# Adaptive dynamic programming-based algorithm for infinite-horizon linear quadratic stochastic optimal control problems

Heng Zhang

School of Control Science and Engineering, Shandong University, Jinan, 250061, China.
E-mail: zhangh2828@163.com

**Abstract:** This paper investigates an infinite-horizon linear quadratic stochastic (LQS) optimal control problem for a class of continuous-time stochastic systems. By employing the technique of adaptive dynamic programming (ADP), we propose a novel model-free policy iteration (PI) algorithm. Without needing all information of the system coefficient matrices, the proposed PI algorithm iterates by using the data of the input and system state collected on a fixed time interval. Finally, a numerical example is presented to demonstrate the feasibility of the obtained algorithm.

**Key Words:** Linear quadratic stochastic optimal control, Policy iteration, Adaptive dynamic programming

## 1 INTRODUCTION

The linear quadratic stochastic (LQS) optimal control problem, initiated by Wonham [15] has been broadly applied in a lot of fields such as engineering. As is known to all, the continuous-time LQS problem in infinite horizon is closely related to the stochastic algebraic Riccati equation (SARE), which is difficult to solve due to its nonlinear structure. With the in-depth study of the LQS optimal control problem, researchers developed some approximation methods to obtain the solution of the SARE. For instance, Ni and Fang [18] proposed a PI algorithm to solve the SARE iteratively. With the help of the positive operators, a Newton's method was proposed by Damm and Hinrichsen [12] to solve the SARE. However, the above methods need all knowledge of the system, i.e., all parameters of the system have to be known beforehand. In fact, the system matrices are difficult to obtain directly in applications such as engineering and finance. The methods mentioned above will become invalid if the system coefficient matrices are unknown. Thus, it is of great importance to propose a model-free strategy to solve LQS optimal control problems, without using the information of system matrices.

For the past decade, adaptive dynamic programming (ADP) (Werbos [7]) and reinforcement learning (RL) (Sutton and Barto [9]) theories have been broadly used to solve optimal control problems with partially model-free or model-free system dynamics. About the development of deterministic system case, see, e.g., Shi and Wang [20], Pang et al. [2], Kiumarsi et al. [1], Vamvoudakis et al. [4], Bian and Jiang [11], Palanisamy et al. [6], Vrabie et al. [3], Wei et al. [8], Jiang and Jiang [16], Mukherjee et al. [10] and the references therein.

Regarding to stochastic optimal control problems, Ge et al. [19] proposed a model-free methodology to get the optimal policy for a kind of mean-field discrete-time stochastic systems by the method of Q-learning. By the technique of ADP, Wang et al. [13] solved a class of discrete-time LQS optimal control problems. Wang et al. [14] developed a model-free Q-learning algorithm to get the optimal control for discrete-time LQS problems. By applying RL techniques, Jiang and Jiang [17] developed an ADP strategy to solve continuous-time optimal control problems where the systems subject to control-dependent noise.

However, to the author's best knowledge, there is no model-free results for continuous-time LQS optimal control problems where drift and diffusion terms contain both control and state variables. The main contribution of this paper is that we propose a model-free algorithm to solve this class of continuous-time LQS problems.

To be specific, we propose a novel data-driven model-free PI algorithm to get the maximal solution to the SARE by using the data of the input and state collected on some time interval. The convergence proof of our model-free strategy is also been provided.

The rest of the paper is organized as follows. In Section 2, the formulation of our problem and some preliminaries are presented. Section 3 develops our data-driven model-free PI algorithm. In Section 4, we provide a simulation example

to illustrate the applicability of the proposed algorithm. In Section 5, some conclusions are presented.

