



基于 Q 学习算法的随机离散时间系统的 随机线性二次最优追踪控制

摘要

针对随机线性离散时间系统,利用 Q 学习算法求解无限时域的随机线性二次最优追踪控制(SLQT)问题.首先,假设通过命令生成器生成追踪所需的参考信号,并建立一个由原随机系统和参考轨迹系统组成的增广系统,把最优追踪问题转化为最优调节问题的形式.其次,为了在线求解随机系统的最优追踪问题,将随机系统转为确定性系统,并根据增广系统定义随机线性二次最优追踪控制的 Q 函数,在无需知道系统模型参数的情况下在线求解增广随机代数方程(GSAE).再次,证明了 Q 学习算法和增广随机代数方程的等价性,给出了 Q 学习算法实现步骤.最后,给出一个仿真实例说明 Q 学习算法的有效性.

关键词

随机系统;Q 学习算法;最优追踪控制;随机代数方程

中图分类号 O232;TP13

文献标志码 A

收稿日期 2021-09-12

资助项目 国家自然科学基金(61873099,62073144);广东省自然科学基金(2020A1515010441);广州市科技计划(202002030158,202002030389)

作者简介

张正义,男,硕士生,研究方向为自适应动态规划、最优控制、强化学习.zhang810024143@163.com

赵学艳(通信作者),女,副教授,硕士生导师,主要从事随机系统和非线性系统的稳定性与镇定、复杂系统的建模、分析和控制的研究.auxyzhao@scut.edu.cn

0 引言

最优控制的目标是找到最优的控制策略,使得被控系统达到指定目标状态的同时,使系统预定义的性能指标为最小.最优控制问题主要有两个研究方向,分别是最优调节问题和最优追踪问题.对于线性系统的二次调节(Linear Quadratic Regulator, LQR)问题,传统方法通常是通过离线求解其对应的代数里卡蒂(Riccati)方程,这种方法需要完全已知系统参数的全部动力学信息^[1-2].但是,在实际情况下,系统动力学信息完全已知的条件难以满足,传统方法不可能得到解析解.所以,通常需要在系统参数未知的情况下在线求解最优控制器,因此利用自适应动态规划(Adaptive Dynamic Programming, ADP)和神经网络方法求解最优控制在近些年备受关注.自适应动态规划^[3]是在系统参数未知或系统参数不确定的情况下设计系统的控制器,不需要提前知道系统动力学信息,充分利用系统的状态信息在线求解最优控制.近些年来,ADP 方法在离散系统和连续系统中有了广泛的应用.文献[4]针对连续时间线性系统提出了自适应动态规划方法,在系统参数矩阵部分未知的情况下得到最优控制器;文献[5]进一步针对连续时间线性系统提出了一种自适应策略迭代方法,在系统参数完全未知的情况下得到最优控制器;文献[6]针对线性离散时间系统的追踪问题使用强化 Q 学习方法,在系统参数完全未知的情况下求解最优控制器.

随机系统控制理论由于其自身的学术难度以及广泛的应用领域,已成为控制理论的重要组成部分与研究热点^[7-8],尤其是随机系统的最优控制问题受到越来越多的关注.与确定性问题相似,随机系统的线性二次最优控制问题(Stochastic Linear Quadratic, SLQ)的可解性等价于随机代数 Riccati 方程的可解性,文献[9]研究了线性终端状态约束下不定随机线性二次最优控制问题,文献[10]研究了具有乘性噪声的随机离散系统的带约束线性二次最优控制问题,但是文献[9-10]需要完全已知的系统参数信息.因此,文献[11]针对随机连续时间系统在系统参数部分未知的情况下提出了策略迭代方法求解随机系统的最优控制问题,文献[12]针对系统参数完全未知的随机线性离散系统提出了使用自适应动态规划的方法求解最优控制问题,

¹ 华南理工大学 自动化科学与工程学院,广州,510640

文献[13]针对模型自由的随机线性离散系统提出了Q学习算法求解最优控制问题.相较于最优调节问题,最优追踪问题在现实中往往有更多的应用,例如文献[14]针对参数未知的随机离散系统提出了基于神经网络的自适应动态规划方法求解最优追踪控制问题.求解系统的最优控制问题,大多需要系统的完全动力学信息,使用Q学习算法的优点是不用直接求解复杂的随机代数方程,而是充分利用系统的状态信息在线求得系统的最优控制.受到文献[13-14]的启发,本文针对离散时间系统的随机线性二次最优追踪控制问题,提出了解决随机线性二次最优追踪控制的Q学习算法,给出算法的具体实现步骤,使用Q学习算法在线解决追踪控制问题而无需系统模型参数,最后给出仿真实例,表明系统输出可以有效地追踪参考轨迹.

本文的结构安排如下:第一节对问题进行描述,定义参考信号系统,将原随机系统和参考信号系统组成增广系统,把最优追踪问题转化为最优调节问题的形式;第二节对随机系统进行了问题转变,将随机系统转化为确定性系统;第三节推导了Q函数;第四节给出算法的具体实现步骤;第五节给出仿真实例;第六节对全文进行了总结.

注1 本文中: \mathbf{R} 表示实数集; \mathbf{R}^n 表示 n 维欧几里得空间; $\mathbf{R}^{n \times m}$ 表示全体 $n \times m$ 维实数矩阵的集合; $\|\cdot\|$ 表示 \mathbf{R}^n 中的欧几里得范数; $\|\mathbf{P}\|_F$ 表示假设 \mathbf{P} 为 $n \times m$ 维实矩阵,则其 F 范数定义为 $\|\mathbf{P}\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2}$; \mathbf{Z}^T 表示矩阵或向量 \mathbf{Z} 的转置; \mathbf{Z}^{-1} 表示方阵 \mathbf{Z} 的逆; \mathbf{I} 表示适当维度的单位矩阵; $\text{tr}(\mathbf{Z})$ 表示方阵 $\mathbf{Z} = (z_{ij})_{n \times n}$ 的迹, $\text{tr}(\mathbf{Z}) = \sum_{i=1}^n z_{ii}$; $\mathbf{E}(\cdot)$ 表示数学期望; $(\Omega, \mathcal{F}, \{\mathcal{F}_k\}_{k \geq k_0}, \mathcal{P})$ 表示具有流 $\{\mathcal{F}_k\}_{k \geq k_0}$ 的完备概率空间,其中 Ω 是样本点的全体, \mathcal{F} 是 Ω 上的一个 σ -代数, \mathcal{P} 是 \mathcal{F} 上的一个概率测度,滤子族 $\{\mathcal{F}_k\}_{k \geq k_0}$ 是一族单增的 \mathcal{F} 子 σ -代数; $\mathbf{w}_k (k = 0, 1, 2, \dots, \mathbf{w}_0 = 0)$ 是定义在完备概率空间上的一维Wiener过程,又称为布朗运动; $\mathbf{P} > 0$ 表示对称矩阵 \mathbf{P} 为正定矩阵.

