

Learning Trust Over Directed Graphs in Multiagent Systems

Orhan Eren Akgün *

Arif Kerem Dayı *

Stephanie Gil *

*SEAS, Harvard University

Angelia Nedić **

**ECE, Arizona State University

ERENAKGUN@G.HARVARD.EDU

KEREMDAYI@COLLEGE.HARVARD.EDU

SGIL@SEAS.HARVARD.EDU

ANGELIA.NEDICH@ASU.EDU

Editors: N. Matni, M. Morari, G. J. Pappas

Abstract

We address the problem of learning the legitimacy of other agents in a multiagent network when an unknown subset is comprised of malicious actors. We specifically derive results for the case of directed graphs and where stochastic side information, or observations of trust, is available. We refer to this as “learning trust” since agents must identify which neighbors in the network are reliable, and we derive a protocol to achieve this. We also provide analytical results showing that under this protocol i) agents can learn the legitimacy of all other agents almost surely, and that ii) the opinions of the agents converge in mean to the true legitimacy of all other agents in the network. Lastly, we provide numerical studies showing that our convergence results hold in practice for various network topologies and variations in the number of malicious agents in the network.

Keywords: Multiagent systems, adversarial learning, directed graphs, networked systems

1. Introduction

Learning the network topology in multiagent systems, what edges exist and are reliable, is critical because of the central role it plays in many multiagent collaboration tasks. This includes a wide range of tasks from estimation, to control, to machine learning, optimization and beyond [Rabbat and Nowak \(2004\)](#); [Olshevsky \(2010\)](#); [Nedić et al. \(2018\)](#). Many times both the coordination protocols and achievable performance of the team is dictated by topology [Olfati-Saber and Murray \(2004\)](#); [Xi et al. \(2018\)](#); [Cai and Ishii \(2012\)](#); [Nedić and Olshevsky \(2014\)](#). Two aspects that can greatly complicate the learning however, are i) directed graphs, and ii) the presence of untrustworthy data. Directed graphs are more common in practice due to heterogeneity in sensing and communication capabilities in multiagent systems, but are often more difficult to analyze due to non-symmetric information flow. On the other hand, the presence of malicious agents are an important real-world consideration but lead to untrustworthy data in the system [Lamport et al. \(2019\)](#); [Sundaram and Hadjicostis \(2010\)](#); [Sundaram and Gharesifard \(2018\)](#); [Fischer et al. \(1985\)](#). Unfortunately, the compounded impact of both of these challenges is a very complex problem with sparse theory to date. *Our objective in this paper is to develop a learning protocol and its related analysis, where agents learn the legitimacy of their neighbors over time in the presence of malicious agents over directed graphs.*

The class of problems over directed graphs pose a particular challenge to achieving resilience: many distributed algorithms on directed graphs require agents to have some information about their out-neighbors, but because of the asymmetric information flow, they cannot sense or obtain information directly from these agents. This makes detection of malicious out-neighbors particularly difficult. For instance, the distributed optimization algorithms presented in [Nedić and Olshevsky](#)

(2014); Tsianos et al. (2012b,a); Makhdoumi and Ozdaglar (2015); Pu et al. (2021) and the distributed consensus algorithms Cai and Ishii (2012); Dominguez-Garcia and Hadjicostis (2012) all require that the agents know the number of out-neighbors they have. This assumption can break if an agent designs the update rule considering an out-neighbor as legitimate, but that agent is malicious in reality. Hence, agents need to have some information about the trustworthiness of their out-neighbors. An interesting concept that has the potential to help this difficult problem is the use of “side information” or data in cyberphysical systems Liu et al. (2019); Xiong and Jamieson (2013); Pasqualetti et al. (2015); Renganathan and Summers (2017); Gil et al. (2017, 2019); Giraldo et al. (2018); Mallmann-Trenn et al. (2021); Cavorsi et al. (2022). Recent work has shown that by leveraging physical channels of information in the system, agents can gain stochastic information about the trustworthiness of the other agents Liu et al. (2019); Giraldo et al. (2018); Xiong and Jamieson (2013); Gil et al. (2017). We call these “stochastic observations of trust.” It has been shown that exploiting these observations leads to stronger results in resilience for multiagent systems Yemini et al. (2022); Gil et al. (2019); Mallmann-Trenn et al. (2021). Unfortunately however, existing results do not immediately extend to the case of directed graphs.

In this work, we are interested in learning a trusted graph topology over a directed graph. Using stochastic information about trustworthy neighbors, agents can decide how they should process information that they receive from their in-neighbors, and with which out-neighbors they should share their information. Since agents cannot necessarily observe their out-neighbors, it is natural to think that they need to get information about their out-neighbors from the other agents. We investigate what sufficient information agents can share and how they should process this information to learn the trustworthiness of the other agents in the system in a robust way. This setup is particularly challenging since there might be malicious agents in the system sharing misinformation during this learning process. We present a learning protocol to enable each agent to learn the trustworthiness of all other agents in the system leveraging the opinion of their neighbors. Agents develop opinions in two ways: For their in-neighbors they can obtain a trust observation, they then use this information to form their own opinions. For the other agents, they use the opinions of their in-neighbors they trust to update their opinions. Under the assumption that the subgraph of legitimate agents is strongly connected and each malicious agent is observed by at least one legitimate agent, we show that all legitimate agents can almost surely learn the trustworthiness of all other agents.

Our contributions can be summarized as follows: i) We present a novel learning protocol that enables the legitimate agents in the system to learn the trustworthiness of the other agents where the underlying communication network is a directed graph; ii) We prove that using our learning protocol, legitimate agents can learn the identities of the other agents almost surely; iii) We show that opinions of the agents converge in mean to the true identity of the agents; iv) We provide extensive numerical studies to show that the convergence results hold in practice for various network topologies and the number of malicious agents.

