# Convergence of policy gradient methods for finite-horizon exploratory linear-quadratic control problems

Michael Giegrich*        Christoph Reisinger*        Yufei Zhang†

**Abstract.** We study the global linear convergence of policy gradient (PG) methods for finite-horizon continuous-time exploratory linear-quadratic control (LQC) problems. The setting includes stochastic LQC problems with indefinite costs and allows additional entropy regularisers in the objective. We consider a continuous-time Gaussian policy whose mean is linear in the state variable and whose covariance is state-independent. Contrary to discrete-time problems, the cost is noncoercive in the policy and not all descent directions lead to bounded iterates. We propose geometry-aware gradient descents for the mean and covariance of the policy using the Fisher geometry and the Bures-Wasserstein geometry, respectively. The policy iterates are shown to satisfy an a-priori bound, and converge globally to the optimal policy with a linear rate. We further propose a novel PG method with discrete-time policies. The algorithm leverages the continuous-time analysis, and achieves a robust linear convergence across different action frequencies. A numerical experiment confirms the convergence and robustness of the proposed algorithm.

**Key words.** Continuous-time linear-quadratic control, policy optimisation, relative entropy, geometry-aware gradient, global linear convergence, mesh-independent convergence

**AMS subject classifications.** 68Q25, 93E20

## 1   Introduction

In recent years, the policy gradient (PG) method and its variants have become an effective tool in seeking optimal polices to control stochastic systems (see e.g., [19, 28, 17, 24, 25]). These algorithms parametrise the policy as a function of the system state, and update the policy parametrisation based on the gradient of the control objective. Most of the progress, especially the convergence analysis of PG methods, has been in discrete-time Markov decision processes (MDPs) (see e.g., [6, 12, 20, 36, 18]). However, most real-world control systems, such as those in aerospace, the automotive industry and robotics, are naturally continuous-time dynamical systems, and hence do not fit in the MDP setting.

One of the most fundamental stochastic control problems is the finite-horizon linear-quadratic control (LQC) problem. It aims to control a linear stochastic differential equation over a given time horizon, subject to a quadratic cost. This problem is important as it provides a reasonable approximation of many nonlinear control problems, and has been used in a wide range of applications, including portfolio optimisation [38, 32], algorithmic trading [5] and production management of exhaustible resources [9]. Moreover, the optimal policy of an LQC problem admits a natural parameterisation as a (time-dependent) linear function of the state, and hence it suffices

---

*Mathematical Institute, University of Oxford, Oxford OX2 6GG, UK (michael.giegrich@maths.ox.ac.uk, christoph.reisinger@maths.ox.ac.uk)

†Department of Mathematics, Imperial College London, London, UK (yufei.zhang@imperial.ac.uk)

to determine the coefficients of this linear function. All these properties make the LQC problem an important theoretical benchmark for studying learning-based control.

**Issues and challenges from continuous-time models.** It is insufficient and improper to rely solely on the analysis and algorithms for discrete-time MDPs to solve continuous-time problems, including LQC problems. There is a mismatch between the algorithm timescale for the former and the underlying systems timescale for the latter. This model mismatch can make conventional discrete-time algorithms very sensitive to the discretisation stepsize. For instance, the empirical studies in [21, 22] suggest that standard PG methods exhibit degraded performance as the agent's action frequency increases (see Section 4 for more details). Similar performance degradation has been observed in [30] for Q-learning methods. Recently, [14] and [15] extend PG and Q-learning methods, respectively, to continuous-time problems without time discretisation, in order to develop algorithms that are robust across different timescales. Nevertheless, the convergence of these algorithms has not been studied, even for LQC problems.

There are technical reasons behind the limited theoretical progress of PG methods for continuous-time LQC problems. The objective of a LQC problem is typically nonconvex with respect to the policies (see Proposition 2.4), analogous to its discrete-time counterpart [6, 36]. This links the convergence analysis of PG methods to the analysis of gradient search for nonconvex objectives, which has always been one of the formidable challenges in optimisation theory. The time-dependent nature of the optimal policy for finite-horizon LQC problems poses new challenges. It requires analysing the optimisation landscape over a suitable infinite-dimensional policy space, instead of in a finite-dimensional parameter space.

One significant new feature of LQC problems with continuous-time policies, in contrast to discrete-time policies, is the *noncoercivity* of the cost function (see Proposition 2.4). Coercivity of the cost means that each sublevel set of the cost is bounded, and this implies that the iterates of a discrete-time algorithm remain bounded as long as the cost decreases along the iteration. This can be ensured by updating the policies along *any descent direction* of the cost with a sufficiently small stepsize. The lack of coercivity of the continuous-time cost function complicates the analysis of PG methods, since for a given descent direction, there may not exist a constant stepsize such that the iterates remain bounded as the algorithm proceeds.

**Our contributions.** This paper proposes convergent PG methods to solve finite-horizon exploratory LQC problems, which generalise classical LQC problems by allowing an entropy regulariser in the objective.

- We reformulate the exploratory LQC problem into a minimisation over Gaussian polices. Each Gaussian policy is parameterised by two time-dependent functions $(K, V)$: the mean is linear in the state with the coefficient $K$, and the covariance is the function $V$. The policy gradient of the cost is characterised by the Pontryagin optimality principle. The cost is shown to satisfy a non-uniform Łojasiewicz condition and a non-uniform smoothness condition (Propositions 2.2 and 2.3). We then prove that the cost is neither coercive nor quasiconvex in $K$, even in a one-dimensional deterministic setting (Proposition 2.4).

- We propose a geometry-aware PG method to solve the LQC problem in continuous time. The gradient for $K$ adapts to the geometry induced by the Fisher information metric (also known as the natural gradient), while the gradient for $V$ adapts to the geometry induced by the Bures-Wasserstein metric. These geometry-aware gradient directions are proved to enjoy an *implicit regularisation* property, i.e., they preserve an $L^2$-bound of $K$, and pointwise

upper and lower bounds of $V$ without an explicit projection step (Proposition 2.5). This allows for exploiting the local regularity of the cost, and proving the PG method converges globally to the optimal policy with a linear rate (Theorem 2.6).

- By leveraging the continuous-time analysis, we propose practically implementable PG methods that take actions at discrete time points, and achieve a linear convergence guarantee independent of the action frequency. Our analysis shows that scaling the discrete-time gradients linearly with respect to action frequency is critical for a robust performance of the algorithm in different timescales (Remark 2.4). The theoretical property is verified through a numerical experiment on an exploratory LQC problem arising from mean-variance portfolio selection problems. This shows that the number of required iterations for conventional PG methods grows linearly in the number of action time points, while the proposed PG methods achieve a robust linear convergence rate over a wide range of action frequencies.

**Our approaches and related works.** Most existing theoretical works of PG methods for LQC problems consider the setting of infinite horizon and deterministic dynamics (see e.g., [6, 3]). For the case with noisy dynamics, existing works focus on discrete-time problems. This includes the setting of infinite horizon and additive noise [16, 36], finite horizon and additive noise [12], and infinite horizon and multiplicative noise [10]. We further refer the reader to [11, 33, 37] for LQ games. In all of these settings, the optimal policy admits a *finite-dimensional* parameterisation.

Compared to existing works, our technical difficulties are three-fold. First, analysing the optimisation landscape over infinite-dimensional continuous-time policies requires continuous-time control theory. For instance, the policy gradient is derived via Pontryagin's maximum principle. The cost regularity (such as Łojasiewicz and smoothness conditions) is proved by using partial differential equation techniques. The lack of cost coercivity also adds complexity to the choice of appropriate descent directions, as discussed in Remark 2.3. Notably, the noncoercivity of the cost function in this context primarily stems from the fact that a policy can have an infinite number of changes in values, occurring at arbitrary time points. This characteristic distinguishes our problem from aforementioned discrete-time scenarios, in which policies change solely at predetermined time points.

Second, the finite-horizon continuous-time setting requires more advanced techniques for the nondegeneracy of the state covariance than the discrete-time setting. In [12, 36], the state covariance is lower bounded by the minimum eigenvalue of the covariance of system noises, *uniformly over all policies*. This bound vanishes as the time discretisation stepsize tends to zero, as the covariance of noise increment typically scales linearly to the stepsize. Moreover, in the present setting, the system noise can degenerate due to a controlled diffusion coefficient. We overcome this difficulty by establishing the positive definiteness of the state covariance *along the policy iterates*. This is possible by a) first estimating the state covariance explicitly using the magnitude of policies, but independent of the system noise (Lemma 3.7), and b) then proving that the geometry-aware gradient directions induce a uniform bound of the iterates. This approach is different from the contraction argument in [23] for problems with uncontrolled diffusion coefficients.

Finally, the possible degeneracy of cost matrices requires sharper estimate of the cost regularity. All existing works assume a running cost of the form $f(x, a) = x^\top Q x + a^\top R a$, with positive definite matrices $Q$ and $R$, and estimate optimisation landscape using minimum eigenvalues of $Q$ and $R$. However, for many applications of stochastic LQC problems, the cost can involve the product of state and control variables [5], or an indefinite weight $R$ [38, 32]. Here, we derive tighter Łojasiewicz and smoothness bounds of the cost using solutions to Lyapunov equations, instead of the cost coefficients. This allows us to consider a general setting where both the drift and diffusion

3

coefficients of the state are controlled, and all cost weights can be negative definite.

**Notation.** For each Euclidean space $E$, we denote by $\langle \cdot, \cdot \rangle$ its usual inner product and $|\cdot|$ the norm induced by $\langle \cdot, \cdot \rangle$. For each $A \in \mathbb{R}^{n \times m}$, we denote by $A^\top$ the transpose of $A$, by $\operatorname{tr}(A)$ the trace of $A$, and by $\|A\|_2$ the spectral norm of $A$. For each $n \in \mathbb{N}$, we denote by $I_n$ the $n \times n$ identity matrix, by $\mathbb{S}^n$, $\overline{\mathbb{S}^n_+}$ and $\mathbb{S}^n_+$ the space of $n \times n$ symmetric, symmetric positive semidefinite, and symmetric positive definite matrices, respectively, and by $\lambda_{\max}(A)$ and $\lambda_{\min}(A)$ the largest and smallest eigenvalues of $A \in \mathbb{S}^n$, respectively. We equip $\mathbb{S}^n$ with the Loewner (partial) order such that for each $A, B \in \mathbb{S}^n$, $A \succeq B$ if $A - B \in \overline{\mathbb{S}^n_+}$. For every measurable functions $F, G : [0, T] \to \mathbb{S}^n$, $F \succeq G$ stands for $F(t) - G(t) \in \overline{\mathbb{S}^n_+}$ for a.e. $t \in [0, T]$.

For each $T > 0$, filtered probability space $(\Omega, \mathcal{F}, \mathbb{F}, \mathbb{P})$ satisfying the usual condition (of right continuity and completeness) and Euclidean space $(E, |\cdot|)$, we introduce the following spaces:

- $\mathcal{B}(0, T; E)$ is the space of Borel measurable functions $\phi : [0, T] \to E$.
- $L^p(0, T; E)$, $p \in [1, \infty]$, is the space of Borel measurable functions $\phi : [0, T] \to E$ satisfying $\|\phi\|_{L^p} = (\int_0^T |\phi_t|^p \, \mathrm{d}t)^{1/p} < \infty$ if $p \in [1, \infty)$ and $\|\phi\|_{L^\infty} = \operatorname{ess\,sup}_{t \in [0,T]} |\phi_t| < \infty$.
- $C([0, T]; E)$ is the space of continuous functions $\phi : [0, T] \to E$ endowed with the norm $\|\cdot\|_{L^\infty}$.
- $\mathcal{S}^2(0, T; E)$ is the space of $\mathbb{F}$-progressively measurable càdlàg processes $X : \Omega \times [0, T] \to E$ satisfying $\|X\|_{\mathcal{S}^2} = \mathbb{E}[\operatorname{ess\,sup}_{t \in [0,T]} |X_t|^2]^{1/2} < \infty$;
- $\mathcal{M}(E)$ is the set of measures on $E$, $\mathcal{P}(E)$ is the set of probability measures on $E$, and $\mathcal{P}_2(E)$ is the set of square integrable probability measures on $E$ endowed with the 2–Wasserstein distance.

For each $\mu \in \mathbb{R}^n$ and $\Sigma \in \overline{\mathbb{S}^n_+}$, we denote by $\mathcal{N}(\mu, \Sigma)$ the Gaussian measure on $\mathbb{R}^n$ with mean $\mu$ and covariance matrix $\Sigma$. We also write $\mathbb{N}_0 = \mathbb{N} \cup \{0\}$ for notation simplicity.

# 2 Problem formulation and main results

This section introduces exploratory LQC problems, proposes a class of geometry-aware PG algorithms to seek the optimal policy, and presents their convergence properties.

## 2.1 Regularised stochastic LQ control problems with indefinite costs

This section recalls the regularised LQC problem introduced in [31, 32] and its optimal feedback controls. Let $T > 0$ be a finite time horizon, $(\Omega, \mathcal{F}, \mathbb{P})$ be a complete filtered probability space on which a $d$-dimensional standard Brownian motion $W = (W_t)_{t \geq 0}$ is defined, and $\mathbb{F} = (\mathcal{F}_t)_{t \geq 0}$ be the natural filtration of $W$ augmented by an independent $\sigma$-algebra $\mathcal{F}_0$.

We first introduce the admissible controls and the associated state dynamics. Let $\mathcal{A}$ be the set of (relaxed) controls $\mathfrak{m} : \Omega \to \mathcal{M}([0, T] \times \mathbb{R}^k)$ such that $\mathfrak{m}_t(\mathrm{d}t, \mathrm{d}a) = \mathfrak{m}_t(\mathrm{d}a)\mathrm{d}t$ for a.e. $t \in [0, T]$, where $\mathfrak{m}_t : \Omega \to \mathcal{P}(\mathbb{R}^k)$ is $\mathcal{F}_t$-measurable for all $t \in [0, T]$ and $\mathbb{E}[\int_0^T \int_{\mathbb{R}^k} |a|^2 \mathfrak{m}_t(\mathrm{d}a)\mathrm{d}t] < \infty$. For each $\mathfrak{m} \in \mathcal{A}$, consider the following controlled dynamics:

$$\mathrm{d}X_t = \Phi_t(X_t, \mathfrak{m}_t) \, \mathrm{d}t + \Gamma_t(X_t, \mathfrak{m}_t) \, \mathrm{d}W_t, \quad t \in [0, T]; \quad X_0 = \xi_0, \tag{2.1}$$

where $\xi_0 \in L^2(\Omega; \mathbb{R}^d)$ is a given $\mathcal{F}_0$-measurable random variable, and the functions $\Phi : [0, T] \times \mathbb{R}^d \times \mathcal{P}_2(\mathbb{R}^k) \to \mathbb{R}^d$ and $\Gamma : [0, T] \times \mathbb{R}^d \times \mathcal{P}_2(\mathbb{R}^k) \to \overline{\mathbb{S}^d_+}$ satisfy for all $(t, x, m) \in [0, T] \times \mathbb{R}^d \times \mathcal{P}_2(\mathbb{R}^k)$,

$$\Phi_t(x, m) = \int_{\mathbb{R}^k} (A_t x + B_t a) \, m(\mathrm{d}a), \quad \Gamma_t(x, m) = \left( \int_{\mathbb{R}^k} (C_t x + D_t a)(C_t x + D_t a)^\top m(\mathrm{d}a) \right)^{\frac{1}{2}}, \tag{2.2}$$

where $(\cdot)^{\frac{1}{2}} : \overline{\mathbb{S}_+^d} \to \overline{\mathbb{S}_+^d}$ is the matrix square root such that $M^{\frac{1}{2}}(M^{\frac{1}{2}})^\top = M$ for all $M \in \overline{\mathbb{S}_+^d}$, and $A, B, C, D$ are measurable functions such that (2.1) admits a unique strong solution $X^{\mathfrak{m}} \in \mathcal{S}^2(0, T; \mathbb{R}^d)$ (see (H.1) for precise conditions).

The state dynamics (2.1) is commonly referred to as an exploratory dynamics (see, e.g., [31, 32, 26]). It models interacting with the system by repeatedly sampling random actions according to a given measure-valued control $\mathfrak{m}$. As a consequence of these random actions, the system's state evolves with the aggregated coefficients (2.2), which indicates that the infinitesimal change of the state at $t$ has a mean and variance integrated with respect to the sampling distribution $\mathfrak{m}_t$. In the special case where $\mathfrak{m}_t(\mathrm{d}t, \mathrm{d}a) = \boldsymbol{\delta}_{\alpha_t}(\mathrm{d}a)\mathrm{d}t$ for some $\alpha_t : \Omega \times [0, T] \to \mathbb{R}^k$, with $\boldsymbol{\delta}_a$ being the Dirac measure on $a \in \mathbb{R}^k$, (2.1) simplifies into

$$\mathrm{d}X_t = (A_t X_t + B_t \alpha_t)\,\mathrm{d}t + (C_t X_t + D_t \alpha_t)\,\mathrm{d}W_t, \quad t \in [0, T]; \quad X_0 = \xi_0, \tag{2.3}$$

which is the dynamics studied in the classical LQC problem [34]. See the end of Section 2.4 for more details on the connection between an exploratory state dynamics and controlling (2.3) with random actions.

We now consider minimising the following cost functional over all $\mathfrak{m} \in \mathcal{A}$, which is known as the exploratory/entropy-regularised control problem [31, 32, 26, 14, 15]:

$$\mathbb{E}\left[\int_0^T \int_{\mathbb{R}^k} \left(\frac{1}{2}\left\langle \begin{pmatrix} Q_t & S_t^\top \\ S_t & R_t \end{pmatrix}\begin{pmatrix} X_t^{\mathfrak{m}} \\ a \end{pmatrix}, \begin{pmatrix} X_t^{\mathfrak{m}} \\ a \end{pmatrix} \right\rangle \mathfrak{m}_t(\mathrm{d}a) + \rho \mathcal{H}(\mathfrak{m}_t \| \overline{\mathfrak{m}}_t)\right)\mathrm{d}t + \frac{1}{2}(X_T^{\mathfrak{m}})^\top G X_T^{\mathfrak{m}}\right], \tag{2.4}$$

where $X^{\mathfrak{m}}$ satisfies the state dynamics (2.1). Here $Q, S, R$ are given matrix-valued functions of proper dimensions, $G \in \mathbb{R}^{d \times d}$ and $\rho \geq 0$ are given constants, $(\overline{\mathfrak{m}}_t)_{t \in [0,T]}$ are given measures on $\mathbb{R}^k$, and for each $t \in [0, T]$, $\mathcal{H}(\cdot \| \overline{\mathfrak{m}}_t) : \mathcal{P}(\mathbb{R}^k) \to [0, \infty]$ is the relative entropy with respect to $\overline{\mathfrak{m}}_t$ such that for all $m \in \mathcal{P}(\mathbb{R}^k)$,

$$\mathcal{H}(m \| \overline{\mathfrak{m}}_t) = \begin{cases} \int_{\mathbb{R}} \ln\left(\frac{m(\mathrm{d}a)}{\overline{\mathfrak{m}}_t(\mathrm{d}a)}\right) m(\mathrm{d}a), & m \text{ is absolutely continuous with respect to } \overline{\mathfrak{m}}_t, \\ \infty, & \text{otherwise.} \end{cases}$$

Note that the cost (2.4) is aggregated with respect to the control distribution $\mathfrak{m}_t$ from which the random actions are sampled. The entropy $\mathcal{H}(\cdot \| \overline{\mathfrak{m}}_t)$ serves as a regularisation term to encourage the minimiser of (2.4) to be close to the provided reference measures $(\overline{\mathfrak{m}}_t)_{t \in [0,T]}$, and the weight parameter $\rho \geq 0$ controls the strength of this regularisation.

The entropy-regularised control problem (2.4), initially introduced in [31], represents a natural extension of the well-established regularised MDPs (see e.g., [8, 20]) into the continuous domain. Common choices of $(\overline{\mathfrak{m}}_t)_{t \in [0,T]}$ in the existing literature include Gibbs measures [26] and the Lebesgue measure [31, 32, 7].

The following assumptions on the coefficients of (2.1)-(2.4) are imposed throughout this paper.

**H.1.** *(1)* $T > 0$, $\xi_0 \in L^2(\Omega; \mathbb{R}^d)$, $A \in L^1(0, T; \mathbb{R}^{d \times d})$, $B \in L^2(0, T; \mathbb{R}^{d \times k})$, $C \in L^2(0, T; \mathbb{R}^{d \times d})$, $D \in L^\infty(0, T; \mathbb{R}^{d \times k})$, $Q \in L^1(0.T; \mathbb{S}^d)$, $S \in L^2(0.T; \mathbb{R}^{k \times d})$, $R \in L^\infty(0, T; \mathbb{S}^k)$ *and* $G \in \mathbb{S}^d$.

*(2)* $\rho > 0$, $\overline{\mathfrak{m}}_t = \mathcal{N}(0, \overline{V}_t)$ *for all* $t \in [0, T]$, $\overline{V} \in L^\infty(0, T; \mathbb{S}_+^k)$ *and* $\overline{V} \succeq \delta I_k$ *for some* $\delta > 0$.