**Notation.** We denote the collections of non-negative integers, positive integers and real numbers by $\mathbb{Z}$, $\mathbb{Z}^+$ and $\mathbb{R}$. $\mathbb{R}^{n \times m}$ represents the collection of all $n \times m$ real matrices. $\mathbb{R}^n$ is the $n$-dimensional Euclidean space and $|\cdot|$ denotes its Euclidean norm for vector or matrix of proper size. Zero matrix (or vector) with appropriate dimension is denoted by $O$. We use $diag\{v\}$ to denote a square diagonal matrix whose main diagonal is the elements of vector $v$. The sets of all symmetric matrices, positive definite matrices and semipositive definite matrices in $\mathbb{R}^{n \times n}$ are represented by $\mathbf{S}^n$, $\mathbf{S}^n_{++}$ and $\mathbf{S}^n_+$, respectively. $w(\cdot)$ is a one-dimensional standard Brownian motion defined on a filtered probability space $(\Omega, \mathcal{F}, \{\mathcal{F}_t\}_{t \geqslant 0}, \mathbb{P})$ that satisfies usual conditions. Moreover, we use $\otimes$ to denote the Kronecker product and for any matrix $B \in \mathbb{R}^{m \times n}$, $vec(B)$ denotes a vectorization map from the matrix $B$ into a column vector of proper size, which stacks the columns of $B$ on top of one another, that is, $vec(B) = [b_1^T, b_2^T, \cdots, b_n^T]^T$, where $b_j \in \mathbb{R}^n$, $j = 1, 2, 3, \cdots, n$, are the columns of $B$. For any $\xi \in \mathbb{R}^n$ and $F \in \mathbf{S}^n$, we define two operators as follows:

$$vecs : \xi \in \mathbb{R}^l \to vecs(\xi) \in \mathbb{R}^{\frac{n(n+1)}{2}},$$
$$\text{and } vech : F \in \mathbf{S}^l \to vech(F) \in \mathbb{R}^{\frac{n(n+1)}{2}},$$

where

$$vecs(\xi) = [\xi_1^2, \xi_1\xi_2, \cdots, \xi_1\xi_n, x_2^2, x_2x_3, \cdots, \xi_{n-1}\xi_n, \xi_n^2]^T,$$
$$vech(F) = [f_{11}, 2f_{12}, \cdots, 2f_{1n}, f_{22}, 2f_{23}, \cdots, 2f_{n-1,n}, f_{nn}]^T,$$

and $\xi_j$, $j = 1, 2, \cdots, n$, is the $j$th element of $\xi$ and $f_{ji}$, $j, i = 1, 2, \cdots, n$, is the $(j,i)$th element of matrix $F$. For simplity, we denote $vecs(\xi)$ by $\overline{\xi}$ in this paper.

## 2 PROBLEM FORMULATION

This section presents the formulation of our LQS optimal control problems.

Consider a continuous-time time-invariant stochastic linear system as follows

$$\begin{cases} dx(s) = [Ax(s) + Bu(s)]ds \\ \qquad\quad + [Cx(s) + Du(s)]dw(s), \\ x(0) = x_0, \end{cases} \tag{1}$$

where $x_0 \in \mathbb{R}^n$ is the initial state. The cost functional is defined as

$$J(u(\cdot)) = \mathbb{E} \int_0^\infty [x(s)^T Q x(s) + u(s)^T R u(s)]ds, \tag{2}$$

where $R > 0$, $Q \geq 0$ and $[A, C|Q]$ is exactly detectable. Now we give the definition of mean-square stabilizability.

**Definition 1.** System (1) is called mean-square stabilizable for any initial state $x_0$, if there exists a matrix $K \in \mathbb{R}^{m \times n}$ such that the solution of

$$\begin{cases} dx(s) = (A + BK)x(s)ds \\ \qquad\quad + (C + DK)x(s)dw(s), \\ x(0) = x_0 \end{cases} \tag{3}$$

satisfies $\lim_{s \to \infty} \mathbb{E}[x(s)^T x(s)] = 0$. In this case, the feedback control $u(\cdot) = Kx(\cdot)$ is called stabilizing and the constant matrix $K$ is called a stabilizer of system (1).

**Assumption 1.** System (1) is mean-square stabilizable.

Under Assumption 1, we define the sets of admissible control as

$$\mathcal{U}_{ad} = \{u(\cdot) \in L^2_\mathcal{F}(\mathbb{R}^m) | u(\cdot) \text{ is stabilizing}\}. \tag{4}$$

Our continuous-time LQS optimal control problems are given as follows:

**Problem (LQS).** For any initial state $x_0 \in \mathbb{R}^n$, we want to find an optimal control $u^*(\cdot) \in \mathcal{U}_{ad}$ such that

$$J(u^*(\cdot)) = \inf_{u(\cdot) \in \mathcal{U}_{ad}} J(u(\cdot)). \tag{5}$$

Ni and Fang [18] shows that the optimal control of Problem (LQS) can be obtained by solving the following stochastic algebraic Riccati equation (SARE)

$$\begin{aligned} PA + A^T P + C^T PC + Q - (PB + C^T PD) \\ \times (R + D^T PD)^{-1}(B^T P + D^T PC) = 0. \end{aligned} \tag{6}$$

Due to the nonlinear structure of SARE (6), the analytical solution of (6) is difficult to obtain. To our best knowledge, there are some iterative algorithms to get the approximation solution of (6), one of which is the PI method developed in Ni and Fang [18]. We summarize the method as the following lemma.

