

欺骗攻击下具备隐私保护的多智能体系统均值趋同控制

应晨铎¹ 伍益明¹ 徐明¹ 郑宁¹ 何熊熊²

摘 要 针对通信网络遭受欺骗攻击的离散时间多智能体系统, 研究其均值趋同和隐私保护问题. 首先, 考虑链路信道存在窃听者的情形, 提出一种基于状态分解思想的分布式网络节点值重构方法, 以阻止系统初始信息的泄露. 其次, 针对所构建的欺骗攻击模型, 利用重构后节点状态信息并结合现有的安全接受广播算法, 提出一种适用于无向通信网络的多智能体系统均值趋同控制方法. 理论分析表明, 该方法能够有效保护节点初始状态信息的隐私, 并能消除链路中欺骗攻击的影响, 实现分布式系统中所有节点以初始值均值趋同. 最后, 通过数值仿真实验验证了该方法的有效性.

关键词 多智能体系统, 均值趋同, 欺骗攻击, 隐私保护, 网络安全

引用格式 应晨铎, 伍益明, 徐明, 郑宁, 何熊熊. 欺骗攻击下具备隐私保护的多智能体系统均值趋同控制. 自动化学报, 2023, 49(2): 425–436

DOI 10.16383/j.aas.c210889

Privacy-preserving Average Consensus Control for Multi-agent Systems Under Deception Attacks

YING Chen-Duo¹ WU Yi-Ming¹ XU Ming¹ ZHENG Ning¹ HE Xiong-Xiong²

Abstract This paper investigates the average consensus and privacy-preserving problem for discrete-time multi-agent systems with deception attacks. First, considering the situation where there are eavesdroppers on the link channel, a distributed node value reconstruction method based on the idea of state decomposition is proposed to prevent the leakage of the initial information of the system. Then, under the constructed deception attack model, a novel average consensus control method for multi-agent systems with undirected communication networks is proposed, which uses the reconstructed node status information and combines with the secure acceptance and broadcast algorithm. Theoretical analysis shows that the proposed method can effectively protect the privacy of the initial state information of the nodes, and eliminate the influence of deception attacks in the links, and reach a consensus. Finally, numerical simulation experiments verify the effectiveness of the proposed method.

Key words Multi-agent systems, average consensus, deception attack, privacy-preserving, cyber security

Citation Ying Chen-Duo, Wu Yi-Ming, Xu Ming, Zheng Ning, He Xiong-Xiong. Privacy-preserving average consensus control for multi-agent systems under deception attacks. *Acta Automatica Sinica*, 2023, 49(2): 425–436

随着技术的进步与时代的发展, 人工智能已然成为当前自动化研究领域的一片广阔热土. 其中分布式人工智能则是未来人工智能发展的趋势之一. 多智能体系统作为分布式人工智能的重要实现应用, 其已成为许多复杂人工智能系统的核心技术^[1]. 多智能体系统是由多个具有一定感知、计算、执行

和通信能力的智能个体组成的网络系统^[2]. 目前, 多智能体系统已经被广泛应用于与日常生活及工业生产息息相关的领域, 例如: 无人机协同编队^[3]、智能城市交通^[1]、智能电网^[4]等. 趋同问题作为多智能体系统分布式协作控制领域中最基础的研究方向之一, 是指在没有控制中心的情况下, 系统中每个智能体 (或节点) 仅使用邻居间相互广播的状态信息, 将智能体动力学方程与通信网络拓扑耦合成复杂网络, 并使用合适的分布式控制算法, 从而在有限时间内实现所有智能体状态值的一致或同步.

然而, 由于多智能体系统所具有的开放式网络环境、通信渠道种类单一、节点同构性高且单个节点资源有限等特性, 使得网络中通讯链路容易被恶意第三方窃听或破坏. 因此如何在恶意网络环境下实现节点之间状态信息的隐私保护和精准趋同, 已成为多智能体系统研究的新挑战. 具体地, 多智能

收稿日期 2021-09-15 录用日期 2022-04-28

Manuscript received September 15, 2021; accepted April 28, 2022

国家自然科学基金 (61803135, 61873239, 62073109), 浙江省公益技术应用研究项目 (LGF21F020011) 资助

Supported by National Natural Science Foundation of China (61803135, 61873239, 62073109) and Zhejiang Provincial Public Welfare Research Project of China (LGF21F020011)

本文责任编辑 秦家虎

Recommended by Associate Editor QIN Jia-Hu

1. 杭州电子科技大学网络空间安全学院 杭州 310018 2. 浙江工业大学信息工程学院 杭州 310023

1. School of Cyberspace, Hangzhou Dianzi University, Hangzhou 310018 2. College of Information Engineering, Zhejiang University of Technology, Hangzhou 310023

体系统趋同控制在实际应用中面临两个关键问题: 1) 节点自身初始状态信息的隐私泄露问题; 2) 节点间的通信链路可能会遭受网络攻击的问题, 如拒绝服务 (Denial-of-service, DoS) 攻击、欺骗攻击等。

在过去 10 年, 已有较多研究人员针对节点初始状态值的隐私问题开展相关的工作。研究的目的是在确保多智能体系统趋同的基础上, 避免网络中节点的状态值隐私泄露。一方面, 有研究人员针对迭代趋同问题提出了一种差分隐私迭代同步趋同机制^[5]。但是采用差分隐私机制所带来的收敛状态与期望状态不精确一致问题无法避免。随后, 使用差分隐私机制并针对不同等级隐私需求的分布式趋同方法被提出^[6-8]。这类方法的基本思想是在信息交互过程中用零和随机噪声掩盖真实状态值, 通过精心设计噪声插入过程实现趋同并保护节点隐私。其次, 有学者针对均值趋同问题提出了一种添加伪随机偏移量的隐私保护方法^[9], 它克服了使用差分隐私机制导致的精度下降的不足。此外, 有学者研究利用可观测性的概念, 结合图论和优化工具来保护节点隐私^[10-12]。这类方法的本质是通过调整网络拓扑中的链路权重以减小窃听者推算被窃听节点的观测能力, 从而保护节点的隐私。另一方面, 部分研究人员开始将目光投向逐步应用的同态加密技术。同态加密技术在文献 [13-14] 中被应用于计算加密域的趋同。网络中每个节点仅能获得其他节点交互的加密值, 因此节点的状态值是保密的。但是同态加密技术也有不足之处, 它的计算复杂度非常大导致资源开销显著增加。为了摆脱使用同态加密技术所带来的限制, 有学者将安全多方计算中的方法融入到分布式趋同控制系统中, 例如: 基于加性秘密共享的隐私保护均值趋同算法^[15], 以及基于 Shamir 秘密共享的隐私保护异步均值算法^[16]等。这类基于秘密共享的隐私保护方法虽然减少了计算和通信消耗, 但仍不适用于单个智能体计算和通信资源有限的分布式多智能体系统。最近, Wang^[17]提出了一种基于状态分解的隐私保护均值趋同机制, 主要思想是将每个节点的初始状态分解为两个随机状态值子状态, 让一个子状态扮演分解前原节点的角色参与邻居节点间的信息交互, 而另一子状态则被隐藏起来仅与第一个子状态通信。该机制能够使系统达成均值趋同的目标, 并且保护每个节点的初始状态信息不被泄露。

上述的研究成果均假设在安全理想的网络环境下, 即系统不存在网络攻击的前提下得出的。然而, 在实际应用场景中, 由于智能体的组成部件众多, 组件之间的通信链路和智能体之间的通信链路皆有可能遭受网络攻击, 导致相关的多智能体系统趋同

控制方法不再适用, 这使得研究多智能体系统在各种类型的网络攻击下的安全趋同发展迅速, 并产出了大量的研究成果。目前, 多智能体系统中常见的网络攻击主要有 DoS 攻击^[18-21]和欺骗攻击^[22-25]两种形式。欺骗攻击作为一种典型的网络攻击类型, 在攻击者精心设计的情况下可以巧妙地绕过攻击检测机制的监测, 造成严重的损失。与 DoS 攻击相比, 欺骗攻击更难发现, 同时严重影响数据的完整性^[26]。多智能体系统分布式网络可能遭受例如数据重放、数据篡改、虚假数据注入等不同形式的欺骗攻击导致系统不能达成共识状态。近年来, 学者们从不同的角度入手对欺骗攻击下的多智能体系统趋同问题开展了相关研究并取得了较多的成果。其中, 有学者提出一种基于后退地平线控制方法的新型分布式弹性算法^[22], 解决了攻击者针对控制器-执行器通信渠道重复传送数据的重放式欺骗攻击。此外, 有学者针对传感器-控制器通信渠道提出了一个新的分布式观测器^[23], 通过使用这个观测器来估计相对完整的状态, 然后在反馈协议中使用估计的状态, 最终实现系统在网络攻击下的共识。同样的, 针对传感器-控制器通道, 有研究者提出了一种分布式安全脉冲控制器^[24], 通过引入与每个通信通道相关的随机变量, 实现了存在虚假数据注入形式的欺骗攻击下的趋同。

