多无人机系统在线强化学习最优 安全跟踪控制



弓镇宇,杨飞生

西北工业大学,陕西西安 710072

摘 要:在无人机(UAV)编队跟踪任务中,虚假数据注入(FDI)攻击者可向控制指令注入误导性数据,导致无人机无法形成 指定的编队构型,故需设计安全编队跟踪控制器。为此,本文利用零和图博弈对攻防过程进行建模,其中FDI攻击者和安全 控制器是博弈的参与者,攻击者的目标是最大化设定的成本函数,而安全控制器的目标与之相反,求解博弈并获得最优安全 控制策略依赖于求取Hamilton-Jacobi-Isaacs(HJI)方程的解。而HJI方程是耦合偏微分方程,难以直接求解,因此结合经验 回放机制引入了有限时间收敛的在线强化学习算法,设计了单评价神经网络近似值函数并获得了最优安全控制策略。最终 利用仿真验证了算法的有效性。

关键词:FDI攻击;多无人机;在线强化学习;优化控制;零和图博弈

中图分类号:V249.1

文献标识码:A

无人机作为典型的无人系统,已经广泛应用于农林作 业、电力巡检、灾后搜救、目标侦察、协同作战等领域^{III}。 相比于单体无人机的机载设备有限、感知范围小、任务执 行容错能力差,整合了信息融合、目标分配、协同控制等技 术的多无人机系统可在复杂周边环境下完成多样化任务。 然而,无人机之间的信息交互依赖于通信网络,因此多无 人机系统面临网络攻击的威胁。攻击者可向无人机传输 信道中注入欺骗型数据以降低系统性能,甚至导致任务失 败,因此设计安全控制方案以抵御虚假数据注入(FDI)攻 击至关重要。

目前,主要有两种安全控制方案以应对FDI攻击,其区 别在于是否引入了攻击检测机制^[2-3]。在无人机系统中, Lin Hong等^[4]设计了攻击检测器和线性二次高斯控制器来 抵抗FDI攻击。Xiao Jiaping等^[5]基于滑动新息序列提出了 一种新型的攻击检测器。为了节约网络通信资源,Yin Tingting等^[6]研究了事件触发机制下的无人机安全编队 控制。

随着人工智能技术的发展,强化学习算法因其良好的 决策优化和实时的策略选择能力而备受关注,越来越多研

DOI: 10.19452/j.issn1007-5453.2024.04.004

究者将其应用于控制问题求解^[7-10]。针对安全控制问题, Wu Chengwei 等^[11]利用Q学习算法研究了控制信号遭受 FDI攻击时的安全控制器设计。Zhou Yuanqiang等^[12]结合 威胁检测水平函数设计了检测机制并利用离策略算法求解 了最优安全控制器。多智能体系统中,Moghadam等^[13]结 合离策略算法求解了非齐次博弈 Riccati 方程并设计了弹性 控制器。考虑事件触发机制,Xu Yuanyuan等^[14]结合预选器 和观测器研究了传感器遭受 FDI 攻击时的最优控制律 设计。

已有文献很少关注强化学习的快速收敛问题, Kokolakis等^[15]设计了有限时间收敛的强化学习算法,但是 其未考虑多无人机情形。本文在领导一跟随多智能体框架 下研究了多无人机遭受网络攻击时的安全编队跟踪控制问 题。考虑到FDI攻击者和安全控制器之间相互对抗的关 系,引入了零和图博弈理论对攻防过程进行建模,最优攻击 和最优安全控制器位于纳什均衡点。为求解博弈并进一步 获取最优安全控制律,引入了有限时间收敛强化学习方法。 通过单评价神经网络架构对值函数进行逼近,采用了经验 回放机制以维持持续激励条件,分析了算法的收敛性以及

收稿日期: 2023-06-27; 退修日期: 2023-11-16; 录用日期: 2023-12-25

基金项目:国家自然科学基金(62073269);航空科学基金(2020Z034053002);陕西省重点研发计划项目(2022GY-244);重庆市自然科学基金 (CSTB2022NSCQ-MSX0963);广东省基础与应用基础研究基金(2023A1515011220)

引用格式: Gong Zhenyu, Yang Feisheng. Optimal secure tracking control in multi-UAVs based on online reinforcement learning [J]. Aeronautical Science & Technology, 2024, 35(04):25-30. 弓镇宇, 杨飞生. 多无人机系统在线强化学习最优安全跟踪控制[J]. 航空 科学技术, 2024, 35(04):25-30. 停息时间上界。

