

An Optimal Control Scheme for a Class of Discrete-time Nonlinear Systems with Time Delays Using Adaptive Dynamic Programming

WEI Qing-Lai¹ ZHANG Hua-Guang² LIU De-Rong¹ ZHAO Yan³

Abstract In this paper, an optimal control scheme for a class of nonlinear systems with time delays in both state and control variables with respect to a quadratic performance index function is proposed using a new iterative adaptive dynamic programming (ADP) algorithm. By introducing a delay matrix function, the explicit expression of the optimal control is obtained using the dynamic programming theory and the optimal control can iteratively be obtained using the adaptive critic technique. Convergence analysis is presented to prove that the performance index function can reach the optimum by the proposed method. Neural networks are used to approximate the performance index function, compute the optimal control policy, solve delay matrix function, and model the nonlinear system, respectively, for facilitating the implementation of the iterative ADP algorithm. Two examples are given to demonstrate the validity of the proposed optimal control scheme.

Key words Adaptive dynamic programming (ADP), approximate dynamic programming, time delay, optimal control, nonlinear system, neural networks

DOI 10.3724/SP.J.1004.2010.00121

The optimal control problem of nonlinear systems has always been a key focus in the control field in the last several decades. Coupled with this is the fact that nothing can happen instantaneously, as is so often presumed in many mathematical models. So strictly speaking, time delays exist in the most practical control systems. Time delays may result in degradation in the control efficiency even instability of the control systems. So there have been many studies on the control systems with time delay in various research fields, such as electrical, chemical engineering, and networked control^[1-2]. The optimal control problem for the time-delay systems always attracts considerable attention of the researchers and many results have been obtained^[3-5]. In general, the optimal control for the time-delay systems is an infinite-dimensional control problem^[3], which is very difficult to solve. So many analysis and applications are limited to a very simple case: the linear systems with only state delays^[6]. For nonlinear case with state delays, the traditional method is to adopt fuzzy method and robust method, which transforms the nonlinear time-delay systems to linear systems^[7]. For systems with time delays both in states and controls, it is still an open problem^[4-5]. The main difficulty lies in the formulation of the optimal controller which must use the information of the delayed control term so as to obtain an efficient control. This makes the analysis of the system much more difficult, and there is no method strictly facing this problem even in the linear cases. This motivates our research.

Adaptive dynamic programming (ADP) is a powerful tool in solving optimal control problems^[8-9] and has attracted considerable attention from many researchers in recent years, such as [10-16]. However, most of the results focus on the optimal control problems without delays. To

the best of our knowledge, there are no results discussing how to use ADP to solve the time-delay optimal control problems. In this paper, the time-delay optimal control problem is solved by the iterative ADP algorithm for the first time. By introducing a delay matrix function, the explicit expression of the optimal control function is obtained. The optimal control can iteratively be obtained using the proposed iterative ADP algorithm which avoids the infinite-dimensional computation. Also, it is proved that the performance index function converges to the optimum using the proposed iterative ADP algorithm.

This paper is organized as follows. Section 1 presents the preliminaries. In Section 2, the time-delay optimal control scheme is proposed based on iterative ADP algorithm. In Section 3, the neural network implementation for the control scheme is discussed. In Section 4, two examples are given to demonstrate the effectiveness of the proposed control scheme. The conclusion is drawn in Section 5.

1 Preliminaries

Basically, we consider the following discrete-time affine nonlinear system with time delays in state and control variables as follows:

$$\begin{aligned} \mathbf{x}(k+1) = & f(\mathbf{x}(k), \mathbf{x}(k-\sigma)) + g_0(\mathbf{x}(k), \mathbf{x}(k-\sigma))\mathbf{u}(k) + \\ & g_1(\mathbf{x}(k), \mathbf{x}(k-\sigma))\mathbf{u}(k-\tau) \end{aligned} \quad (1)$$

with the initial condition given by $\mathbf{x}(s) = \phi(s)$, $s = -\sigma, -\sigma + 1, \dots, 0$, where $\mathbf{x}(k) \in \mathbf{R}^n$ is the state vector, $f: \mathbf{R}^n \times \mathbf{R}^n \rightarrow \mathbf{R}^n$ and $g_0, g_1: \mathbf{R}^n \times \mathbf{R}^n \rightarrow \mathbf{R}^{n \times m}$ are differentiable functions and the control $\mathbf{u}(k) \in \mathbf{R}^m$. The state and control delays σ and τ are both nonnegative integral numbers. Assume that $f(\mathbf{x}(k), \mathbf{x}(k-\sigma)) + g_0(\mathbf{x}(k), \mathbf{x}(k-\sigma))\mathbf{u}(k) + g_1(\mathbf{x}(k), \mathbf{x}(k-\sigma))\mathbf{u}(k-\tau)$ is Lipschitz continuous on a set Ω in \mathbf{R}^n containing the origin, and that system (1) is controllable in the sense that there exists a bounded control on Ω that asymptotically stabilizes the system. In this paper, how to design an optimal state feedback controller for this class of delayed discrete-time systems is mainly discussed. Therefore, it is desired to find the optimal control $\mathbf{u}(\mathbf{x})$ satisfying $\mathbf{u}(\mathbf{x}(k)) = \mathbf{u}(k)$ to minimize the generalized performance functional as follows:

Manuscript received September 5, 2008; accepted March 3, 2009
Supported by National High Technology Research and Development Program of China (863 Program) (2006AA04Z183), National Natural Science Foundation of China (60621001, 60534010, 60572070, 60774048, 60728307), and the Program for Changjiang Scholars and Innovative Research Groups of China (60728307, 4031002)

1. Key Laboratory of Complex Systems and Intelligence Science, Institute of Automation, Chinese Academy of Sciences, Beijing 100190, P. R. China 2. School of Information Science and Engineering, Northeastern University, Shenyang 110004, P. R. China 3. Department of Automatic Control Engineering, Shenyang Institute of Engineering, Shenyang 110136, P. R. China

$$V(\mathbf{x}(0), \mathbf{u}) = \sum_{k=0}^{\infty} \left(\mathbf{x}^T(k) Q_0 \mathbf{x}(k) + 2\mathbf{x}^T(k) Q_1 \mathbf{x}(k - \sigma) + \mathbf{x}^T(k - \sigma) Q_2 \mathbf{x}(k - \sigma) + \mathbf{u}^T(k) R_0 \mathbf{u}(k) + 2\mathbf{u}^T(k) R_1 \mathbf{u}(k - \tau) + \mathbf{u}^T(k - \tau) R_2 \mathbf{u}(k - \tau) \right) \quad (2)$$

where $\begin{bmatrix} Q_0 & Q_1 \\ Q_1^T & Q_2 \end{bmatrix} \geq 0$, $\begin{bmatrix} R_0 & R_1 \\ R_1^T & R_2 \end{bmatrix} > 0$, and $l(\mathbf{x}(k), \mathbf{x}(k - \sigma), \mathbf{u}(k), \mathbf{u}(k - \tau)) = \mathbf{x}^T(k) Q_0 \mathbf{x}(k) + 2\mathbf{x}^T(k) Q_1 \mathbf{x}(k - \sigma) + \mathbf{x}^T(k - \sigma) Q_2 \mathbf{x}(k - \sigma) + \mathbf{u}^T(k) R_0 \mathbf{u}(k) + 2\mathbf{u}^T(k) R_1 \mathbf{u}(k - \tau) + \mathbf{u}^T(k - \tau) R_2 \mathbf{u}(k - \tau)$ is the utility function. Let $V^*(\mathbf{x})$ denote the optimal performance index function which satisfies

$$V^*(\mathbf{x}) = \min_{\mathbf{u}} V(\mathbf{x}, \mathbf{u}) \quad (3)$$

According to the Bellman's optimal principle, we can get the following Hamilton-Jacobi-Bellman (HJB) equation

$$V^*(\mathbf{x}(k)) = \min_{\mathbf{u}(k)} \left\{ \mathbf{x}^T(k) Q_0 \mathbf{x}(k) + 2\mathbf{x}^T(k) Q_1 \mathbf{x}(k - \sigma) + \mathbf{x}^T(k - \sigma) Q_2 \mathbf{x}(k - \sigma) + \mathbf{u}^T(k) R_0 \mathbf{u}(k) + 2\mathbf{u}^T(k) R_1 \mathbf{u}(k - \tau) + \mathbf{u}^T(k - \tau) R_2 \mathbf{u}(k - \tau) + V^*(\mathbf{x}(k + 1)) \right\} \quad (4)$$

For the optimal control problem, the state feedback control $\mathbf{u}(\mathbf{x})$ must not only stabilize the system on Ω but also guarantee that (2) is finite, i.e., $\mathbf{u}(\mathbf{x})$ must be admissible^[17].