*Remark* 2.1. Condition (H.1(1)) ensures that for all $\mathfrak{m} \in \mathcal{A}$, (2.1) admits a unique strong solution in $\mathcal{S}^2(0, T; \mathbb{R}^d)$ (see Proposition A.2), and the associated regularised cost is well-defined. Note that (H.1(1)) allows the coefficients $Q, S, R$ and $G$ to be indefinite or even negative definite (provided that (H.2) holds). Such a control problem is often called indefinite stochastic LQ problem (see e.g. [27] and the references therein) and has important applications in optimal liquidation [5] and mean-variance portfolio selection [38] in finance.

Condition (H.1(2)) assumes that for each $t \in [0,T]$, the reference measure $\overline{\mathfrak{m}}_t$ in (2.4) is a Gaussian measure. This ensures that the optimal strategy of (2.1)-(2.4) is Gaussian (see (2.6)), which in turn implies that (2.1)-(2.4) can be reformulated as an optimisation problem over Gaussian policies. A similar reformulation also holds if $\overline{\mathfrak{m}}_t$ is the Lebesgue measure [31, 32, 7], and our proposed policy descent algorithm and its convergence analysis can be naturally extended to this case.

We also impose the following well-posedness condition of the corresponding Riccati equation for the closed-loop solvability of the (possibly indefinite) control problem (2.1)-(2.4).

**H.2.** *There exists $P^\star \in C([0,T];\mathbb{S}^d)$ satisfying the following Riccati equation: for a.e. $t \in [0,T]$,*

$$
\begin{cases}
(\frac{\mathrm{d}}{\mathrm{d}t}P)_t + A_t^\top P_t + P_t A_t + C_t^\top P_t C_t + Q_t \\
\quad - (B_t^\top P_t + D_t^\top P_t C_t + S_t)^\top (D_t^\top P_t D_t + R_t + \rho \bar{V}_t^{-1})^{-1}(B_t^\top P_t + D_t^\top P_t C_t + S_t) = 0; \quad (2.5) \\
P_T = G,
\end{cases}
$$

*and $D^\top P^\star D + R + \rho \bar{V}^{-1} \succeq \widetilde{\delta} I_k$ for some $\widetilde{\delta} > 0$.*

*Remark* 2.2. Condition (H.2) is called the strongly regular solvability of (2.5) in [27] and ensures that (2.1)-(2.4) admits an optimal feedback control. Note that it suffices to assume the existence of a strongly regular solution, as the uniqueness of a strongly regular solution to (2.5) follows directly from Gronwall's inequality (see [27] and also [34, Proposition 7.1, p. 319]). One can easily show that (H.2) holds if the unregularised (2.5) is strongly regular solvable, i.e., (2.5) with $\rho = 0$ admits a solution $P^{\star,0} \in C([0,T];\mathbb{S}^n)$ and $D^\top P^{\star,0}D + R \succeq \widetilde{\delta} I_k$. This is due to the fact that $P^\star \succeq P^{\star,0}$ (see [27, Theorem 5.3]), and hence $D^\top P^\star D + R + \rho \bar{V}^{-1} \succeq D^\top P^{\star,0}D + R$ by (H.1(2)).

Moreover, by virtue of the regularisation term $\rho \bar{V}^{-1}$, (H.2) may hold even when the unregularised LQ problem (with $\rho = 0$) is not closed-loop solvable. This indicates that the entropy term $\rho \mathcal{H}(\cdot \| \overline{\mathfrak{m}}_t)$ indeed regularises the cost landscape. Such a regularisation effect may not hold if the reference measure $\overline{\mathfrak{m}}_t$, $t \in [0,T]$, is chosen as the Lebesgue measure $\mathcal{L}_k$ on $\mathbb{R}^k$. In fact, as shown in [31, 32, 7], if $\overline{\mathfrak{m}}_t = \mathcal{L}_k$ for all $t \in [0,T]$, then the closed-loop solvability of the regularised problem is equivalent to that of the unregularised problem, and the entropy term will not modify the cost landscape over policies.

Under (H.1) and (H.2), standard verification arguments (see, e.g., [34]) show that the optimal control $\mathfrak{m}^\star \in \mathcal{A}$ of (2.4) is of the form $\mathfrak{m}_t^\star = \nu_t^\star(X_t^{\mathfrak{m}^\star})$, where $\nu^\star : [0,T] \times \mathbb{R}^d \to \mathcal{P}_2(\mathbb{R}^k)$ satisfies for all $(t,x) \in [0,T] \times \mathbb{R}$, $\nu_t^\star(x) = \mathcal{N}(K_t^\star x, V_t^\star)$ and

$$
\begin{aligned}
K_t^\star &= -(D_t^\top P_t^\star D_t + R_t + \rho \bar{V}_t^{-1})^{-1}(B_t^\top P_t^\star + D_t^\top P_t^\star C_t + S_t), \\
V_t^\star &= \rho(D_t^\top P_t^\star D_t + R_t + \rho \bar{V}_t^{-1})^{-1}.
\end{aligned} \quad (2.6)
$$

By (H.1) and (H.2), $K^\star \in L^2(0,T;\mathbb{R}^{k \times d})$, $V^\star \in L^\infty(0,T;\mathbb{S}_+^k)$ and $V^\star \succeq \varepsilon I_k$ for some $\varepsilon > 0$. Note that the optimality of $\mathfrak{m}^\star$ in $\mathcal{A}$ implies that the policy $\nu^\star$ is optimal among all Markovian feedback controls $\nu : [0,T] \times \mathbb{R}^d \to \mathcal{P}_2(\mathbb{R}^k)$ for which the resulting open-loop control $\mathfrak{m}_. = \nu_.(X_.^\nu)$ is square integrable. Here, $X^\nu$ denotes the state dynamics controlled by $\nu$, as defined in (2.8).

## 2.2 Optimisation over Gaussian policies and landscape analysis

Motivated by the optimal Gaussian policy $\nu^\star$ in (2.6), this section reformulates (2.1)-(2.4) as an equivalent minimisation problem over Gaussian policies, and presents key properties of the optimisation landscape $\mathcal{C} : \Theta \to \mathbb{R}$. The proofs of these properties will be given in Section 3.1.

**Policy optimisation.** Let $\Theta$ be the following parameter space

$$\Theta := \left\{ \theta = (K, V) \in \mathcal{B}(0, T; \mathbb{R}^{k \times d} \times \mathbb{S}_+^k) \, \middle| \, \|K\|_{L^2} < \infty, \ \varepsilon I_k \preceq V \preceq \tfrac{1}{\varepsilon} I_k \text{ for some } \varepsilon > 0 \right\},$$

and $\mathcal{V}$ be the space of Gaussian policies parameterised by $\Theta$:

$$\mathcal{V} := \left\{ \nu^\theta : [0, T] \times \mathbb{R}^d \ni (t, x) \mapsto \mathcal{N}(K_t x, V_t) \in \mathcal{P}(\mathbb{R}^k) \, \middle| \, \theta = (K, V) \in \Theta \right\}.^1 \tag{2.7}$$

We shall identify $\nu^\theta \in \mathcal{V}$ with its parameter $\theta = (K, V) \in \Theta$. For each $\nu^\theta \in \mathcal{V}$, consider the associated controlled dynamics (cf. (2.1)):

$$\mathrm{d}X_t = \Phi_t(X_t, \nu_t^\theta(X_t)) \, \mathrm{d}t + \Gamma_t(X_t, \nu_t^\theta(X_t)) \, \mathrm{d}W_t, \quad t \in [0, T]; \quad X_0 = \xi_0, \tag{2.8}$$

with $\Phi$ and $\Gamma$ defined in (2.2), and let $X^\theta \in \mathcal{S}^2(0, T; \mathbb{R}^d)$ be the unique solution to (2.8) (see Proposition A.2). Then we consider minimising the following cost functional:

$$\begin{aligned}
\mathcal{C}(\theta) := \mathbb{E}\bigg[ &\int_0^T \int_{\mathbb{R}^k} \left( \frac{1}{2} \left\langle \begin{pmatrix} Q_t & S_t^\top \\ S_t & R_t \end{pmatrix} \begin{pmatrix} X_t^\theta \\ a \end{pmatrix}, \begin{pmatrix} X_t^\theta \\ a \end{pmatrix} \right\rangle \nu_t^\theta(X_t^\theta; \mathrm{d}a) + \rho \mathcal{H}(\nu_t^\theta(X_t^\theta) \| \overline{\mathrm{m}}_t) \right) \mathrm{d}t \\
&+ \frac{1}{2} (X_T^\theta)^\top G X_T^\theta \bigg]
\end{aligned} \tag{2.9}$$

over all $\theta \in \Theta$, or equivalently all $\nu^\theta \in \mathcal{V}$. It is clear that the cost $\mathcal{C}$ is minimised at $\theta^\star = (K^\star, V^\star)$ defined in (2.6), and the minimum value $\inf_{\theta \in \Theta} \mathcal{C}(\theta)$ is the minimum cost of (2.1)-(2.4).

**Optimisation landscape.** To investigate the regularity of the map $\mathcal{C} : \Theta \to \mathbb{R}$, we introduce two important quantities: for each $\theta = (K, V) \in \Theta$, let $P^\theta \in C([0, T]; \mathbb{S}^d)$ be the solution to following (backward) Lyapunov equation:

$$\begin{aligned}
(\tfrac{\mathrm{d}}{\mathrm{d}t}P)_t &+ (A_t + B_t K_t)^\top P_t + P_t^\top (A_t + B_t K_t) + (C_t + D_t K_t)^\top P_t (C_t + D_t K_t) \\
&+ K_t^\top (R_t + \rho \bar{V}_t^{-1}) K_t + S_t^\top K_t + K_t^\top S_t + Q_t = 0, \quad \text{a.e. } t \in [0, T]; \quad P_T = G,
\end{aligned} \tag{2.10}$$

and let $\Sigma^\theta \in C([0, T]; \overline{\mathbb{S}_+^d})$ be the solution to the following Lyapunov equation: for a.e. $t \in [0, T]$,

$$\begin{aligned}
(\tfrac{\mathrm{d}}{\mathrm{d}t}\Sigma)_t &= (A_t + B_t K_t)\Sigma_t + \Sigma_t(A_t + B_t K_t)^\top + (C_t + D_t K_t)\Sigma_t(C_t + D_t K_t)^\top + D_t V_t D_t^\top, \\
\Sigma_0 &= \mathbb{E}[\xi_0 \xi_0^\top].
\end{aligned} \tag{2.11}$$

Under (H.1), $P^\theta$ and $\Sigma^\theta$ are well-defined by standard well-posedness results of linear differential equations. Note that $P^\theta$ depends only on $K$ and is independent of $V$. Moreover, let $X^\theta$ be the state process governed by (2.8), then $\Sigma_t^\theta = \mathbb{E}[X_t^\theta (X_t^\theta)^\top]$ for all $t \in [0, T]$,[2] due to a straightforward application of Itô's formula to $t \to X_t^\theta (X_t^\theta)^\top$ and the definition (2.2) (see also Lemma 3.1).

Based on the notation $P^\theta$ and $\Sigma^\theta$, the following proposition characterises the Gateaux derivatives of $\mathcal{C}$ at each $\theta \in \Theta$. The proof relies on first reformulating the minimisation problem (2.9) into a deterministic control problem for $\Sigma^\theta$, and then applying the Pontryagin optimality principle.

---

[1] As $\rho > 0$, we require the Gaussian policies in $\mathcal{V}$ to have nondegenerate covariances. If $\rho = 0$, one can restrict admissible policies to be $\nu_t^\theta(x) = \mathcal{N}(K_t x, 0)$. Our analysis and results can be naturally extended to this setting.

[2] Given a state variable $X_t$, the second-moment matrix $\Sigma_t = \mathbb{E}[X_t X_t^\top]$ is often referred to as the state covariance matrix in the reinforcement learning literature (see e.g., [6, 12]). We follow this convention throughout this paper.

**Proposition 2.1.** *Suppose (H.1) holds. For each $\theta \in \Theta$, let $P^\theta \in C([0,T]; \mathbb{S}^d)$ satisfy (2.10), and let $\Sigma^\theta \in C([0,T]; \overline{\mathbb{S}^d_+})$ satisfy (2.11). Then for all $\theta, \theta' \in \Theta$,*

$$\frac{\mathrm{d}}{\mathrm{d}\varepsilon} \mathcal{C}(K + \varepsilon K', V)\Big|_{\varepsilon=0} = \int_0^T \langle \mathcal{D}_K(\theta)_t \Sigma^\theta_t, K'_t \rangle \, \mathrm{d}t,$$

$$\frac{\mathrm{d}}{\mathrm{d}\varepsilon} \mathcal{C}(K, V + \varepsilon(V' - V))\Big|_{\varepsilon=0} = \int_0^T \langle \mathcal{D}_V(\theta)_t, V' - V \rangle \, \mathrm{d}t,$$

*where for a.e. $t \in [0,T]$,*

$$\mathcal{D}_K(\theta)_t := B_t^\top P_t^\theta + D_t^\top P_t^\theta (C_t + D_t K_t) + S_t + (R_t + \rho \bar{V}_t^{-1}) K_t, \tag{2.12}$$

$$\mathcal{D}_V(\theta)_t := \frac{1}{2}(D_t^\top P_t^\theta D_t + R_t + \rho(\bar{V}_t^{-1} - V_t^{-1})). \tag{2.13}$$

We then estimate the regularity of $\mathcal{C} : \Theta \to \mathbb{R}$ by using the gradient terms $\mathcal{D}_K(\theta)$ and $\mathcal{D}_V(\theta)$. The following proposition proves that the functional $\mathcal{C}$ satisfies a non-uniform Łojasiewicz condition in $\theta$. As $\mathcal{C}$ is typically nonconvex in $K$ (see Proposition 2.4), such a Łojasiewicz condition is critical for the global convergence of gradient-based algorithms.

**Proposition 2.2.** *Suppose (H.1) and (H.2) hold. Let $\theta^\star \in \Theta$ be defined by (2.6). For each $\theta \in \Theta$, let $P^\theta \in C([0,T]; \mathbb{S}^d)$ satisfy (2.10), let $\Sigma^\theta \in C([0,T]; \overline{\mathbb{S}^d_+})$ satisfy (2.11), and let $\mathcal{D}_K(\theta)$ and $\mathcal{D}_V(\theta)$ be defined by (2.12) and (2.13), respectively. Then for all $\theta \in \Theta$,*

$$\mathcal{C}(\theta) - \mathcal{C}(\theta^\star) \le \int_0^T \left( \frac{1}{2} \langle (D_t^\top P_t^\theta D_t + R_t + \rho \bar{V}_t^{-1})^{-1} \mathcal{D}_K(\theta)_t, \mathcal{D}_K(\theta)_t \Sigma^{\theta^\star}_t \rangle \right.$$
$$\left. + \frac{1}{\rho} \max(\|V_t^\star\|_2^2, \|V_t\|_2^2) |\mathcal{D}_V(\theta)_t|^2 \right) \mathrm{d}t. \tag{2.14}$$

The next proposition proves that for any $\theta, \theta' \in \Theta$, the cost difference $\mathcal{C}(\theta') - \mathcal{C}(\theta)$ can be upper bounded by the first and second order terms in $\theta' - \theta$. Such a property is often referred to as the "almost smoothness" condition in the literature on PG methods (see e.g., [6, 12, 36]).

**Proposition 2.3.** *Suppose (H.1) holds. For each $\theta \in \Theta$, let $P^\theta \in C([0,T]; \mathbb{S}^d)$ satisfy (2.10), let $\Sigma^\theta \in C([0,T]; \overline{\mathbb{S}^d_+})$ satisfy (2.11), and let $\mathcal{D}_K(\theta)$ and $\mathcal{D}_V(\theta)$ be defined by (2.12) and (2.13), respectively. Then for all $\theta, \theta' \in \Theta$,*

$$\mathcal{C}(\theta') - \mathcal{C}(\theta) \le \int_0^T \left( \langle K'_t - K_t, \mathcal{D}_K(\theta)_t \Sigma^{\theta'}_t \rangle + \frac{1}{2} \langle K'_t - K_t, (D_t^\top P_t^\theta D_t + R_t + \rho \bar{V}_t^{-1})(K'_t - K_t) \Sigma^{\theta'}_t \rangle \right.$$
$$\left. + \langle \mathcal{D}_V(\theta)_t, V'_t - V_t \rangle + \frac{\rho}{4} \frac{|V'_t - V_t|^2}{\min(\lambda^2_{\min}(V_t), \lambda^2_{\min}(V'_t))} \right) \mathrm{d}t.$$

Note that the Łojasiewicz condition in Proposition 2.2 and the smoothness condition in Proposition 2.3 are local properties. The estimates therein depend explicitly on $P^\theta$ and $\Sigma^\theta$, which admit no uniform bound over the unbounded parameter set $\Theta$. For PG methods with finite-dimensional parameter spaces, this difficulty is often overcome by first proving the sublevel set $\{\theta \in \Theta \mid \mathcal{C}(\theta) < \beta\}$ is bounded for any $\beta > 0$, and then designing algorithms whose iterates remain in a fixed sublevel set (see e.g., [6, 10, 12]). However, the following example shows that in the setting with continuous-time policies, the cost is typically noncoercive,[3] and hence the above argument cannot be applied. The proof follows from a straightforward computation, and is given in Appendix A.

---

[3]Let $(X, \|\cdot\|)$ be a normed space. A function $f : X \to \mathbb{R}$ is called coercive if $\lim_{\|x\| \to \infty} f(x) = \infty$.

**Proposition 2.4.** *Let* $\mathcal{C} : L^2(0,1;\mathbb{R}) \to \mathbb{R}$ *be such that for all* $K \in L^2(0,1;\mathbb{R})$,

$$\mathcal{C}(K) := \int_0^1 (K_t X_t)^2 \, \mathrm{d}t, \quad \text{with } X_t = 1 + \int_0^t K_s X_s \, \mathrm{d}s, \ t \in [0,1]. \tag{2.15}$$

*Then* $\mathcal{C} : L^2(0,1;\mathbb{R}) \to \mathbb{R}$ *is neither coercive nor quasiconvex. In particular, let* $K^\varepsilon \in L^2(0,1;\mathbb{R})$, $\varepsilon > 0$, *be such that* $K_t^\varepsilon = -(1 + \varepsilon - t)^{-1}$ *for all* $t \in [0,1]$. *Then* $\lim_{\varepsilon \to 0} \|K^\varepsilon\|_{L^1} = \infty$ *and* $\sup_{\varepsilon > 0} \mathcal{C}(K^\varepsilon) = 1$. *Moreover, there exists* $\varepsilon_0 > 0$ *such that for all* $\varepsilon \in (0, \varepsilon_0]$, $\mathcal{C}(0.5K^\varepsilon) > \max\{\mathcal{C}(\mathbf{0}), \mathcal{C}(K^\varepsilon)\}$, *with* $\mathbf{0}$ *being the zero function.*

## 2.3 Policy gradient method and its convergence analysis

This section proposes a geometry-aware PG method for (2.6) that preserves an a-priori bound, and proves its global linear convergence based on the landscape properties in Section 2.2.

**Geometry-aware policy gradient method.** For each initial guess $\theta^0 = (K^0, V^0) \in \Theta$ and stepsize $\tau > 0$, consider $(\theta^n)_{n \in \mathbb{N}} \subset \mathcal{B}(0, T; \mathbb{R}^{k \times d} \times \mathbb{S}^k)$ such that for all $n \in \mathbb{N}_0$,

$$K_t^{n+1} = K_t^n - \tau \mathcal{D}_K(\theta^n)_t, \quad V_t^{n+1} = V_t^n - \tau \mathcal{D}_V^{\mathrm{bw}}(\theta^n)_t, \quad \text{a.e. } t \in [0,T], \tag{2.16}$$

with

$$\mathcal{D}_V^{\mathrm{bw}}(\theta^n)_t = \mathcal{D}_V(\theta^n)_t V_t^n + V_t^n \mathcal{D}_V(\theta^n)_t,\ ^{4} \tag{2.17}$$

where $\mathcal{D}_K(\theta)$ and $\mathcal{D}_V(\theta)$ are defined by (2.12) and (2.13), respectively. Here we update $K$ and $V$ with the same stepsize $\tau$ for the clarity of presentation, but the results can be naturally extended to the setting where different constant stepsizes are adopted to update $K$ and $V$.

