

离散非线性零和博弈的事件驱动最优控制方案

张欣^{1†}, 薄迎春¹, 崔黎黎²

(1. 中国石油大学(华东) 信息与控制工程学院, 山东 青岛 266580;

2. 沈阳师范大学 科信软件学院, 辽宁 沈阳 110034)

摘要: 在求解离散非线性零和博弈问题时, 为了在有效降低网络通讯和控制器执行次数的同时保证良好的控制效果, 本文提出了一种基于事件驱动机制的最优控制方案. 首先, 设计了一个采用新型事件驱动阈值的事件驱动条件, 并根据贝尔曼最优性原理获得了最优控制对的表达式. 为了求解该表达式中的最优值函数, 提出了一种单网络值迭代算法. 利用一个神经网络构建评价网. 设计了新的评价网权重更新规则. 通过在评价网、控制策略及扰动策略之间不断迭代, 最终获得零和博弈问题的最优值函数和最优控制对. 然后, 利用Lyapunov稳定性理论证明了闭环系统的稳定性. 最后, 将该事件驱动最优控制方案应用到了两个仿真例子中, 验证了所提方法的有效性.

关键词: 博弈论; 事件驱动; 自适应动态规划; 最优控制

引用格式: 张欣, 薄迎春, 崔黎黎. 离散非线性零和博弈的事件驱动最优控制方案. 控制理论与应用, 2018, 35(5): 619 – 626

中图分类号: TP273 文献标识码: A

Event-triggered optimal control scheme for discrete-time nonlinear zero-sum games

ZHANG Xin^{1†}, BO Ying-chun¹, CUI Li-li²

(1. College of Information and Control Engineering, China University of Petroleum, Qingdao Shandong 266580, China;

2. Software College, Shenyang Normal University, Shenyang Liaoning 110034, China)

Abstract: In order to reduce the network communication and controller execution frequency while guarantee a desired control performance, an event-triggered optimal control scheme is proposed for solving the optimal control pair of discrete-time nonlinear zero-sum games in this paper. Firstly, an event-triggered condition with new event-triggered threshold is designed. The expression of the optimal control pair is obtained based on the Bellman optimality principle. Then, a single network value iteration algorithm is proposed to solve the optimal value function in this expression. A neural network is used to construct the critic network. Novel weight update rule of the critic network is derived. Through the iteration between the critic network, the control policy and the disturbance policy, the optimal value function and the optimal control pair can be solved. Further, the Lyapunov theory is used to prove the stability of the event-triggered closed-loop system. Finally, the event-triggered optimal control mechanism is applied to two examples to verify its effectiveness.

Key words: game theory; event-triggered; adaptive dynamic programming; optimal control

Citation: ZHANG Xin, BO Yingchun, CUI Lili. Event-triggered optimal control scheme for discrete-time nonlinear zero-sum games. *Control Theory & Applications*, 2018, 35(5): 619 – 626

1 引言(Introduction)

近年来, 零和博弈问题在博弈论领域和最优控制领域获得了广泛关注^[1-3]. 这是由于零和博弈具有两个决策者, 一方面要求控制输入使性能指标取极小, 而在干扰影响较大时, 又必须考虑干扰信号使性能指标取极大. 这样的对抗性设计既能保证系统在取最优

性的同时又具有较好的抗干扰能力. 然而现有的求解零和博弈问题的方法大都采用时间驱动机制, 即控制器是连续更新的, 在每一个采样时刻系统状态与控制器之间都要进行数据通讯, 控制输入都需要计算并执行. 这就大大增加了通讯网络和执行器的负担.

与传统的采样方法不同, 事件驱动机制采用一种

收稿日期: 2017-11-01; 录用日期: 2018-01-23.

[†]通信作者. E-mail: zhangxin@upc.edu.cn; Tel.: +86 15564879644.

本文责任编辑: 魏鞅.

山东省自然科学基金项目(BS2015DX009), 国家自然科学基金项目(61703289)资助.

Supported by the National Natural Science Foundation of Shandong Province (BS2015DX009) and the National Natural Science Foundation of China (61703289).

非周期采样模式^[4-7]. 文献[4]证明了这种非周期采样比周期采样在计算方面更加有利. 事件驱动机制预先设定了一个事件驱动条件, 只有当该条件不被满足时, 才对系统状态进行采样, 更新系统的控制输入, 在两次更新之间采用零阶保持器保证控制器的输出. 因此, 能够有效地降低网络通讯和控制器执行次数, 同时还能保证系统具有良好的控制性能. 文献[5]研究了线性系统的事件驱动控制. 文献[6]设计了事件驱动光电跟踪系统. Shao等人在文献[7]中研究了连续非线性系统的事件驱动状态反馈控制方案. 文献[8]将事件驱动控制带入了到最优控制领域. 事件驱动控制在求解连续系统的零和博弈问题方面也有了相应的成果, 文献[9]将 H_∞ 问题转化为零和博弈问题, 然后基于事件驱动机制进行求解. 据笔者所知, 目前还没有文献利用事件驱动机制求解离散非线性系统的零和博弈问题.

离散非线性系统的零和博弈问题需要求解离散Hamilton-Jacobi-Isaacs (HJI)方程来获得Nash平衡点, 即最优控制对. 但是对于非线性系统来说, HJI方程的解析解很难获得. Werbos在文献[10]中提出了一种有效的求解最优控制问题的方法——自适应动态规划(adaptive dynamic programming, ADP)算法, 并且得到了广泛应用^[11-13]. 文献[11]利用ADP算法处理鲁棒近似最优跟踪问题. 王鼎等人在文献[12]中综述了连续时间非线性系统的自适应评判鲁棒控制设计的最新研究成果. 文献[13]研究了离散非线性系统的事件驱动控制问题. ADP算法自其诞生之日起产生了一系列的同义词, 例如: 自适应评价设计、启发式动态规划、近似动态规划、神经元动态规划和增强学习等等. 2006年在美国科学基金会组织的“2006 NSF Workshop and Outreach Tutorials on Approximate Dynamic Programming”研讨会上, 建议将该方法统称为“adaptive/approximate dynamic programming (自适应/近似动态规划)”. ADP算法已经在一些文献中被用来处理零和博弈问题, 并取得了一定的理论研究成果^[4-17]. 然而这些研究都是基于时间驱动机制进行的.

本文将事件驱动机制、ADP算法和神经网络各自优势相结合, 提出了一种求解离散非线性零和博弈问题的事件驱动单网络值迭代控制方案. 首先设计了一个新型的事件驱动阈值. 根据贝尔曼最优性原理获得了最优控制对表达式. 然而, 由于HJI固有的非线性其解析解难以获得, 导致该最优控制对无法直接求解. 因此, 一种单网络值迭代算法被提出. 只利用一个神经网络构建评价网, 从而代替了典型ADP算法中的评价——控制双网结构, 有效减少了神经网络的训练次数. 然后, 根据HJI方程和梯度下降法设计了评价网的权值更新规则. 接着, 利用Lyapunov稳定性理论证明了闭环系统的稳定性. 最后, 将事件驱动最优控制方案应用到了两个仿真例子中, 验证了所提方案既能够

有效地降低网络通讯和控制器执行次数, 减少神经网络的训练次数, 又能够保证具有良好的性能.

