

Byzantine-Resilient Multi-Agent Optimization

Lili Su and Nitin H. Vaidya *Fellow, IEEE*

Abstract—We consider the problem of multi-agent optimization wherein an unknown subset of agents suffer Byzantine faults and thus behave adversarially. We assume that each agent i has a local cost function f_i , and the overarching goal of the good agents is to collaboratively minimize a global objective that properly aggregates these local cost functions. To the best of our knowledge, we are among the first to study Byzantine-resilient optimization where no central coordinating agent exists, and we are the first to characterize the structures of the convex coefficients of the achievable global objectives.

Dealing with Byzantine faults is very challenging. For example, in contrast to fault-free networks, reaching Byzantine-resilient agreement even in the simplest setting is far from trivial. We take a step towards solving the proposed Byzantine-resilient multi-agent optimization problem by focusing on scalar local cost functions. Our results might provide useful insights for the general local cost functions.

I. INTRODUCTION

Networked multi-agent systems consist of a group of agents that perform collaborative tasks. The problem of multi-agent optimization typically assumes that each agent in the network has a *local* cost function, and the overarching goal of the networked agents is to minimize a global objective that properly aggregates these local costs. One standard choice of such global objective is the average of the local cost functions [1].

Multi-agent optimization over adversary-free networks is well-studied [1]–[3]. In this paper, we study adversary-prone networks. In particular, we consider the scenario wherein an unknown subset of agents suffer Byzantine faults – a canonical fault model in distributed computing [4] – and thus behave adversarially against the good agents. Unfortunately, having the average of the local cost functions as the global objective makes the multi-agent network extremely vulnerable to Byzantine faults; the average can be completely controlled by even a single adversarial agent. Similar observation is made in [5]. Thus, in our problem, the common goal of the good agents is to collaboratively minimize a global objective that properly aggregates the local cost functions at the *good agents*. The problem formulation (including the network model, the Byzantine fault model, and the

family of desired global objectives) can be found in Section III. To the best of our knowledge, we are among the first to study Byzantine-resilient optimization where no central coordinating agent exists, and we are the first to characterize the structures of the convex coefficients of the achievable global objectives.

Dealing with Byzantine faults is very challenging; in contrast to fault-free networks, reaching Byzantine-resilient agreement even in the simplest setting is far from trivial. We take a step towards solving the proposed Byzantine-resilient multi-agent optimization problem by focusing on scalar local cost functions. Attempts have been made to generalize our results to a broader family of local cost functions [6]. We would like to leave the general local cost functions as one important future direction. Our contributions can be summarized as follows:

- We propose (in Section III-C) a Byzantine-resilient multi-agent optimization problem, wherein the global objective is a convex combination of the local cost functions at the good agents and the structure of the convex coefficients is quantified by a metric named (β, γ) -admissibility (see Definition 1). Notably, our problem formulation is also valid for general local functions.
- We show in Theorem 1 that no algorithms can guarantee (β, γ) -admissibility with $\gamma > n - \phi - b$, where n is the number of agents, ϕ is the actual number of agents suffering Byzantine faults, and b is the maximum number of Byzantine faults.
- We propose an algorithm that is provably resilient to Byzantine faults (i.e., Algorithm 1). We characterize the (β, γ) tuples it can guarantee when the network satisfies certain topological conditions (in Theorem 3). When the network is a complete graph (e.g., logically fully connected), Algorithm 1 can guarantee to achieve (β, γ) for $\beta = \frac{1}{2(n-\phi-b)}$ and $\gamma = n - \phi - b$, matching the bound in the impossibility result in terms of γ . More importantly, given $\gamma = n - \phi - b$, the achieved $\beta = \frac{1}{2(n-\phi-b)}$ matches the optimal β up to a multiplicative factor $\frac{1}{2}$. The results for complete graphs are summarized in Fig. 1. Notably, adversary-prone networks with complete graphs are distributed systems [4].

This paper unifies the results scattered in our unpublished technical reports and the preliminary works [7], [8]. In particular, the impossibility result first appeared in [7] (also in [9]) and the algorithm was presented in

This research is supported in part by NSF awards 1329681 and 1842198.

Lili Su is with the CSAIL, MIT, (lilisu@mit.edu). Nitin H. Vaidya is with the Department of Computer Science, Georgetown University, (nitin.vaidya@georgetown.edu).

[8] and [10], [11]. In contrast to [8] whose focus is on complete graphs, we focus on general networks.

II. RELATED WORK

Byzantine-resilient consensus: There is a significant body of work on Byzantine-resilient consensus [12]–[18]. Readers are referred to [12] for comprehensive survey. Next we give a brief review on this line of work.

The Byzantine fault-tolerance problem was first introduced in [19], and has been one of the most fundamental problems in distributed computing for decades. Fisher, Lynch, and Paterson showed that the fault-tolerant consensus problem cannot be solved in an *asynchronous* system [20]. As one way to circumvent this impossibility result, the notion of *approximate consensus* was introduced [13], which only requires that processes agree with each other approximately rather than exactly in finite time. Reaching approximate consensus is of interest in *synchronous* systems as well [13]–[15]. The discussion in this paper applies to synchronous systems.

For undirected networks, approximate consensus can be achieved if and only if the network connectivity is at least $2b+1$ and $n > 3b$ [21]; recall that b is the maximum number of Byzantine faults. For directed networks, a sufficient and necessary condition on the network structures was characterized in [22]. There has been increasing interest in designing iterative approximate Byzantine consensus algorithms wherein only local communication is allowed [14], [15], [23], [24]. In particular, [24] studied the convergence rate over complete networks, and [14], [15] considered arbitrary directed networks and derived necessary and sufficient topological conditions. Recently, Byzantine consensus subject to differential privacy requirements was considered in [25].

All the above work focuses on scalar inputs. Multi-dimensional inputs have been studied recently [17], [26], [27]. Complete graphs were considered in [26], [27], where tight conditions on the number of agents were identified. Incomplete graphs were studied in [17].

Our work is most relevant to the line of work on *approximate* synchronous Byzantine consensus over *arbitrary graphs* [14], [15]. In our work, in each iteration, each good agent combines an approximate Byzantine consensus update with its local gradient descent update. Intuitively speaking, an approximate consensus update is used as a mechanism for each of the good agents to “robustly collect” information from others in the presence of Byzantine agents. This is in contrast to the exact Byzantine consensus protocols that are involved in Blockchains wherein exact Byzantine consensus is used as a “selection mechanism” to determine among the multiple blocks proposed which one should be appended to the blockchain.

