DATA INJECTION ATTACK ON DECENTRALIZED OPTIMIZATION

Sissi Xiaoxiao Wu[†], Hoi-To Wai[‡], Anna Scaglione[‡], Angelia Nedić[‡], Amir Leshem^{*}

[†]Shenzhen University, College of Information Engineering, Shenzhen, China [‡]School of ECEE, Arizona State Univ., Tempe, AZ, USA *Faculty of Engg., Bar-Ilan Univ., Israel.

ABSTRACT

This paper studies the security aspect of gossip-based decentralized optimization algorithms for multi agent systems against data injection attacks. Our contributions are two-fold. First, we show that the popular distributed projected gradient method (by *Nedić et al.*) can be attacked by *coordinated insider* attacks, in which the attackers are able to steer the final state to a point of their choosing. Second, we propose a metric that can be computed locally by the trustworthy agents processing their own iterates and those of their neighboring agents. This metric can be used by the trustworthy agents to detect and localize the attackers. We conclude the paper by supporting our findings with numerical experiments.

Index Terms— Decentralized optimization, gossip algorithms, data injection attack.

1. INTRODUCTION

Decentralized multi-agent optimization has made significant strides in the past ten years. There is a large literature on multi-agent optimization algorithms, often referred as gossip-based or network diffusion algorithms, that rely on local computations and near neighbors communications to solve iteratively a wide class of constrained optimization problems, e.g., [1-14]. A key advantage of these parallel computation algorithms is the built-in fault tolerance to intermittent computation or communication due to normal failures, as the agents involved can reorganize themselves automatically sailing through these failures. There has been steady progress in expanding the class of problems amenable to a decentralized solution, and a lot of efforts have been made in improving their convergence rate and communication requirements. However, the issue of securing these algorithms against malicious data injection attacks has not received much attention until very recently [15-19]. Naturally, to prevent interference from unauthorized nodes, one can resort to authentication and encryption (see e.g. [20, 21]). However, in the case of an insider attack, gossip-based algorithms are highly vulnerable, even if only one node is compromised. In fact, it is relatively easy to show that gossip-based algorithms are vulnerable to data injection attacks and that such attacks, if coordinated, can steer the network to a final result of the attackers choosing [18, 22]. The flat, self-organizing architecture, which is the selling feature for these algorithms, during an attack becomes a liability.

This paper studies the effect of insider attacks on decentralized optimization algorithms. In particular, we focus on a classical distributed projected gradient (DPG) method introduced by *Nedić* et al. [5]. To this end, we first propose a new, stronger attack model that cannot be detected using the state-of-the-art protection method.

This attack model is shown analytically to be always successful even when the underlying communication graph is time varying. We then propose a new metric that is locally computable by the trustworthy agents to detect and localize the attackers. The latter is shown analytically and empirically to perform the detection and localization tasks successfully. Notice that in [18], the authors proposed a robust consensus based distributed optimization algorithm that is guaranteed to converge to the convex hull of the union of the sets of minimizers of the normal nodes' objectives, which may not include a global optimum in general. On the other hand, with a similar algorithm, [19] proves that an optimal solution can be found under the assumption that the local functions admit a common optimal solution. Both papers did not discuss the task of detection and localization of inside attackers.

2. DISTRIBUTED PROJECTED GRADIENT

We study consensus-based optimization algorithms for tackling the following problem on an *n*-agents network:

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} f(\boldsymbol{\theta}) \coloneqq (1/n) \sum_{i=1}^n f_i(\boldsymbol{\theta}) \text{ s.t. } \boldsymbol{\theta} \in \mathcal{C} , \qquad (1)$$

where $C \subseteq \mathbb{R}^d$ is a closed convex set and $f_i : \mathbb{R}^d \to \mathbb{R}$ is a differentiable function over C. Here, f_i is a *private* function such that it is only known to the *i*th agent, e.g., it may correspond to the measurements made by the sensors of the *i*th agent. We let f^* be the optimal value of problem (1). At time $t \in \mathbb{N}$, the *n* agents are connected via an undirected time varying graph G(t) = (V, E(t)) where V = $[n] := \{1, ..., n\}$ and $E(t) \subseteq [n] \times [n]$ is the edge set. The graph is associated with a weighted adjacency matrix $W(t) \in \mathbb{R}^{n \times n}$ where $[W(t)]_{ij} := W_{ij}(t) = 0$ if $(j, i) \notin E(t)$. The union of these graphs is defined as G := (V, E) with $E := \bigcup_{i=1}^{\infty} E(t)$. We assume the following on the sequence of graphs $\{G(t)\}_{t\geq 1}$ and adjacency matrices $\{W(t)\}_{t>1}$:

H1. If
$$(i, j) \in E(t)$$
, then $W_{ij}(t) \ge \eta$ for some $\eta \in (0, 1)$;

H2. There exists $B_0 < \infty$ such that the graph $(V, \bigcup_{\ell=1}^{B_0} E(t+\ell))$ is connected.

H3. For all $t \ge 1$, it holds that (a) $W(t) \ge 0$, $W(t)\mathbf{1} = \mathbf{1}$ and (b) $W^{\top}(t)\mathbf{1} = \mathbf{1}$.

