

随机线性二次控制的资格迹方法

朱亚楠

上海理工大学, 理学院, 上海

收稿日期: 2023年12月20日; 录用日期: 2023年12月29日; 发布日期: 2024年1月31日

摘要

本文研究了强化学习方法在线性二次控制问题(LQR)中的应用。在LQR问题的研究中, 常见的方法通过求解代数黎卡提方程得到最优控制, 并不直接优化控制增益。本文在策略梯度算法的基础上引入资格迹方法, 直接优化控制增益矩阵。考虑已知和未知参数两种情况下, 资格迹方法的收敛。在有限时域和高斯噪声的条件下, 分别给出了已知和未知参数两种情况下算法的全局收敛保证。参数未知时, 利用零阶优化定理近似梯度项, 这可以将问题扩展至代价函数非凸的情况。数值模拟结果显示资格迹方法与梯度下降算法相比更快收敛, 方差更小。

关键词

线性二次最优控制, 梯度下降, 资格迹

Eligibility Trace Method for Stochastic Linear Quadratic Control

Yanan Zhu

College of Science, University of Shanghai for Science and Technology, Shanghai

Received: Dec. 20th, 2023; accepted: Dec. 29th, 2023; published: Jan. 31st, 2024

Abstract

This paper studies the application of reinforcement learning method to linear quadratic regulator (LQR) problem. For the study of LQR problem, the usual method is to obtain the optimal control by solving the algebraic Riccati equation, but not to optimize the control gain directly. This paper optimizes the control gain directly, proposes the eligibility trace method on the basis of gradient descent algorithm, and produces global convergence guarantee in the case of known and unknown parameters, in the setting of finite time horizon and Gaussian noise. When the parameters are unknown, the zero-order optimization theorem is used to approximate the gradient term, which

can extend the problem to cases where the cost function is not convex. Numerical simulation results show that the eligibility trace method has faster convergence and smaller variance than gradient descent algorithm.

Keywords

Linear Quadratic Optimal Control, Gradient Descent, Eligibility Traces

Copyright © 2024 by author(s) and Hans Publishers Inc.

This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

1. 引言

线性二次型控制(LQR)问题是最优控制理论中最基本的问题之一,系统下一时刻的状态由当前状态和控制决定,代价是关于状态和控制的二次函数[1]。问题的目标是找到最优的系统控制,使得代价函数达到最小值。在系统参数完全已知的情况下,LQR问题已经有广泛的研究,最优控制存在解析表达式[2],被广泛应用在不同领域中。在实际应用中,获取系统所有参数时不现实的,当系统参数未知时,产生许多不确定性,无法直接求解得到最优控制,这促进了对学习方法的探索。

近年来使用强化学习和数据驱动方法解决不确定动力系统问题取得重要进展,强化学习方法成为解决此类不确定问题的有效手段。强化学习(RL)[3]是交互式的学习方法,通过与环境交互积累经验(采样),以最大化数值收益信号为导向,不断从经验中学习,最终得到最优策略(控制)。在已有的研究中,利用RL方法解决参数未知的LQR问题,有基于模型的方法[4][5][6]和无模型的方法[7][8][9]两种。基于模型的RL方法利用经验数据估计系统模型参数,将参数估计值视为真实值进而通过解析表达式求解最优控制[10][11];无模型的RL方法直接与系统环境交互,利用经验数据直接学习最优控制而不估计参数[12][13]。LQR问题的最优控制是状态的线性函数,在学习方法的框架下,相当于优化矩阵形式的参数。自然的想法是在参数空间中使用梯度下降算法,将控制参数化,即RL中的策略梯度算法。这种方法基于代价函数对控制参数的梯度,对参数进行优化。有模型的RL方法在模型结构已知时具有更好的表现,无模型方法则更少依赖于模型假设,具有更好的泛化性能[14]。

RL方法在实践中效果突出,但理论理解较少,相比之下,控制理论已经有较为成熟的研究成果,具有可证明的保证。系统参数已知时,控制优化过程中需要的梯度项有明确的表达式,易于计算;参数未知时,用经验数据对梯度项进行近似,现有的梯度近似有零阶近似、二阶近似以及样条逼近等方法。Dhruv Malik等人给出了零阶近似方法在线性控制问题中的收敛保证[15]。Maryam Fazel等人研究了初始状态有噪声干扰的情况下,不同梯度优化方的全局收敛性[16]。Ben Hambly等人在有限时域和随机状态转移的设置下为零阶近似方法提供了全局收敛保证[17]。Mesbahi A等人将结果推广至连续时间的LQR问题[18][19]。

LQR问题描述

有限时域 T 下的随机离散时间LQR问题,系统方程为:

$$x_{t+1} = Ax_t + Bu_t + \omega_t \quad (1)$$

其中, $x_t \in \mathbb{R}^d$ 和 $u_t \in \mathbb{R}^k$ 分别是系统在 t 时刻的状态向量和控制向量, x_0 从分布 \mathcal{D} 中随机抽样, $\omega_t (t \in [T] = 0, 1, \dots, T-1)$ 是独立同分布的零均值系统噪声序列, 与 x_0 独立。 $A \in \mathbb{R}^{d \times d}$ 和 $B \in \mathbb{R}^{d \times k}$ 是系统转

移矩阵。系统(1)对应的代价函数为:

$$C(\mathbf{K}) = E \left[\sum_{t=0}^{T-1} (x_t^T Q_t x_t + u_t^T R_t u_t) + x_T^T Q_T x_T \right] \quad (2)$$

其中, $Q_t \in \mathbb{R}^{d \times d}$ 是对称正定矩阵, $R_t \in \mathbb{R}^{k \times k}$ 是正定矩阵, $\mathbf{K} = (K_0, K_1, \dots, K_{T-1})$ 是控制序列。

离散时间 LQR 问题的目标是在可行控制空间 \mathcal{U} 中找到使代价函数值最小的最优控制, 在系统参数已知时, 有限时域的离散时间 LQR 问题的最优控制有显式解, 当 (A, B) 可控时, 根据最优控制理论[20], t 时刻的最优控制为当前状态的线性函数:

$$\begin{aligned} u_t^* &= -K_t^* x_t \\ K_t^* &= (B^T P_{t+1}^* B + R_t)^{-1} B^T P_{t+1}^* A \quad \forall t \in [T] \end{aligned} \quad (3)$$

其中, $K_t^* \in \mathbb{R}^{k \times d}$ 是 t 时刻的最优控制, P_t^* 是以下代数黎卡提方程的解:

$$P_t^* = Q_t + A^T P_{t+1}^* A - A^T P_{t+1}^* B (B^T P_{t+1}^* B + R_t)^{-1} B^T P_{t+1}^* A \quad (4)$$

黎卡提方程的求解通常令 $P_T = Q_T$, 通过递推式(6)进行迭代求解, 直到满足精度要求。此时, 代价函数达到最小值

$$C(\mathbf{K}^*) = E \left[\sum_{t=0}^{T-1} x_t^T P_t^* x_t \right].$$

上述通过求解代数黎卡提方程进而得到最优控制的方法有以下缺点: 1) 并不直接优化控制增益 K_t , 也不直接优化代价函数; 2) 在求解黎卡提方程时涉及高维矩阵的求逆运算, 消耗的计算量在系统维度高时过大; 3) 当系统参数 $\theta = (A, B, \{Q_t\}_{t=0}^{T-1}, \{R_t\}_{t=0}^{T-1})$ 未知时, 这种方法并不适用。

本文研究有限时域下的离散时间随机线性二次控制(DLQR)问题, 在策略梯度方法中引入资格迹方法, 不求解代数黎卡提方程, 直接优化控制增益。资格迹方法是强化学习的基本方法之一, 在优化中引入与控制参数向量同维度的资格迹向量, 衡量系统控制参数每个分量的重要性, 在学习过程中影响参数向量的更新[3]。这种方法不仅适用于分幕式的情况, 也适用于持续性问题。在实际效果方面, 资格迹方法比策略梯度算法更快收敛, 且方差更小, 数值模拟结果证明了这一结论。

2. 理论

本文考虑折现的随机 DLQR 问题:

$$\begin{aligned} \min_{\mathbf{K}} \quad & C(\mathbf{K}) = E \left[\sum_{t=0}^{T-1} \gamma^t c_t(x_t, u_t) + \gamma^T c_T(x_T) \right] \\ \text{s.t.} \quad & x_{t+1} = Ax_t + Bu_t + \omega_t \end{aligned} \quad (5)$$