Definition 1. A control $\mathbf{u}(\mathbf{x})$ is defined to be admissible with respect to (3) on Ω if $\mathbf{u}(\mathbf{x})$ is continuous on Ω , $\mathbf{u}(0) = \mathbf{0}$, $\mathbf{u}(\mathbf{x})$ stabilizes (1) on Ω , and $\forall \mathbf{x}(0) \in \Omega$, $V(\mathbf{x}(0))$ is finite.

2 Properties of the iterative ADP approach

Since the nonlinear delayed system (1) is infinite-dimensional^[3] and the control variable $\mathbf{u}(k)$ couples with $\mathbf{u}(k - \tau)$, it is nearly impossible to obtain the expression of the optimal control by solving the HJB equation (1). To overcome the difficulty, a new iterative algorithm is proposed in this paper. The following lemma is necessary to apply the algorithm.

Lemma 1. For the delayed nonlinear system (1) with respect to the performance index function (2), if there exists a control $\mathbf{u}(k) \neq \mathbf{0}$ at time point k , then there exists a bounded matrix function $M(k)$ that makes

$$\mathbf{u}(k - \tau) = M(k) \mathbf{u}(k) \quad (5)$$

hold for $j = 0, 1, \dots, n$.

Proof. As $\mathbf{u}(k)$ and $\mathbf{u}(k - \tau_j)$, $j = 0, 1, \dots, n$ are bounded real vectors, can construct a function that satisfies

$$\mathbf{u}(k - \tau) = h(\mathbf{u}(k)) \quad (6)$$

where $j = 0, 1, \dots, n$. Then, using the method of undetermined coefficients, let $M(\mathbf{u}(k))$ satisfy

$$h(\mathbf{u}(k)) = M(\mathbf{u}(k)) \mathbf{u}(k) \quad (7)$$

Then, we can obtain $M(\mathbf{u}(k))$ expressed as

$$M(\mathbf{u}(k)) = h(\mathbf{u}(k)) \mathbf{u}^T(k) \left(\mathbf{u}(k) \mathbf{u}^T(k) \right)^{-1} \quad (8)$$

where $(\mathbf{u}(k) \mathbf{u}^T(k))^{-1}$ means the generalized inverse matrix of $(\mathbf{u}(k) \mathbf{u}^T(k))$. On the other side, $\mathbf{u}(k)$ and $\mathbf{u}(k - \tau)$ are both bounded real vectors, then we have $h(\mathbf{u}(k))$ and $(\mathbf{u}(k) \mathbf{u}^T(k))^{-1}$ are bounded. So $M(k) = M(\mathbf{u}(k))$ is the solution. \square

According to Lemma 1, the HJB equation becomes

$$V^*(\mathbf{x}(k)) = \mathbf{x}^T(k) Q_0 \mathbf{x}(k) + 2\mathbf{x}^T(k) Q_1 \mathbf{x}(k - \sigma) + \mathbf{x}^T(k - \sigma) Q_2 \mathbf{x}(k - \sigma) + \mathbf{u}^{*\text{T}}(k) R_0 \mathbf{u}^*(k) + 2\mathbf{u}^{*\text{T}}(k) R_1 M^*(k) \mathbf{u}^*(k) + \mathbf{u}^{*\text{T}}(k) M^{*\text{T}}(k) R_2 M^*(k) \mathbf{u}^*(k) + V^*(\mathbf{x}(k + 1)) \quad (9)$$

where $\mathbf{u}^*(k)$ is the optimal control and $\mathbf{u}^*(k - \tau) = M^*(k) \mathbf{u}^*(k)$.

2.1 Derivation of the iterative ADP algorithm

According to the Bellman's principle of optimality, we can obtain the optimal control by differentiating the HJB equation (9) with respect to control \mathbf{u} .

Then, we can obtain the optimal control $\mathbf{u}^*(k)$ formulated as

$$\mathbf{u}^*(k) = -\frac{1}{2} \left(R_0 + 2R_1 M^*(k) + M^{*\text{T}}(k) R_2 M^*(k) \right)^{-1} \times \left(g_0(\mathbf{x}(k), \mathbf{x}(k - \sigma)) + g_1(\mathbf{x}(k), \mathbf{x}(k - \sigma)) M^*(k) \right)^{\text{T}} \times \frac{\partial V^*(\mathbf{x}(k + 1))}{\partial \mathbf{x}(k + 1)} \quad (10)$$

In (10), the inverse of the term $(R_0 + 2R_1 M^*(k) + M^{*\text{T}}(k) R_2 M^*(k))$ should exist and a proof is presented in Appendix to guarantee the existence of the inverse.

From (10), the explicit optimal control expression \mathbf{u}^* is obtained by solving the HJB equation (9). We can see that the optimal control \mathbf{u}^* depends on M^* and $V^*(\mathbf{x})$, where $V^*(\mathbf{x})$ is a solution to the HJB equation (9). While how to solve the HJB equation is still open, there is currently no method for rigorously seeking for this performance index function of this delayed optimal control problem. Furthermore, the optimal delay matrix function M^* is also unknown which makes the optimal control \mathbf{u}^* more difficult to obtain. So an iterative index i is introduced into the ADP approach to obtain the optimal control iteratively.

Firstly, for $i = 0, 1, \dots$, let

$$\mathbf{u}^{(i+1)}(k - \tau) = M^{(i)}(k) \mathbf{u}^{(i+1)}(k) \quad (11)$$

where $M^{(0)}(k) = I$ and $\mathbf{u}^{(0)}(k - \tau) = M^{(0)}(k) \mathbf{u}^{(0)}(k)$. We start with initial performance index $V^{(0)}(\mathbf{x}(k)) = 0$, and the control $\mathbf{u}^{(0)}(k)$ can be computed as follows

$$\mathbf{u}^{(0)}(\mathbf{x}(k)) = \arg \min_{\mathbf{u}} \left\{ \Gamma^0 + V^{(0)}(\mathbf{x}(k + 1)) \right\} \quad (12)$$

where

$$\begin{aligned} \Gamma^0 = & \mathbf{x}^T(k) Q_0 \mathbf{x}(k) + 2\mathbf{x}^T(k) Q_1 \mathbf{x}(k - \sigma) + \\ & \mathbf{x}^T(k - \sigma) Q_2 \mathbf{x}(k - \sigma) + \mathbf{u}^{(0)\text{T}}(k) R_0 \mathbf{u}^{(0)}(k) + \\ & 2\mathbf{u}^{(0)\text{T}}(k) R_1 M^{(0)}(k) \mathbf{u}^{(0)}(k) + \\ & \mathbf{u}^{(0)\text{T}}(k) M^{(0)\text{T}}(k) R_2 M^{(0)}(k) \mathbf{u}^{(0)}(k) \end{aligned}$$

Then, the performance index function is updated as

$$V^{(1)}(\mathbf{x}(k)) = \Gamma^0 + V^{(0)}(\mathbf{x}(k + 1)) \quad (13)$$

Thus, for $i = 1, 2, \dots$, the iterative ADP can be used to implement the iteration between

$$\begin{aligned} \mathbf{u}^{(i)}(\mathbf{x}(k)) = \arg \min_{\mathbf{u}} \left\{ \Gamma^{(i)} + V^{(i)}(\mathbf{x}(k+1)) \right\} = \\ -\frac{1}{2} \left(R_0 + 2R_1 M^{(i-1)}(k) + \right. \\ \left. M^{(i-1)\top}(k) R_2 M^{(i-1)}(k) \right)^{-1} (g_0(\mathbf{x}(k), \mathbf{x}(k-\sigma)) + \\ g_1(\mathbf{x}(k), \mathbf{x}(k-\sigma)) M^{(i-1)}(k)) \frac{\partial V^{(i)}(\mathbf{x}(k+1))}{\partial \mathbf{x}(k+1)} \end{aligned} \quad (14)$$

where

$$\begin{aligned} \Gamma^{(i)} = \mathbf{x}^\top(k) Q_0 \mathbf{x}(k) + 2\mathbf{x}^\top(k) Q_1 \mathbf{x}(k-\sigma) + \\ \mathbf{x}^\top(k-\sigma) Q_2 \mathbf{x}(k-\sigma) + \mathbf{u}^{(i)\top}(k) R_0 \mathbf{u}^{(i)}(k) + \\ 2\mathbf{u}^{(i)\top}(k) R_1 M^{(i-1)}(k) \mathbf{u}^{(i)}(k) + \\ \mathbf{u}^{(i)\top}(k) M^{(i-1)\top}(k) R_2 M^{(i-1)}(k) \mathbf{u}^{(i)}(k) \end{aligned}$$

and

$$V^{(i+1)}(\mathbf{x}(k)) = \Gamma^{(i)} + V^{(i)}(\mathbf{x}(k+1)) \quad (15)$$

Then, the optimal control can be obtained iteratively. From (14) and (15), it can be seen that during the iteration process, the control actions for different control steps obey different control laws. After the iteration number of i , the obtained control laws sequence is $(\mathbf{u}^{(0)}, \mathbf{u}^{(1)}, \dots, \mathbf{u}^{(i)})$. For the infinite-horizon problem, both the optimal performance index function and the optimal control law are unique. Therefore, it is necessary to show that the iterative performance index function $V^{(i)}(\mathbf{x}(k))$ will converge when the iteration number $i \rightarrow \infty$ under the iterative control $\mathbf{u}^{(i)}(k)$ and this will be proved in the following subsection.