The distributed projected gradient method (DPG) method [5] tackles (1) by performing the recursion:

$$\boldsymbol{\theta}_{i}(t+1) = \mathcal{P}_{\mathcal{C}}\left(\boldsymbol{\theta}_{i}(t) - \gamma(t)\nabla f_{i}\left(\boldsymbol{\theta}_{i}(t)\right)\right),$$

$$\bar{\boldsymbol{\theta}}_{i}(t) = \sum_{j=1}^{n} W_{ij}(t)\boldsymbol{\theta}_{j}(t),$$
(2)

for all $i \in [n]$ and $t \ge 1$, where $\gamma(t) > 0$ is a diminishing step size. For convex problems, it was shown in [5] the DPG method converges to an optimal solution of (1):

This work is supported by the National Natural Science Foundation of China under Grant 009989, the US National Science Foundation EAGER CCF 1553746, NSF CCF-BSF 1714672, and BSF Grant 2016660.

Fact 1. Under H1, H2, H3(a) and H3(b). Suppose that each of f_i is convex, $\|\nabla f_i(\theta)\| \leq C$ for some C and for all $\theta \in C$, and the step size satisfies $\sum_{t=1}^{\infty} \gamma(t) = \infty$, $\sum_{t=1}^{\infty} \gamma^2(t) < \infty$, then for all $i, j \in [n]$ we have

$$\lim_{t \to \infty} f(\boldsymbol{\theta}_i(t)) = f^* \text{ and } \lim_{t \to \infty} \|\boldsymbol{\theta}_i(t) - \boldsymbol{\theta}_j(t)\| = 0.$$
 (3)

For non-convex problems, the same DPG method (2) is shown in [14] to converge to a KKT point of (1) under different assumptions on the sequence of adjacency matrices $\{W(t)\}_{t=1}^{\infty}$.

3. COORDINATED ATTACK ON DPG

We next describe an attack scheme on the DPG method and show that a straightforward attack scheme can be successful under mild assumptions on the network topology. To set up the stage, let us define $A \subseteq V$, $A \neq \emptyset$ as the set of attackers and $N := V \setminus A$ as the set of trustworthy agents.

The attackers' goal is to steer the final state $\lim_{t\to\infty} \theta_i(t)$ to a target state $x \in C$ for all agents in V, instead of converging to a stationary point of (1). To do so, the attackers follow a different update rule than (2), *i.e.*, we have

$$\boldsymbol{\theta}_j(t) = \boldsymbol{x} + \boldsymbol{z}_j(t), \ \forall \ j \in A ;$$
(4)

meanwhile the trustworthy agents, *i.e.*, agent *i* with $i \in N$, apply the same DPG rule in Eq. (2). In the above, $z_j(t)$ is an artificial noise introduced to obfuscate the trustworthy agents into believing that these agents are trustworthy as they appear to be converging. The artificial noise satisfies

H4. For all $j \in A$, the artificial noise $z_j(t)$ vanishes almost surely as $t \to \infty$, i.e., $\lim_{t\to\infty} ||z_j(t)|| =^{a.s.} 0$.

Notice that the attackers' states are not affected by the trustworthy agents' states at any time. Equivalently we can model the time varying graph such that for any $i \in N$ and $j \in A$, we have $(j, i) \notin E(t)$, *i.e.*, there is no information flow from any trustworthy agent to an attacker. For the attack to be successful, we require the following assumption which is mild. Let E(N; t) be the edge set of the subgraph of G(t) with only the nodes in N.

H5. There exists $B_1, B_2 < \infty$ such that for all $t \ge 1$, (a) the composite sub-graph $(N, \cup_{\ell=t+1}^{t+B_1} E(N; \ell))$ is connected; (b) there exists a pair $i \in N$, $j \in A$ with $(i, j) \in E(t) \cup \ldots \cup E(t+B_2-1)$.

We show that:

Proposition 1. Under H1, H3(a), H4 and H5. If the gradient is bounded such that $\|\nabla f_i(\theta)\| \leq M$ for some M and for all $\theta \in C$, and $\gamma(t) \to 0$, then:

$$\lim_{t \to \infty} \max_{i \in N} \left\| \boldsymbol{\theta}_i(t) - \boldsymbol{x} \right\| =^{a.s.} 0.$$
 (5)

The proof is relegated to Appendix 5.1 and it involves analyzing the attacked DPG dynamics using a new result on the product of transition matrices of the time varying graphs involved. Notice that our result can be easily specialized to the case with a static transition matrix. Proposition 1 shows that the attackers can always succeed in steering the trustworthy agents to the attackers' desired vector.

Note that the attack strategy in (4) could be driven by the attacker trying to solve another instance of (1) with the DPG method, where the sum of the objective functions only extends to the attacker set A. If the optimum point of this instance of (1) is unique and given by \boldsymbol{x} , assuming that the sub-graph G[A] satisfies H1 to H3, the deviation of the attackers' DPG iterates, $\boldsymbol{\theta}_j(t)$, from \boldsymbol{x} forms a sequence $\boldsymbol{z}_j(t)$ that satisfies H4.

3.1. Detecting and Locating the Attackers

Our next endeavor is to detect the presence of attackers and to localize them. If successful, we can eliminate the influences from the attackers such that the attackers no longer steer the final state away from the optimal solution to (1). In the following, the proposed scheme requires accruing historical information about the iterates from K different instances of the DPG (2).