2. Problem Formulation

We consider a distributed multi-agent system where agents need to collaborate in order to achieve a common task such as solving an optimization problem. We represent the communication graph among agents with a directed graph $G = (V, E)$ where the set V represents the set of agents communicating over G with a set E of directed links. Moreover, we let $N = |V|$ be the number of agents. If there is an edge $(i, j) \in E$, then agent i can send information to j , and we say that j is an out-neighbor of i and i is an in-neighbor of j . We assume every agent i has a self loop $(i, i) \in E$. Moreover, for an agent $i \in V$, we define its in-neighborhood $\mathcal{N}_i^{\text{in}} = \{j \in V \mid (j, i) \in E\}$ and out-

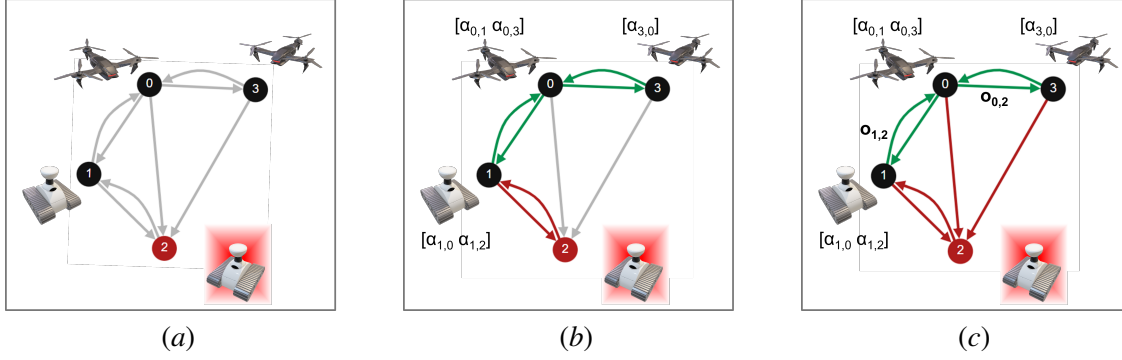


Figure 1: This schematic shows our problem setup with one malicious agent shown as a red node. α_{ij} and o_{ij} are defined in Definitions 1 and 2, respectively. Various stages of learning are depicted: (a) initial state (b) agents use their direct observations to learn the trustworthiness of other agents (c) agents indirectly learn the trustworthiness of the entire network by propagating their opinions.

neighborhood $\mathcal{N}_i^{\text{out}} = \{j \in V \mid (i, j) \in E\}$. We assume that agents in the system communicate at every time step t . Moreover, we assume that there might be a set $\mathcal{M} \subsetneq V$, called *malicious agents*, of non-cooperative agents in the system that are either adversarial or malfunctioning. We assume that malicious agents can act arbitrarily. We call the set of cooperative agents, that is, the set of agents outside the set \mathcal{M} , legitimate agents denoted by \mathcal{L} . We have $\mathcal{L} \cap \mathcal{M} = \emptyset$ and $\mathcal{L} \cup \mathcal{M} = V$. We say that malicious agents are untrustworthy and legitimate agents are trustworthy. We assume that the set of malicious agents \mathcal{M} and the set of legitimate agents \mathcal{L} are unknown. We wish to learn the *trustworthiness* of agents in the network. We are interested in the problems where every agent receives a stochastic observation of trust from an agent that sends information during each communication round. We note that stochastic observations of trust have been developed in previous works [Gil et al. \(2017\)](#); [Yemini et al. \(2022\)](#) and we use a similar definition here:

Definition 1 (Stochastic Observation of Trust α_{ij}) We denote stochastic observations of trust with $\alpha_{ij}(t)$ if agent j sends information to agent i at time t , and we assume that $\alpha_{ij}(t) \in [0, 1]$. Here, $\alpha_{ij}(t)$ represents the stochastic value of trust of agent j as observed by agent i .

Agents can develop opinions about trustworthiness of their in-neighbors using these stochastic trust observations over time. However, it is not straightforward how they can develop opinions about their out-neighbors since they have no direct observations of their trustworthiness. Next, we formalize the notion of opinion and then we discuss how to construct opinions of agents.

Definition 2 (Opinion of Trust) We denote agent i 's opinion of trust about agent j at time t with $o_{ij}(t) \in [0, 1]$. We say agent i trusts agent j at time t if $o_{ij}(t) \geq 1/2$ and does not trust agent j otherwise.

We want to find a learning protocol to enable the legitimate agents to develop accurate opinions $o_{ij}(t)$ about their neighbors in directed graphs, including their out-neighbors. An example case is shown in Figure 1. Next, we state our assumptions under which we develop our protocol.

Assumption 1 (Connectivity of Network)

1. *Sufficiently connected graph:* The subgraph $G_{\mathcal{L}}$ induced by the legitimate agents is strongly connected.
2. *Observation of malicious agents:* For any malicious agent $j \in \mathcal{M}$, there exists some legitimate agent $i \in \mathcal{L}$ that observes j , i.e., $j \in \mathcal{N}_i^{\text{in}}$ for some $i \in \mathcal{L}$.

Assumption 2 (Trust Observations) Suppose that the following hold:

1. *Homogeneity of trust variables:* The expectation of the variables $\alpha_{ij}(t)$ are constant for the case of malicious transmissions and legitimate transmissions, respectively, i.e., for some scalars c, d with $c < 0$ and $d > 0$, $c = \mathbb{E}[\alpha_{ij}(t)] - 1/2$ for all $i \in \mathcal{L}$, $j \in \mathcal{N}_i^{\text{in}} \cap \mathcal{M}$, and $d = \mathbb{E}[\alpha_{ij}(t)] - 1/2$ for all $i \in \mathcal{L}$, $j \in \mathcal{N}_i^{\text{in}} \cap \mathcal{L}$.
2. *Independence of trust observations:* The observations $\alpha_{ij}(t)$ are independent for all t and all pairs of agents i and j , with $i \in \mathcal{L}$, $j \in \mathcal{N}_i^{\text{in}}$. Moreover, for any $i \in \mathcal{L}$ and $j \in \mathcal{N}_i^{\text{in}}$, the observation sequence $\{\alpha_{ij}(t)\}_{t \in \mathbb{N}}$ is identically distributed.

