

Hongyang LI, Qinglai WEI, 2022. Optimal synchronization control for multi-agent systems with input saturation: a nonzero-sum game. *Frontiers of Information Technology & Electronic Engineering*, 23(7):1010-1019.

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

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

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

Corresponding author: Qinglai WEI

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

 ORCID: <https://orcid.org/0000-0001-7002-9800>

□ Research problem

- Optimal synchronization control for multi-agent systems with input saturation

- Control systems: $\dot{\mathbf{x}}_i = \mathbf{A}\mathbf{x}_i + \mathbf{B}\mathbf{u}_i$

$$\mathbf{u}_i \in \mathcal{U}_i \subset \mathbb{R}^m$$

$$\mathcal{U}_i = \{ \mathbf{u}_i \mid \mathbf{u}_i \in \mathbb{R}^m, \|\mathbf{u}_i\|_\infty \leq \lambda_i \}$$

$$\dot{\mathbf{x}}_0 = \mathbf{A}\mathbf{x}_0$$

Leader

- State synchronization condition: $\delta_i \rightarrow 0$ $\delta_i = \sum_{j \in N_i} e_{ij} (\mathbf{x}_i - \mathbf{x}_j) + g_i (\mathbf{x}_i - \mathbf{x}_0)$

□ Research background

- The existing synchronization condition methods for multi-agent systems based on ADP (Vamvoudakis et al., 2012; Wei et al., 2015; Qin et al., 2019)
- **Contribution:** The control constraint and coupled terms in the performance index functions are considered which broaden the application scope of the presented method

□ Multi-agent nonzero-sum game

- Performance index function

$$J_i(\boldsymbol{\delta}_i(0), \mathbf{u}_i, \mathbf{u}_{-i}) = \int_0^\infty \left(\boldsymbol{\delta}_i^\top \mathbf{Q}_{ii} \boldsymbol{\delta}_i + R_i(\mathbf{u}_i) + \sum_{j \in N_i} R_j(\mathbf{u}_j) \right) dt \quad (1)$$

- Nash equilibrium condition

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

- Iterative value function

$$V_i(\boldsymbol{\delta}_i(t)) = \int_t^\infty \left(\boldsymbol{\delta}_i^\top \mathbf{Q}_{ii} \boldsymbol{\delta}_i + R_i(\mathbf{u}_i) + \sum_{j \in N_i} R_j(\mathbf{u}_j) \right) d\tau \quad (3)$$

V_i^* satisfies the HJB equation

$$\mathbf{u}_i^* = -\lambda_i \Psi \left(\frac{1}{2\lambda_i} (d_i + g_i) \mathbf{R}_{ii}^{-1} \mathbf{B}^\top \frac{\partial V_i^*}{\partial \boldsymbol{\delta}_i} \right) \quad (4)$$

Optimal control

$$\left(\frac{\partial V_i^*}{\partial \boldsymbol{\delta}_i} \right)^\top \left(\mathbf{A} \boldsymbol{\delta}_i - \lambda_i (d_i + g_i) \mathbf{B} \Psi(\boldsymbol{\Delta}_i^*) + \sum_{j \in N_i} \lambda_j e_{ij} \mathbf{B} \Psi(\boldsymbol{\Delta}_j^*) \right) + \boldsymbol{\delta}_i^\top \mathbf{Q}_{ii} \boldsymbol{\delta}_i \quad (5)$$

$$+ R_i(-\lambda_i \Psi(\boldsymbol{\Delta}_i^*)) + \sum_{j \in N_i} R_j(-\lambda_j \Psi(\boldsymbol{\Delta}_j^*)) = 0$$

Theorem 1 Assume that the optimal control law is \mathbf{u}_i^* , and that V_i is the positive definite smooth solution to the HJB equation. Then, the tracking error is asymptotically stable. $\mathbf{u}_i^*, i = 1, 2, \dots, N$, constitute the Nash equilibrium, and the solution V_i to the HJB equation is the optimal value of the game, i.e., $J_i^*(\boldsymbol{\delta}_i(0), \mathbf{u}_i^*, \mathbf{u}_{-i}^*) = V_i(\boldsymbol{\delta}_i(0))$.