2 问题描述(Problem descriptions)

考虑如下离散非线性系统的零和博弈问题, 其状态方程描述为

$$\mathbf{x}_{k+1} = f(\mathbf{x}_k) + g(\mathbf{x}_k)\mathbf{u}_k + h(\mathbf{x}_k)\mathbf{w}_k, \quad (1)$$

相应的性能指标函数为普通二次型形式

$$J(\mathbf{x}_0) = \sum_{k=0}^{\infty} U(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k), \quad (2)$$

其中: $\mathbf{x}_k \in \Omega \subseteq \mathbb{R}^n$ 为状态向量; $\mathbf{u}_k \in \mathbb{R}^{m_1}$ 为控制输入, 控制目标是使得性能指标函数最小, 而扰动输入 $\mathbf{w}_k \in \mathbb{R}^{m_2}$ 则希望使得性能指标函数最大; $f(\cdot)$, $g(\cdot)$ 和 $h(\cdot)$ 为光滑可微函数; \mathbf{x}_0 为系统初始状态; $U(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k) = \mathbf{x}_k^T Q \mathbf{x}_k + \mathbf{u}_k^T R \mathbf{u}_k - \mathbf{w}_k^T S \mathbf{w}_k$ 是对应的效用函数, 矩阵 Q , R 和 S 是具有适当维数的对称正定矩阵.

假设 1 系统(1)是可控的, 即存在连续控制策略能够渐近镇定系统(1), $f(0) = 0$, $\mathbf{x}_k = 0$ 是系统(1)唯一的平衡点^[17].

假设 2 $f + gu + hw$ 在紧集 $\Omega \subseteq \mathbb{R}^n$ 上李普希兹连续^[17].

定义 1 容许控制是指控制输入 \mathbf{u}_k 在紧集 $\Omega \subseteq \mathbb{R}^{m_1}$ 上连续且 $\mathbf{u}(0) = 0$, 能够控制系统(1)稳定并且保证性能指标函数(2)有界, $\forall \mathbf{x}_0 \in \Omega$ ^[17].

由容许控制 \mathbf{u}_k 和扰动输入 \mathbf{w}_k 定义值函数

$$V(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k) = \sum_{n=k}^{\infty} \{ \mathbf{x}_n^T Q \mathbf{x}_n + \mathbf{u}_n^T R \mathbf{u}_n - \mathbf{w}_n^T S \mathbf{w}_n \}. \quad (3)$$

求解由式(1)–(2)描述的离散非线性系统的零和博弈问题的最优控制对, 要求最优值函数满足

$$V^*(\mathbf{x}_k) = \min_{\mathbf{u}_k} \max_{\mathbf{w}_k} V(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k) = \max_{\mathbf{w}_k} \min_{\mathbf{u}_k} V(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k). \quad (4)$$

根据Bellman最优性原理, 最优值函数 $V^*(\mathbf{x}_k)$ 满足离散HJI方程^[16]

$$V^*(\mathbf{x}_k) = U(\mathbf{x}_k, \mathbf{u}_k^*, \mathbf{w}_k^*) + V^*(\mathbf{x}_{k+1}), \quad (5)$$

其中最优化控制对 $(\mathbf{u}_k^*, \mathbf{w}_k^*)$ 应该满足

$$\begin{cases} \frac{\partial H(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k)}{\partial \mathbf{u}_k} = 0, \\ \frac{\partial H(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k)}{\partial \mathbf{w}_k} = 0, \end{cases} \quad (6)$$

$H(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k)$ 为汉密尔顿函数

$$H(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k) = U(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k) + \lambda_{k+1}(f(\mathbf{x}_k) + g(\mathbf{x}_k)\mathbf{u}_k + h(\mathbf{x}_k)\mathbf{w}_k), \quad (7)$$

其中协状态 $\lambda_{k+1} = \frac{\partial V^*(\mathbf{x}_{k+1})}{\partial \mathbf{x}_{k+1}}$. 因此,

$$\mathbf{u}^*(\mathbf{x}_k) = -\frac{1}{2}R^{-1}g^T(\mathbf{x}_k)\lambda_{k+1}, \quad (8a)$$

$$\mathbf{w}^*(\mathbf{x}_k) = \frac{1}{2}S^{-1}h^T(\mathbf{x}_k)\lambda_{k+1}. \quad (8b)$$

3 事件驱动最优控制方案 (Event-triggered optimal control mechanism)

3.1 事件驱动条件(Event-triggered condition)

在事件驱动机制中, 定义 $\{k_i\}_0^\infty$ 是一个单调递增序列, k_i 代表第 i 个采样时刻, $i = 0, 1, 2, 3, \dots$. 这个采样系统的输出是由系统(1)在 k_i 时刻的状态 \mathbf{x}_{k_i} 组成的序列. 定义事件驱动误差为

$$\mathbf{e}_k = \mathbf{x}_{k_i} - \mathbf{x}_k, \quad k \in [k_i, k_{i+1}). \quad (9)$$

事件驱动条件为

$$\|\mathbf{e}_k\| \leq e_T, \quad (10)$$

其中 e_T 为事件驱动阈值. 只有当 $\|\mathbf{e}_k\| > e_T$ 时, 驱动条件不再满足, 系统进行采样. 事件驱动误差重置为零, $\mathbf{e}_{k_i} = 0$. 反馈控制输入 $\mathbf{u}(\mathbf{x}_{k_i}) = \boldsymbol{\mu}(\mathbf{x}_{k_i})$ 更新, 并且通过零阶保持器, 该控制输入在 $k \in [k_i, k_{i+1})$ 时间段内保持不变 $\mathbf{u}(\mathbf{x}_k) = \boldsymbol{\mu}(\mathbf{x}_{k_i})$, 直到下一个采样时刻. 需要注意的是, 在本文中假设事件驱动只对控制器 \mathbf{u}_k 有影响, 而对扰动输入 \mathbf{w}_k 没有影响. 根据式(9), 可得

$$\mathbf{u}(\mathbf{x}_k) = \boldsymbol{\mu}(\mathbf{e}_k + \mathbf{x}_k). \quad (11)$$

因此, 系统状态方程(1)重写为

$$\begin{aligned} \mathbf{x}_{k+1} &= f(\mathbf{x}_k) + g(\mathbf{x}_k)\boldsymbol{\mu}(\mathbf{e}_k + \mathbf{x}_k) + \\ &h(\mathbf{x}_k)\mathbf{w}_k. \end{aligned} \quad (12)$$