然而, 上述文献仅考虑隐私保护需求或网络容错功能。例如, Wang^[17]提出的状态分解机制满足了对于节点初始状态值的隐私需求, 但是如果遭受欺骗攻击, 则系统将不能实现预期的均值趋同。文献 [23] 提出的重设计观测器能够抵御网络中可能出现的欺骗攻击, 但是若存在一个只窃听交互信息不产生恶意攻击行为的第三方, 网络中节点的隐私就无法保证。目前, 已有部分同时考虑隐私保护需求和网络容错功能的文献。Li 等^[27]率先开展了分布式多智能体网络在信道攻击下系统全局一致性的工作, 随后在文献 [28] 中更进一步地提出在有限资源条件及隐私保护需求下的高效分布式算法。特别针对分布式系统优化问题, 文献 [29] 提出了一种时变非平衡有向网络环境下差分隐私随机次梯度推送算法。2019 年, Fiore 等^[30]进行了欺骗攻击下满足差分隐私需求的多智能体系统弹性趋同研究工作, 但成果仍存在非精确均值趋同、未考虑节点内部通讯链路的安全状况等可以进一步改进和扩展的地方。随着多智能体系统在关键领域的逐步应用, 如何设计兼顾隐私保护需求和网络容错功能的多智能体系统趋同控制算法成为亟须研究的热点问题, 这也是本文的研究重点。

基于上述研究与总结, 本文主要致力于研究欺

骗攻击下保护节点初始值隐私信息的多智能体系统均值趋同问题, 从而完善和补充现有趋同算法的相关研究成果. 本文围绕系统的容错能力开展研究, 不关注检测和处理网络攻击的能力, 而是关注分布式控制算法在网络攻击下完成预期趋同的鲁棒设计实现. 本文主要贡献包括以下 3 个方面:

1) 不同于文献 [24–25] 对网络攻击的建模, 本文考虑了欺骗攻击在多智能体系统中对不同类型通信链路的攻击特性和发生范围, 提出了广义 f -局部攻击模型. 新攻击模型相比于传统攻击模型考虑的场景更面向实际应用, 具有一定的普适性.

2) 针对提出的广义 f -局部攻击模型欺骗攻击下的无向通信拓扑结构多智能体系统, 提出一种基于状态分解机制的隐私保护与弹性均值趋同控制算法. 相比于文献 [17], 本文提出的算法在保证节点初始状态值的隐私基础上进一步实现了网络容错功能. 理论分析证明在一定网络鲁棒性条件下系统可容忍一定数量的通信链路遭受欺骗攻击破坏, 并最终实现均值趋同.

3) 相比于文献 [24, 26], 本文拓宽了欺骗攻击的抵御范围, 从单一通信链路防御扩展到多类型不同链路的全面防御, 针对不同链路传输的多类型数据的篡改进行相应的处理.

本文内容结构如下: 第 1 节介绍本文需要用到的图论知识和网络鲁棒性等相关预备知识; 第 2 节对广义 f -局部攻击模型欺骗攻击进行建模以及给出相关假设, 随后给出引理及其证明; 第 3 节提出欺骗攻击下具备隐私保护的均值趋同控制算法, 并分别对算法在欺骗攻击下的系统均值趋同以及隐私保护能力进行分析; 第 4 节通过四组数值仿真实验验证所提算法的有效性; 第 5 节是总结与展望.

1 预备知识

1.1 图论知识

考虑一个通信网络拓扑结构为无向加权图 $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathbf{A})$ 的多智能体系统, 系统由 N 个节点组成, 其中节点集和边集分别表示为 $\mathcal{V} = \{v_1, v_2, \dots, v_N\}$ 和 $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$. 两个节点之间的连接关系用邻接矩阵 (权重矩阵) $\mathbf{A} = [a_{ij}] \in \mathbf{R}^{N \times N}$ 表示, 如果 $(v_j, v_i) \in \mathcal{E}$, 表示节点 v_i 和节点 v_j 之间存在着信息交互, 则 $a_{ij} > 0$, 否则 $a_{ij} = 0$. 无向图的定义规定图的邻接矩阵是对称的, 即如果 $(v_j, v_i) \in \mathcal{E}$, 则有 $(v_i, v_j) \in \mathcal{E}$, 且 $a_{ji} = a_{ij}$. 同时, 本文不考虑节点自环情况, 即 $a_{ii} = 0$. $\mathcal{N}_i = \{v_j \in \mathcal{V} | (v_j, v_i) \in \mathcal{E}\}$ 代表节点 v_i 的邻居节点集.

注 1. $\mathcal{A} \setminus \mathcal{B}$ 表示集合 $\{x : x \in \mathcal{A}, x \notin \mathcal{B}\}$.

除了上述提到的图论相关知识以外, 本文涉及到文献 [31] 提出的 r -可达集和 r -鲁棒图概念. 值得一提的是, 其中的 r -鲁棒图概念之后被文献 [32] 加以延伸, 提出了强 r -鲁棒图概念.

定义 1 (r -可达集) [31]. 对于一个图 \mathcal{G} , 其节点的一个子集 \mathcal{S} , 如果 $\exists v_i \in \mathcal{S}$ 满足 $|\mathcal{N}_i \setminus \mathcal{S}| \geq r$, 其中 $r \in \mathbf{Z}_{\geq 1}$, 则称该子集 \mathcal{S} 为 r -可达集合.

定义 2 (强 r -鲁棒图) [32]. 对于图 \mathcal{G} , 如果对任意一个非空子集 $\mathcal{S} \subseteq \mathcal{V}$, \mathcal{S} 是 r -可达集合或者 $\exists v_i \in \mathcal{S}$ 满足 $\mathcal{V} \setminus \mathcal{S} \subseteq \mathcal{N}_i$, 其中 $r \in \mathbf{Z}_{\geq 1}$ 并且 $r \leq \lceil N/2 \rceil$, 则称 \mathcal{G} 为强 r -鲁棒图.

鲁棒性概念是通信网络拓扑图的连通性衡量标准. 根据本文所需, 笔者将 r -可达集和强 r -鲁棒图概念分别修改为如下定义:

定义 3 (r -链路可达集). 对于一个图 \mathcal{G} , 其节点的一个子集 \mathcal{S} , 如果 $\exists v_i \in \mathcal{S}$ 满足 $|\mathcal{C}_i \setminus \mathcal{C}_{\mathcal{S}}| \geq 2r$, 其中 $r \in \mathbf{Z}_{\geq 1}$, \mathcal{C}_i 表示与节点 v_i 相关的所有通信链路, 则称该子集 \mathcal{S} 为 r -链路可达集.

注 2. $\mathcal{C}_{\mathcal{S}}$ 表示与子集 \mathcal{S} 相关的所有通信链路, 当集合中仅包含一个节点时, 该集合相关的通信链路为 0, 即 $|\mathcal{C}_{\mathcal{S}}| = 0$. 当集合中包含多个节点时该集合相关的通信链路为集合中节点相连的通信链路.

定义 4 (强 r -链路鲁棒图). 对于图 \mathcal{G} , 如果对任意一个非空子集 $\mathcal{S} \subseteq \mathcal{V}$, \mathcal{S} 是 r -链路可达集或者 $\exists v_i \in \mathcal{S}$, 满足 $\mathcal{C}_{\mathcal{V}} \setminus \mathcal{C}_{\mathcal{S}} \subseteq \mathcal{C}_i$, 其中 $r \in \mathbf{Z}_{\geq 1}$ 并且 $r \leq \lceil N/2 \rceil$, 则称 \mathcal{G} 为强 r -链路鲁棒图.

1.2 分布式均值趋同

在一个由 N 个智能体组成的多智能体分布式网络 \mathcal{G} 中, 每一个智能体 $v_i \in \mathcal{V}$ 在时刻 k 拥有一个标量状态值 $x_i[k]$. 分布式均值趋同问题的目标即每个智能体节点仅使用从邻居节点接收到的局部信息, 使得网络中智能体节点状态值 $x_i[k]$ 最终收敛于它们的初始平均值 $x_a := \sum_{i=1}^N x_i[0]/N$. 下述定义形式化了分布式均值趋同问题的目标.