其中 $\gamma \in (0, 1)$ 是折现因子, 单步的代价 $c_t(x_t, u_t) = x_t^T Q_t x_t + u_t^T R_t u_t, t \in [T]$, $c_T(x_T) = x_T^T Q_T x_T$ 。控制增益遵循以下法则优化:

$$\begin{aligned} \mathbf{K}^{n+1} &= \mathbf{K}^n - \alpha \delta^n \\ \delta^0 &= \nabla C(\mathbf{K}^0) \quad \forall t \in [T] \\ \delta^n &= \lambda \delta^{n-1} + \nabla C(\mathbf{K}^n), n > 0 \end{aligned} \quad (6)$$

其中, $n \in [N]$ 是迭代次数, α 是步长参数, λ 是衰减参数, $\mathbf{K}^n = (K_0^n, K_1^n, \dots, K_{T-1}^n)$ 是第 n 次迭代时的控制序列, $\delta^n = (\delta_0^n, \delta_1^n, \dots, \delta_{T-1}^n)$ 是与之对应的资格迹序列。

引理 2.1 在策略参数的第 n 次迭代中, 定义状态向量的协方差矩阵为:

$$\Sigma_t = E[x_t x_t^T], \quad \Sigma_{\mathbf{K}} = \sum_{t=0}^{T-1} \Sigma_t \quad (7)$$

则与 \mathbf{K}_t 对应的资格迹由以下表达式给出:

$$\delta_t^n = \begin{cases} 2\gamma^t E_t^n \Sigma_t^n & n=0 \\ \lambda \delta_t^{n-1} + 2\gamma^t E_t^n \Sigma_t^n & n>0 \end{cases}, \forall t \in [T] \quad (8)$$

其中, $E_t^n = (R_t + \gamma B^T P_{t+1}^n B) K_t - \gamma B^T P_{t+1}^n A$ 。

证明: 根据资格迹的定义,

$$\delta_t^n = \lambda \delta_t^{n-1} + \nabla_t C(\mathbf{K}^n), \quad n > 0 \quad (9)$$

随机 DLQR 问题(5)中, 代价函数可拆分为初始状态项与噪声项:

$$\begin{aligned} C(\mathbf{K}) &= E \left[\sum_{t=0}^{T-1} \gamma^t c_t(x_t, u_t) + \gamma^T c_T(x_T) \right] \\ &= E \left[\sum_{t=0}^{T-2} \gamma^t c_t(x_t, u_t) + \gamma^{T-1} c_{T-1}(x_{T-1}, -K_{T-1} x_{T-1}) + \gamma^T c_T(x_T) \right] \\ &= E \left[\sum_{t=0}^{T-2} \gamma^t c_t(x_t, u_t) + \gamma^{T-1} x_{T-1}^T P_{T-1} x_{T-1} + \gamma^T \omega_{T-1}^T P_T \omega_{T-1} \right] \\ &= E \left[x_0^T P_0 x_0 + \sum_{t=0}^{T-1} \gamma^{t+1} \omega_t^T P_{t+1} \omega_t \right] \end{aligned} \quad (10)$$

上式对 K_t 求偏导,

$$\begin{aligned} \nabla_t C(\mathbf{K}) &= \frac{\partial C(\mathbf{K})}{\partial K_t} \\ &= \frac{\partial E \left[\gamma^t x_t^T (Q_t + K_t^T R_t K_t + \gamma (A - BK_t)^T P_{t+1} (A - BK_t)) x_t + \mathbf{K}(-t) \right]}{\partial K_t} \\ &= E \left[2\gamma^t R_t K_t x_t x_t^T - 2\gamma^{t+1} B^T P_{t+1} (A - BK_t) x_t x_t^T \right] \\ &= 2\gamma^t E_t \Sigma_t \end{aligned} \quad (11)$$

其中, $\mathbf{K}(-t) = \sum_{j=0}^{t-1} \gamma^j c_j(x_j, u_j) + \sum_{j=t+1}^{T-1} \gamma^j \omega_j^T P_{j+1} \omega_j$ 证毕。

对衰减参数 λ 的不同设置, 是衡量历史信息对下一步决策的重要程度, 当 $\lambda=0$ 时, 不考虑历史信息, 当 $\lambda=1$ 时, 则认为历史信息与当前信息同样重要。资格迹方法考虑了当前的损失函数和历史策略梯度的关系, 而梯度下降算法只考虑梯度更新的平滑度。资格迹方法能够减少参数更新过程中忽略影响较大的历史动量而造成的错误决策次数。

2.1. 参数已知的资格迹方法

本节讨论系统参数 θ 完全已知时的资格迹方法[3]。资格迹方法结合了蒙特卡洛方法和时序差分方法。算法的核心是定义一个维度与策略参数相同的短时记忆向量 δ , 衡量策略参数分量的重要性。随着迭代次数的增加, 和参与更新的控制参数分量对应的资格迹分量逐渐衰减, 直到这一分量再次参与更新。本文用资格迹代替参数更新中的梯度项, 增强了算法的泛化性能。

在分析算法收敛性之前，首先研究不同控制序列下产生的代价函数的差异以及状态协方差矩阵间的差异上界，定义 $\|S\|$ 为矩阵 S 的谱范数。

引理 2.2 [17] 假设 \mathbf{K}' 与 \mathbf{K} 产生的代价函数均有界， $\{x_t\}_{t=0}^{T-1}$ ， $\{u_t\}_{t=0}^{T-1}$ ， $\{x'_t\}_{t=0}^{T-1}$ ， $\{u'_t\}_{t=0}^{T-1}$ 分别是由 \mathbf{K}' ， \mathbf{K} 生成的序列，令 $x'_0 = x_0 = x$ ，则代价差可表示为：

$$C(\mathbf{K}') - C(\mathbf{K}) = E \left[\sum_{t=0}^{T-1} 2Tr \left(x'_t (x'_t)^T (K'_t - K_t)^T E_t \right) \right] + E \left[\sum_{t=0}^{T-1} Tr \left(x'_t (x'_t)^T (K'_t - K_t)^T (R_t + \gamma B^T P_{t+1} B) (K'_t - K_t) \right) \right] \quad (12)$$

其中， P_t 是下述方程的解：

$$P_t = Q_t + K_t^T R_t K_t + \gamma (A - BK_t)^T P_{t+1} (A - BK_t), \quad \forall t \in [T] \quad (13)$$

引理 2.3 令 $\rho = \max \left\{ \max_t \|A - BK_t\|, \max_t \|A - BK'_t\| \right\}$ ， $\Delta = K_t - K'_t$ ， \mathbf{K}' 与 \mathbf{K} 是任意策略，系统状态向量的协方差满足下面的关系：

$$\|\Sigma_{\mathbf{K}} - \Sigma_{\mathbf{K}'}\| \leq \frac{\rho^{2T} - 1}{\rho^2 - 1} \left(2\rho \|B\| \sum_{t=0}^{T-1} \|\Delta\| + \|B\|^2 \sum_{t=0}^{T-1} \|\Delta\|^2 \right) \left(\frac{C(\mathbf{K})}{\sigma_{\min} \mathbf{Q}} + T \|W\| \right) \quad (14)$$

证明详见附录。

上面的分析为收敛保证奠定了基础，证明算法的收敛性之前，引理 2.4 的论证了控制序列经过一次迭代后对代价函数值的影响。

引理 2.4 设 \mathbf{K}^* 是最优至序列， \mathbf{K}' 由 \mathbf{K} 经一次迭代得到，当

$$\alpha \leq \min \left\{ \alpha_1, \frac{\sigma_{\min} \mathbf{Q}}{2C(\mathbf{K})T \max_t \|R_t + \gamma B^T P_{t+1} B\|} \right\}$$

其中，

$$\alpha_1 = \frac{\rho^2 - 1}{2(\rho^{2T} - 1)(2\rho + 1)\|B\|} \cdot \frac{\sigma_{\min} \mathbf{Q} \sigma_{\min} \Sigma_{\mathbf{K}}}{C(\mathbf{K}) + T \|W\| \sigma_{\min} \mathbf{Q}} \cdot \frac{1}{\max_t \|\delta_t\|}$$

则

$$C(\mathbf{K}') - C(\mathbf{K}^*) \leq \left(1 - \frac{4\alpha \sigma_{\min} \mathbf{R} \sigma_{\min}^2 \Sigma_{\mathbf{K}}}{\|\Sigma_{\mathbf{K}^*}\|} \right) C(\mathbf{K}) - C(\mathbf{K}^*) \quad (15)$$