在事件驱动机制中, 控制输入只在采样时刻更新, 即只在 k_i 时刻生成. 因此, 状态反馈控制策略(8a)应该表示为

$$\boldsymbol{\mu}^*(\mathbf{x}_k) = -\frac{1}{2}R^{-1}g^T(\mathbf{x}_{k_{i+1}})\lambda_{k_{i+1}}, \quad k \in [k_i, k_{i+1}). \quad (13)$$

假设 3 存在正数 L , 满足^[13]

$$\|\mathbf{x}_{k+1}\| \leq L\|\mathbf{e}_k\| + L\|\mathbf{x}_k\|. \quad (14)$$

当最后一次采样时刻为 k_i , $k \in [k_i, k_{i+1})$, 根据式(9), 可得 $\mathbf{e}_{k+1} = \mathbf{x}_{k_i} - \mathbf{x}_{k+1}$. 显然

$$\begin{aligned} \|\mathbf{e}_{k+1}\| &\leq \|\mathbf{x}_{k_i}\| + \|\mathbf{x}_{k+1}\| \leq \\ &\|\mathbf{x}_{k_i}\| + L\|\mathbf{e}_k\| + L\|\mathbf{x}_k\| \leq \\ &\|\mathbf{x}_{k_i}\| + L\|\mathbf{e}_k\| + L(\|\mathbf{x}_{k_i}\| + \|\mathbf{e}_k\|) = \\ &2L\|\mathbf{e}_k\| + (1+L)\|\mathbf{x}_{k_i}\|. \end{aligned} \quad (15)$$

利用其递归性可得

$$\|\mathbf{e}_k\| \leq 2L\|\mathbf{e}_{k-1}\| + (1+L)\|\mathbf{x}_{k_i}\| \leq$$

$$\begin{aligned} &(2L)^2\|\mathbf{e}_{k-2}\| + [1+2L](1+L)\|\mathbf{x}_{k_i}\| \leq \\ &\vdots \\ &(2L)^{k-k_i}\|\mathbf{e}_{k_i}\| + [1+2L+(2L)^2+ \\ &(2L)^3+\dots+(2L)^{k-k_i-1}] \times \\ &(1+L)\|\mathbf{x}_{k_i}\|. \end{aligned} \quad (16)$$

为了确保等比数列收敛, 要求 $2L < 1$, 即 $L < 0.5$. 由于在每一个采样时刻 $\mathbf{e}_{k_i} = 0$, 则式(16)变为

$$\|\mathbf{e}_k\| \leq \frac{(1-(2L)^{k-k_i})(1+L)}{1-2L}\|\mathbf{x}_{k_i}\|. \quad (17)$$

定义事件驱动阈值为

$$e_T = \frac{\alpha(1-(2L)^{k-k_i})(1+L)}{1-2L}\|\mathbf{x}_{k_i}\|, \quad (18)$$

其中 $\alpha \in (0, 1]$ 为常数.

3.2 单网络ADP值迭代算法及神经网络实现(Single network ADP value iteration algorithm and neural network implementation)

对于非线性系统来说, HJI方程(5)的解很难直接求解. 为了获得式(8b)和式(13)中最优值函数的值, 根据贝尔曼最优性原理, 利用ADP值迭代算法来近似求解.

首先, 给定一个初始值函数 $V_0(\mathbf{x}_k)$, 一般情况选择 $V_0(\mathbf{x}_k) = 0$. \mathbf{u}_0 和 \mathbf{w}_0 可以通过下式计算获得:

$$\begin{cases} \mathbf{u}_0 = \arg \min_{\mathbf{u}_k} \{U(\mathbf{x}_k, \mathbf{u}_0, \mathbf{w}_0) + V_0(\mathbf{x}_{k+1})\}, \\ \mathbf{w}_0 = \arg \max_{\mathbf{w}_k} \{U(\mathbf{x}_k, \mathbf{u}_0, \mathbf{w}_0) + V_0(\mathbf{x}_{k+1})\}. \end{cases} \quad (19)$$

那么迭代的值函数 $V_1(\mathbf{x}_k)$ 为

$$V_1(\mathbf{x}_k) = U(\mathbf{x}_k, \mathbf{u}_0, \mathbf{w}_0) + V_0(\mathbf{x}_{k+1}). \quad (20)$$

以此类推, 相应的迭代策略 \mathbf{u}_j 和 \mathbf{w}_j 迭代规则为

$$\mathbf{u}_j = -\frac{1}{2}R^{-1}g^T(\mathbf{x}_k)\lambda_{j(k+1)}, \quad (21a)$$

$$\mathbf{w}_j = \frac{1}{2}S^{-1}h^T(\mathbf{x}_k)\lambda_{j(k+1)}. \quad (21b)$$

值函数 $V_{j+1}(\mathbf{x}_k)$ 的迭代规则为

$$\begin{aligned} V_{j+1}(\mathbf{x}_k) &= (\mathbf{x}_k^T Q \mathbf{x}_k + \mathbf{u}_j^T R \mathbf{u}_j - \mathbf{w}_j^T S \mathbf{w}_j) + \\ &V_j(\mathbf{x}_{k+1}), \end{aligned} \quad (22)$$

其中 j 表示迭代次数.

本文采用的是单神经网络结构, 只利用一个评价网来近似值函数. 该评价网由以下3层神经网络构成:

$$V^*(\mathbf{x}_k) = W_c^{*T} \phi_c(V_c^* \mathbf{x}_k) + \varepsilon_{ck}, \quad (23)$$

其中 $W_c^* \in \mathbb{R}^{N_c \times 1}$ 为未知的隐含层到输出层的理想神经网络权值, $V_c^* \in \mathbb{R}^{N_c \times n}$ 为输入层到隐含层的理想神经网络权值, N_c 是隐含层节点数, $\phi_c(\cdot)$ 为评价网激活函数, $\varepsilon_{ck} \in \mathbb{R}$ 为评价网近似误差.

在评价网训练过程中,输入层到隐含层的权值 V_c^* 保持不变.仅训练隐含层到输出层的权值,定义 $\hat{W}_c(j)$ 为其估计值,则实际的评价网输出为

$$V_j(\mathbf{x}_k) = \hat{W}_c^T(j)\phi_c(\bar{x}_k), \quad (24)$$

其中 $\bar{x}_k = V_c^* \mathbf{x}_k$.

根据值函数的迭代规则(22)和评价网输出(24)以及HJI方程(5),设计评价网的训练误差为

$$\begin{aligned} e_c = & U(\mathbf{x}_k, \mathbf{u}_j, \mathbf{w}_j) + V_j(\mathbf{x}_{k+1}) - V_j(\mathbf{x}_k) = \\ & U(\mathbf{x}_k, \mathbf{u}_j, \mathbf{w}_j) + \\ & \hat{W}_c^T(j)\phi_c(\bar{x}_{k+1}) - \hat{W}_c^T(j)\phi_c(\bar{x}_k) = \\ & \mathbf{x}_k^T Q \mathbf{x}_k + \mathbf{u}_j^T R \mathbf{u}_j - \mathbf{w}_j^T S \mathbf{w}_j + \\ & \hat{W}_c^T(j)\Delta\phi_c(k), \end{aligned} \quad (25)$$

其中: $\Delta\phi_c(k) = \phi_c(\bar{x}_{k+1}) - \phi_c(\bar{x}_k)$, \mathbf{u}_j 和 \mathbf{w}_j 的值由式(21)计算获得.