Distributed optimization: Our work goes beyond Byzantine consensus in that the local inputs are functions, and the goal is to reach an agreement on a value that is a minimizer of some weighted average of these local cost functions.

Distributed optimization has a long history. The seminal works [28], [29] considered separable global objectives for which the local decision variables at different agents are allowed to be different. Nedic and Ozdaglar [1] studied the setting wherein the global objective is the average of these local cost functions, and the local variables at the agents are required to reach consensus asymptotically. Many follow-up works are inspired; see [30], [31] for comprehensive surveys. Nevertheless, little attention has been paid to adversary-prone networks.

Multi-agent optimization over adversary-prone networks, to the best of our knowledge, was first considered by Sundaram and Ghareisifard [5] and our technical reports [9]–[11], [32] with different fault models and global objectives. In contrast to the Byzantine fault model which assumes that a bad agent can send differently-valued messages to different neighbors, in [5] an adversarial agent is only allowed to send identically-valued (broadcast) messages. This difference is significant as in complete graphs, for the faults in [5], there exists a consensus algorithm that can tolerate less than $\frac{1}{2}$ of the agents to be faulty; in contrast, it is well-known that no consensus algorithms can tolerate more than $\frac{1}{3}$ agents to be Byzantine [33]. It might be enough to consider the fault model in [5] when the networked agents communicate with each other via wireless communication. In fact, the more structure on the adversarial behaviors, the easier to secure multi-agent optimization. Additionally, [5] considered the family of global objectives in the form of convex combinations of the local cost functions at the good agents, and no additional structures on the convex coefficients are required. Consequently, the local estimates at the good agent in [5] are only guaranteed to converge to a convex combination of the minima of those local functions. In contrast, in addition to being in the convex hull, we also characterize a structure of the convex coefficients.

Byzantine-tolerant distributed machine learning: Another line of work that is relevant to the discussion of this paper is Byzantine-tolerant distributed machine learning [34], [35]. The main difference between these papers and the problem considered here is that they considered the learning problem under a statistical learning framework. Specifically, they assumed that the training data is *i.i.d.* generated from a unknown distribution, and the learning goal there is to minimize the population loss, which is defined in a form of integration -where the integration is taken over the unknown underlying distribution.

III. PROBLEM FORMULATION

A. Network Model

A multi-agent network consists of a collection of n agents/nodes that interact with each other through a directed communication network $G(\mathcal{V}, \mathcal{E})$, where $\mathcal{V} = \{1, \dots, n\}$ denotes the set of nodes and \mathcal{E} denotes the set of edges. An unknown subset of agents, denoted by \mathcal{A} , are *adversarial*. This set is chosen by the system adversary and $|\mathcal{A}| \leq b$. For ease of exposition, let $\phi \triangleq |\mathcal{A}|$. Clearly, $\phi \leq b$. Additionally, we assume $n \geq 3b + 1$.

Let \mathcal{N}_i be the set of incoming neighbors of agent i with $d_i \triangleq |\mathcal{N}_i|$. Each agent can send messages to itself. However, for convenience, we *exclude self-loops* from set \mathcal{E} , i.e., $(i, i) \notin \mathcal{E}$. The specific network structure under consideration is summarized later in Assumption 1, which ensures Byzantine-resilient consensus.

B. Fault Model

We use Byzantine fault model [4] to capture the system threat. Byzantine fault model assumes that there exists a system adversary that can choose a subset of agents \mathcal{A} to compromise and control. While \mathcal{A} is unknown to the good agents, a standard assumption is that the value of b is common knowledge [4]. We refer to an agent suffering Byzantine fault as Byzantine agent.

The system adversary is powerful in that it has complete knowledge of the network, including the network structures, the local program at each good agent, the current status and running history of the multi-agent network. The Byzantine agents can collude with each other and deviate from their pre-specified local programs to arbitrarily misrepresent information to the good agents [4]. In particular, Byzantine agents can mislead the good agents by sending possibly inconsistent messages: letting $m_{ij}[t]$ be the message sent from a Byzantine agent $i \in \mathcal{A}$ to a good agent $j \in \mathcal{V} \setminus \mathcal{A}$ at iteration t , it is possible that $m_{ij}[t] \neq m_{ij'}[t]$ for $j \neq j' \in \mathcal{V} \setminus \mathcal{A}$.

Due to the freedom given to Byzantine agents as well as the system asymmetry caused by them, dealing with Byzantine faults is challenging. In particular, it is well-known that even for complete graphs and scalar local inputs, no Byzantine-resilient consensus algorithms can tolerate more than $\frac{1}{3}$ of the agents to be Byzantine [33].

C. Byzantine-Resilient Multi-Agent Optimization

Let $\mathcal{X} \subseteq \mathbb{R}$ be a nonempty, closed and convex set. We say a function $f : \mathcal{X} \rightarrow \mathbb{R}$ is *admissible* if (i) f is convex, continuously differentiable, and has L -Lipschitz gradient, (ii) f has bounded gradient, i.e., $|f'(x)| \leq L$ for each $x \in \mathcal{X}$, and (iii) $\arg \min_{x \in \mathcal{X}} f(x)$ is non-empty and compact (i.e., bounded and closed). Notably, the

bounded gradient assumption holds automatically when \mathcal{X} is compact. In our multi-agent optimization problem, we assume each agent i has an *admissible* local cost function f_i . Similar assumptions are adopted in [36].

