



# Optimal synchronization control for multi-agent systems with input saturation: a nonzero-sum game\*

Hongyang LI<sup>1,2</sup>, Qinglai WEI<sup>†1,2,3</sup>

<sup>1</sup>School of Artificial Intelligence, University of Chinese Academy of Sciences, Beijing 100049, China

<sup>2</sup>The State Key Laboratory for Management and Control of Complex Systems, Institute of Automation, Chinese Academy of Sciences, Beijing 100190, China

<sup>3</sup>Institute of Systems Engineering, Macau University of Science and Technology, Macau 999078, China

E-mail: lihongyang2019@ia.ac.cn; qinglai.wei@ia.ac.cn

Received Jan. 14, 2022; Revision accepted Mar. 7, 2022; Crosschecked Apr. 29, 2022; Published online May 20, 2022

**Abstract:** This paper presents a novel optimal synchronization control method for multi-agent systems with input saturation. The multi-agent game theory is introduced to transform the optimal synchronization control problem into a multi-agent nonzero-sum game. Then, the Nash equilibrium can be achieved by solving the coupled Hamilton–Jacobi–Bellman (HJB) equations with nonquadratic input energy terms. A novel off-policy reinforcement learning method is presented to obtain the Nash equilibrium solution without the system models, and the critic neural networks (NNs) and actor NNs are introduced to implement the presented method. Theoretical analysis is provided, which shows that the iterative control laws converge to the Nash equilibrium. Simulation results show the good performance of the presented method.

**Key words:** Optimal synchronization control; Multi-agent systems; Nonzero-sum game; Adaptive dynamic programming; Input saturation; Off-policy reinforcement learning; Policy iteration

<https://doi.org/10.1631/FITEE.2200010>

**CLC number:** TP13

## 1 Introduction

Multi-agent synchronization control has attracted much attention due to its high efficiency and computational performance (Wieland et al., 2011; Wei et al., 2018, 2020, 2021; Li JQ et al., 2021; Rehák and Lynnyk, 2021; Zhang KQ et al., 2021). Generally speaking, synchronization control problems require that the agents synthesize to the same value (Cao et al., 2015; Garcia et al., 2017; Yang JY et al., 2019) or track the trajectories of leaders (Du et al., 2014; Zhao et al., 2014) by designing distributed control laws. Because of the practical signif-

icance of multi-agent systems, many researchers have devoted themselves to tackling various synchronization control problems, including switching topologies (Thunberg et al., 2014), system faults (Ma and Yang, 2016) and so on (Han et al., 2013; Wei et al., 2015; He et al., 2018). Among the research works on synchronization control, optimal synchronization control, which requires each agent to minimize its own local performance index function, is a promising research direction. Multi-agent cooperative games provide an effective tool to study multi-agent optimal control problems, and they rely on solving coupled Hamilton–Jacobi (HJ) equations (Vamvoudakis et al., 2012). However, coupled HJ equations are hard to solve, which limits the applications of cooperative game theory in synchronization control problems.

Reinforcement learning is an effective method

<sup>†</sup> Corresponding author

\* Project supported by the National Key R&D Program of China (No. 2018YFB1702300) and the National Natural Science Foundation of China (Nos. 61722312 and 61533017)

ORCID: Hongyang LI, <https://orcid.org/0000-0001-5891-134X>; Qinglai WEI, <https://orcid.org/0000-0001-7002-9800>

© Zhejiang University Press 2022

for solving coupled HJ equations. The main idea of reinforcement learning is to solve the coupled HJ equations forward by time, which can reduce the computational burden (Wang et al., 2009; Wei and Liu, 2014; Wei et al., 2014, 2016, 2017; Zhang HG et al., 2015; Yang N et al., 2019; Zhang LD et al., 2019). In recent years, reinforcement learning has been further developed to solve multi-agent cooperative game problems. Vamvoudakis et al. (2012) proposed an online policy iteration method for optimal synchronization control problems; however, the external disturbance was not considered. In Jiao et al. (2016), a novel policy iteration method was proposed for the multi-agent zero-sum game problem, and disturbance rejection was achieved. In Wei et al. (2015), the graphical game was studied for heterogeneous multi-agent systems. The off-policy reinforcement learning method was proposed to solve multi-agent synchronization control problems by Li JN et al. (2017), and the input constraint was considered by Qin et al. (2019). However, there are few research results considering cooperative game problems with input saturation, which motivates our study.

In this paper, the multi-agent optimal synchronization control problem with input saturation is studied based on cooperative game theory and reinforcement learning. Compared with Qin et al. (2019), we consider coupled terms with neighboring agents in the performance index functions. The main contributions can be summarized as follows:

1. A novel off-policy reinforcement learning method is presented without the information of system models for cooperative game problems of multi-agent systems. The control constraint and coupled terms in the performance index functions are considered which broaden the application scope of the presented method.

2. The characteristics of the presented model-free off-policy reinforcement learning method, including convergence and optimality, are analyzed, showing that the solutions obtained from the presented method converge to the Nash equilibrium.

3. Critic neural networks (NNs) and actor NNs are used to implement the off-policy reinforcement learning algorithm. Simulation results verify the good performance of the presented method.

## 2 Problem formulation

### 2.1 Graph theory

Let  $G_r = (V, \varepsilon, \mathbf{E})$  be a directed graph, where  $V = \{v_1, v_2, \dots, v_N\}$  denotes the nonempty finite vertex set. Furthermore,  $\varepsilon \subseteq V \times V$  is the set of edges. An edge of graph  $G_r$  is denoted as  $\varepsilon_{ij}$ , which means that agent  $j$  is a neighbor of agent  $i$ .  $\mathbf{E} = [e_{ij}] \in \mathbb{R}^{N \times N}$  is the adjacency matrix, where  $e_{ij}$  represents the weight of edge  $\varepsilon_{ij}$ . If  $\varepsilon_{ij} \in \varepsilon$ ,  $e_{ij} > 0$ ; otherwise,  $e_{ij} = 0$ . Let the set of neighbors of agent  $i$  be  $N_i = \{v_j | (v_j, v_i) \in \varepsilon\}$ . Define  $\mathbf{G} = \text{diag}(g_i) \in \mathbb{R}^{N \times N}$  as the pinning matrix. If agent  $i$  has access to the leader,  $g_i > 0$ ; otherwise,  $g_i = 0$ . Define the Laplacian matrix as  $\mathbf{L} = \mathbf{D} - \mathbf{E}$ , where  $\mathbf{D} = \text{diag}(d_i)$  and  $d_i = \sum_j e_{ij}$ .

### 2.2 Multi-agent synchronization control

For  $i = 1, 2, \dots, N$ , consider the following systems:

$$\dot{\mathbf{x}}_i = \mathbf{A}\mathbf{x}_i + \mathbf{B}\mathbf{u}_i, \quad (1)$$

where  $\mathbf{x}_i \in \mathbb{R}^n$  and  $\mathbf{u}_i \in \mathcal{U}_i \subset \mathbb{R}^m$  are the system state and control, respectively. Here,  $\mathbf{A}$  and  $\mathbf{B}$  are system matrices with suitable dimensions.  $\mathcal{U}_i = \{\mathbf{u}_i | \mathbf{u}_i \in \mathbb{R}^m, \|\mathbf{u}_i\|_\infty \leq \lambda_i\}$  ( $\lambda_i > 0$  is a known constant). Let the leader dynamics be

$$\dot{\mathbf{x}}_0 = \mathbf{A}\mathbf{x}_0, \quad (2)$$

where  $\mathbf{x}_0 \in \mathbb{R}^n$  is the system state. Then, we can define the synchronization error as

$$\boldsymbol{\delta}_i = \sum_{j \in N_i} e_{ij} (\mathbf{x}_i - \mathbf{x}_j) + g_i (\mathbf{x}_i - \mathbf{x}_0). \quad (3)$$

