

引用格式:蒲俊,马清亮,李远冬,等.基于数据驱动自适应动态规划的输入约束的非线性系统 H_∞ 控制[J].电光与控制,2019,26(7):40-45. PU J, MA Q L, LI Y D, et al. H_∞ control of nonlinear systems with input constraints based on data-driven adaptive dynamic programming[J]. Electronics Optics & Control, 2019, 26(7):40-45.

基于数据驱动自适应动态规划的输入约束的非线性系统 H_∞ 控制

蒲俊, 马清亮, 李远冬, 顾凡
(火箭军工程大学, 西安 710025)

摘要:提出了一种包含在线采样、离线学习两个阶段的基于数据驱动的迭代自适应动态规划(ADP)算法,仅通过在线数据,解决输入约束的连续未知模型的非线性系统的 H_∞ 控制问题。通过策略迭代(PI)和迭代强化学习(IRL)方法推导出无模型(HJI)方程。构建3个神经网络,在线采集系统数据结束后,利用离线学习方法,近似求解无模型HJI方程,进而得到值函数、控制策略和扰动策略,神经网络的未知参数通过最小二乘方法求解。仿真结果验证了算法的可行性。

关键词: 自适应动态规划; H_∞ 控制; 输入约束; 最优控制; 神经网络

中图分类号: TP13 文献标志码: A doi:10.3969/j.issn.1671-637X.2019.07.008

H_∞ Control of Nonlinear Systems with Input Constraints Based on Data-Driven Adaptive Dynamic Programming

PU Jun, MA Qing-liang, LI Yuan-dong, GU Fan
(Rocket Force University of Engineering, Xi'an 710025, China)

Abstract: A data-driven iterative Adaptive Dynamic Programming (ADP) algorithm including online sampling and offline learning is proposed. The H_∞ control problem of the nonlinear system with input constraints and unknown models is solved only by online data. The model-free Hamilton-Jacobi-Isaacs (HJI) equation is derived by using the methods of Policy Iteration (PI) and Iteration Reinforcement Learning (IRL). Three neural networks are constructed. After the online acquisition of system data is completed, the off-line learning method is used to approximately solve the model-free HJI equation, and then the value function, control strategy and disturbance strategy are obtained. The parameter of the neural network is solved by the least squares method. Simulation results verify the feasibility of the algorithm.

Key words: adaptive dynamic programming; H_∞ control; input constraints; optimal control; neural network

0 引言

最优控制一直是控制理论非常活跃的分支,其目的就是要找到一个控制策略,使得所定义的性能指标最小^[1],而在各类控制应用中,系统中存在大量扰动,给系统带来负面影响, H_∞ 控制提供了一个有力的工

具减少扰动的影响^[2]。根据博弈论的思想, H_∞ 控制器的设计相当于一个双玩家的零和博弈游戏(Zero-Sum Game, ZSG)^[3],控制器在最坏扰动下最小化性能指标,达到最优控制的目的。对于连续的非线性系统,可以通过求解 Hamilton-Jacobi-Isaacs (HJI) 方程获得 ZSG 解,但是非线性系统 HJI 方程是非线性偏微分方程,求解几乎是不可能的。

迭代自适应动态规划(Adaptive Dynamic Programming, ADP)技术融合了最优控制、自适应控制、强化学习理论,利用函数近似结构估计值函数,这就为近似求解 HJI 方程提供了有效途径^[4]。WU 等^[5]利用强化学

收稿日期:2018-07-11

修回日期:2018-08-29

基金项目:国家自然科学基金(61773387)

作者简介:蒲俊(1994—),男,四川绵阳人,硕士生,研究方向为先导控制理论、非线性系统控制。

习的思想,推导出不依赖系统模型的有效计算同步策略更新算法,设计 H_∞ 反馈控制器;YASINI 等^[6]融合策略迭代(Policy Iteration, PI)和迭代强化学习(Iteration Reinforcement Learning, IRL)算法解决部分系统模型未知的双玩家 ZSG 问题;ZHANG 等^[3]提出迭代自适应动态规划方法解决一类非线性系统零和博弈问题;LUO 等^[7]设计基于数据驱动的 H_∞ 控制器,其设计方法不再需要系统的模型信息,但是 PI 算法在近似求解 HJI 方程时需要初始的容许控制策略,其对于 HJI 方程往往是次优;YANG 等^[8]用 ADP 算法设计输入约束的非线性系统 H_∞ 控制器,但依赖系统的模型信息,而实际系统不可避免地会受到多种因素的影响,很大程度上造成其模型未知或部分未知;ZHU 等^[9]提出迭代 ADP 算法求解无模型 H_∞ 控制问题,但没有考虑输入受约束情况,系统由于工作环境和本身性质,其控制输入受到约束的情况非常常见。

受此启发,本文针对一类输入受约束的非线性系统,基于 ADP 技术,并结合博弈论和数据驱动思想,融合 PI 和 IRL 算法推导出无模型 HJI 方程,求解非线性系统 H_∞ 控制问题。

