULTS

- ★ Byzantine primal-dual method
- ★ Optimal rate

- ONGOING EFFORT

- ★ Nonsmooth optimization problems
- ★ Byzantine constraints

**Thank you !**