Algorithm (2.16) normalises the (Fréchet) derivatives of $\theta^n$ (cf. Proposition 2.1) to incorporate the local geometry of the parameter space. Specifically, it updates $(K^n)_{n \in \mathbb{N}}$ by the steepest descent on the manifold of Gaussian policies endowed with the Fisher information metric (also known as the natural gradient). To see this, for each $n \in \mathbb{N}_0$, consider the following natural gradient update for $K^n$ (see [17]):

$$K_t^{n+1} = K_t^n - \tau \mathcal{I}(\theta^n)_t^{-1} \nabla_K \mathcal{C}(\theta^n)_t,\ ^{5} \quad \text{a.e. } t \in [0,T], \tag{2.18}$$

where $\nabla_K \mathcal{C}(\theta^n) = \mathcal{D}_K(\theta^n) \Sigma^{\theta^n}$ is the derivative in $K^n$, $\mathcal{I}(\theta^n)_t \in \mathbb{R}^{kd \times kd}$ is the Fisher information matrix satisfying for all $i, i' \in \{1, \ldots, k\}$ and $j, j' \in \{1, \ldots, d\}$,

$$(\mathcal{I}(\theta^n)_t)_{ij,i'j'} := \mathbb{E}\left[ \int_{\mathbb{R}^k} \left[ \partial_{(K_t^n)_{ij}} \ln\left( \hat{\nu}_t^{\theta^n}(X_t^{\theta^n}; a) \right) \partial_{(K_t^n)_{i'j'}} \ln\left( \hat{\nu}_t^{\theta^n}(X_t^{\theta^n}; a) \right) \right] \hat{\nu}_t^{\theta^n}(X_t^{\theta^n}; a) \, \mathrm{d}a \right],$$

and $\hat{\nu}_t^{\theta^n}(X_t^{\theta^n}; \cdot)$ is the density of $\mathcal{N}(K_t^n X_t^{\theta^n}, I_k)$. Then by a similar computation as in [6, 11], $\mathcal{I}(\theta^n)_t^{-1} \nabla_K \mathcal{C}(\theta^n)_t = \nabla_K \mathcal{C}(\theta^n)_t (\Sigma_t^{\theta^n})^{-1} = \mathcal{D}_K(\theta^n)_t$.

On the other hand, (2.16) updates $(V^n)_{n \in \mathbb{N}}$ by the steepest descent on the matrix manifold $\mathbb{S}_+^k$ endowed with the Bures-Wasserstein metric [13]. It corresponds to the geometry induced by the 2-Wasserstein metric over the space of centered nondegenerate Gaussian measures. By normalising $\mathcal{D}_V$ according to $V$, the Riemannian gradient $\mathcal{D}_V^{\mathrm{bw}}$ in (2.17) preserves a pointwise upper and lower bound of $(V^n)_{n \in \mathbb{N}}$ without the use of projection (see Remark 2.3).

---

[4]For an arbitrary stepsize $\tau > 0$, $(V^n)_{n \in \mathbb{N}}$ may not be positive definite and hence may not be invertible. In this case, $\mathcal{D}_V$ is defined by replacing $V_t^{-1}$ in (2.13) with the (symmetric) Moore-Penrose inverse of $V_t$. We prove that $(\theta^n)_{n \in \mathbb{N}} \subset \Theta$ for all sufficiently small stepsizes (see Proposition 2.5).

[5]For each $A \in \mathbb{R}^{kd \times kd}$ and $B \in \mathbb{R}^{k \times d}$, indexed by $A_{ij,i'j'}$ and $B_{ij}$ with $i, i' \in \{1, \ldots, k\}$ and $j, j' \in \{1, \ldots, d\}$, we define $AB \in \mathbb{R}^{k \times d}$ with $(AB)_{ij} = \sum_{k,l} A_{ij,kl} B_{kl}$. This is equivalent to reshaping $B$ (with row-major ordering) into a vector, performing the standard matrix-vector multiplication, and reshaping the result into a matrix.

**Convergence analysis.** The key challenge in the convergence analysis of (2.16) is to establish a uniform bound for the corresponding $(P^{\theta^n})_{n\in\mathbb{N}}$ and $(\Sigma^{\theta^n})_{n\in\mathbb{N}}$, as shown in Proposition 2.5. This is achieved by proving a uniform bound of the iterates $(\theta^n)_{n\in\mathbb{N}}$ and quantifying the explicit dependence of $\Sigma^\theta$ on $\theta$. The proof is given in Section 3.2 (Propositions 3.5, 3.6, and 3.8).

**Proposition 2.5.** *Suppose (H.1) and (H.2) hold. For each $\theta \in \Theta$, let $P^\theta \in C([0,T];\mathbb{S}^d)$ satisfy (2.10) and let $\Sigma^\theta \in C([0,T];\overline{\mathbb{S}^d_+})$ satisfy (2.11). Let $\theta^0 \in \Theta$ and $\overline{\lambda}_0 > 0$ such that $\overline{\lambda}_0 I_k \succeq D^\top P^{\theta^0} D + R + \rho\bar{V}^{-1}$. For each $\tau > 0$, let $(\theta^n)_{n\in\mathbb{N}} \subset \mathcal{B}(0,T;\mathbb{R}^{k\times d}\times\mathbb{S}^k)$ be defined in (2.16). Then*

(1) *There exists $\widetilde{C}, \overline{\lambda}_V, \underline{\lambda}_V > 0$ such that for all $\tau \in (0, 1/\overline{\lambda}_0]$, $n \in \mathbb{N}_0$, $\|K^n\|_{L^2} \leq \widetilde{C}$ and $\underline{\lambda}_V I_k \preceq V^n \preceq \overline{\lambda}_V I_k$.*

(2) *For all $\tau \in (0, 2/\overline{\lambda}_0]$, $n \in \mathbb{N}_0$, $P^{\theta^n} \succeq P^{\theta^{n+1}} \succeq P^\star$, with $P^\star \in C([0,T];\mathbb{S}^d)$ in (H.2),*

(3) *Assume further that $\mathbb{E}[\xi_0\xi_0^\top] \succ 0$. Then there exists $\overline{\lambda}_X, \underline{\lambda}_X > 0$ such that for all $\tau \in (0, 1/\overline{\lambda}_0]$ and $n \in \mathbb{N}_0$, $\underline{\lambda}_X I_d \preceq \Sigma^{\theta^n} \preceq \overline{\lambda}_X I_d$.*

*Remark* 2.3 (**Implicit regularisation**). The uniform bounds of $(K^n)_{n\in\mathbb{N}}$ and $(V^n)_{n\in\mathbb{N}}$ are achieved by an *implicit regularisation* feature of the geometry-aware gradient directions $\mathcal{D}_K$ and $\mathcal{D}_V^{\mathrm{bw}}$. Here, "implicit regularisation" means that the iterates preserve certain constraints without an explicit projection step. Note that applying projection to enforce a pointwise lower bound for minimum eigenvalues of $(V^n)_{n\in\mathbb{N}}$ is computationally expensive. It requires performing an eigenvalue decomposition of $V_t^n$ for every time $t \in [0,T]$ and iteration $n \in \mathbb{N}$.

A similar implicit regularisation property holds if $(K^n)_{n\in\mathbb{N}}$ is updated by a preconditioned natural gradient descent: for all $n \in \mathbb{N}_0$,

$$K_t^{n+1} = K_t^n - \tau H_t^n \mathcal{D}_K(\theta^n)_t, \quad \text{with } \tfrac{1}{L}I_k \preceq H^n \preceq LI_k \text{ for some } L > 0 \text{ independent of } n.$$

This includes the Gauss-Newton method with $H^n = \left(D^\top P^n D + R + \rho\bar{V}^{-1}\right)^{-1}$ as a special case (see [6, 10]). However, due to the noncoercivity of $\mathcal{C}$, it is unclear whether an implicit regularisation holds for an arbitrary descent direction of $\mathcal{C}$ in $K$ (e.g., the vanilla gradient direction $\nabla_K\mathcal{C}(\theta) = \mathcal{D}_K(\theta)\Sigma^\theta$), in contrast with PG methods for discrete-time problems [6, 11]; see the discussion above Proposition 2.4.

It is noteworthy that an implicit regularisation feature of natural policy gradient algorithms was observed in [36]. In their study, an agent optimises over stationary linear policies to stablise a linear system with additive noise over an infinite horizon while adhering to robustness constraints on the sup-norm of the input-output transfer matrix. They show that a natural policy gradient algorithm naturally preserves the transfer matrix's sup-norm throughout the iterations, eliminating the need for explicit projection.

The challenges faced in the current setting differ from those in [36]. Firstly, as Proposition 2.4 shows, the cost of a finite-horizon continuous-time LQC problem is already noncoercive without any robustness constraints. This is primarily because a policy can have an infinite number of changes in values, occurring at arbitrary time points. Such a feature is not present in the scenarios studied in [36], where stationary policies are considered. Secondly, instead of optimising a deterministic policy, we optimise both the mean and covariance of a Gaussian policy, for which we derive natural gradient updates with respect to different geometries. Our result implies that Wasserstein gradient descent of negative entropy preserves a-priori bounds on the variance of Gaussian measures, which is novel and of independent interest. Lastly, the possible degeneracy of the system noise and the cost coefficients (Remark 2.1) necessitates a more precise quantification of the desired implicit regularisation within appropriate function spaces.

Proposition 2.5 implies that the functional $\mathcal{C}$ satisfies uniform Łojasiewicz and smoothness conditions *along the iterates* $(\theta^n)_{n \in \mathbb{N}}$. Based on this local regularity, the following theorem establishes the global linear convergence of (2.16) for all sufficiently small stepsizes $\tau$. The proof is given in Section 3.3.

**Theorem 2.6.** *Suppose (H.1) and (H.2) hold, and $\mathbb{E}[\xi_0 \xi_0^\top] \succ 0$. Let $\theta^0 \in \Theta$, and for each $\tau > 0$, let $(\theta^n)_{n \in \mathbb{N}} \subset \mathcal{B}(0, T; \mathbb{R}^{k \times d} \times \mathbb{S}^k)$ be defined in (2.16). Then there exists $\tau_0, C_1, C_2 > 0$ such that for all $\tau \in (0, \tau_0]$ and $n \in \mathbb{N}_0$,*

*(1) $\mathcal{C}(\theta^{n+1}) \le \mathcal{C}(\theta^n)$ and $\mathcal{C}(\theta^{n+1}) - \mathcal{C}(\theta^\star) \le (1 - \tau C_1)(\mathcal{C}(\theta^n) - \mathcal{C}(\theta^\star))$, with $\theta^\star$ defined in (2.6),*

*(2) $\|K^n - K^\star\|_{L^2}^2 + \|V^n - V^\star\|_{L^2}^2 \le C_2 (1 - \tau C_1)^n$.*

The precise expressions of the constants $\tau_0$, $C_1$ and $C_2$ can be found in the proof of the statement. These constants depend on the regularisation weight $\rho$ in (2.4), the constant $\widetilde{\delta}$ in (H.2), the initial guess $\theta^0$, and the a-priori bounds $\underline{\lambda}_X, \overline{\lambda}_X, \underline{\lambda}_V, \overline{\lambda}_V$ in Proposition 2.5. Achieving more precise dependencies in terms of model parameters is challenging. It would entail deriving precise a-priori bounds of solutions to (2.5) and (2.10) in terms of the coefficients given in (H.1(1)). This remains an open problem, particularly when the diffusion coefficient is controlled ($D \ne 0$) and when cost coefficients $Q$, $R$, and $G$ are not positive definite.

## 2.4 Mesh-independent linear convergence with discrete-time policies

By leveraging Theorem 2.6, this section proposes PG methods that take actions at discrete time points and achieve a robust convergence behaviour across different mesh sizes. Our analysis shows that a proper scaling of the discrete-time gradients in terms of mesh size is critical for a robust performance of the algorithm.

More precisely, let $\mathscr{P}_{[0,T]}$ be the collection of all partitions of $[0, T]$. For each $\pi = \{0 = t_0 < \cdots < t_N = T\} \in \mathscr{P}_{[0,T]}$, let $|\pi| = \max_{i=0,\ldots,N-1}(t_{i+1} - t_i)$ be the mesh size of $\pi$, and let $\Theta^\pi \subset \Theta$ be the set of piecewise constant policies on $\pi$:

$$\Theta^\pi = \{\theta \in \Theta \mid \theta_t = \theta_{t_i}, \text{ a.e. } t \in [t_i, t_{i+1}) \text{ and all } i \in \{0, \ldots, N-1\}\}. \tag{2.19}$$

Then define the minimum cost $\mathcal{C}$ over $\Theta^\pi$:

$$\mathcal{C}_\pi^\star = \inf_{\theta \in \Theta^\pi} \mathcal{C}(\theta). \tag{2.20}$$

Note that by $\Theta^\pi \subset \Theta$, $\mathcal{C}_\pi^\star \ge \inf_{\theta \in \Theta} \mathcal{C}(\theta) = \mathcal{C}(\theta^\star) > -\infty$.

We now introduce a family of gradient descent schemes for (2.20). Let $\theta^{\pi,0} \in \Theta^\pi$ be an initial guess and $\tau > 0$ be a stepsize. Consider the following sequence $(\theta^{\pi,n})_{n \in \mathbb{N}_0} \subset \Theta^\pi$ (cf. (2.16)) such that for all $n \in \mathbb{N}_0$,

$$K_t^{\pi,n+1} = K_t^{\pi,n} - \tau \mathcal{D}_K^\pi(\theta^{\pi,n})_t, \quad V_t^{\pi,n+1} = V_t^{\pi,n} - \tau \mathcal{D}_V^\pi(\theta^{\pi,n})_t, \quad \text{a.e. } t \in [0, T], \tag{2.21}$$

where $(\mathcal{D}_K^\pi, \mathcal{D}_V^\pi) : \Theta^\pi \to \Theta^\pi$ approximates the gradient operators $(\mathcal{D}_K, \mathcal{D}_V^{\text{bw}})$ in (2.16) as $|\pi| \to 0$; see (H.3) for the precise condition.

The convergence behaviour of (2.21) is measured by the number of required iterations $N^\pi(\varepsilon)$ to achieve a certain accuracy $\varepsilon > 0$: let $(\theta^{\pi,n})_{n \in \mathbb{N}_0}$ be generated by (2.21) (with some $\theta^{\pi,0} \in \Theta^\pi$ and $\tau > 0$), and for each $\varepsilon > 0$, define

$$N^\pi(\varepsilon) := \min \left\{ n \in \mathbb{N}_0 \,\Big|\, \mathcal{C}(\theta^{\pi,n}) - \inf_{\theta \in \Theta^\pi} \mathcal{C}(\theta) < \varepsilon \right\} \in \mathbb{N}_0 \cup \{\infty\}. \tag{2.22}$$

Note that $N^\pi$ is defined for a fixed mesh $\pi$, and hence the residual is defined using the minimum cost $\mathcal{C}_\pi^\star$ over piecewise constant policies $\Theta^\pi$. Similarly, let $(\theta^n)_{n\in\mathbb{N}_0}$ be a sequence generated by (2.16) (with some $\theta^0 \in \Theta$ and $\tau > 0$), and for each $\varepsilon > 0$, define

$$N(\varepsilon) := \min\left\{n \in \mathbb{N}_0 \,\Big|\, \mathcal{C}(\theta^n) - \inf_{\theta\in\Theta}\mathcal{C}(\theta) < \varepsilon\right\} \in \mathbb{N}_0 \cup \{\infty\}. \tag{2.23}$$

The main result of this section shows that if the gradient operators $(\mathcal{D}_K^\pi, \mathcal{D}_V^\pi)_\pi$ in (2.21) satisfy the consistency condition (H.3), then for all sufficiently fine grids $\pi$, $N^\pi(\varepsilon)$ is essentially equal to $N(\varepsilon)$.

**H.3.** *For every* $\theta \in L^2(0,T;\mathbb{R}^{k\times d}) \times C([0,T];\mathbb{S}_+^k)$, *every sequence* $(\pi_m)_{m\in\mathbb{N}} \subset \mathscr{P}_{[0,T]}$ *such that* $\lim_{m\to\infty} |\pi_m| = 0$, *and every* $(\theta^m)_{m\in\mathbb{N}} \subset \Theta$ *such that* $\theta^m \in \Theta^{\pi_m}$ *for all* $m \in \mathbb{N}$ *and* $\lim_{m\to\infty} \|\theta^m - \theta\|_{L^2\times L^\infty} = 0$, *we have*

$$\lim_{m\to\infty} \|\mathcal{D}_K^{\pi_m}(\theta^m) - \mathcal{D}_K(\theta)\|_{L^2} = 0, \quad \text{and} \quad \lim_{m\to\infty} \|\mathcal{D}_V^{\pi_m}(\theta^m) - \mathcal{D}_V^{\text{bw}}(\theta)\|_{L^\infty} = 0.$$

**Theorem 2.7.** *Suppose (H.1), (H.2) and (H.3) hold,* $\mathbb{E}[\xi_0\xi_0^\top] \succ 0$, $D \in C([0,T];\mathbb{R}^{d\times k})$, $R \in C([0,T];\mathbb{S}^k)$ *and* $\bar{V} \in C([0,T];\mathbb{S}_+^k)$. *Let* $\theta^0 \in L^2(0,T;\mathbb{R}^{k\times d}) \times C([0,T];\mathbb{S}_+^k)$, *let* $(\pi_m)_{m\in\mathbb{N}} \subset \mathscr{P}_{[0,T]}$ *be such that* $\lim_{m\to\infty} |\pi_m| = 0$ *and let* $(\theta^{\pi_m,0})_{m\in\mathbb{N}} \subset \Theta$ *be such that* $\theta^{\pi_m,0} \in \Theta^{\pi_m}$ *for all* $m \in \mathbb{N}$ *and* $\lim_{m\to\infty} \|\theta^{\pi_m,0} - \theta^0\|_{L^2\times L^\infty} = 0$. *Then there exists* $\tau_0 > 0$ *such that for all* $\tau \in (0,\tau_0)$ *and* $\varepsilon > 0$, *there exists* $\overline{m} \in \mathbb{N}$ *such that*

$$N(\varepsilon) - 1 \leq N^{\pi_m}(\varepsilon) \leq N(\varepsilon), \quad \forall m \in \mathbb{N} \cap [\overline{m}, \infty). \tag{2.24}$$

The proof of Theorem 2.7 is given in Section 3.4.

Theorem 2.7 indicates that (2.21) achieves linear convergence uniformly across different timescales. Indeed, by Theorem 2.6, there exists $\tau_0, C_1 > 0$ such that for all $\tau \in (0,\tau_0]$ and $n \in \mathbb{N}_0$, $\mathcal{C}(\theta^{n+1}) - \mathcal{C}(\theta^\star) \leq (1 - \tau C_1)^n(\mathcal{C}(\theta^0) - \mathcal{C}(\theta^\star))$. This implies that $N(\varepsilon) \leq \frac{\ln\left(\frac{\varepsilon}{\mathcal{C}(\theta^0)-\mathcal{C}(\theta^\star)}\right)}{\ln(1-\tau C_1)}$ for all $\varepsilon > 0$. By the identity that $\lim_{x\to 0} \frac{\ln(1+x)}{x} = 1$, $N^{\pi_m}(\varepsilon) \approx \frac{1}{C_1\tau}\ln\left(\frac{\mathcal{C}(\theta^0)-\mathcal{C}(\theta^\star)}{\varepsilon}\right)$ for all $m \geq \overline{m}$ and sufficiently small $\tau$ and $\varepsilon$.

To design a concrete gradient methods satisfying (H.3), for each $\pi = \{0 = t_0 < \cdots < t_N = T\} \in \mathscr{P}_{[0,T]}$, we identify $\Theta^\pi$ with $(\mathbb{R}^{k\times d} \times \mathbb{S}_+^k)^N$ by the natural parameterisation:

$$(\mathbb{R}^{k\times d} \times \mathbb{S}_+^k)^N \ni (K_i, V_i)_{i=1}^{N-1} \mapsto \left(\sum_{i=0}^{N-1} K_i\mathbb{1}_{[t_i,t_{i+1})}(t), \sum_{i=0}^{N-1} V_i\mathbb{1}_{[t_i,t_{i+1})}(t)\right)_{t\in[0,T]} \in \Theta^\pi, \tag{2.25}$$

and by abuse of notation, write $\mathcal{C} : (\mathbb{R}^{k\times d} \times \mathbb{S}_+^k)^N \to \mathbb{R}$ as the cost of a Gaussian policy induced by the parameterisation (2.7) and (2.25). Then for each $\theta^{\pi,0} \in \Theta^\pi$ and $\tau > 0$, consider the following sequence $(\theta^{\pi,n})_{n\in\mathbb{N}_0} \subset \Theta^\pi$ such that for all $n \in \mathbb{N}_0$ and $i \in \{0,\ldots,N-1\}$, $\theta_t^{\pi,n+1} = (K_i^{\pi,n+1}, V_i^{\pi,n+1})$ for all $t \in [t_i, t_{i+1})$, with

$$\begin{aligned}
K_i^{\pi,n+1} &= K_i^{\pi,n} - \frac{\tau}{\Delta_i}\nabla_{K_i}\mathcal{C}(\theta^{\pi,n})\left(\Sigma_{t_i}^{\theta^{\pi,n}}\right)^{-1}, \\
V_i^{\pi,n+1} &= V_i^{\pi,n} - \frac{\tau}{\Delta_i}\left(V_i^{\pi,n}\nabla_{V_i}\mathcal{C}(\theta^{\pi,n}) + \nabla_{V_i}\mathcal{C}(\theta^{\pi,n})V_i^{\pi,n}\right),
\end{aligned} \tag{2.26}$$

where $\Delta_i = t_{i+1} - t_i$, $\Sigma_{t_i}^{\theta^{\pi,n}} = \mathbb{E}[X_{t_i}^{\theta^{\pi,n}}(X_{t_i}^{\theta^{\pi,n}})^\top]$, and $\nabla_{K_i}\mathcal{C}$ (resp. $\nabla_{V_i}\mathcal{C}$) is the partial derivative of $\mathcal{C}$ with respect to the matrix $K_i$ (resp. $V_i$). The practical implementation of the algorithm is further discussed at the end of this section.

