

LOCALIZATION OF DATA INJECTION ATTACKS ON DISTRIBUTED M-ESTIMATION

Or Shalom*, Amir Leshem*, Anna Scaglione†

*Faculty of Engineering, Bar-Ilan University, Ramat Gan, Israel.

†School of ECEE, Arizona State University, Tempe, AZ, USA.

ABSTRACT

This paper describes a distributed statistical estimation problem, corresponding to a network of agents. The network may be vulnerable to data injection attacks, in which attackers control legitimate nodes in the network and use them to inject false data. We have previously shown [1] that the detection metric by *Wu et. al* in [2], is vulnerable to sophisticated attacks where the attacker mixes normal behaviour and false data injection. In this paper we propose a novel metric that can be computed locally by each agent to detect and localize the novel attack in the network in a single instance.

Index Terms—Distributed projected gradient, Decentralized optimization, Data injection attacks, Convex optimization, M-Estimators

I. INTRODUCTION

DECENTRALIZED multi-agent optimization is an important problem in distributed computation. These algorithms rely on local computations as well as in-neighborhood communication to achieve their common goal of minimizing a common cost function or converging to a stable point. As these networks gain popularity [3–11], it has become apparent that they are sensitive to false data injection which can steer the network’s final state, see [2], [12–22] for examples. The structure of an independently -self updating network, which has been the main advantage of these methods, can turn into a vulnerability by allowing an attacker which controls a single node to have a global impact. This type of attack cannot be detected using cryptographic techniques, since the attacker controls a legitimate node in the network. This paper focuses on the problem of localizing attacks on distributed statistical estimation, using M-estimators [23] in general and maximum-likelihood in particular, using the distributed projected gradient (DPG) algorithm. We begin with a novel data injection attack scheme, and its effects on decentralized optimization algorithms, and primarily DPG [3]. We propose a novel, more sophisticated attack scheme which is invisible to all previous detection methods. This attack scheme is shown to be always successful on communication networks, even when the network is dynamically changing over time. We then propose a new

metric, computed locally by each agent over time, to detect and localize an attacker in the network, allowing the users to ignore the attacker and reach convergence to the true optimal state. In contrast to previously proposed techniques, our scheme can detect and mitigate the attack in a single run of the algorithm.

Notations: We use boldfaced letters to denote vectors and boldfaced uppercase letters to denote matrices. For a vector θ , $[\theta]_i$ denote its i -th element, similarly, for a matrix A , A_{ij} denotes its (i, j) -th element.

II. PROBLEM FORMULATION

Consider a grid of sensors measuring independent random processes that depend on a joint parameter. In order to extract this parameter, the sensors solve a M-Estimator problem

$$\arg \min_{\theta} \sum_i \rho(x_i, \theta) \quad (1)$$

where $\rho(x_i, \theta)$ is the i -th agent’s private objective function and θ is an unknown parameter vector. M-Estimators generalize the maximum-likelihood by replacing the likelihood function with a generalized function, $\rho(x_i, \theta)$, for each user $i \in V$. We assume that the process is i.i.d. between sensors. We consider a distributed setup where agents do not share their private information x_i .

II-A. Preliminaries

Consider an undirected, time varying graph $G(t) = (V, E(t))$ defining a network of N agents, where $V = \{1, \dots, N\}$ is a set of N nodes (agents) and $E(t) \subseteq V \times V$ denotes the connections between the nodes for some time $t \in \mathbb{N}$. For each node i , we define $\mathcal{N}_i \subset V$ as the neighborhood set of agent i , as $\mathcal{N}_i := \{j : (j, i) \in E\}$, note that $E = \cup_{t=1}^{\infty} E(t)$. We mark the i -th agent state for some time $t \geq 0$ as $\theta_i(t)$.