1 问题描述

给定随机离散时间线性系统为

$$\begin{cases} \mathbf{x}_{k+1} = [\mathbf{A}\mathbf{x}_k + \mathbf{B}\mathbf{u}_k] + [\mathbf{C}\mathbf{x}_k + \mathbf{D}\mathbf{u}_k]\mathbf{w}_k, \\ \mathbf{y}_k = \mathbf{G}\mathbf{x}_k, \end{cases} \quad (1)$$

式中: $\mathbf{x}_k \in \mathbf{R}^n$ 表示系统的状态变量; \mathbf{x}_0 是系统的初始状态; $\mathbf{u}_k \in \mathbf{R}^m$ 表示系统控制输入; $\mathbf{y}_k \in \mathbf{R}^q$ 表示系统输出变量; $\mathbf{A} \in \mathbf{R}^{n \times n}, \mathbf{B} \in \mathbf{R}^{n \times m}, \mathbf{C} \in \mathbf{R}^{n \times n}, \mathbf{D} \in \mathbf{R}^{n \times m}, \mathbf{G} \in \mathbf{R}^{q \times n}$ 均为常数矩阵; \mathbf{w}_k 是定义在完备概率空间 $(\Omega, \mathcal{F}, \{\mathcal{F}_k\}_{k \geq k_0}, \mathcal{P})$ 上的一维标准布朗运动, $\mathbf{w}_0 = 0$,且满足 $\mathbf{E}(\mathbf{w}_{k+1} | \mathcal{F}_k) = 0, \mathbf{E}(\mathbf{w}_{k+1}^2 | \mathcal{F}_k) = 1$.

首先,对于追踪问题,需要定义追踪问题的参考信号.假设命令生成器生成的随机线性二次最优追踪控制的参考轨迹模型为

$$\mathbf{r}_{k+1} = \mathbf{F}\mathbf{r}_k, \quad (2)$$

式中: $\mathbf{r}_k \in \mathbf{R}^q$ 表示参考轨迹系统的输出变量; $\mathbf{F} \in \mathbf{R}^{q \times q}$ 为常数矩阵.为了生成更为广泛的参考轨迹,假设 \mathbf{F} 不必是赫尔维茨的.

对于追踪控制问题,需要重点关注系统的输出与需要追踪的参考信号输出之间的误差,因此定义系统的追踪误差为

$$\mathbf{e}_k = \mathbf{y}_k - \mathbf{r}_k, \quad (3)$$

式中 \mathbf{r}_k 为参考轨迹模型(2)的输出变量.

随机线性最优追踪问题的目标是设计出最优的控制器,不但能够确保系统的输出变量 \mathbf{y}_k 可以稳定地追踪上给定的参考信号 \mathbf{r}_k ,并且最小化预定义的性能指标函数.定义性能指标函数为

$$J(\mathbf{x}_k, \mathbf{r}_k, \mathbf{u}) = \mathbf{E} \sum_{i=k}^{\infty} U_i = \mathbf{E} \sum_{i=k}^{\infty} [(\mathbf{G}\mathbf{x}_i - \mathbf{r}_i)^T \mathbf{Q}(\mathbf{G}\mathbf{x}_i - \mathbf{r}_i) + \mathbf{u}_i^T \mathbf{R}\mathbf{u}_i], \quad (4)$$

式中: \mathbf{r}_k 为参考轨迹信号, \mathbf{u} 为控制序列, U_i 为 $i = k$ 时对应的性能指标函数, $\mathbf{Q} \in \mathbf{R}^{q \times q} > 0, \mathbf{R} \in \mathbf{R}^{m \times m} > 0. (\mathbf{G}\mathbf{x}_i - \mathbf{r}_i)^T \mathbf{Q}(\mathbf{G}\mathbf{x}_i - \mathbf{r}_i)$ 表示为追踪误差变量的二次型, \mathbf{Q} 为正定的权重系数矩阵,因此 \mathbf{Q} 矩阵越大表示误差变量对性能指标函数的影响越大; $\mathbf{u}_i^T \mathbf{R}\mathbf{u}_i$ 表示为系统控制输入变量的二次型, \mathbf{R} 为正定的权重系数矩阵,因此 \mathbf{R} 矩阵越大表示对系统控制的约束越大.

基于随机系统模型(1)和参考轨迹模型(2),建立一个由原随机系统和参考轨迹系统组成的增广系统,把最优追踪问题转化为最优调节问题的形式.

因此,随机线性离散时间系统(1)和参考轨迹系统(2)构建的增广系统如下:

$$\mathbf{X}_{k+1} = \begin{bmatrix} \mathbf{x}_{k+1} \\ \mathbf{r}_{k+1} \end{bmatrix} = \begin{bmatrix} \mathbf{A} + \mathbf{C}\mathbf{w}_k & \mathbf{O} \\ \mathbf{O} & \mathbf{F} \end{bmatrix} \begin{bmatrix} \mathbf{x}_k \\ \mathbf{r}_k \end{bmatrix} + \begin{bmatrix} \mathbf{B} + \mathbf{D}\mathbf{w}_k \\ \mathbf{O} \end{bmatrix} \mathbf{u}_k \equiv \mathbf{T}\mathbf{X}_k + \mathbf{B}_0\mathbf{u}_k, \quad (5)$$

式中: $T \in \mathbf{R}^{(n+q) \times (n+q)}$, $B_0 \in \mathbf{R}^{(n+q) \times m}$, 增广系统的状态为 $\mathbf{X}_k = \begin{bmatrix} \mathbf{x}_k \\ \mathbf{r}_k \end{bmatrix} \in \mathbf{R}^{(n+q)}$.

显然,只有当 F 为赫尔维茨(Hurwitz)时才可以 使用性能指标函数(4),即要求参考轨迹系统是渐近稳定的,如果参考轨迹随着时间的推迟不趋于零,那么性能指标函数(4)将是无界的.但是在实际应用中并不总是追踪渐近稳定的信号,而通常需要追踪正弦信号等其他复杂信号.为了放宽 F 为赫尔维茨的限制,在性能指标函数(4)中引入折扣系数 γ ,可以使得 F 为非赫尔维茨,因此基于(4)重新定义带折扣系数的性能指标函数为

$$J(\mathbf{x}_k, \mathbf{r}_k, \mathbf{u}) = \mathbb{E} \sum_{i=k}^{\infty} \gamma^{i-k} U_i = \mathbb{E} \sum_{i=k}^{\infty} \gamma^{i-k} [(\mathbf{G}\mathbf{x}_i - \mathbf{r}_i)^T \mathbf{Q}(\mathbf{G}\mathbf{x}_i - \mathbf{r}_i) + \mathbf{u}_i^T \mathbf{R}\mathbf{u}_i], \quad (6)$$

式中: $0 < \gamma \leq 1$ 是折扣系数,注意到只有当 F 为赫尔维茨时,此时参考轨迹为渐近稳定的, γ 可以取值为 1,性能指标函数由式(6)化简为式(4).