**Lemma 1.** Assume $[A, C|Q]$ is exactly detectable. For a given stabilizer $K_0$, let $P_i \in \mathbf{S}^n_+$ be the solution of

$$\begin{aligned} P_i(A + BK_i) + (A + BK_i)^T P_i + Q \\ + (C + DK_i)^T P_i(C + DK_i) + K_i^T RK_i = 0, \end{aligned} \tag{7}$$

where $K_i$ is updated by

$$K_{i+1} = -(R + D^T P_i D)^{-1}(B^T P_i + D^T P_i C). \tag{8}$$

Then $P_i$ and $K_i$, $i = 0, 1, 2, 3, \cdots$ can be uniquely determined at each iteration step, and the following conclusions hold:
(i) $K_i$, $i = 0, 1, 2, \cdots$, are stabilizers.
(ii) $\lim_{i \to \infty} P_i = P^*$, $\lim_{i \to \infty} K_i = K^*$, where $P^*$ is a nonnegative definite solution to SARE (6) and

$K^* = -(R + D^T P^* D)^{-1}(B^T P^* + D^T P^* C)$. In this case, $u^*(\cdot) = K^* x^*(\cdot)$ is an optimal control of Problem (LQS).

Note that the above method needs all knowledge of the system matrices, which are difficult to obtain in the real world. Thus we want to develop a model-free algorithm to solve $P_i$ and $K_i$ without using the information of the coefficient matrices $A$, $B$, $C$, $D$ in system (1).

## 3  MODEL-FREE PI ALGORITHM

In this section, we present our data-driven PI algorithm that does not rely on all knowledge of the coefficient matrices in system (1).

To this end, we first rewrite (7) as

$$
\begin{aligned}
&A^T P_i + P_i A + C^T P_i C \\
&= - P_i B K_i - K_i^T B^T P_i - Q - K_i^T D^T P_i C \\
&\quad - C^T P_i D K_i - K_i^T D^T P_i D K_i - K_i^T R K_i.
\end{aligned}
\tag{9}
$$

Then, by Ito's formula, we know

$$
\begin{aligned}
&d\big(x(s)^T P_i x(s)\big) \\
&= \Big\{ x(s)^T \big(A^T P_i + P_i A + C^T P_i C\big) x(s) \\
&\quad + 2u(s)^T \big(B^T P_i + D^T P_i C\big) x(s) \\
&\quad + u(s)^T D^T P_i D u(s) \Big\} ds + \big\{ \cdots \big\} dw(s).
\end{aligned}
\tag{10}
$$

Combining it with (9), we have

$$
\begin{aligned}
&d\big(x(s)^T P_i x(s)\big) \\
&= \Big\{ - x(s)^T \big(Q + K_i^T R K_i\big) x(s) \\
&\quad + 2\big(u(s) - K_i x(s)\big)^T \big(B^T P_i + D^T P_i C\big) x(s) \\
&\quad + u(s)^T D^T P_i D u(s) \\
&\quad - x(s)^T K_i^T D^T P_i D K_i x(s) \Big\} ds + \big\{ \cdots \big\} dw(s).
\end{aligned}
\tag{11}
$$

Integrating (11) from $t$ to $t + \triangle t$ and taking expection $\mathbb{E}$, we get

$$
\begin{aligned}
&\mathbb{E}\big[x(t + \triangle t)^T P_i x(t + \triangle t) - x(t)^T P_i x(t)\big] \\
&\quad - 2\mathbb{E}\int_t^{t+\triangle t} \big(u(s) - K_i x(s)\big)^T M_i x(s) ds \\
&\quad - \mathbb{E}\int_t^{t+\triangle t} u(s)^T H_i u(s) ds \\
&\quad + \mathbb{E}\int_t^{t+\triangle t} x(s)^T K_i^T H_i K_i x(s) ds \\
&= - \mathbb{E}\int_t^{t+\triangle t} x(s)^T \big(Q + K_i^T R K_i\big) x(s) ds,
\end{aligned}
\tag{12}
$$

where $M_i = B^T P_i + D^T P_i C$, $H_i = D^T P_i D$, $t \geq 0$, $\triangle t$ is any positive real number and $x(\cdot)$ is governed by system (1) with any control $u(\cdot)$.

