

# Safety-Aware Reinforcement Learning Framework with an Actor-Critic-Barrier Structure

Yongliang Yang<sup>1</sup>, Kyriakos G. Vamvoudakis<sup>2</sup>, Hamidreza Modares<sup>3</sup>, Wei He<sup>1</sup>, Yixin Yin<sup>1</sup>, Donald C. Wunsch<sup>4</sup>

**Abstract**—This paper considers the control problem with constraints on full-state and control input simultaneously. First, a novel barrier function based system transformation approach is developed to guarantee the full-state constraints. To deal with the input saturation, the hyperbolic-type penalty function is imposed on the control input. The actor-critic based reinforcement learning technique is combined with the barrier transformation to learn the optimal control policy that considers both the full-state constraints and input saturations. To illustrate the efficacy, a numeric simulation is implemented in the end.

**Index Terms**—reinforcement learning, full-state constraints, input saturation, safe control

## I. INTRODUCTION

Due to the constraints in real-world applications, control problem with constraints has attracted numerous attention recently. In addition to stability and control performance, the constraints on state and/or control is critical for safety purpose. The constrained control problem has wide applications in the field of robotics [1]. This paper presents a novel adaptive optimal learning algorithm to consider the constraints on states and control input simultaneously.

### Related work

To consider the input saturation constraints, [2] employs the hyperbolic function to represent the penalty on the control input. In addition to closed-loop stability, the control performance is considered using optimal control theory [2].

This work was supported in part by the Fundamental Research Funds for the China Central Universities of USTB under grant No. FRF-TP-18-031A1 and No. FRF-BD-17-002A, in part by the China Post-Doctoral Science Foundation under Grant 2018M641197, in part by the National Science Foundation under grant NSF CAREER CPS-1851588, in part by NATO under grant No. SPS G5176, in part by ONR Minerva under grant No. N00014-18-1-2160, in part by the Mary K. Finley Endowment, in part by the Missouri S&T Intelligent Systems Center and in part by the Army Research Laboratory under Cooperative Agreement Number W911NF-18-2-0260. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Army Research Laboratory or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for government purposes notwithstanding any copyright notation herein.

<sup>1</sup>Y. Yang, W. He and Y. Yin are with the School of Automation & Electrical Engineering, University of Science and Technology Beijing, Beijing 100083, China. yangyongliang@ieee.org; weihe@ieee.org; yyx@ies.ustb.edu.cn

<sup>2</sup>K. G. Vamvoudakis is with the Daniel Guggenheim School of Aerospace Engineering, Georgia Tech, GA 30332, USA. kyriakos@gatech.edu

<sup>3</sup>H. Modares is with the Mechanical Engineering Department, Michigan State University, East Lansing, MI 48824, USA. modares@msu.edu

<sup>4</sup>D. C. Wunsch is with the Department of Electrical and Computer Engineering, Missouri University of Science & Technology, Rolla, MO 65401, USA. Wunsch@ieee.org

Another type of input saturation problem is analyzed in [3], where an auxiliary design system is introduced for constrained adaptive control design. However, these results only considers input saturations, whereas the state constraints is not analyzed. In this paper, the control problem with full-state constraints is taken into consideration.

In contrast to the linear quadratic control, the state constraints in control problem is achieved by some penalty on the state variable [4]. The state penalty is extended to the case of Lyapunov stability analysis, where a novel type of Lyapunov function, named barrier Lyapunov function, is proposed to deal with the state constrained control problem [5]. The barrier Lyapunov function based method is initially designed to deal with the output constraints [5], which can be viewed as the case of partial state constraints. Later, it is extended to deal with the full-state constrained control problem [6]. However, these results does not consider the input constraints. In some cases, both the state and input constraints are critical for safety issue. Therefore, we solve the constrained control problem with consideration of the state and input constraints simultaneously.

Adaptive dynamic programming is an efficient algorithm to solve the optimal decision-making problem, such as optimal regulation problem [7], [8], optimal tracking problem [9], multiagent systems [10], [11], robust control problem [12], differential games [13], [14] and intermittent feedback design [15]. ADP has also been successfully applied to the input constrained control problem [16], [17]. In this paper, we extend the ADP method to the full-state constrained optimal control problem.

In this paper, we first propose a novel barrier function based system transformation method. Different from the barrier Lyapunov function method, the full-state constrained system can be transformed to an equivalent system without state constraints. It is guaranteed that if the initial state is within the prescribed bound, the state constraints can be guaranteed. Then, the input saturation is considered for transformed system. The hyperbolic penalty function is designed for the control input. Moreover, we optimize the nonlinear system using an online algorithm.