In fault-free networks (i.e., $b = 0$) one commonly adopted global objective [1], [37], [38] is

$$\frac{1}{n} \sum_{i=1}^n f_i(x). \quad (1)$$

Unfortunately, when $b > 0$, the global objective (1) cannot be minimized. This is because initially the function f_i is the local information known to agent i only, and a Byzantine agent can lie arbitrarily about f_i . Observing this, we formulate a robust multi-agent optimization as follows: For a given tuple (β, γ) , where $\beta \in (0, 1)$ and $\gamma \in \mathbb{N}$, the common goal of the good agents is to identify an output such that

$$\tilde{x} \in \arg \min_{x \in \mathcal{X}} \sum_{i \in \mathcal{V} \setminus \mathcal{A}} \alpha_i f_i(x) \quad (2)$$

where

$$\alpha_i \geq 0, \quad \sum_{i \in \mathcal{V} \setminus \mathcal{A}} \alpha_i = 1, \quad \text{and} \\ \sum_{i \in \mathcal{V} \setminus \mathcal{A}} \mathbf{1}\{\alpha_i \geq \beta\} \geq \gamma.$$

Here $\mathbf{1}\{\alpha_i \geq \beta\} = 1$ if $\alpha_i \geq \beta$, and $\mathbf{1}\{\alpha_i \geq \beta\} = 0$ otherwise. Notably, the above formulation is valid even for the general setting where $\mathcal{X} \subseteq \mathbb{R}^d$ for $d \geq 1$. A similar formulation is *independently* proposed in [39, Eq. (4) on page 12] under a statistical learning setup.

Characterizing the structure of the convex coefficients is important despite the fact that different applications favor different structures on the convex coefficients. In this paper, we require that sufficiently many coefficients be not too small, i.e., are nontrivially lower bounded away from 0. Intuitively, if the coefficient assigned to a local cost function is not too small, then the information encoded is utilized in determining the global decision. Such convex coefficients structure is crucial especially for the setting where f_i s are not *i.i.d.* generated.

Example 1. Suppose that $n = 4$, $\mathcal{A} = \{4\}$, $\mathcal{X} = [-1, 1]$, $f_1(x) = f_4(x) = \frac{1}{2}x^2$,

$$f_2(x) = \begin{cases} \frac{(x-1)^2}{2}, & x > 1; \\ 0, & 0 \leq x \leq 1; \\ \frac{1}{2}x^2, & x < 0, \end{cases} \quad f_3(x) = \begin{cases} \frac{1}{2}x^2, & x > 0; \\ 0, & -1 \leq x \leq 0; \\ \frac{(x+1)^2}{2}, & x < -1. \end{cases}$$

If we can solve (2) for $\gamma \geq 2$ and $\beta > 0$, then $\tilde{x} = 0$.

In Example 1, no matter how the malicious agent 4 lies, eventually each of the good agent agrees on 0. Notably, since the malicious agent 4 can lie arbitrarily about its local function. In particular, it could pretend its

local function is $\frac{1}{2}(x - 100)^2$. The good agents cannot tell whether agent 4 is lying or it is honest but its local function is indeed $\frac{1}{2}(x - 100)^2$. To prevent the malicious agent from arbitrarily controlling the eventual agreement value, in the problem formulation (2), the summation is taken over the good agents only.

Example 2 (Vector example). *Consider the distributed state estimation problem wherein each good agent gets access to some local linear measurement y_i of the true state $\theta^* \in \mathbb{R}^d$, i.e., $y_i = \langle w_i, \theta^* \rangle$, where $w_i \in \mathbb{R}^d$ is the local observation vector. Suppose that $n = 4$, $\mathcal{A} = \{4\}$, $\theta^* \in \mathbb{R}^2$, and the local observation vectors of the good agents are $w_1^\top = [1, 0]$, $w_2^\top = [0, 1]$, $w_3^\top = [1, 1]$, and $w_4^\top = [1, 1]$. Suppose we can solve (2) for $\gamma \geq 2$ and $\beta > 0$. For any such convex coefficients α_i s, it is easy to see that $\sum_{i=1}^3 \alpha_i w_i w_i^\top$ is of full rank. Thus, $\tilde{x} = \theta^*$.*

We use the following notion to describe the structure of the convex coefficients.

Definition 1. (β, γ) -**admissibility**: Given $\beta > 0$ and $\gamma \in \mathbb{N}$, $\alpha^\top = [\alpha_1, \dots, \alpha_n]$ is (β, γ) -admissible if:
(1) $\sum_{i \in \mathcal{V} \setminus \mathcal{A}} \alpha_i = 1$, $\alpha_i \geq 0$, for each $i \in \mathcal{V} \setminus \mathcal{A}$; and
(2) at least γ elements of α are lower bounded by β .

In this paper, we take a step towards solving the Byzantine-resilient multi-agent optimization problem in (2) by focusing on scalar local cost functions. Our results might provide insights for general local cost functions.

Clearly, the problem in (2) cannot be solved for all (β, γ) tuples. The following impossibility result was originally presented in [7] (a preliminary conference version) and the unpublished technical report [9].

Theorem 1. *It is impossible to guarantee that more than $|\mathcal{V} \setminus \mathcal{A}| - b$ entries in vector α are non-zero. That is, for any $\beta > 0$, it is impossible to guarantee that α is (β, γ) -admissible with $\gamma > |\mathcal{V} \setminus \mathcal{A}| - b$.*

Proof of Theorem 1 is omitted here for brevity, and can be found in [9]. This proof is based on an indistinguishability argument that is frequently adopted in distributed computing [4]. Specifically, a collection of admissible functions that are hard to be dealt with were constructed. Notably, if additional structures on the local functions are available, the impossibility results in Theorem 1 might no longer hold.

IV. PRELIMINARIES

A. Projection

Let $\text{Dist}(x, \mathcal{X})$ be the distance of x from set \mathcal{X} , i.e.,

$$\text{Dist}(x, \mathcal{X}) = \inf_{y \in \mathcal{X}} |x - y|. \quad (3)$$

We use $P_{\mathcal{X}}[x]$ to denote the projection of the point x onto the set \mathcal{X} , i.e., $P_{\mathcal{X}}[x] = \arg \min_{z \in \mathcal{X}} |z - x|$. Recall

from Section III-C that $\mathcal{X} \subseteq \mathbb{R}$ is a nonempty, closed and convex set. In our analysis, we use the non-expansiveness property of projection, i.e.,

$$|P_{\mathcal{X}}[x] - P_{\mathcal{X}}[y]| \leq |x - y| \quad \forall x, y \in \mathbb{R}. \quad (4)$$

B. Valid Global Objectives

Definition 2. For a given tuple (β, γ) , let $\mathcal{C}(\beta, \gamma)$ be the collection of functions defined as follows:

$$\mathcal{C}(\beta, \gamma) \triangleq \left\{ p : p = \sum_{i \in \mathcal{V} \setminus \mathcal{A}} \alpha_i f_i, \alpha_i \geq 0, \sum_{i \in \mathcal{V} \setminus \mathcal{A}} \alpha_i = 1 \right. \\ \left. \text{and } \sum_{i \in \mathcal{V} \setminus \mathcal{A}} \mathbf{1}\{\alpha_i \geq \beta\} \geq \gamma \right\}. \quad (5)$$

We refer to the functions in $\mathcal{C}(\beta, \gamma)$ as (β, γ) -valid global objectives. Define $X(\beta, \gamma)$ as

$$X(\beta, \gamma) \triangleq \cup_{p \in \mathcal{C}(\beta, \gamma)} \arg \min_{x \in \mathcal{X}} p(x). \quad (6)$$

Lemma 1. If $\beta \leq \frac{1}{|\mathcal{V} \setminus \mathcal{A}|}$ and $\gamma \leq |\mathcal{V} \setminus \mathcal{A}|$, the set $X(\beta, \gamma)$ is convex and closed.