Next, we give some symbols to develop our data-driven model-free PI algorithm. We define matrices $\eta_{\overline{x}} \in \mathbb{R}^{q \times \frac{n(n+1)}{2}}$, $\eta_{\overline{u}} \in \mathbb{R}^{q \times \frac{m(m+1)}{2}}$, $\eta_{\overline{K_i x}} \in \mathbb{R}^{q \times \frac{m(m+1)}{2}}$, $i = 0, 1, 2, \cdots$, $\eta_{xu} \in \mathbb{R}^{q \times mn}$ and $\eta_{xx} \in \mathbb{R}^{q \times n^2}$, as follows

$$
\eta_{\overline{x}} = \mathbb{E}\left[ \overline{x(t_1)} - \overline{x(t_0)}, \cdots, \overline{x(t_q)} - \overline{x(t_{q-1})} \right]^T,
$$

$$
\eta_{\overline{u}} = \mathbb{E}\left[ \int_{t_0}^{t_1} \overline{u(s)} ds, \cdots, \int_{t_{q-1}}^{t_q} \overline{u(s)} ds \right]^T,
$$

$$
\eta_{\overline{K_i x}} = \mathbb{E}\left[ \int_{t_0}^{t_1} \overline{K_i x(s)} ds, \cdots, \int_{t_{q-1}}^{t_q} \overline{K_i x(s)} ds \right]^T,
$$

$$
\eta_{xx} = \mathbb{E}\left[ \int_{t_0}^{t_1} x(s) \otimes x(s) ds, \cdots, \int_{t_{q-1}}^{t_q} x(s) \otimes x(s) ds \right]^T,
$$

$$
\eta_{xu} = \mathbb{E}\left[ \int_{t_0}^{t_1} x(s) \otimes u(s) ds, \cdots, \int_{t_{q-1}}^{t_q} x(s) \otimes u(s) ds \right]^T,
$$

where $q \in \mathbb{Z}^+$ is any positive integer and $0 \leq t_0 < t_1 < t_2 < \cdots < t_q$.

For any given $K_i$, (12) implies

$$
\Psi_i \begin{bmatrix} vech(P_i) \\ vec(M_i) \\ vech(H_i) \end{bmatrix} = \Theta_i,
\tag{13}
$$

where $\Psi_i \in \mathbb{R}^{q \times (\frac{n(n+1)}{2} + mn + \frac{m(m+1)}{2})}$ and $\Theta_i \in \mathbb{R}^q$ are defined as

$$
\Theta_i = \big[ - \eta_{xx} vec(Q + K_i^T R K_i) \big],
$$

$$
\Psi_i = \big[ \eta_{\overline{x}}, 2\eta_{xx}(I_n \otimes K_i^T) - 2\eta_{xu}, \eta_{\overline{K_i x}} - \eta_{\overline{u}} \big].
$$

If $\Psi_i$ has full column rank for any $i \in \mathbb{Z}$, (13) can be directly transformed to

$$
\begin{bmatrix} vech(P_i) \\ vec(M_i) \\ vech(H_i) \end{bmatrix} = (\Psi_i^T \Psi_i)^{-1} \Psi_i^T \Theta_i.
\tag{14}
$$

Next, we show that, under condition (15) in the following lemma, $\Psi_i$, $i = 0, 1, \cdots$, has full column rank.

**Lemma 3.** If there exists a $q_0 \in \mathbb{Z}^+$, such that, for all $q \geq q_0$,

$$rank([\eta_{xx}, \; \eta_{xu}, \; \eta_{\overline{u}}]) = \frac{n(n+1)}{2} + mn + \frac{m(m+1)}{2}, \tag{15}$$

then, $\Psi_i$, $i = 0, 1, \cdots$, has full column rank.

**Proof.** It is enough to prove that

$$\Psi_i V = O, \quad \forall i \in \mathbb{Z}, \tag{16}$$

has the unique solution $V = O$, where $O$ is a zero matrix (or vector) with appropriate dimension and $V \in \mathbb{R}^{mn + \frac{n(n+1)}{2} + \frac{m(m+1)}{2}}$ .