The following corollary shows that (2.26) satisfies (H.3), whose proof is given in Section 3.4.

**Corollary 2.8.** *Suppose (H.1) and (H.2) hold,* $\mathbb{E}[\xi_0\xi_0^\top] \succ 0$, $D \in C([0,T]; \mathbb{R}^{d\times k})$, $R \in C([0,T]; \mathbb{S}^k)$ *and* $\bar{V} \in C([0,T]; \mathbb{S}_+^k)$. *Then Theorem 2.7 holds for* (2.26).

*Remark* 2.4 (**Scaling hyper-parameters with timescales**). It is critical to scale the stepsize $\tau$ in (2.26) with respect to $\Delta_i$ for the robustness of (2.26) for all small mesh sizes. Indeed, standard discrete-time natural PG methods correspond to setting $\Delta_i = 1$ in (2.26) for all grids. For sufficiently fine grids, this is equivalent to adopting a vanishing stepsize $\tau\Delta_i$ in (2.16), as $\nabla_{K_i}\mathcal{C}(\theta)\left(\Sigma_{t_i}^\theta\right)^{-1} \approx \mathcal{D}_K(\theta)_{t_i}\Delta_i$ and $\nabla_{V_i}\mathcal{C}(\theta^{\pi,n}) \approx \mathcal{D}_V(\theta)_{t_i}\Delta_i$ (see Proposition 2.1). This explains the degraded performance of conventional discrete-time PG methods for small mesh sizes. In contrast, by normalising the stepsize with $\Delta_i$, (2.26) admits a continuous-time limit (2.16) as the time stepsize $|\pi|$ vanishes., and achieves mesh-independent convergence; see Section 4 for more details.

*Remark* 2.5 (**Extensions to discrete-time models**). Corollary 2.8 can be extended to incorporate time discretization of the underlying system. Here we provide a heuristic explanation of such an extension. Consider a sequence of time grids $(\pi_m)_{m\in\mathbb{N}}$ with $\lim_{m\to\infty}|\pi_m| = 0$. For each $m \in \mathbb{N}$, let $X^m$ be the discrete-time state dynamics resulting from the Euler–Maruyama discretization of (2.8) on the grid $\pi_m$, and let $\mathcal{C}^{\pi_m} : \Theta^{\pi_m} \to \mathbb{R}$ be the associated cost functional (2.9). Introduce an analogue of (2.26), where $\nabla_{K_i}\mathcal{C}(\theta^{\pi,n})$ and $\nabla_{V_i}\mathcal{C}(\theta^{\pi,n})$ are replaced by $\nabla_{K_i}\mathcal{C}^{\pi_m}(\theta^{\pi,n})$ and $\nabla_{V_i}\mathcal{C}^{\pi_m}(\theta^{\pi,n})$, respectively, and $\Sigma_{t_i}^{\theta^{\pi,n}}$ is replaced by the covariance of the discrete-time state $X^m$ controlled by $\theta^{\pi_m,n}$. If the coefficients in (H.1(1)) are sufficiently regular in time, one can show that these discrete-time gradients converge to the continuous-time gradients in (2.16) as $m \to \infty$, due to the weak convergence of the Euler–Maruyama scheme. This would verify Condition (H.3), which along with Theorem 2.7 implies that these discrete-time algorithms achieve mesh-independent linear convergence uniformly in $m$.

Similar analyses can be carried out for various time discretizations of the state system. Making these arguments precise for general time discretizations would require accurately quantifying the regularity conditions of the coefficients for the weak convergence of the discretization, and is left for future work.

We end this section by describing a possible practical implementation of the algorithm (2.26) which allows for unknown coefficients in (2.8) and (2.9). Recall that, as shown in [29], for a given Gaussian policy $\nu^\theta$, the aggregated dynamics (2.8) and the associated cost (2.9) can be approximated by interacting with the linear dynamics (2.3) with random actions. More precisely, let $\tilde{\pi} = \{0 = \tilde{t}_0 < \cdots < \tilde{t}_M = T\}$ be a time mesh at which random actions are sampled. Consider $X^{\theta,\zeta}$ governed by the following dynamics:

$$\mathrm{d}X_t = (A_t X_t + B_t\phi_t^\theta(X_t))\,\mathrm{d}t + (C_t X_t + D_t\phi_t^\theta(X_t))\,\mathrm{d}W_t, \quad t \in [0,T]; \quad X_0 = \xi_0, \qquad (2.27)$$

where

$$\phi_t^\theta(x) = K_t x + V_t^{1/2}\vartheta_t, \quad \text{with } \vartheta_t := \sum_{i=0}^{M-1} \zeta_i \mathbb{1}_{[\tilde{t}_i,\tilde{t}_{i+1})}(t),$$

and $(\zeta_i)_{i=0}^{M-1}$ are mutually independent standard normal vectors that are independent of $\xi_0$ and $W$. The associated cost with fixed realisations of $\vartheta$, $\xi_0$ and $W$ is defined as:

$$\hat{\mathcal{C}}(\theta) := \int_0^T \left(\frac{1}{2}\left\langle \begin{pmatrix} Q_t & S_t^\top \\ S_t & R_t \end{pmatrix}\begin{pmatrix} X_t^{\theta,\zeta} \\ \phi_t(X_t^{\theta,\zeta}) \end{pmatrix}, \begin{pmatrix} X_t^{\theta,\zeta} \\ \phi_t(X_t^{\theta,\zeta}) \end{pmatrix}\right\rangle + \rho\mathcal{H}(\nu_t^\theta(X_t^{\theta,\zeta})\|\overline{\mathfrak{m}}_t)\right)\mathrm{d}t + \frac{1}{2}(X_T^{\theta,\zeta})^\top G X_T^{\theta,\zeta},$$

$$(2.28)$$

13

where by $\overline{\mathfrak{m}}_t = \mathcal{N}(0, \bar{V}_t)$ (see Lemma 3.1),

$$\mathcal{H}(\nu_t^\theta(X_t^{\theta,\zeta}) \| \overline{\mathfrak{m}}_t) = \frac{1}{2}\left(\operatorname{tr}\left(K_t^\top \bar{V}_t^{-1} K_t X_t^{\theta,\zeta}(X_t^{\theta,\zeta})^\top + \bar{V}_t^{-1} V_t\right) - k + \ln\left(\frac{\det(\bar{V}_t)}{\det(V_t)}\right)\right).$$

In (2.27), the linear dynamics (2.3) is controlled by sampling actions from $\nu^{\theta^\pi}$ using the injected noises $(\zeta_i)_{i=0}^{M-1}$, and (2.28) is the quadratic cost induced by these random actions. Then, by arguments similar to those in [29, Theorem 2.2], $|\mathbb{E}[X_t^{\theta,\zeta}(X_t^{\theta,\zeta})^\top] - \Sigma_t^\theta| \leq C|\tilde{\pi}|$ for all $t \in [0,T]$, and $|\mathbb{E}[\hat{\mathcal{C}}(\theta)] - \mathcal{C}(\theta)| \leq C|\tilde{\pi}|$, with a constant $C$ independent of $\tilde{\pi}$. One can also establish an error bound of the order $\mathcal{O}(\sqrt{|\tilde{\pi}|})$ in the high-probability sense with respect to the noise process $\vartheta$.

The above observation suggests that, at each iteration of (2.26), the gradients $\nabla_{K_i}\mathcal{C}(\theta^{\pi,n})$, $\nabla_{V_i}\mathcal{C}(\theta^{\pi,n})$ and the state covariance $\Sigma_{t_i}^{\theta^{\pi,n}}$ at all grid points of $\pi$ can be estimated using Monte Carlo methods without relying on knowledge of the coefficients in (2.8) and (2.9). By choosing a sufficiently fine randomisation grid $\tilde{\pi}$, the covariance $\Sigma_{t_i}^{\theta^{\pi,n}}$ can be estimated by the empirical covariance of $X^{\theta^{\pi,n},\zeta}$ corresponding to different realisations of $\vartheta$, $W$ and $\xi_0$. The gradients of the cost $\mathcal{C}(\theta^{\pi,n})$ can be approximated by suitable zero-order optimisation methods based on trajectories of the cost (2.28) (see e.g., [6, 12, 2]). It would be interesting to quantify the precise sample efficiency of such a model-free implementation of (2.26). This would entail estimating the approximation errors of $\nabla_{K_i}\mathcal{C}(\theta^{\pi,n})$, $\nabla_{V_i}\mathcal{C}(\theta^{\pi,n})$ and $\Sigma_{t_i}^{\theta^{\pi,n}}$ in terms of the sample frequency $|\tilde{\pi}|^{-1}$ and the sample size, and quantifying the precise error propagation through the gradient descent iteration. We leave a rigorous analysis of such a model-free algorithm for future research.

## 3 Proofs

### 3.1 Analysis of optimisation landscape

This section proves the regularity of the cost functional $\mathcal{C}$ in (2.9) given in Section 2.2.

We start by proving several technical lemmas. The following lemma expresses the coefficients of (2.8) and the cost function (2.9) in terms of $\theta = (K, V)$. The proof follows from a straightforward computation and hence is omitted.

**Lemma 3.1.** *Suppose (H.1) holds. Then for all $\nu^\theta \in \mathcal{V}$ and $(t,x) \in [0,T] \times \mathbb{R}^d$,*

$$\Phi_t(x, \nu_t^\theta(x)) = (A_t + B_t K_t)x,$$

$$\Gamma_t(x, \nu_t^\theta(x)) = \left((C_t + D_t K_t)xx^\top(C_t + D_t K_t)^\top + D_t V_t D_t^\top\right)^{\frac{1}{2}},$$

$$\int_{\mathbb{R}^k} \left\langle \begin{pmatrix} Q_t & S_t^\top \\ S_t & R_t \end{pmatrix}\begin{pmatrix} x \\ a \end{pmatrix}, \begin{pmatrix} x \\ a \end{pmatrix} \right\rangle \nu_t^\theta(x; \mathrm{d}a) = \left\langle \begin{pmatrix} Q_t & S_t^\top \\ S_t & R_t \end{pmatrix}\begin{pmatrix} x \\ K_t x \end{pmatrix}, \begin{pmatrix} x \\ K_t x \end{pmatrix} \right\rangle + \operatorname{tr}(R_t V_t),$$

$$\mathcal{H}(\nu_t^\theta(x) \| \overline{\mathfrak{m}}_t) = \frac{1}{2}\left((K_t x)^\top \bar{V}_t^{-1} K_t x + \operatorname{tr}(\bar{V}_t^{-1} V_t) - k + \ln\left(\frac{\det(\bar{V}_t)}{\det(V_t)}\right)\right).$$

The next lemma represents the cost $\mathcal{C}(\theta)$ in terms of $P^\theta$ defined in (2.10).

**Lemma 3.2.** *Suppose (H.1) holds. For each $\theta \in \Theta$, let $P^\theta \in C([0,T]; \mathbb{S}^d)$ satisfy (2.10), let $\varphi^\theta \in C([0,T]; \mathbb{R}^d)$ satisfy for a.e. $t \in [0,T]$,*

$$(\tfrac{\mathrm{d}}{\mathrm{d}t}\varphi)_t + \tfrac{1}{2}\operatorname{tr}\left((D_t^\top P_t^\theta D_t + R_t + \rho\bar{V}_t^{-1})V_t\right) + \tfrac{\rho}{2}\left(-k + \ln\left(\tfrac{\det(\bar{V}_t)}{\det(V_t)}\right)\right) = 0; \quad \varphi_T = 0,$$

14

and let $u^\theta : [0,T] \times \mathbb{R}^d \to \mathbb{R}$ be such that $u_t^\theta(x) = \frac{1}{2}x^\top P_t^\theta x + \varphi_t^\theta$ for all $(t,x) \in [0,T] \times \mathbb{R}^d$. Then for all $\theta \in \Theta$ and $x \in \mathbb{R}^d$,

$$
(\tfrac{\mathrm{d}}{\mathrm{d}t}u^\theta)_t + \frac{1}{2}\mathrm{tr}\left(\Gamma_t(x,\nu_t^\theta(x))\Gamma_t(x,\nu_t^\theta(x))^\top (\nabla_x^2 u^\theta)_t(x)\right) + \Phi_t(x,\nu_t^\theta(x))^\top (\nabla_x u^\theta)_t(x)
$$

$$
+ \frac{1}{2}\left(x^\top Q_t x + x^\top S_t^\top K_t x + (K_t x)^\top S_t x + (K_t x)^\top R_t K_t x + \mathrm{tr}(R_t V_t)\right) \tag{3.1}
$$

$$
+ \frac{\rho}{2}\left((K_t x)^\top \bar{V}_t^{-1} K_t x + \mathrm{tr}(\bar{V}_t^{-1} V_t) - k + \ln\left(\frac{\det(\bar{V}_t)}{\det(V_t)}\right)\right) = 0, \quad \text{a.e. } t \in [0,T],
$$

and $u_T^\theta(x) = \frac{1}{2}x^\top G x$, where $\nabla_x u^\theta$ and $\nabla_x^2 u^\theta$ are the gradient and Hessian of $u^\theta$ in $x$, respectively. Moreover, it holds that $\mathcal{C}(\theta) = \mathbb{E}[u_0^\theta(\xi_0)]$.

*Proof.* Let $X^\theta \in \mathcal{S}^2(0,T;\mathbb{R}^d)$ be the solution to (2.8). For notational simplicity, we omit $\theta$ in the superscripts of all variables.

By Lemma 3.1 and the definition of $u$, for all $(t,x) \in [0,T]$,

$$
\Phi_t(x,\nu_t^\theta(x))^\top (\nabla_x u^\theta)_t(x) = \tfrac{1}{2}x^\top\left((A_t + B_t K_t)^\top P_t + P_t(A_t + B_t K_t)\right)x,
$$

$$
\mathrm{tr}\left(\Gamma_t(x,\nu_t^\theta(x))\Gamma_t(x,\nu_t^\theta(x))^\top (\nabla_x^2 u^\theta)_t(x)\right) = \mathrm{tr}\left(\left((C_t + D_t K_t)xx^\top (C_t + D_t K_t)^\top + D V_t D^\top\right)P_t\right).
$$

Then one can easily see from the definitions of $P$ and $\varphi$ that $u$ satisfies (3.1) for a.e. $t \in [0,T]$ and all $x \in \mathbb{R}^d$, and $u_T(x) = \frac{1}{2}x^\top G x$.

Now applying Itô's formula to $t \mapsto u_t(X_t)$ implies that

$$
u_T(X_T) = u_0(X_0) + \int_0^T \left((\tfrac{\mathrm{d}}{\mathrm{d}t}u)_t(X_t) + \frac{1}{2}\mathrm{tr}\left(\Gamma_t(X_t,\nu_t(X_t))\Gamma_t(X_t,\nu_t(X_t))^\top (\nabla_x^2 u)_t(X_t)\right)\right.
$$

$$
\left. + \Phi_t(X_t,\nu_t(X_t))^\top (\nabla_x u)_t(X_t)\right)\mathrm{d}t + \int_0^T (\nabla_x u)_t(X_t)^\top \Gamma_t(X_t,\nu_t(X_t))\,\mathrm{d}W_t. \tag{3.2}
$$

By the identity $\nabla_x u_t = P_t x$ and the integrability of $C, D, \theta$ and $X$, $\int_0^\cdot (\nabla_x u)_t(X_t)^\top \Gamma_t(X_t,\nu_t(X_t))\,\mathrm{d}W_t$ is a martingale. Hence taking expectations on both sides of (3.2) and using (3.1) give that

$$
\mathbb{E}[u_0(\xi_0)] = \mathbb{E}\left[\frac{1}{2}X_T^\top G X_T\right] + \mathbb{E}\left[\int_0^T \left\{\frac{1}{2}\left(\left\langle\begin{pmatrix} Q_t & S_t^\top \\ S_t & R_t \end{pmatrix}\begin{pmatrix} X_t \\ K_t X_t \end{pmatrix}, \begin{pmatrix} X_t \\ K_t X_t \end{pmatrix}\right\rangle + \mathrm{tr}(R_t V_t)\right)\right.\right.
$$

$$
\left.\left. + \frac{\rho}{2}\left((K_t X_t)^\top \bar{V}_t^{-1} K_t X_t + \mathrm{tr}(\bar{V}_t^{-1} V_t) - k + \ln\left(\frac{\det(\bar{V}_t)}{\det(V_t)}\right)\right)\right\}\mathrm{d}t\right], \tag{3.3}
$$

which along with Lemma 3.1 leads to the desired identity $\mathcal{C}(\theta) = \mathbb{E}[u_0(\xi_0)]$. □

The following lemma quantifies the difference of value functions for two policies.

**Lemma 3.3.** *Suppose (H.1) holds. For each $\theta \in \Theta$, let $P^\theta \in C([0,T];\mathbb{S}^d)$ satisfy (2.10), and let $\Sigma^\theta \in C([0,T];\overline{\mathbb{S}_+^d})$ satisfy (2.11). Then for all $\theta, \theta' \in \Theta$,*

$$
\mathcal{C}(\theta') - \mathcal{C}(\theta) = \int_0^T \left(\langle K_t' - K_t, \mathcal{D}_K(\theta)_t \Sigma_t^{\theta'}\rangle + \frac{1}{2}\langle K_t' - K_t, (D_t^\top P_t^\theta D_t + R_t + \rho\bar{V}_t^{-1})(K_t' - K_t)\Sigma_t^{\theta'}\rangle\right.
$$

$$
\left. + \ell_t(V_t', P_t^\theta) - \ell_t(V_t, P_t^\theta)\right)\mathrm{d}t,
$$

*where $\mathcal{D}_K(\theta)_t$ is defined by (2.12), and $\ell : [0,T] \times \mathbb{S}_+^k \times \mathbb{R}^{d\times d} \to \mathbb{R}$ is given by*

$$
\ell_t(V,Z) = \frac{1}{2}\left(\langle D_t^\top Z D_t + R_t + \rho\bar{V}_t^{-1}, V\rangle - \rho\ln(\det(V))\right) \quad \forall(t,V,Z) \in [0,T] \times \mathbb{S}_+^k \times \mathbb{R}^{d\times d}. \tag{3.4}
$$

15

*Proof.* Throughout this proof, let $\theta, \theta' \in \Theta$ be given, let $(P, \Sigma) = (P^\theta, \Sigma^\theta)$, $(P', \Sigma') = (P^{\theta'}, \Sigma^{\theta'})$, $u = u^\theta$ and $u' = u^{\theta'}$, where for each $\theta \in \Theta$, $u^\theta : [0, T] \times \mathbb{R}^d \to \mathbb{R}$ is defined as in Lemma 3.2. By (3.1), for all $x \in \mathbb{R}^d$, $(u' - u)_T(x) = 0$,

$$\begin{aligned}
(\tfrac{\mathrm{d}}{\mathrm{d}t}(u' - u))_t &+ \frac{1}{2}\mathrm{tr}\left(\Gamma_t(x, \nu_t^{\theta'}(x))\Gamma_t(x, \nu_t^{\theta'}(x))^\top (\nabla_x^2(u' - u))_t(x)\right) \\
&+ \Phi_t(x, \nu_t^{\theta'}(x))^\top (\nabla_x(u' - u))_t(x) + F_t(x) = 0, \quad \text{a.e. } t \in [0, T],
\end{aligned} \tag{3.5}$$

where $F : [0, T] \times \mathbb{R}^d \to \mathbb{R}$ is given by

$$\begin{aligned}
F_t(x) = &\frac{1}{2}\mathrm{tr}\left(\Gamma_t(x, \nu_t^{\theta'}(x))\Gamma_t(x, \nu_t^{\theta'}(x))^\top (\nabla_x^2 u)_t(x)\right) + \Phi_t(x, \nu_t^{\theta'}(x))^\top (\nabla_x u)_t(x) \\
&- \frac{1}{2}\mathrm{tr}\left(\Gamma_t(x, \nu_t^{\theta}(x))\Gamma_t(x, \nu_t^{\theta}(x))^\top (\nabla_x^2 u)_t(x)\right) - \Phi_t(x, \nu_t^{\theta}(x))^\top (\nabla_x u)_t(x) \\
&+ \frac{1}{2}\left[\left(x^\top Q_t x + x^\top S_t^\top K_t' x + (K_t' x)^\top S_t x + (K_t' x)^\top R_t K_t' x + \mathrm{tr}(R_t V_t')\right)\right. \\
&\left. - \left(x^\top Q_t x + x^\top S_t^\top K_t x + (K_t x)^\top S_t x + (K_t x)^\top R_t K_t x + \mathrm{tr}(R_t V_t)\right)\right] \\
&+ \frac{\rho}{2}\left[\left((K_t' x)^\top \bar{V}_t^{-1} K_t' x + \mathrm{tr}(\bar{V}_t^{-1} V_t') - \ln\left(\det(V_t')\right)\right)\right. \\
&\left. - \left((K_t x)^\top \bar{V}_t^{-1} K_t x + \mathrm{tr}(\bar{V}_t^{-1} V_t) - \ln\left(\det(V_t)\right)\right)\right].
\end{aligned}$$