II-B. Distributed stochastic M-estimation

The N agents share the common goal of minimizing a joint objective function in a distributed manner; i.e., solve the following optimization problem:

$$\min_{\theta} h(\theta) := \frac{1}{N} \sum_{i=1}^N h_i(\theta), \quad \text{s.t. } \theta \in \mathcal{C} \quad (2)$$

This work is supported by the NSF CCF-BSF 1714672 and grants ISF-1644/18 and ISF-NRF 2277/16.

where $\mathcal{C} \subseteq \mathbb{R}^P$ is a closed, convex, compact set and $h_i : \mathbb{R}^P \rightarrow \mathbb{R}$, $h_i(\boldsymbol{\theta}) = \rho(\mathbf{x}_i, \boldsymbol{\theta})$ is a private differentiable¹ function over \mathcal{C} , known to the i -th agent alone. In our problem, we assume that \mathbf{x}_i are i.i.d. given any value of $\boldsymbol{\theta}$; i.e., the objective function $\rho(\mathbf{x}, \boldsymbol{\theta})$ is the same function, and is known to all nodes. However the specific realization $\rho(\mathbf{x}_i, \boldsymbol{\theta})$ is private since each node has its own data. We mark the optimal solution of the optimization problem as $h^* = h(\boldsymbol{\theta}^*)$, where $\boldsymbol{\theta}^* \in \mathcal{C}$ is the optimal parameter. In this paper we assume that all $h_i(\boldsymbol{\theta})$ are convex. The M-estimation can be solved distributedly using a stochastic distributed projected gradient algorithm.

Let $G(t)$ be the graph associated with a weighted adjacency matrix, $\mathbf{W}(t) \in \mathbb{R}^{N \times N}$, where $\mathbf{W}(t)$ satisfies:

Assumption 1: $\mathbf{W}(t)$ fulfills the next terms for $t \geq 0$:

- $\mathbf{W}(t)$ is a symmetric, nonnegative, bi-stochastic matrix.
- If $(i, j) \in E(t)$ then $\mathbf{W}_{i,j}(t) \geq \xi$ for some $\xi \in (0, 1)$.
- If $(i, j) \notin E(t)$ then $\mathbf{W}_{i,j}(t) = 0$.

Assumption 2: There exists $B < \infty$ such that the graph $(V, \cup_{l=1}^B E(t+l))$ is connected.

The distributed projected gradient (DPG) method [3] solves the optimization problem shown in (2) by performing the recursion:

$$\begin{aligned} \bar{\boldsymbol{\theta}}_i(t) &= \sum_{j=1}^N \mathbf{W}_{ij}(t) \boldsymbol{\theta}_j(t), \quad \forall i \in V, t \geq 0, \\ \boldsymbol{\theta}_i(t+1) &= \mathcal{P}_{\mathcal{C}}(\bar{\boldsymbol{\theta}}_i(t) - \eta(t) \nabla h_i(\bar{\boldsymbol{\theta}}_i(t))) \end{aligned} \quad (3)$$

where $\mathcal{P}_{\mathcal{C}}$ denotes the euclidean projection onto the set \mathcal{C} and $\eta(t)$ satisfies:

Assumption 3: $\eta(t)$ is a time-varying step size satisfying $\sum_{t=1}^{\infty} \eta(t) = \infty$ and $\sum_{t=1}^{\infty} \eta^2(t) < \infty$.

Proposition 1: Under assumptions 1-3, for a compact space, the joint objective function asymptotically reaches a minimum, as seen in [4], [12].

$$\lim_{t \rightarrow \infty} h(\boldsymbol{\theta}(t)) = h^* \quad (4)$$

Our goal is to detect malicious nodes in the network that attempt to destroy the distributed computation by injecting false data.

III. DATA INJECTION ATTACKS

Consider a distributed M-estimation, where some nodes are malicious and inject false data into the network. We divide the set of nodes, V , into two subsets: $R \subset V$ is the set of reliable agents and $A := V \setminus R$, $A \neq \emptyset$ is the set of attackers. Let $n_a = |A|$ be the number of attacking nodes. The attackers' goal is to steer the network's final state $\lim_{t \rightarrow \infty} \boldsymbol{\theta}(t)$ to a target state of their choice, while remaining transparent to the network. To do so the attackers follow a deceiving update rule of their choice while the trustworthy agents follow the DPG update rule as shown

¹If the objective function is non-differentiable, each gradient reference should be considered a sub-gradient.

in (3). A previous work [2] suggested a straightforward attack scheme, as well as a detection method. Unfortunately, the attack scheme can be modified to evade this detection method. In this section we present a novel improved attack method, and later we propose a combined detection and localization scheme, computed locally by each agent in a single instance of the algorithm.