1 问题描述

考虑如下连续非线性系统

$$\begin{cases} \dot{x} = f(x) + g(x)u(t) + k(x)w(t) \\ z = h(x) \end{cases} \quad (1)$$

式中: $x \in \mathbf{R}^n$ 是系统状态; $u(t) \in \mathbf{R}^m$ 是控制输入; $w(t) \in \mathbf{R}^q$ 是扰动输入,满足 $w(t) \in L_2[0, \infty)$; $z \in \mathbf{R}^p$ 是虚拟输出; $f(x) \in \mathbf{R}^n$, $g(x) \in \mathbf{R}^{n \times m}$ 和 $k(x) \in \mathbf{R}^{n \times q}$ 是系统内部模型,假设系统模型都是未知的, $f(x)$, $g(x)$ 和 $k(x)$ 为 Lipschitz 连续且 $f(0) = 0$, 因此 $x = 0$ 是式(1)系统的一个平衡点。在后续的计算式中为了方便表示,用 u 和 w 分别表示 $u(t)$ 和 $w(t)$ 。

H_∞ 控制就是要找到一个控制策略使得如下性能指标,对所有的 $w \in L_2[0, \infty)$ 和 $x(0) = 0$ 都是非正的,即

$$J_1(x(0), u, w) = \int_0^\infty (h(x)^T h(x) + u^T R u - \gamma^2 w^T w) dt \quad (2)$$

$R > 0, \gamma \geq \gamma^* \geq 0$, 如果这个控制策略存在,就称系统具有 $L_2 - g_{\text{gain}} \leq \gamma, \gamma^*$ 表示使该问题有解的最小值。

受文献[10]启发,对于控制输入 u 受约束情况,引用广义的非二次函数 $Y(u) = 2 \int_0^u \lambda \tanh^{-1}(s/\lambda) R ds$, $R > 0$ 为对角矩阵。只关注反馈策略和完整的状态信息,给出一个控制策略 $u(t) = u(x(t))$ 和一个扰动策略 $w(t) = w(x(t))$, 定义值函数为

$$V(x(0)) = \int_0^\infty (h(x)^T h(x) + Y(u) - \gamma^2 w^T w) dt = \int_0^\infty r(x(t), u, w) dt \quad (3)$$

定义 Lyapunov 方程为

$$r(x, u, w) + \nabla V^T (f + gu + kw) = 0 \quad V(0) = 0 \quad (4)$$

式中,符号 ∇ 为求偏导数,即 $\nabla V = \partial V / \partial x$ 。

定义 Hamiltonian 函数为

$$H(x, \nabla V, u, w) \equiv r(x, u, w) + \nabla V^T (f + gu + kw) \quad (5)$$

根据博弈理论,求解 H_∞ 控制问题相当于求解一个双玩家 ZSG 问题。控制策略使性能指标最大化,而扰动策略让性能指标最小化。对于连续非线性 ZSG 就是要找到最优的反馈控制策略和扰动策略,即

$$V^* = \min_u \max_w V(x, u, w) \quad (6)$$

假如满足下列纳什均衡条件

$$\min_u \max_w V(x, u, w) = \min_w \max_u V(x, u, w) \quad (7)$$

ZSG 问题有唯一解,即存在鞍点 (u^*, w^*) , u^* 即为所求系统式(1)的 H_∞ 最优控制策略, w^* 为最坏扰动策略,进一步根据稳定条件, u^* 和 w^* 分别表示为

$$u^* = -\lambda \tanh\left(\frac{1}{2\lambda} R^{-1} g^T \nabla V^*\right) \quad (8)$$

$$w^* = \frac{1}{2} \gamma^{-2} k^T \nabla V^* \quad (9)$$

将式(8)、式(9)代入式(4),得 HJI 方程为

$$(\nabla V^*)^T f + h(x)^T h(x) + \lambda^2 \bar{R} \ln(\underline{1} - \tanh^2(D^*)) + \frac{1}{4\gamma^2} (\nabla V^*)^T k(x) k(x)^T \nabla V^* = 0 \quad V^*(0) = 0 \quad (10)$$

式中: $\bar{R} = (s_1 \ \dots \ s_m) \in \mathbf{R}^m$; $D^* = \frac{1}{2\lambda} R^{-1} g^T \nabla V^*$; $\underline{1}$ 是一个所有元素都等于 1 的列向量。

假设 1 选择 $\gamma > 0$ 且式(1)系统是零状态可观的。在集合 $\Omega \in \mathbf{R}^n$ 上存在一个控制策略 $u(x)$ 使系统渐近稳定且有 $L_2 - g_{\text{gain}} \leq \gamma$, 则式(10) HJI 方程在 Ω 上存在一个光滑解。