### Notation and Background

**Definition 1. (Persistent Excitation)** A vector signal  $y(t) \in \mathbb{R}^p$  is exciting over the interval  $[t, t+T]$  with  $T \in \mathbb{R}^+$  if there exists  $\beta_1 \in \mathbb{R}^+$  and  $\beta_2 \in \mathbb{R}^+$  such that for  $\forall t$ ,

$$\beta_1 I_{p \times p} \leq \int_t^{t+T} y(\tau) y^T(\tau) d\tau \leq \beta_2 I_{p \times p} \quad \square$$

## II. PROBLEM FORMULATION

In this paper, we consider the following continuous-time affine nonlinear dynamical systems

$$\begin{aligned}\dot{x}_1 &= x_2 \\ \dot{x}_2 &= x_3 \\ &\vdots \\ \dot{x}_{n-1} &= x_n \\ \dot{x}_n &= f(x) + g(x)u\end{aligned}\quad (1)$$

where  $x = [x_1 \cdots x_n]^T \in \mathbb{R}^n$  with  $x_i \in \mathbb{R}$  is the system state,  $u \in \mathbb{R}$  is the control input,  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  and  $g: \mathbb{R}^n \rightarrow \mathbb{R}$  are Lipschitz continuous functions. The constrained control problem of system (1) with full state constraints can be formulated as follows.

**Problem 1.** Consider the system (1), find a policy  $u: \mathbb{R}^n \rightarrow \mathbb{R}^m$  such that the closed-loop system has an asymptotic stable equilibrium while the control input satisfies

$$\|u_i\| \leq \lambda, \forall i = 1, \dots, m \quad (2)$$

and the state  $x = [x_1 \cdots x_n]^T$  satisfies

$$\begin{aligned}x_1 &\in (a_1, A_1) \\ &\vdots \\ x_n &\in (a_n, A_n),\end{aligned}\quad (3) \quad \square$$

### A. Barrier Function

**Definition 2.** The function  $b(\cdot): \mathbb{R} \rightarrow \mathbb{R}$  defined on  $(a, A)$  is referred to as barrier function if

$$b(z; a, A) = \log \left( \frac{A - a - z}{a - A - z} \right), \forall z \in (a, A) \quad (4)$$

where  $a$  and  $A$  are two constants satisfying  $a < A$ . Moreover, the barrier function is invertible on interval  $(a, A)$ , i.e.,

$$b^{-1}(y; a, A) = aA \frac{e^{\frac{y}{2}} - e^{-\frac{y}{2}}}{ae^{\frac{y}{2}} - Ae^{-\frac{y}{2}}}, \forall y \in \mathbb{R} \quad (5)$$

with the derivative

$$\frac{db^{-1}(y; a, A)}{dy} = \frac{Aa^2 - aA^2}{a^2e^y - 2aA + A^2e^{-y}} \quad (6)$$

Consider the barrier function based state transformation as

$$\begin{aligned}s_i &= b(x_i; a_i, A_i), \\ x_i &= b^{-1}(s_i; a_i, A_i),\end{aligned} \quad i = 1, \dots, n \quad (7)$$

then,

$$\frac{dx_i}{dt} = \frac{dx_i}{ds_i} \frac{ds_i}{dt} \quad (8)$$

which can yield the following based on (6)

$$\begin{aligned}\dot{s}_i &= \frac{dx_{i+1}(s_{i+1})}{\frac{db^{-1}(y; a_i, A_i)}{dy} \Big|_{y=s_i}} \\ &= F_i(s_i, s_{i+1}), \quad i = 1, \dots, n-1 \\ \dot{s}_n &= \frac{f(x) + g(x)u}{\frac{db^{-1}(y; a_n, A_n)}{dy} \Big|_{y=s_n}} \\ &= F_n(s) + g_n(s)u\end{aligned}\quad (9) \quad (10)$$

where

$$\begin{aligned}F_n(s) &= \frac{A_n^2 e^{-s_n} - 2a_n A_n + a_n^2 e^{s_n}}{A_n a_n^2 - a_n A_n^2} \times \\ &\quad f([b_1^{-1}(s_1) \cdots b_n^{-1}(s_n)]) \\ g_n(s) &= \frac{A_n^2 e^{-s_n} - 2a_n A_n + a_n^2 e^{s_n}}{A_n a_n^2 - a_n A_n^2} \times \\ &\quad g([b_1^{-1}(s_1) \cdots b_n^{-1}(s_n)])\end{aligned}\quad (11)$$

The above system of  $s = [s_1 \cdots s_n]^T$  can be expressed in a compact form as

$$\dot{s} = F(s) + G(s)u \quad (12)$$

with

$$F(s) = \begin{bmatrix} F_1(s_1, s_2) \\ \vdots \\ F_n(s) \end{bmatrix}, G(s) = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ g_n(s) \end{bmatrix} \quad (13)$$

**Assumption 1.** The system dynamics satisfies

- 1)  $F(s)$  is Lipschitz with  $F(0) = 0$ , and there exists a constant  $b_f$  such that, for  $x \in \Omega$ ,  $\|F(s)\| \leq b_f \|s\|$  where  $\Omega$  is a compact set containing the origin.
- 2)  $G(s)$  is bounded on  $\Omega$ , i.e., there exists a constant  $b_g$  such that  $\|G(s)\| \leq b_g$
- 3) The system (5) is controllable over the compact set  $\Omega$ .
- 4) The performance functional (7) satisfies zero-state observability.  $\square$

In the following, in order to solve Problem 1 with input saturation and full-state constraints, a novel problem with constraints of input saturation is presented.