□ Policy iteration method

- Solve V_i^k

$$\left(\frac{\partial V_i^k}{\partial \boldsymbol{\delta}_i} \right)^T \left(\mathbf{A} \boldsymbol{\delta}_i + (d_i + g_i) \mathbf{B} \mathbf{u}_i^k - \sum_{j \in N_i} e_{ij} \mathbf{B} \mathbf{u}_j^k \right) + \boldsymbol{\delta}_i^T \mathbf{Q}_{ii} \boldsymbol{\delta}_i + R_i(\mathbf{u}_i^k) + \sum_{j \in N_i} R_i(\mathbf{u}_j^k) = 0 \quad (6)$$

- Update the iterative control laws

$$\mathbf{u}_i^{k+1} = -\lambda_i \Psi \left(\frac{1}{2\lambda_i} (d_i + g_i) \mathbf{R}_{ii}^{-1} \mathbf{B}^T \frac{\partial V_i^k}{\partial \boldsymbol{\delta}_i} \right) \quad (7)$$

Theorem 2 Assume that agent i and its neighbors update their control policies according to Eqs. (6) and (7). Then, policies $\mathbf{u}_i^k, \mathbf{u}_{-i}^k$ converge to the Nash equilibrium, and V_i^k converges to V_i^* .

□ Model-free off-policy reinforcement learning method

- Rewrite the tracking error dynamics

$$\dot{\delta}_i = \mathbf{A}\delta_i + (d_i + g_i)\mathbf{B}(\mathbf{u}_i - \mathbf{u}_i^k) + (d_i + g_i)\mathbf{B}\mathbf{u}_i^k - \sum_{j \in N_i} e_{ij}\mathbf{B}\mathbf{u}_j^k - \sum_{j \in N_i} e_{ij}\mathbf{B}(\mathbf{u}_j - \mathbf{u}_j^k) \quad (8)$$

$$\dot{V}_i^k = \left(\frac{\partial V_i^k}{\partial \delta_i} \right)^T \left(\mathbf{A}\delta_i - \sum_{j \in N_i} e_{ij}\mathbf{B}\mathbf{u}_j^k + (d_i + g_i)\mathbf{B}\mathbf{u}_i^k \right) \quad (9)$$

$$+ \left(\frac{\partial V_i^k}{\partial \delta_i} \right)^T (d_i + g_i)\mathbf{B}(\mathbf{u}_i - \mathbf{u}_i^k) - \left(\frac{\partial V_i^k}{\partial \delta_i} \right)^T \sum_{j \in N_i} e_{ij}\mathbf{B}(\mathbf{u}_j - \mathbf{u}_j^k) \quad (10)$$

$$V_i^k(\delta_i(t')) - V_i^k(\delta_i(t)) = -\int_t^{t'} \delta_i^T \mathbf{Q}_i \delta_i d\tau - \int_t^{t'} R_i(\mathbf{u}_i^k) d\tau - \int_t^{t'} \sum_{j \in N_i} R_i(\mathbf{u}_j^k) d\tau$$

$$+ \int_t^{t'} 2\lambda_i (\mathbf{A}_i^{k+1})^T \mathbf{R}_i(\mathbf{u}_i - \mathbf{u}_i^k) d\tau - \int_t^{t'} \sum_{j \in N_i} e_{ij} \frac{2\lambda_i}{d_i + g_i} (\mathbf{A}_i^{k+1})^T \mathbf{R}_i(\mathbf{u}_j - \mathbf{u}_j^k) d\tau$$

- The critic and actor NNs can be introduced

$$\hat{V}_i^k = \left(\phi_i(\delta_i) \right)^T \mathbf{W}_{vi}^k \quad (11)$$

$$\hat{\Delta}_{il_1}^k = \left(\varphi_{uil_1}(\delta_i) \right)^T \mathbf{W}_{uil_1}^k \quad (12)$$

- It can be derived that

$$\begin{aligned} \sigma_i^k = & \left(\phi_i(\delta_i) - \phi_i(\delta'_i) \right)^T \mathbf{W}_{vi}^k - \int_t^{t'} \delta_i^T \mathbf{Q}_{ii} \delta_i d\tau - \int_t^{t'} R_i(\hat{u}_i^k) d\tau - \int_t^{t'} \sum_{j \in N_i} R_j(\hat{u}_j^k) d\tau \\ & + 2\lambda_i \sum_{l_1=1}^m \sum_{l_2=1}^m r_{il_1l_2} \int_t^{t'} \left(u_{il_1} + \lambda_i \Psi(\hat{\Delta}_{il_1}^k) \right) \left(\varphi_{uil_2}(\delta_i) \right)^T \mathbf{W}_{uil_2}^{k+1} d\tau \end{aligned} \quad (13)$$