1 代数图论

在通过网络进行信息交互的多无人机系统中,每架无 人机都可看作通信拓扑的节点。通信拓扑图可表示为 \mathcal{G} = ($\mathcal{V}, \mathcal{E}, A$),其中 $\mathcal{V} = \{v_1, v_2, \dots, v_N\}$ 是图的节点集, $\mathcal{E} = \{(v_i, v_j)|v_i, v_j \in \mathcal{V}\}$ 表示图的边集。 $A = \{a_{ij}\}$ 为邻接矩阵, 如果 $\{v_i, v_j \in \mathcal{E}, \prod a_{ij} > 0,$ 否则 $a_{ij} = 0$ 。无人机i的邻居集为 $\mathcal{N}_i = \{j|(v_i, v_j) \in \mathcal{E}\}$ 。定义图的入度矩阵为 $D = \text{diag}\{d_i\},$ 其中 $d_i = \sum_{j \in \mathcal{N}_i} a_{ij^\circ}$ 图的拉普拉斯矩阵可表示为L = D - A。考 虑存在领导者的图 $\overline{\mathcal{G}},$ 则牵引矩阵为 $G = \text{diag}\{g_i\}$ 。当第i架无人机可直接接收领导者信息时, $g_i = 1,$ 否则 $g_i = 0$ 。若 图存在一个根节点,可通过有向路径到达图中其他任意节 点,则该图存在一个有向生成树。

2 问题描述

考虑包含N架无人机的多无人机系统。第*i*架无人机 的动力学方程如下所示

$$\dot{\boldsymbol{p}}_{i}(t) = \boldsymbol{v}_{i}(t)$$

$$\dot{\boldsymbol{v}}_{i}(t) = -\boldsymbol{c}\boldsymbol{g} + \boldsymbol{\bar{\mu}}_{i}(t)$$
(1)

式中, $p_i(t) \in \mathbb{R}^3$ 、 $v_i(t) \in \mathbb{R}^3$ 、 $\bar{u}_i(t) \in \mathbb{R}^3$ 分别代表第*i*架无人 机的位置矢量、速度矢量、控制输入矢量, $c = [0 \ 0 \ 1]^T$, g_e 是 重 力 加 速 度 。 令 系 统 的 状 态 矢 量 为 $x_i(t) = [p_i^T(t) \ v_i^T(t)]^T$,则有

$$\dot{\boldsymbol{x}}_i(t) = \boldsymbol{A}\boldsymbol{x}_i(t) + \boldsymbol{B}\overline{\boldsymbol{u}}_i(t) + \boldsymbol{C}\boldsymbol{g}_e$$
(2)

式中

$$A = \begin{bmatrix} \boldsymbol{0}_{3\times3} & \boldsymbol{I}_3 \\ \boldsymbol{0}_{3\times3} & \boldsymbol{0}_{3\times3} \end{bmatrix}, B = \begin{bmatrix} \boldsymbol{0}_{3\times3} \\ \boldsymbol{I}_3 \end{bmatrix}, C = \begin{bmatrix} \boldsymbol{0}_{1\times5} & -1 \end{bmatrix}^T$$
考虑攻击者会向无人机的控制信号中注入虚假数据,

因此有

$$\bar{\boldsymbol{u}}_i(t) = \check{\boldsymbol{u}}_i(t) + \boldsymbol{w}_i(t)$$
(3)

式中,矢量 $\check{u}_i(t)$ 是安全编队跟踪控制器,矢量 $w_i(t)$ 是攻击 信号。给定控制器为

$$\check{\boldsymbol{u}}_i(t) = \boldsymbol{u}_i(t) + \boldsymbol{c}\boldsymbol{g}_e \tag{4}$$

式中,矢量 $u_i(t)$ 是待设计控制器。将式(3)和式(4)代入式 (2)可得

$$\dot{\boldsymbol{x}}_{i}(t) = \boldsymbol{A}\boldsymbol{x}_{i}(t) + \boldsymbol{B}\boldsymbol{u}_{i}(t) + \boldsymbol{B}\boldsymbol{w}_{i}(t)$$
(5)

考虑虚拟领导者可提供期望位置矢量 $p_0(t)$ 与期望速 度 矢 量 $v_0(t)$,则 虚 拟 领 导 者 的 状 态 矢 量 为 $x_0(t) = [p_0^{\mathsf{T}}(t) v_0^{\mathsf{T}}(t)]^{\mathsf{T}}$ 。令第*i*架无人机的位置跟踪误差矢量和速 度跟踪误差矢量为

$$\begin{aligned} \boldsymbol{\delta}_{\mathrm{p}i}(t) &= \boldsymbol{p}_i(t) - \boldsymbol{p}_0(t) - \boldsymbol{\chi}_i \\ \boldsymbol{\delta}_{\mathrm{v}i}(t) &= \boldsymbol{v}_i(t) - \boldsymbol{v}_0(t) \end{aligned} \tag{6}$$

式中, $\chi_i \in \mathbb{R}^3$ 是第*i*架无人机与领导者之间的相对位置矢量。进而引入以下假设。

假设1:无人机系统的通信拓扑图存在有向生成树,且 树的根节点为领导者。

每架无人机的目标是跟踪上期望轨迹与速度,即满足 lim $\|\delta_{vi}(t)\|_2 = 0$, lim $\|\delta_{vi}(t)\|_2 = 0$, lim $\|\delta_{vi}(t)\|_2 = 0$, i = 1, 2, …, N_o