此处,假设 1 保证了非线性系统 H_∞ 控制问题的有解性。

2 基于 ADP 的无模型策略迭代算法

非线性系统的 HJI 方程是一个非线性的偏微分方程,其解析解不一定存在,可用策略迭代方法得到其近似解,算法步骤如下。

1) 分别给定初始稳定的控制策略 u_1 和扰动策略 w_1 , 此时 $i = 1$ 。

2) 策略评价。依据已知 u_1 和 w_1 , 利用如下 Lyapunov 方程求解 V_i , 即

$$r(\mathbf{x}, \mathbf{u}_i, \mathbf{w}_i) + \nabla V_i^T(\mathbf{f} + \mathbf{g}\mathbf{u}_i + \mathbf{k}\mathbf{w}_i) = 0 \quad V_i(0) = 0 \quad (11)$$

3) 策略提高。\$i = i + 1\$, 更新 \$\mathbf{u}_i\$ 和 \$\mathbf{w}_i\$, 即

$$\begin{cases} \mathbf{u}_{i+1}(\mathbf{x}) = -\lambda \tanh(\mathbf{D}_{i+1}) \\ \mathbf{w}_{i+1}(\mathbf{x}) = \frac{1}{2\gamma^2} \mathbf{R}\mathbf{k}^T \nabla V_i \end{cases} \quad (12)$$

式中: \$\mathbf{D}_{i+1} = \frac{1}{2\lambda} \mathbf{R}^{-1} \mathbf{g}^T \nabla V_i\$。

4) 如果 \$\|V_{i+1} - V_i\| \leq \varepsilon, \varepsilon > 0\$, 则停止计算, 输出 \$V_i\$, 否则 \$i = i + 1\$, 并转至步骤 2)。

定理 1 假设 1 成立且满足 Kantorovich's 收敛条件^[11], 在式(11)和式(12)迭代相当于牛顿方法求解式(10)的 HJI 方程, 当 \$i \to \infty\$ 时, \$\mathbf{u}_i \to \mathbf{u}^*, \mathbf{w}_i \to \mathbf{w}^*, V_i \to V^*\$。\$V^*\$ 为值函数最优解, \$(\mathbf{u}^*, \mathbf{w}^*)\$ 为最优纳什均衡策略。

证明 文献[10]已进行详细证明, 本文不赘述。

式(4)是 \$V_i\$ 的线性偏微分方程。迭代求解式(4)相对于直接求解 HJI 方程更为可行, 但是迭代方法需要系统完整的模型信息, 有的采用 IRL 算法, 可以不依赖 \$\mathbf{f}\$, 但其他信息仍必要。但是在非线性系统考虑输入约束时具体模型往往是难以获得的。

通过模型依赖 PI 算法和 IRL 算法推导无模型的 ADP 算法, 给出两个随机的控制策略 \$\mathbf{u}\$ 和扰动策略 \$\mathbf{w}\$ 使式(1)系统在一个闭区间上稳定, \$\mathbf{u}_i, \mathbf{w}_i\$ 表示式(11)和式(12)第 \$i\$ 次迭代的结果, 进一步计算 \$V_i, \mathbf{u}_{i+1}\$ 和 \$\mathbf{w}_{i+1}\$。\$V_i\$ 对时间的导数 \$\dot{V}_i = \nabla V_i^T(\mathbf{f} + \mathbf{g}\mathbf{u} + \mathbf{k}\mathbf{w})\$, 根据式(11)和式(12)有

$$\begin{aligned} \dot{V}_i &= \nabla V_i^T(\mathbf{g}(\mathbf{u} - \mathbf{u}_i) + \mathbf{k}(\mathbf{w} - \mathbf{w}_i)) - r(\mathbf{x}, \mathbf{u}_i, \mathbf{w}_i) = \\ &2\lambda(\mathbf{D}_i)^T \mathbf{R}(\mathbf{u} - \mathbf{u}_i) + 2\gamma^2 \mathbf{w}_{i+1}^T (\mathbf{w} - \mathbf{w}_i) - r(\mathbf{x}, \mathbf{u}_i, \mathbf{w}_i) \end{aligned} \quad (13)$$

根据 IRL 算法, 同时在区间 \$[t, t + \Delta t]\$ 对式(13)两边积分得到

$$\begin{aligned} V_i(t + \Delta t) - V_i(t) - \int_t^{t+\Delta t} 2\lambda(\mathbf{D}_{i+1})^T \mathbf{R}(\mathbf{u} - \mathbf{u}_i) d\tau - \\ \int_t^{t+\Delta t} 2\gamma^2 \mathbf{w}_{i+1}^T (\mathbf{w} - \mathbf{w}_i) d\tau + \int_t^{t+\Delta t} r(\mathbf{x}, \mathbf{u}_i, \mathbf{w}_i) d\tau = 0 \end{aligned} \quad (14)$$

式中, \$V_i, \mathbf{D}_{i+1}\$ 和 \$\mathbf{w}_{i+1}\$ 是需要求解的未知函数或函数向量。

无模型迭代 ADP 算法主要思想是求解式(14)无模型方程来代替求解式(11)和式(12)的基于模型迭代方程。式(14)方程中不再包含系统的模型信息 \$\mathbf{f}(\mathbf{x}), \mathbf{g}(\mathbf{x})\$ 和 \$\mathbf{k}(\mathbf{x})\$, 而是利用系统的实时数据 \$\mathbf{u}_i\$ 和 \$\mathbf{w}_i\$。事实上系统的模型信息包含于可用的系统数据中。所以, 在用无模型策略学习方法近似求解式(14)方程前, 需要在线采样收集所需要的可用系统数据。