定义最小化目标函数为

$$E_c = \frac{1}{2}e_c^2. \quad (26)$$

利用梯度下降法,可得评价网的权值更新规则为

$$\begin{aligned} \hat{W}_c(j+1) &= \hat{W}_c(j) + \Delta\hat{W}_c(j), \\ \Delta\hat{W}_c(j) &= -\alpha_c \left(\frac{\partial E_c}{\partial \hat{W}_c(j)} \right) = \\ & -\alpha_c \Delta\phi_c(k) (\mathbf{x}_k^T Q \mathbf{x}_k + \mathbf{u}_j^T R \mathbf{u}_j - \\ & \mathbf{w}_j^T S \mathbf{w}_j + \hat{W}_c^T(j)\Delta\phi_c(k))^T, \end{aligned} \quad (27)$$

其中 α_c 为评价网学习率.

假设4 存在常数 θ, α, β 满足

$$0 \leq V^*(\mathbf{x}_{k+1}) \leq \theta U(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k), \quad (28)$$

$$0 \leq \eta_1 V^*(\mathbf{x}_k) \leq V_0(\mathbf{x}_k) \leq \eta_2 V^*(\mathbf{x}_k), \quad (29)$$

其中: $0 < \theta < \infty$, $0 < \eta_1 < 1$, $1 \leq \eta_2 < \infty$, V_0 为任意初始值函数^[17].

若假设4成立,当迭代次数 j 趋于无穷大时, $V_j(\mathbf{x}_k)$ 将收敛到最优值函数 $V^*(\mathbf{x}_k)$,控制对 $(\mathbf{u}_j, \mathbf{w}_j)$ 收敛到最优控制对 $(\mathbf{u}^*, \mathbf{w}^*)$.评价网权值 $\hat{W}_c(j)$ 收敛到 W_c , $W_c = \lim_{j \rightarrow \infty} \hat{W}_c(j)$ ^[17].为了避免神经网络权值在训练过程中陷入到局部极小值,在训练中需要加入持续激励信号.

注1 根据假设2, $f + gu + hw$ 是李普希兹连续的.并且有限的控制输入不可能使得系统状态在一步之内跳跃到无穷大,因此 $f(\mathbf{x}_k) + g(\mathbf{x}_k)\mathbf{u}_k + h(\mathbf{x}_k)\mathbf{w}_k$ 是有限的.考虑到 $V^*(\mathbf{x}_k)$ 对于任意有限的系统状态和控制输入都是有限的,因此一定存在 $0 < \theta < \infty$ 能够保证不等式(28)成立.此外,由于任意的初始值函数 $V_0(\mathbf{x}_k)$ 是有界的,那么不等式(29)也很容易得到满足.

注2 与典型的ADP算法不同,本文采用的是单网络

结构,只利用一个评价网来近似值函数,省略掉了用来近似控制策略和扰动策略的两个控制网.由于本文研究的是模型完全已知仿射非线性系统,因而模型网也被省略.系统状态方程具有的仿射结构保证了控制策略和扰动策略可以根据最优性原理直接通过计算获得.如果系统模型未知或者是非仿射结构,可以通过增加模型网来构建仿射结构的系统状态方程.

单网络ADP值迭代算法具体执行步骤如下:

步骤1 初始化参数 $Q, R, S, \xi, \alpha_c, j_{\max}$,神经网络权值 V_c^* 和 \mathbf{x}_k ;

步骤2 令 $j = 0$, $\hat{W}_c(0) = 0$,使得 $V_0(\mathbf{x}_k) = 0$;

步骤3 根据式(19)计算 \mathbf{u}_0 和 \mathbf{w}_0 ;

步骤4 令 $j = j + 1$;

步骤5 根据式(12)计算 \mathbf{x}_{k+1} ;

步骤6 根据式(27)更新权值 $\hat{W}_c(j + 1)$;

步骤7 根据式(24)计算 $V_{j+1}(\mathbf{x}_k)$;

步骤8 根据式(21)计算 \mathbf{u}_j 和 \mathbf{w}_j ;

步骤9 如果 $\|V_{j+1}(\mathbf{x}_k) - V_j(\mathbf{x}_k)\| < \xi$ 或者迭代次数 $j > j_{\max}$,跳转步骤10,否则跳转步骤4;

步骤10 近似最优的控制对已获得,算法结束.

3.3 事件驱动单网络值迭代算法(Event-triggered single network value iteration algorithm, ET-SNVI)

根据第3.1节可知,事件驱动阈值为 e_T ,事件驱动条件为 $\|e_k\| \leq e_T$.当驱动条件不再满足时,事件驱动误差被重置为零,控制输入 $\mu^*(\mathbf{x}_{k_i})$ 更新.控制输入和扰动输入的计算公式如式(13)和式(8b)所示,其中的最优值函数 $V^*(\mathbf{x}_k)$ 可通过第3.2节中的单网络值迭代算法逼近.因此,最终获得了基于事件驱动的零和博弈问题的近似最优解为

$$\mu(\mathbf{x}_k) = -\frac{1}{2}R^{-1}g^T(\mathbf{x}_{k_i})\lambda_{k_i+1}, \quad k \in [k_i, k_{i+1}), \quad (30a)$$

$$\mathbf{w}(\mathbf{x}_k) = \frac{1}{2}S^{-1}h^T(\mathbf{x}_k)\lambda_{k+1}, \quad (30b)$$

其中协状态 $\lambda_{k_i+1}\mathbf{x}$ 和 λ_{k+1} 中的最优值函数由评价网的输出近似 $V^*(\mathbf{x}_k) = W_c^T\phi_c(V_c^*\mathbf{x}_k)$.

假设5 存在正常数 α, β 和 L_1, K_∞ 类函数 α_1 和 α_2 能够使得下列不等式满足^[13]:

$$\alpha_1(\|\mathbf{x}_k\|) \leq V(\mathbf{x}_k) \leq \alpha_2(\|\mathbf{x}_k\|), \quad (31)$$

$$V(\mathbf{x}_{k+1}) - V(\mathbf{x}_k) \leq -\alpha V(\mathbf{x}_k) + \beta \|e_k\|, \quad (32)$$

$$\alpha_1^{-1}(\|\mathbf{x}_k\|) \leq L_1 \|\mathbf{x}_k\|. \quad (33)$$

定理1 对于离散系统(12),如果假设5成立,对于 $k \in [k_i, k_{i+1})$, $i = 0, 1, \dots$,满足下列不等式:

$$\varphi(k) - \alpha\sigma(k) < 0, \quad (34)$$

其中:

$$\begin{aligned} \sigma(k) &= (1 - \alpha)^{k-k_i} + (1 - (1 - \alpha)^{k-k_i})\varphi(k)/\alpha, \\ \varphi(k) &= (\beta(1 - (2L)^{k-k_i})(1 + L)L_1)/(1 - 2L). \end{aligned}$$

则系统(12)是渐近稳定的.