定义系统的编队误差矢量为

接下来引入博弈思想分析攻防过程。

3 安全零和图博弈

针对第 *i* 架尤人机设计如下成本函数。

$$J_i(\boldsymbol{e}_i(0), \boldsymbol{u}_i(t), \boldsymbol{u}_{-i}(t), \boldsymbol{w}_i(t), \boldsymbol{w}_{-i}(t)) = \int_0^{\infty} U_i(\boldsymbol{e}_i(t), \boldsymbol{u}_i(t), \boldsymbol{u}_{-i}(t), \boldsymbol{w}_i(t), \boldsymbol{w}_{-i}(t)) dt = \int_0^{\infty} (\boldsymbol{e}_i^{\mathrm{T}}(t) \boldsymbol{Q}_{ii} \boldsymbol{e}_i(t) + \boldsymbol{u}_i^{\mathrm{T}}(t) \boldsymbol{R}_{ii} \boldsymbol{u}_i(t) + \sum_{j \in \mathcal{N}_i} \boldsymbol{u}_j^{\mathrm{T}}(t) \boldsymbol{R}_{ij} \boldsymbol{u}_j(t) - \gamma^2 \boldsymbol{w}_i^{\mathrm{T}}(t) \boldsymbol{T}_{ii} \boldsymbol{w}_i(t) - \gamma^2 \sum_{i \in \mathcal{N}_i} \boldsymbol{w}_j^{\mathrm{T}}(t) \boldsymbol{T}_{ij} \boldsymbol{w}_j(t)) dt$$
(9)

式中, $Q_{ii} > 0, R_{ii} > 0, R_{ij} \ge 0, T_{ii} > 0, T_{ij} \ge 0$ 是对称权重矩 阵, $\gamma > 0$ 为攻击抑制水平,下标-*i*代表第*i*架无人机的邻居 为简单起见,在后文中省去时间符号"*t*"。

从攻击者的角度而言,其目标是使得系统性能下降 并影响跟踪效果,而安全控制器的目标是抵御这种负面 作用。可以看出,FDI攻击者与安全控制器形成一种相 互对抗的关系,可引入零和图博弈来描述该攻防过程。 攻击者和安全控制器可看作博弈参与者,攻击者的目的 是最大化式(9),而控制器的目的则与之相反。此过程可 表示为

$$V_i(\boldsymbol{e}_i(0)) = \min \max J_i(\boldsymbol{e}_i, \boldsymbol{u}_i, \boldsymbol{u}_{-i}, \boldsymbol{w}_i, \boldsymbol{w}_{-i})$$

式中, $V_i(e_i(0))$ 为博弈的值。如果存在鞍点 (u_i^*, w_i^*) ,则有

$$V_{i}(e_{i}(0)) = \min_{u_{i}} \max_{w_{i}} J_{i}(e_{i}, u_{i}, u_{-i}^{*}, w_{i}, w_{-i}) = \max_{w_{i}} \min_{u_{i}} J_{i}(e_{i}, u_{i}, u_{-i}^{*}, w_{i}, w_{-i}^{*})$$
其等价于纳什均衡条件
L(u^{*} u^{*} u^{*} u^{*} u^{*}) < L(u^{*} u^{*} u^{*} u^{*} u^{*}) < L(u u^{*} u^{*} u^{*} u^{*}) < L(u u^{*} u^{*} u^{*} u^{*}) < L(u^{*} u^{*} u^{*} u^{*}) < L(u^{*} u^{*} u^{*} u^{*}) < L(u^{*} u^{*} u^{*} u^{*}) < L(u^{*} u^{*} u^{*} u^{*} u^{*} u^{*}) < L(u^{*} u^{*} u^

$$J_{i}(u_{i}, u_{-i}, w_{i}, w_{-i}) \leq J_{i}(u_{i}, u_{-i}, w_{i}, w_{-i}) \leq J_{i}(u_{i}, u_{-i}, w_{i}, w_{-i})$$

(10)

式中,∇为微分算子。根据一阶平稳性条件可得

$$\boldsymbol{u}_{i}^{*} = -\frac{1}{2} \left(\boldsymbol{d}_{i} + \boldsymbol{g}_{i} \right) \boldsymbol{R}_{ii}^{-1} \boldsymbol{B}^{\mathrm{T}} \nabla V_{i}(\boldsymbol{e}_{i})$$
(12)

$$\boldsymbol{w}_{i}^{*} = \frac{1}{2\gamma^{2}} \left(d_{i} + g_{i} \right) \boldsymbol{T}_{ii}^{-1} \boldsymbol{B}^{\mathrm{T}} \nabla V_{i}(\boldsymbol{e}_{i})$$
(13)

将式(12)和式(13)代入式(11)中,可得如下Hamilton-Jacobi-Isaacs(HJI)方程。

$$H_{i}(\boldsymbol{e}_{i},\boldsymbol{u}_{i}^{*},\boldsymbol{u}_{-i}^{*},\boldsymbol{w}_{i}^{*},\boldsymbol{w}_{-i}^{*},\nabla V_{i}(\boldsymbol{e}_{i})) = 0$$
(14)