Taking the derivative of Eq. (3), we have

$$\dot{\boldsymbol{\delta}}_i = \mathbf{A}\boldsymbol{\delta}_i + (d_i + g_i) \mathbf{B}\mathbf{u}_i - \sum_{j \in N_i} e_{ij} \mathbf{B}\mathbf{u}_j. \quad (4)$$

For system (4), the performance index function can be given as

$$\begin{aligned} & J_i(\boldsymbol{\delta}_i(0), \mathbf{u}_i, \mathbf{u}_{-i}) \\ &= \int_0^\infty \left( \boldsymbol{\delta}_i^\top \mathbf{Q}_{ii} \boldsymbol{\delta}_i + R_i(\mathbf{u}_i) + \sum_{j \in N_i} R_i(\mathbf{u}_j) \right) dt, \quad (5) \end{aligned}$$

where the term  $\mathbf{u}_{-i}$  represents the policies of the neighbors of agent  $i$ .  $\mathbf{Q}_{ii} > 0$ ,

$$\begin{aligned} R_i(\mathbf{u}_i) &= 2 \int_0^{\mathbf{u}_i} (\lambda_i \Psi^{-1}(t/\lambda_i))^T \mathbf{R}_{ii} dt, \\ R_i(\mathbf{u}_j) &= 2 \int_0^{\mathbf{u}_j} (\lambda_j \Psi^{-1}(t/\lambda_j))^T \mathbf{R}_{ij} dt, \end{aligned}$$

where

$$\begin{aligned} \mathbf{t} &= [t_1, t_2, \dots, t_m]^T \in \mathbb{R}^m, \\ \mathbf{R}_{ii} &= \text{diag}(r_{ii,1}, r_{ii,2}, \dots, r_{ii,m}) > 0, \\ \mathbf{R}_{ij} &= \text{diag}(r_{ij,1}, r_{ij,2}, \dots, r_{ij,m}) \geq 0, \end{aligned}$$

and  $\Psi^{-1}$  is the inverse function of the hyperbolic tangent function (i.e.,  $\Psi^{-1} = \text{arctanh}(\cdot)$ , or equivalently,  $\Psi(\cdot) \triangleq \tanh(\cdot)$ ). Then,  $R_i(\mathbf{u}_i)$  and  $R_i(\mathbf{u}_j)$  can be written as

$$\begin{aligned} R_i(\mathbf{u}_i) &= 2 \int_0^{\mathbf{u}_i} (\lambda_i \Psi^{-1}(t/\lambda_i))^T \mathbf{R}_{ii} dt \\ &= 2 \sum_{q=1}^m r_{ii,q} \int_0^{\mathbf{u}_{iq}} \lambda_i \Psi^{-1}(t_q/\lambda_i) dt_q, \quad (6) \end{aligned}$$

$$\begin{aligned} R_i(\mathbf{u}_j) &= 2 \int_0^{\mathbf{u}_j} (\lambda_j \Psi^{-1}(t/\lambda_j))^T \mathbf{R}_{ij} dt \\ &= 2 \sum_{q=1}^m r_{ij,q} \int_0^{\mathbf{u}_{jq}} \lambda_j \Psi^{-1}(t_q/\lambda_j) dt_q. \quad (7) \end{aligned}$$

For Eq. (5), the Nash equilibrium condition can be described as

$$J_i(\mathbf{u}_i^*, \mathbf{u}_{-i}^*) \leq J_i(\mathbf{u}_i, \mathbf{u}_{-i}^*), \quad i = 1, 2, \dots, N. \quad (8)$$

For agent  $i$ , we define the iterative value function as

$$\begin{aligned} V_i(\boldsymbol{\delta}_i(t)) &= \int_t^\infty \left( \boldsymbol{\delta}_i^T \mathbf{Q}_{ii} \boldsymbol{\delta}_i + R_i(\mathbf{u}_i) + \sum_{j \in N_i} R_i(\mathbf{u}_j) \right) d\tau. \quad (9) \end{aligned}$$

Then, we can obtain the Bellman equation as

$$\begin{aligned} H_i \left( \boldsymbol{\delta}_i, \frac{\partial V_i}{\partial \boldsymbol{\delta}_i}, \mathbf{u}_i, \mathbf{u}_{-i} \right) &= \left( \frac{\partial V_i}{\partial \boldsymbol{\delta}_i} \right)^T \left( \mathbf{A} \boldsymbol{\delta}_i + (d_i + g_i) \mathbf{B} \mathbf{u}_i - \sum_{j \in N_i} e_{ij} \mathbf{B} \mathbf{u}_j \right) \\ &+ \boldsymbol{\delta}_i^T \mathbf{Q}_{ii} \boldsymbol{\delta}_i + R_i(\mathbf{u}_i) + \sum_{j \in N_i} R_i(\mathbf{u}_j) \quad (10) \end{aligned}$$

with  $V_i(\mathbf{0}) = 0$ . According to the stationary condition (Bertsekas, 2007), it can be derived that

$$\begin{aligned} \frac{\partial H_i}{\partial \mathbf{u}_i} = \mathbf{0} &\Rightarrow \\ \mathbf{u}_i^* &= -\lambda_i \Psi \left( \frac{1}{2\lambda_i} (d_i + g_i) \mathbf{R}_{ii}^{-1} \mathbf{B}^T \frac{\partial V_i^*}{\partial \boldsymbol{\delta}_i} \right). \quad (11) \end{aligned}$$

Substituting Eq. (11) into Eq. (10), we can obtain the Hamilton–Jacobi–Bellman (HJB) equation as

$$\begin{aligned} &\left( \frac{\partial V_i^*}{\partial \boldsymbol{\delta}_i} \right)^T \left( \mathbf{A} \boldsymbol{\delta}_i - \lambda_i (d_i + g_i) \mathbf{B} \Psi(\boldsymbol{\Delta}_i^*) \right) \\ &+ \sum_{j \in N_i} \lambda_j e_{ij} \mathbf{B} \Psi(\boldsymbol{\Delta}_j^*) + \boldsymbol{\delta}_i^T \mathbf{Q}_{ii} \boldsymbol{\delta}_i + R_i(-\lambda_i \Psi(\boldsymbol{\Delta}_i^*)) \\ &+ \sum_{j \in N_i} R_i(-\lambda_j \Psi(\boldsymbol{\Delta}_j^*)) = 0 \quad (12) \end{aligned}$$

with  $V_i(\mathbf{0}) = 0$ , where

$$\boldsymbol{\Delta}_i^* = \frac{1}{2\lambda_i} (d_i + g_i) \mathbf{R}_{ii}^{-1} \mathbf{B}^T \frac{\partial V_i^*}{\partial \boldsymbol{\delta}_i}.$$

We would like to design  $\mathbf{u}_i$ , such that the Nash equilibrium condition represented in inequality (8) and the following state synchronization condition are satisfied:

$$\lim_{t \rightarrow \infty} \|\mathbf{x}_i - \mathbf{x}_0\| = 0, \quad i = 1, 2, \dots, N. \quad (13)$$

**Remark 1** The performance index function considered by Qin et al. (2019) is defined as

$$\begin{aligned} J_i(\boldsymbol{\delta}_i(0), \mathbf{u}_i, \mathbf{u}_{-i}) &= \int_0^\infty (\boldsymbol{\delta}_i^T \mathbf{Q}_{ii} \boldsymbol{\delta}_i + R_i(\mathbf{u}_i)) dt. \quad (14) \end{aligned}$$

Comparing Eqs. (5) and (14), it can be seen that the coupled terms with neighboring agents in the performance index function are considered in this study. In multi-agent systems, the behavior of agent  $i$  may have an impact on its neighboring agents. Therefore, the performance index function (5) is more natural for the optimal synchronization control of multi-agent systems.