证 由式(33)可知

$$\|\mathbf{x}_{k_i}\| \leq \alpha_1(V(\mathbf{x}_{k_i})) \leq L_1V(\mathbf{x}_{k_i}). \quad (35)$$

将式(18)和式(35)代入到式(32)中, 可得

$$V(\mathbf{x}_{k+1}) - V(\mathbf{x}_k) \leq -\alpha V(\mathbf{x}_k) + \varphi(k)V(\mathbf{x}_{k_i}). \quad (36)$$

求解式(36), 可得

$$V(\mathbf{x}_k) \leq \sigma(k)V(\mathbf{x}_{k_i}). \quad (37)$$

将式(37)代入式(36), 可得

$$V(\mathbf{x}_{k+1}) - V(\mathbf{x}_k) \leq (\varphi(k) - \alpha\sigma(k))V(\mathbf{x}_{k_i}). \quad (38)$$

应用式(31), 可得

$$\begin{aligned} \Delta V &= V(\mathbf{x}_{k+1}) - V(\mathbf{x}_k) \leq \\ &(\varphi(k) - \alpha\sigma(k))\alpha_2(\|\mathbf{x}_{k_i}\|). \end{aligned} \quad (39)$$

因此, 当不等式(34)成立时, $\Delta V < 0$. 根据Lyapunov稳定性理论系统(12)渐近稳定. 证毕.

本文提出的事件驱动最优控制方案结构图如图1所示, 其具体步骤如下:

步骤 1 初始化参数 α, L, ϵ 和 i_{\max} . 令 $i = 0, k = 0$;

步骤 2 根据式(9)和式(18)计算事件驱动误差 e_k 和阈值 e_T ;

步骤 3 判断 $\|e_k\|$ 是否大于 e_T , 如果大于执行步骤4, 如果小于等于跳转步骤6;

步骤 4 $i = i + 1, \mathbf{x}_{k_i} = \mathbf{x}_k, e_k = 0$;

步骤 5 根据式(30a)计算 $\mu(\mathbf{x}_k)$;

步骤 6 根据式(30b)计算 $w(\mathbf{x}_k)$;

步骤 7 根据式(12)计算 \mathbf{x}_{k+1} ;

步骤 8 如果 $\|\mathbf{x}_{k+1} - \mathbf{x}_k\| \leq \epsilon$, 或者 $i > i_{\max}$, 跳转步骤9, 否则跳转步骤2;

步骤 9 算法结束.

注 3 将值函数 $V(\mathbf{x}_k)$ 定义为系统的李雅普诺夫函数. 根据HJI方程(5)和公式(22), 值函数 $V(\mathbf{x}_k)$ 可以表述为系统状态 \mathbf{x}_k 的相关函数. 如果系统是一个线性系统, 值函数 $V(\mathbf{x}_k) = \mathbf{x}_k^T P \mathbf{x}_k$, 其中 P 为黎卡提方程的解. 显然, 其满足假设5中的不等式(31). 当系统为一个非线性系统的时候, 用评价网 $W_c^T \phi_c(\bar{\mathbf{x}}_{k+1})$ 来逼近 $V(\mathbf{x}_k)$. 适当的选择激活函数 $\phi_c(\cdot)$ 也能够保证不等式(31)成立.

注 4 本文提出的事件驱动单网络值迭代算法是一种离线的算法, 通过在评价网、控制策略和扰动策略之间的不断迭代, 最终获得全局最优控制对, 该最优控制对可以在线直接应用在每一个事件驱动时刻. 而且该算法一般取初始迭代值函数 $V_0(\mathbf{x}_k) = 0$, 不要求提供一个初始稳定增益. 这对非线性系统来说是非常重要的, 因为非线性系统的初始稳定增益并不容易获得.

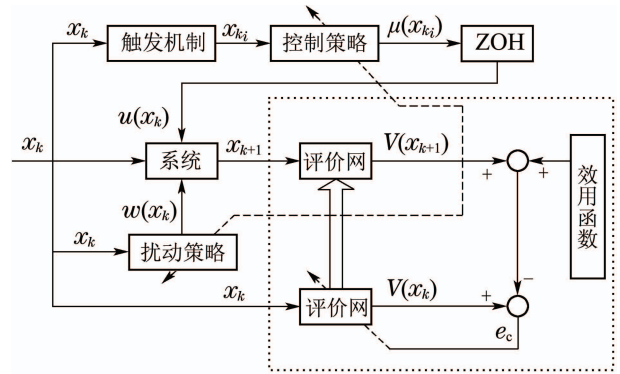


图 1 事件驱动最优控制方案结构图

Fig. 1 The structure of the event-triggered optimal control scheme

4 仿真验证(Simulation)

为验证本文所提的事件驱动最优控制方案的有效性, 本小节将该方案应用到了F-16战斗机的一个非线性系统的仿真例子中.

例 1 F-16战斗机.

考虑如下的F-16战斗机的离散数学模型^[16]:

$$\mathbf{x}_{k+1} = F\mathbf{x}_k + G\mathbf{u}_k + H\mathbf{w}_k, \quad (40)$$

其中: $\mathbf{x}_k = [\alpha_k \ q_k \ \delta_{ek}]^T$, α_k 为攻击角度, q_k 为俯仰角速度, δ_{ek} 为升降舵偏转角, u 为制动器电压, w 为作用到攻击角度上的阵风.

$$F = \begin{bmatrix} 0.906488 & 0.0816012 & -0.0005 \\ 0.0741349 & 0.90121 & -0.000708383 \\ 0 & 0 & 0.132655 \end{bmatrix},$$

$$G = \begin{bmatrix} -0.00150808 \\ -0.0096 \\ 0.867345 \end{bmatrix}, \quad H = \begin{bmatrix} 0.00951892 \\ 0.00038373 \\ 0 \end{bmatrix}.$$

性能指标函数如式(2)所示, 其中: $Q \in \mathbb{R}^{3 \times 3}$, $R \in \mathbb{R}^{1 \times 1}$ 和 $S \in \mathbb{R}^{1 \times 1}$ 为单位阵. 飞行器的初始状态设定为 $\mathbf{x}_0 = [4 \ 2 \ 5]^T$. 采用一个3-8-1的3层神经网络来构成评价网, 评价网的初始权值 V_c 在 $[-1, 1]$ 之间随机生成. $\hat{W}_c(0)$ 设定为零, 从而保证初始迭代值函数 $V_0(\mathbf{x}_k) = 0$. 激活函数 $\phi_c(\cdot)$ 选为tansig函数. 评价网学习率 $\alpha_c = 0.2$. 计算精度为 $\xi = 10^{-5}$. 评价网训练了2000次, 为了避免神经网络权值陷入局部极小值, 在前800迭代步中加入了持续激励. 评价网权值的收敛轨迹如图2所示.