To achieve it, we now prove it by contradiction. We assume $V = [vech(N)^T, vec(F)^T, vech(G)^T]^T \in \mathbb{R}^{mn + \frac{n(n+1)}{2} + \frac{m(m+1)}{2}}$ is a nonzero column vector, where $vech(N) \in \mathbb{R}^{\frac{n(n+1)}{2}}$, $vec(F) \in \mathbb{R}^{mn}$ and $vech(G) \in \mathbb{R}^{\frac{m(m+1)}{2}}$. Then, by the definitions of $vech(\cdot)$ and $vec(\cdot)$, two symmetric matrices $N \in \mathbf{S}^n$, $G \in \mathbf{S}^m$ and a matrix $F \in \mathbb{R}^{m \times n}$ can be uniquely determined by $vech(N)$, $vech(G)$ and $vec(F)$, respectively.

Applying Ito's formula to $x(s)^T N x(s)$, we derive

$$\mathbb{E}\big[x(t + \triangle t)^T N x(t + \triangle t) - X(t)^T N X(t)\big]$$
$$= \mathbb{E} \int_t^{t+\triangle t} x(s)^T \big(A^T N + N A + C^T N C\big) x(s) ds$$
$$+ 2\mathbb{E} \int_t^{t+\triangle t} u(s)^T B^T N x(s) ds \tag{17}$$
$$+ 2\mathbb{E} \int_t^{t+\triangle t} u(s)^T D^T N C x(s) ds$$
$$+ \mathbb{E} \int_t^{t+\triangle t} u(s)^T D^T N D u(s)\big) ds,$$

where $x(\cdot)$ is governed by system (1) with the same input $u(\cdot)$ as in (12).

Using (12), (17) and the definition of $\Psi_i$, we have

$$\Psi_i V = \eta_{xx} vec(\mathcal{T}) + \eta_{xu} vec(\mathcal{J}) + \eta_{\overline{u}} vech(\mathcal{L}), \tag{18}$$

where

$$\mathcal{T} = A^T N + N A + C^T N C + K_i^T G K_i \\ + K_i^T F + F^T K_i \tag{19}$$

$$\mathcal{J} = 2B^T N + 2D^T N C - 2F, \tag{20}$$

$$\mathcal{L} = D^T N D - G. \tag{21}$$

Since $\mathcal{T}$ is a symmetric matrix, we get

$$\eta_{xx} vec(\mathcal{T}) = I_{\overline{x}} vech(\mathcal{T}), \tag{22}$$

where $I_{\overline{x}} \in \mathbb{R}^{q \times \frac{n(n+1)}{2}}$ and

$$I_{\overline{x}} = \mathbb{E}\bigg[ \int_{t_0}^{t_1} \overline{x(s)} ds, \cdots, \int_{t_{q-1}}^{t_q} \overline{x(s)} ds \bigg]^T. \tag{23}$$

Then, (16) and (18) imply

$$[I_{\overline{x}}, \eta_{xu}, \eta_{\overline{u}}] \begin{pmatrix} vech(\mathcal{T}) \\ vec(\mathcal{J}) \\ vech(\mathcal{L}) \end{pmatrix} = O. \tag{24}$$

It is easy to see that $[I_{\overline{x}}, \eta_{xu}, \eta_{\overline{u}}]$ has full column rank under condition (15). Then, the solution to (24) is $vech(\mathcal{T}) = O$, $vec(\mathcal{J}) = O$ and $vech(\mathcal{L}) = O$, and thus $\mathcal{T} = O, \mathcal{J} = O$ and $\mathcal{L} = O$.

Next, since $K_i$ is a stabilizer, by Definition 1, we know the trajectory of

$$\begin{cases} dx(s) = \Big[(A + BK_i)x(s)\Big] ds \\ \qquad\quad + \Big[(C + DK_i)x(s)\Big] dw(s), \\ x(0) = x_0 \in \mathbb{R}^n \end{cases} \tag{25}$$

satisfies $\lim_{s \to +\infty} \mathbb{E}\big[x(s)^T x(s)\big] = 0$.

For any $t > 0$, applying Ito's formula to $d\big(x(s)^T N x(s)\big)$, we get

$$\mathbb{E}\big[x^T(t) N x(t)\big] - x_0^T N x_0$$
$$= \mathbb{E} \int_0^t x^T(s)\big((A + BK_i)^T N + N(A + BK_i) \tag{26}$$
$$+ (C + DK_i)^T N(C + DK_i)\big) x(s) ds,$$

where $x(\cdot)$ is governed by (25).