由于HJI方程是耦合的,因此引入在线强化学习方法 对其进行求解。

4 有限时间收敛在线强化学习

根据 Weierstrass 高阶逼近定理,给定以下评价神经网络对值函数进行逼近

$$V_i(\boldsymbol{e}_i) = \boldsymbol{W}_i^{\mathrm{T}} \boldsymbol{\phi}_i(\boldsymbol{e}_i) + \boldsymbol{\varepsilon}_i(\boldsymbol{e}_i), \ \boldsymbol{e}_i \in \mathbb{E}$$
(15)

式中, W_i , $\phi_i(\cdot)$, $\varepsilon_i(\cdot)$ 和 E分别为神经网络的权重矩阵、激活函数、逼近误差和编队误差集。由于 W_i 是理想权重, 难以直接获得, 所以神经网络的输出可近似为

 $\hat{V}_i(\boldsymbol{e}_i) = \boldsymbol{\hat{W}}_i^{\mathrm{T}} \boldsymbol{\phi}_i(\boldsymbol{e}_i), \ \boldsymbol{e}_i \in \mathbb{E}$ (16)

式中, ŵ, 为估计权重矩阵。因此, 最优策略可近似为

$$\hat{\boldsymbol{u}}_{i} = -\frac{1}{2} \left(d_{i} + g_{i} \right) \boldsymbol{R}_{ii}^{-1} \boldsymbol{B}^{\mathrm{T}} \nabla \boldsymbol{\phi}_{i}^{\mathrm{T}} \left(\boldsymbol{e}_{i} \right) \hat{\boldsymbol{W}}_{i}$$
(17)

$$\hat{\boldsymbol{w}}_{i} = \frac{1}{2\gamma^{2}} \left(d_{i} + g_{i} \right) \boldsymbol{T}_{ii}^{-1} \boldsymbol{B}^{\mathrm{T}} \nabla \boldsymbol{\phi}_{i}^{\mathrm{T}} (\boldsymbol{e}_{i}) \hat{\boldsymbol{W}}_{i}$$
(18)

定义如卜残差函数

$$\xi_{i}(t) = U_{i}(\boldsymbol{e}_{i}(t), \hat{\boldsymbol{u}}_{i}(t), \hat{\boldsymbol{u}}_{-i}(t), \hat{\boldsymbol{w}}_{i}(t), \hat{\boldsymbol{w}}_{-i}(t)) +$$

$$\hat{\boldsymbol{W}}_{i}^{\mathrm{T}}(t) \nabla \boldsymbol{\phi}_{i}(\boldsymbol{e}_{i}(t)) \dot{\boldsymbol{e}}_{i}(t)$$
(19)

在训练神经网络权重的过程中,为了维持持续激励条件,通常需引入探测噪声,但这种做法会影响系统的稳定

性。而经验回放机制弥补了该缺陷,利用了历史数据,因此 定义损失函数为

$$E_{i}(\hat{\boldsymbol{W}}_{i}) = \frac{1}{q+1} \left| \frac{\xi_{i}(t)}{1+\boldsymbol{\psi}_{i}^{\mathrm{T}}(t)\boldsymbol{\psi}_{i}(t)} \right|^{q+1} + \frac{1}{q+1} \sum_{s=t_{0}}^{t_{i}} \left| \frac{\xi_{i}^{[t]}(s)}{1+\boldsymbol{\psi}_{i}^{\mathrm{T}}(s)\boldsymbol{\psi}_{i}(s)} \right|^{q+1}$$
(20)

式中,0 < q < 1, $\psi_i(t) = \nabla \phi_i(e_i(t))\dot{e}_i(t)$,|·|是绝对值函数, 且时间序列0 < t_0 ,..., $t_i < t$ 上的历史残差数据为

$$\xi_{i}^{[i]}(s) = U_{i}(\boldsymbol{e}_{i}(s), \hat{\boldsymbol{u}}_{i}(s), \hat{\boldsymbol{u}}_{-i}(s), \hat{\boldsymbol{w}}_{i}(s), \hat{\boldsymbol{w}}_{-i}(s)) +$$

$$\hat{\boldsymbol{W}}_{i}^{\mathrm{T}}(t)\boldsymbol{\phi}_{i}(\boldsymbol{e}_{i}(s))\dot{\boldsymbol{e}}_{i}(s)$$

