

# Approximate optimal control for a class of nonlinear discrete-time systems with saturating actuators

Yanhong Luo, Huaguang Zhang\*

*School of Information Science and Engineering, Northeastern University, Shenyang 110004, China*

Received 22 December 2007; received in revised form 23 February 2008; accepted 12 March 2008

## Abstract

In this paper, we solve the approximate optimal control problem for a class of nonlinear discrete-time systems with saturating actuators via greedy iterative Heuristic Dynamic Programming (GI-HDP) algorithm. In order to deal with the saturating problem of actuators, a novel nonquadratic functional is developed. Based on the nonquadratic functional, the GI-HDP algorithm is introduced to obtain the optimal saturated controller with a rigorous convergence analysis. For facilitating the implementation of the iterative algorithm, three neural networks are used to approximate the value function, compute the optimal control policy and model the unknown plant, respectively. An example is given to demonstrate the validity of the proposed optimal control scheme.

© 2008 National Natural Science Foundation of China and Chinese Academy of Sciences. Published by Elsevier Limited and Science in China Press. All rights reserved.

*Keywords:* Saturating; GI-HDP algorithm; Nonquadratic functional; Convergence analysis; Neural networks

## 1. Introduction

Saturation, dead-zone, backlash, and hysteresis are the most common actuator nonlinearities in practical control system applications. Saturation nonlinearity is unavoidable in most actuators. Several methods for deriving control laws considering the saturation phenomena were found in Refs. [1–3]. However, most of these methods did not consider optimal control laws for general nonlinear discrete-time systems.

In recent years, in order to obtain the approximate optimal control law, the approximate dynamic programming (ADP) algorithm has been paid much attention by many researchers [4–15]. ADP combines adaptive critic design, reinforcement learning technique with dynamic programming. ADP approaches were classified into four main schemes: Heuristic Dynamic Programming (HDP), Dual

Heuristic Dynamic Programming (DHP), Action Dependent Heuristic Dynamic Programming (ADHDP), also known as Q-learning [5], and Action Dependent Dual Heuristic Dynamic Programming (ADDHP). In Ref. [16], two new ADP schemes known as Globalized-DHP (GDHP) and ADGDHP were developed. In Ref. [14], a greedy HDP iteration scheme was proposed to solve the optimal control problem for nonlinear discrete-time systems with known mathematical model.

Though the ADP algorithm has made great progress in the optimal control field, there are still some problems unsolved, such as if the actuator has saturating characteristic and how to find a constrained optimal control by ADP algorithm. In Refs. [17,18], the nonquadratic functional was used to confront the input constraint. Via the nonquadratic functional, the special HJB equation was formulated and its solution resulted in a smooth saturated controller. However it remains difficult to actually solve for the value function of the HJB equation. Therefore, in this study the constrained optimal control problem is solved by the framework of the Hamilton–Jacobi–Bellman

\* Corresponding author. Tel.: +86 024 83687762; fax: +86 024 83671498.

*E-mail address:* [hg\\_zhang@21cn.com](mailto:hg_zhang@21cn.com) (H. Zhang).

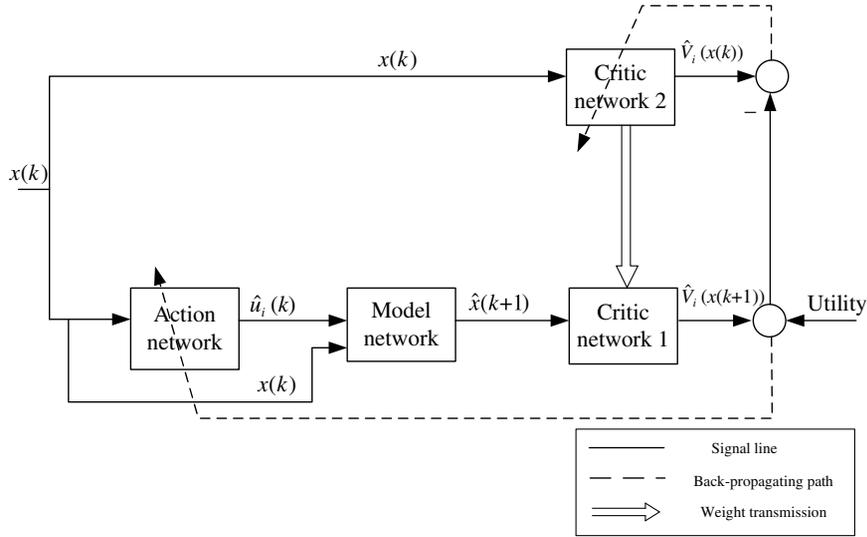


Fig. 1. The structure diagram of the GI-HDP algorithm.

(HJB) equation. After the special HJB equation is derived from the nonquadratic functional, a new ADP algorithm named greedy iterative HDP (GI-HDP) algorithm is proposed with rigorous convergence proof. Furthermore, in order to facilitate the implementation of the GI-HDP algorithm, two neural networks are used to approximate the value function and the corresponding optimal control policy, and a model network is introduced to approximate the unknown nonlinear plant.