Then, by (19), (20), (21), $\mathcal{T} = 0$, $\mathcal{J} = 0$ and $\mathcal{L} = 0$, we can easily see from (26) that $\mathbb{E}\big[x^T(t) N x(t)\big] - x_0^T N x_0 = 0$. Letting $t \to +\infty$, we have $x_0^T N x_0 = \lim_{t \to +\infty} \mathbb{E}\big[x^T(t) N x(t)\big] = 0$. Notice that $x_0$ can be any element in $\mathbb{R}^n$ and $N \in \mathcal{S}^n$, we have $N = 0$. Then it follows from (19), (20), (21), $\mathcal{T} = 0$, $\mathcal{J} = 0$ and $\mathcal{L} = 0$ that $G = 0$ and $F = 0$, which contradicts with $V \neq 0$. The proof is completed. ∎

Using above notations, our model-free algorithm is given in Algorithm 1.

Finally, we show the convergence of our algorithm.

**Algorithm 1**

---

1: Initial $i = 0$ and select $K_0$ as a stabilizer for system (1). Take $u(\cdot) = K_0 x(\cdot) + e(\cdot)$ as the input to system (1) on time interval $[t_0, t_q]$, where $e(\cdot)$ is the exploration noise. Calculate $\eta_{\overline{x}}$, $\eta_{\overline{u}}$, $\eta_{xu}$ and $\eta_{xx}$.
2: **repeat**
3: Compute $\eta_{\overline{K_i x}}$ and solve $P_i$, $M_i$ and $H_i$ from (14).
4: $K_{i+1} = -(R + H_i)^{-1} M_i$.
5: $i \leftarrow i + 1$.
6: **Until** $|P_{i+1} - P_i| < \varepsilon$.

---

**Theorem 1.** Under rank condition (15), starting from a stabilizer $K_0$, the sequences $\{P_i\}_{i=0}^{\infty}$ and $\{K_i\}_{i=1}^{\infty}$ obtained from Algorithm 1 satisfy $\lim_{i \to \infty} P_i = P^*$ and $\lim_{i \to \infty} K_i = K^*$.

**Proof.** Given a stabilizer $K_i$, if $P_i \in \mathbf{S}^{n \times n}$ is the solution of (7), $M_i$ and $H_i$ can be uniquely determined by $M_i = B^T P_i + D^T P_i C$ and $H_i = D^T P_i D$, respectively. Thus, (12) implies that $P_i$, $M_i$ and $H_i$ must satisfy (14).

Moreover, if (15) holds, (14) has the unique solution $(P_i, M_i, H_i)$. Otherwise, (14) has two different solutions and thus contradicts with rank condition (15).

Therefore, under condition (15), $P_i$ and $K_i$, $i = 0, 1, 2, \cdots$, obtained from Algorithm 1 are equivalent to the solution of (7) and (8). Then the convergence of the proposed algorithm can be guaranteed by Lemma 1. ∎

## 4 NUMERICAL EXAMPLE

This section will present a simulation example to illustrate the feasibility of Algorithm 1.

We consider system (1) with $n = 2$ and $m = 1$,

$$A = \begin{bmatrix} 0 & -0.6 \\ 0.6 & -0.3 \end{bmatrix}, B = \begin{bmatrix} 0.05 \\ 0.01 \end{bmatrix},$$

$$C = \begin{bmatrix} -0.02 & 0.03 \\ -0.05 & 0.02 \end{bmatrix}, D = \begin{bmatrix} 0.001 \\ 0.03 \end{bmatrix},$$

and $x_0 = [0.5, -0.1]^T$. The weighting matrices in the cost functional are choosed as $R = 1 > 0$ and $Q = diag(1, 0.5) \geq 0$.

By implementing Algorithm 1, we can obtain

$$\widetilde{P}^* = \begin{bmatrix} 2.9072352 & -0.8296538 \\ -0.8296538 & 2.4975686 \end{bmatrix},$$

$$\widetilde{K}^* = \begin{bmatrix} -0.0669434 & 0.0064058 \end{bmatrix}.$$

Moreover, to check the error of the proposed algorithm, we denote the left sides of (6) and (7) as $\mathcal{R}_1(P)$ and $\mathcal{R}_2(P, K)$. Then we have $|\mathcal{R}_1(\widetilde{P}^*)| = 2.0820041 \times 10^{-3}$ and $|\mathcal{R}_2(\widetilde{P}^*, \widetilde{K}^*)| = 2.0833488 \times 10^{-3}$.