结合梯度下降原理,可得出神经网络权重的更新律为

$$\dot{\hat{\boldsymbol{W}}}_{i} = -\alpha_{i} \bar{\boldsymbol{\psi}}_{i}(t) \left\| \frac{\boldsymbol{\xi}_{i}(t)}{1 + \boldsymbol{\psi}_{i}^{\mathrm{T}}(t) \boldsymbol{\psi}_{i}(t)} \right\|^{q} - \alpha_{i} \sum_{s=t_{0}}^{t_{i}} \bar{\boldsymbol{\psi}}_{i}(s) \left\| \frac{\boldsymbol{\xi}_{i}^{[t]}(s)}{1 + \boldsymbol{\psi}_{i}^{\mathrm{T}}(s) \boldsymbol{\psi}_{i}(s)} \right\|^{q}$$

$$(21)$$

式中, $\bar{\boldsymbol{\psi}}_i(t) = \boldsymbol{\psi}_i(t)/(1 + \boldsymbol{\psi}_i^{\mathrm{T}}(t)\boldsymbol{\psi}_i(t)), \alpha_i > 0$ 是神经网络的学习 率, [[·]] = sgn(·)[·], sgn(·)是符号函数。进一步有以下假设。

假设2:矩阵 $\Psi_i = [\psi_i(t_0), \dots, \psi_i(t_f)]$ 由 k组历史数据组成且是行满秩的。

定义权重误差为 $\tilde{W}_i = \hat{W}_i - W_i$,则权重误差的动力学方程为

$$\begin{split} \hat{\boldsymbol{W}}_{i} &= -\alpha_{i} \bar{\boldsymbol{\psi}}_{i}(t) \left[\left[\frac{\boldsymbol{\psi}_{i}^{\mathrm{T}}(t) \tilde{\boldsymbol{W}}_{i}(t) + \boldsymbol{\epsilon}_{i}(t)}{1 + \boldsymbol{\psi}_{i}^{\mathrm{T}}(t) \boldsymbol{\psi}_{i}(t)} \right]^{q} - \\ &\alpha_{i} \sum_{s=t_{0}}^{t_{i}} \bar{\boldsymbol{\psi}}_{i}(s) \left[\left[\frac{\boldsymbol{\psi}_{i}^{\mathrm{T}}(s) \tilde{\boldsymbol{W}}_{i}(t) + \boldsymbol{\epsilon}_{i}(s)}{1 + \boldsymbol{\psi}_{i}^{\mathrm{T}}(s) \boldsymbol{\psi}_{i}(s)} \right]^{q} \end{split}$$

$$(22)$$

其中

$$\begin{aligned} \boldsymbol{\epsilon}_{i}(t) &= U_{i} \Big(\boldsymbol{e}_{i}(t), \hat{\boldsymbol{u}}_{i}(t), \hat{\boldsymbol{u}}_{-i}(t), \hat{\boldsymbol{w}}_{i}(t), \hat{\boldsymbol{w}}_{-i}(t) \Big) + \\ \boldsymbol{W}_{i}^{\mathrm{T}}(t) \nabla \boldsymbol{\phi}_{i}(\boldsymbol{e}_{i}(t)) \dot{\boldsymbol{e}}_{i}(t) \end{aligned}$$

 $\epsilon_i(t)$ 的有界性与神经网络隐层的神经元个数相关,进而与强化学习算法收敛性相关,同时有以下引理。

引理1^[15]:对于任意的 $\hat{\epsilon}$,存在 $F(\hat{\epsilon}) > 0$ 和 $M_0(\hat{\epsilon}) > 0$,当 $M \rightarrow \infty$ 时满足 sup $|\epsilon_i| < F(\hat{\epsilon}), M_0(\hat{\epsilon}) \leq M, \epsilon_i \equiv 0_\circ$

接下来分析算法的有限时间收敛性,我们有如下定理。

定理 1: 在满足假设2的前提下,令 $\bar{\Psi}_i = [\bar{\Psi}_i(t_0), \cdots, \bar{\Psi}_i(t_f)]$,给定0 < $\theta_i < 1$, $|\bar{\epsilon}_i| < \bar{\epsilon}_m$, $\|\bar{\Psi}_i\|_2 \leq \bar{\Psi}_m$ 。

(1)当 $\epsilon_i \equiv 0$ 时,式(22)的零解 $\tilde{W}_i = 0$ 是全局有限时间 稳定的,且停息时间为

$$T \leq \frac{\left\|\tilde{\boldsymbol{W}}_{i}(0)\right\|_{2}^{1-q}}{\sigma_{\min}^{q+1}(\bar{\boldsymbol{\Psi}}_{i})\alpha_{i}(1-q)}$$

(2)当 ϵ_i ≠0时,式(22)是最终一致有界的,且收敛界为

$$\begin{split} \mu_{i} &= \left(\frac{(k+1)\bar{\epsilon}_{m}^{q}\bar{\Psi}_{m}}{\theta_{i}\sigma_{\min}^{q+1}(\bar{\Psi}_{i})}\right)^{\frac{1}{q}} \\ \bar{P} &\equiv \mathrm{btild} \\ \bar{P} &\equiv \mathrm{btild} \\ T &\leq \frac{\left\|\tilde{\Psi}_{i}(0)\right\|_{2}^{1-q} - \mu_{i}^{1-q}}{\alpha_{i}\sigma_{\min}^{q+1}(\bar{\Psi}_{i})(1-\theta_{i})(1-q)} \\ \mathrm{itfild} &\equiv \mathrm{Lyapunov} \mathrm{ind} \\ \bar{V}_{i} &= \frac{1}{2\alpha_{i}}\tilde{\Psi}_{i}^{\mathsf{T}}\tilde{\Psi}_{i} \\ \bar{V}_{i} &\equiv \mathrm{cl} \\ \bar{\gamma}_{i} &\equiv \mathrm{cl} \\ \bar{\gamma}_{$$