Naturally, we shall focus on the case when H1-H5 hold from the previous section such that the attack is successful. Concretely, at the *k*th instance, we apply the DPG to:

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} f^k(\boldsymbol{\theta}) = (1/n) \sum_{i=1}^n f^k_i(\boldsymbol{\theta}) \text{ s.t. } \boldsymbol{\theta} \in \mathcal{C} , \quad (6)$$

and the iterates of DPG dynamics for the above are denoted using $\{\boldsymbol{\theta}_i^k(t)\}_{i=1}^n$. After running the DPG for K different instances, we consider the following difference vector:

$$\boldsymbol{\eta}_m := (1/K) \sum_{k=1}^K \left(\boldsymbol{\theta}_m^k(\infty) - \boldsymbol{\theta}_m^k(0) \right), \tag{7}$$

which is the difference between the initial value and the steady state value held by node, averaged over all K instances. Note that in practice, it is impossible to compute this metric at $t = \infty$, so we just choose a sufficiently large t at which the system approaches the steady state. It is important to note that for agent i, the vector above can be computed locally as long as $m \in \mathcal{N}_i^{\text{in}}$, *i.e.*, m is in the *inneighbor* set of agent i where $\mathcal{N}_i^{\text{in}} := \{j : (j,i) \in E\}$, note that $E = \bigcup_{i=1}^{\infty} E(t)$ is the union of the edge set over an infinite horizon.

Detection Task. We propose to use the vector η_m for the *detection task*, which is a hypothesis test for:

 \mathcal{H}_{0}^{i} —there is no attacker in $\mathcal{N}_{i}^{\mathsf{in}}$, *i.e.*, $A \cap \mathcal{N}_{i}^{\mathsf{in}} = \emptyset$;

$$\mathcal{H}_1^i$$
 —there exists an attacker in \mathcal{N}_i^{in} , *i.e.*, $A \cap \mathcal{N}_i^{in} \neq \emptyset$.

The detection task corresponds to ----

$$\mathcal{D}^{i} := \frac{1}{|\mathcal{N}_{i}^{\text{in}}|} \sum_{m \in \mathcal{N}_{i}^{\text{in}}} \left| \mathbf{1}^{\top} \left(\boldsymbol{\eta}_{m} - \frac{1}{|\mathcal{N}_{i}^{\text{in}}|} \sum_{j \in \mathcal{N}_{i}^{\text{in}}} \boldsymbol{\eta}_{j} \right) \right| \stackrel{\mathcal{H}_{0}^{i}}{\underset{\mathcal{H}_{1}^{i}}{\overset{\mathcal{S}}{\underset{\mathcal{H}_{1}^{i}}{\underset{\mathcal{H}_{1}^{i}}{\overset{\mathcal{S}}{\underset{\mathcal{H}_{1}^{i}}{\overset{\mathcal{S}}{\underset{\mathcal{H}_{1}^{i}}{\overset{\mathcal{S}}{\underset{\mathcal{H}_{1}^{i}}{\overset{\mathcal{S}}{\underset{\mathcal{H}_{1}^{i}}{\overset{\mathcal{S}}{\underset{\mathcal{H}_{1}^{i}}{\overset{\mathcal{S}}{\underset{\mathcal{H}_{1}^{i}}{\underset{\mathcal{H}_{1}^{i}}{\overset{\mathcal{S}}{\underset{\mathcal{H}_{1}^{i}}{\overset{\mathcal{S}}{\underset{\mathcal{H}_{1}^{i}}{\overset{\mathcal{S}}{\underset{\mathcal{H}_{1}^{i}}{\underset{\mathcal{H}_{1}^{i}}{\overset{\mathcal{S}}{\underset{\mathcal{H}_{1}^{i}}{\underset{\mathcal{H}_{1}^{i}}{\underset{\mathcal{H}_{1}^{i}}{\underset{\mathcal{H}_{1}^{i}}{\underset{\mathcal{H}_{1}^{i}}{\underset{\mathcal{H}_{1}^{i}}{\underset{\mathcal{H}_{1}^{i}}{\underset{\mathcal{H}_{1}^{i}}{\underset{\mathcal{H}_{1}}{\underset{\mathcal{H}_{1}^{i}}{\underset{\mathcal{H}_{1}^{i}}{\underset{\mathcal{H}_{1}^{i}}}{\underset{\mathcal{H}_{1}^{i}}{\underset{\mathcal{H}_{1}^{i}}{\underset{\mathcal{H}_{1}^{i}}{\underset{\mathcal{H}_{1}^{i}}{\underset{\mathcal{H}_{1}^{i}}}{\underset{\mathcal{H}_{1}^{i}}{\underset{\mathcal{H}_{1}^{i}}}{\underset{\mathcal{H}_{1}^{i}}{\underset{\mathcal{H}_{1}^{i}}{\underset{\mathcal{H}_{1}^{i}}{\underset{\mathcal{H}_{1}^{i}}}{\underset{\mathcal{H}_{1}^{i}}{\underset{1}}{\underset{\mathcal{H}_{1}^{i}}{\underset{\mathcal{H}_{1}^{i}}}{\underset{1}}}}}}}}}}}} \right}}$$

where $\delta_I > 0$ is a predefined threshold. To explain the intuition behind the test (8), we observe that under both \mathcal{H}_0^i and \mathcal{H}_1^i ,