引入折扣系数 γ 后,基于增广系统(5)和性能指标函数(6),性能指标函数可进一步重写为

$$J(\mathbf{X}_k, \mathbf{u}) = \mathbb{E} \sum_{i=k}^{\infty} \gamma^{i-k} [\mathbf{X}_i^T \mathbf{Q}_1 \mathbf{X}_i + \mathbf{u}_i^T \mathbf{R}\mathbf{u}_i], \quad (7)$$

式中: $\mathbf{Q}_1 = [\mathbf{G} \quad -\mathbf{I}]^T \mathbf{Q} [\mathbf{G} \quad -\mathbf{I}] \in \mathbf{R}^{(n+q) \times (n+q)}$.

本文的目标是找到最优控制增益矩阵 \mathbf{K}^* 使得性能指标函数(7)达到最小,同时能够追踪参考信号.经过增广原系统,最优追踪问题已经转化为最优调节问题的形式,定义状态反馈线性控制器为

$$\mathbf{u}_k = \mathbf{K}\mathbf{X}_k, \quad (8)$$

式中: $\mathbf{K} \in \mathbf{R}^{m \times (n+q)}$, \mathbf{K} 表示系统的控制增益矩阵.

同时定义性能指标函数的最优值为

$$V^*(\mathbf{X}_0, \mathbf{u}) = \min_{\mathbf{u}} J(\mathbf{X}_0, \mathbf{u}). \quad (9)$$

定义 1^[14] 如果存在线性状态反馈控制器 \mathbf{u}_k , 对于任意初始状态,使得 $\lim_{k \rightarrow \infty} \mathbb{E}(\mathbf{e}_k^T \mathbf{e}_k) = 0$, 那么称 \mathbf{u}_k 是均方稳定的.若系统(5)存在均方稳定的控制器,该系统称为均方可稳的.

定义 2^[10] 如果 $-\infty < V^*(\mathbf{X}_0, \mathbf{u}) < +\infty$, 则称随机线性二次最优控制问题是适定的.

定义 3^[15] 如果线性反馈控制 \mathbf{u}_k 满足:1) \mathbf{u}_k 是关于 \mathcal{F}_k 适定并且可测的随机过程;2) \mathbf{u}_k 能够使得性能指标函数式(6)为最小值;3) \mathbf{u}_k 是均方稳定的.那么称 \mathbf{u}_k 是容许控制.

引理 1 如果存在均方稳定的容许控制 $\mathbf{u}_k = \mathbf{K}\mathbf{X}_k$, 那么随机线性二次最优追踪控制是适定的,并且值函数可重写为

$$V(\mathbf{X}_k, \mathbf{K}) = \mathbb{E}(\mathbf{X}_k^T \mathbf{P}\mathbf{X}_k), \quad (10)$$

式中: $\mathbf{P} \in \mathbf{R}^{(n+q) \times (n+q)}$, \mathbf{P} 为正定矩阵,且为式(11)增广随机代数方程的解.

$$\mathbf{P} = \gamma(\mathbf{A}_1 + \mathbf{B}_1 \mathbf{K})^T \mathbf{P}(\mathbf{A}_1 + \mathbf{B}_1 \mathbf{K}) + \gamma(\mathbf{C}_1 + \mathbf{D}_1 \mathbf{K})^T \mathbf{P}(\mathbf{C}_1 + \mathbf{D}_1 \mathbf{K}) + \mathbf{Q}_1 + \mathbf{K}^T \mathbf{R}\mathbf{K}, \quad (11)$$

式中:

$$\mathbf{A}_1 = \begin{bmatrix} \mathbf{A} & \mathbf{O} \\ \mathbf{O} & \mathbf{F} \end{bmatrix} \in \mathbf{R}^{(n+q) \times (n+q)},$$

$$\mathbf{B}_1 = \begin{bmatrix} \mathbf{B} \\ \mathbf{O} \end{bmatrix} \in \mathbf{R}^{(n+q) \times m},$$

$$\mathbf{C}_1 = \begin{bmatrix} \mathbf{C} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} \end{bmatrix} \in \mathbf{R}^{(n+q) \times (n+q)},$$

$$\mathbf{D}_1 = \begin{bmatrix} \mathbf{D} \\ \mathbf{O} \end{bmatrix} \in \mathbf{R}^{(n+q) \times m}.$$

证明 由式(7)和式(8)有:

$$V(\mathbf{X}_k, \mathbf{K}) = \mathbb{E} \sum_{i=k}^{\infty} \gamma^{i-k} [\mathbf{X}_i^T (\mathbf{Q}_1 + \mathbf{K}^T \mathbf{R}\mathbf{K}) \mathbf{X}_i],$$

拆开化简得:

$$V(\mathbf{X}_k, \mathbf{K}) = \mathbb{E}(\mathbf{X}_k^T (\mathbf{Q}_1 + \mathbf{K}^T \mathbf{R}\mathbf{K}) \mathbf{X}_k) + \gamma V(\mathbf{X}_{k+1}, \mathbf{K}),$$

代入式(10)可得:

$$\mathbb{E}(\mathbf{X}_k^T \mathbf{P}\mathbf{X}_k) = \mathbb{E}(\mathbf{X}_k^T (\mathbf{Q}_1 + \mathbf{K}^T \mathbf{R}\mathbf{K}) \mathbf{X}_k) + \gamma \mathbb{E}(\mathbf{X}_{k+1}^T \mathbf{P}\mathbf{X}_{k+1}),$$

代入增广系统(5),替换 \mathbf{X}_{k+1} 化简后得式(11).

对于随机线性二次最优追踪控制问题:随机线性二次最优追踪控制问题等价于求解其对应的增广随机代数 Riccati 方程的解,并且需要系统模型参数的全部信息.

2 问题转换

目前,确定性系统的最优追踪控制问题有着广泛的研究并且已经得到了很好的解决,随机系统因为随机参数的存在使得系统输出轨迹存在不确定性,且性能指标函数带有期望,在线算法无法实现期望功能.因此本节通过系统转变将随机系统转变为确定性系统,进而将随机系统的最优追踪控制问题转化为确定性的系统最优追踪控制问题.

给定容许控制 $\mathbf{u}_k = \mathbf{K}\mathbf{X}_k$, 令 $\mathbf{T}_k = \mathbb{E}(\mathbf{X}_k \mathbf{X}_k^T)$, 系统(5)转化为

$$\mathbf{T}_{k+1} = \mathbb{E}(\mathbf{X}_{k+1} \mathbf{X}_{k+1}^T) = \mathbb{E}((\mathbf{T} + \mathbf{B}_0 \mathbf{K}) \mathbf{X}_k \mathbf{X}_k^T (\mathbf{T} + \mathbf{B}_0 \mathbf{K})^T), \quad (12)$$

代入增广系统式(5),化简得:

$$T_{k+1} = E(X_{k+1}X_{k+1}^T) = (A_1 + B_1K)T_k(A_1 + B_1K)^T + (C_1 + D_1K)T_k(C_1 + D_1K)^T, \quad (13)$$

式中: $T_k \in \mathbf{R}^{(n+q) \times (n+q)}$ 是确定性系统的状态. 令 $T_0 = E(X_0X_0^T)$ 为初始状态.