2.2 Properties of the iterative ADP algorithm

In this subsection, we focus on the proof of convergence of the iteration between (14) and (15), with the performance index $V^{(i)}(\mathbf{x}(k)) \rightarrow V^*(\mathbf{x}(k)), \forall k$.

Lemma 2^[17]. Let $\tilde{\mathbf{u}}^{(i)}(k), k = 0, 1, \dots$ be any sequence of control, and $\mathbf{u}^{(i)}(k)$ be expressed as (14). Define $V^{(i+1)}(\mathbf{x}(k))$ as (15) and $\Lambda^{(i+1)}(\mathbf{x}(k))$ as

$$\begin{aligned} \Lambda^{(i+1)}(\mathbf{x}(k)) = \mathbf{x}^\top(k) Q_0 \mathbf{x}(k) + 2\mathbf{x}^\top(k) Q_1 \mathbf{x}(k-\sigma) + \\ \mathbf{x}^\top(k-\sigma) Q_2 \mathbf{x}(k-\sigma) + \tilde{\mathbf{u}}^{(i)\top}(k) R_0 \tilde{\mathbf{u}}^{(i)}(k) + \\ 2\tilde{\mathbf{u}}^{(i)\top}(k) R_1 M^{(i-1)}(k) \tilde{\mathbf{u}}^{(i)}(k) + \\ \tilde{\mathbf{u}}^{(i)\top}(k) M^{(i-1)\top}(k) R_2 M^{(i-1)}(k) \tilde{\mathbf{u}}^{(i)}(k) + \\ \Lambda^{(i)}(\mathbf{x}(k+1)) \end{aligned} \quad (16)$$

If $V^{(0)}(\mathbf{x}(k)) = \Lambda^{(0)}(\mathbf{x}(k)) = 0$, then $V^{(i)}(\mathbf{x}(k)) \leq \Lambda^{(i)}(\mathbf{x}(k)), \forall i$. In order to prove the convergence of the performance index function, the following theorem is also necessary.

Theorem 1. Let the performance index function $V^{(i)}(\mathbf{x}(k))$ be defined by (15). If $\mathbf{x}(k)$ for system (1) is controllable, then there exists an upper bound Y such that $0 \leq V^{(i)}(\mathbf{x}(k)) \leq Y, \forall i$.

Proof. As system (1) is Lipschitz, $M^{(i)}(k)$ is a bounded matrix for $i = 0, 1, \dots$. Define a delay matrix function

$\bar{M}(k)$ which makes

$$\begin{aligned} \boldsymbol{\chi}^\top \left(R_0 + 2R_1 \bar{M}(k) + \bar{M}^\top(k) R_2 \bar{M}(k) \right) \boldsymbol{\chi} - \boldsymbol{\chi}^\top (R_0 + \\ 2R_1 M^{(i)}(k) + M^{(i)\top}(k) R_2 M^{(i)}(k)) \boldsymbol{\chi} \geq 0 \end{aligned} \quad (17)$$

hold for $\forall i$, where $\boldsymbol{\chi}$ is any nonzero m -dimensional vector. Let $\bar{\mathbf{u}}(k), k = 0, 1, \dots$ be any admissible control input. Define a new sequence $P^{(i)}(\mathbf{x}(k))$ as follows:

$$\begin{aligned} P^{(i+1)}(\mathbf{x}(k)) = \mathbf{x}^\top(k) Q_0 \mathbf{x}(k) + 2\mathbf{x}^\top(k) Q_1 \mathbf{x}(k-\sigma) + \\ \mathbf{x}^\top(k-\sigma) Q_2 \mathbf{x}(k-\sigma) + \bar{\mathbf{u}}^\top(k) R_0 \bar{\mathbf{u}}(k) + \\ 2\bar{\mathbf{u}}^\top(k) R_1 \bar{M}(k) \bar{\mathbf{u}}(k) + \\ \bar{\mathbf{u}}^\top(k) \bar{M}^\top(k) R_2 \bar{M}(k) \bar{\mathbf{u}}(k) + P^{(i)}(\mathbf{x}(k+1)) \end{aligned} \quad (18)$$

where $P^{(0)}(\mathbf{x}(k)) = V^{(0)}(\mathbf{x}(k)) = 0$ and $\bar{\mathbf{u}}(k-\tau) = \bar{M}(k) \bar{\mathbf{u}}(k)$. $V^{(i)}(\mathbf{x}(k))$ is updated by (15). Thus, we can obtain

$$\begin{aligned} P^{(i+1)}(\mathbf{x}(k)) - P^{(i)}(\mathbf{x}(k)) &= P^{(i)}(\mathbf{x}(k+1)) - P^{(i-1)}(\mathbf{x}(k+1)) \\ &\vdots \\ &= P^{(1)}(\mathbf{x}(k+i)) - P^{(0)}(\mathbf{x}(k+i)) \end{aligned} \quad (19)$$

Because $P^{(0)}(\mathbf{x}(k+i)) = 0$, we have

$$\begin{aligned} P^{(i+1)}(\mathbf{x}(k)) &= P^{(1)}(\mathbf{x}(k+i)) + P^{(i)}(\mathbf{x}(k)) + \\ &\sum_{j=0}^i P^{(1)}(\mathbf{x}(k+j)) \end{aligned} \quad (20)$$

According to (18), (20) can be rewritten as

$$P^{(i+1)}(\mathbf{x}(k)) = \sum_{j=0}^i \Xi(k+j) \leq \sum_{j=0}^{\infty} \Xi(k+j) \quad (21)$$

where

$$\begin{aligned} \Xi(k+j) &= \mathbf{x}^\top(k+j) Q_0 \mathbf{x}(k+j) + \\ &2\mathbf{x}^\top(k+j) Q_1 \mathbf{x}(k+j-\sigma) + \\ &\mathbf{x}^\top(k+j-\sigma) Q_2 \mathbf{x}(k+j-\sigma) + \\ &\bar{\mathbf{u}}^\top(k+j) R_0 \bar{\mathbf{u}}(k+j) + \\ &2\bar{\mathbf{u}}^\top(k+j) R_1 \bar{M}(k+j) \bar{\mathbf{u}}(k+j) + \\ &\bar{\mathbf{u}}^\top(k+j) \bar{M}^\top(k+j) R_2 \bar{M}(k+j) \bar{\mathbf{u}}(k+j) \end{aligned}$$

Noting that the control input $\bar{\mathbf{u}}(k), k = 0, 1, \dots$ is an admissible control, we can obtain

$$P^{(i+1)}(\mathbf{x}(k)) \leq \sum_{j=0}^{\infty} P^{(1)}(\mathbf{x}(k+j)) \leq Y, \quad \forall i \quad (22)$$

From Lemma 1, we have

$$V^{(i+1)}(\mathbf{x}(k)) \leq P^{(i+1)}(\mathbf{x}(k)) \leq Y, \quad \forall i \quad (23)$$

□

With Lemma 1 and Theorem 1, the following main theorem can be derived.

Theorem 2. Define the performance index function $V^{(i)}(\mathbf{x}(k))$ as (15), with $V^{(0)}(\mathbf{x}(k)) = 0$. If $\mathbf{x}(k)$ for system

(1) is controllable, then $V^{(i)}(\mathbf{x}(k))$ is a nondecreasing sequence that is $V^{(i)}(\mathbf{x}(k)) \leq V^{(i+1)}(\mathbf{x}(k))$ and $V^{(i)}(\mathbf{x}(k))$ is convergent as $i \rightarrow \infty$.

Proof. For the convenience of analysis, define a new sequence $\Phi^{(i)}(\mathbf{x}(k))$ as follows:

$$\begin{aligned} \Phi^{(i+1)}(\mathbf{x}(k)) &= \mathbf{x}^T(k)Q_0\mathbf{x}(k) + 2\mathbf{x}^T(k)Q_1\mathbf{x}(k - \sigma) + \\ &\mathbf{x}^T(k - \sigma)Q_2\mathbf{x}(k - \sigma) + \mathbf{u}^{(i+1)\top}(k)R_0\mathbf{u}^{(i+1)}(k) + \\ &2\mathbf{u}^{(i+1)\top}(k)R_1M^{(i)}(k)\mathbf{u}^{(i+1)}(k) + \\ &\mathbf{u}^{(i+1)\top}(k)M^{(i)\top}(k)R_2M^{(i)}(k)\mathbf{u}^{(i+1)}(k) + \\ &\Phi^{(i)}(\mathbf{x}(k + 1)) \end{aligned} \quad (24)$$

with $\mathbf{u}^{(i)}(k)$ obtained by (14) and $\Phi_0(\mathbf{x}(k)) = V_0(\mathbf{x}(k)) = 0$. $V^{(i)}(\mathbf{x}(k))$ is updated by (15).