$$\boldsymbol{\eta}_m - \frac{1}{|\mathcal{N}_i^{\text{in}}|} \sum_{j \in \mathcal{N}_i^{\text{in}}} \boldsymbol{\eta}_j = \left(\frac{1}{|\mathcal{N}_i^{\text{in}}|} \sum_{j \in \mathcal{N}_i^{\text{in}}} \boldsymbol{\theta}_j^k(0)\right) - \boldsymbol{\theta}_m^k(0) , \quad (9)$$

where the equality is due to the fact that $\theta_i^k(\infty) = \theta_j^k(\infty)$ for all i, j regardless of the hypothesis [cf. Fact 1 and Proposition 1]. Now, suppose that the trustworthy agents are *initialized* using a distribution of the same mean $\bar{\theta}$ while the attackers' initial values follow a distribution of the mean \bar{x} , such that $\bar{x} \neq \bar{\theta}$. From (9), it is obvious that $\mathcal{D}^i = 0$ when there is *no attacker* in the neighborhood $\mathcal{N}_i^{\text{in}}$; while $\mathcal{D}^i = \Omega \left(|\mathbf{1}^\top (\bar{x} - \bar{\theta})| \right) \neq 0$ when there is *at least one attacker* in the neighborhood $\mathcal{N}_i^{\text{in}}$.

Formally, to characterize the detection performance, we need the following assumptions on the statistics of the initial values for trustworthy agents and the attackers:

H6. A trustworthy agent $i \in N$ is initialized by $\boldsymbol{\theta}_i^k(0)$, which is a sub-Gaussian random vector with mean $\bar{\boldsymbol{\theta}}$ and covariance $\sigma_{\theta}^2 \boldsymbol{I}$.

H7. An attacker $j \in A$ is initialized by $\theta_j^k(0)$ defined in (4), which is a sub-Gaussian random vector with mean \bar{x} , and covariance $\sigma_x^2 I$, and $z_j^k(0)$ is a sub-Gaussian random vector with zero mean, and covariance $\sigma_z^2 I$.

Under the assumptions above, the following can be established:

Proposition 2. Define
$$\mu_i = \frac{|N \cap \mathcal{N}_i^{\text{in}}|}{|\mathcal{N}_i^{\text{in}}|} \mathbf{1}^\top (\bar{x} - \bar{\theta})$$
 and
$$\sigma_i^2 = \left(\frac{|\mathcal{N}_i^{\text{in}}|^2 - 2|\mathcal{N}_i^{\text{in}}| + |A \cap \mathcal{N}_i^{\text{in}}|}{|\mathcal{N}_i^{\text{in}}|^2} \right) \sigma_z^2 + \frac{|N \cap \mathcal{N}_i^{\text{in}}|^2}{|\mathcal{N}_i^{\text{in}}|^2} (\sigma_x^2 + \sigma_\theta^2)$$

Under the stated assumptions, we have

$$P(\mathcal{D}^{i} > \delta_{I} \mid \mathcal{H}_{0}^{i}) \leq 2|\mathcal{N}_{i}^{\mathsf{in}}| \exp\left(-Kd\frac{\delta_{I}^{2}}{2\sigma_{\theta}^{2}|\mathcal{N}_{i}^{\mathsf{in}}|(|\mathcal{N}_{i}^{\mathsf{in}}|-1)}\right),$$
$$P(\mathcal{D}^{i} > \delta_{I} \mid \mathcal{H}_{1}^{i}) \geq 1 - \exp\left(-Kd\frac{\max\left\{0, -\delta_{I} + |\mu_{i}|\right\}}{2\sigma_{i}^{2}}\right).$$

The proof of Proposition 2 is provided in Appendix 5.3. In general, observe that the performance improves when:

- 1) the number of instances K and dimension d accrued increases;
- 2) $|\mathbf{1}^{\top}(\bar{\boldsymbol{x}}-\bar{\boldsymbol{\theta}})|$ increases and the variances $\sigma_x, \sigma_z, \sigma_{\theta}$ decrease;
- 3) the number of trustworthy agents near agent i increases.

Apparently, as the number of attackers in the network increases, the detection performance will deteriorate.

Localization Task. Suppose that agent i detects the existence or presence of an attacker in his/her neighborhood. Our next focus is on the *localization task*, whose goal is to distinguish between:

$$\mathcal{H}_{0}^{ij} - \text{agent } j \text{ is not an attacker, } i.e., \ j \notin A,$$

$$\mathcal{H}_{1}^{ij} - \text{agent } j \text{ is an attacker, } i.e., \ j \in A,$$
 (10)

for all $j \in \mathcal{N}_i^{\text{in}}$. Similarly, we propose checking the metric:

$$\mathcal{L}^{ij} \coloneqq \left| \mathbf{1}^{\top} \boldsymbol{\eta}_{j} \right| \underset{\mathcal{H}_{0}^{ij}}{\overset{\mathcal{H}_{0}^{ij}}{\overset{\mathcal{H}_{0}^{j}}}{\overset{\mathcal{H}_{0}^{j}}{\overset{\mathcal{H}_{0}^{j}}}{\overset{\mathcal{H}_{0}^{j}}{\overset{\mathcal{H}_{0}^{j}}}{\overset{\mathcal{H}_{0}^{j}}{\overset{\mathcal{H}_{0}^{j}}}{\overset{\mathcal{H}_{0}^{j}}{\overset{\mathcal{H}_{0}^{j}}}{\overset{\mathcal{H}_{0}^{j}}{\overset{\mathcal{H}_{0}^{j}}}{\overset{\mathcal{H}_{0}^{j}}}{\overset{\mathcal{H}_{0}}}{\overset{\mathcal{H}_{0}^{j}}}{\overset{\mathcal{H}_{0}^{j}}}{\overset{\mathcal{H}_{0}^{j}}}{\overset{\mathcal{H}_{0}}}{\overset{\mathcal{H}_{0}^{j}}}{\overset{\mathcal{H}_{0}^{j}}}{\overset{\mathcal{H}_{0}}}{\overset{\mathcal{H}_{0}^{j}}}{\overset{\mathcal{H}_{0}^{j}}}{\overset{\mathcal{H}_{0}^{j}}}{\overset{\mathcal{H}_{0}}}{\overset{\mathcal{H}_{0}^{j}}}{\overset{\mathcal{H}_{0}^{j}}}{\overset{\mathcal{H}_{0}^{j}}}{\overset{\mathcal{H}_{0}^{j}}}{\overset{\mathcal{H}_{0}}}{\overset{$$

for all $j \in \mathcal{N}_i^{\text{in}}$. Consider the same set of assumptions as before – H6, H7 – we observe that under \mathcal{H}_0^{ij} , the metric can be approximated as $\mathcal{L}^{ij} \approx |\mathbf{1}^\top (\bar{\boldsymbol{\theta}} - \bar{\boldsymbol{x}})| > 0$; while under \mathcal{H}_1^{ij} , the metric is approximately zero $\mathcal{L}^{ij} \approx 0$. The performance can be bounded as:

Proposition 3. Under the stated assumptions, we have

$$P(\mathcal{L}^{ij} < \epsilon_I \mid \mathcal{H}_0^{ij}) < \exp\left(-Kd\frac{(\max\left\{0, -\epsilon_I + |\mathbf{1}^\top(\bar{\boldsymbol{\theta}} - \bar{\boldsymbol{x}})|\right\})^2}{2(\sigma_x^2 + \sigma_\theta^2)}\right)$$
$$P(\mathcal{L}^{ij} < \epsilon_I \mid \mathcal{H}_1^{ij}) \ge 1 - 2\exp\left(-Kd\epsilon_I^2/(2\sigma_z^2)\right).$$

The proof of Proposition 3 is provided in Appendix 5.3. We observe that the localization performance improves under the same conditions 1) and 2) for that in detection.

4. SIMULATIONS AND CONCLUSIONS

In this section, we provide numerical experiments to evaluate the performance of the proposed detection and localization methods. We run several trails of the DPG, each of which contains K instances. We calculate the metrics in (8) and (11) based on the accumulated data and test the detection and localization performance. We consider a Manhattan topology with n = 8 normal nodes (c.f., (2)) and one stubborn node (c.f., (4)), as shown in Fig. 3 in [22]. For simplicity, we take an example of the least square problem; i.e.,

$$f^{k}(\boldsymbol{\theta}) = \sum_{i=1}^{n} f_{i}^{k}(\boldsymbol{\theta}) = \sum_{i=1}^{n} |(\boldsymbol{a}_{i}^{k})^{T} \boldsymbol{\theta} - y_{i}^{k}|^{2}, k = 1, ..., K.$$



Fig. 1. ROCs temporal difference detection performance at the neighboring nodes of the attacker.



Fig. 2. ROCs temporal difference localization performance at the neighboring nodes of the attacker.

Herein, f_i^k can be seen as a utility function for instance k. We write the expected transition matrix as $\mathbb{E}[\mathbf{W}(t)] = \mathbf{I} - \frac{1}{2n} \mathbf{\Sigma} + \frac{P+P^{\top}}{2n}$. where $[\mathbf{P}]_{ij} = P_{ij}$ and $\mathbf{\Sigma}$ is defined as a diagonal matrix with $[\mathbf{\Sigma}]_{ii} = \sum_{j=1}^{n} (P_{ij} + P_{ji})$. Therefore, in our setting the randomized gossip-based decentralized protocol is run with probability $P_{ij} = 1/|\mathcal{N}_i^n|$ between node *i* and node *j* at each time *t*, and is terminated at $T_{\infty} = 2000$. The Monte Carlo simulation is run with 300 trials. At each trail under each instance, we set $\mathbf{x}^k \sim \mathcal{U}[-0.5, 0.5]^d$ and $z_j^k(t) \sim \mathcal{U}[-\hat{\lambda}^t, \hat{\lambda}^t]$, where λ is the second largest eigenvalue of $\mathbb{E}[\mathbf{W}(t)]$. Also, we change the specific functions $f_i^k(\boldsymbol{\theta})$ by randomly generating $\mathbf{a}_i^k \sim \mathcal{U}[0.5, 2.5]^d$, $(\boldsymbol{\theta}^\star)^k \sim \mathcal{U}[0, 1]^d$ and then let $y_i^k = (\mathbf{a}_i^k)^T(\boldsymbol{\theta}^\star)^k$ and pass \mathbf{a}_i^k and y_i^k to the decentralized optimization algorithm with the initialization $\boldsymbol{\theta}^k(0) \sim \mathcal{U}[0, 1]^d$. In the simulation, we take d = 2

In Figure 1, we show curves of the attacker detection performance at the nodes next to an attacker. In the legend, K is the number of instances; dim means that among total d dimensions, how many of which is observed to calculate the metric. If dim = 2, we have exactly the same metric as (8) and (11); if dim = 1, we replace the vector 1 in (8) and (11) by the vector [1,0]. We can see that, by employing the temporal difference method, the trustworthy agents next to an attacker can detect the network is under attack. The detection performance improves with K when dim is fixed, and improves with dim when K is fixed. When the attacker is next to the

trustworthy agent, we can detect the attacker when the product of K and dim is 100, but in general localization appears to be more demanding. In Figure 2, we show the localization ROC curves. Notice that the localization task is invoked only when one trustworthy agent has detected an attack in the neighborhood. In the simulations we assume that the neighborhood detection test was completed correctly, which means that the performance shown are optimistic. The numerical results show that the tests improve with K and d. The test is quite reliable when the product of K and dim is 100.