3 算法实现

使用 3 个神经网络(评价网络、控制网络、扰动网络)分别近似值函数 \$V_i\$, 控制策略 \$\mathbf{u}_{i+1}\$ 和扰动策略 \$\mathbf{w}_{i+1}\$。

根据 Weirstrass 高阶近似理论^[12], 式(14)的解 \$V_i\$ 用神经网络表示为

$$V_i^i(\mathbf{x}) = \mathbf{W}_{1,i+1}^T \boldsymbol{\phi}_1(\mathbf{x}) + \varepsilon_{1,i+1}(\mathbf{x}) \quad (15)$$

$$\mathbf{D}_{i+1}(\mathbf{x}) = \mathbf{W}_{2,i+1}^T \boldsymbol{\phi}_2(\mathbf{x}) + \varepsilon_{2,i+1}(\mathbf{x}) \quad (16)$$

$$\mathbf{w}_{i+1}(\mathbf{x}) = \mathbf{W}_{3,i+1}^T \boldsymbol{\phi}_3(\mathbf{x}) + \varepsilon_{3,i+1}(\mathbf{x}) \quad (17)$$

$$\mathbf{D}_i(\mathbf{x}) = \mathbf{W}_{2,i}^T \boldsymbol{\phi}_2(\mathbf{x}) + \varepsilon_{2,i}(\mathbf{x}) \quad (18)$$

$$\mathbf{w}_i(\mathbf{x}) = \mathbf{W}_{3,i}^T \boldsymbol{\phi}_3(\mathbf{x}) + \varepsilon_{3,i}(\mathbf{x}) \quad (19)$$

其中: \$\boldsymbol{\phi}_1: \mathbf{R}^n \to \mathbf{R}^{K_1}, \boldsymbol{\phi}_2: \mathbf{R}^n \to \mathbf{R}^{K_2}, \boldsymbol{\phi}_3: \mathbf{R}^n \to \mathbf{R}^{K_3}\$, 为线性独立基础函数向量; \$\mathbf{W}_{1,\cdot} \in \mathbf{R}^{K_1}, \mathbf{W}_{2,\cdot} \in \mathbf{R}^{K_2 \times m}, \mathbf{W}_{3,\cdot} \in \mathbf{R}^{K_3 \times q}\$ 为权值矩阵; \$\varepsilon_{1,\cdot}, \varepsilon_{2,\cdot}\$ 和 \$\varepsilon_{3,\cdot}\$ 是神经网络重建误差; \$K_1, K_2\$ 和 \$K_3\$ 是隐藏层神经元数。假设基函数参数和重建误差在 \$\Omega\$ 上, 当 \$K_1 \to \infty, K_2 \to \infty, K_3 \to \infty\$ 时, \$\varepsilon_{1,\cdot} \to 0, \varepsilon_{2,\cdot} \to 0, \varepsilon_{3,\cdot} \to 0\$。

为确定 \$\mathbf{W}_{1,i+1}, \mathbf{W}_{2,i+1}, \mathbf{W}_{3,i+1}\$ 的理想参数, 用 \$\hat{\mathbf{W}}_{1,i+1}, \hat{\mathbf{W}}_{2,i+1}, \hat{\mathbf{W}}_{3,i+1}\$ 表示理想权值矩阵, 神经网络表示为

$$\hat{V}_i(\mathbf{x}) = \hat{\mathbf{W}}_{1,i+1}^T \boldsymbol{\phi}_1(\mathbf{x}) \quad (20)$$

$$\hat{\mathbf{D}}_{i+1}(\mathbf{x}) = \hat{\mathbf{W}}_{2,i+1}^T \boldsymbol{\phi}_2(\mathbf{x}) \quad (21)$$

$$\hat{\mathbf{w}}_{i+1}(\mathbf{x}) = \hat{\mathbf{W}}_{3,i+1}^T \boldsymbol{\phi}_3(\mathbf{x}) \quad (22)$$

式中, \$\hat{\mathbf{W}}_{1,i+1}, \hat{\mathbf{W}}_{2,i+1}\$ 和 \$\hat{\mathbf{W}}_{3,i+1}\$ 是当前的估计值。