Note that stochastic observations of trust satisfying Assumption 2.1 were derived in Gil et al. (2017). Additionally, we make the same Assumptions 1.1, 2.1, and 2.2 as in the work Yemini et al. (2022), except for the first assumption, where we require the graph to be strongly connected instead of connected since we deal with directed graphs. Assumption 1.2 is new and needed since it is not possible to learn the legitimacy of an agent if no other agent is observing that agent. This requirement shows up in the analysis later on. We formalize the problem that we are aiming to solve in this paper as follows:

Problem 1 Let $i \in \mathcal{L}$ be a legitimate agent and let $q \in V$ be an arbitrary agent in the system. Assume that stochastic observations of trust are available and Assumption 1 and Assumption 2 hold. We want to find a learning protocol such that for all legitimate agent $i \in \mathcal{L}$ and for all agents $q \in V$, $o_{iq}(t)$ converges to 1 if $q \in \mathcal{L}$ and 0 if $q \in \mathcal{M}$ almost surely.

3. Learning Protocol

In this section we introduce our learning protocol. Let each agent i store a vector of trust $o_i(t)$ at time t , where $o_i(t)$ is an $N \times 1$ column vector. Let $o_{ij}(t)$ denote the j th component of $o_i(t)$. The value $o_{ij}(t)$ represents agent i 's opinion about the node j where a higher $o_{ij}(t)$ indicates that agent i trusts agent j more. Let $\beta_{ij}(t)$ represent an aggregate trust value for the link (j, i) at time t . Following Yemini et al. (2022), we define $\beta_{ij}(t)$ as

$$\beta_{ij}(t) = \sum_{k=0}^t (\alpha_{ij}(k) - 1/2), \quad (1)$$

for all $j \in \mathcal{N}_i^{\text{in}}$ and we define $\beta_{ii}(t) = 1$ for all t . Using the aggregate stochastic trust value $\beta_{ij}(t)$, a legitimate agent i decides on its trusted in-neighbor set at time t by defining $\mathcal{N}_i^{\text{in}}(t) = \{j \in \mathcal{N}_i^{\text{in}} \mid \beta_{ij}(t) \geq 0\}$. In our learning protocol, an agent i shares $o_i(t)$ with its out-neighbors. A legitimate agent i determines its vector of $o_i(t)$ after receiving $o_j(t-1)$ from all of its in-neighbors $j \in \mathcal{N}_i^{\text{in}}$ using the following update rule:

$$o_{iq}(t) = \begin{cases} 1 & \text{if } q \in \mathcal{N}_i^{\text{in}} \text{ and } \beta_{iq}(t) \geq 0 \\ 0 & \text{if } q \in \mathcal{N}_i^{\text{in}} \text{ and } \beta_{iq}(t) < 0. \\ \sum_{j \in \mathcal{N}_i^{\text{in}}(t)} \frac{o_{jq}(t-1)}{|\mathcal{N}_i^{\text{in}}(t)|} & \text{if } q \notin \mathcal{N}_i^{\text{in}} \end{cases} \quad (2)$$

Every legitimate agent i initializes its opinion vector with vector $o_i(0)$ with all ones, meaning that in the beginning, they trust everyone in the network. However, this choice of initialization is arbitrary

and as it does not affect our results. A legitimate agent i decides on its trusted out-neighbor set at time t by defining $\mathcal{N}_i^{\text{out}}(t) = \{j \in \mathcal{N}_i^{\text{out}} \mid o_{ij}(t) \geq 1/2\}$.

Notice that the trust vector $o_i(t)$ is in $[0, 1]^N$ by definition. Since $o_{iq}(t) \in [0, 1]$ for all legitimate agents, any malicious agent that sends an opinion outside this range would reveal itself. Therefore, we assume that malicious agents' opinions are also in the range $[0, 1]$. However, we assume that malicious agents can decide their trust vectors $o_i(t)$ arbitrarily within this range. This assumption captures strong attacks where malicious agents coordinate with each other to choose their opinions knowing the true trustworthiness of everyone. With our protocol, legitimate agents use only the stochastic observations of trust α_{ij} to determine the legitimacy of their in-neighbors. For the other nodes, they use the opinions of their trusted in-neighbors to form their opinion.

4. Analysis

Recall that agents either directly observe an agent and develop their opinions using these observations, or they use the opinions of others to generate an opinion about an agent. In our analysis, we first show that all legitimate agents learn the trustworthiness of their in-neighbors. Then, we analyze the propagation of information thereafter. We show that estimated trust values converge in mean and almost surely to true trust values (1 for legitimate, 0 for malicious agents). **Some of the proofs are provided in our extended technical report due to space limitations** [Akgün et al. \(2022\)](#)

4.1. Notation

Let $[W]_{ij}$ denote entry in row i and column j of matrix W . For some agent j and a set S , define the indicator function $\mathbf{1}_{\{j \in S\}}$ as 1 if $j \in S$ and 0 otherwise. We also use the same notation for indicator vectors when the size of the vector is clear from the context. Finally, let the square matrix $W \in \mathbb{R}^{n \times n}$ be non-negative, i.e. $W_{ij} \geq 0$ for all i, j . Hence, the digraph of W , denoted by $G(W) = (V(W), E(W))$ is the graph such that $V(W) = \{1, \dots, n\}$ and for all $i, j \in \{1, \dots, n\}$, $(i, j) \in E(W)$ if and only if $W_{ij} > 0$.

4.2. Learning Trustworthiness

Since agents use their trusted in-neighbors in their updates, we start by showing that agents learn the trustworthiness of their in-neighbors. This will be useful later to show that the protocol converges to the desired state. The lemma below follows from ([Yemini et al., 2022](#), Proposition 1).

Lemma 3 *There exists a random finite time T_f such that for all $t \geq T_f$ and for all legitimate agents i , the trusted in-neighbor set is $\mathcal{N}_i^{\text{in}}(t) = \mathcal{N}_i^{\text{in}} \cap \mathcal{L}$ almost surely.*

Notice that Lemma 3 shows that every legitimate agent can learn its in-neighbors correctly. Now, let $q \in V$ be an arbitrary but fixed agent in the network. Our goal is to show that all legitimate agents learn whether q is legitimate or not. This process requires information to propagate from agents receiving trust information directly from q to other agents in the network, which motivates the following definition to use in our analysis:

Definition 4 *Let $q \in V$. Define $\mathcal{D}_q \subseteq \mathcal{L}$ to be the subset of legitimate agents directly observing q , i.e. $\mathcal{D}_q \triangleq \mathcal{N}_q^{\text{out}} \cap \mathcal{L}$. Define $\mathcal{C}_q \triangleq \mathcal{L} \setminus \mathcal{D}_q$ as the subset of legitimate agents not observing q .*

These sets are illustrated in Fig. 2. Because of Assumption 1, there is at least one legitimate agent that observes q , so \mathcal{D}_q of observing agents is non-empty. On the other hand, if \mathcal{C}_q is empty, then all agents are directly observing q . In that case, all legitimate agents will eventually learn the identity of q by Lemma 3.