Lemma 1 is proved in our full version [40]. The following quantity is used in proving Lemma 1: Define

$$\bar{X}(\beta, \gamma) \triangleq \cup_{p \in \mathcal{C}(\beta, \gamma)} \arg \min_{x \in \mathbb{R}} p(x). \quad (7)$$

We first show that $\bar{X}(\beta, \gamma)$ is convex and closed. The proof of Lemma 1 then follows easily by carefully examining two cases: (1) $\bar{X}(\beta, \gamma) \cap \mathcal{X} = \emptyset$ and (2) $\bar{X}(\beta, \gamma) \cap \mathcal{X} \neq \emptyset$.

To show an algorithm solves (2) for a given tuple (β, γ) , it is enough to show the convergence of the local estimates at the good agents to set $X(\beta, \gamma)$. Notably, by requiring a common output at the good agents in (2), we implicitly require consensus among the good agents.

C. Byzantine Consensus

In this subsection, we briefly review some relevant definitions and results on Byzantine consensus.

Definition 3. [23] For a given $G(\mathcal{V}, \mathcal{E})$, a reduced graph \mathcal{H} is a subgraph obtained by removing all the Byzantine agents along with their edges; and (ii) by removing any additional up to b incoming edges of each good agent.

Denote the collection of all such reduced graphs by $\mathcal{R}(G)$ with $|\mathcal{R}(G)| \triangleq \tau$.

Definition 4. A source component S of a given graph is the collection of agents each of which has a directed path to every other agent in the graph.

Notably, a graph has at most one source component.

Theorem 2. [23] *Byzantine-resilient consensus with scalar inputs can be achieved on $G(\mathcal{V}, \mathcal{E})$ if and only if every reduced graph has a non-empty source component.*

Throughout this paper, we assume that $G(\mathcal{V}, \mathcal{E})$ satisfies the following assumption so that Byzantine-resilient consensus with scalar inputs can be achieved.

Assumption 1. *Every reduced graph of $G(\mathcal{V}, \mathcal{E})$ has a non-empty source component.*

As stated in Theorem 2, the graph condition in Assumption 1 is both necessary and sufficient for Byzantine-resilient consensus. A similar necessary and sufficient graph condition is assumed by Sundaram and Gharesifard [5] but is stated for the restricted fault model wherein, in each iteration, a malicious agent can send only identically-valued messages to its neighbors. As a result of this difference, the graph condition in Theorem 2 is more restrictive than the graph condition stated in [5]. Note that Assumption 1 immediately implies that $d_i \geq 2b + 1$ for all $i \in \mathcal{V}$ [23]. Otherwise, we might be able to find an isolated node in a reduced graph of $G(\mathcal{V}, \mathcal{E})$, contradicting Assumption 1.

V. BYZANTINE-RESILIENT GRADIENT DESCENT OVER NETWORKS

We present a collaborative gradient descent method, formally described in Algorithm 1 for a good agent; a Byzantine agent may deviate from Algorithm 1 by sending out malicious and possibly inconsistent messages. In Algorithm 1, the good agents exchange both local estimates and local gradients with others. Though exchanging local gradients might not be necessary, it simplifies the exposition of our results.

In Algorithm 1, the maximum number of Byzantine faults b and the sequence of stepsizes $\{\lambda[t]\}_{t=1}^{\infty}$ are known to each good agent a priori. Here, the stepsizes used satisfy the following conditions: (1) $\lambda[t] \geq \lambda[t+1]$ for $t \geq 1$, (2) $\sum_{t=1}^{\infty} \lambda[t] = \infty$, and (3) $\sum_{t=1}^{\infty} \lambda^2[t] < \infty$. In iteration $t \geq 1$, a good agent i sends its local estimate $x_i[t-1]$ and the local gradient $f'_i(x_i[t-1])$ to its outgoing neighbors. Define $\mathcal{D}_i^x[t]$ and $\mathcal{D}_i^g[t]$ as

$$\begin{aligned} \mathcal{D}_i^x[t] &\triangleq \text{set of estimates received from } \mathcal{N}_i, \\ \mathcal{D}_i^g[t] &\triangleq \text{set of received gradients} \cup \{f'_i(x_i[t-1])\}. \end{aligned}$$

Note that the definitions of $\mathcal{D}_i^x[t]$ and $\mathcal{D}_i^g[t]$ are slightly different: $\mathcal{D}_i^x[t]$ does not contain the local estimate, and $\mathcal{D}_i^g[t]$ contains the local gradient. Due to the fact that a Byzantine agent can send differently-valued messages to different outgoing neighbors [4], [21], even in complete graphs it is possible that

$$\mathcal{D}_i^x[t] \neq \mathcal{D}_j^x[t], \text{ and } \mathcal{D}_i^g[t] \neq \mathcal{D}_j^g[t], \quad \forall i \neq j \in \mathcal{V} \setminus \mathcal{A}.$$

Algorithm 1: Byzantine-resilient gradient descent (local program at agent $i \in \mathcal{V} \setminus \mathcal{A}$)