III-A. Novel Attack Scheme

The new attack scheme proposed here is a mixture of two update rules:

- The trustworthy agents' DPG update rule.
- The straightforward attacker's update rule [2].

To combine both update rules we generate a new time-varying proportion coefficient marked as $g(t)$.

Assumption 4: The new proportion coefficient $g(t)$ fulfills the following conditions:

- For all $t \geq 0$, $0 \leq g(t) \leq 1$.
- $g(t)$ decreases over time, i.e. $g(t+1) < g(t)$.
- $g(0) = 1$, $\lim_{t \rightarrow \infty} g(t) = 0$.

Note that the limitations on $g(t)$ is very minor, we have no assumptions neither on the convergence rate to 0, nor on any relation to other components in the network's convergence process.

The new proposed attack scheme is

$$\begin{aligned} \boldsymbol{\theta}_j(t+1) &= g(t) \times DPG(\boldsymbol{\theta}_j(t)) \\ &\quad + (1 - g(t)) \times (\boldsymbol{\alpha}_0 + \mathbf{z}_j(t+1)), \quad \forall j \in A \end{aligned} \quad (5)$$

where $DPG(\boldsymbol{\theta}_j(t))$ refers to (3), $\boldsymbol{\alpha}_0$ is the attacker's desired final state and $\mathbf{z}_j(t)$ is a zero mean and $\sigma_z^2(t) \mathbf{I}_P$ variance random noise, vanishing a.s. over time and satisfying the expected convergence rate of the graph for all $j \in A$.

The result of implementing the new attack scheme on the network forces the initial state of the attackers' nodes to be similar to that of the trustworthy agents. Therefore, the detection scheme in [2] fails. The network's convergence to the attacker's desired state, under the new attack scheme, is demonstrated in Figure 1. Looking at Figure 1, we can see that the entire convergence process under the proposed attack scheme can be divided into 3 time periods. Prior to the attack ($t < T_g$), during the attack ($T_g \leq t < T_\infty$) and post convergence ($t \geq T_\infty$).

Assumption 5: The objective function gradient, ∇h_i , is bounded for each $i \in V$, s.t.

$$|\nabla h_i(\boldsymbol{\theta})| \leq \frac{\mathbf{q}}{2}, \quad \mathbf{q} = [q_1, \dots, q_P] \quad (6)$$

Proposition 2. Under the previous assumption and the proposed attack scheme in (5), the network converges to the attacker's desired state $\boldsymbol{\alpha}_0$.

$$\lim_{t \rightarrow \infty} \|\boldsymbol{\theta}_i(t) - \boldsymbol{\alpha}_0\|_\infty = 0, \quad \forall i \in V, \quad (7)$$

Proof in [1, Appendix A] .

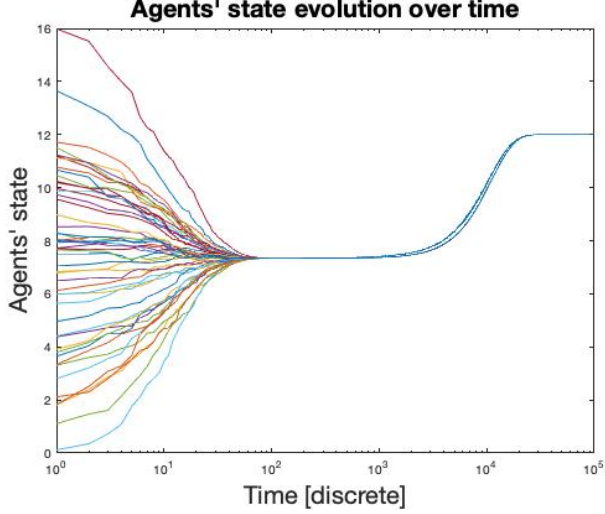


Fig. 1. An example of the novel attack scheme (for $P = 1$). The network reaches convergence to an unstable state that drift over time to α_0 . Note that in this figure a weak attack was drawn in order to emphasize the three periods of the algorithm, generally the attack is controlled by the attackers.