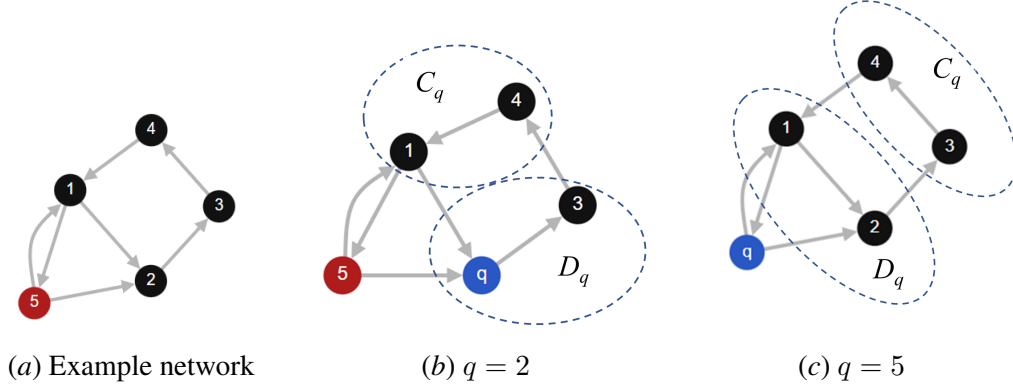


Figure 2: (a) Network with four legitimate (black nodes) and one malicious (red node) agents. (b) and (c) show the sets C_q and D_q for different q nodes based on the directionality of the graph edges.

Now, we analyze the evolution of $o_{iq}(t)$ by writing the evolution of opinions about agent q in matrix form. Let $u_q = |C_q|$. Without loss of generality, reorder the indices of agents such that $C_q = \{1, 2, \dots, u_q\}$, $D_q = \{u_q + 1, \dots, |\mathcal{L}|\}$, and $\mathcal{M} = \{|\mathcal{L}| + 1, \dots, N\}$. Hence, we have:

$$o_{iq}(t) = \sum_{j \in C_q} [W_q(t)]_{ij} o_{jq}(t-1) + \sum_{j \in D_q} [W_q(t)]_{ij} o_{jq}(t-1) + \sum_{j \in \mathcal{M}} [W_q(t)]_{ij} o_{jq}(t-1), \quad (3)$$

where $[W_q(t)]_{ij} = \frac{1}{|\mathcal{N}_i^{\text{in}}(t)|}$ if $j \in \mathcal{N}_i^{\text{in}}(t)$ and $[W_q(t)]_{ij} = 0$ otherwise following from the update rule (2). Here, $W_q(t)$ is a row-stochastic matrix with size $u_q \times N$. Then, we can write $W_q(t) = [W_{C_q}(t) \ W_{D_q}(t) \ W_{\mathcal{M}}(t)]$ where the matrices $W_{C_q}(t)$, $W_{D_q}(t)$, $W_{\mathcal{M}}(t)$ have sizes $u_q \times u_q$, $u_q \times |D_q|$, and $u_q \times |\mathcal{M}|$ respectively. Let the vectors $o_{C_q}(t)$, $o_{D_q}(t)$, and $o_{\mathcal{M}}(t)$ denote the trust estimates of agents in C_q and D_q and \mathcal{M} . So, we can write (3) in the matrix form as:

$$o_{C_q}(t) = [W_{C_q}(t) \ W_{D_q}(t) \ W_{\mathcal{M}}(t)] \begin{bmatrix} o_{C_q}(t-1) \\ o_{D_q}(t-1) \\ o_{\mathcal{M}}(t-1) \end{bmatrix}, \quad (4)$$

Recall that there exists some random finite time T_f such that all legitimate agents learn their in-neighbors correctly. Until the system reaches time T_f , malicious agents can affect the learning dynamics. Nevertheless, we will show that the legitimate agents can recover from that effect after reaching time T_f . Now, we focus our analysis on the system dynamics after time T_f .

Lemma 5 *For $t \geq T_f$, the following hold almost surely (i) $W_{\mathcal{M}}(t) = 0$, (ii) $W_{C_q}(t) = \overline{W}_{C_q}$, (iii) $W_{D_q}(t) = \overline{W}_{D_q}$, (iv) $o_{D_q}(t) = \mathbf{1}_{\{q \in \mathcal{L}\}}$ for some constant matrices \overline{W}_{C_q} and \overline{W}_{D_q} .*

Notice that since $W_q(t)$ is a row-stochastic matrix and that $W_{\mathcal{M}}(t)$ is zero, the matrix $[\overline{W}_{C_q} \ \overline{W}_{D_q}]$ is row stochastic. We now focus on the agents in C_q . For all $t \geq T_f + 1$, we can describe the evolution of $o_{C_q}(t)$ as follows:

$$o_{C_q}(t) = \overline{W}_{C_q} o_{C_q}(t-1) + \overline{W}_{D_q} o_{D_q}(t-1) \quad (5)$$

Now, we define $\Delta_{C_q}(t) = o_{C_q}(t) - \mathbf{1}_{\{q \in \mathcal{L}\}}$. Rearranging 5, and using the fact that $[\overline{W}_{C_q} \ \overline{W}_{D_q}]$ is row-stochastic and $o_{D_q}(t-1) = \mathbf{1}_{\{q \in \mathcal{L}\}}$, we get:

$$\Delta_{C_q}(t) = \overline{W}_{C_q} \Delta_{C_q}(t-1) \quad (6)$$

Now, we can bound the error norm $\|\Delta_{C_q}(t)\| \leq \|\bar{W}_{C_q}^{t-T_f}\| \|\Delta_{C_q}(T_f)\|$. Here, $\|\Delta_{C_q}(T_f)\|$ includes the error introduced by malicious agents before all agents learn their in-neighbors. Since the convergence of the error term $\|\Delta_{C_q}(t)\|$ depends on the convergence of \bar{W}_{C_q} , we analyze the matrix \bar{W}_{C_q} next.