**Remark 2** Based on Eq. (3), it can be derived that

$$\boldsymbol{\delta} = ((\mathbf{L} + \mathbf{G}) \otimes \mathbf{I}_n) (\mathbf{x} - \bar{\mathbf{x}}_0),$$

where

$$\begin{aligned} \boldsymbol{\delta} &= [\boldsymbol{\delta}_1^T, \boldsymbol{\delta}_2^T, \dots, \boldsymbol{\delta}_N^T]^T, \\ \mathbf{x} &= [\mathbf{x}_1^T, \mathbf{x}_2^T, \dots, \mathbf{x}_N^T]^T, \\ \bar{\mathbf{x}}_0 &= [\mathbf{x}_0^T, \mathbf{x}_0^T, \dots, \mathbf{x}_0^T]^T, \end{aligned}$$

“ $\otimes$ ” is the Kronecker product, and  $\mathbf{I}_n$  is an identity matrix with dimension  $n$ . According to Vamvoudakis et al. (2012), we have

$$\|\mathbf{x} - \bar{\mathbf{x}}_0\| \leq \|\delta\|/\sigma_{\min}(\mathbf{L} + \mathbf{G}),$$

where  $\sigma_{\min}(\cdot)$  represents the minimum singular value of the matrix. Therefore, the stability of the tracking error dynamics represented in Eq. (4) guarantees the state synchronization condition denoted by Eq. (13).

### 3 Main results

#### 3.1 Multi-agent nonzero-sum game

A theorem is provided, which shows that the solution to the HJB equation (i.e., Eq. (12)) satisfies the Nash equilibrium condition (i.e., inequality (8)) under certain conditions.

**Theorem 1** Assume that the optimal control law  $\mathbf{u}_i^*$  is given as shown in Eq. (11), and that  $V_i$  is the positive definite smooth solution to the HJB equation (12). Then, system (4) is asymptotically stable. The optimal control laws  $\mathbf{u}_i^*$  ( $i = 1, 2, \dots, N$ ) constitute the Nash equilibrium, and the solution  $V_i$  to the HJB equation (12) is the optimal value of the game, i.e.,

$$J_i^*(\delta_i(0), \mathbf{u}_i^*, \mathbf{u}_{-i}^*) = V_i(\delta_i(0)). \quad (15)$$

**Proof** Choosing the iterative value function  $V_i$  as the Lyapunov function, it can be derived that

$$\frac{dV_i}{dt} = \left(\frac{\partial V_i}{\partial \delta_i}\right)^T \left(\mathbf{A}\delta_i + (d_i + g_i)\mathbf{B}\mathbf{u}_i - \sum_{j \in N_i} e_{ij}\mathbf{B}\mathbf{u}_j\right). \quad (16)$$

For the right-hand side of the Bellman equation, i.e., Eq. (10), adding and subtracting  $R_i(\mathbf{u}_i^*) + \sum_{j \in N_i} R_i(\mathbf{u}_j^*)$  and  $\left(\frac{\partial V_i}{\partial \delta_i}\right)^T \left(\mathbf{A}\delta_i + (d_i + g_i)\mathbf{B}\mathbf{u}_i^* - \sum_{j \in N_i} e_{ij}\mathbf{B}\mathbf{u}_j^*\right)$ , it can be derived that

$$\begin{aligned} & H_i\left(\delta_i, \frac{\partial V_i}{\partial \delta_i}, \mathbf{u}_i, \mathbf{u}_{-i}\right) \\ &= H_i\left(\delta_i, \frac{\partial V_i}{\partial \delta_i}, \mathbf{u}_i^*, \mathbf{u}_{-i}^*\right) + R_i(\mathbf{u}_i) - R_i(\mathbf{u}_i^*) \\ & \quad + \sum_{j \in N_i} (R_i(\mathbf{u}_j) - R_i(\mathbf{u}_j^*)) \end{aligned}$$

$$\begin{aligned} & + \left(\frac{\partial V_i}{\partial \delta_i}\right)^T (d_i + g_i)\mathbf{B}(\mathbf{u}_i - \mathbf{u}_i^*) \\ & - \left(\frac{\partial V_i}{\partial \delta_i}\right)^T \sum_{j \in N_i} e_{ij}\mathbf{B}(\mathbf{u}_j - \mathbf{u}_j^*). \quad (17) \end{aligned}$$

According to Eq. (12), we have

$$H_i\left(\delta_i, \frac{\partial V_i}{\partial \delta_i}, \mathbf{u}_i^*, \mathbf{u}_{-i}^*\right) = 0.$$

Letting  $\mathbf{u}_i = \mathbf{u}_i^*$  and  $\mathbf{u}_j = \mathbf{u}_j^*$  ( $j \in N_i$ ), Eq. (17) can be written as

$$\frac{dV_i}{dt} + \delta_i^T \mathbf{Q}_{ii} \delta_i + R_i(\mathbf{u}_i) + \sum_{j \in N_i} R_i(\mathbf{u}_j) = 0. \quad (18)$$

It can be derived that  $\frac{dV_i}{dt} < 0$ , and thus system (4) is asymptotically stable.

Based on Eqs. (5) and (9), it can be derived that

$$\begin{aligned} & J_i(\delta_i(0), \mathbf{u}_i, \mathbf{u}_{-i}) \\ &= \int_0^\infty \left( \delta_i^T \mathbf{Q}_{ii} \delta_i + R_i(\mathbf{u}_i) + \sum_{j \in N_i} R_i(\mathbf{u}_j) \right) dt \\ & \quad + V_i(\delta_i(0)) - V_i(\delta_i(\infty)) \\ & \quad + \int_0^\infty \left(\frac{\partial V_i}{\partial \delta_i}\right)^T \left( \mathbf{A}\delta_i + (d_i + g_i)\mathbf{B}\mathbf{u}_i \right. \\ & \quad \left. - \sum_{j \in N_i} e_{ij}\mathbf{B}\mathbf{u}_j \right) dt. \quad (19) \end{aligned}$$

Because of the asymptotic stability of system (4) and the boundary condition  $V_i(\mathbf{0}) = 0$ , it can be derived that  $V_i(\delta_i(\infty)) = 0$ . Then, substituting Eq. (12) into Eq. (19) and completing the squares, we can obtain

$$\begin{aligned} & J_i(\delta_i(0), \mathbf{u}_i, \mathbf{u}_{-i}) \\ &= V_i(\delta_i(0)) + \int_0^\infty \left( 2 \int_{\mathbf{u}_i^*/\lambda_i}^{\mathbf{u}_i/\lambda_i} \lambda_i^2 (\Psi^{-1}(\mathbf{s})) \right. \\ & \quad \left. - \Psi^{-1}(\mathbf{u}_i^*/\lambda_i) \right)^T \mathbf{R}_{ii} d\mathbf{s} \\ & \quad + 2 \sum_{j \in N_i} \int_{\mathbf{u}_j^*/\lambda_j}^{\mathbf{u}_j/\lambda_j} \lambda_j^2 (\Psi^{-1}(\mathbf{s}))^T \mathbf{R}_{ij} d\mathbf{s} \\ & \quad - \left(\frac{\partial V_i}{\partial \delta_i}\right)^T \sum_{j \in N_i} e_{ij}\mathbf{B}(\mathbf{u}_j - \mathbf{u}_j^*) \right) dt. \quad (20) \end{aligned}$$

For  $\mathbf{u}_{-i} = \mathbf{u}_{-i}^*$ , it can be obtained that

$$\begin{aligned} & J_i(\delta_i(0), \mathbf{u}_i, \mathbf{u}_{-i}^*) \\ &= V_i(\delta_i(0)) + \int_0^\infty 2\lambda_i^2 \int_{\mathbf{u}_i^*/\lambda_i}^{\mathbf{u}_i/\lambda_i} (\Psi^{-1}(\mathbf{s})) \\ & \quad - \Psi^{-1}(\mathbf{u}_i^*/\lambda_i) \right)^T \mathbf{R}_{ii} d\mathbf{s} dt. \quad (21) \end{aligned}$$