Applying Itô's formula to $t \mapsto (u' - u)_t(X_t^{\theta'})$ (recall the definition of $u^\theta$ in Lemma 3.2) and using (3.5) yield that

$$\mathbb{E}[(u' - u)_T(X_T^{\theta'})] - \mathbb{E}[(u' - u)_0(X_0^{\theta'})] = \mathbb{E}\left[\int_0^T -F_t(X_t^{\theta'})\,\mathrm{d}t\right],$$

which along with $\mathcal{C}(\theta) = \mathbb{E}[u_0^\theta(\xi_0)]$ (see Lemma 3.2) and $(u' - u)_T = 0$ implies that

$$\mathcal{C}(\theta') - \mathcal{C}(\theta) = \mathbb{E}\left[\int_0^T F_t(X_t^{\theta'})\,\mathrm{d}t\right]. \tag{3.6}$$

We now simplify the expression of $F_t(x)$ for any given $(t, x) \in [0, T] \times \mathbb{R}^d$. To this end, let $\overline{H} : [0, T] \times \mathbb{R}^d \times \mathbb{R}^k \times \mathbb{R}^d \times \mathbb{R}^{d \times d} \to \mathbb{R}$ be a modified Hamiltonian such that $(t, x, a, y, z) \in [0, T] \times \mathbb{R}^d \times \mathbb{R}^k \times \mathbb{R}^d \times \mathbb{R}^{d \times d}$,

$$\begin{aligned}
\overline{H}_t(x, a, y, z) = &\tfrac{1}{2}\mathrm{tr}\left((C_t x + D_t a)(C_t x + D_t a)^\top z\right) + \langle A_t x + B_t a, y\rangle \\
&+ \tfrac{1}{2}\left(x^\top Q_t x + x^\top S_t^\top a + a^\top S_t x + a^\top(R_t + \rho \bar{V}_t^{-1})a\right),
\end{aligned}$$

and let $\ell : [0, T] \times \mathbb{S}_+^k \times \mathbb{R}^{d \times d} \to \mathbb{R}$ be defined as in (3.4). Recall that $(\nabla_x u)_t(x) = P_t x$ and $(\nabla_x^2 u)_t(x) = P_t$. Hence by Lemma 3.1,

$$F_t(x) = \overline{H}_t(x, K_t' x, P_t x, P_t) - \overline{H}_t(x, K_t x, P_t x, P_t) + \ell_t(V_t', P_t) - \ell_t(V_t, P_t). \tag{3.7}$$

Observe that for all $(t, x, y, z) \in [0, T] \times \mathbb{R}^k \times \mathbb{R}^d \times \mathbb{S}^d$, $a \mapsto \overline{H}_t(x, a, y, z)$ is a quadratic function, and hence Taylor's expansion shows that for all $a, a \in \mathbb{R}^k$,

$$\overline{H}_t(x, a', y, z) - \overline{H}_t(x, a, y, z) = \langle a' - a, \partial_a \overline{H}_t(x, a, y, z)\rangle + \frac{1}{2}\langle a' - a, \partial_a^2 \overline{H}_t(x, a, y, z)(a' - a)\rangle,$$

16

where $\partial_a \overline{H}_t(x, a, y, z)$ and $\partial_a^2 \overline{H}_t(x, a, y, z)$ are given by

$$\partial_a \overline{H}_t(x, a, y, z) = D_t^\top z (C_t x + D_t a) + B_t^\top y + S_t x + (R_t + \rho \bar{V}_t^{-1}) a,$$
$$\partial_a^2 \overline{H}_t(x, a, y, z) = D_t^\top z D_t + R_t + \rho \bar{V}_t^{-1}.$$

Substituting the above identities into (3.7) yields

$$\begin{aligned}
F_t(x) = &\ \ell_t(V_t', P_t) - \ell_t(V_t, P_t) \\
&+ \langle (K_t' - K_t) x, D_t^\top P_t (C_t x + D_t K_t x) + B_t^\top P_t x + S_t x + (R_t + \rho \bar{V}_t^{-1}) K_t x \rangle \\
&+ \frac{1}{2} \langle (K_t' - K_t) x, (D_t^\top P_t D_t + R_t + \rho \bar{V}_t^{-1})(K_t' - K_t) x \rangle,
\end{aligned}$$

which along with (3.6), the definition of $\mathcal{D}_K(\theta)$ in (2.12), and $\Sigma_t' = \mathbb{E}[X_t^{\theta'}(X_t^{\theta'})^\top]$ leads to the desired conclusion. $\qquad \square$

*Proof of Proposition 2.1.* For each $\theta \in \Theta$, by (3.3),

$$\begin{aligned}
\mathcal{C}(\theta) = &\frac{1}{2} \int_0^T \left( \operatorname{tr} \left( (Q_t + K_t^\top S_t + S_t^\top K_t + K_t^\top (R_t + \rho \bar{V}_t^{-1}) K_t) \Sigma_t^\theta \right) \right. \\
&\left. + \operatorname{tr}(R_t V_t) + \rho \left( \operatorname{tr}(\bar{V}_t^{-1} V_t) - k + \ln \left( \frac{\det(\bar{V}_t)}{\det(V_t)} \right) \right) \right) \mathrm{d}t + \frac{1}{2} \operatorname{tr}(G \Sigma_T^\theta),
\end{aligned} \tag{3.8}$$

where $\Sigma^\theta \in C([0, T]; \overline{\mathbb{S}_+^d})$ satisfies (2.11). We then apply [4, Corollary 4.11] to characterise the Gateaux derivatives. Let $H : [0, T] \times \overline{\mathbb{S}_+^d} \times \mathbb{R}^{k \times d} \times \mathbb{S}_+^k \times \mathbb{R}^d \to \mathbb{R}$ be the Hamiltonian of (3.8)-(2.11) such that for all $(t, \Sigma, K, V, Y) \in [0, T] \times \overline{\mathbb{S}_+^d} \times \mathbb{R}^{k \times d} \times \mathbb{S}_+^k \times \mathbb{R}^{d \times d}$,

$$\begin{aligned}
H_t(\Sigma, K, V, Y) = &\langle (A_t + B_t K) \Sigma + \Sigma (A_t + B_t K)^\top + (C_t + D_t K) \Sigma (C_t + D_t K)^\top + D_t V D_t^\top, Y \rangle \\
&+ \frac{1}{2} \left\{ \operatorname{tr} \left( (Q_t + K^\top S_t + S_t^\top K + K^\top (R_t + \rho \bar{V}_t^{-1}) K) \Sigma \right) + \operatorname{tr}(R_t V) \right. \\
&\left. + \rho \left( \operatorname{tr}(\bar{V}_t^{-1} V) - k + \ln \left( \frac{\det(\bar{V}_t)}{\det(V)} \right) \right) \right\},
\end{aligned}$$

and for each $\theta \in \Theta$, let $Y^\theta \in C([0, T]; \mathbb{R}^{d \times d})$ be the adjoint process satisfying

$$(\tfrac{\mathrm{d}}{\mathrm{d}t} Y)_t = -\partial_\Sigma H_t(\Sigma_t^\theta, K_t, V_t, Y_t), \quad \text{a.e. } t \in [0, T]; \quad Y_T = \tfrac{1}{2} G.$$

Then by [4, Corollary 4.11], for all $\theta, \theta \in \Theta$,

$$\frac{\mathrm{d}}{\mathrm{d}\varepsilon} \mathcal{C}(K + \varepsilon K', V) \Big|_{\varepsilon=0} = \int_0^T \langle \partial_K H_t(\Sigma_t^\theta, K_t, V_t, Y_t^\theta), K_t' \rangle \, \mathrm{d}t,$$
$$\frac{\mathrm{d}}{\mathrm{d}\varepsilon} \mathcal{C}(K, V + \varepsilon(V' - V)) \Big|_{\varepsilon=0} = \int_0^T \langle \partial_V H_t(\Sigma_t^\theta, K_t, V_t, Y_t^\theta), V_t' - V_t \rangle \, \mathrm{d}t.$$

Observe that $Y^\theta = \frac{1}{2} P^\theta \in C([0, T]; \mathbb{S}^d)$, and for all $(t, \Sigma, K, V, Y) \in [0, T] \times \overline{\mathbb{S}_+^d} \times \mathbb{R}^{k \times d} \times \mathbb{S}_+^k \times \mathbb{S}^d$,

$$\partial_K H_t(\Sigma, K, V, Y) = \left( 2 B_t^\top Y + 2 D_t^\top Y(C_t + D_t K) + S_t + (R_t + \rho \bar{V}_t^{-1}) K \right) \Sigma,$$
$$\partial_V H_t(\Sigma, K, V, Y) = D_t^\top Y D_t + \tfrac{1}{2} (R_t + \rho(\bar{V}_t^{-1} - V^{-1})).$$

This proves the desired claims. $\qquad \square$

*Proof of Proposition 2.2.* Observe from a direct computation that for all $Z, \Gamma \in \mathbb{R}^{k \times d}$, $\Sigma \in \overline{\mathbb{S}^k_+}$ and $M \in \mathbb{S}^k_+$,

$$\langle Z, \Gamma\Sigma \rangle + \frac{1}{2}\langle Z, MZ\Sigma \rangle = \frac{1}{2}\left\langle Z + M^{-1}\Gamma, M\left(Z + M^{-1}\Gamma\right)\Sigma \right\rangle - \frac{1}{2}\langle M^{-1}\Gamma, \Gamma\Sigma \rangle$$

$$\geq -\frac{1}{2}\langle M^{-1}\Gamma, \Gamma\Sigma \rangle, \tag{3.9}$$

where the last inequality uses the fact that $\mathrm{tr}(AB) \geq 0$ if $A, B \in \overline{\mathbb{S}^d_+}$. Hence for all $\theta \in \Theta$ and $t \in [0, T]$, substituting (3.9) with $Z = K_t^\star - K_t$, $\Gamma = \mathcal{D}_K(\theta)_t$, $\Sigma = \Sigma_t^{\theta^\star}$ and $M = D_t^\top P_t^\theta D_t + R_t + \rho \bar{V}_t^{-1}$ yields that

$$\int_0^T \left( \langle K_t^\star - K_t, \mathcal{D}_K(\theta)_t \Sigma_t^{\theta^\star} \rangle + \frac{1}{2}\langle K_t^\star - K_t, (D_t^\top P_t^\theta D_t + R_t + \rho \bar{V}_t^{-1})(K_t^\star - K_t)\Sigma_t^{\theta^\star} \rangle \right) \mathrm{d}t$$

$$\geq -\frac{1}{2}\int_0^T \langle (D_t^\top P_t^\theta D_t + R_t + \rho \bar{V}_t^{-1})^{-1}\mathcal{D}_K(\theta)_t, \mathcal{D}_K(\theta)_t \Sigma_t^{\theta^\star} \rangle \, \mathrm{d}t. \tag{3.10}$$

Then by Lemma 3.3 and (3.10):

$$\mathcal{C}(\theta^\star) - \mathcal{C}(\theta)$$

$$= \int_0^T \left( \langle K_t^\star - K_t, \mathcal{D}_K(\theta)_t \Sigma_t^{\theta^\star} \rangle + \frac{1}{2}\langle K_t^\star - K_t, (D_t^\top P_t^\theta D_t + R_t + \rho \bar{V}_t^{-1})(K_t^\star - K_t)\Sigma_t^{\theta^\star} \rangle \right.$$

$$\left. + \ell_t(V_t^\star, P_t^\theta) - \ell_t(V_t, P_t^\theta) \right) \mathrm{d}t$$

$$\geq \int_0^T \left( -\frac{1}{2}\langle (D_t^\top P_t^\theta D_t + R_t + \rho \bar{V}_t^{-1})^{-1}\mathcal{D}_K(\theta)_t, \mathcal{D}_K(\theta)_t \Sigma_t^{\theta^\star} \rangle + \ell_t(V_t^\star, P_t^\theta) - \ell_t(V_t, P_t^\theta) \right) \mathrm{d}t. \tag{3.11}$$

Now by (3.4), for all $(t, Z) \in [0, T] \times \mathbb{R}^{d \times d}$ and $V, V' \in \mathbb{S}^k_+$,

$$\ell_t(V', Z) - \ell_t(V, Z)$$

$$= \langle \partial_V \ell_t(V, Z), V' - V \rangle + \int_0^1 \left( \frac{\mathrm{d}}{\mathrm{d}s}\ell_t(V + s(V' - V), Z) - \langle \partial_V \ell_t(V, Z), V' - V \rangle \right) \mathrm{d}s$$

$$= \langle \partial_V \ell_t(V, Z), V' - V \rangle + \int_0^1 \langle \partial_V \ell_t(V + s(V' - V), Z) - \partial_V \ell_t(V, Z), V' - V \rangle \, \mathrm{d}s.$$

Recall that $\partial_V \ell_t(V, Z) = \frac{1}{2}(D_t^\top Z D_t + R_t + \rho \bar{V}_t^{-1} - \rho V^{-1})$, and $A^{-1} - B^{-1} = B^{-1}(B - A)A^{-1}$ for all $A, B \in \mathbb{S}^k_+$. Then for all $(t, Z) \in [0, T] \times \mathbb{R}^{d \times d}$ and $V, V' \in \mathbb{S}^k_+$,

$$\ell_t(V', Z) - \ell_t(V, Z)$$

$$= \langle \partial_V \ell_t(V, Z), V' - V \rangle + \frac{\rho}{2}\int_0^1 \langle V^{-1}\big(s(V' - V)\big)(V + s(V' - V))^{-1}, V' - V \rangle \, \mathrm{d}s. \tag{3.12}$$

Hence for all $\theta, \theta' \in \Theta$ and $t \in [0, T]$, by using (2.13), the fact that $\mathrm{tr}(AB) \geq 0$ for all $A, B \in \overline{\mathbb{S}^d_+}$, and (3.9) (with $Z = V_t' - V_t$, $\Gamma = \mathcal{D}_V(\theta)_t$, $\Sigma = I_k$, $M = \frac{\rho}{2}\Lambda(V_t', V_t)^2 I_k$),

$$\ell_t(V_t', P_t^\theta) - \ell_t(V_t, P_t^\theta)$$

$$= \langle \partial_V \ell_t(V_t, P_t^\theta), V_t' - V_t \rangle + \frac{\rho}{2}\int_0^1 \langle V_t^{-1}\big(s(V_t' - V_t)\big)(V_t + s(V_t' - V_t))^{-1}, V_t' - V_t \rangle \, \mathrm{d}s \tag{3.13}$$

$$\geq \langle \mathcal{D}_V(\theta)_t, V_t' - V_t \rangle + \frac{\rho}{4}\Lambda(V_t', V_t)^2 \langle V_t' - V_t, V_t' - V_t \rangle \geq -\frac{1}{\rho \Lambda(V_t', V_t)^2}|\mathcal{D}_V(\theta)_t|^2,$$

18

with $\Lambda(V_t', V_t) > 0$ defined as

$$\Lambda(V_t', V_t) := \min_{s \in [0,1]} \lambda_{\min} \left( (V_t + s(V_t' - V_t))^{-1} \right) = \frac{1}{\max_{s \in [0,1]} \lambda_{\max}(V_t + s(V_t' - V_t))}$$

$$= \frac{1}{\max(\lambda_{\max}(V_t), \lambda_{\max}(V_t'))} = \frac{1}{\max(\|V_t\|_2, \|V_t'\|_2)},$$

due to the convexity of $[0,1] \ni s \mapsto \lambda_{\max}(V_t + s(V_t' - V_t)) \in \mathbb{R}$, and $\|V\|_2 = \lambda_{\max}(V)$ for all $V \in \overline{\mathbb{S}_+^k}$. Substituting (3.13) with $V' = V^\star$ and using (3.11) yield the desired estimate (2.14). $\qquad\square$

*Proof of Proposition 2.3.* By (2.13) and (3.12), for all $\theta, \theta' \in \Theta$ and $t \in [0, T]$,

$$\ell_t(V_t', P_t^\theta) - \ell_t(V_t, P_t^\theta)$$

$$= \langle \partial_V \ell_t(V_t, P_t^\theta), V_t' - V_t \rangle + \frac{\rho}{2} \int_0^1 \langle V_t^{-1} \big( s(V_t' - V_t) \big)(V_t + s(V_t' - V_t))^{-1}, V_t' - V_t \rangle \, \mathrm{d}s$$

$$\leq \langle \mathcal{D}_V(\theta)_t, V_t' - V_t \rangle + \frac{\rho}{4} \overline{\Lambda}(V_t', V_t)^2 \langle V_t' - V_t, V_t' - V_t \rangle,$$

where $\overline{\Lambda}(V_t', V_t) > 0$ is given by

$$\overline{\Lambda}(V_t', V_t) := \max_{s \in [0,1]} \lambda_{\max} \left( (V_t + s(V_t' - V_t))^{-1} \right) = \frac{1}{\min_{s \in [0,1]} \lambda_{\min}(V_t + s(V_t' - V_t))}$$

$$= \frac{1}{\min(\lambda_{\min}(V_t), \lambda_{\min}(V_t'))}.$$

Combining this and Lemma 3.3 yields the desired estimate. $\qquad\square$

## 3.2 Proof of Proposition 2.5

The following lemma compares solutions to (2.10) for different $\theta, \theta' \in \Theta$.

**Lemma 3.4.** *Suppose (H.1) holds. For each $\theta \in \Theta$, let $P^\theta \in C([0,T]; \mathbb{S}^d)$ satisfy (2.10). Then for all $\theta, \theta' \in \Theta$, $\Delta P := P^{\theta'} - P^\theta$ satisfies for a.e. $t \in [0, T]$,*

$$(\tfrac{\mathrm{d}}{\mathrm{d}t} \Delta P)_t + (A_t + B_t K_t')^\top \Delta P_t + \Delta P_t^\top (A_t + B_t K_t') + (C_t + D_t K_t')^\top \Delta P_t (C_t + D_t K_t')$$
$$+ (K_t' - K_t)^\top \mathcal{D}_K(\theta)_t + \mathcal{D}_K(\theta)_t^\top (K_t' - K_t)$$
$$+ (K_t' - K_t)^\top (D_t^\top P_t^\theta D_t + R_t + \rho \bar{V}_t^{-1})(K_t' - K_t), = 0; \quad \Delta P_T = 0,$$

*where $\mathcal{D}_K(\theta)_t$ is defined in (2.12).*

*Proof.* By (2.10), $\Delta P_T = 0$ and for a.e. $t \in [0, T]$,

$$(\tfrac{\mathrm{d}}{\mathrm{d}t} \Delta P)_t + (A_t + B_t K_t')^\top \Delta P_t + \Delta P_t^\top (A_t + B_t K_t') + (C_t + D_t K_t')^\top \Delta P_t (C_t + D_t K_t')$$
$$+ q_t(K_t') - q_t(K_t) = 0,$$

where for all $K \in \mathbb{R}^{k \times d}$,

$$q_t(K) := (A_t + B_t K)^\top P_t^\theta + (P_t^\theta)^\top (A_t + B_t K) + (C_t + D_t K)^\top P_t^\theta (C_t + D_t K)$$
$$+ S_t^\top K + K^\top S_t + K^\top (R_t + \rho \bar{V}_t^{-1}) K.$$

19

Observe that for any $K_1, K_2 \in \mathbb{R}^{k \times d}$ and $P \in \mathbb{S}^d$,

$$K_1^\top P K_1 - K_2^\top P K_2 = (K_1 - K_2)^\top P K_2 + K_2^\top P (K_1 - K_2) + (K_1 - K_2)^\top P (K_1 - K_2).$$

Thus for a.e. $t \in [0, T]$,

$$q_t(K_t') - q_t(K_t) = (K_t' - K_t)^\top \left( B_t^\top P_t^\theta + D_t^\top P_t^\theta (C_t + D_t K_t) + S_t + (R_t + \rho \bar{V}_t^{-1}) K_t \right)$$
$$+ \left( B_t^\top P_t^\theta + D_t^\top P_t^\theta (C_t + D_t K_t) + S_t + (R_t + \rho \bar{V}_t^{-1}) K_t \right)^\top (K_t' - K_t)$$
$$+ (K_t' - K_t)^\top (D_t^\top P_t^\theta D_t + R_t + \rho \bar{V}_t^{-1})(K_t' - K_t),$$

which along with the definition of $\mathcal{D}_K(\theta)_t$ leads to the desired identity. $\qquad\square$

Based on Lemma 3.4, we establish a uniform bound of $(P^{\theta^n})_{n \in \mathbb{N}}$ and $(K^n)_{n \in \mathbb{N}}$.