In the following part, we prove $\Phi^{(i)}(\mathbf{x}(k)) \leq V^{(i+1)}(\mathbf{x}(k))$ by mathematical induction.

First, we prove it holds for $i = 0$. Note that

$$\begin{aligned} V^{(1)}(\mathbf{x}(k)) - \Phi^{(0)}(\mathbf{x}(k)) &= \\ \mathbf{x}^T(k)Q_0\mathbf{x}(k) + 2\mathbf{x}^T(k)Q_1\mathbf{x}(k - \sigma) + \\ \mathbf{x}^T(k - \sigma)Q_2\mathbf{x}(k - \sigma) &\geq 0 \end{aligned} \quad (25)$$

Thus for $i = 0$, we can get

$$V^{(1)}(\mathbf{x}(k)) \geq \Phi^{(0)}(\mathbf{x}(k)) \quad (26)$$

Second, we assume it holds for $i - 1$, i.e., $V^{(i)}(\mathbf{x}(k)) - \Phi^{(i-1)}(\mathbf{x}(k)) \geq 0, \forall \mathbf{x}(k)$. Then, for i , from (15) and (24), we can obtain

$$\begin{aligned} V^{(i+1)}(\mathbf{x}(k)) - \Phi^{(i)}(\mathbf{x}(k)) &= \\ V^{(i)}(\mathbf{x}(k + 1)) - \Phi^{(i-1)}(\mathbf{x}(k + 1)) &\geq 0 \end{aligned} \quad (27)$$

i.e.,

$$\Phi^{(i)}(\mathbf{x}(k)) \leq V^{(i+1)}(\mathbf{x}(k)) \quad (28)$$

Therefore, the mathematical induction proof is completed.

Moreover, from Lemma 1, we know that $V^{(i)}(\mathbf{x}(k)) \leq \Phi^{(i)}(\mathbf{x}(k))$ and therefore we can obtain

$$V^{(i)}(\mathbf{x}(k)) \leq \Phi^{(i)}(\mathbf{x}(k)) \leq V^{(i+1)}(\mathbf{x}(k)) \quad (29)$$

which proves that $V^{(i)}(\mathbf{x}(k))$ is a nondecreasing sequence bounded by (23). Hence, we conclude that $V^{(i)}(\mathbf{x}(k))$ is a nondecreasing convergent sequence as $i \rightarrow \infty$. \square

We note the obvious corollary.

Corollary 1. If Theorem 2 holds, then the delay matrix function $M^{(i)}(k)$ is a convergent sequence, as $i \rightarrow \infty$.

According to Corollary 1, we define

$$M^{(\infty)}(k) = \lim_{i \rightarrow \infty} M^{(i)}(k) \quad (30)$$

Next, we will prove that the performance index function sequence $V^{(i)}(\mathbf{x}(k))$ converges to $V^*(\mathbf{x}(k))$ as $i \rightarrow \infty$. As $V^{(i)}(\mathbf{x}(k))$ is a convergent sequence as $i \rightarrow \infty$, we define

$$V^{(\infty)}(\mathbf{x}(k)) = \lim_{i \rightarrow \infty} V^{(i)}(\mathbf{x}(k)) \quad (31)$$

Let $\bar{\mathbf{u}}_l$ be the l -th admissible control. Similar to the proof of Theorem 1, we can construct the performance index function sequence $P_l^{(i)}(\mathbf{x})$ as follows:

$$\begin{aligned} P_l^{(i+1)}(\mathbf{x}(k)) &= \mathbf{x}^T(k)Q_0\mathbf{x}(k) + 2\mathbf{x}^T(k)Q_1\mathbf{x}(k - \sigma) + \\ &\mathbf{x}^T(k - \sigma)Q_2\mathbf{x}(k - \sigma) + \\ &\bar{\mathbf{u}}_l^T(k)R_0\bar{\mathbf{u}}_l(k) + 2\bar{\mathbf{u}}_l(k)R_1M^{(\infty)}(k)\bar{\mathbf{u}}_l(k) + \\ &\bar{\mathbf{u}}_l(k)M^{(\infty)\top}(k)R_2M^{(\infty)}(k)\bar{\mathbf{u}}_l(k) + \\ &P_l^{(i)}(\mathbf{x}(k + 1)) \end{aligned} \quad (32)$$

with $P_l^{(0)}(\cdot) = 0$ and $\bar{\mathbf{u}}_l(k) = M^{(\infty)}(k)\bar{\mathbf{u}}_l(k - \tau)$. According to Theorem 1, we have

$$\begin{aligned} P_l^{(i+1)}(\mathbf{x}(k)) &= \sum_{j=0}^i \left(\mathbf{x}^T(k + j)Q_0\mathbf{x}(k + j) + \right. \\ &2\mathbf{x}^T(k + j)Q_1\mathbf{x}(k + j - \sigma) + \\ &\mathbf{x}^T(k + j - \sigma)Q_2\mathbf{x}(k + j - \sigma) + \\ &\bar{\mathbf{u}}_l^T(k + j)R_0\bar{\mathbf{u}}_l(k + j) + \\ &2\bar{\mathbf{u}}_l^T(k + j)R_1M^{(\infty)}(k + j)\bar{\mathbf{u}}_l(k + j) + \\ &\bar{\mathbf{u}}_l^T(k + j)M^{(\infty)\top}(k + j)R_2 \times \\ &\left. M^{(\infty)}(k + j)\bar{\mathbf{u}}_l(k + j) \right) \end{aligned} \quad (33)$$

Let

$$P_l^{(\infty)}(\mathbf{x}(k)) = \lim_{i \rightarrow \infty} P_l^{(i+1)}(\mathbf{x}(k)) \quad (34)$$

So, we have

$$P_l^{(i)}(\mathbf{x}(k)) \leq P_l^{(\infty)}(\mathbf{x}(k)) \quad (35)$$

Theorem 3. Define $P_l^{(\infty)}(\mathbf{x}(k))$ as in (34), and define the performance index function $V^{(i)}(\mathbf{x}(k))$ as in (15) with $V^{(0)}(\cdot) = 0$. For any state vector $\mathbf{x}(k)$, define $V^*(\mathbf{x}(k)) = \min_l \{P_l^{(\infty)}(\mathbf{x}(k))\}$ starting from $\mathbf{x}(k)$ for all admissible control sequences. Then, we can conclude that $V^*(\mathbf{x}(k))$ is the limit of the performance index function $V^{(i)}(\mathbf{x}(k))$ as $i \rightarrow \infty$.

Proof. For any l , there exists an upper bound Y_l , such that

$$P_l^{(i+1)}(\mathbf{x}(k)) \leq P_l^{(\infty)}(\mathbf{x}(k)) \leq Y_l \quad (36)$$

According to (23), for $\forall l$, we have

$$V^{(\infty)}(\mathbf{x}(k)) \leq P_l^{(\infty)}(\mathbf{x}(k)) \leq Y_l \quad (37)$$

Since $V^*(\mathbf{x}(k)) = \min_l \{P_l^{(\infty)}(\mathbf{x}(k))\}$, for any $\epsilon > 0$, there exists an admissible control $\bar{\mathbf{u}}_K$, where K is a nonnegative number such that the associated performance index function satisfies $P_K^{(\infty)}(\mathbf{x}(k)) \leq V^*(\mathbf{x}(k)) + \epsilon$. According to (23), we have $V^{(\infty)}(\mathbf{x}(k)) \leq P_l^{(\infty)}(\mathbf{x}(k))$ for any l . Thus, we can obtain $V^{(\infty)}(\mathbf{x}(k)) \leq P_K^{(\infty)}(\mathbf{x}(k)) \leq V^*(\mathbf{x}(k)) + \epsilon$. Noting that ϵ is chosen arbitrarily, we have

$$V^{(\infty)}(\mathbf{x}(k)) \leq V^*(\mathbf{x}(k)) \quad (38)$$

On the other hand, since $V^{(i)}(\mathbf{x}(k))$ is bounded for $\forall i$, according to the definition of admissible control, the control sequence associated with the performance index function $V^{(\infty)}(\mathbf{x}(k))$ must be an admissible control, i.e., there

exists an admissible control $\bar{\mathbf{u}}_N^{(i)}$ such that $V^{(\infty)}(\mathbf{x}(k)) = P_N^{(\infty)}(\mathbf{x}(k))$. Combining with the definition $V^*(\mathbf{x}(k)) = \min_l \{P_l^{(\infty)}(\mathbf{x}(k))\}$, we can obtain

$$V^{(\infty)}(\mathbf{x}(k)) \geq V^*(\mathbf{x}(k)) \quad (39)$$

Therefore, combining (38) and (39), we can conclude that

$$V^{(\infty)}(\mathbf{x}(k)) = \lim_{i \rightarrow \infty} V^{(i)}(\mathbf{x}(k)) = V^*(\mathbf{x}(k)) \quad (40)$$

namely, $V^*(\mathbf{x}(k))$ is the limit of the performance index function $V^{(i)}(\mathbf{x}(k))$, as $i \rightarrow \infty$. \square

Based on Theorem 3, we will prove that the performance index function $V^*(\mathbf{x}(k))$ satisfies the principle of optimality, which shows that $V^{(i)}(\mathbf{x}(k))$ can reach the optimum as $i \rightarrow \infty$.