To conclude, in this paper we have studied how to detect the malicious node in a network running an instance of the DPG optimization algorithm. Considering the nature of the problem, we first prove that under a coordinated attack, the attackers can steer the network towards their desired optimum point. We therefore proposed a attack detection and a localization method and studied its performance analytically and by simulations. The efficacy of the method is demonstrated by numerical experiments. Future work include the extension to other distributed optimization algorithms, the analysis of other detection strategies and the Byzantine resilience of our detection methods.

5. APPENDIX

5.1. Proof of Proposition 1

Define the following product of matrices:

$$\boldsymbol{\Phi}(t,s) := \boldsymbol{W}(t)\boldsymbol{W}(t-1)\cdots\boldsymbol{W}(s), \ t \ge s , \qquad (12)$$

and $\Phi_{ij}(t,s) := [\mathbf{\Phi}(t,s)]_{ij}$. For $i \in N$, we observe the inequality:

$$\begin{aligned} \|\boldsymbol{\theta}_{i}(t+1) - \boldsymbol{x}\| &= \left\| \mathcal{P}_{\mathcal{C}} \left(\bar{\boldsymbol{\theta}}_{i}(t) - \gamma(t) \nabla f_{i} \left(\bar{\boldsymbol{\theta}}_{i}(t) \right) \right) - \boldsymbol{x} \right\| \\ &\leq \left\| \bar{\boldsymbol{\theta}}_{i}(t) - \boldsymbol{x} \right\| + \gamma(t) \| \nabla f_{i}(\bar{\boldsymbol{\theta}}_{i}(t)) \| \\ &\leq \left\| \sum_{j \in N} W_{ij}(t) (\boldsymbol{\theta}_{j}(t) - \boldsymbol{x}) \right\| + \sum_{j \in A} W_{ij}(t) \| \boldsymbol{z}_{j}(t) \| + \gamma(t) M. \end{aligned}$$

Here, the first inequality is due to the projection inequality. Using the fact that $\sum_{j \in A} W_{ij}(t) \leq 1$, the second last term can be bounded by $\overline{z}(t) := \max_{j \in A} \| \boldsymbol{z}_j(t) \|.$

Now define $B := (n-1)B_1 + B_2$, we can proceed with the recursion above to get:

$$\begin{aligned} \|\boldsymbol{\theta}_{i}((k+1)B) - \boldsymbol{x}\| &\leq \Big\| \sum_{j \in N} \Phi_{ij}((k+1)B - 1, kB)(\boldsymbol{\theta}_{j}(kB) - \boldsymbol{x}) \Big\| \\ &+ \sum_{\ell=kB}^{(k+1)B-1} \left(\gamma(\ell)M + \bar{z}(\ell) \right) . \end{aligned}$$
(13)

Using the triangular inequality, it is easy to bound the first term on the right hand side above as:

$$\Big(\sum_{j\in N}\Phi_{ij}((k+1)B-1,kB)\Big)\cdot\max_{j\in N}\big\|\boldsymbol{\theta}_j(kB)-\boldsymbol{x}\big\|.$$

Let us invoke the following lemma, whose proof can be found in Appendix 5.2:

Lemma 1. Under H1, H3(a), H5. It holds that

$$\max_{i \in N} \sum_{j \in N} \phi_{ij}(t+B,t+1) \le 1 - \eta^B < 1, \ \forall \ t \ge 0 \ . \ (14)$$

Substituting the above into (13) and (14) yields the inequality:

$$\max_{i \in N} \|\boldsymbol{\theta}_i((k+1)B) - \boldsymbol{x}\| \le (1 - \eta^B) \cdot \max_{i \in N} \|\boldsymbol{\theta}_i(kB) - \boldsymbol{x}\| + \sum_{\ell=kB}^{(k+1)B-1} (\gamma(\ell)M + \bar{z}(\ell)),$$

for all k > 1. As $B < \infty$, the latter term vanishes as $k \to \infty$. Combining these observations and applying Corollary 3 in [23] implies that $\max_{i \in N} \|\boldsymbol{\theta}_i(kB) - \boldsymbol{x}\| \to 0$. Consequently, (5) holds since for all $s \in [1, B-1], \theta_i(kB+s) - \theta_i(kB)$ can also be bounded by a vanishingly small quantity.