**2. Discrete-time HJB equation for constrained nonlinear systems**

Consider a class of discrete-time nonlinear systems as follows:

$$x(k + 1) = f(x(k)) + g(x(k))u(k) \tag{1}$$

where  $x(k) \in \mathbb{R}^n$  is the state vector,  $f(\cdot)$  and  $g(\cdot)$  are differentiable in their argument with  $f(0) = 0, g(0) = 0$ . Assume that  $f + gu$  is Lipschitz continuous on a set  $\Omega$  in  $\mathbb{R}^n$  containing the origin, and that the system (1) is controllable in the sense that there exists a continuous control on  $\Omega$  that asymptotically stabilizes the system. And the control  $u(k) \in \Omega_u, \Omega_u = \{u(k) = [u_1(k), u_2(k), \dots, u_m(k)]^T \in \mathbb{R}^m : |u_i(k)| \leq \bar{u}_i, i = 1, \dots, m\}$ , here  $\bar{u}_i$  denotes the saturating bound for the  $i$ th actuator. Let  $\bar{U} \in \mathbb{R}^{m \times m}$  be the constant diagonal matrix described as  $\bar{U} = \text{diag}\{\bar{u}_1, \bar{u}_2, \dots, \bar{u}_m\}$ .

In this paper, we mainly discuss how to design an optimal state feedback controller for this class of constrained discrete-time systems. Therefore, it is desired to find the control policy  $u(x)$  which minimizes a generalized performance functional as follows:

$$J(x(k), u) = \sum_{i=k}^{\infty} \{x(i)^T Qx(i) + W(u(i))\} \tag{2}$$

where both  $W(u(i))$  and  $Q$  are positive definite.

For optimal control problem, the state feedback control  $u(x)$  must not only stabilize the system on  $\Omega$  but also guarantee that (2) is finite, i.e., admissible control, see Ref. [19]. From now on, we let  $V^*(x(k))$  denote the minimum value of the performance functional  $J(x(k), u)$ , which is called value function or the optimal cost function in the later parts.

**Definition 1.** A control  $u(x)$  is defined to be admissible control with respect to (2) on  $\Omega$  if  $u$  stabilizes (1) on  $\Omega, u(0) = 0$ , and for all  $x(0) \in \Omega, J(x(0), u)$  is finite.

According to Bellman optimality principle, we can obtain

$$\begin{aligned} V^*(x(k)) &= \min_{u(i)} \sum_{i=k}^{\infty} \{x(i)^T Qx(i) + W(u(i))\} \\ &= \min_{u(k)} \{x(k)^T Qx(k) + W(u(k)) + V^*(x(k + 1))\} \end{aligned} \tag{3}$$

For unconstrained control problem, a common choice for  $W(u(i))$  is  $W(u(i)) = u(i)^T R u(i)$ , where  $R \in \mathbb{R}^{m \times m}$  is positive definite. With the first order necessity condition, we compute the gradient of right-hand side of (3) with respect to  $u$  as

$$\begin{aligned} \frac{\partial V^*(x(k))}{\partial u(k)} &= \frac{\partial (x(k)^T Qx(k) + u(k)^T R u(k))}{\partial u(k)} \\ &+ \left( \frac{\partial x(k + 1)}{\partial u(k)} \right)^T \frac{\partial V^*(x(k + 1))}{\partial x(k + 1)} = 0 \end{aligned} \tag{4}$$

Therefore we can obtain

$$u^*(k) = -\frac{1}{2} R^{-1} g^T(x(k)) \frac{\partial V^*(x(k + 1))}{\partial x(k + 1)} \tag{5}$$

where  $V^*$  is the value function corresponding to the optimal control policy  $u^*$ .

However, for constrained control problem, the above derivation is unfeasible. To confront this bounded control problem, we introduce a nonquadratic functional motivated by Refs. [12,18].

$$W(\mathbf{u}(i)) = 2 \int_0^{u(i)} \boldsymbol{\varphi}^{-T}(\bar{\mathbf{U}}^{-1}\mathbf{s})\bar{\mathbf{U}}\mathbf{R}d\mathbf{s} \tag{6}$$

$$\boldsymbol{\varphi}^{-1}(\mathbf{u}(i)) = [\phi^{-1}(u_1(i)), \phi^{-1}(u_2(i)), \dots, \phi^{-1}(u_m(i))]^T$$

where  $\mathbf{R}$  is positive definite and assumed to be diagonal for simplicity of analysis,  $\mathbf{s} \in \mathbb{R}^m$ ,  $\boldsymbol{\varphi} \in \mathbb{R}^m$ ,  $\phi(\cdot)$  is a bounded one-to-one function satisfying  $|\phi(\cdot)| \leq 1$  and belonging to  $C^p(p \geq 1)$  and  $L_2(\Omega)$ . Moreover, it is a monotonic increasing odd function with its first derivative bounded by a constant  $M$ . Such function is easy to find, one example is the hyperbolic tangent function  $\phi(\cdot) = \tanh(\cdot)$ . It should be noticed that by the definition above,  $W(\mathbf{u}(i))$  is assured to be positive definite because  $\boldsymbol{\varphi}^{-1}(\mathbf{u}(i))$  is monotonic odd function and  $R$  is positive definite.