**Proposition 3.5.** *Suppose (H.1) and (H.2) hold. For each $\theta \in \Theta$, let $P^\theta \in C([0, T]; \mathbb{S}^d)$ satisfy (2.10). Let $\theta^0 \in \Theta$, $\bar{\lambda}_0 > 0$ be such that $\bar{\lambda}_0 I_k \succeq D^\top P^{\theta^0} D + R + \rho \bar{V}^{-1}$, and for each $\tau > 0$, let $(K^n)_{n \in \mathbb{N}} \subset \mathcal{B}(0, T; \mathbb{R}^{d \times k})$ be defined in (2.16). Then*

*(1) for all $\tau \in (0, 2/\bar{\lambda}_0]$ and $n \in \mathbb{N}_0$, $P^{\theta^n} \succeq P^{\theta^{n+1}} \succeq P^\star$, and $\widetilde{\delta} I_k \preceq D^\top P^{\theta^n} D + R + \rho \bar{V}^{-1} \preceq \bar{\lambda}_0 I_k$, with $P^\star \in C([0, T]; \mathbb{S}^d)$ and $\widetilde{\delta} > 0$ in (H.2);*

*(2) there exists $\widetilde{C}_{(\theta^0)} \geq 0$ such that for all $\tau \in (0, 1/\bar{\lambda}_0]$ and $n \in \mathbb{N}_0$, $\|K^n\|_{L^2} \leq \widetilde{C}_{(\theta^0)}$.*

*Proof.* We write $P^n = P^{\theta^n}$ for notational simplicity. For each $n \in \mathbb{N}$, applying (2.16) and Lemma 3.4 with $\theta' = \theta^n$ and $\theta = \theta^{n-1}$, $\Delta P := P^n - P^{n-1} \in C([0, T]; \mathbb{S}^d)$ satisfies $\Delta P_T = 0$, and for a.e. $t \in [0, T]$,

$$(\tfrac{\mathrm{d}}{\mathrm{d}t} \Delta P)_t + (A_t + B_t K_t^{n+1})^\top \Delta P_t + \Delta P_t^\top (A_t + B_t K_t^{n+1}) + (C_t + D_t K_t^{n+1})^\top \Delta P_t (C_t + D_t K_t^{n+1})$$
$$= -(K_t^{n+1} - K_t^n)^\top \mathcal{D}_K(\theta^n)_t - \mathcal{D}_K(\theta^n)_t^\top (K_t^{n+1} - K_t^n)$$
$$\quad - (K_t^{n+1} - K_t^n)^\top (D_t^\top P_t^n D_t + R_t + \rho \bar{V}_t^{-1})(K_t^{n+1} - K_t^n)$$
$$= 2\tau \mathcal{D}_K(\theta^n)_t^\top \left( I_k - \tfrac{\tau}{2}(D_t^\top P_t^n D_t + R_t + \rho \bar{V}_t^{-1}) \right) \mathcal{D}_K(\theta^n)_t.$$

Now suppose that $\tau \in (0, 2/\bar{\lambda}_0]$, then $I_k - \tfrac{\tau}{2}(D_t^\top P_t^0 D_t + R_t + \rho \bar{V}_t^{-1}) \succeq 0$, which implies that $P^1 \preceq P^0$ (see e.g., [34, Lemma 7.3, p. 320]), and hence

$$I_k - \tfrac{\tau}{2}(D^\top P^1 D + R + \rho \bar{V}^{-1}) \succeq I_k - \tfrac{\tau}{2}(D^\top P^0 D + R + \rho \bar{V}^{-1}) \succeq 0.$$

An induction argument shows that $P^n \succeq P^{n+1}$ for all $n \in \mathbb{N}_0$. Moreover, observe from (2.5) and (2.6) that $\mathcal{D}_K(\theta^\star) = 0$ and $P^\star = P^{\theta^\star}$. By applying Lemma 3.4 with $\theta' = \theta^n$ and $\theta = \theta^\star$, one can deduce from similar arguments that $P^n \succeq P^{\theta^\star}$ for all $n \in \mathbb{N}_0$. Consequently, by (H.2),

$$\bar{\lambda}_0 I_k \succeq D^\top P^0 D + R + \rho \bar{V}^{-1} \succeq D^\top P^n D + R + \rho \bar{V}^{-1} \succeq D^\top P^\star D + R + \rho \bar{V}^{-1} \succeq \widetilde{\delta} I_k.$$

This proves Item (1).

Item (1) implies that there exists $\widetilde{C}_{(\theta^0)} > 0$ such that $\|P^n\|_{L^\infty} \leq \widetilde{C}_{(\theta^0)}$ for all $n \in \mathbb{N}_0$. Then for all $n \in \mathbb{N}_0$, by (2.12) and (2.16),

$$|K_t^{n+1}| = \left| K_t^n - \tau \left( (B_t^\top P_t^n + D_t^\top P_t^n C_t + S_t) + (D_t^\top P_t^n D_t + R_t + \rho \bar{V}_t^{-1}) K_t^n \right) \right|$$
$$\leq |I_k - \tau(D_t^\top P_t^n D_t + R_t + \rho \bar{V}_t^{-1})||K_t^n| + \tau |B_t^\top P_t^n + D_t^\top P_t^n C_t + S_t|.$$

20

Thus for all $\tau \in (0, 1/\overline{\lambda}_0]$ and $n \in \mathbb{N}_0$,

$$\|K^{n+1}\|_{L^2} \leq (1 - \tau\overline{\lambda}_0)\|K^n\|_{L^2} + \tau\|B^\top P^n + D^\top P^n C + S\|_{L^2}$$

$$\leq \|K^0\|_{L^2} + \sup_{n \in \mathbb{N}_0} \frac{1}{\overline{\lambda}_0}\|B^\top P^n + D^\top P^n C + S\|_{L^2} < \infty,$$

where the last inequality follows from a straightforward induction argument. $\square$

The next proposition proves a uniform upper and lower bound of $(V^n)_{n \in \mathbb{N}}$.

**Proposition 3.6.** *Suppose (H.1) and (H.2) hold. Let $\theta^0 \in \Theta$, and for each $\tau > 0$, let $(\theta^n)_{n \in \mathbb{N}} \subset \mathcal{B}(0, T; \mathbb{R}^{k \times d} \times \mathbb{S}^k)$ be defined in (2.16). Let $\overline{\lambda}_0 > 0$ be such that $\overline{\lambda}_0 I_k \succeq D^\top P^{\theta^0} D + R + \rho\overline{V}^{-1}$ with $P^{\theta_0} \in C([0, T]; \mathbb{S}^d)$ defined in (2.10), let $\underline{\lambda}_V = \min\left(\min_{t \in [0,T]} \lambda_{\min}(V_t^0), \frac{\rho}{\overline{\lambda}_0}\right)$, and let $\overline{\lambda}_V = \max\left(\max_{t \in [0,T]} \lambda_{\max}(V_t^0), \frac{\rho}{\delta}\right)$. Then for all $\tau \in (0, 1/\overline{\lambda}_0]$ and $n \in \mathbb{N}_0$, $\underline{\lambda}_V I_k \preceq V^n \preceq \overline{\lambda}_V I_k$.*

*Proof.* For each $n \in \mathbb{N}_0$, let $M^n = D^\top P^{\theta^n} D + R + \rho\overline{V}^{-1}$. By (2.13), for each $n \in \mathbb{N}_0$ and a.e. $t \in [0, T]$,

$$V_t^{n+1} = V_t^n - \tau\left(\frac{1}{2}\left(M_t^n - \rho(V_t^n)^{-1}\right)V_t^n + V_t^n\frac{1}{2}\left(M_t^n - \rho(V_t^n)^{-1}\right)\right)$$

$$= \frac{1}{2}\left(I_k - \tau M_t^n\right)V_t^n + \frac{1}{2}V_t^n\left(I_k - \tau M_t^n\right) + \rho\tau I_k.$$

Let $\tau \in (0, 1/\overline{\lambda}_0]$. By Proposition 3.5 Item (1), for all $n \in \mathbb{N}_0$, $\widetilde{\delta}I_k \preceq M^n \preceq \overline{\lambda}_0 I_k$, and hence $0 \preceq (1 - \tau\overline{\lambda}_0)I_k \preceq I_k - \tau M^n \preceq (1 - \tau\widetilde{\delta})I_k$. Thus for all $n \in \mathbb{N}_0$ and a.e. $t \in [0, T]$,

$$\lambda_{\min}(V_t^{n+1}) \geq \lambda_{\min}\left(I_k - \tau M_t^n\right)\lambda_{\min}(V_t^n) + \rho\tau \geq \left(1 - \tau\overline{\lambda}_0\right)\lambda_{\min}(V_t^n) + \rho\tau.$$

Setting $v_t^n = \lambda_{\min}(V_t^n)$ for all $n \in \mathbb{N}_0$. An induction argument shows that

$$v_t^n \geq \left(1 - \tau\overline{\lambda}_0\right)^n v_t^0 + \rho\tau \sum_{i=0}^{n-1}\left(1 - \tau\overline{\lambda}_0\right)^i = \left(v_t^0 - \frac{\rho}{\overline{\lambda}_0}\right)\left(1 - \tau\overline{\lambda}_0\right)^n + \frac{\rho}{\overline{\lambda}_0} \geq \min\left(v_t^0, \frac{\rho}{\overline{\lambda}_0}\right).$$

Similarly, for all $n \in \mathbb{N}_0$ and a.e. $t \in [0, T]$,

$$\lambda_{\max}(V_t^{n+1}) \leq \lambda_{\max}\left(I_k - \tau M_t^n\right)\lambda_{\max}(V_t^n) + \rho\tau \leq \left(1 - \tau\widetilde{\delta}\right)\lambda_{\max}(V_t^n) + \rho\tau,$$

which implies that $\lambda_{\max}(V_t^n) \leq \max\left(\lambda_{\max}(V_t^0), \frac{\rho}{\delta}\right)$. $\square$

The following lemma establishes an upper and lower bounds of the state covariance matrices for any $\theta \in \Theta$, which is crucial for the convergence analysis of (2.16).

**Lemma 3.7.** *Suppose (H.1) and (H.2) hold. For each $\theta \in \Theta$, let $\Sigma^\theta \in C([0, T]; \overline{\mathbb{S}_+^d})$ satisfy (2.11). Then there exists $\widetilde{C} > 0$ such that for all $\theta \in \Theta$,*

$$\lambda_{\min}(\mathbb{E}[\xi_0\xi_0^\top])\exp\left(-\widetilde{C}(1 + \|K\|_{L^2}^2)\right)I_d \preceq \Sigma^\theta \preceq \widetilde{C}\left(|\Sigma_0| + \|V\|_{L^1}\right)\exp\left(\widetilde{C}(1 + \|K\|_{L^2}^2)\right)I_d.$$

*Proof.* Let $\theta \in \Theta$ be fixed. We omit the superscript of $\Sigma^\theta$ to simplify the notation. To estimate $\lambda_{\max}(\Sigma_t)$, by (2.11), for all $t \in [0, T]$,

$$\|\Sigma_t\|_2 \leq \|\Sigma_0\|_2 + \int_0^t \left( (2\|\widetilde{A}_s\|_2 + \|\widetilde{C}_s\|_2^2)\|\Sigma_s\|_2 + \|D_s\|_2^2\|V_s\|_2 \right) \mathrm{d}s,$$

where $\widetilde{A}_t = A_t + B_t K_t$ and $\widetilde{C}_t = C_t + D_t K_t$. Then by (H.1(1)) and Gronwall's inequality, $\|\Sigma\|_{L^\infty} \leq \widetilde{C}\left(|\Sigma_0| + \|V\|_{L^1}\right) \exp\left(\widetilde{C}(1 + \|K\|_{L^2}^2)\right)$ for some $\widetilde{C}_1 > 0$.

Now we obtain a lower bound of $\lambda_{\min}(\Sigma_t)$. As $(C + DK)\Sigma(C + DK)^\top + DVD^\top \succeq 0$, by (2.11), $\Sigma \succeq \widetilde{\Sigma}$, where $\widetilde{\Sigma} \in C([0, T]; \mathbb{S}_+^d)$ satisfies for a.e. $t \in [0, T]$,

$$\left(\tfrac{\mathrm{d}}{\mathrm{d}t}\Sigma\right)_t = (A_t + B_t K_t)\Sigma_t + \Sigma_t(A_t + B_t K_t)^\top; \quad \Sigma_0 = \mathbb{E}[\xi_0\xi_0^\top]. \tag{3.14}$$

Note that for all $t \in [0, T]$, $\widetilde{\Sigma}_t = \Psi_t^\top \mathbb{E}[\xi_0\xi_0^\top]\Psi_t$, where $\Psi \in C([0, T]; \mathbb{R}^{d \times d})$ satisfies $\Psi_0 = I_d$ and for a.e. $t \in [0, T]$, $\mathrm{d}\Psi_t = \Psi_t \widetilde{A}_t^\top \mathrm{d}t$, with $\widetilde{A} = A + BK \in L^1(0, T; \mathbb{R}^{d \times d})$. For each $t \in [0, T]$, let $x_t \in \mathbb{R}^d$ be such that $|x_t| = 1$ and $\lambda_{\min}(\widetilde{\Sigma}_t) = x_t^\top \widetilde{\Sigma}_t x_t$, and let $y_t = \Psi_t x_t$. Then

$$\lambda_{\min}(\Sigma_t) \geq \lambda_{\min}(\widetilde{\Sigma}_t) = x_t^\top \left((\Psi_t)^\top \mathbb{E}[\xi_0\xi_0^\top]\Psi_t\right) x_t = \frac{y_t^\top \mathbb{E}[\xi_0\xi_0^\top]y_t}{|y_t|^2}|y_t|^2 \geq \frac{\lambda_{\min}(\mathbb{E}[\xi_0\xi_0^\top])}{\left\|\Psi_t^{-1}\right\|_2^2},$$

where the last inequality uses $1 = |x_t| \leq \|(\Psi_t)^{-1}\|_2|y_t|$, with the spectral norm $\|\cdot\|_2$. Observe that $\Psi^{-1} \in C([0, T]; \mathbb{R}^{d \times d})$ be such that $\Psi_0^{-1} = I_d$ and for a.e. $t \in [0, T]$, $\mathrm{d}\Psi_t^{-1} = -\widetilde{A}_t^\top \Psi_t^{-1}\mathrm{d}t$. Hence for all $t \in [0, T]$,

$$\|\Psi_t^{-1}\|_2 \leq 1 + \int_0^t \|\widetilde{A}_s\|_2\|\Psi_s^{-1}\|_2\,\mathrm{d}s \leq 1 + \int_0^t |\widetilde{A}_s|\|\Psi_s^{-1}\|_2\,\mathrm{d}s,$$

which along with Gronwall's inequality shows that $\|\Psi_t^{-1}\|_{L^\infty} \leq \exp\left(\|\widetilde{A}\|_{L^1}\right)$. Consequently, $\inf_{t \in [0, T]} \lambda_{\min}(\Sigma_t) \geq \lambda_{\min}(\mathbb{E}[\xi_0\xi_0^\top]) \exp\left(-2\|\widetilde{A}\|_{L^1}\right)$, which along with (H.1(1)) leads to the desired lower bound of $\lambda_{\min}(\Sigma_t)$. $\qquad\square$

A direct consequence of Proposition 3.6 and Lemma 3.7 are the following uniform bounds of the state covariance matrices along the iterates $(\theta^n)_{n \in \mathbb{N}}$ generated by (2.16).

**Proposition 3.8.** *Suppose (H.1) and (H.2) hold, and $\mathbb{E}[\xi_0\xi_0^\top] \succ 0$. For each $\theta \in \Theta$, let $P^\theta \in C([0, T]; \mathbb{S}^d)$ satisfy (2.10), and let $\Sigma^\theta \in C([0, T]; \overline{\mathbb{S}_+^d})$ satisfy (2.11). Let $\theta^0 \in \Theta$, let $\overline{\lambda}_0 > 0$ be such that $\overline{\lambda}_0 I_k \succeq D^\top P^{\theta^0}D + R + \rho\overline{V}^{-1}$, and for each $\tau \in (0, 1/\overline{\lambda}_0)$, let $(\theta^n)_{n \in \mathbb{N}} \subset \Theta$ be defined in (2.16). Then there exists $\overline{\lambda}_X, \underline{\lambda}_X > 0$, depending on $\theta_0$, such that for all $\tau \in (0, 1/\overline{\lambda}_0]$ and $n \in \mathbb{N}_0$, $\underline{\lambda}_X I_d \preceq \Sigma^{\theta^n} \preceq \overline{\lambda}_X I_d$.*

*Proof.* By Proposition 3.5, for all $\tau \in (0, 1/\overline{\lambda}_0]$, $\sup_{n \in \mathbb{N}_0} \|K^n\|_{L^2} \leq \widetilde{C}_{(\theta^0)}$ for some $\widetilde{C}_{(\theta^0)} > 0$. The uniform lower and upper bounds of $(\Sigma^{\theta^n})_{n \in \mathbb{N}_0}$ follow from Proposition 3.6 and Lemma 3.7. $\quad\square$

## 3.3 Proof of Theorem 2.6

The following proposition compares the value functions of two consecutive iterates.

**Proposition 3.9.** *Suppose (H.1) and (H.2) hold, and $\mathbb{E}[\xi_0\xi_0^\top] \succ 0$. Let $\theta^0 \in \Theta$, and $\overline{\lambda}_0 > 0$ be such that $\overline{\lambda}_0 I_k \succeq D^\top P^{\theta^0} D + R + \rho \bar{V}^{-1}$ with $P^{\theta_0} \in C([0,T];\mathbb{S}^d)$ defined in (2.10). For each $\tau \in (0, 1/\overline{\lambda}_0]$, let $(\theta^n)_{n\in\mathbb{N}} \subset \Theta$ be defined in (2.16), let $\underline{\lambda}_V, \overline{\lambda}_V > 0$ be such that $\underline{\lambda}_V I_k \preceq V^n \preceq \overline{\lambda}_V I_k$ for all $n \in \mathbb{N}_0$ (cf. Proposition 3.6), and let $\underline{\lambda}_X, \overline{\lambda}_X > 0$ be such that $\underline{\lambda}_X I_k \preceq \Sigma^{\theta^n} \preceq \overline{\lambda}_X I_k$ for all $n \in \mathbb{N}_0$ (cf. Proposition 3.8). Then for all $\tau \in (0, 1/\overline{\lambda}_0]$ and $n \in \mathbb{N}_0$,*

$$\mathcal{C}(\theta^{n+1}) - \mathcal{C}(\theta^n) \leq -\tau \int_0^T \left( \left( \underline{\lambda}_X - \frac{\tau}{2}\overline{\lambda}_0\overline{\lambda}_X \right) |\mathcal{D}_K(\theta^n)_t|^2 + \left( 2\underline{\lambda}_V - \frac{\rho\tau\overline{\lambda}_V^2}{\underline{\lambda}_V^2} \right) |\mathcal{D}_V^n(\theta^n)_t|^2 \right) \mathrm{d}t.$$

*Proof.* For each $n \in \mathbb{N}_0$, let $\Sigma^n = \Sigma^{\theta^n}$, $P^n = P^{\theta^n}$, $\Delta K^n = K^{n+1} - K^n$, $\Delta V^n = V^{n+1} - V^n$, $\mathcal{D}_K^n = \mathcal{D}_K(\theta^n)$, and $\mathcal{D}_V^n = \mathcal{D}_V(\theta^n)$. By using Proposition 2.3 and the fact that $\underline{\lambda}_V I_k \preceq V^n \preceq \overline{\lambda}_V I_k$ for all $n \in \mathbb{N}_0$,

$$\begin{aligned}
\mathcal{C}(\theta^{n+1}) - \mathcal{C}(\theta^n) &\leq \int_0^T \left( \langle \Delta K_t^n, \mathcal{D}_{K,t}^n \Sigma_t^{n+1} \rangle + \frac{1}{2} \langle \Delta K_t^n, (D_t^\top P_t^n D_t + R_t + \rho \bar{V}_t^{-1})(\Delta K_t^n)\Sigma_t^{n+1} \rangle \right. \\
&\qquad \left. + \langle \mathcal{D}_{V,t}^n, \Delta V_t^n \rangle + \frac{\rho}{4\underline{\lambda}_V^2} |\Delta V_t^n|^2 \right) \mathrm{d}t \\
&\leq \int_0^T \left( \langle -\tau \mathcal{D}_{K,t}^n, \mathcal{D}_{K,t}^n \Sigma_t^{n+1} \rangle + \frac{\tau^2}{2} \langle \mathcal{D}_{K,t}^n, (D_t^\top P_t^n D_t + R_t + \rho \bar{V}_t^{-1})\mathcal{D}_{K,t}^n \Sigma_t^{n+1} \rangle \right. \\
&\qquad \left. - \tau \left[ \langle \mathcal{D}_{V,t}^n, \{\mathcal{D}_{V,t}^n V_t^n\}_S \rangle - \frac{\rho\tau}{4\underline{\lambda}_V^2} |\{\mathcal{D}_{V,t}^n V_t^n\}_S|^2 \right] \right) \mathrm{d}t
\end{aligned}$$

with $\{\mathcal{D}_{V,t}^n V_t^n\}_S := \mathcal{D}_{V,t}^n V_t^n + V_t^n \mathcal{D}_{V,t}^n$, where the last inequality used (2.16). Recall that for all $S_1, S_2 \in \overline{\mathbb{S}_+^k}$, $\lambda_{\min}(S_1)\mathrm{tr}(S_2) \leq \mathrm{tr}(S_1 S_2) \leq \lambda_{\max}(S_1)\mathrm{tr}(S_2)$. Hence $\langle \mathcal{D}_{V,t}^n, \{\mathcal{D}_{V,t}^n V_t^n\}_S \rangle \geq 2\underline{\lambda}_V |\mathcal{D}_{V,t}^n|^2$, and $|\{\mathcal{D}_{V,t}^n V_t^n\}_S|^2 \leq 4\overline{\lambda}_V^2 |\mathcal{D}_{V,t}^n|^2$. Hence for all $n \in \mathbb{N}_0$,

$$\begin{aligned}
\mathcal{C}(\theta^{n+1}) - \mathcal{C}(\theta^n) &\leq \int_0^T \left( -\tau \left( \lambda_{\min}(\Sigma_t^{n+1}) - \frac{\tau}{2}\lambda_{\max}((D_t^\top P_t^n D_t + R_t + \rho \bar{V}_t^{-1}))\lambda_{\max}(\Sigma_t^{n+1}) \right) |\mathcal{D}_{K,t}^n|^2 \right. \\
&\qquad \left. - \tau \left( 2\underline{\lambda}_V - \frac{\rho\tau\overline{\lambda}_V^2}{\underline{\lambda}_V^2} \right) |\mathcal{D}_{V,t}^n|^2 \right) \mathrm{d}t.
\end{aligned}$$

The desired inequality then follows from Propositions 3.5 and 3.8. $\qquad\square$

The next proposition establishes a uniform Łojasiewicz property of the cost $\mathcal{C} : \Theta \to \mathbb{R}$ along the iterates (2.16).