$$\dot{\tilde{V}}_{i} = -\tilde{\boldsymbol{W}}_{i}^{\mathrm{T}}(t)\bar{\boldsymbol{\psi}}_{i}(t)\left[\!\left[\bar{\boldsymbol{\psi}}_{i}^{\mathrm{T}}(t)\tilde{\boldsymbol{W}}_{i}(t) + \bar{\boldsymbol{\epsilon}}_{i}(t)\right]\!\right]^{q} - \sum_{s=t_{0}}^{t_{i}}\tilde{\boldsymbol{W}}_{i}^{\mathrm{T}}(t)\bar{\boldsymbol{\psi}}_{i}(s)\left[\!\left[\bar{\boldsymbol{\psi}}_{i}^{\mathrm{T}}(s)\tilde{\boldsymbol{W}}_{i}(t) + \bar{\boldsymbol{\epsilon}}_{i}(s)\right]\!\right]^{q}$$
(24)

式中, $\bar{\boldsymbol{\epsilon}}_i(t) = \boldsymbol{\epsilon}_i(t)/(1 + \boldsymbol{\psi}_i^{\mathrm{T}}(t)\boldsymbol{\psi}_i(t))_{\circ}$

$$(1) \stackrel{\text{def}}{=} \epsilon_{i} \equiv 0 \text{ fr}, \exists \stackrel{\text{def}}{=} \tilde{\boldsymbol{V}}_{i}^{\mathrm{T}}(t) \bar{\boldsymbol{\psi}}_{i}(t) \left[\!\left[\boldsymbol{\bar{\psi}}_{i}^{\mathrm{T}}(t) \tilde{\boldsymbol{W}}_{i}(t) \right]\!\right]^{q} - \sum_{s=t_{0}}^{t_{t}} \left[\boldsymbol{\tilde{\psi}}_{i}^{\mathrm{T}}(t) \bar{\boldsymbol{\psi}}_{i}(s) \left[\!\left[\boldsymbol{\bar{\psi}}_{i}^{\mathrm{T}}(s) \tilde{\boldsymbol{W}}_{i}(t) \right]\!\right]^{q} \leq -\sum_{s=t_{0}}^{t_{t}} \left| \boldsymbol{\bar{\psi}}_{i}^{\mathrm{T}}(s) \tilde{\boldsymbol{W}}_{i}(t) \right|^{q+1}$$

$$(25)$$

$$\begin{split} &-\sum_{s=t_0}^{t_i} \left| \bar{\boldsymbol{\psi}}_i^{\mathrm{T}}(s) \tilde{\boldsymbol{W}}_i(t) \right|^{q+1} = - \left\| \bar{\boldsymbol{\Psi}}_i^{\mathrm{T}} \tilde{\boldsymbol{W}}_i \right\|_{q+1}^{q+1} \\ & \text{根据范数的单调性可得} \end{split}$$

式中, $\sigma_{\min}(\cdot)$ 代表最小奇异值。进而有

$$\dot{\tilde{V}}_{i} \leq -\sigma_{\min}^{q+1}(\bar{\Psi}_{i})(2\alpha_{i})^{\frac{q+1}{2}}\tilde{V}_{i}^{\frac{q+1}{2}}$$

$$\tag{27}$$

相应地,停息时间为

$$T \leq \frac{\left\|\tilde{\boldsymbol{W}}_{i}(0)\right\|_{2}^{1-q}}{\sigma_{\min}^{q+1}(\bar{\boldsymbol{\Psi}}_{i})\alpha_{i}(1-q)}$$
(28)

(2)当
$$\epsilon_i \neq 0$$
时,根据引理1可得
sm($\bar{\mu}^{T}\tilde{\Psi} + \bar{c}$)= sm($\bar{\mu}^{T}\tilde{\Psi}$)

$$\operatorname{sgn}(\bar{\boldsymbol{\psi}}_{i}^{\mathrm{T}}\tilde{\boldsymbol{W}}_{i} + \bar{\boldsymbol{\epsilon}}_{i}) = \operatorname{sgn}(\bar{\boldsymbol{\psi}}_{i}^{\mathrm{T}}\tilde{\boldsymbol{W}}_{i})$$

$${}^{\underline{\mathrm{W}}}\overline{\mathbf{n}}\overline{\mathbf{n}} \, {}^{\underline{\mathrm{M}}} \, {}^{\underline{\mathrm{M}}}$$