Substituting (6) into (3), we can obtain

$$\begin{aligned} V^*(\mathbf{x}(k)) &= \min_{\mathbf{u}(i)} \sum_{i=k}^{\infty} \left\{ \mathbf{x}(i)^T \mathbf{Q}\mathbf{x}(i) + 2 \int_0^{u(i)} \boldsymbol{\varphi}^{-T}(\bar{\mathbf{U}}^{-1}\mathbf{s})\bar{\mathbf{U}}\mathbf{R}d\mathbf{s} \right\} \\ &= \min_{\mathbf{u}(k)} \left\{ \mathbf{x}(k)^T \mathbf{Q}\mathbf{x}(k) + 2 \int_0^{u(k)} \boldsymbol{\varphi}^{-T}(\bar{\mathbf{U}}^{-1}\mathbf{s})\bar{\mathbf{U}}\mathbf{R}d\mathbf{s} \right. \\ &\quad \left. + V^*(\mathbf{x}(k+1)) \right\} \end{aligned} \tag{7}$$

According to the first order necessary condition of the optimal control, the following equation holds:

$$\begin{aligned} \frac{\partial V^*(\mathbf{x}(k))}{\partial \mathbf{u}(k)} &= 2\bar{\mathbf{U}}\mathbf{R}\boldsymbol{\varphi}^{-1}(\bar{\mathbf{U}}^{-1}\mathbf{u}(k)) + \left( \frac{\partial \mathbf{x}(k+1)}{\partial \mathbf{u}(k)} \right)^T \\ &\quad \times \frac{\partial V^*(\mathbf{x}(k+1))}{\partial \mathbf{x}(k+1)} = 0 \end{aligned} \tag{8}$$

Therefore, the following formulation can be obtained:

$$\mathbf{u}^*(k) = \bar{\mathbf{U}}\boldsymbol{\varphi} \left( -\frac{1}{2}(\bar{\mathbf{U}}\mathbf{R})^{-1}\mathbf{g}^T(\mathbf{x}(k))V_x^*(\mathbf{x}(k+1)) \right) \tag{9}$$

Substituting (9) into (7), we obtain the special discrete-time HJB equation as follows:

$$\begin{aligned} V^*(\mathbf{x}(k)) &= \mathbf{x}(k)^T \mathbf{Q}\mathbf{x}(k) \\ &\quad + 2 \int_0^{\bar{\mathbf{U}}\boldsymbol{\varphi} \left( -\frac{1}{2}(\bar{\mathbf{U}}\mathbf{R})^{-1}\mathbf{g}^T(\mathbf{x}(k))V_x^*(\mathbf{x}(k+1)) \right)} \boldsymbol{\varphi}^{-T}(\bar{\mathbf{U}}^{-1}\mathbf{s})\bar{\mathbf{U}}\mathbf{R}d\mathbf{s} \\ &\quad + V^* \left( \mathbf{f}(\mathbf{x}(k)) + \mathbf{g}(\mathbf{x}(k))\bar{\mathbf{U}} \right. \\ &\quad \left. \times \boldsymbol{\varphi} \left( -\frac{1}{2}(\bar{\mathbf{U}}\mathbf{R})^{-1}\mathbf{g}^T(\mathbf{x}(k))V_x^*(\mathbf{x}(k+1)) \right) \right) \end{aligned} \tag{10}$$

The optimal control  $\mathbf{u}^*(k)$  can be computed if the value function  $V^*(\mathbf{x}(k+1))$  can be solved from HJB equation (10). However, there is currently no method for rigorously solving the value function of this constrained optimal control problem. Therefore, in the next section we will discuss

how to use the approximate dynamic programming algorithm named GI-HDP algorithm to solve the near-optimal control solution.

### 3. The derivation and implementation of near-optimal control scheme based on GI-HDP algorithm

#### 3.1. Derivation of the GI-HDP algorithm and the corresponding convergence analysis

In the GI-HDP algorithm, for dealing with the control constraint, we choose  $W(\mathbf{u}(k)) = 2 \int_0^{u(k)} \boldsymbol{\varphi}^{-T}(\bar{\mathbf{U}}^{-1}\mathbf{s})\bar{\mathbf{U}}\mathbf{R}d\mathbf{s}$ . First, we start with initial cost function  $V_0(\mathbf{x}(k)) = 0$  which is not necessarily the value function, and then find the control  $\mathbf{u}_0(\mathbf{x}(k))$  as follows:

$$\mathbf{u}_0(\mathbf{x}(k)) = \arg \min_{\mathbf{u}(k)} \left( \mathbf{x}(k)^T \mathbf{Q}\mathbf{x}(k) + W(\mathbf{u}(k)) + V_0(\mathbf{x}(k+1)) \right) \tag{11}$$

then update the cost function as

$$\begin{aligned} V_1(\mathbf{x}(k)) &= \mathbf{x}(k)^T \mathbf{Q}\mathbf{x}(k) + W(\mathbf{u}_0(k)) + V_0(\mathbf{f}(\mathbf{x}(k))) \\ &\quad + \mathbf{g}(\mathbf{x}(k))\mathbf{u}_0(k) \end{aligned} \tag{12}$$

The GI-HDP algorithm therefore iterates between

$$\begin{aligned} \mathbf{u}_i(\mathbf{x}(k)) &= \arg \min_{\mathbf{u}(k)} \left( \mathbf{x}(k)^T \mathbf{Q}\mathbf{x}(k) + W(\mathbf{u}(k)) \right. \\ &\quad \left. + V_i(\mathbf{f}(\mathbf{x}(k)) + \mathbf{g}(\mathbf{x}(k))\mathbf{u}(k)) \right) \\ &= \bar{\mathbf{U}}\boldsymbol{\varphi} \left( -\frac{1}{2}(\bar{\mathbf{U}}\mathbf{R})^{-1}\mathbf{g}^T(\mathbf{x}(k)) \frac{\partial V_i(\mathbf{x}(k+1))}{\partial \mathbf{x}(k+1)} \right) \end{aligned} \tag{13}$$

and

$$\begin{aligned} V_{i+1}(\mathbf{x}(k)) &= \min_{\mathbf{u}(k)} \left( \mathbf{x}(k)^T \mathbf{Q}\mathbf{x}(k) + W(\mathbf{u}(k)) \right. \\ &\quad \left. + V_i(\mathbf{f}(\mathbf{x}(k)) + \mathbf{g}(\mathbf{x}(k))\mathbf{u}(k)) \right) \\ &= \mathbf{x}(k)^T \mathbf{Q}\mathbf{x}(k) + W(\mathbf{u}_i(k)) \\ &\quad + V_i(\mathbf{f}(\mathbf{x}(k)) + \mathbf{g}(\mathbf{x}(k))\mathbf{u}_i(k)) \end{aligned} \tag{14}$$