## 5 CONCLUSION

This paper has developed a model-free PI algorithm to solve infinite-horizon LQS problems, i.e., Problem (LQS). By applying ADP techniques, the solution of Problem (LQS) can be learned from the collected data. Moreover, an example is given to show the applicability of the obtained algorithm.

## REFERENCES

[1] B. Kiumarsi, K.G. Vamvoudakis, H. Modares, F.L. Lewis, Optimal and autonomous control using reinforcement learning: A survey, IEEE Trans. Neural Netw. Learn. Syst. 29 (6) (2017) 2042-2062.

[2] B. Pang, Z. Jiang, I. Mareels, Reinforcement learning for adaptive optimal control of continuous-time linear periodic systems, Automatica 118 (2020) 1-9.

[3] D. Vrabie, O. Pastravanu, M. Abu-Khalaf, F.L. Lewis, Adaptive optimal control for continuous-time linear systems based on policy iteration, Automatica 45 (2) (2009) 477-484.

[4] K.G. Vamvoudakis, Q-learning for continuous-time linear systems: A model-free infinite horizon optimal control approach, Systems Control Lett. 100 (2017) 14–20.

[5] M. Ait Rami, X. Zhou, Linear matrix inequalities, riccati equations, and indefinite stochastic linear quadratic controls, IEEE Trans. Automat. Control 45 (6) (2000) 1131-1143.

[6] M. Palanisamy, H. Modares, F.L.Lewis, M. Aurangzeb, Continuous-time q-learning for infinite-horizon discounted cost linear quadratic regulator problems, IEEE Trans. Cybern. 45 (2) (2015) 165–176.

[7] P.J. Werbos, Beyond regression: new tools for prediction and analysis in the behavioural sciences, Ph.D. Thesis, Harvard University, 1974.

[8] Q. Wei, H. Zhang, J. Dai, Model-free multiobjective approximate dynamic programming for discrete-time nonlinear systems with general performance index functions, Neurocomputing 72 (7) (2009) 1839–1848.

[9] R.S. Sutton, A.G. Barto, Reinforcement learning: an introduction, MIT Press, 1998.

[10] S. Mukherjee, H. Bai, A. Chakrabortty, Model-based and model-free designs for an extended continuous-time LQR with exogenous inputs, Systems Control Lett. 154 (2021) 1-9.

[11] T. Bian, Z. Jiang, Value iteration and adaptive dynamic programming for data-driven adaptive optimal control design, Automatica 71 (2016) 348-360.

[12] T. Damm, D. Hinrichsen, Newton's method for a rational matrix equation occuring in stochastic control, Linear Algebra Appl. 332–334 (2001) 81–109.

[13] T. Wang, H. Zhang, Y. Luo, Infinite-time stochastic linear quadratic optimal control for unknown discrete-time systems using adaptive dynamic programming approach, Neurocomputing 171 (2016) 379-386.

[14] T. Wang, H. Zhang, Y. Luo, Stochastic linear quadratic optimal control for model-free discrete-time systems based on Q-learning algorithm, Neurocomputing 312 (2018) 1-8.

[15] W.M. Wonham, On a matrix Riccati equation of stochastic control, SIAM J. Control 6 (4) (1968) 681-697.

[16] Y. Jiang, Z. Jiang, Computational adaptive optimal control for continuous-time linear systems with completely unknown dynamics, Automatica 48 (10) (2012) 2699-2704.

[17] Y. Jiang, Z. Jiang, Approximate dynamic programming for optimal stationary control with control-dependent noise, IEEE Trans. Neural Netw. 22 (12) (2011) 2392–2398.

[18] Y. Ni, H. Fang, Policy iteration algorithm for singular controlled diffusion processes, SIAM J. control optim. 51 (5) (2013) 3844-3862.

[19] Y. Ge, X. Liu, Y. Li, Optimal control for unknown mean-field discrete-time system based on Q-Learning, Int. J. Syst. Sci. 52 (15) (2021) 1-15.

[20] Z. Shi, Z. Wang, Adaptive output-feedback optimal control for continuous-time linear systems based on adaptive dynamic programming approach, Neurocomputing 438 (2021) 334–344.