**Problem 2. (Optimal Control Problem)** find a policy  $u(\cdot): \mathbb{R}^n \rightarrow \mathbb{R}^m$  such that the performance

$$V(s_0) = \int_{t_0}^{\infty} [Q(s(t)) + \Theta(u(t))] dt, \quad (14)$$

is minimized with the constraints of input saturation,

$$\|u_i\| \leq \lambda, \forall i = 1, \dots, m \quad (15)$$

where  $\lambda > 0$  is the bound of control input,  $Q(s)$  is a positive definite monotonically increasing function and  $\Theta(u)$  is a positive definite integrand function. For simplicity, we denote the reward function  $U(s, u) = Q(s) + \Theta(u)$ .  $\square$

**Lemma 1.** Suppose that  $u^*(\cdot)$  solves Problem 2 for system (12). Then,  $u^*(\cdot)$  can also solve Problem 1 provided that the initial state  $x_0$  of system (1) satisfies the constraints in (3).

*Proof.* This proof follows from [5, Lemma 1].  $\square$

To deal with input saturation, the nonquadratic penalty function in [2] is adopted,

$$\Theta(u) = 2 \int_0^u [\theta^{-1}(v)]^T R dv, \quad (16)$$

with  $R = \text{diag} \left( \begin{bmatrix} r_1 & \cdots & r_m \end{bmatrix} \right)$  with the positive penalty weight  $r_i \in \mathbb{R}^+$  for  $\forall i \in \{1, \dots, m\}$ , and  $\theta(\cdot)$  is the hyperbolic tangent, i.e.,

$$\theta(v) = \lambda \tanh\left(\frac{v}{\lambda}\right) \quad (17)$$

$$\Theta(u) = 2 \int_0^u \left[ \lambda \tanh^{-1} \left( \frac{v}{\lambda} \right) \right]^T R dv \quad (18)$$

Given an admissible policy  $u$ , the Hamiltonian is defined as

with

$$D = \frac{1}{2\lambda} R^{-1} G^T(s) (\nabla \phi)^T W \quad (30)$$

It is assumed that the ideal value function approximation (24) guarantees  $\|\varepsilon_{hjb}\| \leq b_{hjb}$ .

#### A. Value Function Approximation

The ideal weight,  $W$  in (24), provides the best approximation to the optimal value function  $V^*(s)$  on the compact set  $\Omega$  and is unknown. Therefore, the estimation of  $W$  is implemented by the critic network with the approximations of the value function and value gradient

$$\hat{V}(s) = W_c^T \phi_c(s) \quad (31)$$

$$\nabla \hat{V}(s) = [\phi_c(s)]^T W_c \quad (32)$$

Then, for a given policy  $u(\cdot)$ , the residual of Bellman equation approximation using the identifier NN and the critic NN, can be determined as

$$e_c(t) := U(s(t), u(t)) + \hat{W}_c^T(t) \sigma(t) \quad (33)$$

Define the critic weight approximation error

$$\tilde{W}_c = W - W_c \quad (34)$$

Then, from (26), the relation between Bellman residual  $e_c$  and the Bellman equation approximation error  $\varepsilon_B$  can be written in terms of the critic weight error  $\tilde{W}_c$  as

$$e_c = \varepsilon_B - \tilde{W}_c^T \sigma. \quad (35)$$

As discussed in [7], [16], the Bellman residual  $e_c$  using critic NN (31) is the counterpart of the temporal difference for continuous-time system. The policy evaluation for an admissible control policy  $u(\cdot)$  can be formulated as adapting the critic weight  $\tilde{W}_c$  to minimize the objective function [8]

$$E_c = \frac{1}{2} \frac{[e_c(t)]^2}{(1 + \sigma^T \sigma)^2} \quad (36)$$

Then  $e_c \rightarrow \varepsilon_B$  as  $W_c \rightarrow W$ . Using the chain rule yields the gradient descent algorithm for minimizing  $E_c$  given by [8]

$$\dot{W}_c = -\alpha_c \frac{\partial E_c}{\partial W_c} = -\alpha_c \frac{\sigma}{(1 + \sigma^T \sigma)^2} [\sigma^T W_c + U(s, u)]. \quad (37)$$

**Theorem 1.** *Let  $u$  be any admissible control policy. Let the critic NN (31) with the adaptive tuning law (37) be used to evaluate the given control policy. Suppose that the signal  $\frac{\sigma(t)}{1 + \sigma^T(t) \sigma(t)}$  satisfies the PE condition. Then,  $\tilde{W}_c$  is uniformly ultimately bounded.*  $\square$

*Proof.* Then, from (34) and (35), the dynamics of  $\tilde{W}_c$  is

$$\begin{aligned} \dot{\tilde{W}}_c = & -\alpha_c \left[ \frac{\sigma(t) \sigma^T(t)}{[1 + \sigma^T(t) \sigma(t)]^2} \right] \tilde{W}_c \\ & + \alpha_c \left[ \frac{\sigma(t)}{[1 + \sigma^T(t) \sigma(t)]^2} \right] \varepsilon_B(t), \end{aligned} \quad (38)$$

which can be viewed as a linear time-varying system with the control input  $\varepsilon_B(t)$ . Note that the signal  $\frac{\sigma(t)}{1 + \sigma^T(t) \sigma(t)}$  satisfies the PE condition. Then,  $\tilde{W}_c$  can be expressed as