Theorem 4. For any state vector $\mathbf{x}(k)$, the ‘‘optimal’’ performance index function $V^*(\mathbf{x}(k))$ satisfies $V^*(\mathbf{x}(k)) = \min_{\mathbf{u}(k)} \{\mathbf{x}^T(k)Q_0\mathbf{x}(k) + 2\mathbf{x}^T(k)Q_1\mathbf{x}(k - \sigma) + \mathbf{x}^T(k - \sigma)Q_2\mathbf{x}(k - \sigma) + \mathbf{u}^T(k)R_0\mathbf{u}(k) + 2\mathbf{u}^T(k)R_1M(k)\mathbf{u}(k) + \mathbf{u}^T(k)M(k)R_2M(k)\mathbf{u}(k) + V^*(\mathbf{x}(k + 1))\}$, where $\mathbf{u}(k - \tau) = M(k)\mathbf{u}(k)$.

Proof. For any $\mathbf{u}(k)$ and i , based on Bellman’s optimality principle, we have

$$V^{(i)}(\mathbf{x}(k)) \leq \Upsilon^{(i-1)} + V^{(i-1)}(\mathbf{x}(k + 1)) \quad (41)$$

where

$$\begin{aligned} \Upsilon^{(i-1)} = & \mathbf{x}^T(k)Q_0\mathbf{x}(k) + 2\mathbf{x}^T(k)Q_1\mathbf{x}(k - \sigma) + \\ & \mathbf{x}^T(k - \sigma)Q_2\mathbf{x}(k - \sigma) + \\ & \mathbf{u}^T(k)R_0\mathbf{u}(k) + 2\mathbf{u}^T(k)R_1M^{(i-1)}(k)\mathbf{u}(k) + \\ & \mathbf{u}^T(k)M^{(i-1)T}(k)R_2M^{(i-1)}(k)\mathbf{u}(k) \end{aligned}$$

As $V^{(i)}(\mathbf{x}(k)) \leq V^{(i+1)}(\mathbf{x}(k)) \leq V^{(\infty)}(\mathbf{x}(k))$ and $V^{(\infty)}(\mathbf{x}(k)) = V^*(\mathbf{x}(k))$, we can obtain

$$V^{(i)}(\mathbf{x}(k)) \leq \Upsilon^{(i-1)} + V^*(\mathbf{x}(k + 1)) \quad (42)$$

If $i \rightarrow \infty$, then we have

$$V^*(\mathbf{x}(k)) \leq \Upsilon^{(\infty)} + V^*(\mathbf{x}(k + 1)) \quad (43)$$

Since $\mathbf{u}(k)$ in the above equation is chosen arbitrarily, the following equation holds

$$V^*(\mathbf{x}(k)) \leq \min_{\mathbf{u}(k)} \left\{ \Upsilon^{(\infty)} + V^*(\mathbf{x}(k + 1)) \right\} \quad (44)$$

On the other hand, for any i , the performance index function satisfies

$$V^{(i)}(\mathbf{x}(k)) = \Omega^{(i-1)} + V^{(i-1)}(\mathbf{x}(k + 1)) \quad (45)$$

where

$$\begin{aligned} \Omega^{(i-1)} = & \mathbf{x}^T(k)Q_0\mathbf{x}(k) + 2\mathbf{x}^T(k)Q_1\mathbf{x}(k - \sigma) + \\ & \mathbf{x}^T(k - \sigma)Q_2\mathbf{x}(k - \sigma) + \mathbf{u}^{(i-1)T}(k)R_0\mathbf{u}^{(i-1)}(k) + \\ & 2\mathbf{u}^{(i)T}(k)R_1M^{(i-2)}(k)\mathbf{u}^{(i-1)T}(k) + \\ & \mathbf{u}^{(i-1)T}(k)M^{(i-2)T}(k)R_2M^{(i-2)}(k)\mathbf{u}^{(i-1)T}(k) \end{aligned}$$

Combining with $V^{(i)}(\mathbf{x}(k)) \leq V^*(\mathbf{x}(k))$, $\forall i$, we have

$$V^*(\mathbf{x}(k)) \geq \Omega^{(i-1)} + V^{(i-1)}(\mathbf{x}(k + 1)) \quad (46)$$

Let $i \rightarrow \infty$, then

$$\begin{aligned} V^*(\mathbf{x}(k)) \geq & \lim_{i \rightarrow \infty} \left\{ \Omega^{(i-1)} + V^{(i-1)}(\mathbf{x}(k + 1)) \right\} \geq \\ & \min_{\mathbf{u}(k)} \left\{ \Omega^{(\infty)} + V^*(\mathbf{x}(k + 1)) \right\} \end{aligned} \quad (47)$$

Combining (44) with (47), we have

$$\begin{aligned} V^*(\mathbf{x}(k)) = & \min_{\mathbf{u}(k)} \left\{ \Omega^{(\infty)} + V^*(\mathbf{x}(k + 1)) \right\} = \\ & \mathbf{x}^T(k)Q_0\mathbf{x}(k) + 2\mathbf{x}^T(k)Q_1\mathbf{x}(k - \sigma) + \\ & \mathbf{x}^T(k - \sigma)Q_2\mathbf{x}(k - \sigma) + \mathbf{u}^{*T}(k)R_0\mathbf{u}^*(k) + \\ & 2\mathbf{u}^{*T}(k)R_1M^{(\infty)}(k)\mathbf{u}^*(k) + \\ & \mathbf{u}^{*T}(k)M^{(\infty)T}(k)R_2M^{(\infty)}(k)\mathbf{u}^*(k) + \\ & V^*(\mathbf{x}(k + 1)) \end{aligned} \quad (48)$$

Thus, we have that $\mathbf{u}^{(i)}(k) \rightarrow \mathbf{u}^*(k)$ as $i \rightarrow \infty$ so does $\mathbf{u}^{(i)}(k - \tau)$. On the other hand, we also have $M^{(i)}(k) \rightarrow M^{(\infty)}(k)$ and $\mathbf{u}^{(i)}(k - \tau) = M^{(i-1)}(k)\mathbf{u}^{(i)}(k)$. Letting $i \rightarrow \infty$, we get

$$\mathbf{u}^*(k - \tau) = M^{(\infty)}(k)\mathbf{u}^*(k) \quad (49)$$

Therefore, we have $M^{(\infty)}(k) = M^*(k)$ and (48) can be written as

$$\begin{aligned} V^*(\mathbf{x}(k)) = & \mathbf{x}^T(k)Q_0\mathbf{x}(k) + 2\mathbf{x}^T(k)Q_1\mathbf{x}(k - \sigma) + \\ & \mathbf{x}^T(k - \sigma)Q_2\mathbf{x}(k - \sigma) + \mathbf{u}^{*T}(k)R_0\mathbf{u}^*(k) + \\ & 2\mathbf{u}^{*T}(k)R_1M^*(k)\mathbf{u}^*(k) + \\ & \mathbf{u}^{*T}(k)M^{*T}(k)R_2M^*(k)\mathbf{u}^*(k) + \\ & V^*(\mathbf{x}(k + 1)) \end{aligned} \quad (50)$$

where $\mathbf{u}^*(k - \tau) = M^*(k)\mathbf{u}^*(k)$. \square

Therefore, we can conclude that the performance index function $V^{(i)}(\mathbf{x}(k))$ converges to the optimum $V^*(\mathbf{x}(k))$ as $i \rightarrow \infty$.

2.3 Implementation of iterative ADP algorithm

Given the above preparation, we may formulate the desired iterative ADP approach to nonlinear systems with delays as follows.

Step 1. Give initial state $\mathbf{x}(s) = \phi(s)$, $s = -\sigma, -\sigma + 1, \dots, 0$, initial control $\mathbf{u}(\rho)$, $\rho = 0, 1, \dots, k - 1$; give i_{\max} and computation accuracy ε .

Step 2. Set the iterative step $i = 0$, $M^{(0)}(k) = I$, and $V^{(0)}(\cdot) = 0$.

Step 3. Compute $\mathbf{u}^{(0)}(k)$ by (12) and the performance index function $V^{(1)}(\mathbf{x}(k))$ by (13).

Step 4. For the iterative step $i \geq 1$, compute $\mathbf{u}^{(i)}(k)$ by (14).

Step 5. Compute the performance index function $V^{(i)}(\mathbf{x}(k))$ by (15).

Step 6. If

$$\left[V^{(i)}(\mathbf{x}(k)) - V^{(i-1)}(\mathbf{x}(k)) \right]^2 < \varepsilon \quad (51)$$

go to Step 9; otherwise, go to Step 7.