证明详见附录。

经过以上分析，下面给出参数已知时，资格迹算法在 DLQR 问题中的全局收敛性保证。

定理 2.4 (收敛性定理) 假设 $C(\mathbf{K}^0)$ 有界，步长 α 满足引理 2.4 的约束，对 $\forall \varepsilon > 0$ ，当迭代次数 N 满足下述条件：

$$N \geq \frac{\|\Sigma_{\mathbf{K}^*}\|}{4\alpha \sigma_{\min}^2 \Sigma_{\mathbf{K}} \sigma_{\min} \mathbf{R}} \log \frac{C(\mathbf{K}^0) - C(\mathbf{K}^*)}{\varepsilon}$$

代价函数值收敛至最优值，即：

$$C(\mathbf{K}^N) - C(\mathbf{K}^*) \leq \varepsilon \quad (16)$$

证明：令 $\mathbf{K}^1 = \mathbf{K}^0 - \alpha \delta^0$ ，根据引理 2.4 的结论，

$$C(\mathbf{K}^1) - C(\mathbf{K}^*) \leq \left(1 - \frac{4\alpha\sigma_{\min}\mathbf{R}\sigma_{\min}^2\Sigma_{\mathbf{K}}}{\|\Sigma_{\mathbf{K}^*}\|} \right) (C(\mathbf{K}^0) - C(\mathbf{K}^*))$$

假设经 $n+1$ 次迭代后， $C(\mathbf{K}^{n+1}) \leq C(\mathbf{K}^0)$ ，此时 $K_t^{n+1} = K_t^n - \alpha\delta_t^n$ ，根据 Cauchy-Schwarz 不等式，

$$\begin{aligned} \sum_{t=0}^{T-1} \delta_t^n &= \sum_{t=0}^{T-1} \sum_{i=0}^n \lambda^{n-i} \nabla_i C(\mathbf{K}^n) \leq \sum_{t=0}^{T-1} \sqrt{n \sum_{i=0}^n \|\nabla_i C(\mathbf{K}^n)\|^2} \\ &\leq \sum_{t=0}^{T-1} \sqrt{4n \sum_{i=0}^n \text{Tr}(\Sigma_t^i (E_t^i)^T E_t^i \Sigma_t^i)} \\ &\leq \sqrt{T \cdot \sum_{t=0}^{T-1} 4n \sum_{i=0}^n \|\Sigma_t^i\|^2 \text{Tr}((E_t^i)^T E_t^i)} \\ &\leq \left(\frac{2C(\mathbf{K})}{\sigma_{\min}\mathbf{Q}} \right) \sqrt{nT \frac{\max_t (R_t + \gamma B^T P_t B)}{\sigma_{\min}\Sigma_{\mathbf{K}}} (C(\mathbf{K}) - C(\mathbf{K}^*))} \end{aligned}$$

最后一个不等式成立基于下面的结果：

$$\begin{aligned} C(\mathbf{K}) - C(\mathbf{K}^*) &\geq C(\mathbf{K}) - C(\mathbf{K}') \\ &= E \left[\sum_{t=0}^{T-1} \text{Tr} \left(E_t^T (R_t + \gamma B^T P_{t+1} B)^{-1} E_t \right) \right] \\ &\geq \frac{\sigma_{\min}\Sigma_{\mathbf{K}}}{\max_t (R_t + \gamma B^T P_t B)} \sum_{t=0}^{T-1} \text{Tr}(E_t^T E_t) \end{aligned}$$

结合引理 2.3 的分析，引理 2.3 中的结论仍然成立，即：

$$C(\mathbf{K}^{n+1}) - C(\mathbf{K}^*) \leq \left(1 - \frac{4\alpha\sigma_{\min}\mathbf{R}\sigma_{\min}^2\Sigma_{\mathbf{K}}}{\|\Sigma_{\mathbf{K}^*}\|} \right) (C(\mathbf{K}^n) - C(\mathbf{K}^*)) \quad (17)$$

将 $n+1$ 次的结果进行累积，

$$C(\mathbf{K}^{n+1}) - C(\mathbf{K}^*) \leq \left(1 - \frac{4\alpha\sigma_{\min}\mathbf{R}\sigma_{\min}^2\Sigma_{\mathbf{K}}}{\|\Sigma_{\mathbf{K}^*}\|} \right)^{n+1} (C(\mathbf{K}^0) - C(\mathbf{K}^*))$$

对 $\forall \varepsilon > 0$ ，当

$$N \geq \frac{\|\Sigma_{\mathbf{K}^*}\|}{4\alpha\sigma_{\min}^2\Sigma_{\mathbf{K}}\sigma_{\min}\mathbf{R}} \log \frac{C(\mathbf{K}^0) - C(\mathbf{K}^*)}{\varepsilon}$$

时， $C(\mathbf{K}^N) - C(\mathbf{K}^*) < \varepsilon$ 。证毕。

2.2. 参数未知的资格迹方法

当系统参数 θ 未知时，式(13)中对梯度项的表达无效，在控制增益的更新中，性能指标的资格迹无法以解析式精确表达，只能对模型进行仿真得到训练数据。本文使用零阶优化方法[6]近似资格迹，无需估计系统参数，直接利用目标函数的值优化控制增益。并且零阶优化方法对目标函数的凸性没有要求，直接用函数值估计函数梯度，这大大降低了算法的复杂度。本节首先分析零阶优化方法在随机最优优化问题上的拓展[11] [18]，在 DLQR 问题中，参数未知时，在每一步的控制上加入随机噪声进行采样来估计

代价函数值。目标函数可表示为

$$C(\mathbf{K}) = E_{\omega} [C(\mathbf{K}; \omega)] \quad (18)$$

这里利用带噪声的代价函数值构造梯度的近似无偏估计。令 $\mathcal{U}^r = \{\mathbf{U} \in \mathbb{R}^{k \times d} : \|\mathbf{U}\|_F = r\}$ ，设 \mathcal{P}_0 是 \mathcal{U}^r 上的均匀分布。任意给一个度量 $r > 0$ ，以及 $\mathbf{U} \sim \mathcal{P}_0$ 与 ω 独立，则零阶优化方法对 $C(\mathbf{K})$ 的梯度估计为：

$$\nabla C(\mathbf{K}) = \frac{k \times d}{r} C(\mathbf{K} + r\mathbf{U})\mathbf{U} \quad (19)$$

随着 r 越来越小，近似值越来越精确[11]，但 r 过小可能导致估计值有较大的方差。表 1 给出了算法的框架。

Table 1. Eligibility track-gradient algorithm for DLQR problem

表 1. DLQR 的资格迹 - 梯度算法

算法：资格迹 - 梯度算法

1. 输入： \mathbf{K} ，最大迭代次数 M ，采样轨迹数 N ，幕长 T ，随机项参数 r ，维度 D
2. for $n=0, 1, \dots, N-1$:
3. for $i=0, 1, \dots, I-1$:
4. for $t=0, 1, \dots, T-1$:
 - 1) 从 $x_0^i \in \mathcal{D}$ 开始，根据 $(\mathbf{K}_{-t}, \hat{K}_t^i) = (K_0, \dots, K_{t-1}, \hat{K}_t^i, K_{t+1}, \dots, K_{T-1})$ 采样，其中 $\hat{K}_t^i = K_t + U_t^{ni}$ $\|U_t^{ni}\|_F = r$.
 - 2) 记录单幕代价 \hat{c}_t^i .
5. 计算资格迹的估计值：

$$\hat{\delta}_t^n = \begin{cases} \frac{1}{I} \sum_{i=0}^{I-1} \frac{D}{r^2} \hat{c}_t^i U_t^i & n=0 \\ \lambda \hat{\delta}_t^{n-1} + \frac{1}{I} \sum_{i=0}^{I-1} \frac{D}{r^2} \hat{c}_t^i U_t^i & n>0 \end{cases}$$

6. $\mathbf{K}^{n+1} = \mathbf{K}^n - \alpha \hat{\delta}^n$

引理 2.5 对给定的 $r > 0$ 以及从 $\mathcal{U}^r = \{\mathbf{U} \in \mathbb{R}^d : \|\mathbf{U}\|_F = r\}$ 中随机抽取的随机向量 \mathbf{U} ， I 为采样幕数， λ 是折扣系数，资格迹的经验近似为：

$$\hat{\delta}_t^n = \begin{cases} \frac{1}{I} \sum_{i=0}^{I-1} \frac{D}{r^2} \hat{c}_t^i U_t^i & n=0 \\ \lambda \hat{\delta}_t^{n-1} + \frac{1}{I} \sum_{i=0}^{I-1} \frac{D}{r^2} \hat{c}_t^i U_t^i & n>0 \end{cases} \quad (20)$$