For Eq. (21), it can be derived that

$$\begin{aligned} & \int_{\mathbf{u}_i^*/\lambda_i}^{\mathbf{u}_i/\lambda_i} (\Psi^{-1}(\mathbf{s}) - \Psi^{-1}(\mathbf{u}_i^*/\lambda_i))^T \mathbf{R}_{ii} d\mathbf{s} \\ &= \sum_{q=1}^m r_{ii,q} \int_{u_{iq}^*/\lambda_i}^{u_{iq}/\lambda_i} (\Psi^{-1}(s_q) - \Psi^{-1}(u_{iq}^*/\lambda_i)) ds_q, \end{aligned} \quad (22)$$

where  $\Psi^{-1} = \operatorname{arctanh}(\cdot)$  is monotonically increasing, i.e.,  $(\Psi^{-1})' > 0$ . Therefore, based on the mean value theorem for integrals, it can be derived that

$$\begin{aligned} & \int_{u_{iq}^*/\lambda_i}^{u_{iq}/\lambda_i} (\Psi^{-1}(s_q) - \Psi^{-1}(u_{iq}^*/\lambda_i)) ds_q \\ &= (\Psi^{-1}(\bar{s}_q) - \Psi^{-1}(u_{iq}^*/\lambda_i)) (u_{iq}/\lambda_i - u_{iq}^*/\lambda_i) \\ &> 0, \end{aligned} \quad (23)$$

where  $\bar{s}_q$  is between  $u_{iq}^*/\lambda_i$  and  $u_{iq}/\lambda_i$ . Therefore, the second term of the right-hand side of Eq. (21) is positive. If  $\mathbf{u}_i = \mathbf{u}_i^*$ , it can be derived that the solution  $V_i$  to Eq. (12) is the optimal value of the game. Comparing Eqs. (21) and (15), it can be obtained that the condition expressed in inequality (8) is satisfied. The proof is completed.

### 3.2 Policy iteration method for solving the HJB equation

In the previous subsection, it was derived that the optimal control, represented in Eq. (11), can be calculated to construct the Nash equilibrium represented in inequality (8). However, the optimal control in Eq. (11) requires the information of  $V_i^*$ , which can be calculated from the HJB equation (12). The HJB equation (12) is a nonlinear partial differential equation, which is hard to solve analytically. Therefore, a policy iteration method is provided (Algorithm 1) to solve the HJB equation (12) numerically. Then, a theorem can be provided, which shows the convergence of the presented policy iteration method.

**Theorem 2** Assume that agent  $i$  and its neighbors update their control policies according to Algorithm 1. Then, policies  $\mathbf{u}_i^k$  and  $\mathbf{u}_{-i}^k$  converge to the Nash equilibrium, and  $V_i^k$  converges to  $V_i^*$ , where  $V_i^*$  is the solution to the HJB equation (12).

**Proof** Integrating  $\dot{V}_i^k - \dot{V}_i^{k+1}$  along the system

$$\dot{\delta}_i = \mathbf{A}\delta_i + (d_i + g_i) \mathbf{B}\mathbf{u}_i^{k+1} - \sum_{j \in N_i} e_{ij} \mathbf{B}\mathbf{u}_j^k, \quad (26)$$

---

#### Algorithm 1 Policy iteration method

---

**Initialization:** Give the initial stabilizing control laws  $\mathbf{u}_i^0, i = 1, 2, \dots, N$

**Iteration:**

- 1: Let  $k = 0$
- 2: For  $i = 1, 2, \dots, N$ , solve for  $V_i^k$  from the following equation:

$$\begin{aligned} & \left( \frac{\partial V_i^k}{\partial \delta_i} \right)^T \left( \mathbf{A}\delta_i + (d_i + g_i) \mathbf{B}\mathbf{u}_i^k - \sum_{j \in N_i} e_{ij} \mathbf{B}\mathbf{u}_j^k \right) \\ & + \delta_i^T \mathbf{Q}_{ii} \delta_i + R_i(\mathbf{u}_i^k) + \sum_{j \in N_i} R_i(\mathbf{u}_j^k) = 0 \end{aligned} \quad (24)$$

- 3: For  $i = 1, 2, \dots, N$ , update the control laws as

$$\mathbf{u}_i^{k+1} = -\lambda_i \Psi \left( \frac{1}{2\lambda_i} (d_i + g_i) \mathbf{R}_{ii}^{-1} \mathbf{B}^T \frac{\partial V_i^k}{\partial \delta_i} \right) \quad (25)$$

- 4: Let  $k = k + 1$ ; repeat steps 2 and 3 until convergence

- 5: **return**  $V_i^k, \mathbf{u}_i^k, i = 1, 2, \dots, N$
- 

we have

$$\begin{aligned} & V_i^k - V_i^{k+1} \\ &= \int_t^\infty \left( \left( \frac{\partial V_i^{k+1}}{\partial \delta_i} \right)^T \dot{\delta}_i - \left( \frac{\partial V_i^k}{\partial \delta_i} \right)^T \dot{\delta}_i \right) d\tau. \end{aligned} \quad (27)$$

Based on Eq. (24), we have

$$\begin{aligned} & \left( \frac{\partial V_i^{k+1}}{\partial \delta_i} \right)^T \left( \mathbf{A}\delta_i + (d_i + g_i) \mathbf{B}\mathbf{u}_i^{k+1} \right. \\ & \left. - \sum_{j \in N_i} e_{ij} \mathbf{B}\mathbf{u}_j^{k+1} \right) + \delta_i^T \mathbf{Q}_{ii} \delta_i \\ & + R_i(\mathbf{u}_i^{k+1}) + \sum_{j \in N_i} R_i(\mathbf{u}_j^{k+1}) = 0. \end{aligned} \quad (28)$$

Subtracting and adding Eqs. (24) and (28) to the right-hand side of Eq. (27), we have

$$\begin{aligned} & \left( \frac{\partial V_i^{k+1}}{\partial \delta_i} \right)^T \dot{\delta}_i - \left( \frac{\partial V_i^k}{\partial \delta_i} \right)^T \dot{\delta}_i \\ &= \left( \frac{\partial V_i^{k+1}}{\partial \delta_i} \right)^T \sum_{j \in N_i} e_{ij} \mathbf{B}(\mathbf{u}_j^{k+1} - \mathbf{u}_j^k) - R_i(\mathbf{u}_i^{k+1}) \\ & \quad - \sum_{j \in N_i} R_i(\mathbf{u}_j^{k+1}) + R_i(\mathbf{u}_i^k) + \sum_{j \in N_i} R_i(\mathbf{u}_j^k) \\ & \quad - \left( \frac{\partial V_i^k}{\partial \delta_i} \right)^T (d_i + g_i) \mathbf{B}(\mathbf{u}_i^{k+1} - \mathbf{u}_i^k). \end{aligned} \quad (29)$$