4.3. Convergence of Weakly Chained Substochastic Matrices

Now, we aim to show that \bar{W}_{C_q} is *convergent*, i.e. $\|\bar{W}_{C_q}^t\| \rightarrow 0$ as $t \rightarrow \infty$. In this part, we will show that \bar{W}_{C_q} belongs to a family of *convergent* substochastic matrices called *weakly chained substochastic* matrices. To analyze the convergence properties of \bar{W}_{C_q} , we define the index of contraction following [Azimzadeh \(2019\)](#)

Definition 6 *Index of contraction:* Let the matrix $W \in \mathbb{R}^{n \times n}$ be substochastic. Define the set $\hat{J}(W) \triangleq \{1 \leq i \leq n : \sum_{j=1}^n W_{ij} < 1\}$, and let the set $\hat{K}_i(W)$ be the set of all paths¹ in the digraph of W from i to all $j \in \hat{J}(W)$. The index of contraction $\widehat{con}W$ associated with matrix W is defined as:

$$\widehat{con}W \triangleq \max \left\{ 0, \sup_{i \notin \hat{J}(W)} \left\{ \inf_{\omega \in \hat{K}_i(W)} \{|\omega|\} \right\} \right\}, \quad (7)$$

where $|\omega|$ denotes the length of the path ω . Also, we follow the conventions that $\inf \emptyset = \infty$ and $\sup \emptyset = -\infty$.

([Azimzadeh, 2019](#), Corollary 2.6) shows that a square substochastic matrix W is convergent if and only if $\widehat{con}W$ is finite. We call a substochastic matrix with finite contraction index *weakly chained substochastic matrix*.

Remark 7 Matrix W is a weakly chained substochastic matrix if and only if for all rows i that are not in the set $\hat{J}(W)$, set $\hat{K}_i(W)$ is non-empty, i.e there is a path $i \rightarrow i_1 \rightarrow \dots \rightarrow i_j$ in $G(W)$ such that row i_j sums to less than one. Moreover, a weakly chained substochastic matrix is convergent.

This remark follows directly from the definition of the index of contraction and ([Azimzadeh, 2019](#), Corollary 2.6). The following sequence of results will show that \bar{W}_{C_q} is weakly chained substochastic. We will establish that the links $E(\bar{W}_{C_q})$ in the digraph of \bar{W}_{C_q} are the inversion of links in the original graph. Then use assumptions of strong connectivity and existence of a directly observing agent to conclude that \bar{W}_{C_q} is weakly chained substochastic.

Lemma 8 Let $\bar{W}_{C_q} \in \mathbb{R}^{u_q \times u_q}$ be defined as before, and let G_{C_q} be the subgraph of $G_{\mathcal{L}}$ induced by the set of agents C_q . Then, $(i, j) \in G(\bar{W}_{C_q})$ if and only if $(j, i) \in G_{C_q}$.

Corollary 9 If there is a path $v_1 \rightarrow v_2 \rightarrow \dots \rightarrow v_l$ in G_{C_q} , then there is a path $v_l \rightarrow v_{l-1} \rightarrow \dots \rightarrow v_1$ in $G(\bar{W}_{C_q})$.

Theorem 10 For all agents q , given that the set C_q is non-empty, the update matrix \bar{W}_{C_q} is a weakly chained substochastic matrix. Moreover, \bar{W}_{C_q} is convergent.

Proof Let $i \in C_q$. If agent i has a neighbor $d \in \mathcal{D}_q$ directly observing agent q , row i must sum up to less than one since agent i receives information from d and $d \notin C_q$. So, $i \in \hat{J}(\bar{W}_{C_q})$.

1. We use path instead of walk in contrast to [Azimzadeh \(2019\)](#) in our definition, however these definitions are equivalent.

Assume agent i doesn't have a directly observing neighbor, i.e. $i \notin \hat{J}(\overline{W}_{C_q})$. We know that there exists some agent $d \in \mathcal{D}_q$ that directly observes agent q by Assumption 1.1 and Assumption 1.2. By Assumption 1.1, the subgraph induced by legitimate agents are strongly connected, so there exists a path $d = i_0 \rightarrow i_1 \rightarrow i_2 \rightarrow \dots \rightarrow i_l \rightarrow i$ in $G_{\mathcal{L}}$ where each arrow denotes a directed edge. $l \geq 1$ since agent i does not have a directly observing neighbor. Now, choose the largest j such that $i_j \in \mathcal{D}_q$, and consider the path $i_j \rightarrow i_{j+1} \rightarrow \dots \rightarrow i_l \rightarrow i$. Here, since j is chosen as the largest j s.t. $i_j \in \mathcal{D}_q$, we have that $i_{j+1}, \dots, i_l, i \in C_q$. Moreover, we assumed $i \notin \hat{J}(\overline{W}_{C_q})$, so $j < l$ since i does not have a directly observing neighbor.

Now, we know, i_{j+1} has a neighbor directly observing q , i.e. i_j . Therefore, row i_{j+1} of \overline{W}_{C_q} sums to less than 1, meaning that $i_{j+1} \in \hat{J}(\overline{W}_{C_q})$. From Corollary 9, there exists a path $i \rightarrow i_l \rightarrow i_{l-1} \rightarrow \dots \rightarrow i_{j+2} \rightarrow i_{j+1}$ in the graph $G(\overline{W}_{C_q})$. Hence, $\hat{K}_i(\overline{W}_{C_q})$ is non-empty for $i \notin \hat{J}(\overline{W}_{C_q})$. Therefore, \overline{W}_{C_q} is weakly chained substochastic and convergent by Remark 7. ■

Corollary 11 *For all agents $q \in V$ where the set C_q is non-empty, $o_{C_q}(t)$ almost surely converges to $\mathbf{1}_{\{q \in \mathcal{L}\}}$ where $\mathbf{1}_{\{q \in \mathcal{L}\}}$ is a vector with all values equal to 1 if $q \in \mathcal{L}$ and to 0 if $q \in \mathcal{M}$.*