IV. LOCALIZING THE ATTACKERS

In this section we propose a novel low-complexity metric, computed over time by each agent, to detect and localize the attackers instantaneously, running the algorithm for a single instance (as opposed to previous works, including [1]). Once we find the attackers we can ignore their data and have a trustworthy network solving (2), reaching an optimal state. We run the recursive DPG algorithm as seen in (3), where the attackers are following the proposed attack scheme shown in (5). The algorithm runs for some time index marked as T_∞ , sufficient for convergence.

In this method, each agent compares the state updates received from each agent in its neighborhood with the rest of the agents in the neighborhood over time, after reaching convergence. The agents are assumed to be identically distributed and therefore if an agent is malicious and updates differently, it will stand-out and be considered as an outlier. Denote the two hypotheses:

$\mathcal{H}_{i,j}^0$ – Agent $j \in \mathcal{N}_i$ is not an attacker; i.e., $j \notin A$.

$\mathcal{H}_{i,j}^1$ – Agent $j \in \mathcal{N}_i$ is an attacker; i.e., $j \in A$.

The proposed metric, computed over time by each agent is given by

$$\Delta U_{i,j} = \left| \frac{1}{\Delta T} \sum_{\Delta T} U_{i,j}(t) \right|_{\mathcal{H}_{i,j}^0}^{\mathcal{H}_{i,j}^1} \lesseqgtr \delta_u \quad (8)$$

where

$$U_{i,j}(t) = u_{i,j}(t) - \text{median}\{u_{i,l}(t) : l \in \mathcal{N}_i \setminus j\} \quad (9)$$

$$u_{i,j}(t) = \frac{\|\theta_{i,j}(t+1) - \theta_{i,j}(t)\|}{\eta(t)}$$

$\theta_{i,j}$ is the data that agent i receives from agent j and δ_u is a predefined threshold. An example of the proposed detection and localization scheme in a single instance for different integration time, ΔT , can be seen in Figure 2.

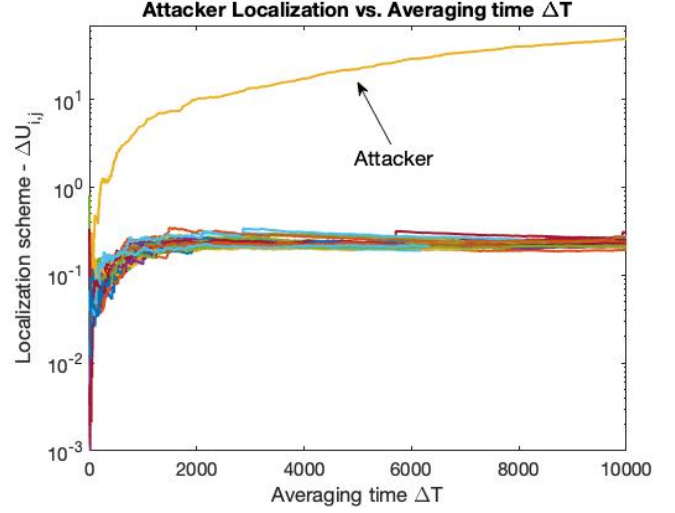


Fig. 2. An example of the new proposed detection and localization scheme, $\Delta U_{i,j}$, computed for different averaging time, ΔT , in a dynamically changing random graph ($P = 1$). It is easy to notice that the attacker is exceptional.

V. LOCALIZATION SCHEME ANALYSIS

In the proposed localization scheme, we look at the tail of the algorithm (post convergence, at times $t \geq T_\infty$). We can show analytically, that the attacker will behave differently than other agents in most cases. Due to the nature of the algorithm, the attacker constantly oppose the trustworthy agents' gradient update and therefore will stand off other agents.