In this way, the cost function and control policy are updated by recurrent iteration until they converge to the optimal ones, with the iteration number  $i$  increasing from 0 to  $\infty$ . In the following part, we shall present a proof of convergence of the iteration between (13) and (14) with the cost function  $V_i \Rightarrow V^*$  and the control policy  $\mathbf{u}_i \Rightarrow \mathbf{u}^*$  as  $i \Rightarrow \infty$ .

**Lemma 1.** Let  $\boldsymbol{\mu}_i$  be any arbitrary sequence of control policies and  $\mathbf{u}_i$  be the policies as (13). Let  $V_i$  be as (14) and  $A_i$  as

$$\begin{aligned} A_{i+1}(\mathbf{x}(k)) &= \mathbf{x}(k)^T \mathbf{Q}\mathbf{x}(k) + W(\boldsymbol{\mu}_i(k)) + A_i(\mathbf{f}(\mathbf{x}(k))) \\ &\quad + \mathbf{g}(\mathbf{x}(k))\boldsymbol{\mu}_i(k) \end{aligned} \tag{15}$$

If  $V_0 = A_0 = 0$ , then  $V_i \leq A_i, \forall i$ .

**Proof.** It is quite clear from the fact that  $V_{i+1}$  is a result of minimizing the right-hand side of (14) with respect to the control input  $\mathbf{u}$ , while  $A_{i+1}$  is a result of any arbitrary control input.  $\square$

**Lemma 2.** Let the sequence  $\{V_i\}$  be defined as (14). If the system is controllable, there is an upper bound  $Y$  such that  $0 \leq V_i \leq Y, \forall i$ .

**Proof.** Let  $\eta(\mathbf{x}(k))$  be any stabilizing and admissible control input, and let  $V_0(\cdot) = Z_0(\cdot) = 0$ , where  $V_i$  is updated as (14) and  $Z_i$  is updated by

$$Z_{i+1}(\mathbf{x}(k)) = \mathbf{x}(k)^T \mathbf{Q}\mathbf{x}(k) + W(\eta(k)) + Z_i(\mathbf{x}(k+1)) \quad (16)$$

It follows that the difference

$$\begin{aligned} Z_{i+1}(\mathbf{x}(k)) - Z_i(\mathbf{x}(k)) &= Z_i(\mathbf{x}(k+1)) - Z_{i-1}(\mathbf{x}(k+1)) \\ &= Z_{i-1}(\mathbf{x}(k+2)) - Z_{i-2}(\mathbf{x}(k+2)) \\ &= Z_{i-2}(\mathbf{x}(k+3)) - Z_{i-3}(\mathbf{x}(k+3)) \dots \\ &= Z_1(\mathbf{x}(k+i)) - Z_0(\mathbf{x}(k+i)) \end{aligned} \quad (17)$$

Then the following relation can be obtained

$$Z_{i+1}(\mathbf{x}(k)) - Z_i(\mathbf{x}(k)) = Z_1(\mathbf{x}(k+i)) - Z_0(\mathbf{x}(k+i)) \quad (18)$$

Since  $Z_0(\cdot) = 0$ , so we have

$$\begin{aligned} Z_{i+1}(\mathbf{x}(k)) &= Z_1(\mathbf{x}(k+i)) + Z_i(\mathbf{x}(k)) \\ &= Z_1(\mathbf{x}(k+i)) + Z_1(\mathbf{x}(k+i-1)) + Z_{i-1}(\mathbf{x}(k)) \\ &= Z_1(\mathbf{x}(k+i)) + Z_1(\mathbf{x}(k+i-1)) \\ &\quad + Z_1(\mathbf{x}(k+i-2)) + Z_{i-2}(\mathbf{x}(k)) \\ &= Z_1(\mathbf{x}(k+i)) + Z_1(\mathbf{x}(k+i-1)) \\ &\quad + Z_1(\mathbf{x}(k+i-2)) + \dots + Z_1(\mathbf{x}(k)) \end{aligned} \quad (19)$$

So (19) can be written as

$$\begin{aligned} Z_{i+1}(\mathbf{x}(k)) &= \sum_{j=0}^i Z_1(\mathbf{x}(k+j)) \\ &= \sum_{j=0}^i (\mathbf{x}(k+j)^T \mathbf{Q}\mathbf{x}(k+j) + W(\eta(\mathbf{x}(k+j)))) \\ &\leq \sum_{j=0}^{\infty} (\mathbf{x}(k+j)^T \mathbf{Q}\mathbf{x}(k+j) + W(\eta(\mathbf{x}(k+j)))) \end{aligned} \quad (20)$$

Note that  $\eta(\mathbf{x}(k))$  is a stabilizing and admissible control input, i.e.,  $\mathbf{x}(k) \rightarrow 0$  as  $k \rightarrow \infty$ , therefore we have

$$\forall i: Z_{i+1}(\mathbf{x}(k)) \leq \sum_{i=0}^{\infty} Z_1(\mathbf{x}(k+i)) \leq Y \quad (21)$$

Combining with Lemma 1, we can obtain

$$\forall i: V_{i+1}(\mathbf{x}(k)) \leq Z_{i+1}(\mathbf{x}(k)) \leq Y \quad (22)$$

This completes the proof.  $\square$

With Lemma 1 and Lemma 2, the next main theorem can be derived.