Input: b and $\{\lambda[t]\}_{t=1}^{\infty}$;
Initialization: Set $x_i[0]$ to an arbitrary value in \mathcal{X} .
for $t = 1, \dots$ **do**
 - Compute the local gradient $f'_i(x_i[t-1])$;
 - Send $(x_i[t-1], f'_i(x_i[t-1]))$ to its outgoing neighbors;
 - $\mathcal{D}_i^x[t] \leftarrow$ set of estimates received from incoming neighbors;
 - $\mathcal{D}_i^g[t] \leftarrow$ set of received gradients $\cup \{f'_i(x_i[t-1])\}$;
 - $\tilde{x}_i[t] \leftarrow \text{Average}(\text{Trim}(\mathcal{D}_i^x[t]) \cup \{x_i[t-1]\})$;
 - $\tilde{g}_i[t] \leftarrow \frac{1}{2} (\max\{\text{Trim}(\mathcal{D}_i^g[t])\} + \min\{\text{Trim}(\mathcal{D}_i^g[t])\})$;
 - Compute $x_i[t] \leftarrow P_{\mathcal{X}}[\tilde{x}_i[t] - \lambda[t]\tilde{g}_i[t]]$;
end

Function Trim(\cdot)

Input: Multi-set \mathcal{D} of size $\geq 2b + 1$;
 - Sort the elements in \mathcal{D} in a non-decreasing order (breaking ties arbitrarily);
 - Remove the smallest b elements and the largest b elements;
Output: Trimmed set \mathcal{D} (after $2b$ elements were removed in the previous step).

If agent i does not receive a message from an incoming neighbor in iteration t , it must be true that this incoming neighbor is Byzantine. Thus, a default value is assumed for such missing messages. An algorithmic function, named Trim(\cdot), is used in our gradient descent method, and its description can be found in Algorithm 1. Trim(\cdot) takes a multi-set of size at least $2b + 1$ as input, and removes the b largest and the b smallest elements (breaking ties arbitrarily). Note that in Algorithm 1, the averaging strategies used to compute \tilde{x}_i and \tilde{g}_i are also slightly different. In particular, \tilde{x}_i is the average of $d_i - 2b + 1$ elements in $\text{Trim}(\mathcal{D}_i^x[t]) \cup \{x_i[t-1]\}$, and \tilde{g}_i is the average of the max and min elements in $\text{Trim}(\mathcal{D}_i^g[t])$, where the local gradient may possibly be trimmed away. As can be seen later, these averaging strategies ensure some desired structures on \tilde{x}_i and \tilde{g}_i , respectively.

With the computed $\tilde{x}_i[t]$ and $\tilde{g}_i[t]$, the local estimate x_i is updated via a projected gradient descent update in the last step of the **for** loop in Algorithm 1:

$$x_i[t] \leftarrow P_{\mathcal{X}}[\tilde{x}_i[t] - \lambda[t]\tilde{g}_i[t]].$$

Theorem 3 (main result). *Suppose Assumption 1 holds. Let $|\mathcal{N}_i \cap \mathcal{A}| \triangleq \phi_i$, $\tilde{\gamma} \triangleq \min_{i \in \mathcal{V} \setminus \mathcal{A}} (d_i + 1 - \phi_i - b)$, and $\tilde{\beta} \triangleq \min \left\{ \frac{1}{2 \max_{i \in \mathcal{V} \setminus \mathcal{A}} (d_i + 1 - \phi_i - b)}, \frac{1}{n - \phi} \right\}$. Then $\lim_{t \rightarrow \infty} \text{Dist} \left(x_i[t], X(\tilde{\beta}, \tilde{\gamma}) \right) = 0$ for $i \in \mathcal{V} \setminus \mathcal{A}$.*

Notably, Assumption 1 implies that $\max_{i \in \mathcal{V} \setminus \mathcal{A}} (d_i + 1 - \phi_i - b) > 0$.

VI. MAIN ANALYSIS

In this section, we provide a proof sketch of Theorem 3. All the missing detailed proofs can be found in the full version [40]. The proof of Theorem 3 consists of two parts: achievability of consensus (Section VI-A) and convergence to set $X(\tilde{\beta}, \tilde{\gamma})$ (Section VI-B).

A. Achievability of Consensus

Define the projection error as

$$e_i[t] = P_{\mathcal{X}} [\tilde{x}_i[t] - \lambda[t] \tilde{g}_i[t]] - (\tilde{x}_i[t] - \lambda[t] \tilde{g}_i[t]). \quad (8)$$

Thus, the update of x_i in Algorithm 1 can be written as

$$x_i[t] = \tilde{x}_i[t] - \lambda[t] \tilde{g}_i[t] + e_i[t]. \quad (9)$$

It turns out that the projection error diminishes over time.

Proposition 1. *For each $i \in \mathcal{V} \setminus \mathcal{A}$ and each $t \geq 1$, the projection error $e_i[t]$ satisfies $|e_i[t]| \leq \lambda[t] L$.*

Without loss of generality, assume that agents indexed from 1 through $n - \phi$ are good agents, and agents indexed from $n - \phi + 1$ to n are Byzantine faulty. Let $\mathbf{x} \in \mathbb{R}^{n - \phi}$ be the vector that stacks the local estimates of the good agents with $\mathbf{x}[t]$ being these estimates at the end of iteration t , and $\mathbf{x}_i[t] = x_i[t]$ for $i \in \mathcal{V} \setminus \mathcal{A}$. Similarly, let $\tilde{\mathbf{g}}[t] \in \mathbb{R}^{n - \phi}$ be the vector that stacks the aggregated gradients adopted by the good agents at iteration t (i.e., $\tilde{g}_i[t], \forall i \in \mathcal{V} \setminus \mathcal{A}$), and let $\mathbf{e}[t] \in \mathbb{R}^{n - \phi}$ be the vector that stacks the projection errors (defined in (8)) at the good agents at iteration t .