$$(29)$$

$$\begin{split} & |\bar{\boldsymbol{\psi}}_{i}^{\mathrm{T}}\tilde{\boldsymbol{W}}_{i}|^{q} - |\bar{\boldsymbol{\epsilon}}_{i}|^{q} \leq |\bar{\boldsymbol{\psi}}_{i}^{\mathrm{T}}\tilde{\boldsymbol{W}}_{i} + \bar{\boldsymbol{\epsilon}}_{i}|^{q} \end{split} \tag{30} \\ & \text{R}\text{Bz}(29) \text{Az}(30) \overline{0}$$

$$\dot{\tilde{V}}_{i} \leq -\sigma_{\min}^{q+1}(\bar{\boldsymbol{\Psi}}_{i}) \left\| \tilde{\boldsymbol{W}}_{i} \right\|_{2}^{q+1} + \sum_{s=t_{0}}^{t_{i}} \left| \bar{\boldsymbol{\psi}}_{i}^{\mathrm{T}}(s) \tilde{\boldsymbol{W}}_{i}(t) \right\| \left| \bar{\boldsymbol{\epsilon}}_{i}(s) \right|^{q} + \left\| \bar{\boldsymbol{\psi}}_{i}^{\mathrm{T}}(t) \tilde{\boldsymbol{W}}_{i}(t) \right\| \left| \bar{\boldsymbol{\epsilon}}_{i}(t) \right|^{q}$$

$$(31)$$

进而有

$$\begin{split} \dot{\tilde{V}}_{i} &\leq -(1-\theta_{i})\sigma_{\min}^{q+1}(\bar{\Psi}_{i}) \left\| \tilde{\boldsymbol{W}}_{i} \right\|_{2}^{q+1} - \theta_{i}\sigma_{\min}^{q+1}(\bar{\Psi}_{i}) \left\| \tilde{\boldsymbol{W}}_{i} \right\|_{2}^{q+1} + (k+1)\bar{\epsilon}_{m}^{q}\bar{\Psi}_{m} \left\| \tilde{\boldsymbol{W}}_{i} \right\|_{2} \\ \end{split}$$
因此有

$$\begin{split} \dot{\tilde{V}}_{i} &\leq -(1-\theta_{i})\sigma_{\min}^{q+1}(\bar{\Psi}_{i}) \left\| \tilde{\boldsymbol{W}}_{i} \right\|_{2}^{q+1}, \left\| \tilde{\boldsymbol{W}}_{i} \right\|_{2} \geq \mu_{i} \\ \end{cases} (32)$$
根据比较引理可得出停息时间为

$$T \leq \frac{\left\|\tilde{\boldsymbol{W}}_{i}(0)\right\|_{2}^{1-q} - \mu_{i}^{1-q}}{\alpha_{i}\sigma_{\min}^{q+1}(\bar{\boldsymbol{\Psi}}_{i})(1-\theta_{i})(1-q)}$$

$$\mathcal{U}\hat{\boldsymbol{z}}^{\text{H}}\boldsymbol{u}\boldsymbol{\tilde{\boldsymbol{v}}}_{\circ}$$
(33)

5 仿真结果

假设系统的通信拓扑图如图1所示。图1中,1、2、3、4 为无人机,0为虚拟领导者,对应的拉普拉斯矩阵和牵引矩 阵分别为

<i>L</i> =	1	-1	0	0]	<i>G</i> =	1	0	0	0]
	-1	1	0	0		0	0	0	0
	0	0	1	-1		0	0	1	0
	0	0	-1	1		0	0	0	0

给定攻击抑制水平为8,参数q=0.1,所有神经网络的学习 率为0.1,激活函数为多项式函数。虚拟领导者的期望轨迹和 期望速度(矢量)为 p_0 = [15sin(0.1t), 10cos(0.1t), 10t]^T, v_0 = [1.5cos(0.1t), -sin(0.1t), 10]^T。仿真结果如图2~图7所示。 由图2~图4可得,所有无人机均能到达预定的编队位



Fig.2 Components of UAVs' positions in x direction

置并跟踪上期望轨迹。由图 5~图 7 可得,所有无人机都能 达到期望速度。因此,文中给出的最优安全控制方案可抵 御FDI攻击并实现编队跟踪。



图3 无人机位置在y方向的分量





图4 无人机位置在z方向的分量





图5 无人机速度在x方向的分量





图6 无人机速度在y方向的分量





图 7 无人机速度在z方向的分量 Fig.7 Components of UAVs'velocities in z direction

6 结束语

本文研究了无人机控制信号遭受FDI攻击时的最优安 全跟踪控制问题,利用零和图博弈理论对攻击者与安全控 制器之间的攻防进行了建模,引入了有限时间收敛的强化 学习方法和单评价神经网络架构,在线求解了最优安全控 制器。在后续工作中,可基于无模型在线强化学习方法实 现最优安全编队跟踪控制器设计。