相应地,性能指标函数式(7) 重写为

$$J(T_k, K) = \text{tr} \left\{ \sum_{i=k}^{\infty} \gamma^{i-k} (Q_1 + K^T R K) T_i \right\}. \quad (14)$$

注 2 经系统转变后,随机系统转变为确定性系统.由式(13)可知,系统转变后随机参数 w_k 对系统的影响包含在系统状态中,这为下文 Q 学习算法的推导和应用做出了准备.

3 Q 学习算法的推导

为了求解随机系统最优追踪问题而无需系统参数信息,本节推导出在线求解随机系统追踪问题的 Q 学习算法.

引理 2^[14] 给定容许控制 $u_k = KX_k$,由引理 1 可知系统的最优值函数为

$$V^*(X_k) = E(X_k^T P^* X_k) = \text{tr}(P^* T_k), \quad (15)$$

最优控制增益矩阵为

$$u_k^* = K^* X_k = -(R + \gamma B_1^T P^* B_1 + \gamma D_1^T P^* D_1)^{-1} \gamma (B_1^T P^* A_1 + D_1^T P^* C_1) X_k, \quad (16)$$

式中:矩阵 P^* 是式(17) 增广随机代数方程的解.

$$P = Q_1 + \gamma(A_1^T P A_1 + C_1^T P C_1) - \gamma(A_1^T P B_1 + C_1^T P D_1) \times (R + \gamma B_1^T P B_1 + \gamma D_1^T P D_1)^{-1} \times \gamma(B_1^T P A_1 + D_1^T P C_1), \quad (17)$$

式中: $R + \gamma B_1^T P B_1 + \gamma D_1^T P D_1 > 0$.

注 3 由引理 2 可知,求解 SLQ 最优追踪控制的有效方法之一就是求解增广随机代数方程.然而,求解式(17) 需要完全已知系统的全部动力学信息,因此,当系统动力学参数未知时,式(17) 就无法求解.

为了求解最优追踪控制无需系统的动力学信息,下文给出 Q 函数的定义和对应的 H 矩阵,基于式(7) 得:

$$J(X_k, u) = E(X_k^T Q_1 X_k + u_k^T R u_k) + \gamma \sum_{i=k+1}^{\infty} \gamma^{i-(k+1)} [X_i^T Q_1 X_i + u_i^T R u_i], \quad (18)$$

由式(18)和贝尔曼最优性原理得:

$$V^*(X_k) = \min_{u_k} \{ E(X_k^T Q_1 X_k + u_k^T R u_k) + \gamma V^*(X_{k+1}) \}, \quad (19)$$

定义 Q 函数为

$$Q(X_k, u_k) = E(X_k^T Q_1 X_k + u_k^T R u_k) + \gamma E(X_{k+1}^T P X_{k+1}), \quad (20)$$

代入增广系统(5),式(20)变为

$$Q(X_k, u_k) = E(X_k^T Q_1 X_k + u_k^T R u_k) + \gamma E((TX_k + B_0 u_k)^T P (TX_k + B_0 u_k)), \quad (21)$$

化简式(21)可得:

$$Q(X_k, u_k) = E \left\{ \begin{bmatrix} X_k \\ u_k \end{bmatrix}^T \begin{bmatrix} H_{XX} & H_{uX} \\ H_{Xu} & H_{uu} \end{bmatrix} \begin{bmatrix} X_k \\ u_k \end{bmatrix} \right\}. \quad (22)$$

进一步化简式(22):

$$Q(X_k, u_k) = E \left\{ \begin{bmatrix} X_k \\ u_k \end{bmatrix}^T \begin{bmatrix} H_{XX} & H_{uX} \\ H_{Xu} & H_{uu} \end{bmatrix} \begin{bmatrix} X_k \\ u_k \end{bmatrix} \right\} = E \left\{ \begin{bmatrix} X_k \\ u_k \end{bmatrix}^T H \begin{bmatrix} X_k \\ u_k \end{bmatrix} \right\}. \quad (23)$$

式中: $H = H^T \in \mathbf{R}^{(n+q+m) \times (n+q+m)}$,

$$\begin{aligned} H_{XX} &= \gamma A_1^T P A_1 + \gamma C_1^T P C_1 + Q_1 \in \mathbf{R}^{(n+q) \times (n+q)}, \\ H_{uX} &= \gamma B_1^T P A_1 + \gamma D_1^T P C_1 \in \mathbf{R}^{m \times (n+q)}, \\ H_{Xu} &= \gamma A_1^T P B_1 + \gamma C_1^T P D_1 \in \mathbf{R}^{(n+q) \times m}, \\ H_{uu} &= R + \gamma B_1^T P B_1 + \gamma D_1^T P D_1 \in \mathbf{R}^{m \times m}. \end{aligned}$$

由式(16)可知: $u_k^* = K^* X_k = -(R + \gamma B_1^T P^* B_1 + \gamma D_1^T P^* D_1)^{-1} \gamma (B_1^T P^* A_1 + D_1^T P^* C_1) X_k$. 当控制策略 u_k^* 为最优的控制策略时, $Q^*(X_k, u_k)$ 和 $V^*(X_k)$ 是等价的,由最优化一阶必要条件,通过微分 $\frac{\partial}{\partial u_k} Q(X_k, u_k) = 0$,可以得到:

$$K^* = -H_{uu}^{-1} H_{uX}. \quad (24)$$

并且 Q 函数可重写为

$$Q(X_k, u_k) = E \{ [X_k^T \ u_k^T] H [X_k^T \ u_k^T]^T \}. \quad (25)$$

由引理 1 和式(25)可知 P 和 H 的联系为

$$P = [I \ K^T] H [I \ K^T]^T. \quad (26)$$

通过式(24)可知,求解随机线性二次最优追踪控制增益矩阵摆脱了系统模型参数的限制,系统控制的最优解仅依赖于 H 矩阵, H 矩阵包含了系统的动态特性信息,如果 H 矩阵已知,就能够实现在无需系统参数的情况下在线求解系统的最优控制增益矩阵 K^* .

文献[16] 使用 Q 学习算法解决确定性离散系统的 H_∞ 最优控制问题,文献[13] 使用 Q 学习算法解决随机离散系统的最优调节控制问题,并且均证明了 Q 学习算法的收敛性.本文的 Q 学习算法推导过程基于文献[13] 和文献[16]. 基于值迭代方法^[13],由式(25) 可得:

$$\begin{aligned} E \{ [X_k^T \ u_k^T] H_{i+1} [X_k^T \ u_k^T]^T \} &= \\ E \{ X_k^T [I \ K_i^T] H_{i+1} [I \ K_i^T]^T X_k \} &= \\ \text{tr} \{ [I \ K_i^T] H_{i+1} [I \ K_i^T]^T T_k \}. & \quad (27) \end{aligned}$$