定义 5 (均值趋同) [32]. 对于一个由 N 个智能体节点组成的拓扑图 \mathcal{G} 所代表的多智能体分布式网络, 如果对每一个节点的初始状态值 $x_i[0]$, $i = 1, 2, \dots, N$, 满足 $\lim_{k \rightarrow +\infty} x_i[k] = x_a$, $\forall i \in \{1, 2, \dots, N\}$, 则称该网络达成均值趋同.

1.3 状态分解

状态分解方法由文献 [17] 在 2019 年提出, 是一种噪声避免的隐私保护方法. 该方法的主要思想是将每个节点的状态值 x_i 分解成两个子状态, 分别用 x_i^α 和 x_i^β 表示. 具体来说, 节点初始状态的子状态

值 $x_i^\alpha[0]$ 和 $x_i^\beta[0]$ 在满足条件: $x_i^\alpha[0] + x_i^\beta[0] = 2x_i[0]$ 的前提下可以取为任意实数.

为便于对状态分解方法的理解以及后续工作的展开, 本文以 N 个节点的抽象拓扑连通图为例, 阐述状态分解方法的主要思想, 抽象拓扑图如图 1 所示. 由图 1 可以看出, 图中节点 v_i 经过状态分解方法处理后生成两个子状态, 即 x_i^α 和 x_i^β . 其中, 第一个子状态 x_i^α 的作用与原节点相同, 即接收并广播邻居节点和节点自身的状态值信息, 同时对于节点 v_i 的邻居节点来说 x_i^α 是唯一可见的. 第二个子状态 x_i^β 的作用是与第一个子状态 x_i^α 进行交互, 但是对于节点 v_i 的邻居节点来说 x_i^β 是不可见的. 本文假设两个子状态 x_i^α 和 x_i^β 之间的权重称为耦合权重, 表示为 $a_{i,\alpha\beta}$ 以及 $a_{i,\beta\alpha}$, 并且两者是相同的, 即 $a_{i,\alpha\beta} = a_{i,\beta\alpha}$.

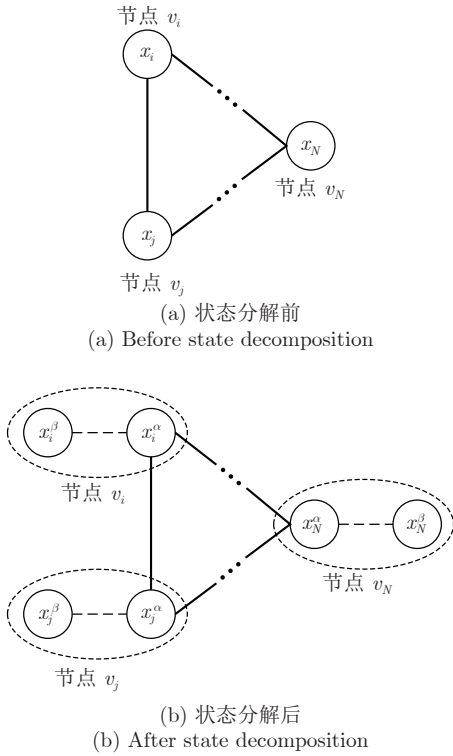


图 1 状态分解方法示例图

Fig.1 Example diagram of state decomposition method

2 问题描述

本文考虑的研究对象为由 N 个智能体组成的一阶离散多智能体系统, 系统中的智能体遵循预先设定的动力学方程:

$$x_i[k] = \begin{cases} \theta_i[k], & k = 1 \\ \varepsilon_i x_i[k-1] + (1 - \varepsilon_i) u_i[k], & k > 1 \end{cases} \quad (1)$$

式中, $x_i[k]$ 表示智能体 v_i 在 k 时刻的状态值, θ_i 是下文将要设计的隐私保护状态值, ε_i 为控制增益, u_i 为下文将要设计的控制输入. 需要注意的是, 本文所考虑的智能体节点由控制器、执行器、传感器三部分构成, 控制器负责接收由传感器发送来的邻居节点数据并经过处理发送给执行器, 执行器收到控制器发送的数据对其进行处理并最终更新节点的状态值, u_i 在数据交互过程中可能受欺骗攻击的影响而无法传输原始的正确状态值.

2.1 欺骗攻击模型

近年来随着云技术的发展, 大量的分布式多智能体系统开始采用基于云辅助的协同控制^[33]. 基于云辅助的协同控制通过一个开放的网络环境进行数据传输. 本文考虑针对通信链路传输数据的数据篡改欺骗攻击. 基于云辅助的协同控制中通信链路可分为两类, 一类为智能体组件 (如控制器、执行器等) 之间的节点内部通信链路, 另一类为节点间信息交互的外部通信链路. 在多智能体系统分布式网络中, 不同类型的通信链路遭遇欺骗攻击的示意图如图 2 所示.

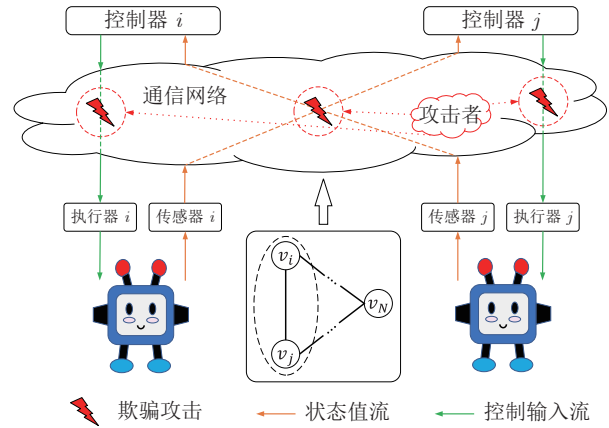


图 2 欺骗攻击下多智能体系统分布式网络示意图

Fig.2 The diagram of the multi-agent system distributed network under deception attacks

图 2 中, 节点 v_i 和节点 v_j 为多智能体系统网络中相邻的一对邻居节点. 对于系统中任一智能体, 一条智能体组件之间的内部通信链路 (控制器至执行器) 对应一条节点间信息交互的外部通信链路 (传感器至控制器).

当多智能体系统启动并完成初始化阶段后, 攻击者探测到通信链路中的数据传输行为, 开始针对不同的通信链路发动攻击. 在 k 时刻欺骗攻击对节点 v_j 至节点 v_i 的外部通信链路中传输的数据进行篡改, 其数学模型可以表示为:

$$\tilde{x}_j[k] = x_j[k] + p_{ij}[k] x_j^\alpha[k] \quad (2)$$

式中, $p_{ij}[k]$ 为攻击参数, 当欺骗攻击在节点 v_j 至节点 v_i 的外部通信链路中发生时 $p_{ij}[k] = 1$, 否则 $p_{ij}[k] = 0$, $x_j^a[k]$ 表示欺骗攻击篡改的数值. 在 k 时刻欺骗攻击对节点 v_i 的内部通信链路中传输的数据进行篡改, 其数学模型可以表示为:

$$\tilde{x}_i[k] = x_i[k] + p_i[k]x_i^a[k] \quad (3)$$

式中, $p_i[k]$ 为攻击参数, 当欺骗攻击在节点 v_i 的内部通信链路中发生时 $p_i[k] = 1$, 否则 $p_i[k] = 0$, $x_i^a[k]$ 表示欺骗攻击篡改的数值. 攻击者的目标是通过改变精心设计的 $x_i^a[k]$ 得到任意篡改后的数值 $\tilde{x}_i[k]$, 使得系统中各个节点接收到篡改后的错误状态值信息, 进而达到破坏多智能体系统分布式网络趋同的目的.

考虑到实际应用环境中恶意攻击者自身资源的有限性以及经济效益的回报率, 通常真实环境下的网络攻击存在着一定的约束. 因此, 本文将首先介绍 f -局部攻击的模型定义, 该攻击模型在现有文献 [23, 25, 34] 中已被广泛采用. 其概念具体定义如下:

定义 6 (f -局部攻击)^[34]. 对于分布式网络中的任一智能体, 在任意时刻其邻居节点中恶意节点的数量小于等于 f , 则称此类攻击模型为 f -局部攻击.