**Proposition 3.10.** *Suppose (H.1) and (H.2) hold, and $\mathbb{E}[\xi_0\xi_0^\top] \succ 0$. Let $\theta^\star \in \Theta$ be defined in (2.6). For each $\theta \in \Theta$, let $P^\theta \in C([0,T];\mathbb{S}^d)$ satisfy (2.10), and let $\Sigma^\theta \in C([0,T];\overline{\mathbb{S}_+^d})$ satisfy (2.11). Let $\theta^0 \in \Theta$ and $\overline{\lambda}_0 > 0$ such that $\overline{\lambda}_0 I_k \succeq D^\top P^{\theta^0} D + R + \rho \bar{V}^{-1}$. For each $\tau \in (0, 1/\overline{\lambda}_0]$, let $(\theta^n)_{n\in\mathbb{N}_0} \subset \Theta$ be defined in (2.16). Then for all $\tau \in (0, 1/\overline{\lambda}_0]$ and $n \in \mathbb{N}_0$,*

$$\mathcal{C}(\theta^n) - \mathcal{C}(\theta^\star) \leq \max \left( \frac{\overline{\lambda}_X^\star}{2\widetilde{\delta}}, \frac{\max(\overline{\lambda}_V, \overline{\lambda}_V^\star)^2}{\rho} \right) \int_0^T \left( |\mathcal{D}_K(\theta^n)_t|^2 + |\mathcal{D}_V(\theta^n)_t|^2 \right) \mathrm{d}t,$$

*where $\widetilde{\delta} > 0$ is the same as in (H.2), $\overline{\lambda}_X^\star > 0$ satisfies $\Sigma^{\theta^\star} \preceq \overline{\lambda}_X^\star I_d$, $\overline{\lambda}_V^\star > 0$ satisfies $V^\star \preceq \overline{\lambda}_V^\star I_k$, and $\overline{\lambda}_V > 0$ satisfies $V^n \preceq \overline{\lambda}_V I_k$ for all $n \in \mathbb{N}_0$.*

23

*Proof.* Let $\overline{\lambda}_V^\star > 0$ be such that $V^\star \preceq \overline{\lambda}_V^\star I_k$. For each $n \in \mathbb{N}_0$, let $\Sigma^n = \Sigma^{\theta^n}$, $P^n = P^{\theta^n}$, $\mathcal{D}_K^n = \mathcal{D}_K(\theta^n)$, and $\mathcal{D}_V^n = \mathcal{D}_V(\theta^n)$. Recall that there exists $\underline{\lambda}_V, \overline{\lambda}_V > 0$ such that $\underline{\lambda}_V I_k \preceq V^n \preceq \overline{\lambda}_V I_k$ for all $n \in \mathbb{N}_0$. Then for all $n \in \mathbb{N}_0$, Proposition 3.5 shows that $D_t^\top P_t^n D_t + R_t + \rho \bar{V}_t^{-1} \succeq \widetilde{\delta} I_k$, which along with Proposition 2.2 shows that

$$\mathcal{C}(\theta^n) - \mathcal{C}(\theta^\star) \leq \int_0^T \left( \frac{1}{2} \langle (D_t^\top P_t^n D_t + R_t + \rho \bar{V}_t^{-1})^{-1} \mathcal{D}_{K,t}^n, \mathcal{D}_{K,t}^n \Sigma_t^{\theta^\star} \rangle + \frac{\max(\overline{\lambda}_V, \overline{\lambda}_V^\star)^2}{\rho} |\mathcal{D}_{V,t}^n|^2 \right) dt$$

$$\leq \int_0^T \left( \frac{\overline{\lambda}_X^\star}{2\widetilde{\delta}} |\mathcal{D}_{K,t}^n|^2 + \frac{\max(\overline{\lambda}_V, \overline{\lambda}_V^\star)^2}{\rho} |\mathcal{D}_{V,t}^n|^2 \right) dt,$$

with $\overline{\lambda}_X^\star > 0$ such that $\Sigma^{\theta^\star} \preceq \overline{\lambda}_X^\star I_d$ (cf. Lemma 3.7). This proves the desired estimate. $\qquad\square$

*Proof of Theorem 2.6.* Let $\overline{\lambda}_0 > 0$ be such that $\overline{\lambda}_0 I_k \succeq D^\top P^{\theta^0} D + R + \rho \bar{V}^{-1}$, where $P^{\theta_0} \in C([0,T]; \mathbb{S}^d)$ satisfies (2.10) with $\theta = \theta_0$. Then by Proposition 3.9, for all $\tau \in (0, 1/\overline{\lambda}_0]$ and $n \in \mathbb{N}_0$,

$$\mathcal{C}(\theta^{n+1}) - \mathcal{C}(\theta^n) \leq -\tau \int_0^T \left( \left( \underline{\lambda}_X - \frac{\tau}{2} \overline{\lambda}_0 \overline{\lambda}_X \right) |\mathcal{D}_K(\theta^n)_t|^2 + \left( 2\underline{\lambda}_V - \frac{\rho \tau \overline{\lambda}_V^2}{\underline{\lambda}_V^2} \right) |\mathcal{D}_V(\theta^n)_t|^2 \right) dt,$$

with the constants $\underline{\lambda}_X, \overline{\lambda}_X > 0$ in Proposition 3.8. Hence by setting $\widetilde{C}_1 = \max(\overline{\lambda}_0, \frac{2\rho \overline{\lambda}_V^2}{3\underline{\lambda}_V^3}, \frac{\overline{\lambda}_0 \overline{\lambda}_X}{\underline{\lambda}_X})$, it holds for all $\tau \in (0, 1/\widetilde{C}_1]$ and $n \in \mathbb{N}_0$,

$$\mathcal{C}(\theta^{n+1}) - \mathcal{C}(\theta^n) \leq -\tau \int_0^T \left( \frac{\underline{\lambda}_X}{2} |\mathcal{D}_K(\theta^n)_t|^2 + \frac{\underline{\lambda}_V}{2} |\mathcal{D}_V(\theta^n)_t|^2 \right) dt$$

$$\leq -\tau \frac{1}{2} \min(\underline{\lambda}_X, \underline{\lambda}_V) \int_0^T \left( |\mathcal{D}_K(\theta^n)_t|^2 + |\mathcal{D}_V(\theta^n)_t|^2 \right) dt$$

$$\leq -\tau C_1 \left( \mathcal{C}(\theta^n) - \mathcal{C}(\theta^\star) \right), \quad \text{with } C_1 := \frac{\min(\underline{\lambda}_X, \underline{\lambda}_V)}{2 \max \left( \frac{\overline{\lambda}_X^\star}{2\widetilde{\delta}}, \frac{\max(\overline{\lambda}_V, \overline{\lambda}_V^\star)^2}{\rho} \right)},$$

where the last inequality used Proposition 3.10. Thus, for all $\tau \in (0, \tau_0]$ with $\tau_0 > 0$ satisfying

$$\frac{1}{\tau_0} \geq \max \left( \overline{\lambda}_0, \frac{2\rho \overline{\lambda}_V^2}{3\underline{\lambda}_V^3}, \frac{\overline{\lambda}_0 \overline{\lambda}_X}{\underline{\lambda}_X}, \frac{\min(\underline{\lambda}_X, \underline{\lambda}_V)}{2 \max \left( \frac{\overline{\lambda}_X^\star}{2\widetilde{\delta}}, \frac{\max(\overline{\lambda}_V, \overline{\lambda}_V^\star)^2}{\rho} \right)} \right),$$

we have for all $n \in \mathbb{N}_0$, $\mathcal{C}(\theta^{n+1}) \leq \mathcal{C}(\theta^n)$ and

$$\mathcal{C}(\theta^{n+1}) - \mathcal{C}(\theta^\star) \leq \mathcal{C}(\theta^{n+1}) - \mathcal{C}(\theta^n) + \mathcal{C}(\theta^n) - \mathcal{C}(\theta^\star) \leq (1 - \tau C_1) \left( \mathcal{C}(\theta^n) - \mathcal{C}(\theta^\star) \right). \tag{3.15}$$

To prove Item (2), observe that $\mathcal{D}_K(\theta^\star) = 0$ and $\mathcal{D}_V(\theta^\star) = 0$. Hence by Lemma 3.3 and (3.13), for all $n \in \mathbb{N}_0$,

$$\mathcal{C}(\theta^n) - \mathcal{C}(\theta^\star)$$
$$\geq \int_0^T \left( \frac{1}{2} \langle K_t^n - K_t^\star, (D_t^\top P_t^\star D_t + R_t + \rho \bar{V}_t^{-1})(K_t^n - K_t^\star) \Sigma_t^{\theta^n} \rangle + \frac{\rho}{4} \frac{|V_t^n - V_t^\star|^2}{\max(\|V_t^n\|_2^2, \|V_t^\star\|_2^2)} \right) dt$$
$$\geq \int_0^T \left( \frac{1}{2} \widetilde{\delta} \underline{\lambda}_X |K_t^n - K_t^\star|^2 + \frac{\rho}{4 \overline{\lambda}_V^2} |V_t^n - V_t^\star|^2 \right) dt,$$

where the last inequality used (H.2), Proposition 3.8 and $V^\star, V^n \preceq \overline{\lambda}_V I_k$. This along with Item (1) proves Item (2) with $C_2 = 1 / \min \left( \frac{1}{2} \widetilde{\delta} \underline{\lambda}_X, \frac{\rho}{4 \overline{\lambda}_V^2} \right)$. $\qquad\square$

## 3.4 Proofs of Theorem 2.7 and Corollary 2.8

The following lemma proves the optimal costs of piecewise constant policies converges to the optimal cost of continuous-time policies as $|\pi| \to 0$.

**Lemma 3.11.** *Suppose (H.1) and (H.2) hold. Let $(\pi_m)_{m \in \mathbb{N}} \subset \mathscr{P}_{[0,T]}$ be such that $\lim_{m \to \infty} |\pi_m| = 0$. Then $\lim_{m \to \infty} \mathcal{C}^\star_{\pi_m} = \inf_{\theta \in \Theta} \mathcal{C}(\theta)$.*

*Proof.* For each $m \in \mathbb{N}$, by $\Theta^{\pi_m} \subset \Theta$, $\mathcal{C}^\star_{\pi_m} = \inf_{\theta \in \Theta^{\pi_m}} \mathcal{C}(\theta) \geq \inf_{\theta \in \Theta} \mathcal{C}(\theta)$, which implies that $\liminf_{m \to \infty} \mathcal{C}^\star_{\pi_m} \geq \inf_{\theta \in \Theta} \mathcal{C}(\theta)$. On the other hand, let $\theta^\star = (K^\star, V^\star)$ be defined in (2.6), and for each $m \in \mathbb{N}$, let $\theta^{m,\star} = (K^{m,\star}, V^{m,\star})$ be the $L^2$ projection of $\theta^\star$ onto $\Theta^m$ such that $K_t^{m,\star} = \sum_{i=0}^{N_m-1} \overline{K}^\star_{t_i} \mathbb{1}_{[t_i,t_{i+1})}(t)$ and $V_t^{m,\star} = \sum_{i=0}^{N_m-1} \overline{V}^\star_{t_i} \mathbb{1}_{[t_i,t_{i+1})}(t)$ for a.e. $t \in [0,T]$, where

$$\overline{K}^\star_{t_i} = \frac{1}{t_{i+1} - t_i} \int_{t_i}^{t_{i+1}} K_t^\star \, \mathrm{d}t, \quad \overline{V}^\star_{t_i} = \frac{1}{t_{i+1} - t_i} \int_{t_i}^{t_{i+1}} V_t^\star \, \mathrm{d}t, \quad \forall i \in \{0, \ldots, N_m - 1\}.$$

A standard mollification argument shows that $\lim_{m \to \infty} \|\theta^{m,\star} - \theta^\star\|_{L^2} = 0$. Moreover, the fact that $\varepsilon I_k \preceq V^\star \preceq \frac{1}{\varepsilon} I_k$ for some $\varepsilon > 0$ implies that $\varepsilon I_k \preceq V^{m,\star} \preceq \frac{1}{\varepsilon} I_k$ for all $m \in \mathbb{N}$. By the uniform $L^2$-bound of $(K^{m,\star})_{m \in \mathbb{N}}$ and the $L^\infty$-bound of $(V^{m,\star})_{m \in \mathbb{N}}$, there exists $C \geq 0$ such that $\Sigma^{\theta^{m,\star}} \preceq C I_d$ for all $m \in \mathbb{N}$ due to Lemma 3.7. Then by Proposition 2.3, for all $m \in \mathbb{N}$,

$$\mathcal{C}(\theta^{m,\star}) - \mathcal{C}(\theta^\star) \leq \int_0^T \left( \frac{1}{2} \langle K_t^{m,\star} - K_t^\star, (D_t^\top P_t^{\theta^\star} D_t + R_t + \rho \bar{V}_t^{-1})(K_t^{m,\star} - K_t^\star) \Sigma_t^{\theta^{m,\star}} \rangle \right.$$
$$\left. + \frac{\rho}{4} \frac{|V_t^{m,\star} - V_t^\star|^2}{\min(\lambda_{\min}^2(V_t^\star), \lambda_{\min}^2(V_t^{m,\star}))} \right) \mathrm{d}t,$$

which along with $\lim_{m \to \infty} \|\theta^{m,\star} - \theta^\star\|_{L^2} = 0$ and $V^{m,\star} \succeq \varepsilon I_k$, $\Sigma^{\theta^{m,\star}} \preceq C I_d$ for all $m \in \mathbb{N}$ implies that $\lim_{m \to \infty} \mathcal{C}(\theta^{m,\star}) = \inf_{\theta \in \Theta} \mathcal{C}(\theta)$. As $\mathcal{C}^\star_{\pi_m} \leq \mathcal{C}(\theta^{m,\star})$ for all $m \in \mathbb{N}$,

$$\inf_{\theta \in \Theta} \mathcal{C}(\theta) \leq \lim_{m \to \infty} \inf \mathcal{C}^\star_{\pi_m} \leq \lim_{m \to \infty} \sup \mathcal{C}^\star_{\pi_m} \leq \lim_{m \to \infty} \sup \mathcal{C}(\theta^{m,\star}) = \inf_{\theta \in \Theta} \mathcal{C}(\theta).$$

This leads to the desired convergence result. $\qquad \square$

The following proposition proves that when the mesh size $|\pi|$ are sufficiently small, the policies from (2.21) have similar costs as those from (2.16).

**Proposition 3.12.** *Suppose (H.1), (H.2) and (H.3) hold. Assume further that $D \in C([0,T]; \mathbb{R}^{d \times k})$, $R \in C([0,T]; \mathbb{S}^k)$ and $\bar{V} \in C([0,T]; \mathbb{S}^k_+)$. Let $\theta^0 \in L^2(0,T; \mathbb{R}^{k \times d}) \times C([0,T]; \mathbb{S}^k_+)$, let $(\pi_m)_{m \in \mathbb{N}} \subset \mathscr{P}_{[0,T]}$ be such that $\lim_{m \to \infty} |\pi_m| = 0$, and let $(\theta^{\pi_m,0})_{m \in \mathbb{N}} \subset \Theta$ be such that $\theta^{\pi_m,0} \in \Theta^{\pi_m}$ for all $m \in \mathbb{N}$, $\lim_{m \to \infty} \|\theta^{\pi_m,0} - \theta^0\|_{L^2 \times L^\infty} = 0$. Let $\bar{\lambda}_0 > 0$ be such that $\bar{\lambda}_0 I_k \succeq D^\top P^{\theta^0} D + R + \rho \bar{V}^{-1}$, with $P^{\theta^0} \in C([0,T]; \mathbb{S}^d)$ defined in (2.10), and for each $\tau > 0$, let $(\theta^n)_{n \in \mathbb{N}}$ and $(\theta^{\pi_m,n})_{m,n \in \mathbb{N}}$ be defined in (2.16) and (2.21), respectively. Then for all $\tau \in (0, 1/\bar{\lambda}_0]$ and $N \in \mathbb{N}_0$,*

$$\lim_{m \to \infty} \sup_{n=0,\ldots,N} |\mathcal{C}(\theta^{\pi_m,n}) - \mathcal{C}(\theta^n)| = 0.$$

*Proof.* For each $L > 0$, define $\Theta_L = \left\{ \theta = (K,V) \in \Theta \,\middle|\, \frac{1}{L} I_k \preceq V \preceq L I_k \right\}$. Let $\tau \in (0, 1/\bar{\lambda}_0]$ be fixed. By Proposition 3.6, there exists $\underline{\lambda}_V, \bar{\lambda}_V > 0$ such that $\underline{\lambda}_V I_k \preceq V^n \preceq \bar{\lambda}_V I_k$ for all $n \in \mathbb{N}_0$. Moreover, by the continuity of $D$, $R$ and $\bar{V}$, and the expressions (2.13) and (2.16), a straightforward induction argument shows that $V^n \in C([0,T]; \mathbb{S}^k_+)$ for all $n \in \mathbb{N}_0$.

We first prove by induction that for all $n \in \mathbb{N}_0$, there exists $L > 0, m_0 \in \mathbb{N}$ such that

$$\lim_{m \to \infty} \|\theta^{\pi_m,n} - \theta^n\|_{L^2 \times L^\infty} = 0, \text{ and } \theta^{\pi_m,n} \in \Theta_L \cap \Theta^{\pi_m}, \ \forall m \geq m_0. \tag{3.16}$$

Note that as $\theta^0 \in \Theta$ and $\lim_{m \to \infty} \|V^{\pi_m,0} - V^0\|_{L^\infty} = 0$, there exists $L > 0$ such that $\frac{1}{L} I_k \preceq V^{\pi_m,0} \preceq L I_k$ for all large $m \in \mathbb{N}$. This proves (3.16) for $n = 0$. Now suppose that the induction statement (3.16) holds for some $n \in \mathbb{N}_0$. As $V^n \in C([0,T]; \mathbb{S}_+^k)$, by (2.16) and (H.3), the triangle inequality shows that $\lim_{m \to \infty} \|\theta^{\pi_m,n+1} - \theta^{n+1}\|_{L^2 \times L^\infty} = 0$, which subsequently implies that there exists $L > 0$ such that $\frac{1}{L} I_k \preceq V^{\pi_m,n+1} \preceq L I_k$ for all sufficiently large $m$. This proves the statement (3.16) for $n + 1$.