It was shown in [23] that when $G(\mathcal{V}, \mathcal{E})$ satisfies the condition in Assumption 1, the update of $\mathbf{x} \in \mathbb{R}^{n - \phi}$ can be written compactly in a matrix form:

$$\mathbf{x}[t] = \mathbf{M}[t] \mathbf{x}[t - 1] - \lambda[t] \tilde{\mathbf{g}}[t] + \mathbf{e}[t], \quad (10)$$

where $\mathbf{M}[t]$ is a row-stochastic matrix, and its dependency on t arises from the fact that the Byzantine agents can behave differently from iterations to iterations. In particular, $\tilde{x}_i[t]$ can be written as a convex combination of $x_j[t - 1]$ s at agent i 's good incoming neighbors, i.e.,

$$\tilde{x}_i[t] = \sum_{j=1}^{n - \phi} M_{ij}[t] x_j[t - 1]. \quad (11)$$

The adjacency matrix \mathbf{H} of a reduced graph $\mathcal{H} \in \mathcal{R}(G)$ is defined as

$$\mathbf{H}_{ij} = \begin{cases} 1, & \text{if } i = j \text{ or edge } (i, j) \text{ is contained in } \mathcal{H}; \\ 0, & \text{otherwise.} \end{cases}$$

Henceforth, with a little abuse of terminology and notation, we do not distinguish a reduced graph and its adjacency matrix. The matrix $\mathbf{M}[t]$ in (11) satisfies: for every $t \geq 1$, there exists a reduced graph $\mathcal{H}[t]$ such that

$$\mathbf{M}[t] \geq \xi \mathbf{H}[t], \quad (12)$$

where $\xi \triangleq \frac{1}{2(d_{\max} + 1 - 2b)}$, and $d_{\max} \triangleq \max_{i \in \mathcal{V}} d_i$.

Let $\Phi(t, r) \triangleq \mathbf{M}[t] \mathbf{M}[t - 1] \dots \mathbf{M}[r]$ be a backward product, $\Phi(t, t) \triangleq \mathbf{M}[t]$ and $\Phi(t, t + 1) \triangleq \mathbf{I}$. Equation (10) can be expanded out as

$$\begin{aligned} \mathbf{x}[t] &= \Phi(t, 1) \mathbf{x}[0] - \sum_{r=1}^t \lambda[r] \Phi(t, r + 1) \tilde{\mathbf{g}}[r] \\ &\quad + \sum_{r=1}^t \Phi(t, r + 1) \mathbf{e}[r]. \end{aligned} \quad (13)$$

Proposition 2. [15, Proof of Theorem 1] *Suppose Assumption 1 holds. Then for $i, j, \ell \in \mathcal{V} \setminus \mathcal{A}$, and $r \leq t$, it holds that $|\Phi_{i\ell}(t, r) - \Phi_{j\ell}(t, r)| \leq \theta^{\lceil \frac{t-r+1}{\nu} \rceil}$, where $\nu \triangleq \tau(n - \phi)$ and $\theta \triangleq 1 - \xi^\nu$.*

Recall that $|\mathcal{R}(G)| = \tau$ and $|\mathcal{A}| = \phi \leq b$.

Lemma 2. *For all $i, j \in \mathcal{V} \setminus \mathcal{A}$ and for each $t \geq 1$,*

$$\begin{aligned} |x_i[t] - x_j[t]| &\leq (n - \phi) \theta^{\lceil \frac{t}{\nu} \rceil} \max_{\ell \in \mathcal{V} \setminus \mathcal{A}} |x_\ell[0]| \\ &\quad + 2L(n - \phi) \sum_{r=1}^t \lambda[r] \theta^{\lceil \frac{t-r}{\nu} \rceil}. \end{aligned}$$

The second term in the above upper bound concerns the convolution of an exponentially convergent sequence with a convergent sequence. It is well-known that

$$L(n - \phi) \sum_{r=1}^t \lambda[r] \theta^{\lceil \frac{t-r}{\nu} \rceil} \rightarrow 0 \text{ as } t \rightarrow \infty.$$

A proof was given in our technical report [11] for completeness. Thus, an immediate consequence of Lemma 2 is the achievability of asymptotic consensus.

Corollary 1. *For all $i, j \in \mathcal{V} \setminus \mathcal{A}$, $|x_i[t] - x_j[t]| \rightarrow 0$, as $t \rightarrow \infty$.*

In addition, a convergence rate of $|x_i[t] - x_j[t]|$ for $i, j \in \mathcal{V} \setminus \mathcal{A}$ can be obtained by plugging in the specific choice of stepsizes $\{\lambda[t]\}_{t=1}^\infty$, such as $\lambda[t] = \frac{1}{t} \forall t \geq 1$.

B. Convergence Analysis

Lemma 3. At each iteration t , for each $i \in \mathcal{V} \setminus \mathcal{A}$, there exists a valid function $p^{it} = \sum_{j \in \mathcal{V} \setminus \mathcal{A}} \alpha_j^{it} f_j$ in $\mathcal{C}(\tilde{\beta}, \tilde{\gamma})$, where $\{\alpha_j^{it}\}_{j \in \mathcal{V} \setminus \mathcal{A}}$ are the $(\tilde{\beta}, \tilde{\gamma})$ -admissible convex coefficients, such that

$$\tilde{g}_i[t] = \sum_{j \in \mathcal{V} \setminus \mathcal{A}} \alpha_j^{it} f'_j(x_j[t-1]). \quad (14)$$

In (14), the local cost functions are evaluated at different estimates. Moreover, the valid function involved in Equation (14) is time-varying, and the local gradients at different good agent i s might correspond to different valid functions p s. This is because that the Byzantine agents might behave maliciously.

To show Algorithm 1 solves (2), it is enough to show the local estimates at the good agents converge to set $X(\tilde{\beta}, \tilde{\gamma})$. Towards this goal, we introduce an auxiliary estimate sequence $\{z[t]\}_{t=1}^\infty$, defined as follows:

$$z[t] = x_{j_t}[t], \text{ where } j_t \in \arg \max_{j \in \mathcal{V} \setminus \mathcal{A}} \text{Dist}(x_j[t], X(\tilde{\beta}, \tilde{\gamma})).$$

By this definition, there is a sequence of agents $\{j_t\}_{t=1}^\infty$ associated with the auxiliary sequence $\{z[t]\}_{t=1}^\infty$.

Lemma 4. The auxiliary estimate sequence $\{z[t]\}_{t=1}^\infty$ is asymptotically trapped in set $X(\tilde{\beta}, \tilde{\gamma})$, i.e., $\lim_{t \rightarrow \infty} \text{Dist}(z[t], X(\tilde{\beta}, \tilde{\gamma})) = 0$.

Lemma 4 immediately implies Theorem 3.