$$\begin{aligned} \tilde{W}_c(t) = & \varphi(t, t_0) \tilde{W}_c(t_0) + \int_{t_0}^t \varphi(t, \tau) \frac{\alpha_c \sigma(\tau) \varepsilon_B(\tau)}{[1 + \sigma^T(\tau) \sigma(\tau)]^2} d\tau \\ \leq & \kappa_1 e^{-\kappa_2 t} + \frac{\alpha_c}{\kappa_2 [1 + \sigma^T(\tau) \sigma(\tau)]} \|\varepsilon_B(t)\| \end{aligned}$$

where

$$\frac{\partial \varphi(t, t_0)}{\partial t} = -\alpha_c \frac{\sigma(\tau) \sigma^T(\tau)}{[1 + \sigma^T(\tau) \sigma(\tau)]^2} \varphi(t, t_0), \varphi(t, t) = I$$

Note that  $\|\varepsilon_B\| \leq b_B$ . Then, according to [8], the critic weight error  $\tilde{W}_c$  is uniformly ultimately bounded. This completes the proof.  $\blacksquare$

#### B. Actor Learning: Online Synchronous Policy Iteration

As shown in (22), the optimal control policy depends on the optimal value gradient  $\frac{\partial V^*(s)}{\partial s}$ . Therefore, consider the value gradient approximation using the critic weight  $W_c$  in (31), the policy can be determined as

$$u_c(s) = -\lambda \tanh(D_c) \quad (39)$$

$$D_c = \frac{1}{2\lambda} R^{-1} G^T (\nabla \phi)^T W_c \quad (40)$$

However, this policy improvement does not guarantee the stability of the system [8], [16]. Therefore, to ensure stability in a Lyapunov sense (as discussed later), the policy applied to the system is implemented by the actor network as

$$u_a(s) = -\lambda \tanh(D_a) \quad (41)$$

$$D_a = \frac{1}{2\lambda} R^{-1} G^T (\nabla \phi)^T W_a \quad (42)$$

with the actor network weight  $W_a$ . The actor-critic adaptive learning is discussed as the following.

**Theorem 2.** *Consider the dynamical system (1). Let the control law  $u_a$  be given by (41). Let the experience replay based critic learning be*

$$\dot{W}_c = -\alpha_c \frac{\sigma_a(t)}{[1 + \sigma_a^T(t) \sigma_a(t)]^2} [U_a(t) + W_c^T(t) \sigma_a(t)] \quad (43)$$

where

$$\begin{aligned} \sigma_a(t) &:= \nabla \phi[F(s(t)) + G(s(t)) u_a(t)], \\ U_a(t) &= U(s(t), u_a(t)), \end{aligned} \quad (44)$$

Let the gradient-descent based actor learning be

$$\dot{W}_a = -\alpha_a [\nabla \phi G e_u + \nabla \phi G \tanh^2(D_a) e_u + Y W_a] \quad (45)$$

which is designed to minimize the objective function

$$E_u = \frac{1}{2} e_u^T R e_u \quad (46)$$

where

$$e_u = u_c - u_a = \lambda [\tanh(D_a) - \tanh(D_c)] \quad (47)$$

denotes the difference between the actor  $u_a$  (41) applied to the system and the control input  $u_c$  (39) as an approximation

of the optimal control policy (22) with value function approximation by (31). Then, under Assumptions 1 - 3 and suppose that the signal  $\frac{\sigma_a(t)}{1+\sigma_a^T(t)\sigma_a(t)}$  satisfies the PE condition, the closed-loop system states and the critic and actor NN errors are uniformly ultimately bounded for a sufficiently large number of NN basis, provided that

$$Y > \frac{M_a M_a^T}{2}$$

where  $M_a = \nabla \phi G \lambda [\tanh(\kappa D_a) - \tanh(D_a)]$ .  $\square$

*Proof.* When the policy  $u_a$  is applied to the system (12), we consider the following Lyapunov candidate

$$V = V^*(s) + \underbrace{\frac{1}{2} \tilde{W}_c^T \alpha_1^{-1} \tilde{W}_c}_{V_c} + \underbrace{\frac{1}{2} \tilde{W}_a^T \alpha_2^{-1} \tilde{W}_a}_{V_a} \quad (48)$$

where  $V^*(s)$  is the optimal value function,  $\tilde{W}_c$  and  $\tilde{W}_a$  are the critic and actor weight errors, respectively. Differentiating  $V$  yields,

$$\dot{V} = \dot{V}^* + \dot{V}_c + \dot{V}_a. \quad (49)$$

The first term in (49) can be expressed as

$$\begin{aligned} \dot{V}^* &= \dot{V}^*(s) = \langle \nabla V^*(x), F + G u_a \rangle \\ &= \langle (\nabla \phi)^T W + \nabla \varepsilon(x), F - G \lambda \tanh(D_a) \rangle \\ &= W^T \nabla \phi F - W^T \nabla \phi G \lambda \tanh(D_a) + \varepsilon_1 \end{aligned} \quad (50)$$

where  $\varepsilon_1 = (\nabla \varepsilon)^T [F - G \lambda \tanh(D_a)]$  satisfying  $\|\varepsilon_1\| \leq b_{d\varepsilon} b_f \|x\| + \lambda b_{d\varepsilon} b_g$ . The HJB equation approximation (29) can be equivalently rewritten as

$$\varepsilon_{hjb} = W^T \nabla \phi [F - G \lambda \tanh(D)] + Q(s) + \Theta, \quad (51)$$

where

$$\Theta = W^T \nabla \phi g \lambda \tanh(D) + \lambda^2 R \log(1 - \tanh^2(D)), \quad (52)$$

with  $D$  as defined in (30). Then

$$W^T \nabla \phi F = W^T \nabla \phi G \lambda \tanh(D) - Q(s) - \Theta + \varepsilon_{hjb}, \quad (53)$$