For a trustworthy agent, $i \in R$, and an attacker, $j \in A$, the states after convergence are

$$\begin{aligned} \theta_i(t+1) &= \alpha_0 + \epsilon_i(t) - \eta(t) \nabla h_i(\alpha_0) \\ \theta_j(t+1) &= g(t) (\alpha_0 + \epsilon_j(t) - \eta(t) \nabla h_j(\alpha_0)) \\ &\quad + (1 - g(t)) (\alpha_0 + z_j(t)) \end{aligned} \quad (10)$$

where $\forall i \in V$, ϵ_i is a random noise, vanishing according to the convergence rate of the network. To compute $u_{i,j}(t)$, we have to compute the following

$$\begin{aligned} [\theta_i(t+1) - \theta_i(t)] / \eta(t) &= \hat{\epsilon}_i(t) / \eta(t) - \nabla h_i(\alpha_0) \\ [\theta_j(t+1) - \theta_j(t)] / \eta(t) &= [z_j(t) - \epsilon_j(t-1)] / \eta(t) \end{aligned} \quad (11)$$

where $\hat{\epsilon}_i(t) = \epsilon_i(t) - \epsilon_i(t-1)$.

It is easy to see that from the definition of $\epsilon_i(t)$, $\forall i \in V$, and $z_j(t)$, they both rely on the convergence rate of the network; i.e. rely on $\eta(t)$. That means that the failure of the post-convergence localization scheme, relies solely on the objective function h_i . Therefore, in order for the attacker to

bypass the given localization test, the attacker has to steer the network to a state that satisfies

$$\nabla h_i(\alpha_0) \simeq \lim_{t \rightarrow \infty} \frac{\epsilon_i(t) - z_j(t)}{\eta(t)} \quad (12)$$

VI. SIMULATIONS

This section presents the simulations conducted and the results are shown in the figures below. In the simulations we generated an "Erdos-Renyi" random-graph, consisting of N agents ($N = 50, 100, 500$) with random edge probability, $0 \leq p \leq 1$. We generated the adjacency matrix $\mathbf{W}(t) = \mathbf{I} - \frac{1}{2N}\mathbf{S} + \frac{1}{2N}(\mathbf{P} + \mathbf{P}^T)$, where \mathbf{P} is a random $N \times N$ matrix and \mathbf{S} is a diagonal matrix consisting of the column sum of $(\mathbf{P} + \mathbf{P}^T)$.

We assume that there is a majority of trustworthy agents in each neighborhood. In a case that the majority of neighbors in a trustworthy agent's neighborhood are attackers, this agent is likely to consider the attackers as trustworthy agents.

VI-A. Example: Detecting, localizing and eliminating 5 coordinated attackers while estimating logistic distribution mean:

In this example we present detection and localization in a single instance. After the network converges, each agent periodically look for attackers using (8). If a trustworthy agent suspect another agent in his neighborhood, the suspicious agent's data will be ignored for a given time period. When the given time period expires, the suspicious agent is being examined again with a more rigorous threshold, if the agent remain suspicious his data will still be ignored. If not, the trustworthy agent will use the data again (and so on). By ignoring the data from suspicious agents, we make sure that the network is reaching convergence to the true optimal state.

We assume that each agent holds a single measurement, (in this example $P = 1, 3, 5$ and 7) consisting of the desired signal with zero mean noise. Our goal is to extract the desired signal by eliminating the noise from the given measurement in a distributed manner. To simulate the problem we initialize the agents' state with values generated from a logistic distribution with parameters (μ, Σ) where $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_P)$. The agents solve a distributed problem with the following private objective function:

$$h_i(\theta) = 2 \sum_{p=1}^P \log \left(2 \cosh \left(\frac{[\mathbf{x}_i]_p - [\theta]_p}{2\sigma_p} \right) \right) + C \quad (13)$$

for all $i \in V$, where C is a constant number known to all agents, \mathbf{x}_i is the measured signal for some agent i and $\theta = \mu$, the desired variable. An example of the network convergence to the true optimal state after attackers elimination can be see in Figure 3. The localization scheme ROC, for different P s, is depicted in Figure 4. Looking at the simulations, we see that it is easier to perform a successful yet transparent attack on high dimensional problems.