同时由式(21)、(25)以及式(20)的右半部分化简为

$$\begin{aligned} E\{(\mathbf{X}_k^T \mathbf{Q}_1 \mathbf{X}_k + \mathbf{u}_k^T \mathbf{R} \mathbf{u}_k) + \gamma [\mathbf{X}_{k+1}^T \quad \mathbf{u}_{k+1}^T] \mathbf{H}_i [\mathbf{X}_{k+1}^T \quad \mathbf{u}_{k+1}^T]^T\} = \\ E\{[\mathbf{X}_k^T \quad \mathbf{u}_k^T] \begin{bmatrix} \mathbf{Q}_1 & \mathbf{O} \\ \mathbf{O} & \mathbf{R} \end{bmatrix} [\mathbf{X}_k^T \quad \mathbf{u}_k^T]^T + \\ \gamma [\mathbf{X}_{k+1}^T \quad \mathbf{u}_{k+1}^T] \mathbf{H}_i [\mathbf{X}_{k+1}^T \quad \mathbf{u}_{k+1}^T]^T\} = \\ E\{\mathbf{X}_k^T [\mathbf{I} \quad \mathbf{K}_i^T] \begin{bmatrix} \mathbf{Q}_1 & \mathbf{O} \\ \mathbf{O} & \mathbf{R} \end{bmatrix} [\mathbf{I} \quad \mathbf{K}_i^T]^T \mathbf{X}_k + \\ \gamma \mathbf{X}_{k+1}^T [\mathbf{I} \quad \mathbf{K}_i^T] \mathbf{H}_i [\mathbf{I} \quad \mathbf{K}_i^T]^T \mathbf{X}_{k+1}\} = \\ \text{tr}\{[\mathbf{I} \quad \mathbf{K}_i^T] \begin{bmatrix} \mathbf{Q}_1 & \mathbf{O} \\ \mathbf{O} & \mathbf{R} \end{bmatrix} [\mathbf{I} \quad \mathbf{K}_i^T]^T \mathbf{T}_k + \\ \gamma [\mathbf{I} \quad \mathbf{K}_i^T] \mathbf{H}_i [\mathbf{I} \quad \mathbf{K}_i^T]^T \mathbf{T}_{k+1}\}. \end{aligned} \quad (28)$$

由式(27)、(28),可以得到:

$$\begin{aligned} \text{tr}\{[\mathbf{I} \quad \mathbf{K}_i^T] \mathbf{H}_{i+1} [\mathbf{I} \quad \mathbf{K}_i^T]^T \mathbf{T}_k\} = \\ \text{tr}\{[\mathbf{I} \quad \mathbf{K}_i^T] \begin{bmatrix} \mathbf{Q}_1 & \mathbf{O} \\ \mathbf{O} & \mathbf{R} \end{bmatrix} [\mathbf{I} \quad \mathbf{K}_i^T]^T \mathbf{T}_k + \\ \gamma [\mathbf{I} \quad \mathbf{K}_i^T] \mathbf{H}_i [\mathbf{I} \quad \mathbf{K}_i^T]^T \mathbf{T}_{k+1}\}. \end{aligned} \quad (29)$$

由式(24)得控制迭代式:

$$\mathbf{K}_i = -(\mathbf{H}_{uu}^{-1} \mathbf{H}_{uX})^T. \quad (30)$$

式(29)、(30)构成了 Q 学习算法,可知系统参数信息 $\mathbf{A} \in \mathbf{R}^{n \times n}$, $\mathbf{B} \in \mathbf{R}^{n \times m}$, $\mathbf{C} \in \mathbf{R}^{n \times n}$, $\mathbf{D} \in \mathbf{R}^{n \times m}$ 均包含在系统状态 \mathbf{T}_k 和 \mathbf{T}_{k+1} 之中,因此系统矩阵信息 $\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}$ 矩阵可以是未知的.接下来说明 Q 学习算法的等价性和收敛性.

定理 1 Q 学习算法式(29)——(30) 等价于

$$\begin{aligned} \mathbf{P}_{i+1} = \mathbf{Q}_1 + \gamma (\mathbf{A}_1^T \mathbf{P}_i \mathbf{A}_1 + \mathbf{C}_1^T \mathbf{P}_i \mathbf{C}_1) - \\ \gamma (\mathbf{A}_1^T \mathbf{P}_i \mathbf{B}_1 + \mathbf{C}_1^T \mathbf{P}_i \mathbf{D}_1) \times \\ (\mathbf{R} + \gamma \mathbf{B}_1^T \mathbf{P}_i \mathbf{B}_1 + \gamma \mathbf{D}_1^T \mathbf{P}_i \mathbf{D}_1)^{-1} \times \\ \gamma (\mathbf{B}_1^T \mathbf{P}_i \mathbf{A}_1 + \mathbf{D}_1^T \mathbf{P}_i \mathbf{C}_1). \end{aligned} \quad (31)$$

证明 由于

$$\begin{aligned} E\{[\mathbf{X}_{k+1}^T \quad \mathbf{u}_i^T(\mathbf{X}_{k+1})] \mathbf{H}_i [\mathbf{X}_{k+1}^T \quad \mathbf{u}_i^T(\mathbf{X}_{k+1})]^T\} = \\ E\{\mathbf{X}_{k+1}^T [\mathbf{I} \quad \mathbf{K}_i^T] \mathbf{H}_i [\mathbf{I} \quad \mathbf{K}_i^T]^T \mathbf{X}_{k+1}\} = \\ E\{(\mathbf{T} \mathbf{X}_k + \mathbf{B}_0 \mathbf{u}_k)^T [\mathbf{I} \quad \mathbf{K}_i^T] \mathbf{H}_i [\mathbf{I} \quad \mathbf{K}_i^T]^T (\mathbf{T} \mathbf{X}_k + \mathbf{B}_0 \mathbf{u}_k)\}, \end{aligned} \quad (32)$$

代入式(26),化简式(27)和式(28)的左半部分后,得:

$$\begin{aligned} \mathbf{H}_{i+1} = \begin{bmatrix} \mathbf{Q}_1 & \mathbf{O} \\ \mathbf{O} & \mathbf{R} \end{bmatrix} + \begin{bmatrix} \gamma \mathbf{A}_1^T \mathbf{P}_i \mathbf{A}_1 & \gamma \mathbf{A}_1^T \mathbf{P}_i \mathbf{B}_1 \\ \gamma \mathbf{B}_1^T \mathbf{P}_i \mathbf{A}_1 & \gamma \mathbf{B}_1^T \mathbf{P}_i \mathbf{B}_1 \end{bmatrix} + \\ \begin{bmatrix} \gamma \mathbf{C}_1^T \mathbf{P}_i \mathbf{C}_1 & \gamma \mathbf{C}_1^T \mathbf{P}_i \mathbf{D}_1 \\ \gamma \mathbf{D}_1^T \mathbf{P}_i \mathbf{C}_1 & \gamma \mathbf{D}_1^T \mathbf{P}_i \mathbf{D}_1 \end{bmatrix}, \end{aligned} \quad (33)$$

又因为:

$$\mathbf{P}_{i+1} = [\mathbf{I} \quad \mathbf{K}_{i+1}^T] \mathbf{H}_{i+1} [\mathbf{I} \quad \mathbf{K}_{i+1}^T]^T. \quad (34)$$

结合式(30)、(33)和(34)可得:

$$\begin{aligned} \mathbf{P}_{i+1} = \mathbf{Q}_1 + \gamma (\mathbf{A}_1^T \mathbf{P}_i \mathbf{A}_1 + \mathbf{C}_1^T \mathbf{P}_i \mathbf{C}_1) - \\ \gamma (\mathbf{A}_1^T \mathbf{P}_i \mathbf{B}_1 + \mathbf{C}_1^T \mathbf{P}_i \mathbf{D}_1) \times \\ (\mathbf{R} + \gamma \mathbf{B}_1^T \mathbf{P}_i \mathbf{B}_1 + \gamma \mathbf{D}_1^T \mathbf{P}_i \mathbf{D}_1)^{-1} \times \\ \gamma (\mathbf{B}_1^T \mathbf{P}_i \mathbf{A}_1 + \mathbf{D}_1^T \mathbf{P}_i \mathbf{C}_1). \end{aligned} \quad (35)$$

式中: $\mathbf{R} + \gamma \mathbf{B}_1^T \mathbf{P}_i \mathbf{B}_1 + \gamma \mathbf{D}_1^T \mathbf{P}_i \mathbf{D}_1 > 0$.