Step 7. If $i > i_{\max}$, go to Step 9; otherwise, compute $M^{(i)}(k)$ by

$$M^{(i)}(k) = \mathbf{u}^{(i)}(k - \tau) \mathbf{u}^{(i)T}(k) \left(\mathbf{u}^{(i)}(k) \mathbf{u}^{(i)T}(k) \right)^{-1} \quad (52)$$

Step 8. Set $i = i + 1$ and go to Step 4.

Step 9. Stop.

In (52) of the above algorithm, the term $(\mathbf{u}^{(i)}(k)\mathbf{u}^{(i)\top}(k))^{-1}$ can be obtained by the Moore-Penrose pseudoinverse technique to compute the delay matrix function $M^{(i)}(k)$. There are two other methods to compute $M^{(i)}(k)$. One choice is to introduce a small zero-mean Gaussian noise with variances γ^2 denoted by $\delta(0, \gamma^2)$ into the control $\mathbf{u}(k - \tau)$ ^[18]. The other choice is to use a neural network to approximate delay matrix function $M^{(i)}(k)$. In this paper, we use the neural network approximation method and the details will be shown in the next section.

3 Neural network implementation

In the case of linear systems, the performance index function is quadratic and the control policy is linear. In the nonlinear case, this is not necessarily true and therefore we use neural networks to approximate $\mathbf{u}^{(i)}(k)$ and $V^{(i)}(\mathbf{x}(k))$.

Assume the number of hidden layer neurons is denoted by l , that the weight matrix between the input layer and hidden layer is denoted by V , and that the weight matrix between the hidden layer and output layer is denoted by W . Then the output of three-layer neural network (NN) is represented by

$$\hat{F}(\mathbf{X}, V, W) = W^T \sigma(V^T \mathbf{X}) \quad (53)$$

where $\sigma(V^T \mathbf{X}) \in \mathbf{R}^l$, $[\sigma(z)]_i = \frac{e^{z_i} - e^{-z_i}}{e^{z_i} + e^{-z_i}}$, $i = 1, \dots, l$, are the activation functions.

The NN estimation error can be expressed by

$$F(\mathbf{X}) = F(\mathbf{X}, V^*, W^*) + \varepsilon(\mathbf{X}) \quad (54)$$

where V^* and W^* are the ideal weight parameters, and $\varepsilon(\mathbf{X})$ is the reconstruction error.

Here, there are four neural networks, which are critic network, model network, action network, and delay matrix function network (M network), respectively. All the neural networks are chosen as three-layer feedforward network. The whole structure diagram is shown in Fig. 1. The utility term in the figure denotes $\mathbf{x}^T(k)Q_0\mathbf{x}(k) + 2\mathbf{x}^T(k)Q_1\mathbf{x}(k - \sigma) + \mathbf{x}^T(k - \sigma)Q_2\mathbf{x}(k - \sigma) + \mathbf{u}^T(k)R_0\mathbf{u}(k) + 2\mathbf{u}^T(k)R_1\mathbf{u}(k - \tau) + \mathbf{u}^T(k - \tau)R_2\mathbf{u}(k - \tau)$.

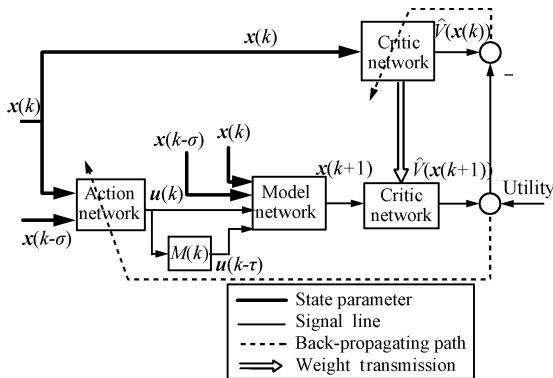


Fig. 1 The structure diagram of the algorithm

3.1 The model network

The model network is to approximate the system dynamic and it should be trained before the implementation

of the iterative ADP algorithm. The update rule of the model network is adopted as the gradient decent method. The training process is simple and general. The details can be seen in [13, 19] and it is omitted here.

After the model network is trained, its weights are kept unchanged.

3.2 The M network

The M network is to approximate the delay matrix function $M(k)$. The output of the M network is denoted as

$$\hat{\mathbf{u}}(k - \tau) = W_M^T \sigma(V_M^T \mathbf{u}(k)) \quad (55)$$

We define the error function of the model network as

$$\mathbf{e}_M(k) = \hat{\mathbf{u}}(k - \tau) - \mathbf{u}(k - \tau) \quad (56)$$

Define the performance error measure as

$$E_M(k) = \frac{1}{2} \mathbf{e}_M^T(k) \mathbf{e}_M(k) \quad (57)$$

Then, the gradient-based weight update rule for the critic network can be described by

$$w_M(k+1) = w_M(k) + \Delta w_M(k) \quad (58)$$

$$\Delta w_M(k) = \alpha_M \left[-\frac{\partial E_M(k)}{\partial w_M(k)} \right] \quad (59)$$

where α_M is the learning rate of the M network.

3.3 The critic network

The critic network is used to approximate the performance index function $V^{(i)}(\mathbf{x}(k))$. The output of the critic network is denoted as

$$\hat{V}^{(i)}(\mathbf{x}(k)) = W_{ci}^T \sigma(V_{ci}^T \mathbf{z}(k)) \quad (60)$$

The target function can be written as

$$V^{(i+1)}(\mathbf{x}(k)) = \Gamma^{(i)} + \hat{V}^{(i)}(\mathbf{x}(k+1)) \quad (61)$$

Then, we define the error function for the critic network as

$$e_{ci}(k) = \hat{V}^{(i+1)}(\mathbf{x}(k)) - V^{(i+1)}(\mathbf{x}(k)) \quad (62)$$

And the objective function to be minimized in the critic network is

$$E_{ci}(k) = \frac{1}{2} e_{ci}^2(k) \quad (63)$$

So the gradient-based weight update rule for the critic network is given by

$$w_{c(i+1)}(k) = w_{ci}(k) + \Delta w_{ci}(k) \quad (64)$$

$$\Delta w_{ci}(k) = \alpha_c \left[-\frac{\partial E_{ci}(k)}{\partial w_{ci}(k)} \right] \quad (65)$$

$$\frac{\partial E_{ci}(k)}{\partial w_{ci}(k)} = \frac{\partial E_{ci}(k)}{\partial \hat{V}^{(i)}(\mathbf{x}(k))} \frac{\partial \hat{V}^{(i)}(\mathbf{x}(k))}{\partial w_{ci}(k)} \quad (66)$$

where $\alpha_c > 0$ is the learning rate of critic network and $w_c(k)$ is the weight vector in the critic network.

3.4 The action network

In the action network, the state $\mathbf{x}(k)$ is used as input to create the optimal control as the output of the network. The output can be formulated as

$$\hat{\mathbf{u}}^{(i)}(k) = W_{ai}^T \sigma(V_{ai}^T \mathbf{x}(k)) \quad (67)$$

The target of the output of the action network is given by (14). So, we can define the output error of the action network as

$$\mathbf{e}_{ai}(k) = \hat{\mathbf{u}}^{(i)}(k) - \mathbf{u}^{(i)}(k) \quad (68)$$

where $\mathbf{u}^{(i)}(k)$ is the target function which can be described by

$$\begin{aligned} \mathbf{u}^{(i)}(k) = & -\frac{1}{2} \left(R_0 + 2R_1 M^{(i-1)}(k) + \right. \\ & \left. M^{(i-1)T}(k) R_2 M^{(i-1)}(k) \right)^{-1} \times (g_0(\mathbf{x}(k), \mathbf{x}(k-\sigma)) + \\ & g_1(\mathbf{x}(k), \mathbf{x}(k-\sigma)) M^{(i-1)}(k))^T \frac{\partial \hat{V}^{(i)}(\mathbf{x}(k+1))}{\partial \mathbf{x}(k+1)} \end{aligned}$$

As $\mathbf{u}^{(i)}(k-\tau) = M^{(i-1)}(k) \mathbf{u}^{(i)}(k)$, we have $\frac{\partial \mathbf{u}^{(i)}(k-\tau)}{\partial \mathbf{u}^{(i)}(k)} = M^{(i-1)}(k)$. Then, according to (55), $M^{(i-1)}(k)$ can be expressed as

$$M_{ij}^{(i-1)}(k) = V_{Mi}^T \left[1 - \left(\sigma(V_M^T \mathbf{u}(k)) \right)_i \right]^2 W_{Mj} \quad (69)$$

for $i, j = 1, 2, \dots, m$. $M_{ij}^{(i-1)}(k)$ denotes the element at row i , column j of matrix $M^{(i-1)}(k)$; V_{Mi} and W_{Mj} mean the column i and column j of the weight matrices V_M and W_M , respectively; $(\sigma(V_M^T \mathbf{u}(k)))_i$ is the i th element of the vector $\sigma(V_M^T \mathbf{u}(k))$.