其中，

$$\hat{c}_t^i = \sum_{s=0}^{T-1} \gamma^s \left((x_t^i)^T Q_t x_t^i + (u_t^i)^T R_t u_t^i \right) + (x_t^i)^T Q_T x_t^i$$

引理 2.6 假设任意不同控制 \mathbf{K}' 与 \mathbf{K} 的分量满足：

$$\|K'_i - K_i\| \leq \min \left\{ \|K_i\|, \frac{(\rho^2 - 1) \sigma_{\min} \mathbf{Q} \sigma_{\min} \Sigma_{\mathbf{K}}}{2T(\rho^{2T} - 1)(2\rho + 1)(C(\mathbf{K}) + \sigma_{\min} \mathbf{Q} T \|W\|) \|B\|} \right\} \quad (21)$$

则存在

$$h_c \leq \left\{ \frac{\rho^2 - 1}{(2\rho + 1)(\rho^{2T} - 1)} \cdot \frac{1}{\|B\|} \cdot \frac{1}{\|W\|} \cdot \frac{1}{C(\mathbf{K}^0)} \right\}, \quad h_g \leq \left\{ \frac{\rho^2 - 1}{(2\rho + 1)(\rho^{2T} - 1)} \cdot \frac{1}{\|B\|} \cdot \frac{1}{\|W\|} \cdot \frac{\sigma_{\min} \mathbf{Q}}{C(\mathbf{K}^0)} \cdot \frac{1}{\|\Sigma_{\mathbf{K}^0}\|} \right\}$$

使得,

$$\|C(\mathbf{K}') - C(\mathbf{K})\| \leq h_c \sum_{t=0}^{T-1} \|K'_t - K_t\|, \quad \|\nabla_t C(\mathbf{K}') - \nabla_t C(\mathbf{K})\| \leq h_g \sum_{t=0}^{T-1} \|K'_t - K_t\|$$

证明详见附录。

引理 2.7 对 $\forall \varepsilon > 0$, 当 $r \leq \frac{\varepsilon}{4h_g}$ 时, 资格迹的经验近似满足:

$$\|\hat{\delta}_t^n - \delta_t^n\|_F \leq \varepsilon, \quad \forall t \in [T]$$

证明详见附录。

定理 2.8 (收敛性定理) 假设 $C(\mathbf{K}^0)$ 有界, 步长 α 满足引理 2.3 的约束, 对 $\forall \varepsilon > 0$, 当迭代次数 N 满足

$$N \geq \frac{\|\Sigma_{\mathbf{K}^*}\|}{2\alpha\sigma_{\min}^2 \Sigma_{\mathbf{K}} \sigma_{\min} \mathbf{R}} \log \frac{C(\mathbf{K}^0) - C(\mathbf{K}^*)}{\varepsilon}$$

代价函数值收敛至最优值, 即:

$$C(\mathbf{K}^N) - C(\mathbf{K}^*) \leq \varepsilon \quad (22)$$

证明: 令 $K'_t = K_t - \alpha\delta_t$, 根据引理 2.4, 有

$$C(\mathbf{K}') - C(\mathbf{K}^*) \leq \left(1 - \frac{4\alpha\sigma_{\min} \mathbf{R} \sigma_{\min}^2 \Sigma_{\mathbf{K}}}{\|\Sigma_{\mathbf{K}^*}\|} \right) (C(\mathbf{K}) - C(\mathbf{K}^*)) \quad (23)$$

令 $\hat{K}'_t = K_t - \alpha\hat{\delta}_t$, 当样本数足够多时, 只要 $C(\mathbf{K}) - C(\mathbf{K}^*) \geq \varepsilon$ 时,

$$C(\hat{\mathbf{K}}') - C(\mathbf{K}^*) \leq \left(1 - \frac{4\alpha\sigma_{\min} \mathbf{R} \sigma_{\min}^2 \Sigma_{\mathbf{K}}}{\|\Sigma_{\mathbf{K}^*}\|} \right) (C(\mathbf{K}) - C(\mathbf{K}^*)) \quad (24)$$

因为

$$C(\hat{\mathbf{K}}') - C(\mathbf{K}^*) = C(\hat{\mathbf{K}}') - C(\mathbf{K}') + C(\mathbf{K}') - C(\mathbf{K}^*)$$

首先讨论右式第一项, 根据引理 2.6, 当 $\forall t \in T$, 若

$$\|\hat{K}'_t - K'_t\| \leq \frac{\alpha\sigma_{\min} \mathbf{R} \sigma_{\min}^2 \Sigma_{\mathbf{K}} \varepsilon}{T \|\Sigma_{\mathbf{K}^*}\| h_c}$$

成立, 则有

$$\|C(\hat{\mathbf{K}}') - C(\mathbf{K}')\| \leq \frac{\alpha\sigma_{\min} \mathbf{R} \sigma_{\min}^2 \Sigma_{\mathbf{K}} \varepsilon}{\|\Sigma_{\mathbf{K}^*}\|}$$

因为 $\hat{K}'_t - K'_t = \alpha(\delta_t - \hat{\delta}_t)$, 所以只需证明

$$\|\delta_t - \hat{\delta}_t\| \leq \frac{\sigma_{\min} \mathbf{R} \sigma_{\min}^2 \Sigma_{\mathbf{K}} \varepsilon}{T \|\Sigma_{\mathbf{K}^*}\| h_c} \quad (25)$$

根据引理 2.7, 上式成立。综上所述

$$C(\hat{\mathbf{K}}') - C(\mathbf{K}^*) \leq \left(1 - \frac{3\alpha\sigma_{\min} \mathbf{R}\sigma_{\min}^2 \Sigma_{\mathbf{K}}}{\|\Sigma_{\mathbf{K}^*}\|} \right) (C(\mathbf{K}) - C(\mathbf{K}^*)) \quad (26)$$

基于以上结论，与定理 2.4 相同的证明保证了收敛性。

3. 数值模拟

当状态维度 $d=2$ 时，设定系统参数为

$$A = \begin{pmatrix} 0.9 & -1.1 \\ 1.0 & 0.8 \end{pmatrix}, B = \begin{pmatrix} 0.8 \\ 1 \end{pmatrix}, Q = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, R = 1$$

折现因子 $\gamma=0.99$ ，比较资格迹方法与梯度下降算法的收敛情况。在折扣系数 $\lambda=0.1$ 的条件下，设定指数衰减的步长参数，分别令时域 $T=20$ ，迭代次数 $N=40$ ；时域 $T=50$ ，迭代次数 $N=70$ ，代价函数的收敛情况见图 1 和图 2：

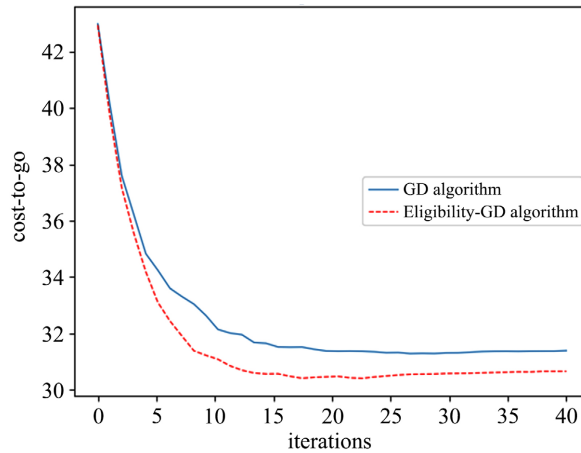


Figure 1. The convergence of $C(K)$ when $d=2, T=20$

图 1. $d=2, T=20$ ，代价函数的收敛情况

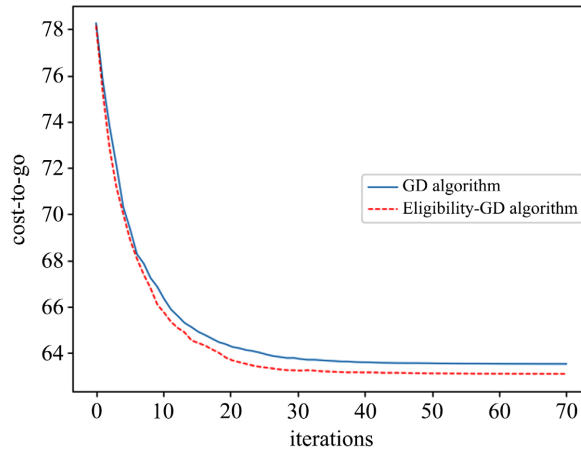


Figure 2. The convergence of $C(K)$ when $d=2, T=50$

图 2. $d=2, T=50$ ，代价函数的收敛情况

图 1 和图 2 的结果说明，资格迹算法比梯度下降算法具有更快的收敛速度。资格迹方法中折扣系数的取值对最终结果有显著影响，图 3 展示了 $T=50, N=70$ 时不同的折扣系数对算法性能的影响：