为多次迭代, 定义一个严格增的时间序列 \$\{t_k\}_{k=0}^l\$, \$l > 0\$ 为采集数据点的数量。由于 \$\hat{V}_i(\mathbf{x}), \hat{\mathbf{D}}_{i+1}(\mathbf{x}), \hat{\mathbf{w}}_{i+1}(\mathbf{x})\$ 替代 \$V_i(\mathbf{x}), \mathbf{D}_{i+1}(\mathbf{x}), \mathbf{w}_{i+1}(\mathbf{x})\$ 代入方程(14), 因此存在残差表示为

$$\begin{aligned} e_k &= \hat{V}_i(t_{k+1}) - \hat{V}_i(t_k) + \int_{t_k}^{t_{k+1}} r(\mathbf{x}, \mathbf{u}_i, \mathbf{w}_i) d\tau - \\ &\int_{t_k}^{t_{k+1}} 2\lambda(\hat{\mathbf{D}}_{i+1})^T \mathbf{R}(\mathbf{u} - \mathbf{u}_i) d\tau - \int_{t_k}^{t_{k+1}} 2\gamma^2 \mathbf{w}_{i+1}^T (\mathbf{w} - \mathbf{w}_i) d\tau = \\ &(\boldsymbol{\phi}_1(\mathbf{x}(t_{k+1})) - \boldsymbol{\phi}_1(\mathbf{x}(t_k)))^T \hat{\mathbf{W}}_{1,i+1} - 2\lambda \int_{t_k}^{t_{k+1}} [\boldsymbol{\phi}_2(\mathbf{x})^T \cdot \\ &\hat{\mathbf{W}}_{2,i+1} \mathbf{R}(\mathbf{u} + \lambda \tanh(\hat{\mathbf{W}}_{2,i}^T \boldsymbol{\phi}_2(\mathbf{x}))) - \int_{t_k}^{t_{k+1}} 2\gamma^2 \boldsymbol{\phi}_3(\mathbf{x})^T \cdot \\ &\hat{\mathbf{W}}_{3,i+1} (\mathbf{w} - \hat{\mathbf{W}}_{3,i+1}^T \boldsymbol{\phi}_3(\mathbf{x})) d\tau + \int_{t_k}^{t_{k+1}} [\mathbf{h}(\mathbf{x})^T \mathbf{h}(\mathbf{x}) - \\ &\boldsymbol{\phi}_3(\mathbf{x})^T \hat{\mathbf{W}}_{3,i+1} \hat{\mathbf{W}}_{3,i+1}^T \boldsymbol{\phi}_3(\mathbf{x}) + \\ &\int_0^{-\lambda \tanh(\mathbf{W}_{2,i}^T \boldsymbol{\phi}_2(\mathbf{x}))} (\lambda \tanh^{-1}(s/\lambda))^T \mathbf{R} ds] d\tau \end{aligned} \quad (23)$$

根据 Kronecker 积 \$\otimes\$, 式(23)可改写为

$$e_k = \boldsymbol{\theta}_k(\bar{\mathbf{W}}_i)^T \bar{\mathbf{W}}_{i+1} + \boldsymbol{\xi}(\bar{\mathbf{W}}_i) \quad (24)$$

式中: $\bar{W}_{i+1} = (\hat{W}_{1,i+1}^T, \text{vec } \hat{W}_{2,i+1}^T, \text{vec } \hat{W}_{3,i+1}^T)^T \in \mathbf{R}^{\bar{K}}$, $\bar{K} = K_1 + mK_2 + qK_3$, \bar{W}_k 可用相同方式表达, 符号 $\text{vec}(\cdot)$ 表示矩阵向量化, 迭代标志 $i \in \{0, 1, \dots\}$, 时间序列标志 $k \in \{0, \dots, l\}$, $\theta_k(\bar{W}_i)$ 和 $\xi_k(\bar{W}_i)$ 可以分别定义为

$$\theta_k(\bar{W}_i) = \begin{pmatrix} \phi_1(x(t_{k+1})) - \phi_1(x(t_k)) \\ - \int_{t_k}^{t_{k+1}} 2\lambda R(u + \lambda \tanh(\hat{W}_{2,i}^T \phi_2)) \otimes \phi_2 d\tau \\ - \int_{t_k}^{t_{k+1}} 2\gamma^2(w - \hat{W}_{3,i}^T \phi_3) \otimes \phi_3 d\tau \end{pmatrix} \quad (25)$$

$$\xi_k(\bar{W}_i) = \int_{t_k}^{t_{k+1}} [h(x)^T h(x) - \phi_3(x)^T \hat{W}_{3,i+1} \hat{W}_{3,i+1}^T \phi_3(x) + \int_0^{-\lambda \tanh(\hat{W}_{2,i}^T \phi_2(x))} (\lambda \tanh^{-1}(s/\lambda))^T R ds] d\tau \quad (26)$$

根据最小二乘 (Least-Squares, LS) 原理找到一组权值使残差最小, 即

$$\min_{\bar{w}_{i+1}} \sum_{k=0}^{l-1} e_k^2 \quad (27)$$

假设 2 持续激励 (Persistence of Excitation, PE) 对于 $i \geq 0$ 存在 $l_0 > 0$ 和 $\delta > 0$, 因此对所有的 $l \geq l_0$ 有

$$\frac{1}{l} \sum_{k=0}^{l-1} \theta_k(\bar{W}_i) \theta_k(\bar{W}_i)^T \geq \delta I_{\bar{K}} \quad (28)$$