For Eq. (29), we have

$$\begin{aligned}
 & R_i(\mathbf{u}_i^k) - R_i(\mathbf{u}_i^{k+1}) + \left(\frac{\partial V_i^k}{\partial \delta_i}\right)^T (d_i + g_i) \mathbf{B}(\mathbf{u}_i^k - \mathbf{u}_i^{k+1}) \\
 &= R_i(\mathbf{u}_i^k) - R_i(\mathbf{u}_i^{k+1}) \\
 &\quad - 2\lambda_i \left(\Psi^{-1}\left(\frac{\mathbf{u}_i^{k+1}}{\lambda_i}\right)\right)^T \mathbf{R}_{ii}(\mathbf{u}_i^k - \mathbf{u}_i^{k+1}) \\
 &= 2\lambda_i^2 \int_{\mathbf{u}_i^{k+1}/\lambda_i}^{\mathbf{u}_i^k/\lambda_i} (\Psi^{-1}(s))^T \mathbf{R}_{ii} ds \\
 &\quad - 2\lambda_i^2 \int_{\mathbf{u}_i^{k+1}/\lambda_i}^{\mathbf{u}_i^k/\lambda_i} \left(\Psi^{-1}\left(\frac{\mathbf{u}_i^{k+1}}{\lambda_i}\right)\right)^T \mathbf{R}_{ii} ds \\
 &= 2\lambda_i^2 \int_{\mathbf{u}_i^{k+1}/\lambda_i}^{\mathbf{u}_i^k/\lambda_i} \left(\Psi^{-1}(s) - \Psi^{-1}\left(\frac{\mathbf{u}_i^{k+1}}{\lambda_i}\right)\right)^T \mathbf{R}_{ii} ds
 \end{aligned} \tag{30}$$

and

$$\begin{aligned}
 & \sum_{j \in N_i} R_i(\mathbf{u}_j^k) - \sum_{j \in N_i} R_i(\mathbf{u}_j^{k+1}) \\
 &= \sum_{j \in N_i} 2\lambda_j^2 \int_{\mathbf{u}_j^{k+1}/\lambda_j}^{\mathbf{u}_j^k/\lambda_j} (\Psi^{-1}(s))^T \mathbf{R}_{ij} ds.
 \end{aligned} \tag{31}$$

Therefore, we can rewrite Eq. (27) as

$$\begin{aligned}
 & V_i^k - V_i^{k+1} \\
 &= \int_t^\infty \left( \sum_{j \in N_i} 2\lambda_j^2 \int_{\mathbf{u}_j^{k+1}/\lambda_j}^{\mathbf{u}_j^k/\lambda_j} (\Psi^{-1}(s))^T \mathbf{R}_{ij} ds \right. \\
 &\quad \left. + 2\lambda_i^2 \int_{\mathbf{u}_i^{k+1}/\lambda_i}^{\mathbf{u}_i^k/\lambda_i} \left(\Psi^{-1}(s) - \Psi^{-1}\left(\frac{\mathbf{u}_i^{k+1}}{\lambda_i}\right)\right)^T \mathbf{R}_{ii} ds \right. \\
 &\quad \left. + \left(\frac{\partial V_i^{k+1}}{\partial \delta_i}\right)^T \sum_{j \in N_i} e_{ij} \mathbf{B}(\mathbf{u}_j^{k+1} - \mathbf{u}_j^k) \right) d\tau.
 \end{aligned} \tag{32}$$

Based on Eq. (22) and inequality (23), it can be derived that

$$\int_{\mathbf{u}_i^{k+1}/\lambda_i}^{\mathbf{u}_i^k/\lambda_i} \left(\Psi^{-1}(s) - \Psi^{-1}\left(\frac{\mathbf{u}_i^{k+1}}{\lambda_i}\right)\right)^T \mathbf{R}_{ii} ds > 0. \tag{33}$$

For Eq. (32), if  $e_{ij}$  and  $\sigma_{\max}(\mathbf{R}_{ij})$  are small enough, where  $\sigma_{\max}(\cdot)$  represents the maximum singular value of the matrix, it can be derived that  $V_i^k \geq V_i^{k+1}$ . According to Eq. (9), we have  $V_i^k \geq 0$ . Therefore,  $V_i^k$  is convergent as  $k \rightarrow \infty$ . As  $V_i^*$  is the optimal iterative value function obtained by minimizing Eq. (9), we have  $\lim_{k \rightarrow \infty} V_i^k = V_i^\infty \geq V_i^*$ . Therefore, it

can be derived that  $V_i^k$  converges to  $V_i^*$ , where  $V_i^*$  is the solution to Eq. (12). The proof is completed.

However, system matrices are still required to solve Eqs. (24) and (25). A novel off-policy reinforcement learning method will be presented in the next subsection without the information of system matrices.

### 3.3 Model-free off-policy reinforcement learning method for solving the HJB equation

We can rewrite the tracking error dynamics, represented in Eq. (4), as follows:

$$\begin{aligned}
 \dot{\delta}_i &= \mathbf{A}\delta_i + (d_i + g_i) \mathbf{B}(\mathbf{u}_i - \mathbf{u}_i^k) + (d_i + g_i) \mathbf{B}\mathbf{u}_i^k \\
 &\quad - \sum_{j \in N_i} e_{ij} \mathbf{B}\mathbf{u}_j^k - \sum_{j \in N_i} e_{ij} \mathbf{B}(\mathbf{u}_j - \mathbf{u}_j^k).
 \end{aligned} \tag{34}$$

Taking the derivative of  $V_i^k$  along system (34), we have

$$\begin{aligned}
 \dot{V}_i^k &= \left(\frac{\partial V_i^k}{\partial \delta_i}\right)^T \dot{\delta}_i \\
 &= \left(\frac{\partial V_i^k}{\partial \delta_i}\right)^T \left( \mathbf{A}\delta_i - \sum_{j \in N_i} e_{ij} \mathbf{B}\mathbf{u}_j^k + (d_i + g_i) \mathbf{B}\mathbf{u}_i^k \right) \\
 &\quad + \left(\frac{\partial V_i^k}{\partial \delta_i}\right)^T (d_i + g_i) \mathbf{B}(\mathbf{u}_i - \mathbf{u}_i^k) \\
 &\quad - \left(\frac{\partial V_i^k}{\partial \delta_i}\right)^T \sum_{j \in N_i} e_{ij} \mathbf{B}(\mathbf{u}_j - \mathbf{u}_j^k).
 \end{aligned} \tag{35}$$

According to Eq. (25), we have

$$\Delta_i^{k+1} = \frac{1}{2\lambda_i} (d_i + g_i) \mathbf{R}_{ii}^{-1} \mathbf{B}^T \frac{\partial V_i^k}{\partial \delta_i}. \tag{36}$$

Then, substituting Eqs. (24) and (36) into Eq. (35), it can be derived that

$$\begin{aligned}
 & V_i^k(\delta_i(t')) - V_i^k(\delta_i(t)) \\
 &= \int_t^{t'} \dot{V}_i^k d\tau \\
 &= - \int_t^{t'} \delta_i^T \mathbf{Q}_{ii} \delta_i d\tau - \int_t^{t'} R_i(\mathbf{u}_i^k) d\tau - \int_t^{t'} \sum_{j \in N_i} R_i(\mathbf{u}_j^k) d\tau \\
 &\quad + \int_t^{t'} 2\lambda_i (\Delta_i^{k+1})^T \mathbf{R}_{ii}(\mathbf{u}_i - \mathbf{u}_i^k) d\tau \\
 &\quad - \int_t^{t'} \sum_{j \in N_i} e_{ij} \frac{2\lambda_i}{d_i + g_i} (\Delta_i^{k+1})^T \mathbf{R}_{ii}(\mathbf{u}_j - \mathbf{u}_j^k) d\tau.
 \end{aligned} \tag{37}$$

Therefore, it can be seen that the system matrices are not included in Eq. (37). Based on the Weierstrass high-order approximation theorem (Abu-Khalaf and Lewis, 2005), the critic and actor NNs can be introduced as

$$\hat{V}_i^k = (\phi_i(\delta_i))^T \mathbf{W}_{vi}^k, \tag{38}$$

$$\hat{\Delta}_{il_1}^k = (\varphi_{uil_1}(\delta_i))^T \mathbf{W}_{uil_1}^k, \tag{39}$$

where  $\phi_i \in \mathbb{R}^{h_v}$  and  $\varphi_{uil_1} \in \mathbb{R}^{h_{u_{l_1}}}$  ( $l_1 = 1, 2, \dots, m$ ) are activation functions.  $\mathbf{W}_{vi}^k$  and  $\mathbf{W}_{uil_1}^k$  are constant weights. Eq. (39) can be written in the following compact form:

$$\begin{aligned} \hat{\Delta}_i^k &= \begin{bmatrix} \hat{\Delta}_{i1}^k & \hat{\Delta}_{i2}^k & \dots & \hat{\Delta}_{im}^k \end{bmatrix}^T \\ &= \begin{bmatrix} (\varphi_{ui1}(\delta_i))^T \mathbf{W}_{ui1}^k & (\varphi_{ui2}(\delta_i))^T \mathbf{W}_{ui2}^k & \dots & (\varphi_{uim}(\delta_i))^T \mathbf{W}_{uim}^k \end{bmatrix}^T. \end{aligned} \tag{40}$$