Q 学习算法基于值迭代算法,值迭代算法的收敛性在文献[14]已经有详细的证明,并且 Q 学习算法的迭代和增广系统的随机代数方程自身的迭代是等价的.

4 Q 学习算法的实现

本节给出 Q 学习算法的具体实现步骤.

式(29)可进一步化简为

$$\begin{aligned} \text{tr}\{\mathbf{H}_{i+1} [\mathbf{I} \quad \mathbf{K}_i^T]^T \mathbf{T}_k [\mathbf{I} \quad \mathbf{K}_i^T]\} = \\ \text{tr}\left\{\begin{bmatrix} \mathbf{Q}_1 & \mathbf{O} \\ \mathbf{O} & \mathbf{R} \end{bmatrix} [\mathbf{I} \quad \mathbf{K}_i^T]^T \mathbf{T}_k [\mathbf{I} \quad \mathbf{K}_i^T] + \right. \\ \left. \gamma \mathbf{H}_i [\mathbf{I} \quad \mathbf{K}_i^T]^T \mathbf{T}_{k+1} [\mathbf{I} \quad \mathbf{K}_i^T]\right\}. \end{aligned} \quad (36)$$

为了方便计算式(36),另记

$$\mathbf{Z}_k = [\mathbf{I} \quad \mathbf{K}_i^T]^T \mathbf{T}_k [\mathbf{I} \quad \mathbf{K}_i^T], \mathbf{G} = \begin{bmatrix} \mathbf{Q}_1 & \mathbf{O} \\ \mathbf{O} & \mathbf{R} \end{bmatrix},$$

使用最小二乘法求解式(36)之前,首先引入三个算子:

$$\begin{aligned} \mathbf{H} \in \mathbf{R}^{(n+q+m) \times (n+q+m)} \rightarrow \hat{\mathbf{H}} \in \mathbf{R}^{\frac{1}{2}(n+q+m) \times (n+q+m+1)}, \\ \mathbf{Z}_k \in \mathbf{R}^{(n+q+m) \times (n+q+m)} \rightarrow \hat{\mathbf{Z}}_k \in \mathbf{R}^{\frac{1}{2}(n+q+m) \times (n+q+m+1)}, \\ \mathbf{G} \in \mathbf{R}^{(n+q+m) \times (n+q+m)} \rightarrow \hat{\mathbf{G}} \in \mathbf{R}^{\frac{1}{2}(n+q+m) \times (n+q+m+1)}. \end{aligned}$$

式中:

$$\hat{\mathbf{H}} = [h_{11}, 2h_{12}, 2h_{13}, \dots, 2h_{1n}, h_{22}, 2h_{23}, \dots, 2h_{(n-1)n}, h_{nn}]^T,$$

$$\hat{\mathbf{Z}}_k = [z_{11}, z_{12}, z_{13}, \dots, z_{1n}, z_{22}, z_{23}, \dots, z_{(n-1)n}, z_{nn}]^T,$$

$$\hat{\mathbf{G}} = [g_{11}, 2g_{12}, 2g_{13}, \dots, 2g_{1n}, g_{22}, 2g_{23}, \dots, 2g_{(n-1)n}, g_{nn}]^T.$$

因此,式(36)可化简为

$$\begin{aligned} \text{tr}\{\mathbf{H}_{i+1} [\mathbf{I} \quad \mathbf{K}_i^T]^T \mathbf{T}_k [\mathbf{I} \quad \mathbf{K}_i^T]\} = \hat{\mathbf{Z}}_k^T \hat{\mathbf{H}}_{i+1}, \\ \text{tr}\left\{\begin{bmatrix} \mathbf{Q}_1 & \mathbf{O} \\ \mathbf{O} & \mathbf{R} \end{bmatrix} [\mathbf{I} \quad \mathbf{K}_i^T]^T \mathbf{T}_k [\mathbf{I} \quad \mathbf{K}_i^T] + \right. \\ \left. \gamma \mathbf{H}_i [\mathbf{I} \quad \mathbf{K}_i^T]^T \mathbf{T}_{k+1} [\mathbf{I} \quad \mathbf{K}_i^T]\right\} = \hat{\mathbf{G}}^T \hat{\mathbf{Z}}_k + \gamma \hat{\mathbf{Z}}_{k+1}^T \hat{\mathbf{H}}_i. \end{aligned}$$

由最小二乘法可知,至少需要收集 N 组数据,即

$N \geq \frac{1}{2}(n+q+m)(n+q+m+1)$. 记 $\hat{\mathbf{Z}}_{k+N}^T$ 为收集的

N 个状态信息构成的矩阵, $\hat{\mathbf{Z}}_{k+N}^T =$

$[\hat{\mathbf{Z}}_1, \hat{\mathbf{Z}}_2, \hat{\mathbf{Z}}_3, \dots, \hat{\mathbf{Z}}_N]^T$, 如果矩阵 $\hat{\mathbf{Z}}_{k+N}^T$ 是列满秩的,则

式(36)可以由下式直接求解:

$$\hat{H}_{i+1} = (\hat{Z}_{k+\lambda} \hat{Z}_{k+\lambda}^T)^{-1} \hat{Z}_{k+\lambda} (\hat{G}^T \hat{Z}_{k+\lambda} + \gamma \hat{Z}_{k+1+\lambda}^T \hat{H}_i), \quad (37)$$

Q 学习算法的流程如图 1 所示.

注 4 因为 H 是维度为 $(n+q+m) \times (n+q+m)$ 的对称矩阵, 因此可知 H 有 $\frac{1}{2}(n+q+m) \times (n+q+m+1)$ 个独立元素, 所以求解式 (37):

1) 至少需要与 H 独立元素相同数量的方程, 也就是必须保证 $\text{rank}(\hat{Z}_{k+\lambda}^T) = \frac{1}{2}(n+q+m) \times (n+q+m+1)$, 收集足够多的数据以保证式 (37) 可解;

2) 同时在系统 (13) 需要引入探测噪声, 以保证能够充分探测系统参数信息.

5 系统仿真

本节给出一个仿真实例来说明 Q 学习算法的有效性.