传统攻击模型可以通过节点的度等网络特性确定网络攻击约束条件 f 的数值. 然而, 本文所考虑的欺骗攻击发生范围已不再局限于恶意邻居节点数量, 而是同时考虑了节点间信息交互的外部通信链路以及智能体组件内部通信链路遭受欺骗攻击的数量. 节点的度等网络特性通常仅能表示节点间信息交互的外部通信链路数量, 使得传统攻击模型不再适用. 同时, 本文攻击模型相比较于传统攻击模型适用范围更广泛, 从而具有一定的普适性以及实际应用价值. 因此, 本文拓展了传统的 f -局部攻击模型, 提出了广义 f -局部攻击模型, 具体定义如下:

定义 7 (广义 f -局部攻击). 对于分布式网络中的任一智能体, 如果在该节点相关的通信链路中, 任意时刻遭受欺骗攻击的通信链路数量之和小于等于 f , 则称此类攻击模型为广义 f -局部攻击. 数学模型表示为:

$$\sum_{v_j \in \mathcal{N}_i} p_{ij}[k] + \sum p_i[k] \leq f \quad (4)$$

注 3. 本文欺骗攻击针对的对象不同于传统模型, 不考虑邻居恶意节点数量而是关注相关通信链路被破坏数量, 因此与节点的度等网络特性无关.

2.2 系统假设

通过结合上述给出的广义 f -局部攻击模型以及欺骗攻击针对多智能体系统分布式网络的发生特

性, 本文对所研究的基于云辅助协同控制的一阶离散多智能体系统做出以下假设:

假设 1. 系统中与任意一个智能体相关的所有通信链路中在任意 $k > 2$ 时刻至多有 f 条通信链路同时遭受欺骗攻击, 即满足定义 7 攻击模型. 具体来说, 即对于任意 $v_i \in \mathcal{V}$, 在任意 $k > 2$ 时刻, 都有下式成立:

$$|C_i \cap A_i| \leq f, \forall v_i \in \mathcal{V} \quad (5)$$

式中, A_i 表示与节点 v_i 相关的遭受欺骗攻击的通信链路.

注 4. 不同于单一链路攻击仅考虑对某一种数据 (如节点状态值或控制输入等) 进行篡改, 本文攻击模型考虑与节点相关的所有被攻击通信链路, 并在后续控制算法设计中处理不同类型被篡改的数据.

假设 2. 系统在 $k = 1$ 与 $k = 2$ 时刻的网络环境是安全的, 即没有恶意攻击的存在. 并且对于所有 $v_i \in \mathcal{V}$, $a_{i, \alpha\beta}$ 的值是未知的.

注 5. 在真实应用环境中, 攻击者对目标发动的攻击往往需要一定的部署和准备时间, 因此本文合理假设在系统初始运行阶段网络环境暂时是安全的.

根据上述假设, 可以得出系统具备如下属性:

引理 1. 多智能体系统 (1) 在 $k = 1$ 时刻执行状态分解方法更新后重组各节点子状态, 其均值不变.

证明. 由第 1.3 节可知, 节点初始状态值 $x_i[0]$ 经过状态分解后产生的子状态值 $x_i^\alpha[0]$ 和 $x_i^\beta[0]$ 需要满足条件: $x_i^\alpha[0] + x_i^\beta[0] = 2x_i[0]$. 本文引用文献 [17] 的状态更新方程并加以修改, 得到系统中节点在 $k = 1$ 时刻的动力学方程为:

$$\begin{cases} \theta_i^\alpha[1] = x_i^\alpha[0] + \sum_{v_j \in \mathcal{N}_i} a_{ij}(x_j^\alpha[0] - x_i^\alpha[0]) + \\ \quad a_{i, \alpha\beta}(x_i^\beta[0] - x_i^\alpha[0]) \\ \theta_i^\beta[1] = x_i^\beta[0] + a_{i, \alpha\beta}(x_i^\alpha[0] - x_i^\beta[0]) \end{cases} \quad (6)$$

此时, 得到重组后的节点状态值为:

$$x_i[1] = \theta_i[1] = \frac{1}{2}(\theta_i^\alpha[1] + \theta_i^\beta[1]) =$$

$$\frac{1}{2} \left\{ x_i^\alpha[0] + x_i^\beta[0] + \sum_{v_j \in \mathcal{N}_i} a_{ij}(x_j^\alpha[0] - x_i^\alpha[0]) + \right. \\ \left. a_{i, \alpha\beta}(x_i^\beta[0] - x_i^\alpha[0]) + \right. \\ \left. a_{i, \alpha\beta}(x_i^\alpha[0] - x_i^\beta[0]) \right\} =$$

$$\frac{1}{2} \left\{ x_i^\alpha[0] + x_i^\beta[0] + \sum_{v_j \in \mathcal{N}_i} a_{ij}(x_j^\alpha[0] - x_i^\alpha[0]) \right\} \quad (7)$$

由此,可以得到:

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N x_i[1] &= \frac{1}{2N} \sum_{i=1}^N (x_i^\alpha[0] + x_i^\beta[0]) + \\ &\quad \frac{1}{2N} \sum_{i=1}^N \left\{ \sum_{v_j \in \mathcal{N}_i} a_{ij} (x_j^\alpha[0] - x_i^\alpha[0]) \right\} \end{aligned} \quad (8)$$

根据无向拓扑图的特性可以得到 $a_{ij} = a_{ji}$, 对于任意 $v_i, v_j \in \mathcal{V}$, 有:

$$a_{ij} (x_j^\alpha[0] - x_i^\alpha[0]) = -a_{ji} (x_i^\alpha[0] - x_j^\alpha[0]) \quad (9)$$

将式 (9) 代入式 (8), 易得:

$$\sum_{i=1}^N \left\{ \sum_{v_j \in \mathcal{N}_i} a_{ij} (x_j^\alpha[0] - x_i^\alpha[0]) \right\} = 0 \quad (10)$$

最后, 将式 (10) 代入式 (8), 可得:

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N x_i[1] &= \frac{1}{2N} \sum_{i=1}^N (x_i^\alpha[0] + x_i^\beta[0]) = \\ &\quad \frac{1}{N} \sum_{i=1}^N x_i[0] \end{aligned} \quad (11)$$

综上, 由式 (11) 不难看出, 在经过一次状态分解方法后重构各节点子状态, 系统节点状态均值仍维持不变. \square

根据上述提出的欺骗攻击模型以及相关的系统假设, 本文的研究目标是设计一种控制算法, 使得: 1) 多智能体系统分布式网络达到均值趋同的同时实现智能体初始状态值的隐私保护; 2) 算法无论是面对特定通信链路的攻击抑或是面对针对不同类型的通信链路同时发动攻击均能够拥有一定的弹性.

3 算法设计

3.1 弹性分布式检索

为便于后续控制算法的描述与理解, 首先对弹性分布式检索的概念进行定义.

定义 8 (弹性分布式检索)^[32]. 如果系统中任意节点 $v_i \in \mathcal{V}$ 可以检索到所有其他节点的初始状态值, 即 $x_j[0]$, $v_j \in \mathcal{V} \setminus v_i$, 则称该欺骗攻击下的多智能体系统分布式网络中节点实现了弹性分布式检索.

3.2 基于状态分解的弹性均值趋同算法

为了实现欺骗攻击下具备隐私保护的多智能体系统均值趋同, 本文提出了基于状态分解的弹性均值趋同算法, 该算法的伪代码如算法 1 所示.

具体地, 多智能体系统分布式网络中每个节点

v_i 会使用一个永久存储向量 $s^i[k] = [s_1^i[k], \dots, s_{\hat{N}}^i[k]]$, $\hat{N} \geq N$, 该存储向量的用处是记录下节点从邻居节点 $v_j \in \mathcal{N}_i$ 接收到的与最终接受的状态值. 存储向量中元素 $s_n^i[k]$, $n \in \{1, 2, \dots, \hat{N}\}$ 表示节点 v_n 在节点 v_i 内存中所记录的状态值. 在 $k = 1$ 时刻, 存储向量被创建为 $s^i[1] = [\]_{1 \times \hat{N}}$, 其中 $[\]$ 代表一个空向量. 需要注意的是, 本文假设节点知道但不限制网络中节点数量的上限 $\hat{N} \geq N$, 这样做的目的是更便于真实环境下的实际应用^[28].

算法 1. 基于状态分解的弹性均值趋同算法

输入. 节点初始状态值 $x_i[0]$.

输出. 节点更新后状态值 $x_i[k]$, 存储向量 $s^i[k] = [s_1^i[k], \dots, s_{\hat{N}}^i[k]]$, $\hat{N} \geq N$.

步骤 1. 初始化节点 v_i , 随后根据状态分解方法分解为两个子状态.

步骤 2.

1) if $k = 1$ then;

2) for $v_i \in \mathcal{V}$ do;

3) 节点 v_i 接收 $x_j^\alpha[0]$, $j \in \mathcal{N}_i$ 并且使用式 (3) 更新子状态值, 随后使用式 (4) 重组两个子状态得到 $x_i[1]$;

4) end;

5) 节点 v_i 创建一个永久存储向量 $s^i[1] = [\]_{1 \times \hat{N}}$, 其中 $\hat{N} \geq N$, 并且设置 $s_1^i[1] = x_i[1]$;

6) end.