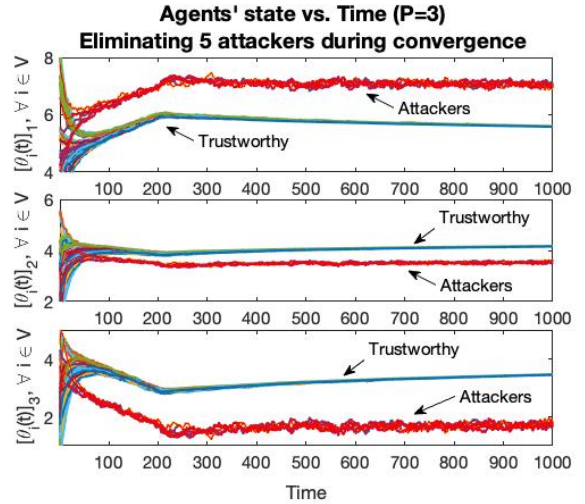


Fig. 3. Detection, localization and elimination of 5 coordinated attackers in a single instance, as explained in VI-A. In each subplot, presented different dimension of $\theta(t)$ ($P = 3$). We see that the trustworthy agents converge to the true optimal state, θ^* , while the attackers converge to α_0 .

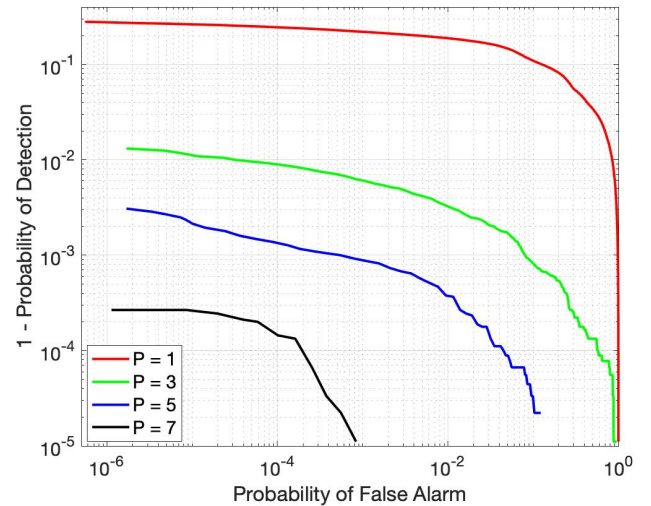


Fig. 4. ROCs temporal difference localization performance at the neighboring agents of the 5 coordinated attackers. $\theta \in \mathbb{C}^P$, $P = 1, 3, 5, 7$.

VII. CONCLUSIONS

In this paper we presented a novel attack on distributed multi-agent optimization. We then presented a combined detection and localization method in the case of distributed M-Estimators with i.i.d agents data. In an extension of this work we present detailed proofs of the exponential bounds for P_{FA} and P_D , as well as the propositions presented in this paper.