Inserting (53) into (50) yields

$$\begin{aligned} \dot{V}^* &= -Q - \Theta + W^T \nabla \phi(s) G \lambda \tanh(D) \\ &\quad - W^T \nabla \phi G \lambda \tanh(D_a) + \varepsilon_1 + \varepsilon_{hjb}. \end{aligned} \quad (54)$$

Note that  $\Theta$  and  $Q$  are positive definite, there exists a positive constant  $q$  such that

$$-Q - \Theta \leq -Q \leq -q s^T s. \quad (55)$$

The third term in (54) can be upper bounded by

$$W^T \nabla \phi(x) G \lambda \tanh(D) \leq \lambda b_g b_{d\phi} \|W\| \quad (56)$$

Taking  $W = W_a + \tilde{W}_a$  into (54) yields

$$\begin{aligned} &-W^T \nabla \phi G \lambda \tanh(D_a) \\ &= -\tilde{W}_a^T \nabla \phi G \lambda \tanh(D_a) - \underbrace{W_a^T \nabla \phi G \lambda \tanh(D_a)}_{=2\lambda^2 R [D_a \tanh(D_a)] \geq 0} \\ &\leq -\tilde{W}_a^T \nabla \phi G \lambda \tanh(D_a) \end{aligned} \quad (57)$$

where the above inequality results from the fact that  $x^T \tanh(x) \geq 0, \forall x$ . Collecting (55) - (57), one has

$$\begin{aligned} \dot{V}^* &\leq -s^T q s + b_{d\varepsilon} b_f \|s\| + \lambda b_g b_{d\phi} \|W\| + \lambda b_{d\varepsilon} b_g + b_h \\ &\quad - \tilde{W}_a^T \nabla \phi G \lambda \tanh(D_a) \\ &= -s^T q s + M_s \|s\| + N_s - \tilde{W}_a^T \nabla \phi G \lambda \tanh(D_a), \end{aligned} \quad (58)$$

where  $M_s = b_{d\varepsilon} b_f$  and  $N_s = \lambda b_g b_{d\phi} \|W\| + \lambda b_{d\varepsilon} b_g + b_h$  are bounded.

Second, we analyze the term  $\dot{V}_c$  in (49). From (35) and (43), one has

$$\dot{\tilde{W}}_c^T = \alpha_c \frac{\sigma_a}{[1 + \sigma_a^T \sigma_a]^2} e_c. \quad (59)$$

From (44),  $e_c$  in the above equation can be further written as

$$\begin{aligned} e_c &= Q + \Theta_a + W_c^T \sigma_a \\ &= Q + \Theta_a + W_c^T \sigma_a \\ &\quad - Q - \Theta - W^T \nabla \phi [F - G \lambda \tanh(D)] + \varepsilon_{hjb}, \end{aligned}$$

where the second equality results from (29) and  $\Theta$  is defined in (52). Then,  $e_c$  can be equivalently written as

$$e_c = \Theta_a - \Theta + W_c^T \sigma_a - W^T \nabla \phi [F - G \lambda \tanh(D)] + \varepsilon_{hjb}, \quad (60)$$

where

$$\begin{aligned} \Theta_a - \Theta &= W_a^T \nabla \phi(s) g \lambda \tanh(D_a) - W^T \nabla \phi(s) g \lambda \tanh(D) \\ &\quad + \lambda^2 R \log(1 - \tanh^2(D_a)) - \lambda^2 R \log(1 - \tanh^2(D)). \end{aligned} \quad (61)$$

Note that the term  $\lambda^2 R \log(1 - \tanh^2(D))$  in (61) can be rewritten as

$$\begin{aligned} &\lambda^2 R \log(1 - \tanh^2(D)) \\ &= \lambda^2 R [\log 4 - 2D - 2 \log(1 + e^{-2D})], \end{aligned}$$

where  $-2 \log(1 + e^{-2D})$  can be approximated as [16]

$$-2 \log(1 + e^{-2D}) = 2D - 2D \text{sgn}(D) + \varepsilon_D,$$

where  $\|\varepsilon_D\| \leq \log 4$ . Then,

$$\lambda^2 R \log(1 - \tanh^2(D)) = \lambda^2 R [\log 4 - 2D \text{sgn}(D) + \varepsilon_D]. \quad (62)$$