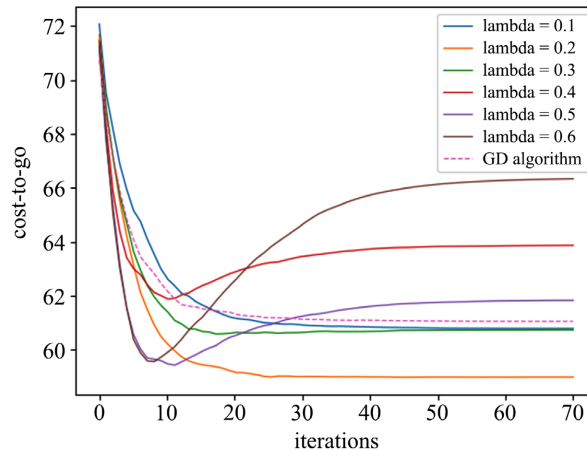


Figure 3. The convergence of $C(K)$ with different λ

图 3. 不同 λ , 代价函数的收敛情况

在本节设定的系统参数下，图 3 的结果显示，折扣系数 $\lambda < 0.3$ 时，资格迹算法表现优于策略梯度算法，当 $\lambda > 0.3$ 后，结果出现不收敛的情况，随 λ 的增大，结果偏离越大。 λ 是过去梯度信息的权重，说明在这一数值范例中，过去梯度信息对问题求解只能提供少量信息。

当状态维度 $d = 4$ 时，设定系统参数为

$$A = \begin{pmatrix} 0.3 & 0.3 & 0.1 & 0.2 \\ 0.2 & 0.2 & 0.3 & 0.1 \\ 0.3 & 0.1 & 0.3 & 0.4 \\ 0.3 & 0.2 & 0.15 & 0.1 \end{pmatrix}, B = \begin{pmatrix} 0.4 & 0.3 \\ -0.5 & 0.1 \\ 1.0 & 0.2 \\ 0.2 & 0.9 \end{pmatrix}$$

$$Q = \begin{pmatrix} 0.9 & 0.2 & -0.5 & 0.15 \\ 0.2 & 1.1 & 0.15 & 0.1 \\ -0.5 & 0.15 & 0.9 & -0.8 \\ 0.15 & 0.1 & -0.8 & 0.88 \end{pmatrix}, R = \begin{pmatrix} 0.9 & -0.5 \\ -0.5 & 0.8 \end{pmatrix}$$

设定 $T = 50$, $N = 100$, 其余参数与二维系统相同，代价函数收敛情况如图 4 所示：

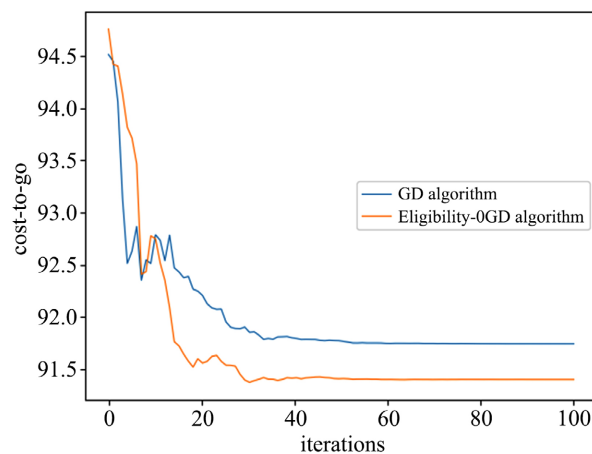


Figure 4. The convergence of $C(K)$ when $d = 4$, $T = 50$

图 4. $d = 4$, $T = 50$, 代价函数的收敛情况

上述数值结果说明资格迹方法的收敛效果优于传统的资格梯度方法，并且在算法初期方差更小。

4. 结论

本文研究了无模型强化学习方法在有限时域离散 LQR 问题中的应用，不同于通过解代数黎卡提方程得到最优控制的方法，本文直接优化控制增益，在梯度下降算法的基础上引入资格迹方法，并给出在参数已知和参数未知两种情况下算法的收敛保证。在初始代价函数有界的条件下，算法可以扩展至无限时域。数值模拟验证了算法的收敛性，展示了不同参数设置对结果的影响。另一个方向是基于有模型的强化学习方法，在更少样本量的基础上，进一步达到更好的收敛结果。

致 谢

感谢张老师在论文写作过程中给出的指导和建议。

参考文献

- [1] Birge, J. and Louveaux, F. (2011) Introduction to Stochastic Programming. Springer Science & Business Media, Heidelberg. <https://doi.org/10.1007/978-1-4614-0237-4>
- [2] Kučera, V. (1992) Optimal Control: Linear Quadratic Methods: Brian D. O. Anderson and John B. Moore. *Automatica*, **28**, 1068-1069. [https://doi.org/10.1016/0005-1098\(92\)90166-D](https://doi.org/10.1016/0005-1098(92)90166-D)
- [3] Sutton, R.S. and Barto, A.G. (2018) Reinforcement Learning: An Introduction. 2nd ed., the MIT Press, Cambridge.
- [4] Basei, M., Guo, X., Hu, A. and Zhang, Y. (2020) Logarithmic Regret for Episodic Continuous-Time Linear-Quadratic Reinforcement Learning over a Finite-Time Horizon. *Computation Theory eJournal*. <https://doi.org/10.2139/ssrn.3848428>
- [5] Dean, S., Mania, H., Matni, N., Recht, B. and Tu, S. (2017) On the Sample Complexity of the Linear Quadratic Regulator. *Foundations of Computational Mathematics*, **20**, 633-679. <https://doi.org/10.1007/s10208-019-09426-y>
- [6] Ren, Z., Zhong, A. and Li, N. (2021) LQR with Tracking: A Zeroth-Order Approach and Its Global Convergence. 2021 *American Control Conference (ACC)*, New Orleans, LA, 25-28 May 2021, 2562-2568. <https://doi.org/10.23919/ACC50511.2021.9483417>
- [7] Bertsekas, D.P. (2011) Approximate Policy Iteration: A Survey and Some New Methods. *Journal of Control Theory and Applications*, **9**, 310-335. <https://doi.org/10.1007/s11768-011-1005-3>
- [8] Mania, H., Guy, A. and Recht, B. (2018) Simple Random Search Provides a Competitive Approach to Reinforcement Learning. arXiv preprint arXiv:1803.07055
- [9] Abbasi-Yadkori, Y., Lázic, N. and Szepesvári, C. (2019) Model-Free Linear Quadratic Control via Reduction to Expert Prediction. *The 22nd International Conference on Artificial Intelligence and Statistics*, Naha, 16-18 April 2019, 3108-3117.
- [10] Mahdi, I. and Braga-Neto, U.M. (2018) Finite-Horizon lqr Controller for Partially-Observed Boolean Dynamical Systems. *Automatica*, **95**, 172-179. <https://doi.org/10.1016/j.automatica.2018.05.028>
- [11] Zhang, H. and Li, N. (2022) Data-Driven Policy Iteration Algorithm for Continuous-Time Stochastic Linear-Quadratic Optimal Control Problems. *Asian Journal of Control*, **26**, 481-489. <https://doi.org/10.1002/asjc.3223>
- [12] Farjadnasab, M. and Babazadeh, M. (2022) Model-Free LQR Design by Q-Function Learning. *Automatica*, **137**, Article ID: 110060. <https://doi.org/10.1016/j.automatica.2021.110060>
- [13] Yaghmaie, F.A., Gustafsson, F.K. and Ljung, L. (2023) Linear Quadratic Control Using Model-Free Reinforcement Learning. *IEEE Transactions on Automatic Control*, **68**, 737-752. <https://doi.org/10.1109/TAC.2022.3145632>
- [14] Tu, S. and Recht, B. (2019) The Gap between Model-Based and Model-Free Methods on the Linear Quadratic Regulator: An Asymptotic Viewpoint. *Conference on Learning Theory*, USA, 9 December 2019, 3036-3083.
- [15] Malik, D., Pananjady, A., Bhatia, K., Khamaru, K., Bartlett, P.L. and Wainwright, M.J. (2018) Derivative-Free Methods for Policy Optimization: Guarantees for Linear Quadratic Systems. *Journal of Machine Learning Research*, **21**, 1-21.
- [16] Fazel, M., Ge, R., Kakade, S.M. and Mesbahi, M. (2018) Global Convergence of Policy Gradient Methods for the Linear Quadratic Regulator. *International Conference on Machine Learning*, Stockholm, 10-15 July 2018, 1467-1476.
- [17] Hambly, B.M., Xu, R. and Yang, H. (2021) Policy Gradient Methods for the Noisy Linear Quadratic Regulator over a Finite Horizon. *SIAM Journal on Control and Optimization*, **59**, 3359-3391. <https://doi.org/10.1137/20M1382386>
- [18] Shamir, O. (2017) An Optimal Algorithm for Bandit and Zero-Order Convex Optimization with Two-Point Feedback.