Next we provide a high level sketch of the proof of Lemma 4. It can be shown that

$$\begin{aligned} & \text{Dist}(z[t], X(\tilde{\beta}, \tilde{\gamma})) \\ & \leq \max_{i \in \mathcal{V} \setminus \mathcal{A}} \text{Dist}(x_i[t-1] - \lambda[t] \tilde{g}_{j_t}[t], X(\tilde{\beta}, \tilde{\gamma})). \end{aligned}$$

By Lemma 3, there exists a valid global objective function $p^{j_t t} = \sum_{k \in \mathcal{V} \setminus \mathcal{A}} \alpha_k f_k \in \mathcal{C}(\tilde{\beta}, \tilde{\gamma})$ such that

$$\tilde{g}_{j_t}[t] = \sum_{k \in \mathcal{V} \setminus \mathcal{A}} \alpha_k f'_k(x_k[t-1]). \quad (15)$$

Let $j'_t \in \arg \max_{i \in \mathcal{V} \setminus \mathcal{A}} \text{Dist}(x_i[t-1] - \lambda[t] \tilde{g}_{j_t}[t], X(\tilde{\beta}, \tilde{\gamma}))$. We have

$$\begin{aligned} & \text{Dist}(z[t], X(\tilde{\beta}, \tilde{\gamma})) \\ & \leq \inf_{y \in X(\tilde{\beta}, \tilde{\gamma})} \left| x_{j'_t}[t-1] - \lambda[t] (p^{j_t t})'(x_{j'_t}[t-1]) - y \right| \\ & \quad + \lambda[t] L \max_{i,j \in \mathcal{V} \setminus \mathcal{A}} |x_i[t-1] - x_j[t-1]|. \end{aligned} \quad (16)$$

By Lemma 2, we know the second term in the right-hand side of (16) is controllable. It remains to push a recursion out of (16). Towards that, we need to consider several different cases. In particular, to ease exposition, a notion of resilient points is introduced. Detailed proof can be found in our full version [40].

VII. DISCUSSION

In this paper, we studied multi-agent optimization over adversary-prone networks. We introduced a metric, named (β, γ) -admissibility, to quantify the “quality” of a global objective. We took a step towards solving the proposed Byzantine-resilient multi-agent optimization problem by focusing on scalar local cost functions. Unfortunately, even for this simple setting, an impossibility result (i.e., Theorem 1) was shown. We presented an algorithm that is provably robust to Byzantine faults for a certain range of (β, γ) assuming that the communication network satisfies Assumption 1, which ensures Byzantine-resilient consensus can be achieved. For the special case when the network is complete, our results achieved the optimality implied by the impossibility result in Theorem 1, summarized in Fig. 1.

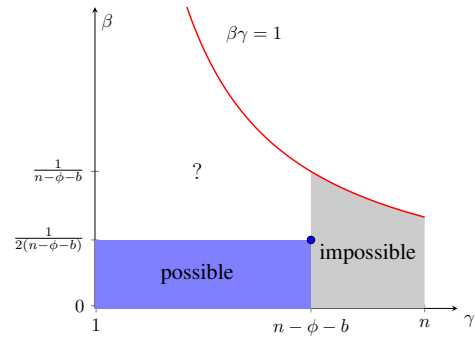


Fig. 1: Achievability of (β, γ) for complete graphs

Open problems:

- In this paper, we assumed Assumption 1. It would be interesting to investigate the minimal conditions on the communication graphs, and to study the trade-off between the (β, γ) -admissibility and the graph structures.
- We focused on asymptotic convergence. A finite-time convergence rate would be of highly practical interests.
- In our algorithm, for ease of analysis, both local estimates and local gradients are exchanged. Directly exchanging gradients might not be necessary because this gradient information, in a sense, is contained in agents' local estimates. It would be interesting to see whether it is possible to relax the requirement on gradient exchanges, and to quantify the gain in exchanging gradients, if any.

REFERENCES

- [1] A. Nedic and A. Ozdaglar, “Distributed subgradient methods for multi-agent optimization,” *Automatic Control, IEEE Transactions on*, vol. 54, no. 1, pp. 48–61, Jan 2009. 1, 2, 3