Substituting Eqs. (38) and (39) into Eq. (37), we can obtain Eq. (41) (on the top of the next page), where  $\hat{\mathbf{u}}_i^k = -\lambda_i \Psi(\hat{\Delta}_i^k)$ ,  $\sigma_i^k$  represents the residual error, and  $\delta'_i = \delta_i(t')$ . Then, Eq. (41) can be written in a simplified form as follows:

$$\sigma_i^k = \rho_{i,[t,t']}^k \mathbf{W}_i^{k+1} - \pi_{i,[t,t']}^k, \tag{42}$$

where

$$\begin{aligned} \pi_{i,[t,t']}^k &= \int_t^{t'} \delta_i^T \mathbf{Q}_{ii} \delta_i d\tau + \int_t^{t'} R_i(\hat{\mathbf{u}}_i^k) d\tau \\ &+ \int_t^{t'} \sum_{j \in N_i} R_i(\hat{\mathbf{u}}_j^k) d\tau, \end{aligned}$$

$$\begin{aligned} \mathbf{W}_i^{k+1} &= \begin{bmatrix} (\mathbf{W}_{vi}^k)^T & (\mathbf{W}_{ui1}^{k+1})^T & \dots & (\mathbf{W}_{uim}^{k+1})^T \end{bmatrix}^T, \\ \rho_{i,[t,t']}^k &= \begin{bmatrix} \rho_{vi}^k & \rho_{ui1}^k & \dots & \rho_{uim}^k \end{bmatrix}, \end{aligned}$$

with

$$\begin{aligned} \rho_{vi}^k &= (\phi_i(\delta_i) - \phi_i(\delta'_i))^T, \\ \rho_{uil_2}^k &= 2\lambda_i \sum_{l_1=1}^m r_{ii,l_1 l_2} \\ &\cdot \int_t^{t'} \left( u_{il_1} + \lambda_i \Psi(\hat{\Delta}_{il_1}^k) \right) (\varphi_{uil_2}(\delta_i))^T d\tau \\ &- \frac{2\lambda_i}{d_i + g_i} \sum_{l_1=1}^m r_{ii,l_1 l_2} \\ &\cdot \int_t^{t'} \sum_{j \in N_i} e_{ij} \left( u_{jl_1} + \lambda_j \Psi(\hat{\Delta}_{jl_1}^k) \right) (\varphi_{uil_2}(\delta_i))^T d\tau, \end{aligned}$$

for  $l_2 = 1, 2, \dots, m$ .

We would like to determine matrix  $\mathbf{W}_i^{k+1}$  that can minimize the residual error  $\sigma_i^k$ . For a positive integer  $v$ , we can define  $\Gamma_i^k$  and  $\Pi_i^k$  as

$$\begin{aligned} \Gamma_i^k &= \left[ \left( \rho_{i,[t_0,t_1]}^k \right)^T \left( \rho_{i,[t_1,t_2]}^k \right)^T \dots \left( \rho_{i,[t_{v-1},t_v]}^k \right)^T \right]^T, \\ \Pi_i^k &= \begin{bmatrix} \pi_{i,[t_0,t_1]}^k & \pi_{i,[t_1,t_2]}^k & \dots & \pi_{i,[t_{v-1},t_v]}^k \end{bmatrix}^T. \end{aligned}$$

Based on the least square approach, it can be obtained that

$$\mathbf{W}_i^{k+1} = \left( (\Gamma_i^k)^T \Gamma_i^k \right)^{-1} (\Gamma_i^k)^T \Pi_i^k. \tag{43}$$

A model-free off-policy reinforcement learning algorithm can be provided in Algorithm 2.

---

**Algorithm 2** Model-free off-policy reinforcement learning

---

**Initialization:** Choose the initial stabilizing network weights  $\mathbf{W}_{uil_1}^0$  for  $i = 1, 2, \dots, N$ ,  $l_1 = 1, 2, \dots, m$ ; choose the computation precision  $\varepsilon$

**Iteration:**

- 1: Employ the control laws  $\mathbf{u}_i$  ( $i = 1, 2, \dots, N$ ) with exploration noises to system (34) on the time interval  $[t_0, t_v]$ . Collect the system data  $\{\delta_i, \mathbf{u}_i\}$  for  $i = 1, 2, \dots, N$
  - 2: Let iteration index  $k = 0$
  - 3: Compute  $\rho_{i,[t_0,t_1]}^k, \rho_{i,[t_1,t_2]}^k, \dots, \rho_{i,[t_{v-1},t_v]}^k$  and  $\pi_{i,[t_0,t_1]}^k, \pi_{i,[t_1,t_2]}^k, \dots, \pi_{i,[t_{v-1},t_v]}^k$
  - 4: Compute  $\mathbf{W}_i^{k+1}$  according to Eq. (43)
  - 5: Let iteration index  $k = k + 1$ ; repeat steps 2 and 3 until  $\|\mathbf{W}_i^{k+1} - \mathbf{W}_i^k\| \leq \varepsilon$  for  $k \geq 1$
  - 6: **return**  $\mathbf{W}_i^k, i = 1, 2, \dots, N$
- 

**Lemma 1** For system (4), if the iterative value functions and iterative control laws are designed as Eqs. (38) and (39), respectively, where the weights  $\mathbf{W}_{vi}^k$  and  $\mathbf{W}_{uil_1}^k$  ( $l_1 = 1, 2, \dots, m$ ) are updated as Algorithm 2,  $\lim_{k \rightarrow \infty} \hat{\mathbf{u}}_i^k = \mathbf{u}_i^*$  ( $i = 1, 2, \dots, N$ ).

The proof process can be found in Li JN et al. (2017) and Qin et al. (2019), and thus is omitted here.

**Remark 3** In Algorithm 2, the selection of the control laws  $\mathbf{u}_i$  ( $i = 1, 2, \dots, N$ ) is the key to the convergence of the algorithm. Generally, the control laws are selected as  $\mathbf{u}_i = -\mathbf{K}_i \delta_i + \xi_i$  ( $i = 1, 2, \dots, N$ ), where  $\xi_i$  is the exploration noise and  $\mathbf{K}_i$  is the stabilizing gain matrix.

**Remark 4** For traditional on-policy integral reinforcement learning methods (Vrabie and Lewis, 2011;

$$\begin{aligned} \sigma_i^k = & (\phi_i(\delta_i) - \phi_i(\delta'_i))^T \mathbf{W}_{vi}^k - \int_t^{t'} \delta_i^T \mathbf{Q}_{ii} \delta_i d\tau - \int_t^{t'} R_i(\hat{\mathbf{u}}_i^k) d\tau - \int_t^{t'} \sum_{j \in N_i} R_i(\hat{\mathbf{u}}_j^k) d\tau \\ & + 2\lambda_i \sum_{l_1=1}^m \sum_{l_2=1}^m r_{ii,l_1 l_2} \int_t^{t'} (u_{il_1} + \lambda_i \Psi(\hat{\Delta}_{il_1}^k)) (\varphi_{uil_2}(\delta_i))^T \mathbf{W}_{uil_2}^{k+1} d\tau \\ & - \frac{2\lambda_i}{d_i + g_i} \sum_{l_1=1}^m \sum_{l_2=1}^m r_{ii,l_1 l_2} \int_t^{t'} \sum_{j \in N_i} e_{ij} (u_{jl_1} + \lambda_j \Psi(\hat{\Delta}_{jl_1}^k)) (\varphi_{uil_2}(\delta_i))^T \mathbf{W}_{uil_2}^{k+1} d\tau. \end{aligned} \quad (41)$$

Liu et al., 2021), the performance index function is evaluated using the inaccurate data, which cause a biased estimation. The presented off-policy reinforcement learning method can avoid this problem and thus obtain results with higher accuracy.