式中, $I_{\bar{K}}$ 是适当维数的单位矩阵。

根据 $\theta_k(\bar{W}_i)$ 的定义, 为了保证 PE 条件, u 和 w 需要与 \hat{u}_{i+1} 和 \hat{w}_{i+1} 大不相同而且状态需要被持续激励, 所以 u 和 w 设计为深度探测信号。LS 表达如下

$$\bar{W}_{i+1} = -(\Theta(\bar{W}_i)^T \Theta(\bar{W}_i))^{-1} \Theta(\bar{W}_i)^T \Xi(\bar{W}_i) \quad (29)$$

其中,

$$\Theta(\bar{W}_i) = (\theta_0(\bar{W}_i) \quad \dots \quad \theta_{l-1}(\bar{W}_i))^T; \quad (30)$$

$$\Xi(\bar{W}_i) = (\xi_0(\bar{W}_i) \quad \dots \quad \xi_{l-1}(\bar{W}_i))^T \quad (31)$$

当计算得到 \bar{W}_{i+1} , 它就替换 \bar{W}_i 开始下面的迭代。

所提出的 ADP 算法实际上是无模型策略学习方法, 通过给定适当的初始策略权值和在线采集系统数据计算出 $\Theta(\bar{W}_k)$ 和 $\Xi(\bar{W}_k)$, 而后用式 (25) 迭代求解 \bar{W}_{k+1} , 最后利用式 (20) ~ 式 (22) 近似求解 $\hat{V}_i(x)$, $\hat{D}_{i+1}(x)$ 和 $\hat{w}_{i+1}(x)$, 进而得到 H_∞ 控制器。

为了式 (29) 中矩阵 $\Theta(\bar{W}_i)^T \Theta(\bar{W}_i)$ 的逆存在, 矩阵 $\Theta(\bar{W}_i)$ 是列满秩的, 在实际计算中, 采样次数 l 需要满足 $l \geq \text{rank}(\Theta(\bar{W}_i))$, 即 $l \geq K_1 + mK_2 + qK_3$ 。

图 1 为 ADP 算法的流程图, 包括两个阶段, 在采集阶段时, 输入深度控制信号和扰动信号, 在线采集系统数据。充分时间后, 算法转到学习阶段, 训练神经网络权值, 采用离线迭代收敛方式使权值收敛到固定值,

如果没有收敛, 则需要重新返回第一阶段采集更多数据, 最终, 控制网络输出 H_∞ 控制器, 其中 $u' \in \Omega$ 为随机容许控制策略。

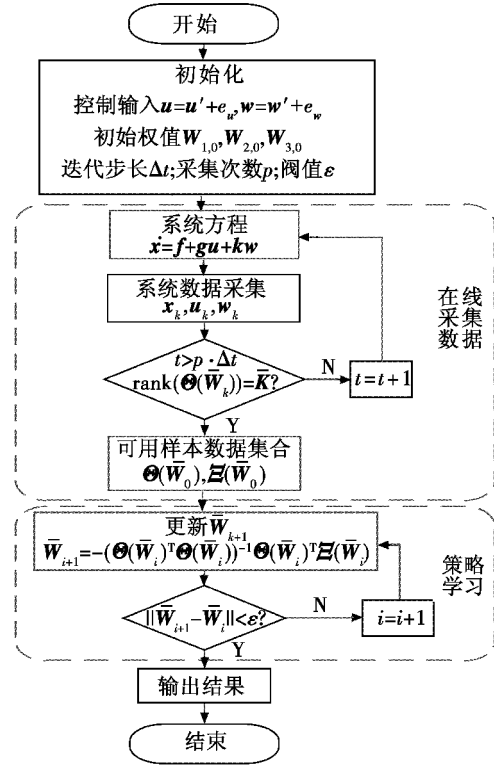


图 1 ADP 算法流程图

Fig. 1 Flow chart of ADP algorithm

定理 2 如果假设 2 成立, 对 $\forall \epsilon > 0$ 存在 $i^* > 0$, $K_1^* > 0, K_2^* > 0$ 和 $K_3^* > 0$, 那么如果 $i^* > i, K_1 > K_1^*, K_2 > K_2^*$ 和 $K_3 > K_3^*$, 则

$$\begin{cases} |\hat{V}_i(x) - V_i(x)| \leq \epsilon \\ \|\hat{D}_i - D_i\| \leq \epsilon \\ \|\hat{w}_i - w_i\| \leq \epsilon \\ |\hat{V}_i(x) - V^*(x)| \leq \epsilon \\ \|\hat{D}_i - D^*\| \leq \epsilon \\ \|\hat{w}_i - w^*\| \leq \epsilon \end{cases} \quad (32)$$