Similarly,

$$\begin{aligned} &\lambda^2 R \log(1 - \tanh^2(D_a)) \\ &= \lambda^2 R [\log 4 - 2D_a \text{sgn}(D_a) + \varepsilon_{D_a}], \end{aligned} \quad (63)$$

where  $\|\varepsilon_{D_a}\| \leq \log 4$ . Consider (61) - (63), one has

$$\begin{aligned} \Theta_a - \Theta &= W_a^T \nabla \phi(s) g \lambda \tanh(D_a) - W^T \nabla \phi(s) g \lambda \tanh(D) \\ &\quad + \lambda^2 R [2D \text{sgn}(D) - 2D_a \text{sgn}(D_a) + \varepsilon_{D_a} - \varepsilon_D] \end{aligned} \quad (64)$$

The approximation of function  $x \text{sgn}(x)$  is investigated in [18] with

$$0 \leq x \text{sgn}(x) - x \tanh(\kappa x) \leq \frac{3.59}{\kappa} \quad (65)$$

Then, based on (30) and (41), one has

$$\begin{aligned} & \lambda^2 R [2D \text{sgn}(D) - 2D_a \text{sgn}(D_a)] \\ &= W^T \nabla \phi(x) g \lambda [\tanh(\kappa D) - \tanh(\kappa D_a)] \\ &+ \tilde{W}_a^T \nabla \phi(x) g \lambda \tanh(\kappa D_a) + \varepsilon_\kappa \end{aligned} \quad (66)$$

with approximation error satisfying  $0 \leq \varepsilon_\kappa \leq \frac{7.18}{\kappa}$ . Based on (64), (66), by adding and subtracting  $W^T \nabla \phi G \lambda \tanh(D_a)$  in (60), one has

$$\begin{aligned} e_c(t) &= -\tilde{W}_c^T(t) \sigma_a(t) + \tilde{W}_a^T(t) M_a(t) + N_a(t) \\ M_a(t) &= \nabla \phi(t) G(t) \lambda [\tanh(\kappa D_a(t)) - \tanh(D_a(t))] \\ N_a(t) &= W^T \nabla \phi(t) G(t) \lambda [\tanh(\kappa D(t)) - \tanh(\kappa D_a(t))] \\ &+ \lambda^2 R [\varepsilon_{D_a}(t) - \varepsilon_D(t)] + \varepsilon_{hjb}(t) + \varepsilon_\kappa(t) \end{aligned} \quad (67)$$

Based on Assumptions 1 and 3, both  $M_a$  and  $N_a$  are bounded. Substituting (67) into (59) yields,

$$\dot{\tilde{W}}_c = \alpha_c \frac{\sigma_a}{[1 + \sigma_a^T \sigma_a]^2} \left[ -\tilde{W}_c^T(t) \sigma_a + \tilde{W}_a^T(t) M_a + N_a \right]$$

Therefore,  $\dot{\tilde{W}}_c$  can be written as in (68), where the matrix  $Q_c > 0$  because the signal  $\frac{\sigma_a(t)}{1 + \sigma_a^T(t) \sigma_a(t)}$  satisfies the PE condition. Also, note that  $N_a$ ,  $M_{aj}$ , and  $N_{aj}$  are bounded. Then,  $M_1$  in (68) is also bounded.

Next, we give the upper bound of  $\dot{V}_a$ . Based on (45), differentiating  $V_a$  yields (69). Based on Assumption 1 - 3 and the definition of the actor learning error  $e_u$  in (47),  $M_2$  is also bounded.

Finally, collecting (58), (68) and (69), one can obtain

$$\begin{aligned} \dot{V}^* &\leq -s^T q s + M_s \|s\| + N_s \\ &- \tilde{W}_c^T(t) Q_c \tilde{W}_c(t) + \tilde{W}_c^T M_1 \\ &+ \tilde{W}_a^T(t) \left[ \frac{M_a M_a^T}{2} - Y \right] \tilde{W}_a(t) + \tilde{W}_a^T M_2. \end{aligned} \quad (70)$$

Using the matrix theory, if the parameter  $Y$  is sufficiently large such that  $Y - \frac{M_a M_a^T}{2} > 0$ . Then,  $\dot{V}$  is negative, provided that

$$\begin{aligned} \|s\| &\geq \frac{M_s}{2q} + \sqrt{\frac{M_s^2}{4q^2} + \frac{N_s}{q}} \\ \|\tilde{W}_c\| &\geq \frac{M_1}{\lambda_{\min}(Q_c)} \\ \|\tilde{W}_a\| &\geq \frac{M_2}{\lambda_{\min}\left(Y - \frac{M_a M_a^T}{2}\right)} \end{aligned}$$

This completes the proof.  $\blacksquare$

#### IV. SIMULATION STUDY

In order to validate the effectiveness of the presented actor-critic-barrier structure for solving Problem 1, the van der pol oscillator will be used in this section. Consider the controlled Van der Pol oscillator with the dynamics

$$\dot{x} = \begin{bmatrix} x_2 \\ -x_1 + 0.5(1 - x_2^2)x_2 \end{bmatrix} + \begin{bmatrix} 0 \\ x_1 \end{bmatrix} u. \quad (71)$$

In this scenario, it is desired that the state  $x = [x_1 \ x_2]^T$  satisfies the following constraints,

$$x_1 \in (-0.6, 0.2), x_2 \in (-0.2, 0.5) \quad (72)$$

According to the converse HJB method [19], when the performance parameters are selected as  $Q = I_{2 \times 2}$  and  $R = 1$ , the optimal controller would be  $u^*(x) = -x_1 x_2$ . When applying this optimal control policy to the system (71), the evolution is shown in Figure 2, where the solid lines represent the state evolutions and the dashed lines denote the asymmetric bounds for the states. The system state trajectories in two-dimensional space is shown in Figure 4, where the black box denotes the safety region. It is desired to drive the system states to the origin without violating the safety constraints. Even the states are regulated to the origin asymptotically, the full state constraints can not be guaranteed.