5.2. Proof of Lemma 1

Define $D(t) := [W_{ij}(t)]_{i,j \in N}$ and $\Phi_D(t,s) := [\Phi_{ij}(t,s)]_{i,j \in N}$, *i.e.*, the lower right sub-block of D(t), $\Phi(t, s)$. It can be shown that

$$\boldsymbol{\Phi}_D(t,s) = \boldsymbol{D}(t)\boldsymbol{D}(t-1)\cdots\boldsymbol{D}(s), \ t \ge s \ . \tag{15}$$

We are interested in the matrix-vector product $\Phi_D(t+B,t+1)\mathbf{1}$. Notice that under H1, H5(b), there exists $t^* \in [t + 1, t + B_2]$ such that:

 $[\mathbf{\Phi}_D(t^{\star}, t+1)\mathbf{1}]_{j^{\star}} \leq 1-\eta$ for some $j^{\star} \in N$. (16)Since $t + B - t^* \ge (n - 1)B_1$, under H5(a) and using similar arguments as in the proof of Lemma 2 in [4], we can show

$$\Phi_{ij}(t+B, t^*+1) \ge \eta^{t+B-t^*}, \,\forall \, i, j \in N .$$
(17)

Consequently, for all $i \in N$, we have:

$$[\mathbf{\Phi}_{D}(t+B,t+1)\mathbf{1}]_{i} = \sum_{\ell \in N} \Phi_{i\ell}(t+B,t^{*}+1)[\mathbf{\Phi}_{D}(t^{*},t+1)\mathbf{1}]_{\ell}$$

$$\leq 1 - \eta \cdot \Phi_{i,j^{*}}(t+B,t^{*}+1) \leq 1 - \eta^{B},$$

where the first equality is obtained by simply expanding the matrixmatrix product; and the second inequality is due to (16). This concludes the proof of Lemma 1.

5.3. Proof of Propositions 2 and 3

The proof is similar to that in [22, Theorem 1]. It is easy to check that under \mathcal{H}_0^i , \mathcal{D}_i is a zero mean r.v. with sub-Gaussian parameter $\sigma_{\theta}^2(|\mathcal{N}_i^{\text{in}}| - 1)/(Kd|\mathcal{N}_i^{\text{in}}|)$. Applying the union bound gives

$$\begin{split} & \mathbf{P}(\mathcal{D}^{i} > \delta_{I} \mid \mathcal{H}_{0}^{i}) \\ \leq & |\mathcal{N}_{i}^{\mathsf{in}}| \mathbf{P}\left(\left| \mathbf{1}^{\top} \left(\boldsymbol{\eta}_{m} - \sum_{j \in \mathcal{N}_{i}^{\mathsf{in}}} \frac{\boldsymbol{\eta}_{j}}{|\mathcal{N}_{i}^{\mathsf{in}}|} \right) \right| > \frac{\delta_{I}}{|\mathcal{N}_{i}^{\mathsf{in}}|} \mid \mathcal{H}_{0}^{i} \right), \exists m \in \mathcal{N}_{i}^{\mathsf{in}}, \end{split}$$

By using Hoeffding's inequality [24] on the above, we obtain the

desired result for the false alarm rate. On the other hand, under \mathcal{H}_1^i , we have $\mathbf{1}^{\top}(\boldsymbol{\eta}_m - \frac{1}{|\mathcal{N}_i^{\text{in}}|}\sum_{j \in \mathcal{N}_i^{\text{in}}} \boldsymbol{\eta}_j)$ being a r.v. with mean μ_i and sub-Gaussian parameter $\sigma_i^2/(Kd)$ for $m \in A \cap \mathcal{N}_i^{\text{in}}$. Then,

$$\begin{split} & \mathbf{P}(\mathcal{D}^{i} < \delta_{I} \mid \mathcal{H}_{1}^{i}) \leq \mathbf{P}\left(\left|\boldsymbol{e}_{j}^{\top}\left(\boldsymbol{\eta}_{m} - \sum_{j \in \mathcal{N}_{i}^{\mathsf{in}}} \frac{\boldsymbol{\eta}_{j}}{|\mathcal{N}_{i}^{\mathsf{in}}|}\right)\right| < \delta_{I} \mid \mathcal{H}_{1}^{i}\right), \forall m \in \mathcal{N}_{i}^{\mathsf{in}} \\ \leq & \mathbf{P}\left(\underbrace{\boldsymbol{e}_{j}^{\top}\left(\boldsymbol{\eta}_{m} - \sum_{j \in \mathcal{N}_{i}^{\mathsf{in}}} \frac{\boldsymbol{\eta}_{j}}{|\mathcal{N}_{i}^{\mathsf{in}}|}\right) - \mu_{i}}_{\hat{\eta}} > -\delta_{I} + |\mu_{i}| \mid \mathcal{H}_{1}^{i}\right), \forall m \in \mathcal{N}_{i}^{\mathsf{in}}. \end{split}$$

Apparently, $\hat{\eta}$ is a sub-Gaussian r.v. with mean zero and parameter $\sigma_i^2/(Kd)$. By using Hoeffding's inequality we have the desired lower bound for the detected rate.

For the localization task, we have $\mathbf{1}^{\top} \boldsymbol{\eta}_j$ being a zero mean r.v. with sub-Gaussian parameter $\sigma_z^2/(Kd)$ under \mathcal{H}_1^{ij} , and being a sub-Gaussian r.v. with mean $\mathbf{1}^{\top}(\bar{\boldsymbol{x}}-\bar{\boldsymbol{\theta}})$ and variance $(\sigma_x^2+\sigma_{\theta}^2)/(Kd)$ under \mathcal{H}_0^{ij} . Then, Hoeffding's inequality gives the desired bounds for localization performance. This completes the proofs.