The weights in the action network are updated to minimize the following performance error measure:

$$E_{ai}(k) = \frac{1}{2} \mathbf{e}_{ai}^T(k) \mathbf{e}_{ai}(k) \quad (70)$$

The weights updating algorithm is similar to the one for the critic network. By the gradient descent rule, we can obtain

$$w_{a(i+1)}(k) = w_{ai}(k) + \Delta w_{ai}(k) \quad (71)$$

$$\Delta w_{ai}(k) = \beta_a \left[-\frac{\partial E_{ai}(k)}{\partial w_{ai}(k)} \right] \quad (72)$$

$$\frac{\partial E_{ai}(k)}{\partial w_{ai}(k)} = \frac{\partial E_{ai}(k)}{\partial e_{ai}(k)} \frac{\partial e_{ai}(k)}{\partial \mathbf{u}^{(i)}(k)} \frac{\partial \mathbf{u}^{(i)}(k)}{\partial w_{ai}(k)} \quad (73)$$

where $\beta_a > 0$ is the learning rate of then action network.

4 Simulation

In this section, two examples are provided to demonstrate the effectiveness of the control scheme proposed in this paper.

4.1 Optimal control for state delayed system

For the first example, the nonlinear system is a modification of Example 1 in [13], which introduces state delays into the system.

Consider the following affine nonlinear system:

$$\mathbf{x}(k+1) = f(\mathbf{x}(k), \mathbf{x}(k-\sigma)) + g(\mathbf{x}(k), \mathbf{x}(k-\sigma)) \mathbf{u}(k) \quad (74)$$

where $\mathbf{x}(k) = [\mathbf{x}_1(k) \ \mathbf{x}_2(k)]^T$, $\mathbf{u}(k) = [u_1(k) \ u_2(k)]^T$, and $f(\mathbf{x}(k), \mathbf{x}(k-\sigma)) = \begin{bmatrix} \mathbf{x}_1(k) \exp(\mathbf{x}_2^3(k)) \mathbf{x}_2(k-2) \\ \mathbf{x}_2^3(k) \mathbf{x}_1(k-2) \end{bmatrix}$, $g(\mathbf{x}(k), \mathbf{x}(k-\sigma)) = \begin{bmatrix} -0.2 & 0 \\ 0 & -0.2 \end{bmatrix}$. The time delay in the

state is $\sigma = 2$ and the initial condition is $\mathbf{x}(k) = [1 \ -1]^T$ for $-2 \leq k \leq 0$. The performance index function is defined as (2), where $Q_0 = Q_2 = R_0 = I$ and $Q_1 = R_1 = R_2 = 0$.

We implement the algorithm at the time instant $k = 5$. We choose three-layer neural networks as the critic network, the action network, and the model network with the structures 4-10-2, 2-10-1, and 6-10-2, respectively. The initial weights of the action network, critic network, and model network are all set to be random in $[-0.5, 0.5]$. It should be mentioned that the model network should be trained first. For the given initial state, we train the model network for 3000 steps under the learning rate $\alpha_m = 0.05$. After the training of the model network, the weights keep unchanged. Then, the critic network and the action network are trained for 3000 steps so that the given accuracy $\varepsilon = 10^{-6}$ is reached. In the training process, the learning rate $\beta_a = \alpha_c = 0.05$. The convergence curve of the performance index function is shown in Fig. 2. Then, we apply the optimal control to the system for $T_f = 30$ time steps and obtain the following results. The state trajectories are given as Fig. 3 and the corresponding control curves are given as Fig. 4.

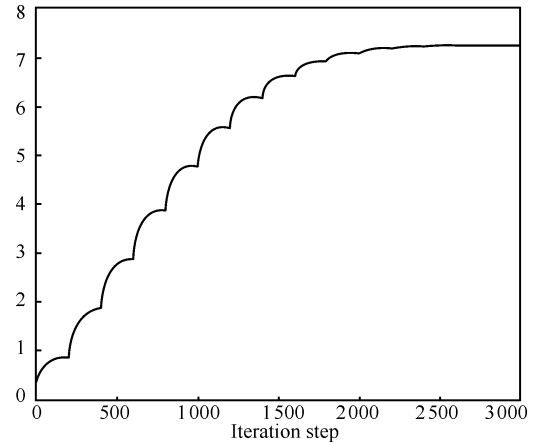


Fig. 2 The convergence of performance index function

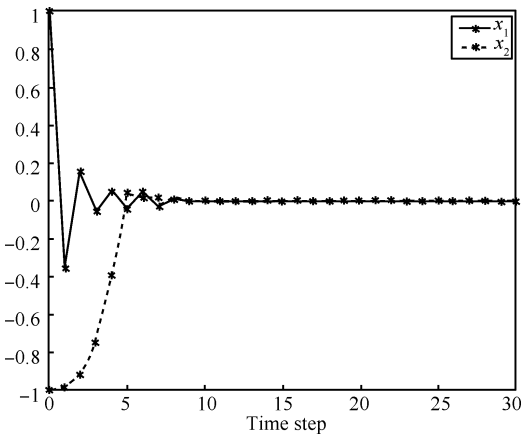


Fig. 3 The state variable trajectories

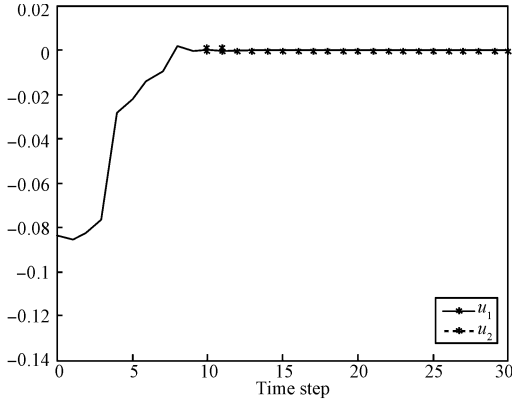


Fig. 4 The optimal control trajectories

4.2 Optimal control for nonlinear system with state and control delays

For the second example, the control time delay is added into the system given in the first example and the system becomes

$$\mathbf{x}(k+1) = f(\mathbf{x}(k), \mathbf{x}(k-\sigma)) + g_0(\mathbf{x}(k), \mathbf{x}(k-\sigma))\mathbf{u}(k) + g_1(\mathbf{x}(k), \mathbf{x}(k-\sigma))\mathbf{u}(k-\tau) \quad (75)$$

where $\mathbf{x}(k) = [\mathbf{x}_1(k) \ \mathbf{x}_2(k)]^T$, $\mathbf{u}(k) = [u_1(k) \ u_2(k)]^T$, and $f(\mathbf{x}(k), \mathbf{x}(k-\sigma))$ is the same as the first example, $g_0(\mathbf{x}(k), \mathbf{x}(k-\sigma)) = g_1(\mathbf{x}(k), \mathbf{x}(k-\sigma)) = \begin{bmatrix} -0.2 & 0 \\ 0 & -0.2 \end{bmatrix}$. The state time delay $\sigma = 2$ and the control time delay $\tau = 1$. The initial condition is $\mathbf{x}(k) = [-1 \ -1]^T$ and $\mathbf{u}(k) = 0$ for $-2 \leq k \leq 0$. The performance index function is defined as (2), where $Q_0 = Q_2 = R_0 = R_2 = I$ and $Q_1 = R_1 = 0$.

We also implement the algorithm at the time instant $k = 5$. We choose three-layer neural networks as the critic network, the action network, the model network, and the M network with the structures 4-10-2, 2-10-1, 8-10-2, and 2-8-2, respectively. All the other parameters are set the same as the first example. The initial weights of action network, critic network, model network, and the M network are all set to be random in $[-0.5, 0.5]$. For the given initial state, we train the model network for 4000 steps. After the training of the model network, the weights keep unchanged. Then, the critic network, the action network, and the M network are trained for 3000 steps to reach the given accuracy $\varepsilon = 10^{-6}$. The convergence curve of the performance index function is shown in Fig. 5. Then, we apply the optimal control to the system for $T_f = 30$ time steps and obtain the following results. The state trajectories are given as Fig. 6 and the corresponding control curves are given as Fig. 7.

5 Conclusion

In this paper, we proposed an effective algorithm to find the optimal infinite-time controller for a class of discrete-time nonlinear systems with time delays in state and control variables. By introducing a delay matrix function, the explicit expression of the optimal control was obtained. Then, the iterative ADP algorithm was implemented to deal with the time delay problem with rigorous convergence analysis. Four neural networks were used as parametric structures to approximate the performance index function, compute the optimal control policy, model the unknown system, and solve delay matrix function, respectively, i.e., the critic network, the action network, the model

network, and the M network. The simulation studies have successfully demonstrated the outstanding performance of the proposed time-delay optimal control scheme for various discrete-time nonlinear systems.