Proof Recall that the error is defined as $\Delta_{C_q}(t) = o_{C_q}(t) - \mathbf{1}_{\{q \in \mathcal{L}\}}$. By Lemma 3 we know that there exists a finite time T_f such that for all $t \geq T_f + 1$ we have $\|\Delta_{C_q}(t)\| \leq \|\overline{W}_{C_q}^{t-T_f}\| \|\Delta_{C_q}(T_f)\|$. Since both $o_{C_q}(t)$ and $\mathbf{1}_{\{q \in \mathcal{L}\}}$ are in $[0, 1]^{u_q}$, we have $\|\Delta_{C_q}(T_f)\| \leq \sqrt{u_q}$. By Theorem 10, we have that $\|\overline{W}_{C_q}^{t-T_f}\| \rightarrow 0$. Therefore, $\|\Delta_{C_q}(t)\| \rightarrow 0$ almost surely. ■

4.4. Main Results

In this part, we present our main results which show that the trust vector of legitimate agents $o_{\mathcal{L}_q}(t)$ converges to the true vector $\mathbf{1}_{\{q \in \mathcal{L}\}}$. **Proofs of these results are provided in our extended technical report Akgün et al. (2022)**

Theorem 12 (Convergence to the true trust vector almost surely) *For all agents $q \in V$, $o_{\mathcal{L}_q}(t)$ converges almost surely to the true trust vector $\mathbf{1}_{\{q \in \mathcal{L}\}}$, where $\mathbf{1}_{\{q \in \mathcal{L}\}}$ is an $|L| \times 1$ vector with all of its values equal to 1 if $q \in \mathcal{L}$ and equal to 0 if $q \in \mathcal{M}$.*

The proof idea is that $o_{\mathcal{D}_q}(t)$ converges to the true vector because agents in \mathcal{D}_q directly observe q , and $o_{C_q}(t)$ converges to the true value by Corollary 11.

Theorem 13 (Convergence in mean to the true trust vector) *For all agents $q \in V$ and $r \geq 1$, $o_{\mathcal{L}_q}(t)$ converges in mean to the true trust vector $\mathbf{1}_{\{q \in \mathcal{L}\}}$. That is,*

$$\lim_{t \rightarrow \infty} E[\|o_{\mathcal{L}_q}(t) - \mathbf{1}_{\{q \in \mathcal{L}\}}\|^r] = 0. \quad (8)$$

Finally, the following Corollary shows that following this protocol, every legitimate agent can learn the trustworthiness of *all agents* in the network, including their in- and out-neighbors.

Corollary 14 (Learning the Trustworthiness of All Agents) *All legitimate agents $i \in \mathcal{L}$ can learn the trustworthiness of all agents in the network correctly. That is, there exists a finite time T_{max} such that for all $t \geq T_{max}$ and for all $q \in V$, $o_{iq}(t) \geq 1/2$ if $q \in \mathcal{L}$ and $o_{iq}(t) < 1/2$ if $q \in \mathcal{M}$ almost surely.*

5. Numerical Studies

Now, we verify our theoretical results and provide more insight into our protocol with numerical studies.

	$ \mathcal{L} = 20, \mathcal{M} = 30$		$ \mathcal{L} = 40, \mathcal{M} = 60$		$ \mathcal{L} = 80, \mathcal{M} = 120$	
	\hat{T}_{max}	\widehat{con}_{max}	\hat{T}_{max}	\widehat{con}_{max}	\hat{T}_{max}	\widehat{con}_{max}
Cyclic	66	19	109	38	192	78
Erdős–Rényi	49	3	64	3	76	3

Table 1: This table shows \hat{T}_{max} and \widehat{con}_{max} for 8 different setups. Large \widehat{con}_{max} usually corresponds to a large \hat{T}_{max} since \widehat{con}_{max} is an approximation of how long it takes for information to propagate from observing agents to non-observing agents.

5.1. Experimental Setup

We generate the graph of legitimate agents with cyclic graphs where the contraction index grows linearly with $|\mathcal{L}|$, and random graphs generated using the Erdős–Rényi model where each edge in the graph is either included or not with probability $\frac{2 \log |\mathcal{L}|}{|\mathcal{L}|}$ Erdős et al. (1960). The probability $\frac{2 \log |\mathcal{L}|}{|\mathcal{L}|}$ is chosen to have a high probability of generating a strongly connected graph Graham and Pike (2008), and we repeat the process to ensure strong connectivity as in Assumption 1.1. The Erdős–Rényi graphs we generated are likely to have stronger connectivity and lower contraction indices compared to cyclic graphs. In this aspect, the cyclic graph represents a difficult case where trust information propagates slowly. Then, we add malicious agents randomly, and they send the exact opposite of the true legitimacy values to their neighbors. This ensures a strong attack. Following the previous work Yemini et al. (2022), we model the trust observations $\alpha_{ij}(t)$ as follows: At each time step t we sample $\alpha_{ij}(t)$ uniformly from the interval $[0.35, 0.75]$ if $j \in \mathcal{L}$ and from $[0.25, 0.65]$ if $j \in \mathcal{M}$. This way, $\mathbb{E}[\alpha_{ij}(t)] = 0.55$ if j is a legitimate agent and $\mathbb{E}[\alpha_{ij}(t)] = 0.45$ otherwise. With this setup, Assumption 2 is satisfied.

5.2. Numerical Results

We evaluate the protocol performance based on mean squared error (MSE) of vectors of trust and time to learn \hat{T}_{max} , defined as the first time agents classify others correctly for N consequent steps, which is a proxy to T_{max} as defined in Corollary 14. Here, we present the results for three different setups with $|\mathcal{L}| \in \{20, 40, 80\}$ and $|\mathcal{M}| = 1.5 \times |\mathcal{L}|$. For each $|\mathcal{L}|$, we test over cyclic graphs and an Erdős–Rényi graph over legitimate agents which are illustrated in Figure 3. The results are presented in Figure 4. It can be seen that both MSE and maximum error converge to 0 in all setups. For each setup, we present the maximum contraction index, denoted by \widehat{con}_{max} and \hat{T}_{max} in Table 1. We define the maximum contraction index as $\widehat{con}_{max} = \max_{q \in V} \widehat{con} \overline{W}_{C_q}$, where $\widehat{con} \overline{W}_{C_q}$ is defined in (7).