VIII. REFERENCES

- [1] Or Shalom, Amir Leshem, Anna Scaglione, and Angelia Nedić, “Detection of data injection attacks on decentralized statistical estimation,” in *2018 IEEE International Conference on the Science of Electrical Engineering (ICSEE)*. IEEE, 2018.
- [2] Sissi Xiaoxiao Wu, Hoi-To Wai, Anna Scaglione, Angelia Nedić, and Amir Leshem, “Data injection attack on decentralized optimization,” in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 3644–3648.
- [3] S Sundhar Ram, Angelia Nedić, and Venugopal V Veeravalli, “Distributed stochastic subgradient projection algorithms for convex optimization,” *Journal of optimization theory and applications*, vol. 147, no. 3, pp. 516–545, 2010.
- [4] Angelia Nedic and Asuman Ozdaglar, “Distributed subgradient methods for multi-agent optimization,” *IEEE Transactions on Automatic Control*, vol. 54, no. 1, pp. 48–61, 2009.
- [5] John Nikolas Tsitsiklis, “Problems in decentralized decision making and computation,” Tech. Rep., Massachusetts Inst of Tech Cambridge lab for information and decision systems, 1984.
- [6] John C Duchi, Alekh Agarwal, and Martin J Wainwright, “Dual averaging for distributed optimization: Convergence analysis and network scaling,” *IEEE Transactions on Automatic control*, vol. 57, no. 3, pp. 592–606, 2012.
- [7] Ali H Sayed et al., “Adaptation, learning, and optimization over networks,” *Foundations and Trends® in Machine Learning*, vol. 7, no. 4-5, pp. 311–801, 2014.
- [8] Stephen Boyd, Arpita Ghosh, Balaji Prabhakar, and Devavrat Shah, “Randomized gossip algorithms,” *IEEE transactions on information theory*, vol. 52, no. 6, pp. 2508–2530, 2006.
- [9] Alexandros G Dimakis, Soumya Kar, José MF Moura, Michael G Rabbat, and Anna Scaglione, “Gossip algorithms for distributed signal processing,” *Proceedings of the IEEE*, vol. 98, no. 11, pp. 1847–1864, 2010.
- [10] Dušan Jakovetić, Joao Xavier, and José MF Moura, “Fast distributed gradient methods,” *IEEE Transactions on Automatic Control*, vol. 59, no. 5, pp. 1131–1146, 2014.
- [11] Pascal Bianchi and Jérémie Jakubowicz, “Convergence of a multi-agent projected stochastic gradient algorithm for non-convex optimization,” *IEEE Transactions on Automatic Control*, vol. 58, no. 2, pp. 391–405, 2013.
- [12] Reinhard Gentz, Sissi Xiaoxiao Wu, Hoi-To Wai, Anna Scaglione, and Amir Leshem, “Data injection attacks in randomized gossiping,” *IEEE Transactions on Signal and Information Processing over Networks*, vol. 2, no. 4, pp. 523–538, 2016.
- [13] Yanpeng Guan and Xiaohua Ge, “Distributed attack detection and secure estimation of networked cyber-physical systems against false data injection attacks and jamming attacks,” *IEEE Transactions on Signal and Information Processing over Networks*, vol. 4, no. 1, pp. 48–59, 2018.
- [14] Edmond Nurellari, Des McLernon, Mounir Ghogho, et al., *Distributed detection and estimation in wireless sensor networks: resource allocation, fusion rules, and network security*, Ph.D. thesis, University of Leeds, 2017.
- [15] Chengcheng Zhao, Jianping He, and Jiming Chen, “Resilient consensus with mobile detectors against malicious attacks,” *IEEE Transactions on Signal and Information Processing over Networks*, vol. 4, no. 1, pp. 60–69, 2018.
- [16] Yanpeng Guan and Xiaohua Ge, “Distributed secure estimation over wireless sensor networks against random multichannel jamming attacks,” *IEEE Access*, vol. 5, pp. 10858–10870, 2017.
- [17] Michel Toulouse and Phuong Khanh Nguyen, “Protecting consensus seeking nids modules against multiple attackers,” in *Proceedings of the Eighth International Symposium on Information and Communication Technology*. ACM, 2017, pp. 226–233.
- [18] Rachad Atat, Lingjia Liu, Hao Chen, Jinsong Wu, Hongxiang Li, and Yang Yi, “Enabling cyber-physical communication in 5g cellular networks: challenges, spatial spectrum sensing, and cyber-security,” *IET Cyber-Physical Systems: Theory & Applications*, vol. 2, no. 1, pp. 49–54, 2017.
- [19] Rachad Atat, Lingjia Liu, Jonathan Ashdown, Michael J Medley, John D Matyjas, and Yang Yi, “A physical layer security scheme for mobile health cyber-physical systems,” *IEEE Internet of Things Journal*, vol. 5, no. 1, pp. 295–309, 2018.
- [20] Edmond Nurellari, Des McLernon, and Mounir Ghogho, “A secure optimum distributed detection scheme in under-attack wireless sensor networks,” *IEEE Transactions on Signal and Information Processing over Networks*, vol. 4, no. 2, pp. 325–337, 2018.
- [21] Reinhard Gentz, *Wireless Sensor Data Transport, Aggregation and Security*, Ph.D. thesis, Arizona State University, 2017.
- [22] Wei Xu, Zhengqing Li, and Qing Ling, “Robust decentralized dynamic optimization at presence of malfunctioning agents,” .
- [23] Sara A Van de Geer and Sara van de Geer, *Empirical Processes in M-estimation*, vol. 6, Cambridge university press, 2000.