Next, we apply the safe reinforcement learning algorithm with the actor-critic-barrier structure developed in Theorem 2, where the corresponding results are given in Figure 3. In this case, we start the system from the same initial condition. The system state trajectories in two-dimensional space is shown in Figure 4. In contrast to the previous case, one can observe that the state approach to the origin without violating the full state constraints. Both the full-state constraints and the closed-loop stability can be guaranteed when using the actor-critic-barrier learning algorithm.

#### V. CONCLUSION

This paper presents a novel barrier function based system transformation to deal with both symmetric and asymmetric constraints on full-state. With the barrier function, the regulation problem with full-state constraints is transformed into an unconstrained optimal control problem. It is shown that the solution of the unconstrained optimal control problem is guaranteed to tackle the full-state constrained regulation problem. Then, the barrier function is combined with the actor-critic structure to solve the unconstrained optimal control problem in an online fashion. It is shown that additional barrier function to the actor-critic structure guarantees the constraints would not be violated by the adaptive optimal controller design. Boundedness and stability of signals in the closed-loop system is analyzed, with the convergence to the optimal control policy. The efficacy of the presented approach is demonstrated using a simulation example in the end.

#### REFERENCES

- [1] W. He, Z. Li, and C. L. P. Chen, "A survey of human-centered intelligent robots: issues and challenges," *IEEE/CAA Journal of Automatica Sinica*, vol. 4, no. 4, pp. 602–609, 2017.
- [2] M. Abu-Khalaf and F. L. Lewis, "Nearly optimal control laws for nonlinear systems with saturating actuators using a neural network HJB approach," *Automatica*, vol. 41, no. 5, pp. 779 – 791, 2005.
- [3] M. Chen, S. S. Ge, and B. Ren, "Adaptive tracking control of uncertain mimo nonlinear systems with input constraints," *Automatica*, vol. 47, no. 3, pp. 452 – 465, 2011.
- [4] M. Abu-Khalaf and F. L. Lewis, "Nearly optimal state feedback control of constrained nonlinear systems using a neural networks HJB approach," *Annual Reviews in Control*, vol. 28, no. 2, pp. 239 – 251, 2004.
- [5] K. P. Tee, S. S. Ge, and E. H. Tay, "Barrier lyapunov functions for the control of output-constrained nonlinear systems," *Automatica*, vol. 45, no. 4, pp. 918 – 927, 2009.

$$\begin{aligned}\dot{V}_c &= -\tilde{W}_c^T(t) \frac{\sigma_a \sigma_a^T}{(1 + \sigma_a^T \sigma_a)^2} \tilde{W}_c(t) + \tilde{W}_c^T(t) \frac{\sigma_a}{[1 + \sigma_a^T \sigma_a]^2} M_a^T \tilde{W}_a(t) + \tilde{W}_c^T \frac{\sigma_a}{[1 + \sigma_a^T \sigma_a]^2} N_a \\ &\leq -\tilde{W}_c^T(t) \underbrace{\left[ \frac{\sigma_a \sigma_a^T}{(1 + \sigma_a^T \sigma_a)^2} - \frac{\sigma_a \sigma_a^T}{2(1 + \sigma_a^T \sigma_a)^4} \right]}_{Q_c} \tilde{W}_c(t) + \frac{1}{2} \tilde{W}_a^T(t) M_a M_a^T \tilde{W}_a(t) + \tilde{W}_c^T \underbrace{\frac{\sigma_a}{(1 + \sigma_a^T \sigma_a)^2} N_a}_{:=M_1}\end{aligned}\quad (68)$$

$$\begin{aligned}\dot{V}_a &= \tilde{W}_a^T [\nabla \phi G \lambda \tanh(D_a) - \nabla \phi G \lambda \tanh(D_c) + \nabla \phi G \tanh^2(D_a) e_u + YW - Y\tilde{W}_a] \\ &= -\tilde{W}_a^T Y \tilde{W}_a + \tilde{W}_a^T \nabla \phi G \lambda \tanh(D_a) + \tilde{W}_a^T \underbrace{\left[ -\nabla \phi G \lambda \tanh(D_c) + \tilde{W}_a^T \nabla \phi G \tanh^2(D_a) e_u + \tilde{W}_a^T YW \right]}_{:=M_2}\end{aligned}\quad (69)$$

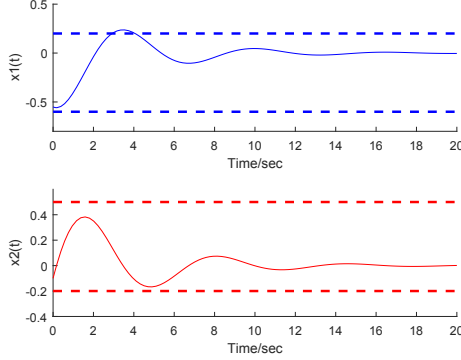


Fig. 2. State trajectories when using the converse HJB approach [19].