步骤 3.

1) if $k = 2$ then;

2) for $v_i \in \mathcal{V}$ do;

3) 节点 v_i 广播 $s^i[1]$ 至节点 $v_j \in \mathcal{N}_i$, 同时接收 $s^j[1]$ 并且更新自身存储向量 $s_j^i[2] = s_j^j[1]$, $v_j \in \mathcal{N}_i$;

4) end;

5) end.

步骤 4.

1) if $k > 2$ then;

2) 节点 v_i 向所有邻居节点广播存储向量 $s^i[k-1]$;

3) for $n \in \{1, 2, \dots, \hat{N}\}$ do;

4) 若节点 v_i 从至少 $f+1$ 条内部通信链路接收到一个相同值 $s_n^j[k-1]$, $v_j \in \mathcal{N}_i$, 则接受这个新的值并且将其更新至存储向量 $s_n^i[k]$;

5) end;

6) end.

步骤 5. 节点 v_i 根据式 (1) 得到状态值 $x_i[k]$.

网络中的每个节点 v_i 在 $k = 1$ 时刻将会执行状态分解方法下的更新方程. 在 $k = 2$ 时刻节点广播 $s^i[1]$ 给所有的邻居节点 $v_j \in \mathcal{N}_i$, 同时接收 $s^j[1]$ 并且更新其自身存储向量, 即 $s_j^i[2] = s_j^j[1]$, $v_j \in \mathcal{N}_i$. 这里本文假设通信网络中没有通信延迟存在, 也就是

每个节点 v_i 在 k 时刻同时发送自身的存储向量 $s^i[k-1]$ 与接收邻居节点发送来的存储向量信息 $s^j[k-1]$. 为了抵抗 $k > 2$ 时刻发生的欺骗攻击, 引入了一种少数服从多数的更新机制, 即对于节点 v_i 在 $k > 2$ 时刻, 仅接受由超过 $f+1$ 条内部通信链路传输的相同的邻居节点状态值, 并将这些值更新至存储向量 $s_n^i[k]$ 中, 否则不变. 最后, 节点 v_i 根据当前时刻的存储向量 $s_n^i[k]$ 以及动力学方程 (1) 得到节点更新后状态值 $x_i[k]$.

3.3 均值趋同分析

接下来将要给出本文的主要结论, 在此之前还需要给出下述引理知识.

引理 2. 对于一个多智能体系统分布式网络 \mathcal{G} , 如果该网络满足强 $(2f+1)$ -链路鲁棒图, 则在网络中存在符合广义 f -局部攻击模型的欺骗攻击时, 通过执行至少 $\hat{K} \geq N-1$ 次算法迭代, 网络中的任意节点 $v_i \in \mathcal{V}$ 都能够实现弹性分布式检索.

证明. 上述假设的广义 f -局部攻击模型根据图 2 具体可分为 2 种情况: 1) 所有 f 条遭受欺骗攻击的通信链路是同一类型, 即全部攻击发生在图中靠外两线或者靠内两线处; 2) 遭受欺骗攻击的通信链路是不同类型, 即一部分攻击发生在图中靠内两线处, 另一部分攻击发生在靠外两线处. 针对第 1 种情况, 无论是发生在传感器至控制器的外部通信链路或控制器至执行器的内部通信链路, 根据定义 2, 可以得出此时节点 v_i 至少拥有 $2f+1$ 个邻居节点所对应的 $2f+1$ 对通信链路传送过来的信息, 其中最多存在 f 对通信链路遭受欺骗攻击导致虚假错误数据被节点 v_i 接收. 但是根据本文算法, 至多 f 个虚假错误数据并不会被节点 v_i 所接受, 其所对应的存储向量不会发生改变. 并且, 最终在某一时刻 $k \leq \hat{K}$ 将会接收到由至少 $f+1$ 条内部通信链路发送的正确存储向量, 将其正确的状态值更新到自己的存储向量中. 针对第 2 种情况, 本文假设节点内部一条传感器至控制器外部通信链路对应一条控制器至执行器内部通信链路. 因此, 如果传感器至控制器通信链路或控制器至执行器通信链路遭受攻击时, 可以看作与之对应的通信链路同时遭受了攻击. 此时, 节点不同类型的通信链路中遭受攻击链路总数之和仍然至多为 f 对, 后续证明过程与第 1 种情况相同. 迭代次数 \hat{K} 的证明与文献 [32] 相同. \square

注 6. 相比于文献 [32] 证明的最少迭代次数 $\hat{K} \geq N-2$, 本文所提出的算法 1 在设计过程中进一步考虑了保护系统中节点的初始状态值隐私, 因此引入了状态分解的方法, 算法迭代所需的最少次数变为 $\hat{K} \geq N-1$.

下面给出本文第 1 个主要结果.

定理 1. 考虑存在定义 7 欺骗攻击下的多智能体系统网络 (1), 在满足假设 1 和假设 2 条件下, 若其网络拓扑结构满足强 $(2f+1)$ -链路鲁棒图, 则网络中每个节点在算法 1 下进行不少于 $\hat{K} \geq N-1$ 次更新后, 所有节点实现初始值均值趋同.

证明. 本文假设多智能体系统分布式网络同步进行状态值更新. 网络中节点 v_i 在时刻 $k > 1$, 通过执行本文提出的算法 1 更新其存储向量 $s^i[k]$, 并且随着时间的推移从其邻居节点那里接收并接受更多的网络中其他节点的初始状态值, 本文将节点 v_i 在 $k > 1$ 时刻得到的控制输入定义如下:

$$u_i[k] = \frac{\sum (s_n^i[k])}{\lambda[k]}, \quad n \in \mathbf{S}^i[k] \quad (12)$$

式中, $\mathbf{S}^i[k]$ 表示存储向量 $s^i[k]$ 中非空元素的索引集, 索引集的基数由 $\lambda[k] = |\mathbf{S}^i[k]|$ 给出. 在这里, 笔者重写动力学方程 (1) 为:

$$\begin{cases} x_i[k] = \varepsilon_i x_i[k-1] + (1 - \varepsilon_i) u_i[k], & k > 1 \\ u_i[k] = \frac{\sum s_n^i[k]}{\lambda[k]}, & n \in \mathbf{S}^i[k] \end{cases} \quad (13)$$

式中, 控制增益的取值范围为 $0 \leq \varepsilon_i < 1$. 当 $\varepsilon_i = 0$ 时, 节点更新后的状态值 $x_i[k]$ 与控制输入 $u_i[k]$ 相等. 需要说明的是, 控制增益 ε_i 的取值在取值范围内可以是任意的, 但是它的大小可能决定了节点的动力学方程对状态值瞬时变化的敏感性. 根据引理 2, 一个满足强 $(2f+1)$ -链路鲁棒图的多智能体系统分布式网络 \mathcal{G} , 通过执行至少 \hat{K} 次算法迭代, 网络中的任意节点 $v_i \in \mathcal{V}$ 都能够实现弹性分布式检索. 同时, 根据式 (12) 易得控制输入 $u_i[k]$ 是节点在时刻 k 接收并接受的初始状态值 $x_i[0]$, $v_i \in \mathcal{V}$ 的线性组合. 并且, 因为每个节点 v_i 都将在 $t = \hat{K}$ 前接收到所有其他节点的初始状态值, 因此, 控制输入 $u_i[k]$ 最终将会渐近收敛到 x_a , 即:

$$\lim_{k \rightarrow +\infty} u_i[k] = x_a \quad (14)$$

此时, 在式 (13) 等号左边加上和减去 x_a , 等号右边加上和减去 $\varepsilon_i x_a$, 可得:

$$(x_i[k] - x_a) - (u_i[k] - x_a) = \varepsilon_i (x_i[k-1] - x_a) - \varepsilon_i (u_i[k] - x_a) \quad (15)$$

最后, 当 $k \rightarrow +\infty$ 时, 根据式 (14), 可得:

$$x_i[k] - x_a = \varepsilon_i (x_i[k-1] - x_a) \quad (16)$$

因此, 当 $0 \leq \varepsilon_i < 1$ 时上述系统满足舒尔稳定, 所有节点实现初始值均值趋同. \square

3.4 隐私保护分析

本节将对本文所提算法执行过程中单个节点初始状态值的隐私保护进行分析. 本文考虑两种针对节点初始状态值的隐私窃听者: 好奇窃听者和第三方窃听者. 好奇窃听者指的是网络中遵循预先规定的控制协议算法执行状态更新, 但是想要通过接收到的信息计算推测邻居节点的状态值的一类好奇节点. 第三方窃听者指的是能够获取整个网络拓扑结构信息 (包括节点间链路权重) 的外部窃听者, 并且这些第三方窃听者能够抓取网络中节点之间通信链路所传输的信息数据.