The Journal of Machine Learning Research, **18**, 1703-1713.

- [19] Bu, J., Mesbahi, A. and Mesbahi, M. (2020) Policy Gradient-Based Algorithms for Continuous-Time Linear Quadratic Control. arXiv: 2006.09178.
- [20] Bertsekas, D.P. (1995) Dynamic Programming and Optimal Control. 3rd Edition, Athena Scientific, Nashua, NH.

附录

引理 2.3 证明:

在对 $\Sigma_{\mathbf{k}}$ 进行分析前, 定义以下线性算子以方便证明。令

$$\begin{aligned}
\mathcal{F}_{K_t}(\Sigma) &= (A - BK_t)\Sigma(A - BK_t)^T, \quad \mathcal{G}_t(\Sigma) = \mathcal{F}_{K_t} \circ \mathcal{F}_{K_{t-1}} \circ \cdots \circ \mathcal{F}_{K_0} \\
\Sigma_{t+1} &= E\left(\left((A - BK_t)x_t + \omega_t\right)\left((A - BK_t)x_t + \omega_t\right)^T\right) \\
&= (A - BK_t)\Sigma_t(A - BK_t)^T + W = \mathcal{F}_{K_t}(\Sigma_t) + W \\
&= (A - BK_t)\left(\mathcal{F}_{K_{t-1}}(\Sigma_{t-1}) + W\right)(A - BK_t)^T + W \\
&= \mathcal{F}_{K_t} \circ \mathcal{F}_{K_{t-1}}(\Sigma_{t-1}) + \mathcal{F}_{K_t}(W) + W \\
&= \mathcal{F}_{K_t} \circ \mathcal{F}_{K_{t-1}} \circ \mathcal{F}_{K_{t-2}}(\Sigma_{t-2}) + \mathcal{F}_{K_t} \circ \mathcal{F}_{K_{t-1}}(W) + \mathcal{F}_{K_t}(W) + W \\
&= \mathcal{G}_t(\Sigma_0) + \sum_{s=0}^t \prod_{u=0}^s \mathcal{F}_{K_t} \circ \cdots \circ \mathcal{F}_{K_{t-u}}(W) + W \\
\Sigma_{\mathbf{k}} &= \Sigma_0 + \sum_{t=0}^{T-1} \left[\mathcal{G}_t(\Sigma_0) + \sum_{s=0}^t \prod_{u=0}^s \mathcal{F}_{K_t} \circ \cdots \circ \mathcal{F}_{K_{t-u}}(W) \right] + TW \\
\sum_{t=0}^{T-1} \left\| (\mathcal{G}'_t - \mathcal{G}_t)(\Sigma_0) \right\| &\leq \sum_{t=0}^{T-1} \left\| (\mathcal{F}_{K'_t} \circ \mathcal{G}'_{t-1} - \mathcal{F}_{K'_t} \circ \mathcal{G}_{t-1} + \mathcal{F}_{K'_t} \circ \mathcal{G}_{t-1} - \mathcal{F}_{K_t} \circ \mathcal{G}_{t-1})(\Sigma_0) \right\| \\
&\leq \sum_{t=0}^{T-1} \left\| \mathcal{F}_{K'_t} \right\| \left\| (\mathcal{G}'_{t-1} - \mathcal{G}_{t-1})(\Sigma_0) \right\| + \left\| \mathcal{G}_{t-1} \right\| \left\| (\mathcal{F}_{K'_t} - \mathcal{F}_{K_t})(\Sigma_0) \right\| \\
&\leq \sum_{t=0}^{T-1} \rho^2 \left\| (\mathcal{G}'_{t-1} - \mathcal{G}_{t-1})(\Sigma_0) \right\| + \rho^{2(t+1)} \left\| (\mathcal{F}_{K'_t} - \mathcal{F}_{K_t})(\Sigma_0) \right\| \\
&\leq \frac{\rho^{2T} - 1}{\rho^2 - 1} \sum_{t=0}^{T-1} \left\| (\mathcal{F}_{K'_t} - \mathcal{F}_{K_t})(\Sigma_0) \right\| \\
\sum_{t=0}^{T-1} \left\| (\mathcal{F}_{K_t} - \mathcal{F}_{K'_t})(\Sigma_0) \right\| &= \sum_{t=0}^{T-1} \left\| (A - BK_t)\Sigma_0(A - BK_t)^T - (A - BK'_t)\Sigma_0(A - BK'_t)^T \right\| \\
&= \sum_{t=0}^{T-1} \left\| (A - BK_t)\Sigma_0(B\Delta)^T + (B\Delta)\Sigma_0(A - BK_t)^T - (B\Delta)\Sigma_0(B\Delta)^T \right\| \\
&\leq \sum_{t=0}^{T-1} \left\| \Sigma_0 \left(2\|A - BK_t\| \|B\| \|K_t - K'_t\| + \|B\|^2 \|K_t - K'_t\|^2 \right) \right\| \\
&\leq \left(2\rho \|B\| \sum_{t=0}^{T-1} \|K_t - K'_t\| + \|B\|^2 \sum_{t=0}^{T-1} \|K_t - K'_t\|^2 \right) \left\| \Sigma_0 \right\|
\end{aligned}$$

同理可得,

$$\begin{aligned}
&\sum_{t=0}^{T-1} \left\| \sum_{s=0}^t \prod_{u=0}^s (\mathcal{F}_{K'_t} \circ \cdots \circ \mathcal{F}_{K'_{t-u}} - \mathcal{F}_{K_t} \circ \cdots \circ \mathcal{F}_{K_{t-u}})(W) \right\| \\
&\leq \frac{\rho^{2T} - 1}{\rho^2 - 1} \sum_{t=0}^{T-1} \left\| \sum_{s=0}^t (\mathcal{F}_{K'_t} - \mathcal{F}_{K_t})(W) \right\|
\end{aligned}$$

综上所述,

$$\begin{aligned}
\|\Sigma_{\mathbf{K}} - \Sigma_{\mathbf{K}'}\| &\leq \sum_{t=0}^{T-1} \left\| (\mathcal{G}'_t - \mathcal{G}_t)(\Sigma_0) \right\| + \left\| \sum_{s=0}^t \prod_{u=0}^s (\mathcal{F}_{K'_t} \circ \dots \circ \mathcal{F}_{K'_{t-u}} - \mathcal{F}_{K_t} \circ \dots \circ \mathcal{F}_{K_{t-u}})(W) \right\| \\
&\leq \frac{\rho^{2T} - 1}{\rho^2 - 1} \left(\sum_{t=0}^{T-1} \left\| (\mathcal{F}_{K'_t} - \mathcal{F}_{K_t})(\Sigma_0) \right\| + \sum_{t=0}^{T-1} \left\| \sum_{s=0}^t (\mathcal{F}_{K'_t} - \mathcal{F}_{K_t})(W) \right\| \right) \\
&\leq \frac{\rho^{2T} - 1}{\rho^2 - 1} \left(\sum_{t=0}^{T-1} \left\| \mathcal{F}_{K'_t} - \mathcal{F}_{K_t} \right\| \right) (\|\Sigma_0\| + T\|W\|) \\
&\leq \frac{\rho^{2T} - 1}{\rho^2 - 1} \left(\sum_{t=0}^{T-1} \left\| \mathcal{F}_{K'_t} - \mathcal{F}_{K_t} \right\| \right) (\|\Sigma_0\| + T\|W\|) \\
&\leq \frac{\rho^{2T} - 1}{\rho^2 - 1} \left(2\rho \|B\| \sum_{t=0}^{T-1} \|\Delta\| + \|B\|^2 \sum_{t=0}^{T-1} \|\Delta\|^2 \right) \left(\frac{C(\mathbf{K})}{\sigma_{\min} \mathbf{Q}} + T\|W\| \right)
\end{aligned}$$