对所有 $x \in \Omega$ 成立, 类似的结论证明详见文献 [12]。

4 算例仿真

本章通过 Matlab 软件仿真结果说明本文算法的有效性。

考虑非线性系统

$$\dot{x} = f(x) + g(x)u + k(x)w \quad (33)$$

式中: $f(x) = \begin{pmatrix} -0.25x_1 \\ 0.5x_2^3 - 0.25\gamma^{-2}x_1^2x_2 - 0.25x_2 \end{pmatrix}; g(x) = (0 \quad x_2)^T; k(x) = (0 \quad x_1)^T; x = (x_1 \quad x_2)^T \in \mathbf{R}^2$ 为系统

状态, $u \in \mathbb{R}^m$ 和 $w \in \mathbb{R}^r$ 分别为控制输入和扰动输入。

选择 $h(x) = (x_1 \ x_2)^T$, $R = I$, $\gamma = 4$, I 是单位矩阵。控制输入 u 被约束为 $\|u\| \leq 1$, 那么 $Y(u)$ 可以定义为

$$Y(u) = 2 \int_0^u \tanh^{-1}(s) R ds \quad (34)$$

根据近似式(21)和式(22), 选择3个3层的前馈神经网络, 其基函数向量定义为 $\phi_1(1) = \phi_2(1) = \phi_3(1) = (x_1^2 \ x_1 x_2 \ x_2^2 \ x_1^4 \ x_1^4)^T$ 。

初始状态设置为 $x_0 = (0.5 \ -0.5)^T$, 选择深度控制输入 u 和扰动输入 w , 即

$$u = -\tanh(\sin 1.1\pi t + \sin 1.4\pi t + \sin 1.2\pi t + \sin 2.9\pi t - 3.2\sin 3.2\pi t + 1), \quad (35)$$

$$w = 1.1\sin \pi t - \sin 1.5\pi t - \sin 1.8\pi t - \sin 2\pi t + 0.2 \quad (36)$$

收敛阈值为 $\varepsilon = 10^{-6}$, 积分计算步长为 0.1 s, 采样次数 $l = 50$, 所以, 无模型学习阶段在 5 s 时, 通过采样数据获得最优控制策略。神经网络初始迭代权值 $W_1^0 = W_2^0 = W_3^0 = (0 \ 0 \ 0 \ 0 \ 0)^T$ 。

图2为评价网络的权值 W_1 随迭代次数的变化; 图3为控制网络的权值 W_2 随迭代次数的变化; 图4为扰动网络的权值 W_3 随迭代次数的变化。5 s 后, 最优控制策略和最优扰动策略代替深度控制信号加在被控系统, 全局的系统状态曲线如图5所示, 图6为系统的控制输入曲线, 图7为系统的全部扰动输入曲线。

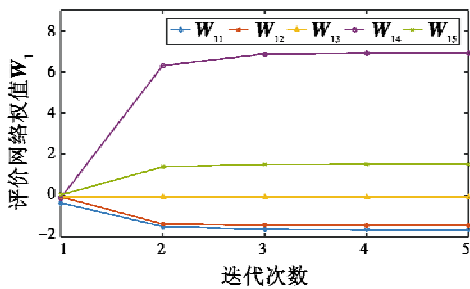


图2 W_1 参数更新

Fig.2 W_1 parameter updating

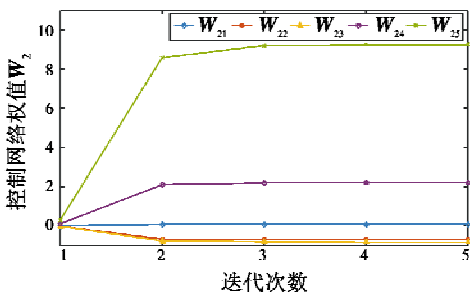


图3 W_2 参数更新

Fig.3 W_2 parameter updating

由图2~图4可以看出, 评价网络、控制网络和扰

动网络的权值 W_1, W_2, W_3 在 4 次迭代后收敛情况分别为 $W_1 = (-1.632 \ -1.424 \ -0.065 \ 6.919 \ 1.537)$, $W_2 = (0.071 \ -0.692 \ -0.838 \ 2.197 \ 9.272)$, $W_3 = (-0.043 \ 0.005 \ 0.005 \ 0.018 \ -0.147)$ 。