具体地, 第三方窃听者能够同时抓取多个节点间通信链路上所传输的数据, 而好奇窃听者只能获得与该节点有信息交互的邻居节点的信息, 因此, 第三方窃听者通常比好奇窃听者更值得关注对其隐私泄露的防御. 相比较于第三方窃听者, 好奇窃听者也具有自己的优势, 即该好奇窃听者自身的初始状态值是已知的.

本文采用与文献 [17] 一致的隐私定义, 具体描述如下:

定义 9 (节点初始状态值隐私)^[17]. 如果窃听者无法以任何确定的精度准确估计节点初始状态信息 $x_i[0]$ 的值, 则称节点 v_i 的初始状态值 $x_i[0]$ 的隐私得到了保护.

下面给出本文第 2 个结果:

定理 2. 考虑欺骗攻击下多智能体系统分布式网络 (1), 在满足假设 1 和假设 2 条件下, 若网络中节点采用算法 1 进行状态值更新, 则该网络中所有节点的初始状态值信息具备隐私保护.

证明. 首先, 分析系统中存在好奇窃听者 v_j 的情况. 本文使用 $x_i[0]$ 表示好奇窃听者尝试去推测的节点 v_i 的初始状态值. 根据状态分解方法, 节点 v_i 的初始状态值 $x_i[0]$ 可以由 $x_i^\alpha[0] + x_i^\beta[0] = 2x_i[0]$ 推得. 其中, $x_i^\alpha[0]$ 将会被好奇窃听者获取. 因此, 估计节点初始状态信息 $x_i[0]$ 的值等同于估计 $x_i^\beta[0]$ 的值. 对于好奇窃听者 v_j 来说, $x_i^\alpha[0]$ 与 $x_j^\alpha[0]$ 的值是已知的并且可以被看作是一个常量. 因此, 节点 v_i 初始状态值的隐私泄露可以被定义为:

$$I(x_i^\alpha[1]; x_i^\beta[0] | x_i^\alpha[0], x_j^\alpha[0])$$

根据式 (3), 条件互信息可以表示为:

$$I\left(\sum_{v_j \in \mathcal{N}_i} a_{ij}(x_j^\alpha[0] - x_i^\alpha[0]) + a_{i, \alpha\beta}(x_i^\beta[0] - x_i^\alpha[0]); x_i^\beta[0] | x_i^\alpha[0], x_j^\alpha[0]\right)$$

此时, 由于好奇窃听者 v_j 对于系统整体网络拓扑并不清楚且无法确定节点 v_i 的内部权重 $a_{i, \alpha\beta}$, 可得:

$$I(x_i^\alpha[1]; x_i^\beta[0] | x_i^\alpha[0], x_j^\alpha[0]) = 0$$

根据上述推导可知, 该情况下好奇窃听者无法准确估计节点 v_i 的初始状态值, 因此节点 v_i 的初始状态值 $x_i[0]$ 得到了隐私保护.

接着, 分析系统存在第三方窃听者的情况. 相比较于好奇窃听者, 第三方窃听者可以获得包括节点间链路权重与链路所对应的邻居节点交互的状态值数据等更多信息. 同样的, 节点 v_i 初始状态值的隐私泄露可以被定义为:

$$I(x_i^\alpha[1]; x_i^\beta[0] | x_i^\alpha[0], x_j^\alpha[0], a_{ij})$$

根据式 (3), 条件互信息可以表示为:

$$I(a_{i, \alpha\beta}(x_i^\beta[0] - x_i^\alpha[0]); x_i^\beta[0] | x_i^\alpha[0])$$

此时, 由于第三方窃听者无法确定节点 v_i 的内部权重 $a_{i, \alpha\beta}$, 可得:

$$I(x_i^\alpha[1]; x_i^\beta[0] | x_i^\alpha[0], x_j^\alpha[0], a_{ij}) = 0$$

根据上述推导可知, 该情况下第三方窃听者无法准确估计节点 v_i 的初始状态值, 因此节点 v_i 的初始状态值 $x_i[0]$ 同样得到了隐私保护. \square

4 数值仿真实验

本文实验采用了与文献 [17] 相同的模拟仿真实验平台进行对比说明, 同时参考文献 [24, 25, 32] 构建了简单抽象化通信拓扑结构, 通过模拟一些数值仿真实验来验证本文算法的有效性以及特性.

考虑一个由 6 个节点组成的多智能体系统分布式网络, 其抽象通信拓扑如图 3 所示. 欺骗攻击将会对节点 v_2 与节点 v_4 间的外部通信链路进行攻击. 图中每个节点分别指定初始状态值 $x_1[0] = 2$, $x_2[0] = 4$, $x_3[0] = 6$, $x_4[0] = 8$, $x_5[0] = 10$, $x_6[0] = 12$.

网络中每个节点在 $k = 0$ 时刻开始执行本文提出的算法 1, 这里不妨假设控制增益 $\varepsilon_i = 0$. 欺骗攻击将会在 $k > 2$ 时刻对被攻击的通信链路上传输的数据进行篡改. 具体来说, 当 k 为偶数时, 将链路上数据篡改为 $\tilde{x}_2[k] = 12$; 当 k 为奇数时, 将链路上数据篡改为 $\tilde{x}_2[k] = 2$. 此时, 网络中除节点 v_4 以外均满足广义 1-局部攻击模型.

首先, 考虑节点 v_4 和节点 v_5 之间的通信链路 (见图 3 中虚线) 被移除. 此时网络拓扑图鲁棒性不满足强 3-链路鲁棒图. 本文的模拟仿真实验数据采集的是节点 v_4 未被攻击通信链路上的各个节点状态量测值 $y_i[k]$. 网络中各个节点的状态量测值变化轨迹如图 4 所示. 图 4 中实线表示各个节点状态量测值

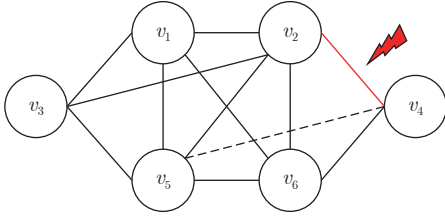
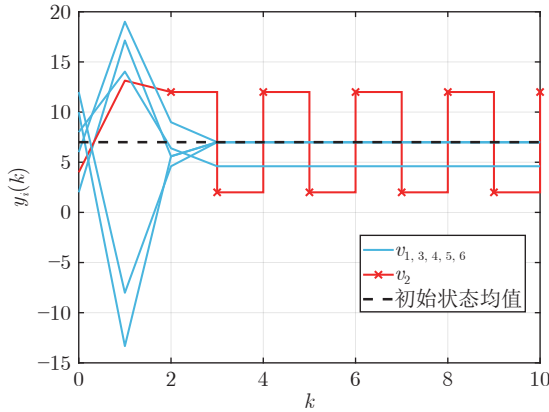


图 3 6 个节点组成的多智能体系统通信拓扑图

Fig.3 Network topology of multi-agent system with 6 nodes

图 4 系统不满足强 $(2f+1)$ -链路鲁棒图下各节点的状态量测值变化轨迹Fig.4 State trajectory of each node with system that does not meet the strong $(2f+1)$ -links robustness

变化, 虚线表示网络中节点的初始状态均值. 可以看出, 在针对节点 v_2 与节点 v_4 间的外部通信链路的欺骗攻击影响下, 当系统网络拓扑图不满足一定鲁棒性条件时, 系统整体无法达成均值趋同.