- [2] S. Shahrampour and A. Jadbabaie, "Distributed online optimization in dynamic environments using mirror descent," *IEEE Transactions on Automatic Control*, vol. 63, no. 3, pp. 714–725, 2018. 1
- [3] C. A. Uribe, S. Lee, A. Gasniov, and A. Nedić, "Optimal algorithms for distributed optimization," *arXiv preprint arXiv:1712.00232*, 2017. 1
- [4] N. A. Lynch, *Distributed Algorithms*. Morgan Kaufmann, 1996. 1, 3, 4, 5
- [5] S. Sundaram and B. Ghahesifard, "Distributed optimization under adversarial nodes," *IEEE Transactions on Automatic Control*, 2018. 1, 2, 5
- [6] Z. Yang and W. U. Bajwa, "Byrdie: Byzantine-resilient distributed coordinate descent for decentralized learning," *arXiv preprint arXiv:1708.08155*, 2017. 1
- [7] L. Su and N. H. Vaidya, "Multi-agent optimization in the presence of byzantine adversaries: Fundamental limits," in *Proceedings of IEEE American Control Conference (ACC)*, July, 2016. 1, 4
- [8] —, "Fault-tolerant multi-agent optimization: Optimal distributed algorithms," in *Proceedings of ACM Symposium on Principles of Distributed Computing*, July, 2016, pp. 425–434. 1, 2
- [9] —, "Byzantine multi-agent optimization: Part I," *arXiv preprint arXiv:1506.04681*, 2015. 1, 2, 4
- [10] —, "Fault-tolerant distributed optimization (Part IV): Constrained optimization with arbitrary directed networks," *arXiv preprint arXiv:1511.01821*, 2015. 2
- [11] —, "Byzantine multi-agent optimization: Part III," *Technical Report*, 2015. 2, 6
- [12] R. Wattenhofer, *Blockchain Science: Distributed Ledger Technology*. Inverted Forest Publishing, 2019. 2
- [13] D. Dolev, N. A. Lynch, S. S. Pinter, E. W. Stark, and W. E. Weihl, "Reaching approximate agreement in the presence of faults," *J. ACM*, vol. 33, no. 3, pp. 499–516, May 1986. [Online]. Available: <http://doi.acm.org/10.1145/5925.5931> 2
- [14] H. J. LeBlanc, H. Zhang, S. Sundaram, and X. Koutsoukos, "Consensus of multi-agent networks in the presence of adversaries using only local information," in *Proceedings of the 1st International Conference on High Confidence Networked Systems*, ser. HiCoNS '12. New York, NY, USA: ACM, 2012, pp. 1–10. [Online]. Available: <http://doi.acm.org/10.1145/2185505.2185507> 2
- [15] N. H. Vaidya, L. Tseng, and G. Liang, "Iterative approximate byzantine consensus in arbitrary directed graphs," in *Proceedings of the 2012 ACM symposium on Principles of distributed computing*. ACM, 2012, pp. 365–374. 2, 6
- [16] R. Friedman, A. Mostefaoui, S. Rajsbaum, and M. Raynal, "Asynchronous agreement and its relation with error-correcting codes," *Computers, IEEE Transactions on*, vol. 56, no. 7, pp. 865–875, 2007. 2
- [17] N. H. Vaidya, "Iterative byzantine vector consensus in incomplete graphs," in *Distributed Computing and Networking*. Springer, 2014, pp. 14–28. 2
- [18] L. Tseng and N. H. Vaidya, "Iterative approximate consensus in the presence of byzantine link failures," in *Networked Systems*. Springer International Publishing, 2014, pp. 84–98. 2
- [19] M. Pease, R. Shostak, and L. Lamport, "Reaching agreement in the presence of faults," *J. ACM*, vol. 27, no. 2, pp. 228–234, Apr. 1980. [Online]. Available: <http://doi.acm.org/10.1145/322186.322188> 2
- [20] M. J. Fischer, N. A. Lynch, and M. S. Paterson, "Impossibility of distributed consensus with one faulty process," *J. ACM*, vol. 32, pp. 374–382, April 1985. [Online]. Available: <http://doi.acm.org/10.1145/3149.214121> 2
- [21] M. J. Fischer, N. A. Lynch, and M. Merritt, "Easy impossibility proofs for distributed consensus problems," in *Proceedings of the fourth annual ACM symposium on Principles of distributed computing*, ser. PODC '85. New York, NY, USA: ACM, 1985, pp. 59–70. [Online]. Available: <http://doi.acm.org/10.1145/323596.323602> 2, 5
- [22] L. Tseng and N. H. Vaidya, "Fault-tolerant consensus in directed graphs," in *Proceedings of the 2015 ACM Symposium on Principles of Distributed Computing*. ACM, 2015, pp. 451–460. 2
- [23] N. H. Vaidya, "Matrix representation of iterative approximate byzantine consensus in directed graphs," *arXiv preprint arXiv:1203.1888*, 2012. 2, 4, 5, 6
- [24] A. D. Fekete, "Asymptotically optimal algorithms for approximate agreement," *Distributed Computing*, vol. 4, no. 1, pp. 9–29, 1990. 2
- [25] D. Fiore and G. Russo, "Resilient consensus for multi-agent systems subject to differential privacy requirements," *Automatica*, vol. 106, pp. 18–26, 2019. 2
- [26] H. Mendes and M. Herlihy, "Multidimensional approximate agreement in byzantine asynchronous systems," in *Proceedings of the Forty-fifth Annual ACM Symposium on Theory of Computing*, ser. STOC '13. New York, NY, USA: ACM, 2013, pp. 391–400. [Online]. Available: <http://doi.acm.org/10.1145/2488608.2488657> 2
- [27] N. H. Vaidya and V. K. Garg, "Byzantine vector consensus in complete graphs," in *Proceedings of the 2013 ACM symposium on Principles of distributed computing*. ACM, 2013, pp. 65–73. 2
- [28] J. N. Tsitsiklis, "Problems in decentralized decision making and computation," DTIC Document, Tech. Rep., 1984. 2
- [29] J. Tsitsiklis, D. Bertsekas, and M. Athans, "Distributed asynchronous deterministic and stochastic gradient optimization algorithms," *Automatic Control, IEEE Transactions on*, vol. 31, no. 9, pp. 803–812, Sep 1986. 2
- [30] S. Boyd, N. Parikh, E. Chu, B. Peleato, J. Eckstein *et al.*, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Foundations and Trends® in Machine learning*, vol. 3, no. 1, pp. 1–122, 2011. 2
- [31] A. Nedić and J. Liu, "Distributed optimization for control," *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 1, no. 1, pp. 77–103, 2018. [Online]. Available: <https://doi.org/10.1146/annurev-control-060117-105131> 2
- [32] L. Su and N. H. Vaidya, "Byzantine multi-agent optimization: Part II," *CoRR*, vol. abs/1507.01845, 2015. [Online]. Available: <http://arxiv.org/abs/1507.01845> 2
- [33] L. Lamport, R. Shostak, and M. Pease, "The byzantine generals problem," *ACM Transactions on Programming Languages and Systems (TOPLAS)*, vol. 4, no. 3, pp. 382–401, 1982. 2, 3
- [34] Y. Chen, L. Su, and J. Xu, "Distributed statistical machine learning in adversarial settings: Byzantine gradient descent," *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, vol. 1, no. 2, p. 44, 2017. 2
- [35] D. Yin, Y. Chen, K. Ramchandran, and P. Bartlett, "Byzantine-robust distributed learning: Towards optimal statistical rates," *arXiv preprint arXiv:1803.01498*, 2018. 2
- [36] A. Nedic, A. Olshevsky, A. Ozdaglar, and J. N. Tsitsiklis, "On distributed averaging algorithms and quantization effects," *IEEE Transactions on Automatic Control*, vol. 54, no. 11, pp. 2506–2517, 2009. 3
- [37] A. Nedic and A. Olshevsky, "Distributed optimization over time-varying directed graphs," *Automatic Control, IEEE Transactions on*, vol. 60, no. 3, pp. 601–615, 2015. 3
- [38] K. I. Tsianos, S. Lawlor, and M. G. Rabbat, "Push-sum distributed dual averaging for convex optimization," in *Decision and Control (CDC), 2012 IEEE 51st Annual Conference on*, Dec 2012, pp. 5453–5458. 3
- [39] M. Charikar, J. Steinhardt, and G. Valiant, "Learning from untrusted data," in *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing*. ACM, 2017, pp. 47–60. 3
- [40] L. Su and N. H. Vaidya, "Byzantine-resilient multi-agent optimization," *Technical report*. [Online]. Available: <https://sites.google.com/site/lilisuece/home/byzantine-resilient-multi-agent-optimization> 4, 6, 7