引理 2.6 证明:

令 $\phi_t = Q_t + K_t^T R_t K_t$,

$$\begin{aligned}
C(\mathbf{K}') - C(\mathbf{K}) &= \sum_{t=0}^{T-1} \gamma^t \left[\text{Tr}(\Sigma'_t (Q_t + (K'_t)^T R_t K'_t)) - \text{Tr}(\Sigma_t (Q_t + K_t^T R_t K_t)) \right] \\
&= \sum_{t=0}^{T-1} \gamma^t \left[\text{Tr}(\mathcal{G}'_t(\Sigma_0) \phi'_t - \mathcal{G}_t(\Sigma_0) \phi_t) + \text{Tr}((\phi'_t - \phi_t)W) \right] \\
&\quad + \sum_{t=0}^{T-1} \gamma^t \left[\text{Tr} \left(\sum_{s=0}^t \prod_{u=0}^s \mathcal{F}_{K'_t} \circ \dots \circ \mathcal{F}_{K'_{t-u}}(W) \phi'_t - \sum_{s=0}^t \prod_{u=0}^s \mathcal{F}_{K_t} \circ \dots \circ \mathcal{F}_{K_{t-u}}(W) \phi_t \right) \right]
\end{aligned}$$

首先看右边第一项,

$$\left\| \sum_{t=0}^{T-1} \gamma^t \text{Tr}(\mathcal{G}'_t(\Sigma_0) \phi'_t - \mathcal{G}_t(\Sigma_0) \phi_t) \right\| \leq \text{Tr}(\Sigma_0) \sum_{t=0}^{T-1} \gamma^t \left\| \text{Tr}(\mathcal{G}'_t(\phi'_t) - \mathcal{G}_t(\phi_t)) \right\|$$

根据引理 2.3 证明中的定义,

$$\begin{aligned}
\sum_{t=0}^{T-1} \gamma^t \left\| \text{Tr}(\mathcal{G}'_t(\phi'_t) - \mathcal{G}_t(\phi_t)) \right\| &\leq \sum_{t=0}^{T-1} \gamma^t \left\| \text{Tr}(\mathcal{G}'_t(\phi'_t) - \mathcal{G}_t(\phi'_t) + \mathcal{G}_t(\phi'_t) - \mathcal{G}_t(\phi_t)) \right\| \\
&\leq \frac{\rho^{2T} - 1}{\rho^2 - 1} \left[(2\rho + 1) \|B\| \sum_{t=0}^{T-1} \gamma^t \|K_t - K'_t\| \|\phi'_t\| \right] + \sum_{t=0}^{T-1} \gamma^t \|\mathcal{G}_t\| \left\| (K'_t)^T R_t K'_t - K_t^T R_t K_t \right\| \\
&\leq \frac{\rho^{2T} - 1}{\rho^2 - 1} \left[(2\rho + 1) \|B\| \sum_{t=0}^{T-1} \gamma^t \|K_t - K'_t\| \cdot \sum_{t=0}^{T-1} \|\phi'_t\| + \sum_{t=0}^{T-1} \gamma^t \left\| (K'_t)^T R_t K'_t - K_t^T R_t K_t \right\| \right] \\
&\leq \frac{\rho^{2T} - 1}{\rho^2 - 1} (2\rho + 1) \|B\| \left(\sum_{t=0}^{T-1} \|Q_t + K_t^T R_t K_t\| \right) \sum_{t=0}^{T-1} \gamma^t \|K_t - K'_t\| \\
&\quad + \frac{\rho^{2T} - 1}{\rho^2 - 1} \left[(2\rho + 1) \|B\| \sum_{t=0}^{T-1} \|K_t - K'_t\| + 1 \right] \sum_{t=0}^{T-1} \gamma^t \left\| (K'_t)^T R_t K'_t - K_t^T R_t K_t \right\|
\end{aligned}$$

第二项,

$$\sum_{t=0}^{T-1} \gamma^t \text{Tr}((\phi'_t - \phi_t)W) \leq \text{Tr}(W) \sum_{t=0}^{T-1} \gamma^t \left\| (K'_t)^T R_t K'_t - K_t^T R_t K_t \right\|$$

第三项,

$$\begin{aligned}
& \sum_{t=0}^{T-1} \gamma^t \text{Tr} \left[\sum_{s=0}^t \prod_{u=0}^s \mathcal{F}'_{K_t} \circ \dots \circ \mathcal{F}'_{K_{t-u}}(W) \phi'_t - \sum_{s=0}^t \prod_{u=0}^s \mathcal{F}_{K_t} \circ \dots \circ \mathcal{F}_{K_{t-u}}(W) \phi_t \right] \\
&= \sum_{t=0}^{T-1} \gamma^t \text{Tr} \left[\sum_{s=0}^t \left[\prod_{u=0}^s \mathcal{F}'_{K_t} \circ \dots \circ \mathcal{F}'_{K_{t-u}}(W) (\phi'_t - \phi_t + \phi_t) - \prod_{u=0}^s \mathcal{F}_{K_t} \circ \dots \circ \mathcal{F}_{K_{t-u}}(W) \phi_t \right] \right] \\
&\leq \sum_{t=0}^{T-1} \sum_{s=0}^t \text{Tr}(W) \left\| \prod_{u=0}^s \mathcal{F}'_{K_t} \circ \dots \circ \mathcal{F}'_{K_{t-u}} \right\|^2 \sum_{t=0}^{T-1} \gamma^t \left\| (K'_t)^T R_t K'_t - K_t^T R_t K_t \right\| \\
&\quad + \sum_{t=0}^{T-1} \gamma^t \left\| \sum_{s=0}^t \prod_{u=0}^s \mathcal{F}'_{K_t} \circ \dots \circ \mathcal{F}'_{K_{t-u}}(W) - \prod_{u=0}^s \mathcal{F}_{K_t} \circ \dots \circ \mathcal{F}_{K_{t-u}}(W) \right\| \left\| \sum_{t=0}^{T-1} \left[Q_t + K_t^T R_t K_t \right] \right\| \\
&\leq \frac{(T-1)(\rho^{4(T-1)} - 1)}{\rho^4 - 1} \text{Tr}(W) \sum_{t=0}^{T-1} \gamma^t \left\| (K'_t)^T R_t K'_t - K_t^T R_t K_t \right\| \\
&\quad + \frac{T(\rho^{2T} - 1)}{\rho^2 - 1} \left[(2\rho + 1) \|W\| \|B\| \sum_{t=0}^{T-1} \|K_t - K'_t\| \right] \left[\sum_{t=0}^{T-1} \|Q_t + K_t^T R_t K_t\| \right]
\end{aligned}$$

其中,

$$\begin{aligned}
\sum_{t=0}^{T-1} \gamma^t \left\| (K'_t)^T R_t K'_t - K_t^T R_t K_t \right\| &= \sum_{t=0}^{T-1} \gamma^t \left\| (K'_t - K_t + K_t)^T R_t (K'_t - K_t + K_t) - K_t^T R_t K_t \right\| \\
&\leq \sum_{t=0}^{T-1} \gamma^t \|K_t - K'_t\|^2 \|R_t\| + 2\gamma^t \|K_t\| \|R_t\| \|K_t - K'_t\| \\
&\leq \sum_{t=0}^{T-1} 3\gamma^t \|K_t\| \|R_t\| \|K_t - K'_t\|
\end{aligned}$$

根据假设 $\|K_t - K'_t\| \leq \|K_t\|$, 倒数第二个不等式成立, 综上所述,

$$\begin{aligned}
C(\mathbf{K}') - C(\mathbf{K}) &\leq \left\{ \frac{\rho^{2T} - 1}{\rho^2 - 1} \left[\text{Tr}(\Sigma_0) (2\rho + 1) \|B\| \sum_{t=0}^{T-1} \|K_t - K'_t\| + 1 + \text{Tr}(W) \right] + \text{Tr}(W) \right\} \sum_{t=0}^{T-1} 3\gamma^t \|K_t\| \|R_t\| \|K_t - K'_t\| \\
&\quad + \frac{(T-1)(\rho^{4(T-1)} - 1)}{\rho^4 - 1} \text{Tr}(W) \sum_{t=0}^{T-1} 3\gamma^t \|K_t\| \|R_t\| \|K_t - K'_t\| \\
&\quad + \frac{\rho^{2T} - 1}{\rho^2 - 1} (2\rho + 1) \|B\| \left(T \|W\| + \sum_{t=0}^{T-1} \|Q_t + K_t^T R_t K_t\| \right) \sum_{t=0}^{T-1} \gamma^t \|K_t - K'_t\| \\
&\leq h_1 \sum_{t=0}^{T-1} \|K_t - K'_t\| + h_2 \left[\sum_{t=0}^{T-1} \|K_t - K'_t\| \right]^2 \\
&\leq h_c \sum_{t=0}^{T-1} \|K_t - K'_t\|
\end{aligned}$$