**Theorem 1.** Define the sequence  $\{V_i\}$  as (14), with  $V_0(\cdot) = 0$ . Then  $\{V_i\}$  is a nondecreasing sequence satisfying  $V_{i+1}(\mathbf{x}(k)) \geq V_i(\mathbf{x}(k)), \forall i$  and converging to the value function of the discrete-time HJB equation (10), i.e.,  $V_i \Rightarrow V^*$  as  $i \Rightarrow \infty$ . Meanwhile, the control policy also converges to the optimal control policy (9), i.e.,  $\mathbf{u}_i \Rightarrow \mathbf{u}^*$  as  $i \Rightarrow \infty$ .

**Proof.** For the convenience of analysis, define a new sequence  $\Phi_i$  as follows:

$$\Phi_{i+1}(\mathbf{x}(k)) = \mathbf{x}(k)^T \mathbf{Q}\mathbf{x}(k) + W(\mathbf{u}_{i+1}(k)) + \Phi_i(\mathbf{x}(k+1)) \quad (23)$$

with  $\Phi_0 = V_0 = 0$  and the policies  $\mathbf{u}_i$  defined as (13), the cost function  $V_i$  is updated by (14).

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

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

$$V_1(\mathbf{x}(k)) - \Phi_0(\mathbf{x}(k)) = \mathbf{x}(k)^T \mathbf{Q}\mathbf{x}(k) + W(\mathbf{u}_0(k)) \geq 0 \quad (24)$$

thus for  $i = 0$ , we get

$$V_1(\mathbf{x}(k)) \geq \Phi_0(\mathbf{x}(k)) \quad (25)$$

Second, we assume that it holds for  $i - 1$ , i.e.,  $V_i(\mathbf{x}(k)) \geq \Phi_{i-1}(\mathbf{x}(k)), \forall \mathbf{x}(k)$ . Then for  $i$ , since

$$\Phi_i(\mathbf{x}(k)) = \mathbf{x}(k)^T \mathbf{Q}\mathbf{x}(k) + W(\mathbf{u}_i(k)) + \Phi_{i-1}(\mathbf{x}(k+1)) \quad (26)$$

and

$$V_{i+1}(\mathbf{x}(k)) = \mathbf{x}(k)^T \mathbf{Q}\mathbf{x}(k) + W(\mathbf{u}_i(k)) + V_i(\mathbf{x}(k+1)) \quad (27)$$

hold, then we obtain

$$V_{i+1}(\mathbf{x}(k)) - \Phi_i(\mathbf{x}(k)) = V_i(\mathbf{x}(k+1)) - \Phi_{i-1}(\mathbf{x}(k+1)) \geq 0 \quad (28)$$

i.e., the following equation holds

$$\Phi_i(\mathbf{x}(k)) \leq V_{i+1}(\mathbf{x}(k)) \quad (29)$$

Therefore, the mathematical induction proof is completed.

Furthermore, from Lemma 1 we know that  $V_i(\mathbf{x}(k)) \leq \Phi_i(\mathbf{x}(k))$ , therefore we have

$$V_i(\mathbf{x}(k)) \leq \Phi_i(\mathbf{x}(k)) \leq V_{i+1}(\mathbf{x}(k)) \quad (30)$$

hence we can draw the conclusion that  $\{V_i\}$  is a nondecreasing sequence with upper bound as shown in Lemma 2, i.e.,  $V_i \rightarrow V^*$  as  $i \rightarrow \infty$ .

We have just proved that the cost function converges to the value function of the discrete-time HJB equation, and according to (9) and (13), we can conclude that the corresponding control policy converges to the optimal one as well.

This completes the proof.  $\square$

### 3.2. Neural network implementation for GI-HDP algorithm

In the case of linear systems the value function and the control policy are quadratic and linear, respectively. In the nonlinear case, this is not necessarily true and therefore we need to use parametric structures or neural networks to approximate both  $\mathbf{u}_i(\mathbf{x}(k))$  and  $V_i(\mathbf{x}(k))$ .

Assume that the number of hidden layer neurons is denoted by  $l$ , the weight matrix between the input layer and the hidden layer denoted by  $\mathbf{v}$ , the weight matrix between the hidden layer and the output layer denoted by  $\mathbf{w}$ , then output of three-layer NN is represented by:

$$\hat{f}(X, \mathbf{v}, \mathbf{w}) = \mathbf{w}^T \boldsymbol{\sigma}(\mathbf{v}^T X) \tag{31}$$

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

In order to implement the GI-HDP algorithm on (13) and (14), we now employ neural networks to approximate the value function and the corresponding optimal control policy. In the GI-HDP algorithm, there are three networks, which are critic network, model network and action network. All the neural networks are chosen as the three-layer feedforward networks. The inputs of the critic network and action network are  $\mathbf{x}(k)$ , and the inputs of the model network are  $\mathbf{x}(k)$  and  $\hat{\mathbf{u}}_i(k)$ . The whole structure diagram is shown in Fig. 1, where the utility term denotes  $\mathbf{x}(k)^T \mathbf{Q}\mathbf{x}(k) + W(\mathbf{u}_i(k))$ .