参考文献

- Wang Haijun, Zhao Haitao, Zhang Jiao, et al. Survey on unmanned aerial vehicle networks: A cyber physical system perspective [J]. IEEE Communications Surveys & Tutorial, 2020, 22(2): 1027-1070.
- [2] Li Xiaomeng, Zhou Qi, Li Panshuo, et al. Event-triggered consensus control for multi-agent systems against false datainjection attacks [J]. IEEE Transactions on Cybernetics, 2020, 50(5): 1856-1866.
- [3] Tan Yushun, Liu Qingyi, Liu Jinliang, et al. Observer-based security control for interconnected semi-Markovian jump systems with unknown transition probabilities [J]. IEEE Transactions on Cybernetics, 2022, 52(9): 9013-9025.
- [4] Lin Hong, Sun Pei, Cai Chenxiao, et al. Secure LQG control for a quadrotor under false data injection attacks [J]. IET Control Theory & Applications, 2022, 16(9): 925-934.
- [5] Xiao Jiaping, Feroskhan M. Cyber attack detection and isolation for a quadrotor UAV with modified sliding innovation sequences [J]. IEEE Transactions on Vehicular Technology, 2022, 71(7): 7202-7214.
- [6] Yin Tingting, Gu Zhou, Park J H, et al. Event-based intermittent formation control of multi-UAV systems under

deception attacks [J]. IEEE Transactions on Neural Networks and Learning Systems, 2022, 12: 1-12.

- [7] Lewis F L, Vrabie J, Vamvoudakis K G, et al. Reinforcement learning and feedback control: Using natural decision methods to design optimal adaptive controllers[J]. IEEE Control Systems Magazine, 2012, 32(6): 76-105.
- [8] Peng Zhinan, Luo Rui, Hu Jiangping, et al. Optimal tracking control of nonlinear multiagent systems using internal reinforce Q-learning [J]. IEEE Transactions on Neural Networks and Learning Systems, 2022, 33(8): 4043-4055.
- [9] Xie Kedi, Yu Xiao, Lan Weiyao. Optimal output regulation for unknown continuous-time linear systems by internal model and adaptive dynamic programming [J]. Automatica, 2022, 146: 1-7.
- [10] Wei Qinglai, Zhu Liao, Li Tao, et al. A new approach to finitehorizon optimal control for discrete-time affine nonlinear systems via a pseudolinear method [J]. IEEE Transactions on Automatic Control, 2022, 67(5): 2610-2617.
- [11] Wu Chengwei, Li Xiaolei, Pan Wei, et al. Zero-sum game

based optimal secure control under actuator attacks[J]. IEEE Transactions on Automatic Control, 2021, 66(8): 3773-3780.

- [12] Zhou Yuanqiang, Vamvoudakis K G, Haddad W M, et al. A secure control learning framework for cyber-physical systems under sensor and actuator attacks [J]. IEEE Transactions on Cybernetics, 2021, 51(9): 4648-4660.
- [13] Moghadam R, Modares H. Resilient autonomous control of distributed multiagent systems in contested environments [J].
 IEEE Transactions on Cybernetics, 2019, 49(11): 3957-3967.
- [14] Xu Yuanyuan, Li Tieshan, Yang Yue, et al. Simplified ADP for event-triggered control of multiagent systems against FDI attacks[J]. IEEE Transactions on Systems, Man, and Cybernetics: System, 2023, 53(8): 4672-4683.
- [15] Kokolakis N-M T, Vamvoudakis K G. Safety-aware pursuitevasion games in unknown environments using Gaussian processes and finite-time convergent reinforcement learning[J]. IEEE Transactions on Neural Networks and Learning Systems, 2024,35(3): 3130-3143.

Optimal Secure Tracking Control in Multi-UAVs Based on Online Reinforcement Learning

Gong Zhenyu, Yang Feisheng

Northwestern Polytechnical University, Xi'an 710072, China

Abstract: In Unmanned Aerial Vehicle (UAV) formation tracking missions, False Data Injection (FDI) attackers can inject misleading data into the control commands, resulting in the fact that UAVs can not form the specified formation configuration, so there is a need to design a secure formation tracking controller. The attack-defense process was modeled as a zero-sum graphical game, in which the FDI attacker and the secure controller were viewed as game players. The attacker aims to maximize the cost function yet the secure controller serves a contrary purpose. Solving the game and acquiring the optimal secure control policy rely on solving the Hamilton-Jacobi-Isaacs (HJI) equation. The HJI equation is a coupled partial differential equation, which is difficult to solve directly. Therefore, the finite-time convergent online reinforcement learning algorithm that combines the experience replay mechanism was introduced and the critic-only neural network was utilized to approximate the value function for obtaining the optimal secure control policy. A numerical simulation was given to show the effectiveness of the raised scheme.

Key Words: FDI attack; multi-UAVs; online reinforcement learning; optimal control; zero-sum graphical game

Received: 2023-06-27; Revised: 2023-11-16; Accepted: 2023-12-25

Foundation item: National Natural Science Foundation of China (62073269); Aeronautical Science Foundation of China (2020Z034053002); Key Research and Development Program of Shaanxi (2022GY-244); Natural Science Foundation of Chongqing (CSTB2022NSCQ-MSX0963); Guangdong Basic and Applied Basic Research Foundation(2023A1515011220)