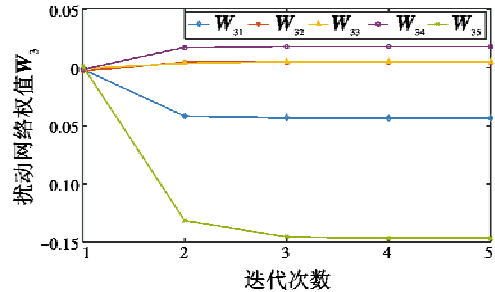


图4 W_3 参数更新

Fig.4 W_3 parameter updating

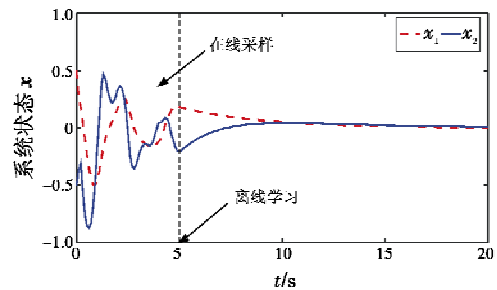


图5 系统状态

Fig.5 Curve of system state

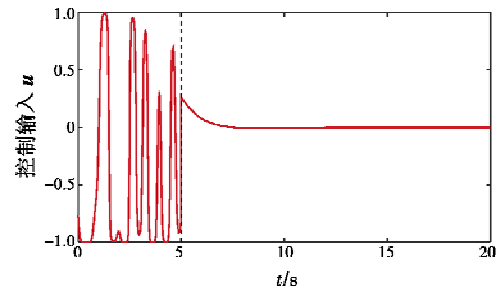


图6 控制输入 u

Fig.6 Control input u

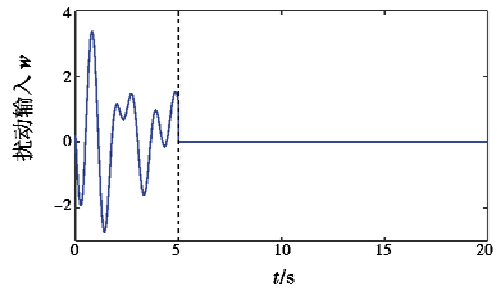


图7 扰动输入 w

Fig.7 Disturbance input w

进而可以得出鞍点为

$$u^*(x) = 0.071x_1^2 - 0.692x_1x_2 - 0.838x_2^2 + 2.197x_1^4 + 9.272x_1^4, \quad (37)$$

$$w^*(x) = -0.043x_1^2 + 0.005x_1x_2 + 0.005x_2^2 + 0.018x_1^4 - 0.147x_1^4. \quad (38)$$

由图5可以看出,5 s时,系统由采样阶段进入无模型策略学习阶段,5 s后,系统在所设计的控制器和最坏扰动条件下,逐渐收敛到零,验证了所得控制策略的合理性。从图6、图7可以看出,控制输入被约束在 $\|u\| \leq 1$,控制和扰动输入最终收敛到零。

5 结束语

本文针对一类内部动力模型未知的非线性系统,基于数据驱动的自适应动态规划思想,提出了一种通过在线采集系统数据,迭代求解输入约束 H_∞ 最优控制问题的无模型算法。该算法分为在线采集数据和离线学习两个部分,在离线学习时,构建3个神经网络,利用实时系统数据,近似求解无模型 HJI 方程,最后,算例仿真验证了算法的有效性。

参考文献

- [1] BERTSEKAS D P. Dynamic programming and optimal control[M]. Belmont: Athena Scientific, 1995.
- [2] 顾凡,马清亮,岳瑞华,等. 基于遗传算法的多项式非线性系统静态输出反馈 H_∞ 控制[J]. 电光与控制, 2017,24(10):55-58,89.
- [3] ZHANG H, WEI Q, LIU D. An iterative adaptive dynamic programming method for solving a class of nonlinear zero-sum differential game[J]. Automatica, 2011, 47(1):207-214.
- [4] LIU D R, YANG X, WANG D, et al. Reinforcement-learning-based robust controller design for continuous-time uncertain nonlinear systems subject to input constraints[J]. IEEE Transactions on Cybernetics, 2015, 45(7):1372-1385.
- [5] LUO B, WU H N. Computationally efficient simultaneous policy update algorithm for nonlinear H_∞ state feedback control with Galerkin's method[J]. International Journal of Robust and Nonlinear Control, 2013, 23(7):991-1012.
- [6] YASINI S, KARIMPOUR A, SISTANI M B, et al. Online concurrent reinforcement learning algorithm to solve two-player zero-sum games for partially unknown nonlinear continuous-time systems[J]. International Journal of Adaptive Control and Signal Processing, 2015, 29(4):473-493.
- [7] LUO B, HUANG T, WU H N, et al. Data-driven H_∞ control for nonlinear distributed parameter systems[J]. IEEE Transactions on Neural Networks & Learning Systems, 2015, 26(11):2949-2961.
- [8] YANG X, LIU D, WEI Q, et al. Adaptive dynamic programming for H_∞ control of constrained-input nonlinear systems[C]//Proceedings of the 34th Chinese Control Conference, IEEE, 2015:3027-3032.
- [9] ZHU Y H, ZHAO D B, LI X J, et al. Iterative adaptive dynamic programming for solving unknown nonlinear zero-sum game based on online data[J]. IEEE Transactions on Neural Networks & Learning Systems, 2017, 28(3):714-725.
- [10] ZHANG Q C, ZHAO D B, ZHU Y H. Data-driven adaptive dynamic programming for continuous-time fully cooperative games with partially constrained inputs[J]. Neurocomputing, 2017, 238:377-386.
- [11] TAPIA R A. The Kantorovich theorem for Newton's method[J]. The American Mathematical Monthly, 1971, 78(4):389-392.
- [12] FINLAYSON B A. The method of weighted residuals and variational principles[M]. New York: Academic Press, Inc., 2014.