For an unknown plant, before carrying out the GI-HDP algorithm, we should first train the model network. For given  $\mathbf{x}(k)$  and  $\hat{\mathbf{u}}_i(k)$ , we can obtain  $\hat{\mathbf{x}}(k+1)$ , and the output of the model network is denoted as

$$\hat{\mathbf{x}}(k+1) = \mathbf{w}_m^T \boldsymbol{\sigma}(\mathbf{v}_m^T I_m(k)) \tag{32}$$

where  $I_m(k) = [\mathbf{x}(k) \hat{\mathbf{u}}_i(k)]$  is the input vector of the model network.

We define the error function of the model network as

$$e_m(k) = \hat{\mathbf{x}}(k+1) - \mathbf{x}(k+1) \tag{33}$$

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

$$E_m(k) = \frac{1}{2} e_m^T e_m(k) \tag{34}$$

The weight update rule for model network is chosen as a gradient-based adaptation rule

$$\mathbf{w}_m(k+1) = \mathbf{w}_m(k) - \alpha_m \left[ \frac{\partial E_m(k)}{\partial \mathbf{w}_m(k)} \right] \tag{35}$$

$$\mathbf{v}_m(k+1) = \mathbf{v}_m(k) - \alpha_m \left[ \frac{\partial E_m(k)}{\partial \mathbf{v}_m(k)} \right] \tag{36}$$

where  $\alpha_m$  is the learning rate of the model network.

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

The critic network is used to approximate the cost function  $V_i(k)$ . The output of the critic network is denoted as

$$\hat{V}_i(k) = \mathbf{w}_{ci}^T \boldsymbol{\sigma}(\mathbf{v}_{ci}^T \mathbf{x}(k)) \tag{37}$$

The target function can be written as

$$V_{i+1}(k) = \mathbf{x}^T(k) \mathbf{Q}\mathbf{x}(k) + W(\mathbf{u}_i(k)) + \hat{V}_i(k+1) \tag{38}$$

Then we define the error function for the critic network as

$$e_{ci}(k) = \hat{V}_i(k) - V_{i+1}(k) \tag{39}$$

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

$$E_{ci}(k) = \frac{1}{2} e_{ci}^2(k) \tag{40}$$

The weight update rule for the critic network is a gradient-based adaptation given by

$$\mathbf{w}_{ci(i+1)}(k) = \mathbf{w}_{ci}(k) - \alpha_c \left[ \frac{\partial E_{ci}(k)}{\partial \mathbf{w}_{ci}(k)} \right] \tag{41}$$

$$\mathbf{v}_{ci(i+1)}(k) = \mathbf{v}_{ci}(k) - \alpha_c \left[ \frac{\partial E_{ci}(k)}{\partial \mathbf{v}_{ci}(k)} \right] \tag{42}$$

where  $\alpha_c > 0$  is the learning rate of the critic network.

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

$$\hat{\mathbf{u}}_i(k) = \mathbf{w}_{ai}^T \boldsymbol{\sigma}(\mathbf{v}_{ai}^T \mathbf{x}(k)) \tag{43}$$

And the target control input is given as

$$\mathbf{u}_i(k) = \bar{\mathbf{U}} \boldsymbol{\varphi} \left( -\frac{1}{2} (\bar{\mathbf{U}} \mathbf{R})^{-1} \mathbf{g}^T(\mathbf{x}(k)) \frac{\partial V_i(\mathbf{x}(k+1))}{\partial \mathbf{x}(k+1)} \right) \tag{44}$$

So we can define the output error of the action network as

$$e_{ai}(k) = \hat{\mathbf{u}}_i(k) - \mathbf{u}_i(k) \tag{45}$$

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

$$E_{ai}(k) = \frac{1}{2} e_{ai}^T e_{ai}(k) \tag{46}$$

The update algorithm is similar to that in the critic network. By the gradient descent rule

$$\mathbf{w}_{ai(i+1)}(k) = \mathbf{w}_{ai}(k) - \beta_a \left[ \frac{\partial E_{ai}(k)}{\partial \mathbf{w}_{ai}(k)} \right] \tag{47}$$

$$\mathbf{v}_{ai(i+1)}(k) = \mathbf{v}_{ai}(k) - \beta_a \left[ \frac{\partial E_{ai}(k)}{\partial \mathbf{v}_{ai}(k)} \right] \tag{48}$$

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

#### 4. Simulation study

In this section, an example is provided to demonstrate the effectiveness of the control scheme proposed in this paper.

Consider the following nonlinear system in Ref. [20]:

$$\mathbf{x}(k+1) = \mathbf{f}(\mathbf{x}(k)) + \mathbf{g}(\mathbf{x}(k))\mathbf{u}(k) \tag{49}$$

where  $\mathbf{f}(\mathbf{x}(k)) = \begin{bmatrix} -0.8\mathbf{x}_2(k) \\ \sin(0.8\mathbf{x}_1(k) - \mathbf{x}_2(k)) + 1.8\mathbf{x}_2(k) \end{bmatrix}$ ,  $\mathbf{g}(\mathbf{x}(k)) = \begin{bmatrix} 0 \\ -\mathbf{x}_2(k) \end{bmatrix}$ , and the control constraint is set to  $\|\mathbf{u}\| \leq 0.3$ .