假设乘性噪声随机线性系统动态模型如下:

$$\begin{cases} \mathbf{x}_{k+1} = \begin{bmatrix} 0.2 & -0.8 \\ 0.5 & -0.7 \end{bmatrix} \mathbf{x}_k + \begin{bmatrix} 0.03 \\ -0.5 \end{bmatrix} \mathbf{u}_k + \\ \left(\begin{bmatrix} -0.04 & 0.4 \\ -0.3 & 0.13 \end{bmatrix} \mathbf{x}_k + \begin{bmatrix} 0.05 \\ -0.3 \end{bmatrix} \mathbf{u}_k \right) \mathbf{w}_k, \\ \mathbf{y}_k = \begin{bmatrix} 3 & 3 \end{bmatrix} \mathbf{x}_k. \end{cases} \quad (38)$$

假设参考轨迹的动态模型如下:

$$\mathbf{r}_{k+1} = -\mathbf{r}_k. \quad (39)$$

设置 $Q = 10, R = 1, \gamma = 0.9$, 增广系统的初始状态 $\mathbf{X}_0 = \begin{bmatrix} 10 \\ 10 \\ 1 \end{bmatrix}$, 参考轨迹的初始状态 $\mathbf{r}_0 = 1$. 更具一般

性, 初始化 $H = I$, 此时 $K = [0 \ 0 \ 0]$, 在每一次迭代过程中, 加入探测噪声并收集 20 组数据去求解更新控制增益矩阵 K . 经过 Q 学习算法迭代后, 控制增益矩阵曲线如图 2 所示.

为了验证 Q 学习算法迭代求解的 K 矩阵的准确性, 需要把迭代计算的 K 矩阵同 SAE 方程求解出的最优控制增益矩阵 K^* 比较. 首先通过式 (16)、(17) 的 SAE 方程求得系统的最优控制矩阵 K^* . K 矩阵收敛到最优控制矩阵 K^* 的过程如图 3 所示. 由图 3 可以看出, 随着迭代次数的增加, 控制矩阵 K 收敛到最优控制矩阵 K^* , 并且 P 收敛到代数方程的解 P^* . P 收敛到代数方程的解 P^* 的过程如图 4 所示, 验证了 Q 学习算法的有效性.

追踪问题的目标是追踪到参考信号轨迹, 结果如图 5 所示, 系统的期望输出 $E(\mathbf{y})$ 追踪到了参考轨迹 \mathbf{r} .