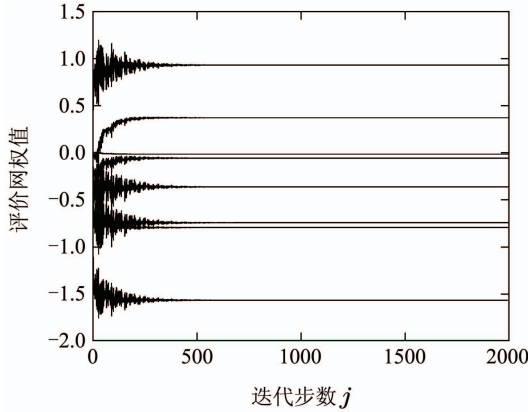


图2 评价网权值收敛轨迹

Fig. 2 The convergent trajectories of critic network weights

由式(18)可知, 事件驱动阈值 e_T 与 α 和 L 的值有关. 为了选择适当的 α 和 L , 作者进行了一系列的试验. 当 $L = 0.2$ 时, α 取不同的值时, 累计采样次数和系统状态曲线如图3所示. 图中箭头指向的方向为 α 增大的方向. 从图3中可以看出, 随着 α 的增大, 累计采样次数逐渐减少, 系统状态 x_1 和 x_2 逐渐接近最优状态轨迹. 但是系统状态 x_3 随着 α 的增大, 距离最优状态轨迹越来越远. 在综合考虑了累计采样次数和系统性能之后, 最终选择 $\alpha = 0.1$. 同理, 当 $\alpha = 0.1$ 时, 选取不同的 L 进行了一系列的仿真, 发现随着 L 的增大, 累计采样次数逐渐减少, 但是对系统状态的影响不大. 最终, 本文选取了 $\alpha = 0.1, L = 0.1$ 来确定事件驱动阈值.

$$e_T = 0.1375(1 - (0.2)^{k-k_i})\|\mathbf{x}_{k_i}\|. \quad (41)$$

当 $\alpha = 0.1, L = 0.1$ 时, 系统的状态轨迹如图4所示. 从图4可以看出, 系统在796步之后能够达到精度 $\epsilon = 10^{-5}$. 事件驱动误差的范数 $\|e_k\|$ 和阈值 e_T 的变化情况如图5所示.

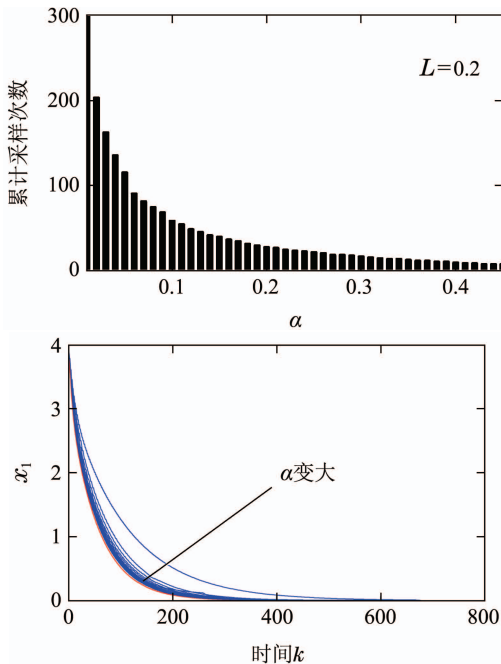


图4 系统状态轨迹

Fig. 4 The trajectories of system states

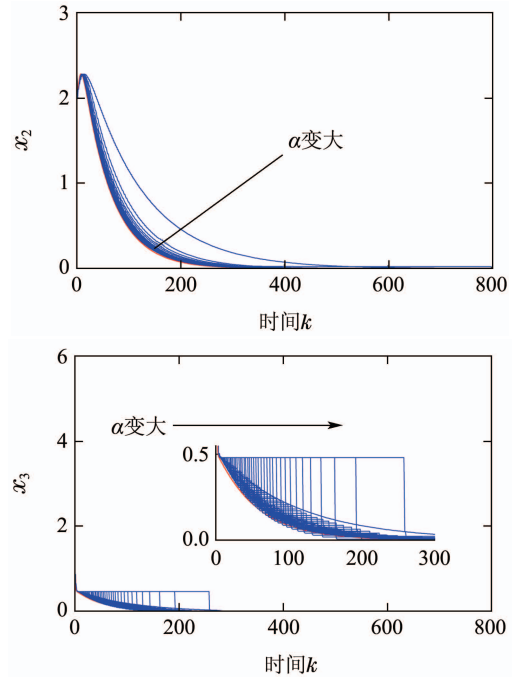


图3 α 取不同值时累计采样次数和系统状态轨迹

Fig. 3 The number of cumulative samples and the trajectories of system states with different α

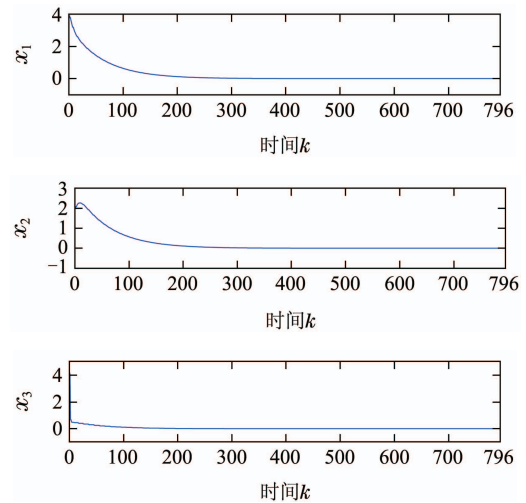


图4 系统状态轨迹

Fig. 4 The trajectories of system states

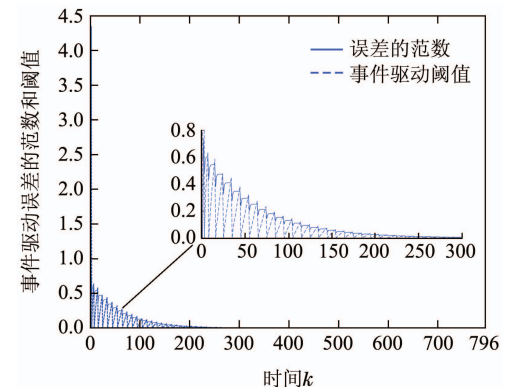


图5 事件驱动误差的范数和事件驱动阈值轨迹

Fig. 5 The trajectories of the norm of event-triggered error and event-triggered threshold

由于事件驱动条件在前300步变化明显,所以在图5中给出了前300步的局部放大图. 控制输入和扰动输入的变化轨迹如图6所示. 图7给出了典型ADP算法和事件驱动单网络值迭代算法的累计采样次数对比图.