接着, 考虑节点 v_4 和节点 v_5 之间的通信链路未被移除的情况, 此时该网络拓扑图是一个强 3-链路鲁棒图. 同时, 考虑存在第三方窃听者试图通过收集到的数据推测节点 v_1 的初始状态值 $x_1[0]$. 窃听者对节点 v_1 的初始状态值计算推测公式为:

$$z[1] = z[0] + x_1^\alpha[0] + \sum_{v_j \in \mathcal{N}_1} a_{1j}(x_j^\alpha[0] - x_1^\alpha[0]) \quad (17)$$

式中, $z[1]$ 表示窃听者在 $k=1$ 时刻计算出的推测值. 不妨设窃听者初始赋值为 $z[0]=0$, 并假设窃听者除了节点 v_1 分解后的内部权重 $a_{1,\alpha\beta}$ 未知, 拥有整个分布式网络其他的数据信息. 系统满足强 3-链路鲁棒图条件下各节点的状态量测值变化轨迹如图 5 所示. 其中方形点线表示窃听者对节点 v_1 在 $k=1$ 时刻状态的推测值, 圆形点线表示节点 v_1 希望得到保护的隐私状态值. 从图 5 可以看出, 系统在欺骗攻击下仍然可以精确收敛到节点初始状态值的平均值 7, 且窃听者无法准确估计节点 v_1 在 $k=1$

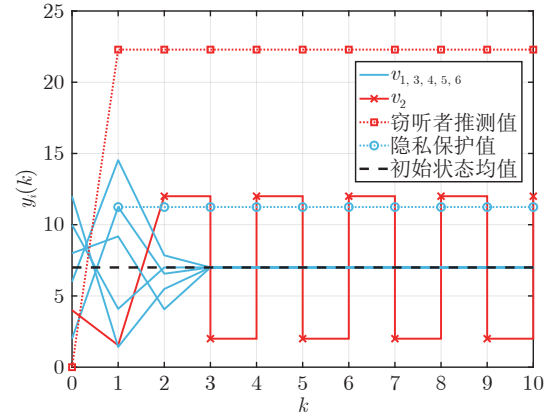
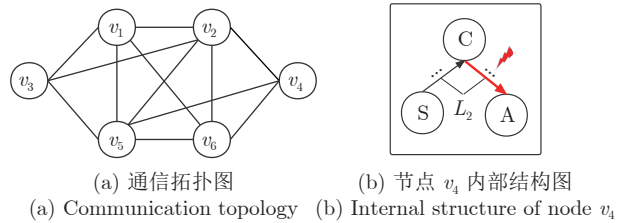


图 5 系统外部通信链路遭受欺骗攻击下各节点的状态量测值变化轨迹

Fig.5 State trajectory of each node under deception attack on the external communication link of the system

时刻的状态值.

同样考虑一个由 6 个节点组成的多智能体系统分布式网络, 其通信拓扑及针对节点 v_4 的攻击示意图如图 6 所示. 图 6(b) 中 C、A、S 分别表示控制器、执行器、传感器, L_2 为节点 v_4 中第二对内部通信链路. 各节点初始状态值设置与图 3 相同. 欺骗攻击将在 $k > 2$ 时刻对节点 v_4 中控制器与执行器的内部通信链路进行攻击. 具体来说, 当 k 为偶数时将链路上数据篡改为 $\tilde{x}_4[k] = 12$, 当 k 为奇数时将链路上数据篡改为 $\tilde{x}_4[k] = 2$. 此时, 多智能体系统满足广义 1-局部攻击模型.

图 6 节点 v_4 内部遭受欺骗攻击的通信拓扑及攻击示意图
Fig.6 Communication topology and attack diagram of the deception attack inside node v_4

首先, 考虑网络中各节点使用文献 [17] 提出的状态分解算法进行实验. 网络中各个节点的量测值变化轨迹如图 7 所示. 可以看出, 在针对节点 v_4 内部通信链路的欺骗攻击影响下, 当系统网络拓扑图鲁棒性符合条件时, 系统不能够达成均值趋同. 网络中节点虽然通过状态分解方法实现了隐私保护, 但多智能体系统最基本的趋同目标没有实现.

最后, 考虑网络中各节点使用本文提出的算法进行实验. 网络中各个节点的量测值变化轨迹如图 8

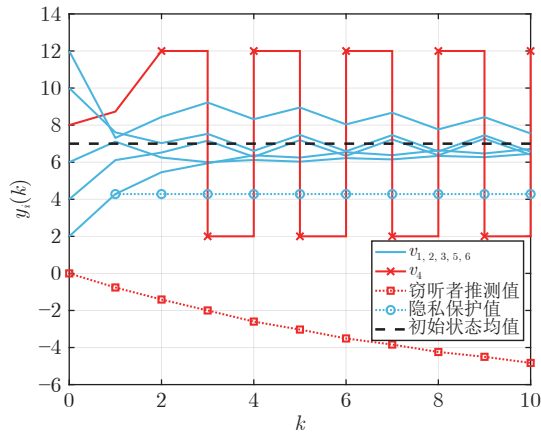


图 7 系统内部通信链路遭受欺骗攻击下使用状态分解算法各节点的状态量测值变化轨迹

Fig.7 State trajectory of each node under deception attack on the internal communication link of the system by using the state decomposition algorithm

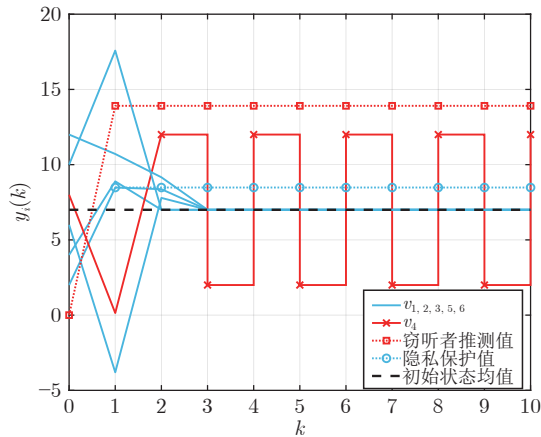


图 8 系统内部通信链路遭受欺骗攻击下使用本文算法各节点的状态量测值变化轨迹

Fig.8 State trajectory of each node under deception attack on the internal communication link of the system by using the proposed algorithm

所示。可以看出,在针对节点 v_4 内部通信链路的欺骗攻击影响下,当系统网络拓扑图鲁棒性符合条件时,系统不仅能够达成均值趋同的目标而且节点的初始状态值隐私得到了保护。

5 结束语

本文针对欺骗攻击下无向多智能体分布式网络均值趋同问题,提出了一种具备隐私保护能力的趋同控制算法,实现了欺骗攻击下多智能体系统分布式网络的均值趋同控制。首先,本文对传统的攻击模型进行了拓展提出了广义 f -局部攻击模型。其

次,使用一种改进后的状态分解方法对分布式网络中节点的初始状态值信息进行处理,经过处理后的系统中所有节点初始状态值得到了隐私保护。然后,利用已具备隐私保护特性的节点状态值加以使用安全接受广播算法实现欺骗攻击下的均值趋同,综合上述方法构造出一种适用于广义 f -局部欺骗攻击下无向通信拓扑的多智能体系统均值趋同控制算法。最后,通过 4 组数值仿真实验验证了所提算法的有效性。

然而,本文提出的控制方法仍存在着不足:1) 所提方法目前仅适用于无向网络,在实际应用中,有向网络则更为普遍,因此接下来的研究将针对有向拓扑图网络对具备隐私保护的均值趋同控制器进行设计;2) 为了抵消欺骗攻击的影响,本文要求系统网络拓扑图符合强 $(2f+1)$ -链路鲁棒图条件,这在实际应用中实现具有一定难度,因此如何在放宽通信网络拓扑图鲁棒性要求下仍能有效抵御欺骗攻击的控制器设计将是接下来另一个值得研究的方向。这些不足,需要在未来的工作中进一步研究。

References

- Ding Li-Fu, Yan Gang-Feng. A survey of the security issues and defense mechanisms of multi-agent systems. *CAAI Transactions on Intelligent Systems*, 2020, **15**(3): 425-434 (丁俐夫, 颜钢锋. 多智能体系统安全性问题及防御机制综述. 智能系统学报, 2020, **15**(3): 425-434)
- Li Tao, Meng Yang, Zhang Ji-Feng. An overview on quantized consensus and consensus with limited data rate of multi-agent systems. *Acta Automatica Sinica*, 2013, **39**(11): 1805-1811 (李韬, 孟杨, 张纪锋. 多自主体量化趋同与有限数据率趋同综述. 自动化学报, 2013, **39**(11): 1805-1811)
- Wang Xiang-Ke, Li Xun, Zheng Zhi-Qiang. Survey of developments on multi-agent formation control related problems. *Control and Decision*, 2013, **28**(11): 1601-1613 (王祥科, 李迅, 郑志强. 多智能体系统编队控制相关问题研究综述. 控制与决策, 2013, **28**(11): 1601-1613)
- Sun Qiu-Ye, Teng Fei, Zhang Hua-Guang. Energy internet and its key control issues. *Acta Automatica Sinica*, 2017, **43**(2): 176-194 (孙秋野, 滕菲, 张化光. 能源互联网及其关键控制问题. 自动化学报, 2017, **43**(2): 176-194)
- Huang Z, Mitra S, Dullerud G. Differentially private iterative synchronous consensus. In: *Proceedings of the ACM Workshop on Privacy in the Electronic Society*. New York, USA: 2012. 81-90
- Huang Z, Mitra S, Vaidya N. Differentially private distributed optimization. In: *Proceedings of the International Conference on Distributed Computing and Networking*. New York, USA: 2015. 1-10
- Nozari E, Tallapragada P, Cortes J. Differentially private average consensus: Obstructions, trade-offs, and optimal algorithm design. *Automatica*, 2017, **81**(7): 221-231
- Katewa V, Pasqualetti F, Gupta V. On privacy vs cooperation in multi-agent systems. *International Journal of Control*, 2018, **91**(7): 1693-1707
- Manitara N, Hadjicostis C. Privacy-preserving asymptotic average consensus. In: *Proceedings of the European Control Conference*. Zurich, Switzerland: IEEE, 2013. 760-765
- Kia S, Cortes J, Martinez S. Dynamic average consensus under