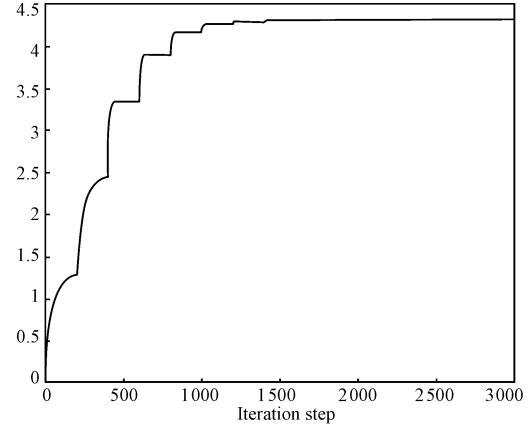


Fig. 5 The convergence of performance index function

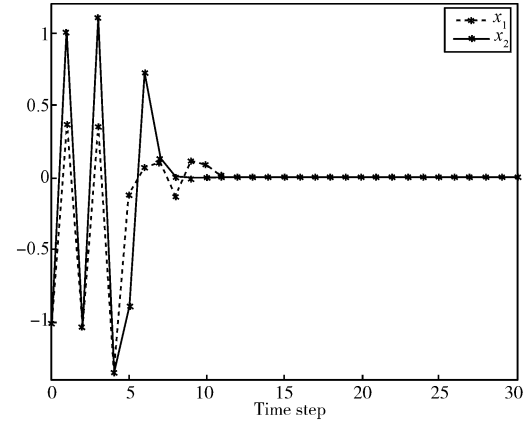


Fig. 6 The state variable trajectories

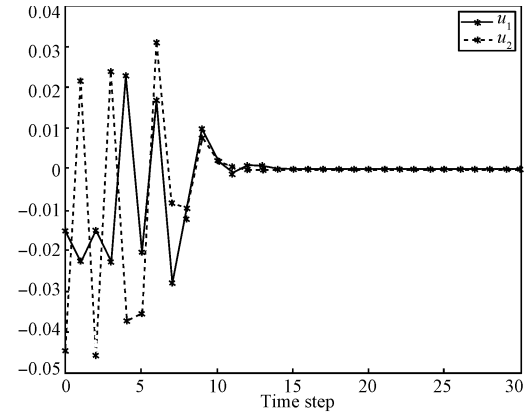


Fig. 7 The optimal control trajectories

Appendix

Lemma A1. If $\begin{bmatrix} R_0 & R_1 \\ R_1^T & R_2 \end{bmatrix}$ is a positive definite matrix, where $R_0, R_1, R_2 \in \mathbf{R}^{n \times n}$, then for any nonsingular matrix $M \in \mathbf{R}^{n \times n}$, $\begin{bmatrix} R_0 & R_1 M \\ M^T R_1^T & M^T R_2 M \end{bmatrix} > 0$.

Proof. Since $\begin{bmatrix} R_0 & R_1 \\ R_1^T & R_2 \end{bmatrix}$ is a positive definite matrix, accord-

ing to Schur complement^[20], we have

$$R_2 - R_1^T R_0^{-1} R_1 > 0 \quad (\text{A1})$$

As the matrix $M \in \mathbf{R}^{n \times n}$ is nonsingular, let M^{-1} denote the inverse matrix of M . Then, (A1) can be written as

$$R_2 - R_1^T M^T M^{-T} R_0^{-1} M^{-1} M R_1 > 0 \quad (\text{A2})$$

where $M^{-T} = (M^{-1})^T$. Again, using Schur complement, we can obtain

$$\begin{bmatrix} R_0 & R_1 M \\ M^T R_1^T & M^T R_2 M \end{bmatrix} > 0 \quad (\text{A3})$$

□

References

- Zhang H G, Yang D D, Chai T Y. Guaranteed cost networked control for T-S fuzzy systems with time delay. *IEEE Transactions on Systems, Man, and Cybernetics, Part C: Applications and Reviews*, 2007, **37**(2): 160–172
- Kong Shu-Lan, Zhang Huan-Shui, Zhang Zhao-Sheng, Zhang Cheng-Hui. Joint predictive control of power and rate for wireless networks. *Acta Automatica Sinica*, 2007, **33**(7): 761–764
- Malek-Zavarei M, Jashmidi M. *Time-Delay Systems: Analysis, Optimization and Applications*. North-Holland: The Netherlands, 1987. 80–96
- Basin M, Rodriguez-Gonzalez J. Optimal control for linear systems with multiple time delays in control input. *IEEE Transactions on Automatic Control*, 2006, **51**(1): 91–97
- Richard J P. Time-delay systems: an overview of some recent advances and open problems. *Automatica*, 2003, **39**(10): 1667–1694
- Wu Zheng-Guang, Zhou Wu-Neng. Delay-dependent robust stabilization for uncertain singular systems with state delay. *Acta Automatica Sinica*, 2007, **33**(7): 714–718
- Zhang H G, Wang Y C, Liu D R. Delay-dependent guaranteed cost control for uncertain stochastic fuzzy systems with multiple time delays. *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics*, 2008, **38**(1): 126–140
- Hou Zeng-Guang, Wu Cang-Pu. A dynamic programming neural network for large-scale optimization problems. *Acta Automatica Sinica*, 1999, **25**(1): 45–51 (in Chinese)
- Bellman R E. *Dynamic Programming*. Princeton: Princeton University Press, 1957. 150–155
- Tang Hao, Yuan Ji-Bin, Lu Yang, Cheng Wen-Juan. Performance potential-based neuro-dynamic programming for SMDPs. *Acta Automatica Sinica*, 2005, **31**(4): 642–645
- Werbos P J. Approximate dynamic programming for real-time control and neural modeling. *Handbook of Intelligent Control: Neural, Fuzzy, and Adaptive Approaches*. New York: Van Nostrand Reinhold, 1992
- Prokhorov D V, Wunsch D C. Adaptive critic designs. *IEEE Transactions on Neural Networks*, 1997, **8**(5): 997–1007
- Zhang H G, Wei Q L, Luo Y H. A novel infinite-time optimal tracking control scheme for a class of discrete-time nonlinear systems via the greedy HDP iteration algorithm. *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics*, 2008, **38**(4): 937–942
- Liu D R, Javaherian H, Kovalenko O, Huang T. Adaptive critic learning techniques for engine torque and air-fuel ratio control. *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics*, 2008, **38**(4): 988–993
- Liu De-Rong. Approximate dynamic programming for self-learning control. *Acta Automatica Sinica*, 2005, **31**(1): 13–18
- Ray S, Venayagamoorthy G K, Chaudhuri B, Majumder R. Comparison of adaptive critic-based and classical wide-area controllers for power systems. *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics*, 2008, **38**(4): 1002–1007
- Al-Tamimi A, Lewis F L, Abu-Khalaf M. Discrete-time nonlinear HJB solution using approximate dynamic programming-convergence proof. *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics*, 2008, **38**(4): 943–949
- Al-Tamimi A, Abu-Khalaf M, Lewis F L. Adaptive critic designs for discrete-time zero-sum games with application to H_∞ control. *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics*, 2007, **37**(1): 240–247
- Si J, Wang Y T. On-line learning control by association and reinforcement. *IEEE Transactions on Neural Networks*, 2001, **12**(2): 264–276
- Balakrishnam V, Feron E, Ghaoui L E. *Linear Matrix Inequalities in System and Control Theory*. Philadelphia: Society for Industrial and Applied Mathematics, 1994



WEI Qing-Lai Received his B.Sc. degree in automation control and M.Sc. degree in control theory and control engineering and the Ph.D. degree in control theory and control engineering from Northeastern University, in 2002, 2005, and 2008, respectively. He is currently a postdoctoral fellow with the Key Laboratory of Complex Systems and Intelligence Science, Institute of Automation, Chinese Academy of Sciences. His research interest covers neural-networks-based control, non-linear control, adaptive dynamic programming, and their industrial application. Corresponding author of this paper. E-mail: qinglaiwei@gmail.com



ZHANG Hua-Guang Received his Ph.D. degree from Southeast University, in 1991. He is currently as a full professor. His research interest covers fuzzy system theory, fuzzy control, neural network-based control, adaptive control, chaotic control, complex industry process automation, electric power system automation, and motor driving system automation. E-mail: hg Zhang@ieee.org



LIU De-Rong Received his Ph.D. degree in electrical engineering from the University of Notre Dame, USA in 1994. In 1999, he joined the University of Illinois at Chicago, Chicago, as a full professor of electrical and computer engineering and of computer science. He is now a professor at the Institute of Automation, Chinese Academy of Sciences. His research interest covers intelligent control, neural networks, fuzzy systems, computational neuroscience, power systems, and automotive engine. E-mail: derongliu@gmail.com



ZHAO Yan Received his Ph.D. degree from Northeastern University in 2008. He is a lecturer in the Department of Automatic Control Engineering, Shenyang Institute of Engineering. His research interest covers intelligent control and applications. E-mail: zhaoyan@sie.edu.cn