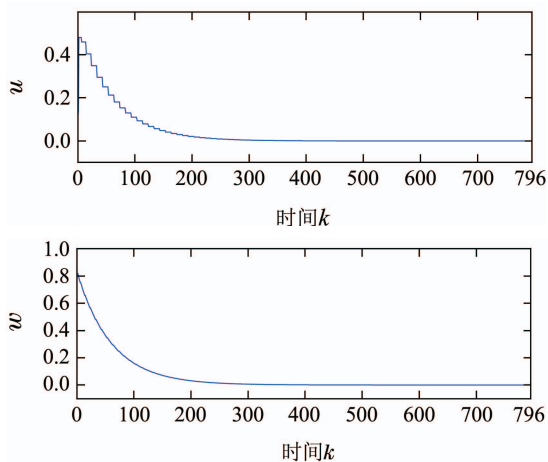


图 6 控制输入和扰动输入轨迹

Fig. 6 The trajectories of control input and distribute input

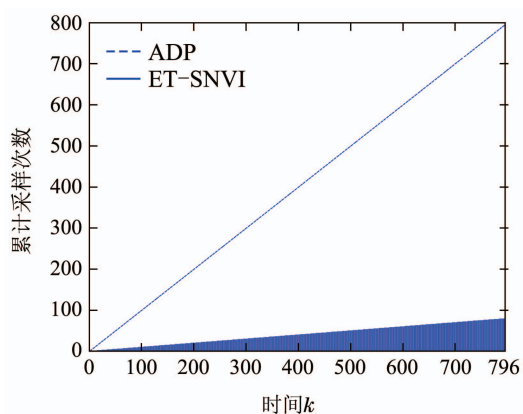


图 7 累计采样次数

Fig. 7 The cumulative samples

如图7所示, 本文所提出的事件驱动单网络值迭代算法只需要进行80次采样, 而典型的时间驱动的ADP算法则需要进行796次采样. 本文所提算法能够减少近90%的通讯次数和计算量. 同时, 由于只采用了一个神经网络, 省略了用来近似控制策略和扰动策略的两个控制网, 所以减少了近67%的神经网络权值训练量.

例 2 离散非线性系统.

考虑如下的离散非线性零和博弈问题, 其状态方程为

$$\mathbf{x}_{k+1} = f(\mathbf{x}_k) + g(\mathbf{x}_k)\mathbf{u}_k + h(\mathbf{x}_k)\mathbf{w}_k, \quad (42)$$

其中:

$$f(\mathbf{x}_k) = [f_1(\mathbf{x}_k) \ f_2(\mathbf{x}_k)]^T, \\ f_1(\mathbf{x}_k) = 0.9x_{1k} - 0.1x_{2k},$$

$$f_2(\mathbf{x}_k) = 0.9x_{1k} - 0.5x_{2k} + 0.025x_{2k}(\cos x_{1k})^2, \\ g(\mathbf{x}_k) = [-0.01x_{1k} \ 0.5 \cos(2x_{1k})]^T, \\ h(\mathbf{x}_k) = [0.1 \ 0.1 \sin(4x_{1k}) + 2]^T.$$

性能指标函数如式(2)所示, 其中 Q, R 和 S 为具有适当维数的单位阵. 初始状态设定为 $\mathbf{x}_0 = [4 \ 2]^T$. 采用一个2-8-1的3层神经网络来构成评价网, 评价网的初始权值 V_c 在 $[-1, 1]$ 之间随机生成. \hat{W}_c 设定为零. 激活函数 $\phi_c(\cdot)$ 选为tansig函数. 评价网学习率 $\alpha_c = 0.1$. 选取 $\alpha = 0.1, L = 0.2$ 来确定事件驱动阈值.

$$e_T = 0.2(1 - (0.4)^{k-k_i})\|\mathbf{x}_{k_i}\|. \quad (43)$$

系统的状态轨迹如图8所示. 从图8可以看出, 系统在125步之后能够达到精度 $\epsilon = 10^{-5}$. 图9给出了控制输入和扰动输入的变化轨迹. 事件驱动误差的范数 $\|\mathbf{e}_k\|$ 和事件驱动阈值 e_T 的变化情况如图10所示. 与典型的时间驱动的ADP算法需要进行125次采样相比, 本文所提的事件驱动最优控制方法只进行了63次采样, 减少了近50%的网络通讯量和控制器计算以及执行次数.

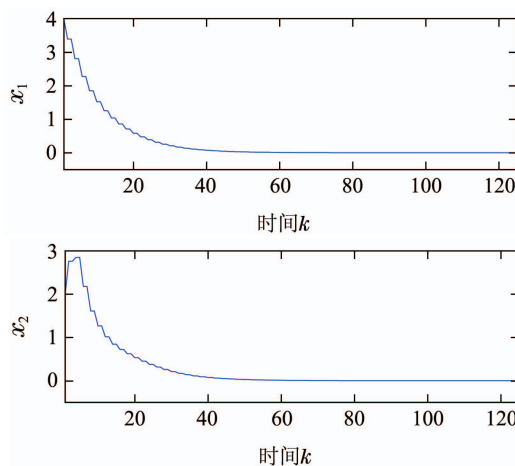


图 8 系统状态轨迹

Fig. 8 The trajectories of system states

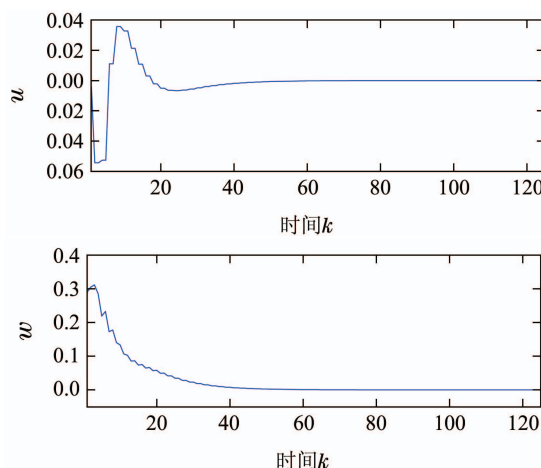


图 9 控制输入和扰动输入轨迹

Fig. 9 The trajectories of control input and distribute input

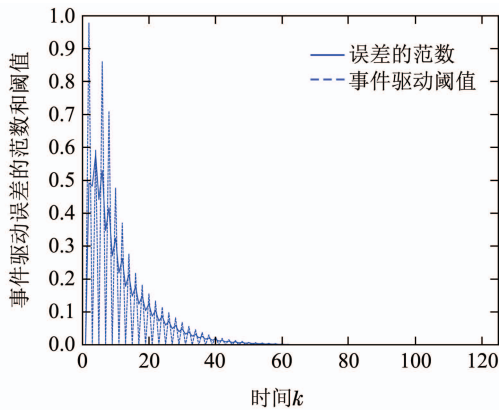


图 10 事件驱动误差的范数和事件驱动阈值的轨迹

Fig. 10 The trajectories of the norm of event-triggered error and event-triggered threshold

从上述仿真结果中可以看出, 本文提出的零和博弈问题的事件驱动最优控制方案, 能够很好的镇定系统, 并且获得零和博弈问题的近似最优控制对. 通过事件驱动机制, 能够有效减少控制输入与系统之间的数据传输次数、控制器计算次数以及执行器变动次数. 并且单网络值迭代算法能够有效降低神经网络权值的训练量.