Simplified form

$$\begin{aligned} - \frac{2\lambda_i}{d_i + g_i} \sum_{l_1=1}^m \sum_{l_2=1}^m r_{il_1l_2} \int_t^{t'} \sum_{j \in N_i} e_{ij} \left(u_{jl_1} + \lambda_j \Psi(\hat{\Delta}_{jl_1}^k) \right) \left(\varphi_{uil_2}(\delta_i) \right)^T \mathbf{W}_{uil_2}^{k+1} d\tau \\ \sigma_i^k = \boldsymbol{\rho}_{i,[t,t']}^k \mathbf{W}_i^{k+1} - \boldsymbol{\pi}_{i,[t,t']}^k \end{aligned} \quad (14)$$

$$\mathbf{W}_i^{k+1} = \left(\left(\boldsymbol{\Gamma}_i^k \right)^T \boldsymbol{\Gamma}_i^k \right)^{-1} \left(\boldsymbol{\Gamma}_i^k \right)^T \boldsymbol{\Pi}_i^k \quad (15)$$

□ Algorithm 1 Model-free off-policy reinforcement learning

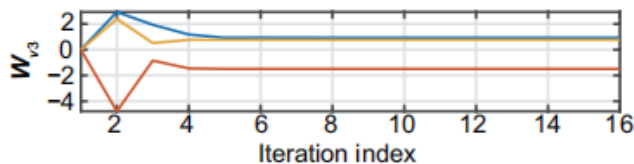
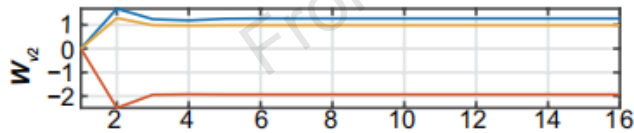
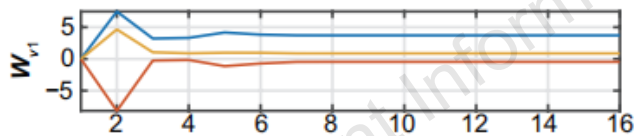
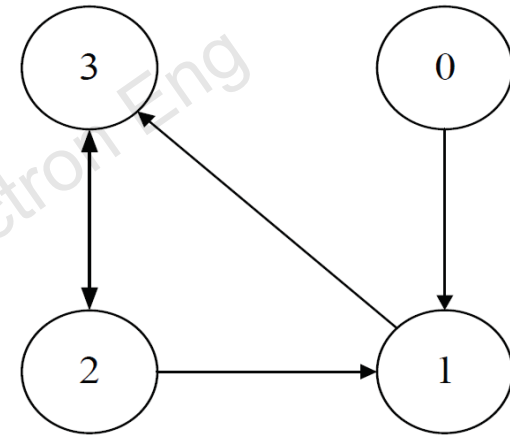
1. Choose the initial stabilizing network weights $\mathbf{W}_{uil_1}^0, i = 1, 2, \dots, N, l_1 = 1, 2, \dots, m$
2. Employ the control law $\mathbf{u}_i, i = 1, 2, \dots, N$, to the system on the time interval $[t_0, t_v]$. Collect the system data $\{\delta_i, \mathbf{u}_i\}, i = 1, 2, \dots, N$. Let $k = 0$
3. Compute $\rho_{i,[t_0,t_1]}^k, \dots, \rho_{i,[t_{v-1},t_v]}^k$ and $\pi_{i,[t_0,t_1]}^k, \dots, \pi_{i,[t_{v-1},t_v]}^k$, and compute \mathbf{W}_i^{k+1}
4. Let iteration index $k = k + 1$; repeat step 3 until $\|\mathbf{W}_i^{k+1} - \mathbf{W}_i^k\| \leq \varepsilon$
5. Return $\mathbf{W}_i^k, i = 1, 2, \dots, N$

Lemma 1 If the iterative value functions and iterative control laws are designed as Eqs. (11) and (12), where \mathbf{W}_{vi}^k and $\mathbf{W}_{uil_1}^k, l_1 = 1, 2, \dots, m$ are updated as Algorithm 1, then $\lim_{k \rightarrow \infty} \hat{\mathbf{u}}_i^k = \mathbf{u}_i^*, i = 1, 2, \dots, N$.