### 4 Numerical analysis

In this section, a simulation example is provided to show the good performance of the presented method. The structure of the multi-agent systems is shown in Fig. 1, with the following dynamics:

$$\mathbf{A} = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}, \mathbf{B} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}.$$

The Laplacian matrix and the pinning matrix are given as

$$\mathbf{L} = \begin{bmatrix} 1 & -1 & 0 \\ 0 & 1 & -1 \\ -1 & -1 & 2 \end{bmatrix}, \mathbf{G} = \begin{bmatrix} 1 & & \\ & 0 & \\ & & 0 \end{bmatrix}.$$

We define the weight matrices of the performance index function, represented by Eq. (5), as follows:

$$\mathbf{Q}_{11} = \mathbf{I}_2, \mathbf{Q}_{22} = 2\mathbf{I}_2, \mathbf{Q}_{33} = 1.5\mathbf{I}_2,$$

$$R_{11} = 2, R_{22} = 1, R_{33} = 3,$$

$$R_{12} = 0.02, R_{23} = 0.01, R_{31} = 0.03, R_{32} = 0.01.$$

The simulation is performed with  $\mathbf{x}_0(0) = [1 \ 1]^T$ ,  $\mathbf{x}_1(0) = [0.5 \ -0.5]^T$ ,  $\mathbf{x}_2(0) = [1 \ -0.5]^T$ ,  $\mathbf{x}_3(0) = [2 \ -1]^T$ ,  $\lambda_1 = 2$ ,  $\lambda_2 = 1.5$ ,  $\lambda_3 = 3$ . First, we collect the system data  $\{\delta_i, \mathbf{u}_i\}$  every 0.01 s for  $i = 1, 2, 3$ . Then, we solve Eq. (43) iteratively based on the collected system data. The activation functions  $\phi_i(\delta_i)$  and  $\varphi_{ui}(\delta_i)$  are chosen as

$$\begin{aligned} \phi_i(\delta_i) &= [\delta_{i1}^2 \ \delta_{i1}\delta_{i2} \ \delta_{i2}^2]^T, \\ \varphi_{ui}(\delta_i) &= [\delta_{i1} \ \delta_{i2} \ \delta_{i1}^2 \ \delta_{i1}\delta_{i2} \ \delta_{i2}^2]^T. \end{aligned}$$

The simulation results are shown in Figs. 2–6. The weights of critic and actor NNs are shown in Figs. 2 and 3, respectively, which show the stability of Algorithm 2. The synchronization error curves are provided in Fig. 4, and the three-dimensional curves are provided in Fig. 5. From Figs. 4 and 5, it can be derived that the optimal synchronization control is achieved. The control curves are shown in Fig. 6, verifying that the control constraint is satisfied.

### 5 Conclusions

The nonzero-sum game problem of multi-agent systems with input saturation has been studied

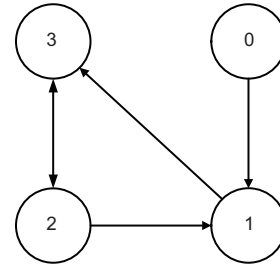


Fig. 1 Structure of the multi-agent systems

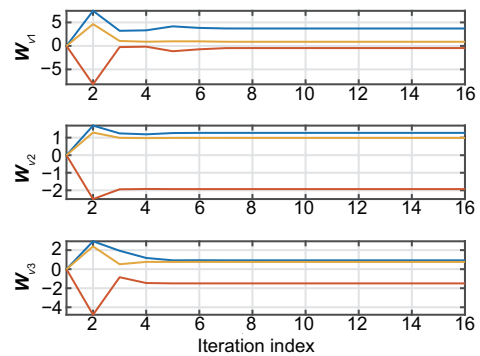


Fig. 2 Weights of critic neural networks of the multi-agent systems

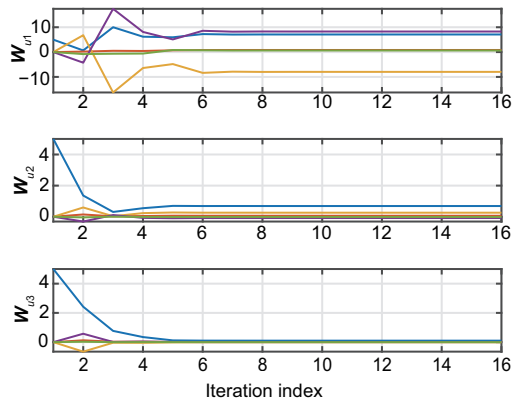


Fig. 3 Weights of actor neural networks of the multi-agent systems

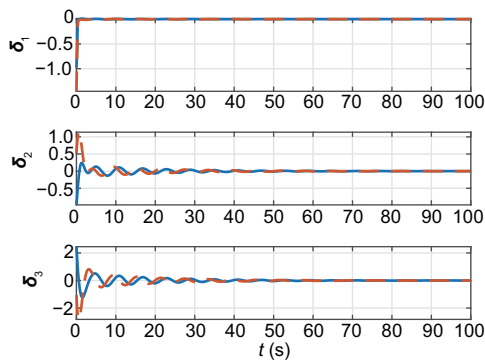


Fig. 4 Synchronization errors of the multi-agent systems

based on the model-free off-policy reinforcement learning method. It is shown that the presented off-policy reinforcement learning algorithm can make the iterative control laws converge to the Nash equilibrium without the information of system models. The simulation results showed the good performance of the presented method.

### Contributors

Hongyang LI designed the method, conducted the simulation, and drafted the paper. Qinglai WEI revised and finalized the paper.

### Compliance with ethics guidelines

Hongyang LI and Qinglai WEI declare that they have no conflict of interest.

### References

Abu-Khalaf M, Lewis FL, 2005. Nearly optimal control laws for nonlinear systems with saturating actuators using a neural network HJB approach. *Automatica*, 41(5):779-791. <https://doi.org/10.1016/j.automatica.2004.11.034>

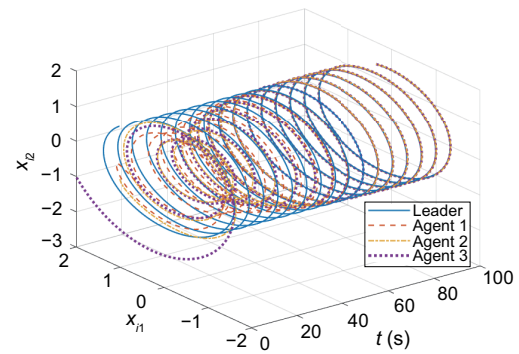


Fig. 5 Three-dimensional curves of the multi-agent systems

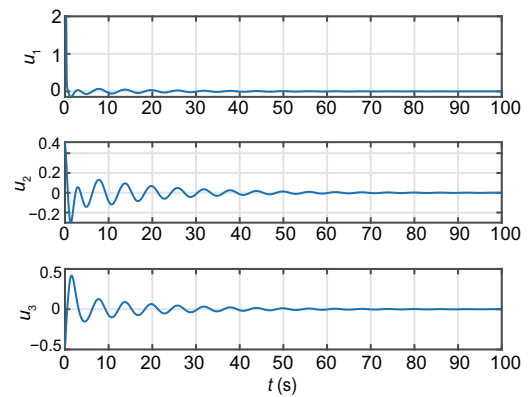


Fig. 6 Control laws of the multi-agent systems

- Bertsekas DP, 2007. *Dynamic Programming and Optimal Control*. Athena Scientific, Belmont, USA.
- Cao MT, Xiao F, Wang L, 2015. Event-based second-order consensus control for multi-agent systems via synchronous periodic event detection. *IEEE Trans Autom Contr*, 60(9):2452-2457. <https://doi.org/10.1109/TAC.2015.2390553>
- Du HB, He YG, Cheng YY, 2014. Finite-time synchronization of a class of second-order nonlinear multi-agent systems using output feedback control. *IEEE Trans Circ Syst I*, 61(6):1778-1788. <https://doi.org/10.1109/TCSI.2013.2295012>
- Garcia E, Cao YC, Casbeer D, 2017. Periodic event-triggered synchronization of linear multi-agent systems with communication delays. *IEEE Trans Autom Contr*, 62(1):366-371. <https://doi.org/10.1109/TAC.2016.2555484>
- Han YJ, Lu WL, Chen TP, 2013. Cluster consensus in discrete-time networks of multi-agents with inter-cluster nonidentical inputs. *IEEE Trans Neur Netw Learn Syst*, 24(4):566-578. <https://doi.org/10.1109/TNNLS.2013.2237786>
- He WL, Gao XY, Zhong WM, et al., 2018. Secure impulsive synchronization control of multi-agent systems under deception attacks. *Inform Sci*, 459:354-368. <https://doi.org/10.1016/j.ins.2018.04.020>