5 结论(Conclusions)

本文研究了博弈论中常见的零和博弈问题. 为了降低数据传输和计算次数, 获得最优控制对, 提出了一种基于事件驱动的单网络值迭代算法. 将事件驱动控制应用到零和博弈问题求解中, 设计新型事件驱动阈值. 采用单网络值迭代算法, 利用一个神经网络构建评价网, 根据Bellman最优性原理直接计算控制对, 通过在评价网、控制策略和扰动策略之间进行迭代, 获得最优值函数. 给出了神经网络权训练步骤. 接着, 利用Lyapunov理论证明了闭环系统的稳定性, 并给出了事件驱动最优控制方案的执行步骤. 最后, 将该方案应用于F-16战斗机和—个非线性系统的零和博弈问题仿真实验中, 仿真结果表明所提方法能够获得近似最优控制对, 并且成功地降低了网络通信频率, 控制输入的执行次数以及神经网络权值的训练次数.

参考文献(References):

- [1] FU Yue, CHAI Tianyou. Online solution of two-player zero-sum games for linear systems with unknown dynamics [J]. *Control Theory & Applications*, 2015, 32(2): 196 – 201.
(富月, 柴天佑. 具有未知动态的线性系统二人零和博弈问题在线学习方案 [J]. 控制理论与应用, 2015, 32(2): 196 – 201.)
- [2] YVES A, PEREZ V. Iterative strategies for solving linearized discrete mean field games systems [J]. *Netw Heterog Media*, 2012, 7(2): 197 – 217.
- [3] FU Y, FU J, CHAI T. Robust adaptive dynamic programming of two-player zero-sum games for continuous-time linear systems [J]. *IEEE Transactions on Neural Networks and Learning Systems*, 2015, 26(12): 3314 – 3319.
- [4] ASTROM K J, BERNHARDSSON B M. Comparison of Riemann and Lebesgue sampling for first order stochastic systems [C] // *Proceedings of the 41st IEEE Conference on Decision Control*. Las Vegas: IEEE, 2002, 2: 2011 – 2016.

- [5] HEEMELES W, DONKERS M, TEEL A. Periodic event-triggered control for linear systems [J]. *IEEE Transactions on Automatic Control*, 2013, 58(4): 847 – 861.
- [6] LIANG Yuan, QI Guoqing, LI Yinya, et al. Design and application of event-triggered mechanism for a kind of optical-electronic tracking system [J]. *Control Theory & Applications*, 2017, 34(10): 1328 – 1338.
(梁苑, 戚国庆, 李银份, 等. 一类光电跟踪系统中事件触发机制的设计及应用 [J]. 控制理论与应用, 2017, 34(10): 1328 – 1338.)
- [7] SAHOO A, XU H, JAGANNATHAN S. Neural network-based event-triggered state feedback control of nonlinear continuous-time systems [J]. *IEEE Transactions on Neural Networks and Learning Systems*, 2016, 27(3): 497 – 509.
- [8] VAMVOUDAKIS K G. Event-triggered optimal adaptive control algorithm for continuous-time nonlinear systems [J]. *IEEE/CAA Journal of Automatica Sinica*, 2014, 1(3): 282 – 293.
- [9] ZHANG Q, ZHAO D, ZHU Y. Event-triggered H_∞ control for continuous-time nonlinear system via concurrent learning [J]. *IEEE Transactions on Systems, Man, and Cybernetics*, 2017, 47(7): 1071 – 1081.
- [10] WERBOS P J. Approximate dynamic programming for real-time control and neural modeling [M] // *Handbook of Intelligent Control: Neural, Fuzzy and Adaptive Approaches*. New York: Van Nostrand Reinhold, 1992.
- [11] QU Qiuxia, LUO Yanhong, ZHANG Huaguang. Robust approximate optimal tracking control of time-varying trajectory for nonlinear affine systems [J]. *Control Theory & Applications*, 2016, 33(1): 77 – 84.
(屈秋霞, 罗艳红, 张化光. 针对时变轨迹的非线性仿射系统的鲁棒近似最优跟踪控制 [J]. 控制理论与应用, 2016, 33(1): 77 – 84.)
- [12] WANG D, HE H, LIU D. Adaptive critic nonlinear robust control: a survey [J]. *IEEE Transactions on Cybernetics*, 2017, 47(10): 3429 – 3451.
- [13] DONG L, ZHONG X N, SUN C Y, et al. Adaptive event-triggered control based on heuristic dynamic programming for nonlinear discrete-time systems [J]. *IEEE Transactions on Neural Networks and Learning Systems*, 2017, 28(7): 1594 – 1605.
- [14] LUO B, WU H N, HUANG T. Off-policy reinforcement learning for H_∞ control design [J]. *IEEE Transactions on Cybernetics*, 2015, 45(1): 65 – 76.
- [15] ZHANG X, ZHANG H G, WANG F Y. A new iteration approach to solve a class of Finite-horizon continuous-time nonaffine nonlinear zero-sum game [J]. *International Journal of Innovative, Computing, Information and Control*, 2011, 7(2): 597 – 608.
- [16] AL-TAMIMI A, KHALAF M, LEWIS F L. Adaptive critic designs for discrete-time zero-sum games with application to H_∞ control [J]. *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics*, 2007, 37(1): 240 – 247.
- [17] LIU D, LI H, WANG D. Neural-network-based zero-sum game for discrete-time nonlinear systems via iterative adaptive dynamic programming algorithm [J]. *Neurocomputing*, 2013, 110(8): 92 – 100.
- [18] JIANG Z P, WANG Y. Input-to-state stability for discretetime nonlinear systems [J]. *Automatica*, 2001, 37(6): 857 – 869.

作者简介:

张欣 (1982–), 女, 讲师, 博士, 目前研究方向为自适应动态规划、最优控制、神经网络等, E-mail: zhangxin@upc.edu.cn;

薄迎春 (1977–), 男, 讲师, 博士, 目前研究方向为非线性系统控制、动态规划, E-mail: boyingchun@sina.com;

崔黎黎 (1983–), 女, 讲师, 博士, 目前研究方向为自适应动态规划、非线性系统、微分对策、最优控制等, E-mail: cuilili8396@163.com.