Define the performance functional as

$$J(\mathbf{x}(k), \mathbf{u}) = \sum_{i=k}^{\infty} \left\{ \mathbf{x}(i)^T \mathbf{Q}\mathbf{x}(i) + 2 \int_0^{\mathbf{u}(i)} \tanh^{-T}(\bar{\mathbf{U}}^{-1}\mathbf{s}) \bar{\mathbf{U}} \mathbf{R} d\mathbf{s} \right\} \tag{50}$$

where the weight matrix is chosen as  $\mathbf{Q} = \begin{bmatrix} 0.1 & 0 \\ 0 & 0.1 \end{bmatrix}$  and

$$\mathbf{R} = \begin{bmatrix} 0.1 & 0 \\ 0 & 0.1 \end{bmatrix}.$$

For this plant, the control constraint is represented as  $\bar{\mathbf{U}} = 0.3$ . We choose three-layer feedforward neural net-

works as the critic network, the action network and the model network with the structures 2-8-1, 2-8-1 and 3-8-2, respectively. The initial weights of action network, critic network and model network are all set to be random in  $[-1,1]$ . It should be mentioned that the model network should be trained first. We train the model network for 1000 steps under the learning rate  $\alpha_m = 0.1$ . After the training of the model network is completed, the weights keep unchanged. Then the critic network and the action network are trained for 100 training cycles with each cycle of 1000 steps. In the training process, the learning rate  $\beta_a = \alpha_c = 0.1$ . Then for the given initial states  $x_1(0) = 0.5$ ,  $x_2(0) = 0.5$ , we apply the saturated optimal control policy to the system for 100 time steps and obtain the state curves shown in Fig. 2, the control curve shown in Fig. 3 and the convergence process of the cost function shown in Fig. 4.

Moreover, in order to make comparison with the controller without considering the actuator saturation, we also present the controller designed by GI-HDP algorithm regardless of the saturation of the actuators. The state curves are shown in Fig. 5 and the control curve is shown in Fig. 6.

From the simulation results, we can see that the iterative cost function sequence does converge to the optimal value quite rapidly, which also indicates the validity of the GI-HDP algorithm for dealing with constrained nonlinear system. Comparing Fig. 3 with Fig. 6, we can find that in Fig. 3 the saturation of the actuator has been overcome successfully.

### 5. Conclusion

In this paper an effective algorithm is proposed to find the approximate optimal controller for a class of discrete-time constrained systems. First a new type of nonquadratic functional is defined to deal with the control constraint, and then the GI-HDP algorithm is introduced to solve the value function of the HJB equation with rigorous convergence analysis. Three neural networks are used as parametric structures to approximate the value function, compute the optimal control policy and model the unknown system, respectively. The simulation study has demonstrated the effectiveness of the proposed optimal control scheme.

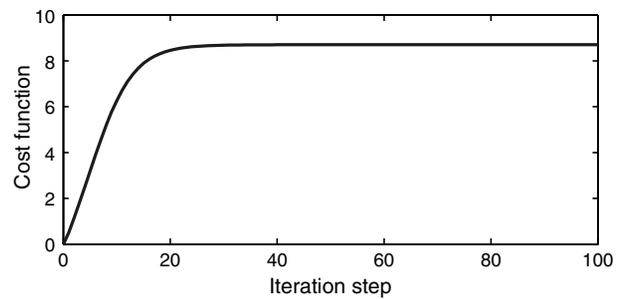


Fig. 4. The convergence process of the cost function.

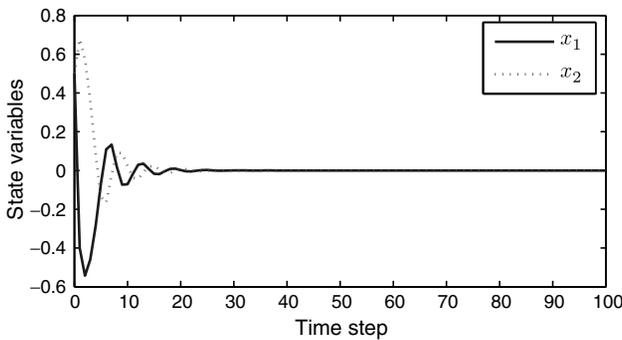


Fig. 2. The state variables curves.

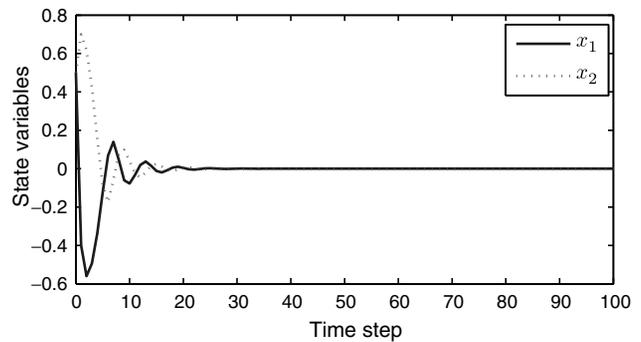


Fig. 5. The state variables curves.