Now, we test the effect of malicious agents in the system to the learning protocol. We use the Erdős–Rényi graph setup as before with 40 legitimate agents. First, we fix the number of malicious agents to 60 and change the probability of connection between the malicious agents and legitimate agents. Then, we change the number of malicious agents. The MSE graphs are shown in Fig 5

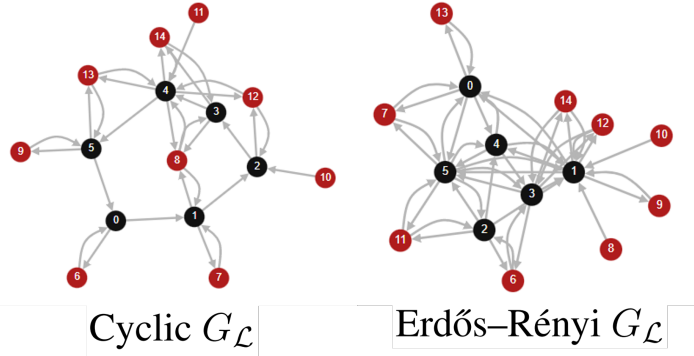


Figure 3: Example graph topologies with $|\mathcal{L}| = 6, |\mathcal{M}| = 9$ nodes.

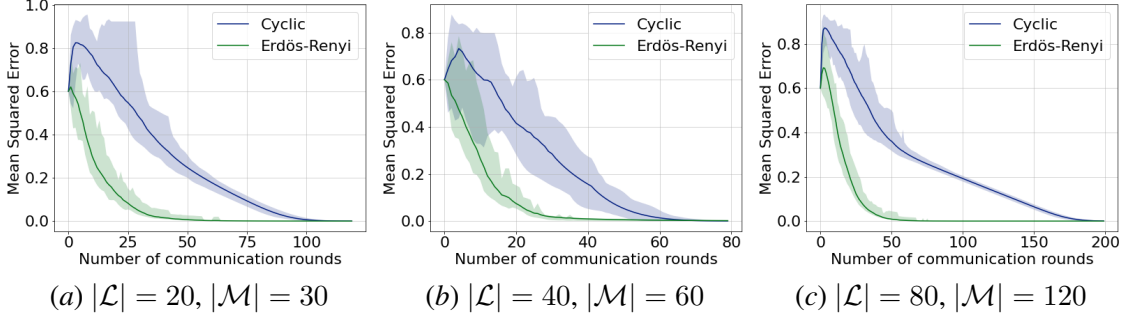


Figure 4: Aggregate MSE plots for three cases where $|\mathcal{M}| = 1.5 \times |\mathcal{L}|$. Error converges to zero as predicted by our theory. Initial disturbance by malicious agents is higher with cyclic graphs since information propagates more slowly. Convergence time increases as size of graph increases, but Erdős–Rényi graphs are less sensitive to this change because of their good connectivity.

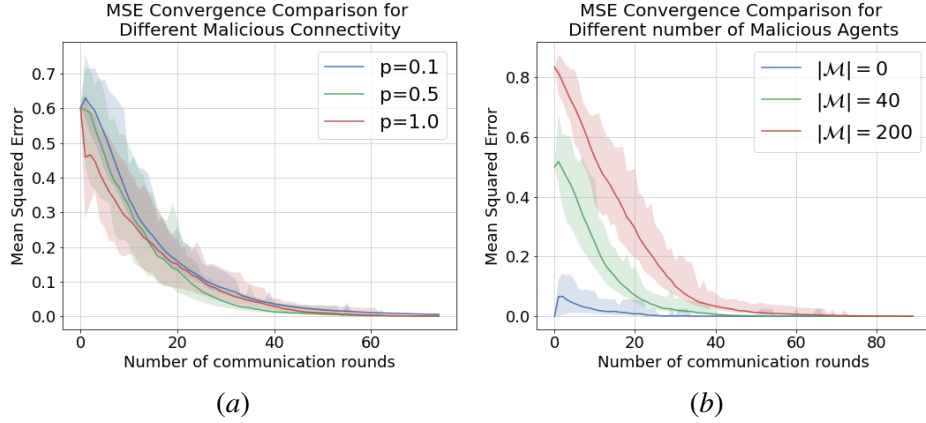


Figure 5: (a) The effect of connectivity of malicious agents with legitimate agents on MSE. Here, p denotes the probability that the edge (m, i) is present in the system for $i \in \mathcal{L}$ and $m \in \mathcal{M}$. As malicious agents are being observed by more agents, disturbance in the initial period decreases. (b) The effect of increasing the number of malicious agents in the system. Even though increasing number of malicious agents slows down convergence, the system still converges to the true values as shown in Theorems 12 and 13.

6. Conclusion

This paper presents a protocol for learning which agents to trust, and the accompanying analysis, for directed multiagent graphs with stochastic observations of trust. Here, the directed nature of the graph presents an important challenge where the out-neighbors of a node cannot directly observe or receive information from it; this leads to a learning dynamic that makes accurate assessment of malicious agents in the network particularly elusive. The learning protocol developed herein specifically addresses this challenge of learning trust in directed graphs and constitutes the novelty of this paper. Since directed graphs often arise in practical multiagent systems due to heterogeneity in sensing and communication, we believe that the learning protocol and theory presented here can support many optimization, estimation, and learning tasks for general multiagent systems.

Acknowledgments

The authors gratefully acknowledge partial support through NSF awards CNS 2147641, CNS-2147694, AFOSR grant number FA9550-22-1-0223, and ONR YIP grant number N00014-21-1-2714.