- limited control authority and privacy requirements. *International Journal of Robust and Nonlinear Control*, 2015, **25**(13): 1941–1966
- 11 Pequito S, Kar S, Sundaram S, Aguiar A P. Design of communication networks for distributed computation with privacy guarantees. In: Proceedings of the 54th IEEE Conference on Decision and Control. Osaka, Japan: IEEE, 2015. 1370–1376
 - 12 Alaeddini A, Morgansen K, Mesbahi M. Adaptive communication networks with privacy guarantees. In: Proceedings of the American Control Conference. Seattle, USA: IEEE, 2017. 4460–4465
 - 13 Hendriks R C, Erkin Z, Gerkmann T. Privacy preserving distributed beamforming based on homomorphic encryption. In: Proceedings of the 21st European Signal Processing Conference. Marrakech, Morocco: IEEE, 2013. 1–5
 - 14 Hendriks R C, Erkin Z, Gerkmann T. Privacy-preserving distributed speech enhancement for wireless sensor networks by processing in the encrypted domain. In: Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing. Vancouver, Canada: IEEE, 2013. 7005–7009
 - 15 Li Q, Cascudo I, Christensen M G. Privacy-preserving distributed average consensus based on additive secret sharing. In: Proceedings of the 27th European Signal Processing Conference. Coruna, Spain: IEEE, 2019. 1–5
 - 16 Li Q, Christensen M G. A privacy-preserving asynchronous averaging algorithm based on shamir's secret sharing. In: Proceedings of the 27th European Signal Processing Conference. Coruna, Spain: IEEE, 2019. 1–5
 - 17 Wang Y. Privacy-preserving average consensus via state decomposition. *IEEE Transactions on Automatic Control*, 2019, **64**(11): 4711–4716
 - 18 Zhang D, Liu L, Feng G. Consensus of heterogeneous linear multiagent systems subject to aperiodic sampled-data and DOS attack. *IEEE Transactions on Cybernetics*, 2019, **49**(4): 1501–1511
 - 19 Feng Z, Hu G. Secure cooperative event-triggered control of linear multiagent systems under DoS attacks. *IEEE Transactions on Control Systems Technology*, 2020, **28**(3): 741–752
 - 20 Yang Y, Xu H, Yue D. Observer-based distributed secure consensus control of a class of linear multi-agent systems subject to random attacks. *IEEE Transactions on Circuits and Systems I: Regular Papers*, 2019, **66**(8): 3089–3099
 - 21 Xu W, Hu G, Ho D W C, Feng Z. Distributed secure cooperative control under denial-of-service attacks from multiple adversaries. *IEEE Transactions on Cybernetics*, 2020, **50**(8): 3458–3467
 - 22 Zhu M, Martínez S. On distributed constrained formation control in operator-vehicle adversarial networks. *Automatica*, 2013, **49**(12): 3571–3582
 - 23 Ding D, Wang Z, Ho D W C, Wei G. Observer-based event-triggering consensus control for multi-agent systems with lossy sensors and cyber-attacks. *IEEE Transactions on Cybernetics*, 2016, **47**(8): 1936–1947
 - 24 He W, Gao X, Zhong W, Qian F. Secure impulsive synchronization control of multi-agent systems under deception attacks. *Information Sciences*, 2018, **459**: 354–368
 - 25 Fu W, Qin J, Shi Y, Zheng W X, Kang Y. Resilient consensus of discrete-time complex cyber-physical networks under deception attacks. *IEEE Transactions on Industrial Informatics*, 2019, **16**(7): 4868–4877
 - 26 He W, Mo Z, Han Q L, Qian F. Secure impulsive synchronization in Lipschitz-type multi-agent systems subject to deception attacks. *IEEE/CAA Journal of Automatica Sinica*, 2020, **7**(5): 1326–1334
 - 27 Li H, Liao X, Huang T, Zhu W, Liu Y. Second-order global consensus in multiagent networks with random directional link failure. *IEEE Transactions on Neural Networks and Learning Systems*, 2014, **26**(3): 565–575
 - 28 Li H, Chen G, Huang T, Dong Z. High-performance consensus control in networked systems with limited bandwidth communication and time-varying directed topologies. *IEEE Transactions on Neural Networks and Learning Systems*, 2016, **28**(5): 1043–1054
 - 29 Lu Q, Liao X, Xiang T, Li H, Huang T. Privacy masking stochastic subgradient-push algorithm for distributed online optimization. *IEEE Transactions on Cybernetics*, 2020, **51**(6): 3224–3237
 - 30 Fiore D, Russo G. Resilient consensus for multi-agent systems subject to differential privacy requirements. *Automatica*, 2019, **106**: 18–26
 - 31 Zhang H, Sundaram S. Robustness of information diffusion algorithms to locally bounded adversaries. In: Proceedings of the American Control Conference. Montréal, Canada: IEEE, 2012. 5855–5861
 - 32 Dibaji S M, Safi M, Ishii H. Resilient distributed averaging. In: Proceedings of the American Control Conference. Philadelphia, USA: 2019. 96–101
 - 33 Hale M T, Egerstedty M. Differentially private cloud-based multi-agent optimization with constraints. In: Proceedings of the American Control Conference. Chicago, USA: 2015. 1235–1240
 - 34 LeBlanc H J, Zhang H, Koutsoukos X, Sundaram S. Resilient asymptotic consensus in robust networks. *IEEE Journal on Selected Areas in Communications*, 2013, **31**(4): 766–781



应晨铎 杭州电子科技大学网络空间安全学院硕士研究生。2020 年获得浙大宁波理工学院软件工程学士学位。主要研究方向为弹性趋同, 隐私保护和分布式系统安全。

E-mail: cdying@hdu.edu.cn

(YING Chen-Duo Master student

at the School of Cyberspace, Hangzhou Dianzi University. He received his bachelor degree in software engineering from NingboTech University in 2020. His research interest covers resilient consensus, privacy preservation, and distributed system security.)



伍益明 杭州电子科技大学网络空间安全学院副教授。2016 年获得浙江工业大学控制科学与工程博士学位。主要研究方向为分布式系统安全控制, 多智能体系统网络安全和迭代学习控制。本文通信作者。

E-mail: ymwu@hdu.edu.cn

(WU Yi-Ming Associate professor at the School of Cyberspace, Hangzhou Dianzi University. He received his Ph.D. degree in control science and engineering from Zhejiang University of Technology in 2016. His research interest covers distributed system secure control, cyber-security for multi-agent systems, and iterative learning control. Corresponding author of this paper.)



徐 明 杭州电子科技大学网络空间安全学院教授. 2004 年获得浙江大学博士学位. 主要研究方向为网络信息安全, 数字取证.

E-mail: mxu@hdu.edu.cn

(XU Ming Professor at the School of Cyberspace, Hangzhou Dianzi

University. He received his Ph.D. degree from Zhejiang University in 2004. His research interest covers network security and digital forensics.)



郑 宁 杭州电子科技大学网络空间安全学院研究员. 1987 年获得浙江大学硕士学位. 主要研究方向为信息安全, 信息管理系统和多智能体系统.

E-mail: nzheng@hdu.edu.cn

(ZHENG Ning Professor at the School of Cyberspace, Hangzhou Di-

anzi University. He received his master degree from Zhejiang University in 1987. His research interest covers information security, information management system, and multi-agent systems.)



何熊熊 浙江工业大学信息工程学院教授. 1997 年获得浙江大学博士学位. 主要研究方向为迭代学习控制, 智能控制及其在多智能体系统和传感器网络中的应用.

E-mail: hxx@zjut.edu.cn

(HE Xiong-Xiong Professor at the

College of Information Engineering, Zhejiang University of Technology. He received his Ph.D. degree from Zhejiang University in 1997. His research interest covers iterative learning control, intelligent control and its applications in multi-agent systems and sensor networks.)