- Jiao Q, Modares H, Xu SY, et al., 2016. Multi-agent zero-sum differential graphical games for disturbance rejection in distributed control. *Automatica*, 69:24-34. <https://doi.org/10.1016/j.automatica.2016.02.002>
- Li JN, Modares H, Chai TY, et al., 2017. Off-policy reinforcement learning for synchronization in multiagent graphical games. *IEEE Trans Neur Netw Learn Syst*, 28(10):2434-2445. <https://doi.org/10.1109/TNNLS.2016.2609500>
- Li JQ, Wang QL, Su YX, et al., 2021. Robust distributed model predictive consensus of discrete-time multi-agent systems: a self-triggered approach. *Front Inform Technol Electron Eng*, 22(8):1068-1079. <https://doi.org/10.1631/FITEE.2000182>
- Liu DR, Xue S, Zhao B, et al., 2021. Adaptive dynamic programming for control: a survey and recent advances. *IEEE Trans Syst Man Cybern Syst*, 51(1):142-160. <https://doi.org/10.1109/TSMC.2020.3042876>
- Ma HJ, Yang GH, 2016. Adaptive fault tolerant control of cooperative heterogeneous systems with actuator faults and unreliable interconnections. *IEEE Trans Autom Contr*, 61(11):3240-3255. <https://doi.org/10.1109/TAC.2015.2507864>
- Qin JH, Li M, Shi Y, et al., 2019. Optimal synchronization control of multiagent systems with input saturation via off-policy reinforcement learning. *IEEE Trans Neur Netw Learn Syst*, 30(1):85-96. <https://doi.org/10.1109/TNNLS.2018.2832025>
- Rehák B, Lynnyk V, 2021. Leader-following synchronization of a multi-agent system with heterogeneous delays. *Front Inform Technol Electron Eng*, 22(1):97-106. <https://doi.org/10.1631/FITEE.2000207>
- Thunberg J, Song W, Monitijano E, et al., 2014. Distributed attitude synchronization control of multi-agent systems with switching topologies. *Automatica*, 50(3):832-840. <https://doi.org/10.1016/j.automatica.2014.02.002>
- Vamvoudakis KG, Lewis FL, Hudas GR, 2012. Multi-agent differential graphical games: online adaptive learning solution for synchronization with optimality. *Automatica*, 48(8):1598-1611. <https://doi.org/10.1016/j.automatica.2012.05.074>
- Vrabie D, Lewis F, 2011. Adaptive dynamic programming for online solution of a zero-sum differential game. *J Contr Theory Appl*, 9(3):353-360. <https://doi.org/10.1007/s11768-011-0166-4>
- Wang FY, Zhang HG, Liu DR, 2009. Adaptive dynamic programming: an introduction. *IEEE Comput Intell Mag*, 4(2):39-47. <https://doi.org/10.1109/MCI.2009.932261>
- Wei QL, Liu DR, 2014. Adaptive dynamic programming for optimal tracking control of unknown nonlinear systems with application to coal gasification. *IEEE Trans Autom Sci Eng*, 11(4):1020-1036. <https://doi.org/10.1109/TASE.2013.2284545>
- Wei QL, Wang FY, Liu DR, et al., 2014. Finite-approximation-error-based discrete-time iterative adaptive dynamic programming. *IEEE Trans Cybern*, 44(12):2820-2833. <https://doi.org/10.1109/TCYB.2014.2354377>
- Wei QL, Liu DR, Lewis FL, 2015. Optimal distributed synchronization control for continuous-time heterogeneous multi-agent differential graphical games. *Inform Sci*, 317:96-113. <https://doi.org/10.1016/j.ins.2015.04.044>
- Wei QL, Liu DR, Lin HQ, 2016. Value iteration adaptive dynamic programming for optimal control of discrete-time nonlinear systems. *IEEE Trans Cybern*, 46(3):840-853. <https://doi.org/10.1109/TCYB.2015.2492242>
- Wei QL, Lewis FL, Sun QY, et al., 2017. Discrete-time deterministic Q-learning: a novel convergence analysis. *IEEE Trans Cybern*, 47(5):1224-1237. <https://doi.org/10.1109/TCYB.2016.2542923>
- Wei QL, Lewis FL, Liu DR, et al., 2018. Discrete-time local value iteration adaptive dynamic programming: convergence analysis. *IEEE Trans Syst Man Cybern Syst*, 48(6):875-891. <https://doi.org/10.1109/TSMC.2016.2623766>
- Wei QL, Li HY, Wang FY, 2020. Parallel control for continuous-time linear systems: a case study. *IEEE/CAA J Autom Sin*, 7(4):919-928. <https://doi.org/10.1109/JAS.2020.1003216>
- Wei QL, Wang X, Zhong XN, et al., 2021. Consensus control of leader-following multi-agent systems in directed topology with heterogeneous disturbances. *IEEE/CAA J Autom Sin*, 8(2):423-431. <https://doi.org/10.1109/JAS.2021.1003838>
- Wieland P, Sepulchre R, Allgöwer F, 2011. An internal model principle is necessary and sufficient for linear output synchronization. *Automatica*, 47(5):1068-1074. <https://doi.org/10.1016/j.automatica.2011.01.081>
- Yang JY, Xi F, Ma J, 2019. Model-based edge-event-triggered containment control under directed topologies. *IEEE Trans Cybern*, 49(7):2556-2567. <https://doi.org/10.1109/TCYB.2018.2828645>
- Yang N, Xiao JW, Xiao L, et al., 2019. Non-zero sum differential graphical game: cluster synchronisation for multi-agents with partially unknown dynamics. *Int J Contr*, 92(10):2408-2419. <https://doi.org/10.1080/00207179.2018.1441550>
- Zhang HG, Zhang JL, Yang GH, et al., 2015. Leader-based optimal coordination control for the consensus problem of multiagent differential games via fuzzy adaptive dynamic programming. *IEEE Trans Fuzzy Syst*, 23(1):152-163. <https://doi.org/10.1109/TFUZZ.2014.2310238>
- Zhang KQ, Yang ZR, Başar T, 2021. Decentralized multi-agent reinforcement learning with networked agents: recent advances. *Front Inform Technol Electron Eng*, 22(6):802-814. <https://doi.org/10.1631/FITEE.1900661>
- Zhang LD, Wang B, Liu ZX, et al., 2019. Motion planning of a quadrotor robot game using a simulation-based projected policy iteration method. *Front Inform Technol Electron Eng*, 20(4):525-537. <https://doi.org/10.1631/FITEE.1800571>
- Zhao DY, Zhu QM, Li N, et al., 2014. Synchronized control with neuro-agents for leader-follower based multiple robotic manipulators. *Neurocomputing*, 124:149-161. <https://doi.org/10.1016/j.neucom.2013.07.016>