References

- Orhan Eren Akgün, Arif Kerem Dayı, Stephanie Gil, and Angelia Nedić. Learning trust over directed graphs in multiagent systems (extended version). *ArXiv*, 2022.
- Parsiad Azimzadeh. A fast and stable test to check if a weakly diagonally dominant matrix is a nonsingular m-matrix. *Mathematics of Computation*, 88(316):783–800, 2019.
- Kai Cai and Hideaki Ishii. Average consensus on general strongly connected digraphs. *Automatica*, 48(11):2750–2761, 2012.
- Matthew Cavorsi, Orhan Eren Akgün, Michal Yemini, Andrea Goldsmith, and Stephanie Gil. Exploiting trust for resilient hypothesis testing with malicious robots. *arXiv preprint arXiv:2209.12285*, 2022.
- Alejandro D Dominguez-Garcia and Christoforos N Hadjicostis. Distributed matrix scaling and application to average consensus in directed graphs. *IEEE Transactions on Automatic Control*, 58(3):667–681, 2012.
- Paul Erdős, Alfréd Rényi, et al. On the evolution of random graphs. *Publ. Math. Inst. Hung. Acad. Sci*, 5(1):17–60, 1960.
- Michael J Fischer, Nancy A Lynch, and Michael S Paterson. Impossibility of distributed consensus with one faulty process. *Journal of the ACM (JACM)*, 32(2):374–382, 1985.
- Stephanie Gil, Swarun Kumar, Mark Mazumder, Dina Katabi, and Daniela Rus. Guaranteeing spoof-resilient multi-robot networks. *Autonomous Robots*, 41(6):1383–1400, 2017.
- Stephanie Gil, Cenk Baykal, and Daniela Rus. Resilient multi-agent consensus using wi-fi signals. *IEEE Control Systems Letters*, 3(1):126–131, 2019. doi: 10.1109/LCSYS.2018.2853641.
- Jairo Giraldo, David Urbina, Alvaro Cardenas, Junia Valente, Mustafa Faisal, Justin Ruths, Nils Ole Tippenhauer, Henrik Sandberg, and Richard Candell. A survey of physics-based attack detection in cyber-physical systems. *ACM Computing Surveys (CSUR)*, 51(4):1–36, 2018.
- Alasdair J Graham and David A Pike. A note on thresholds and connectivity in random directed graphs. *Atl. Electron. J. Math*, 3(1):1–5, 2008.
- Leslie Lamport, Robert Shostak, and Marshall Pease. The byzantine generals problem. In *Concurrency: the works of leslie lamport*, pages 203–226. 2019.
- Rui Liu, Fan Jia, Wenhao Luo, Meghan Chandarana, Changjoo Nam, Michael Lewis, and Katia P Sycara. Trust-aware behavior reflection for robot swarm self-healing. In *Aamas*, pages 122–130, 2019.

- Ali Makhdoumi and Asuman Ozdaglar. Graph balancing for distributed subgradient methods over directed graphs. In *2015 54th IEEE Conference on Decision and Control (CDC)*, pages 1364–1371. IEEE, 2015.
- Frederik Mallmann-Trenn, Matthew Cavorsi, and Stephanie Gil. Crowd vetting: Rejecting adversaries via collaboration with application to multirobot flocking. *IEEE Transactions on Robotics*, 38(1):5–24, 2021.
- Angelia Nedić and Alex Olshevsky. Distributed optimization over time-varying directed graphs. *IEEE Transactions on Automatic Control*, 60(3):601–615, 2014.
- Angelia Nedić, Alex Olshevsky, and Michael G Rabbat. Network topology and communication-computation tradeoffs in decentralized optimization. *Proceedings of the IEEE*, 106(5):953–976, 2018.
- Reza Olfati-Saber and Richard M Murray. Consensus problems in networks of agents with switching topology and time-delays. *IEEE Transactions on automatic control*, 49(9):1520–1533, 2004.
- Alex Olshevsky. Efficient information aggregation strategies for distributed control and signal processing. *arXiv preprint arXiv:1009.6036*, 2010.
- Fabio Pasqualetti, Florian Dorfler, and Francesco Bullo. Control-theoretic methods for cyberphysical security: Geometric principles for optimal cross-layer resilient control systems. *IEEE Control Systems Magazine*, 35(1):110–127, 2015.
- S. Pu, W. Shi, J. Xu, and A. Nedić. Push-pull gradient methods for distributed optimization in networks. *IEEE Transactions on Automatic Control*, 66(1):1–16, 2021.
- Michael Rabbat and Robert Nowak. Distributed optimization in sensor networks. In *Proceedings of the 3rd international symposium on Information processing in sensor networks*, pages 20–27, 2004.
- Venkatraman Renganathan and Tyler Summers. Spoof resilient coordination for distributed multi-robot systems. In *2017 International Symposium on Multi-Robot and Multi-Agent Systems (MRS)*, pages 135–141. IEEE, 2017.
- Shreyas Sundaram and Bahman Gharesifard. Distributed optimization under adversarial nodes. *IEEE Transactions on Automatic Control*, 64(3):1063–1076, 2018.
- Shreyas Sundaram and Christoforos N Hadjicostis. Distributed function calculation via linear iterative strategies in the presence of malicious agents. *IEEE Transactions on Automatic Control*, 56(7):1495–1508, 2010.
- Konstantinos I Tsianos, Sean Lawlor, and Michael G Rabbat. Consensus-based distributed optimization: Practical issues and applications in large-scale machine learning. In *2012 50th annual allerton conference on communication, control, and computing (allerton)*, pages 1543–1550. IEEE, 2012a.
- Konstantinos I Tsianos, Sean Lawlor, and Michael G Rabbat. Push-sum distributed dual averaging for convex optimization. In *2012 IEEE 51st IEEE conference on decision and control (cdc)*, pages 5453–5458. IEEE, 2012b.

Chenguang Xi, Van Sy Mai, Ran Xin, Eyad H Abed, and Usman A Khan. Linear convergence in optimization over directed graphs with row-stochastic matrices. *IEEE Transactions on Automatic Control*, 63(10):3558–3565, 2018.

Jie Xiong and Kyle Jamieson. Securearray: Improving wifi security with fine-grained physical-layer information. In *Proceedings of the 19th annual international conference on Mobile computing & networking*, pages 441–452, 2013.

Michal Yemini, Angelia Nedić, Andrea J. Goldsmith, and Stephanie Gil. Characterizing trust and resilience in distributed consensus for cyberphysical systems. *IEEE Trans. Robot.*, 38(1):71–91, 2022.