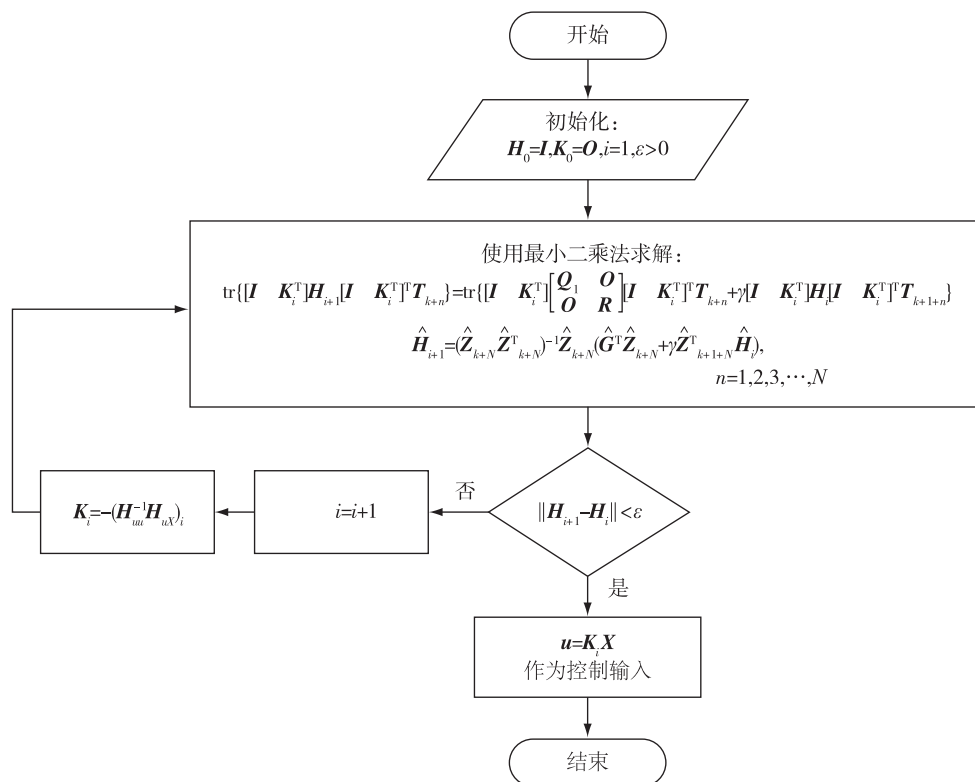


图 1 Q 学习算法流程

Fig. 1 Flowchart of Q-learning

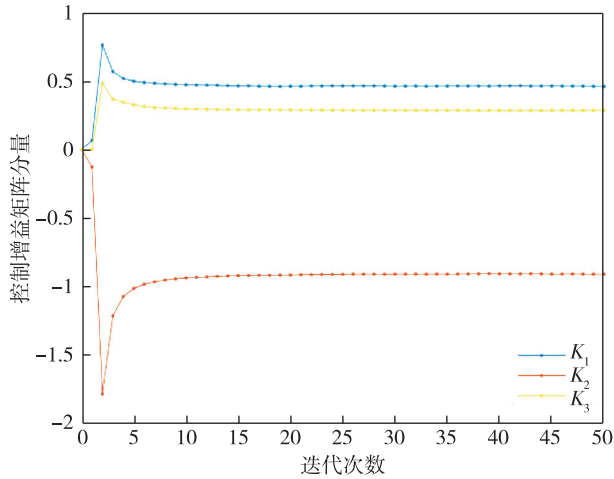


图2 控制增益矩阵 K 曲线
Fig. 2 Curves of control gain matrix K

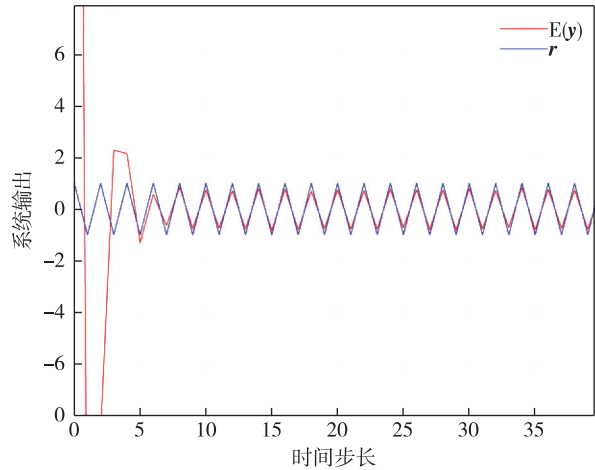


图5 系统的期望输出曲线
Fig. 5 Curves of system output $E(y)$

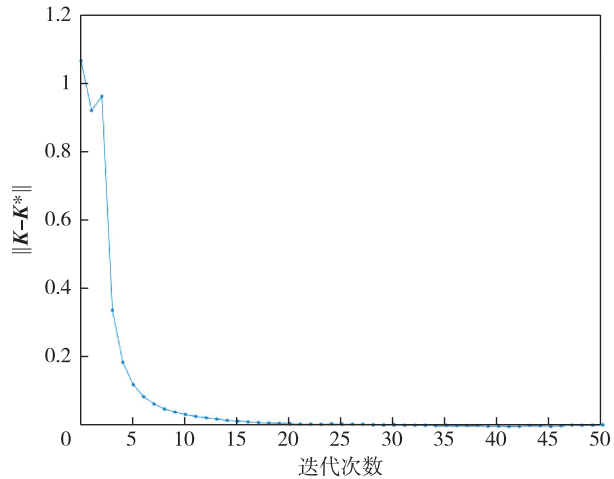


图3 控制矩阵 K 的收敛
Fig. 3 Convergence of control gain matrix K

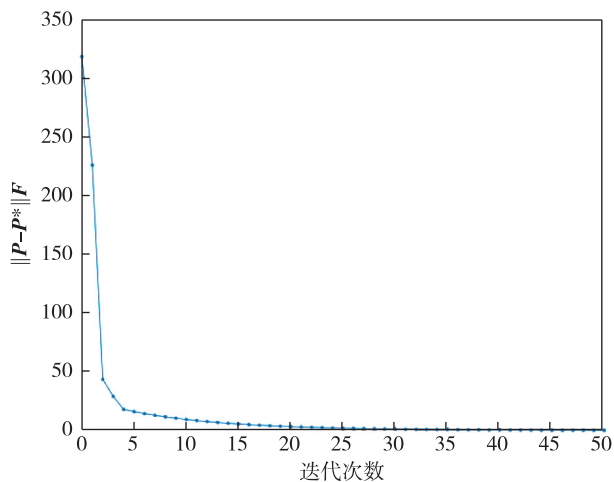


图4 矩阵 P 的收敛
Fig. 4 Convergence of matrix P

6 结论

通常来说,求解随机最优追踪控制问题需要完全的系统参数信息,本文针对离散时间系统的随机线性二次最优追踪控制问题,推导了 Q 学习算法,给出算法的具体实现步骤,使用 Q 学习算法在线解决追踪控制问题而无需系统模型参数,最后给出仿真结果表明系统输出可以有效地追踪参考轨迹.

参考文献

References

- [1] Byers R. Solving the algebraic Riccati equation with the matrix sign function [J]. Linear Algebra and Its Applications, 1987, 85: 267-279
- [2] Kleinman, D. On an iterative technique for Riccati equation computations [J]. IEEE Transactions on Automatic Control, 1968, 13(1): 114-115
- [3] Wang F Y, Zhang H G, Liu D R. Adaptive dynamic programming: an introduction [J]. IEEE Computational Intelligence Magazine, 2009, 4(2): 39-47
- [4] Vrabie D, Pastravanu O, Abu-Khalaf M, et al. Adaptive optimal control for continuous-time linear systems based on policy iteration [J]. Automatica, 2009, 45(2): 477-484
- [5] Jiang Y, Jiang Z P. Computational adaptive optimal control for continuous-time linear systems with completely unknown dynamics [J]. Automatica, 2012, 48(10): 2699-2704
- [6] Kiumarsi B, Lewis F L, Modares H, et al. Reinforcement Q-learning for optimal tracking control of linear discrete-time systems with unknown dynamics [J]. Automatica, 2014, 50(4): 1167-1175
- [7] Zhao X Y, Deng F Q. Divided state feedback control of stochastic systems [J]. IEEE Transactions on Automatic Control, 2015, 60(7): 1870-1885
- [8] Zhao X Y, Deng F Q. A new type of stability theorem for

- stochastic systems with application to stochastic stabilization[J].IEEE Transactions on Automatic Control,2016,61(1):240-245
- [9] 黄玉林,张维海.约束随机线性二次最优控制的研究[J].自动化学报,2006,32(2):246-254
HUANG Yulin, ZHANG Weihai. Study on stochastic linear quadratic optimal control with constraint[J].Acta Automatica Sinica,2006,32(2):246-254
- [10] Liu X K, Li Y, Zhang W H. Stochastic linear quadratic optimal control with constraint for discrete-time systems[J].Applied Mathematics and Computation,2014,228:264-270
- [11] 王涛,张化光.基于策略迭代的连续时间系统的随机线性二次最优控制[J].控制与决策,2015,30(9):1674-1678
WANG Tao, ZHANG Huaguang. Stochastic linear quadratic optimal control for continuous-time systems based on policy iteration[J].Control and Decision,2015,30(9):1674-1678
- [12] Wang T, Zhang H G, Luo Y H. Infinite-time stochastic linear quadratic optimal control for unknown discrete-time systems using adaptive dynamic programming approach[J].Neurocomputing,2016,171:379-386
- [13] Wang T, Zhang H G, Luo Y H. Stochastic linear quadratic optimal control for model-free discrete-time systems based on Q-learning algorithm [J]. Neurocomputing,2018,312:1-8
- [14] Chen X, Wang F. Neural-network-based stochastic linear quadratic optimal tracking control scheme for unknown discrete-time systems using adaptive dynamic programming[J].Control Theory and Technology,2021,19(3):315-327
- [15] Xiao G Y, Zhang H G, Luo Y H, et al. Data-driven optimal tracking control for a class of affine non-linear continuous-time systems with completely unknown dynamics[J].IET Control Theory & Applications,2016,10(6):700-710
- [16] Al-Tamimi A, Lewis F L, Abu-Khalaf M. Model-free Q-learning designs for linear discrete-time zero-sum games with application to H-infinity control [J]. Automatica,2007,43(3):473-481

Stochastic linear quadratic optimal tracking control for stochastic discrete time systems based on Q-learning

ZHANG Zhengyi¹ ZHAO Xueyan¹

¹ School of Automation Science and Engineering, South China University of Technology, Guangzhou 510640

Abstract For stochastic linear discrete time systems, a Q-learning algorithm is proposed in this paper to solve the stochastic linear quadratic optimal tracking control problem in the infinite time domain. First, it is assumed that the reference signal required for tracking is generated by the command generator, and an augmented system consisting of the original stochastic system and the reference trajectory system is established, then the optimal tracking problem is transformed into an optimal regulation problem. Second, in order to solve the optimal tracking problem online, the stochastic system is transformed into a deterministic one, the Q function of stochastic linear quadratic optimal tracking control is defined according to the augmented system, and the augmented stochastic algebraic equation is solved online without knowing the parameters of the system model. Third, the equivalence between the Q-learning algorithm and the augmented stochastic algebraic equation is proved, and the implementation steps of the Q-learning algorithm are given. Finally, a simulation example is given to illustrate the effectiveness of the proposed Q-learning algorithm.

Key words stochastic systems; Q-learning algorithm; optimal tracking control; stochastic algebraic equation