6. REFERENCES

- J. Tsitsiklis, "Problems in decentralized decision making and computation," Ph.D. dissertation, Dept. of Electrical Engineering and Computer Science, M.I.T., Boston, MA, 1984.
- [2] J. Duchi, A. Agarwal, and M. J. Wainwright, "Dual averaging for distributed optimization: Convergence analysis and network scaling," *IEEE Trans. Autom. Control*, vol. 57, no. 3, pp. 592–606, March 2012.
- [3] A. H. Sayed, Adaptation, Learning, and Optimization over Networks. Foundations and Trends in Machine Learning, 2014, vol. 7.
- [4] A. Nedić and A. Ozdaglar, "Distributed subgradient methods for multi-agent optimization," *IEEE Transactions on Automatic Control*, vol. 54, no. 1, pp. 48–61, 2009.
- [5] S. Ram, A. Nedić, and V. Veeravalli, "Distributed stochastic subgradient projection algorithms for convex optimization," *Journal of Optimization Theory and Applications*, vol. 147, no. 3, pp. 516–545, 2010.
- [6] S. Boyd, A. Ghosh, B. Prabhakar, and D. Shah, "Randomized gossip algorithms," *IEEE Trans. Inf. Theory*, vol. 52, no. 6, pp. 2508–2530, Jun. 2006.
- [7] A. G. Dimakis, S. Kar, J. M. F. Moura, M. G. Rabbat, and A. Scaglione, "Gossip Algorithms for Distributed Signal Processing," *Proc. IEEE*, vol. 98, no. 11, pp. 1847–1864, Nov. 2010.
- [8] M. Rabbat and R. Nowak, "Distributed optimization in sensor networks," in *Proceedings of the 3rd International Symposium on Information Processing in Sensor Networks*, ser. IPSN '04. New York, NY, USA: ACM, 2004, pp. 20–27. [Online]. Available: http://doi.acm.org/10.1145/984622.984626
- [9] D. Jakovetic, J. Xavier, and J. M. F. Moura, "Fast distributed gradient methods," *IEEE Trans. Autom. Control*, vol. 59, no. 5, pp. 1131–1146, May 2014.
- [10] E. Wei and A. Ozdaglar, "On the o(1/k) convergence of asynchronous distributed alternating direction method of multipliers," *CoRR*, 2013.
- [11] Y. Yang, G. Scutari, D. P. Palomar, and M. Pesavento, "A parallel stochastic approximation method for nonconvex multiagent optimization problems," *CoRR*, vol. abs/1410.5076, Oct 2014.
- [12] T.-H. Chang, A. Nedic, and A. Scaglione, "Distributed constrained optimization by consensus-based primal-dual perturbation method," *IEEE Trans. Autom. Control*, vol. 59, no. 6, pp. 1524–1538, June 2014.
- [13] H.-T. Wai, J. Lafond, A. Scaglione, and E. Moulines, "Decentralized frank-wolfe algorithm for convex and non-convex problems," *IEEE Transactions on Automatic Control*, 2017.
- [14] P. Bianchi and J. Jakubowicz, "Convergence of a multi-agent projected stochastic gradient algorithm for non-convex optimization," *IEEE Trans. Autom. Control*, vol. 58, no. 2, pp. 391– 405, Feb 2013.
- [15] M. Ozay, I. Esnaola, F. T. Y. Vural, S. R. Kulkarni, and H. V. Poor, "Sparse attack construction and state estimation in the smart grid: Centralized and distributed models," *IEEE Journal* on Selected Areas in Communications, vol. 31, no. 7, pp. 1306– 1318, 2013.

- [16] O. Vuković and G. Dán, "Security of fully distributed power system state estimation: Detection and mitigation of data integrity attacks," *IEEE Journal on Selected Areas in Communications*, vol. 32, no. 7, pp. 1500–1508, July 2014.
- [17] X. Liu, Z. Bao, D. Lu, and Z. Li, "Modeling of local false data injection attacks with reduced network information," *IEEE Transactions on Smart Grid*, vol. 6, no. 4, pp. 1686–1696, 2015.
- [18] S. Sundaram and B. Gharesifard, "Consensus-based distributed optimization with malicious nodes," in 2015 53rd Annual Allerton Conference on Communication, Control, and Computing (Allerton). IEEE, 2015, pp. 244–249.
- [19] L. Su and N. Vaidya, "Byzantine multi-agent optimization: Part II," CoRR, July 2015.
- [20] A. Perrig, R. Szewczyk, J. D. Tygar, V. Wen, and D. E. Culler, "SPINS: security protocols for sensor networks," *Wireless Networks*, vol. 8, no. 5, pp. 521–534, Sep. 2002.
- [21] S. Zhu, S. Setia, and S. Jajodia, "LEAP: efficient security mechanisms for large-scale distributed sensor networks," in *Proc CCS* '03, 2003, pp. 62–72.
- [22] R. Gentz, S. X. Wu, H. T. Wai, A. Scaglione, and A. Leshem, "Data injection attacks in randomized gossiping," *IEEE Transactions on Signal and Information Processing over Networks*, vol. PP, no. 99, pp. 1–1, 2016.
- [23] B. P. Polyak, *Introduction to Optimization*. Optimization Software, Inc., 1987.
- [24] P. Massart, Concentration Inequalities and Model Selection. Springer, 2003.