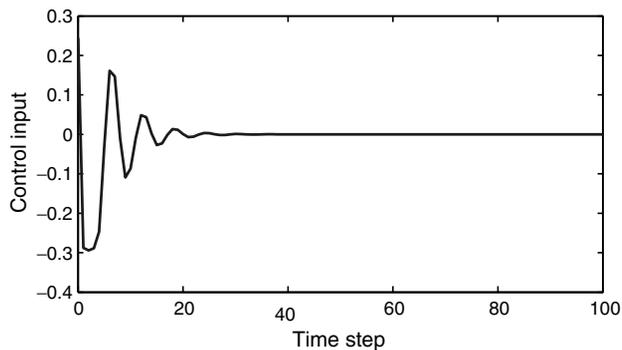


Fig. 3. The control input curve.

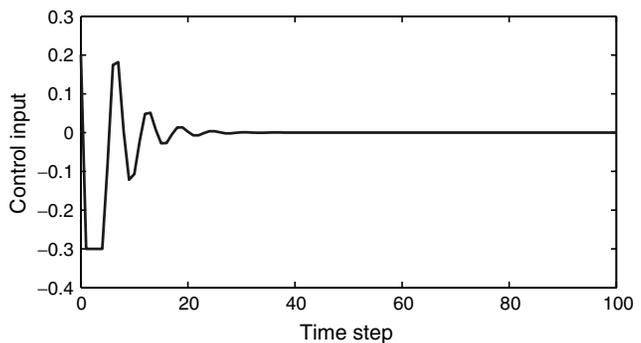


Fig. 6. The control input curve.

## Acknowledgment

This work was supported by the National Natural Science Foundation of China (Grant Nos. 60534010, 60572070, 60774048, 60728307), the Program for Changjiang Scholars and Innovative Research Groups of China (60521003), the Research Fund for the Doctoral Program of China Higher Education (20070145015), and the National High Technology Research and Development Program of China (2006AA04Z183).

## References

- [1] Sussmann H, Sontag ED, Yang Y. A general result on the stabilization of linear systems using bounded controls. *IEEE Trans Automatic Control* 1994;39(12):2411–25.
- [2] Bernstein DS. Optimal nonlinear, but continuous, feedback control of systems with saturating actuators. *Int J Control* 1995;62(5):1209–16.
- [3] Saberi A, Lin Z, Teel A. Control of linear systems with saturating actuators. *IEEE Trans Automatic Control* 1996;41(3):368–78.
- [4] Widrow B, Gupta N, Maitra S. Punish/reward: learning with a critic in adaptive threshold systems. *IEEE Trans Syst Man Cybern* 1973;3(5):455–65.
- [5] Watkins C. Learning from delayed rewards. Ph.D. thesis, Cambridge University, Cambridge, England; 1989.
- [6] Bertsekas DP, Tsitsiklis JN. *Neuro-dynamic programming*. Massachusetts: Athena Scientific; 1996.
- [7] Liu X, Balakrishnan SN. Convergence analysis of adaptive critic based optimal control. In: *Proceedings of American control conference*, Chicago, IL, USA; 2000. p. 1929–33.
- [8] Si J, Wang YT. On-line learning control by association and reinforcement. *IEEE Trans Neural Netw* 2001;12(2):264–76.
- [9] Liu D, Xiong X, Zhang Y. Action-dependent adaptive critic designs. In: *Proceedings of the INNS-IEEE international joint conference on neural networks*, Washington, DC; 2001. p. 990–5.
- [10] Liu D, Zhang HG. A neural dynamic programming approach for learning control of failure avoidance problems. *Int J Intell Control Syst* 2005;10(1):21–32.
- [11] Murray JJ, Cox CJ, Lendaris GG, et al. Adaptive dynamic programming. *IEEE Trans Syst Man Cybern C Appl Rev* 2002;32(2):140–53.
- [12] Abu-Khalaf M, Lewis FL. Nearly optimal control laws for nonlinear systems with saturating actuators using a neural network HJB approach. *Automatica* 2005;41(5):779–91.
- [13] Liu D. Approximate dynamic programming for self-learning control. *Acta Automatica Sin* 2005;31(1):13–8.
- [14] Al-Tamimi A, Lewis FL. Discrete-time nonlinear HJB solution using approximate dynamic programming: convergence proof. In: *Proceedings of the IEEE symposium on approximate dynamic programming and reinforcement learning*, Honolulu, HI, USA; 2007. p. 38–43.
- [15] Zhang HG, Wei QL, Liu D. On-line learning control for discrete nonlinear systems via an improved ADDHP method. In: *Proceedings of the 4th International Symposium on Neural Networks*. Nanjing, China; 2007. p. 387–96.
- [16] Prokhorov DV, Wunsch DC. Adaptive critic designs. *IEEE Trans Neural Netw* 1997;8(5):997–1007.
- [17] Lyshevski SE. Optimal control of nonlinear continuous-time systems: design of bounded controllers via generalized nonquadratic functionals. In: *Proceedings of the American control conference*, Philadelphia, USA; 1998. p. 205–9.
- [18] Lyshevski SE. Nonlinear discrete-time systems: constrained optimization and application of nonquadratic costs. In: *Proceedings of the American control conference*, Philadelphia, USA; 1998. p. 3699–703.
- [19] Beard R. Improving the closed-loop performance of nonlinear systems. Ph.D. thesis, Rensselaer Polytechnic Institute, Troy, New York, USA; 1995.
- [20] Chen Z, Jagannathan S. Generalized Hamilton–Jacobi–Bellman formulation-based neural network control of affine nonlinear discrete-time systems. *IEEE Trans Neural Netw* 2008;19(1):90–106.