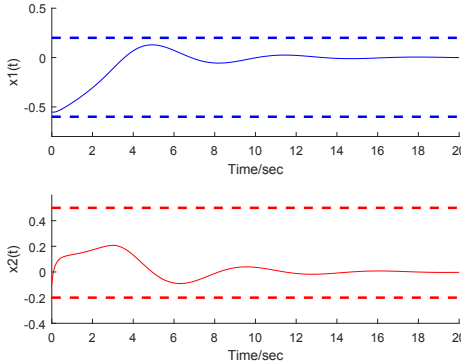


Fig. 3. State trajectories when applying actor-critic-barrier algorithm.

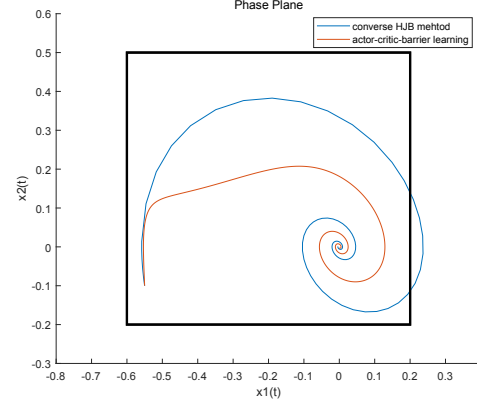


Fig. 4. Two-dimensional phase plot of state trajectories using the converse HJB approach [19] and the actor-critic-barrier algorithm. The black box denotes the safety region.

- [6] W. He, Y. Chen, and Z. Yin, "Adaptive neural network control of an uncertain robot with full-state constraints," *IEEE Transactions on Cybernetics*, vol. 46, no. 3, pp. 620–629, March 2016.
- [7] Y. Yang, D. Wunsch, and Y. Yin, "Hamiltonian-driven adaptive dynamic programming for continuous nonlinear dynamical systems," *IEEE transactions on neural networks and learning systems*, vol. 28, no. 8, pp. 1929–1940, 2017.
- [8] K. G. Vamvoudakis and F. L. Lewis, "Online actor-critic algorithm to solve the continuous-time infinite horizon optimal control problem," *Automatica*, vol. 46, no. 5, pp. 878 – 888, 2010.
- [9] H. Modares and F. L. Lewis, "Optimal tracking control of nonlinear partially-unknown constrained-input systems using integral reinforcement learning," *Automatica*, vol. 50, no. 7, pp. 1780 – 1792, 2014.
- [10] Y. Yang, H. Modares, D. C. Wunsch, and Y. Yin, "Leader-follower output synchronization of linear heterogeneous systems with active leader using reinforcement learning," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 29, no. 6, pp. 2139–2153, June 2018.
- [11] Y. Yang, H. Modares, D. C. Wunsch, and Y. Yin, "Optimal containment control of unknown heterogeneous systems with active leaders," *IEEE Transactions on Control Systems Technology*, 2018, to Appear.
- [12] Y. Yang, Z. Guo, H. Xiong, D. Ding, Y. Yin, and D. C. Wunsch, "Data-driven robust control of discrete-time uncertain linear systems via off-policy reinforcement learning," *IEEE Transactions on Neural Networks and Learning Systems*, 2019, to Appear.
- [13] D. Vrabie, K. G. Vamvoudakis, and F. L. Lewis, *Optimal Adaptive Control and Differential Games by Reinforcement Learning Principles*. IET Press, 2012.
- [14] K. G. Vamvoudakis, H. Modares, B. Kiumarsi, and F. L. Lewis, "Game theory-based control system algorithms with real-time reinforcement learning: How to solve multiplayer games online," *IEEE Control Systems Magazine*, vol. 37, no. 1, pp. 33–52, Feb 2017.
- [15] Y. Yang, K. G. Vamvoudakis, H. Ferraz, and H. Modares, "Dynamic intermittent Q-learning-based model-free suboptimal co-design of  $\mathcal{L}_2$ -stabilization," *International Journal of Robust and Nonlinear Control*, 2019, to Appear.
- [16] H. Modares, F. L. Lewis, and M.-B. Naghibi-Sistani, "Adaptive optimal control of unknown constrained-input systems using policy iteration and neural networks," *IEEE transactions on neural networks and learning systems*, vol. 24, no. 10, pp. 1513–1525, 2013.
- [17] K. G. Vamvoudakis, M. F. Miranda, and J. P. Hespanha, "Asymptotically stable adaptive-optimal control algorithm with saturating actuators and relaxed persistence of excitation," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 27, no. 11, pp. 2386–2398, Nov 2016.
- [18] M. M. Polycarpou and P. A. Ioannou, "A robust adaptive nonlinear control design," in *1993 American Control Conference*, June 1993, pp. 1365–1369.
- [19] V. Nevisti and J. A. Primbs, "Constrained nonlinear optimal control: a converse HJB approach," *California Institute of Technology*, 1996.