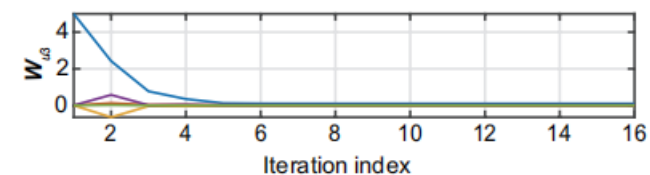
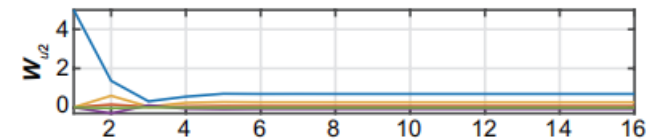
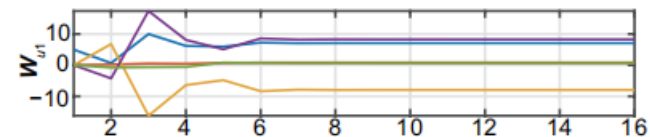
□ Numerical analysis

- System matrices $A = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}$, $B = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

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



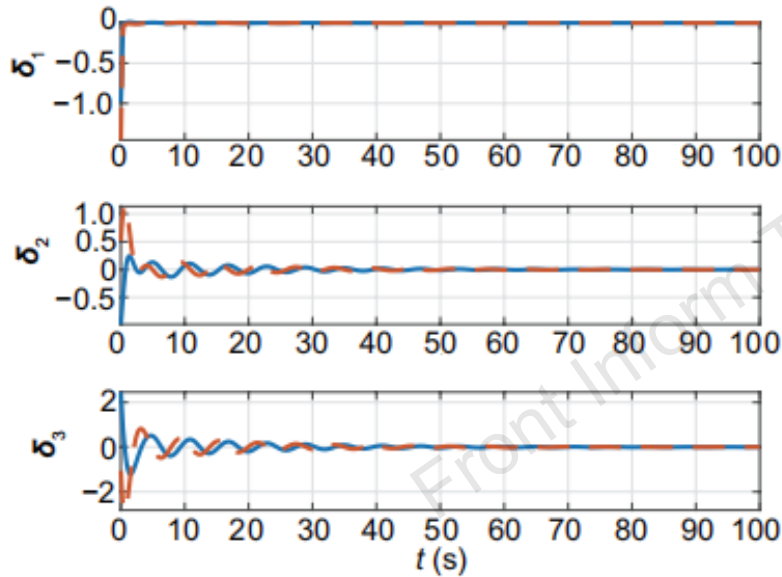
Critic NNs



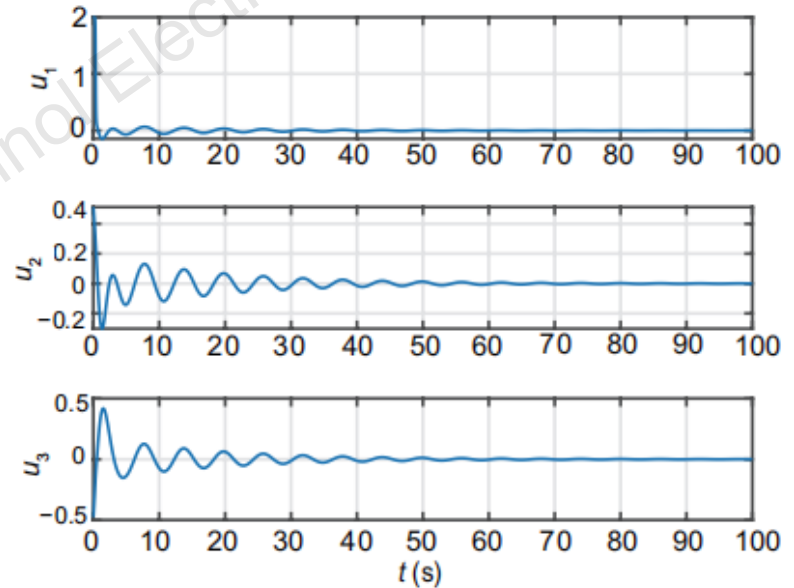
Actor NNs

□ Numerical analysis (Cont'd)

- $\lambda_1 = 2, \lambda_2 = 1.5, \lambda_3 = 3$



Synchronization errors



Control laws

□ Conclusions

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



Qinglai WEI received his BS degree in automation, MS degree in control theory and control engineering, and PhD degree in control theory and control engineering, from Northeastern University, Shenyang, China, in 2002, 2005, and 2008, respectively. He is an Associate Editor-in-Chief of *Neurocomputing*. He has been a vice president of the IEEE Computational Intelligence Society Beijing Chapter since 2021. His research interests include adaptive dynamic programming, neural-network-based control, optimal control, nonlinear systems, and their industrial applications.