下面证明梯度项,

$$\|\nabla_i C(\mathbf{K}') - \nabla_i C(\mathbf{K})\| = 2 \|E'_i \Sigma'_i - E_i \Sigma_i\| \leq 2 \|E'_i - E_i\| \|\Sigma'_i\| + 2 \|E_i\| \|\Sigma'_i - \Sigma_i\|$$

$$\|\Sigma'_i\| \leq \|\Sigma_{\mathbf{K}'}\| \leq \|\Sigma_{\mathbf{K}'} - \Sigma_{\mathbf{K}}\| + \|\Sigma_{\mathbf{K}}\| \leq \frac{C(\mathbf{K})}{\sigma_{\min} \mathbf{Q}} + \|\Sigma_{\mathbf{K}}\|$$

$$\begin{aligned}
\|E'_i - E_i\| &= \|R_i (K'_i - K_i) - B^T (P'_{i+1} - P_{i+1}) A + B^T (P'_{i+1} - P_{i+1}) B K'_i + B^T P_{i+1} B (K'_i - K_i)\| \\
&\leq (\|R_i\| + \|B\|^2 \|P_0\|) \sum_{t=0}^{T-1} \|K_t - K'_t\| + \|B\| \|P'_0 - P_0\| \|A\| + 2 \|B\|^2 \|P'_0 - P_0\| \sum_{t=0}^{T-1} \|K_t\|
\end{aligned}$$

$$\begin{aligned}
\|P'_0 - P_0\| &\leq 3\|K_0\|\|R_0\|\|K'_0 - K_0\| + \left\| \sum_{t=0}^{T-1} \gamma^t \text{Tr}(\mathcal{G}'_t(\Sigma_0)\phi'_t - \mathcal{G}_t(\Sigma_0)\phi_t) \right\| \\
&\quad + \frac{\rho^{2T} - 1}{\rho^2 - 1} \left[(2\rho + 1)\|B\|\|Q_T\| \left\| \sum_{t=0}^{T-1} K_t - K'_t \right\| \right] \\
\|E_t\| &\leq \sum_{t=0}^{T-1} \|E_t\| \leq \sum_{t=0}^{T-1} \sqrt{\text{Tr}(E_t^T E_t)} \leq \sqrt{T \cdot \frac{\max_t \|R_t + \gamma B^T P_{t+1} B\|}{\sigma_{\min} \Sigma_{\mathbf{K}}} (C(\mathbf{K}) - C(\mathbf{K}^*))} \\
\|\Sigma'_t - \Sigma_t\| &\leq \|(\mathcal{G}'_t - \mathcal{G}_t)(\Sigma_0)\| + \left\| \sum_{s=0}^t \prod_{u=0}^s (\mathcal{F}_{K'_t} \circ \dots \circ \mathcal{F}_{K'_{t-u}} - \mathcal{F}_{K_t} \circ \dots \circ \mathcal{F}_{K_{t-u}})(W) \right\| \\
&\leq \left[\rho^{2t} (2\rho + 1)\|B\|\|\Sigma_0\| + \frac{\rho^{2T} - 1}{\rho^2 - 1} (2\rho + 1)\|B\|\|W\| \right] \sum_{t=0}^{T-1} \|K_t - K'_t\| \\
\|\nabla_t C(\mathbf{K}') - \nabla_t C(\mathbf{K})\| &\leq h_g \sum_{t=0}^{T-1} \|K'_t - K_t\|
\end{aligned}$$

引理 2.7 证明过程

定义

$$\tilde{\nabla}_t C(\mathbf{K}) \triangleq \frac{1}{I} \sum_{i=0}^{I-1} \frac{D}{r^2} (C(\mathbf{K} + \mathbf{U}^i) U_i^i)$$

则

$$\hat{\nabla}_t C(\mathbf{K}) - \nabla_t C(\mathbf{K}) = \hat{\nabla}_t C(\mathbf{K}) - \tilde{\nabla}_t C(\mathbf{K}) + \tilde{\nabla}_t C(\mathbf{K}) - \nabla_t C(\mathbf{K})$$

因为 $E_{x_0, \omega} [\hat{\nabla}_t C(\mathbf{K})] = \tilde{\nabla}_t C(\mathbf{K})$, 根据霍夫丁不等式, 对 $\forall \zeta > 0$,

$$P\left\{ \|\hat{\nabla}_t C(\mathbf{K}) - \tilde{\nabla}_t C(\mathbf{K})\| \leq \zeta \right\} \geq 1 - 2 \exp\left(-\frac{2I\zeta^2}{L^2}\right)$$

因为

$$\begin{aligned}
\|\hat{c}_t^i\| &\leq \sum_{r=0}^{T-1} \gamma^r \left\| \left((x_r^i)^T Q_r x_r^i + (u_r^i)^T R_r u_r^i \right) + (x_r^i)^T Q_r x_r^i \right\| \\
&\leq \sum_{r=0}^{T-1} \gamma^r \|x_r^i\|^2 \left\| Q_r + (K_r^i)^T R_r K_r^i \right\| \\
&\leq \max_t \left\| Q_t + (K_t^i)^T R_t K_t^i \right\| \sum_{r=0}^{T-1} \gamma^r \|\Sigma_r^i\| \\
&\leq \max_t \left\| Q_t + (K_t^i)^T R_t K_t^i \right\| \left(\frac{C(\mathbf{K})}{\sigma_{\min} \mathbf{Q}} + \|\Sigma_{\mathbf{K}^0}\| \right) \triangleq L
\end{aligned}$$

所以, $\|\hat{\nabla}_t C(\mathbf{K}) - \tilde{\nabla}_t C(\mathbf{K})\|_F \leq \frac{\varepsilon}{2}$.

又 $\nabla_t C^r(\mathbf{K}) = \frac{D}{r^2} E_{\mathbf{U} \sim \mathbb{B}_r} [C(\mathbf{K} + \mathbf{U}_t) U_t]$, $\tilde{\nabla}_t C(\mathbf{K}) \triangleq \frac{1}{I} \sum_{i=0}^{I-1} \frac{D}{r^2} (C(\mathbf{K} + \mathbf{U}^i) U_i^i)$

所以, $E[\tilde{\nabla}_t C(\mathbf{K})] = \nabla_t C^r(\mathbf{K})$, $\|\tilde{\nabla}_t C(\mathbf{K}) - \nabla_t C^r(\mathbf{K})\|_F \leq \frac{\varepsilon}{4}$

$\|\nabla_t C^r(\mathbf{K}) - \nabla_t C(\mathbf{K})\|_F \leq \|\nabla_t C^r(\mathbf{K}) - \nabla_t C(\mathbf{K} + \mathbf{U}_t)\|_F + \|\nabla_t C(\mathbf{K} + \mathbf{U}_t) - \nabla_t C(\mathbf{K})\|_F$

根据引理 2.9, 当 $r \leq \frac{\varepsilon}{4h_g}$ 时, $\|\nabla_i C(\mathbf{K} + \mathbf{U}_i) - \nabla_i C(\mathbf{K})\|_F \leq h_g \|\mathbf{U}_i\|_F \leq \frac{\varepsilon}{8}$

所以, $\|\nabla_i C^r(\mathbf{K}) - \nabla_i C(\mathbf{K})\|_F \leq \frac{\varepsilon}{4}$

综上所述, $\|\hat{\nabla}_i C(\mathbf{K}) - \nabla_i C(\mathbf{K})\| \leq \varepsilon$ 。根据资格迹与梯度的关系,

$$\|\hat{\delta}_i^n - \delta_i^n\|_F = \sum_{i=0}^n \lambda^{n-i} \|\hat{\nabla}_i C(\mathbf{K}^n) - \nabla_i C(\mathbf{K}^n)\|_F \leq \sum_{i=0}^n \lambda^{n-i} \varepsilon$$

因为 $\lambda \in (0,1)$, 由 ε 的任意性可知结论成立。