By (3.16), for each $n \in \mathbb{N}$, $\sup_{m \in \mathbb{N}} \|K^{\pi_m,n}\|_{L^2} < \infty$ and $\limsup_{m \in \mathbb{N}} \|V^{\pi_m,n}\|_{L^\infty} < \infty$. Thus by Lemma 3.7, there exists $C \geq 0$ such that $0 \preceq \Sigma^{\theta^{\pi_m,n}} \preceq C I_d$ for all $m \in \mathbb{N}$. Then $\lim_{m \to \infty} |\mathcal{C}(\theta^{\pi_m,n}) - \mathcal{C}(\theta^n)| = 0$ follows from Proposition 2.3 and $\lim_{m \to \infty} \|\theta^{\pi_m,n} - \theta^n\|_{L^2 \times L^\infty} = 0$. This implies the desired convergence result for any given $N \in \mathbb{N}$. $\qquad \square$

*Proof of Theorem 2.7.* Let $C^\star = \inf_{\theta \in \Theta} \mathcal{C}(\theta) = C(\theta^\star)$, and for each $\tau > 0$ and $m \in \mathbb{N}$, let $(\theta^n)_{n \in \mathbb{N}}$ and $(\theta^{\pi_m,n})_{n \in \mathbb{N}}$ be defined by (2.16) and (2.21) with stepsize $\tau$, respectively. Then by Theorem 2.6 and Proposition 3.12, there exists $\tau_0 > 0$ such that for all $\tau \in (0, \tau_0]$ and $n \in \mathbb{N}_0$, $\mathcal{C}(\theta^{n+1}) \leq \mathcal{C}(\theta^n)$, $\mathcal{C}(\theta^{n+1}) - C^\star \leq \eta(\mathcal{C}(\theta^n) - C^\star)$ for some $\eta \in [0,1)$ (independent of $n$), and $\lim_{m \to \infty} |\mathcal{C}(\theta^{\pi_m,n}) - \mathcal{C}(\theta^n)| = 0$. Moreover, for all $\varepsilon > 0$, $N(\varepsilon) = \frac{\widetilde{C}}{\tau} \log(\frac{\widetilde{C}}{\varepsilon})$ for some $\widetilde{C} > 0$ independent of $\tau$ and $\varepsilon$.

We first prove for all $\tau \in (0, \tau_0]$ and all $\varepsilon, \gamma > 0$, there exists $\mathfrak{m}_{\varepsilon,\gamma} \in \mathbb{N}$ such that for all $m \geq \mathfrak{m}_{\varepsilon,\gamma}$,

$$N(\varepsilon + \gamma) \leq N^{\pi_m}(\varepsilon) \leq N(\varepsilon). \tag{3.17}$$

To prove $N^{\pi_m}(\varepsilon) \leq N(\varepsilon)$, by Lemma 3.11 and the choice of $\tau_0$, for all $n \in \mathbb{N}_0$, $\lim_{m \to \infty} (\mathcal{C}(\theta^{\pi_m,n}) - \mathcal{C}^\star_{\pi_m}) = \mathcal{C}(\theta^n) - C^\star$. Hence, for all $\varepsilon > 0$ and $n \in \mathbb{N}_0$, if $\mathcal{C}(\theta^n) - C^\star < \varepsilon$, then there exists $\mathfrak{m}_\varepsilon \in \mathbb{N}$ such that for all $m \geq \mathfrak{m}_\varepsilon$, $\mathcal{C}(K^{\pi_m,n}) - \mathcal{C}^\star_{\pi_m} < \varepsilon$, which implies $N^{\pi_m}(\varepsilon) \leq N(\varepsilon)$ for all $m \geq \mathfrak{m}_\varepsilon$. We then prove $N(\varepsilon + \gamma) \leq N^{\pi_m}(\varepsilon)$ with a given $\gamma > 0$. The convergence of $(\mathcal{C}(\theta^n))_{n \in \mathbb{N}}$ implies that $N(\varepsilon + \gamma) \in \mathbb{N}_0$, which along with Lemma 3.11 and Proposition 3.12 shows that

$$\lim_{m \to \infty} \max_{0 \leq n \leq N(\varepsilon + \gamma)} \left| (\mathcal{C}(\theta^{\pi_m,n}) - \mathcal{C}^\star_{\pi_m}) - (\mathcal{C}(\theta^n) - C^\star) \right| = 0. \tag{3.18}$$

The definition of $N(\varepsilon + \gamma)$ implies that $\mathcal{C}(\theta^n) - C^\star \geq \varepsilon + \gamma$ for all $n < N(\varepsilon + \gamma)$. Moreover, by (3.18), there exists $\mathfrak{m}_\gamma \in \mathbb{N}$ such that for all $m \geq \mathfrak{m}_\gamma$,

$$\max_{0 \leq n < N(\varepsilon + \gamma)} \left| (\mathcal{C}(\theta^{\pi_m,n}) - \mathcal{C}^\star_{\pi_m}) - (\mathcal{C}(\theta^n) - C^\star) \right| \leq \gamma.$$

Hence for all $m \geq \mathfrak{m}_\gamma$ and $n < N(\varepsilon + \gamma)$,

$$
\begin{aligned}
\mathcal{C}(\theta^{\pi_m,n}) - \mathcal{C}^\star_{\pi_m} &= (\mathcal{C}(\theta^{\pi_m,n}) - \mathcal{C}^\star_{\pi_m}) - (\mathcal{C}(\theta^n) - C^\star) + (\mathcal{C}(\theta^n) - C^\star) \\
&\geq (\mathcal{C}(\theta^n) - C^\star) - \left| (\mathcal{C}(\theta^{\pi_m,n}) - \mathcal{C}^\star_{\pi_m}) - (\mathcal{C}(\theta^n) - C^\star) \right| \\
&\geq (\mathcal{C}(\theta^n) - C^\star) - \max_{0 \leq n < N(\varepsilon + \gamma)} \left| (\mathcal{C}(\theta^{\pi_m,n}) - \mathcal{C}^\star_{\pi_m}) - (\mathcal{C}(\theta^n) - C^\star) \right| \geq \varepsilon.
\end{aligned}
$$

This implies that $N^{\pi_m}(\varepsilon) \geq N(\varepsilon + \gamma)$ for all $m \geq \mathfrak{m}_\gamma$. Taking $\mathfrak{m}_{\varepsilon,\gamma} = \max(\mathfrak{m}_\varepsilon, \mathfrak{m}_\gamma)$ completes the proof of (3.17).

Now we are ready to establish (2.24) for fixed $\tau \in (0, \tau_0]$ and $\varepsilon > 0$. By the choice of $\tau_0$, there exists $\eta \in [0,1)$, independent of $\varepsilon$, such that for all $n \in \mathbb{N}_0$, $\mathcal{C}(\theta^{n+1}) - C^\star \leq \eta(\mathcal{C}(\theta^n) - C^\star)$. Then, by the definition of $N(\varepsilon)$, $\mathcal{C}(\theta^n) - C^\star \geq \varepsilon$ for all $n < N(\varepsilon)$, which yields the estimate

$$\eta^{N(\varepsilon)-1-n}(\mathcal{C}(\theta^n) - C^\star) \geq \mathcal{C}(\theta^{N(\varepsilon)-1}) - C^\star \geq \varepsilon, \quad \forall n < N(\varepsilon) - 1.$$

26

This implies that $\mathcal{C}(\theta^n) - \mathcal{C}^\star \geq \frac{\varepsilon}{\eta} > \varepsilon$ for all $n < N(\varepsilon) - 1$. Now let $\gamma_\varepsilon := \min\{\mathcal{C}(\theta^n) - \mathcal{C}^\star - \varepsilon \mid n < N(\varepsilon) - 1\}$. Note that $\gamma_\varepsilon > 0$ as $N(\varepsilon) < \infty$. By the definition of $\gamma_\varepsilon$, for all $n < N(\varepsilon) - 1$, $\mathcal{C}(\theta^n) - \mathcal{C}^\star \geq \varepsilon + \gamma_\varepsilon$, which implies that $N(\varepsilon + \gamma_\varepsilon) \geq N(\varepsilon) - 1$. Hence, by (3.17), there exists $\mathfrak{m}_\varepsilon \in \mathbb{N}$ such that

$$N(\varepsilon) - 1 \leq N(\varepsilon + \gamma_\varepsilon) \leq N^{\pi_m}(\varepsilon) \leq N(\varepsilon), \quad \forall m \geq \mathfrak{m}_\varepsilon.$$

This proves the desired estimate (2.24). $\qquad\square$

*Proof of Corollary 2.8.* By Proposition 2.1 and (2.25), for all $\pi \in \mathscr{P}_{[0,T]}$, $\theta \in \Theta^\pi$ and $i \in \{0, \ldots, N-1\}$,

$$\nabla_{K_i} \mathcal{C}(\theta) = \int_{t_i}^{t_{i+1}} \mathcal{D}_K(\theta)_t \Sigma_t^\theta \, dt, \quad \nabla_{V_i} \mathcal{C}(\theta) = \int_{t_i}^{t_{i+1}} \mathcal{D}_V(\theta)_t \, dt,$$

where $\mathcal{D}_K(\theta)$ and $\mathcal{D}_V(\theta)$ are defined by (2.12) and (2.13), respectively. Hence $(\mathcal{D}_K^\pi, \mathcal{D}_V^\pi) : \Theta^\pi \to \Theta^\pi$ in (2.26) satisfies for all $\theta \in \Theta^\pi$, and a.e. $t \in [0, T]$,

$$\begin{aligned}
\mathcal{D}_K^\pi(\theta)_t &= \sum_{i=0}^{N-1} \left( \frac{1}{t_{i+1} - t_i} \int_{t_i}^{t_{i+1}} \mathcal{D}_K(\theta)_t \Sigma_t^\theta \, dt \right) \left( \Sigma_{t_i}^\theta \right)^{-1} \mathbb{1}_{[t_i, t_{i+1})}(t), \\
\mathcal{D}_V^\pi(\theta)_t &= \sum_{i=0}^{N-1} \left( \frac{1}{t_{i+1} - t_i} \int_{t_i}^{t_{i+1}} (V_{t_i} \mathcal{D}_V(\theta)_t + \mathcal{D}_V(\theta)_t V_{t_i}) \, dt \right) \mathbb{1}_{[t_i, t_{i+1})}(t).
\end{aligned} \tag{3.19}$$

To simplify the notation, for each Euclidean space $E$, let $\mathcal{PC}_\pi(E)$ be the space of piecewise constant functions $f : [0, T] \to E$ on $\pi$, let $\Pi^\pi : L^2(0, T; E) \to \mathcal{PC}_\pi(E)$ be such that for all $f \in L^2(0, T; E)$, $\Pi^\pi(f)_t := \sum_{i=0}^{N-1} \left( \frac{1}{t_{i+1} - t_i} \int_{t_i}^{t_{i+1}} f_t \, dt \right) \mathbb{1}_{[t_i, t_{i+1})}(t)$ for all $t \in [0, T]$, and let $\mathcal{T}^\pi : C([0, T]; E) \to \mathcal{PC}_\pi(E)$ be such that for all $f \in C([0, T]; E)$, $\mathcal{T}^\pi(f)_t := \sum_{i=0}^{N-1} f_{t_i} \mathbb{1}_{[t_i, t_{i+1})}(t)$ for all $t \in [0, T]$. Note that $\Pi^{\pi_m}$ is the orthogonal projection with respect to the $\|\cdot\|_{L^2}$ norm, and hence is 1-Lipschitz continuous with respect to the $\|\cdot\|_{L^2}$ norm. Moreover, by (3.19), for all $\theta \in \Theta^\pi$,

$$\mathcal{D}_K^\pi(\theta) = \Pi^\pi \left( \mathcal{D}_K(\theta) \Sigma^\theta \left( \mathcal{T}^\pi(\Sigma^\theta) \right)^{-1} \right), \quad \mathcal{D}_V^\pi(\theta) = \Pi^\pi \left( \mathcal{T}^\pi(V) \mathcal{D}_V(\theta) + \mathcal{D}_V(\theta) \mathcal{T}^\pi(V) \right). \tag{3.20}$$

The definition of $(\mathcal{D}_K^\pi, \mathcal{D}_V^\pi)$ in (3.20) can be naturally extended to all $\theta \in L^2(0, T; \mathbb{R}^{k \times d}) \times C([0, T]; \mathbb{S}_+^k)$. Note that $\Sigma^\theta$ is pointwise invertible due to $\mathbb{E}[\xi_0 \xi_0^\top] \succ 0$ (see Lemma 3.7).

We are now ready to verify (H.3) for (3.20). Let $\theta \in L^2(0, T; \mathbb{R}^{k \times d}) \times C([0, T]; \mathbb{S}_+^k)$, $(\pi_m)_{m \in \mathbb{N}} \subset \mathscr{P}_{[0,T]}$ be such that $\lim_{m \to \infty} |\pi_m| = 0$, and $(\theta^m)_{m \in \mathbb{N}} \subset \Theta$ be such that $\theta^m \in \Theta^{\pi_m}$ for all $m \in \mathbb{N}$ and $\lim_{m \to \infty} \|\theta^m - \theta\|_{L^2 \times L^\infty} = 0$. Then for all $m \in \mathbb{N}$, by the Lipschitz continuity of $\Pi^{\pi_m}$,

$$\begin{aligned}
&\|\mathcal{D}_K^{\pi_m}(\theta^m) - \mathcal{D}_K(\theta)\|_{L^2} \\
&\leq \|\mathcal{D}_K^{\pi_m}(\theta^m) - \Pi^{\pi_m}(\mathcal{D}_K(\theta))\|_{L^2} + \|\Pi^{\pi_m}(\mathcal{D}_K(\theta)) - \mathcal{D}_K(\theta)\|_{L^2} \\
&\leq \left\| \mathcal{D}_K(\theta^m) \Sigma^{\theta^m} \left( \mathcal{T}^{\pi_m}(\Sigma^{\theta^m}) \right)^{-1} - \mathcal{D}_K(\theta) \right\|_{L^2} + \|\Pi^{\pi_m}(\mathcal{D}_K(\theta)) - \mathcal{D}_K(\theta)\|_{L^2}. \tag{3.21}
\end{aligned}$$

The density of $(\mathcal{PC}_\pi(\mathbb{R}^{k \times d}))_{m \in \mathbb{N}}$ in $L^2(0, T; \mathbb{R}^{k \times d})$ shows that the second term of (3.21) tends to zero as $m \to \infty$. Standard stability results of (2.10) and (2.11) (see, e.g., Lemma 3.4) show that $\lim_{m \to \infty} \|P^{\theta^m} - P^\theta\|_{L^\infty} = 0$ and $\lim_{m \to \infty} \|\Sigma^{\theta^m} - \Sigma^\theta\|_{L^\infty} = 0$. Thus by (H.1) and (2.12), $\lim_{m \to \infty} \|\mathcal{D}_K(\theta^m) - \mathcal{D}_K(\theta)\|_{L^2} = 0$. Moreover, as $\inf_{t \in [0,T]} \lambda_{\min}(\Sigma_t^\theta) > 0$ (see Lemma 3.7),

$\Sigma^{\theta^m}\left(\mathcal{T}^{\pi_m}\left(\Sigma^{\theta^m}\right)\right)^{-1}$ tends to the identity function in $L^\infty$ as $m \to \infty$. Consequently, the first term of (3.21) tends to zero as $m \to \infty$, which proves $\lim_{m\to\infty} \|\mathcal{D}_K^{\pi_m}(\theta^m) - \mathcal{D}_K(\theta)\|_{L^2} = 0$.

We then prove the convergence of $(\mathcal{D}_V^{\pi_m}(\theta^m))_{m\in\mathbb{N}}$. Note that for each $m \in \mathbb{N}$ and Euclidean space $E$, $\|\Pi^{\pi_m}(f)\|_{L^\infty} \le \|f\|_{L^\infty}$ if $f \in L^\infty(0,T;E)$, and $\lim_{m\to\infty} \|\Pi^{\pi_m}(f) - f\|_{L^\infty} = 0$ if $f \in C([0,T];E)$. The same property also holds for the operator $\mathcal{T}^{\pi_m}$. Then for all $m \in \mathbb{N}$,

$$
\begin{aligned}
&\|\mathcal{D}_V^{\pi_m}(\theta^m) - \mathcal{D}^{\mathrm{bw}}(\theta)\|_{L^\infty} \\
&\qquad \le \|\mathcal{D}_V^{\pi_m}(\theta^m) - \Pi^{\pi_m}(\mathcal{D}_V^{\mathrm{bw}}(\theta))\|_{L^\infty} + \|\Pi^{\pi_m}(\mathcal{D}_V^{\mathrm{bw}}(\theta)) - \mathcal{D}_V^{\mathrm{bw}}(\theta)\|_{L^\infty}.
\end{aligned}
\tag{3.22}
$$

By the continuity of $D$, $R$, $\bar{V}$ and $V$, $\mathcal{D}_V^{\mathrm{bw}}(\theta) \in C([0,T];\mathbb{S}_+^k)$ (cf. (2.13) and (2.17)), and hence the second term in (3.22) tends to zero as $m \to \infty$. To show the first term tends to zero, by (2.17) and (3.20), it suffices to prove $\lim_{m\to\infty} \|\mathcal{D}_V(\theta^m) - \mathcal{D}_V(\theta)\|_{L^\infty} = 0$. This follows directly from the facts that $\lim_{m\to\infty} \|P^{\theta^m} - P^\theta\|_{L^\infty} = 0$, $\lim_{m\to\infty} \|V^m - V\|_{L^\infty} = 0$ and $V \in C([0,T];\mathbb{S}_+^k)$. This verifies (H.3) for (3.20). $\qquad\square$

# 4   Numerical experiments

In this section, we test the theoretical findings through a numerical experiment on an exploratory LQC problem arising from mean-variance portfolio selection. Our experiments confirm that the proposed iteration (2.26) converges linearly to the optimal policy. They also show that conventional PG methods exhibit a degraded performance for small timesteps in the policy updates, while our algorithm demonstrates robustness across different step sizes.

**Problem setup.**   We minimise the following cost $\mathcal{C} : \Theta \to \mathbb{R}$ (cf. (2.9)):

$$
\mathcal{C}(\theta) = \mathbb{E}\left[\frac{1}{2}\mu\,(X_T^\theta)^2 + \rho \int_0^T \mathcal{H}(\nu_t^\theta(X_t^\theta)\|\overline{\mathfrak{m}}_t)\,\mathrm{d}t\right],
\tag{4.1}
$$

where $\overline{\mathfrak{m}}_t = \mathcal{N}(0,\bar{V})$ with $\bar{V} \in \mathbb{S}_+^3$, and for each $\theta \in \Theta$, $X^\theta \in \mathcal{S}^2(0,T;\mathbb{R})$ satisfies for all $t \in [0,T]$,

$$
\mathrm{d}X_t = \int_{\mathbb{R}^3}\left(B_t a\,\nu_t^\theta(X_t;\mathrm{d}a)\right)\mathrm{d}t + \left(\int_{\mathbb{R}^3}\sum_{j=1}^3\left(D^{(j)}a\right)^2\nu_t^\theta(X_t;\mathrm{d}a)\right)^{\frac{1}{2}}\mathrm{d}W_t, \quad X_0 = \xi_0,
\tag{4.2}
$$

for some $B : [0,T] \to \mathbb{R}^{1\times 3}$ and $D^{(j)} \in \mathbb{R}^{1\times 3}$, $j = 1,2,3$. The coefficients are chosen as follows: $T = 1$, $\mu = 0.5$, $\rho = 0.01$, $\bar{V} = 0.1I_3$, $\xi_0 \sim \mathcal{N}(0.5,0.01)$, $B_t = (0.4,0.8,0.4) + 0.2\sin(2\pi t)\mathbf{1}_3$ for all $t \in [0,T]$, and $D = \begin{pmatrix} D^{(1)} \\ D^{(2)} \\ D^{(3)} \end{pmatrix}$ with $D^\top D = \begin{pmatrix} 0.5 & 0.25 & -0.125 \\ 0.25 & 1 & -0.25 \\ -0.125 & -0.25 & 0.5 \end{pmatrix}$. Note that $D^\top D \in \mathbb{S}_+^3$, and hence (H.2) holds for all $\rho \ge 0$ (see [38] and Remark 2.2).

The problem (4.1)-(4.2) arises from an exploratory mean-variance portfolio selection problem, where the agent allocates their wealth among three risky assets by sampling from the policy $\nu^\theta$ (see [32]). Indeed, as illustrated at the end of Section 2.4, for each $\theta = (K,V) \in \Theta$, $\mathcal{C}(\theta)$ can be approximated by replacing (4.2) with the following dynamics: $X_0 = \xi_0$, and for all $t \in [0,T]$,

$$
\mathrm{d}X_t = B_t\left(K_tX_t + V_t^{\frac{1}{2}}\xi_t\right)\mathrm{d}t + \sum_{j=1}^3 D^{(j)}\left(K_tX_t + V_t^{\frac{1}{2}}\xi_t\right)\mathrm{d}W_t^{(j)}
\tag{4.3}
$$

with $\xi_t = \sum_{i=1}^n \zeta_i \mathbb{1}_{[t_i,t_{i+1})}(t)$, where $(W^{(j)})_{j=1}^3$ are independent Brownian motions, $(\zeta_i)_{i=1}^n$ are independent standard normal random vectors, and $(t_i)_{i=1}^n$ is a sufficiently fine time mesh.

**Linear convergence.** We first implement (2.26) on the uniform time mesh $\pi_c$ with mesh size $1/128$, and examine its convergence. The scheme is initialised with $K^0 \equiv (1/3, 1/3, 1/3)$ and $V^0 \equiv 0.1 D^\top D$. For each $n \in \mathbb{N}_0$, given $\theta^n \subset \Theta^{\pi_c}$, we simulate $10^5$ independent trajectories of (4.3) (with $\theta = \theta^n$) using the Euler–Maruyama method on the mesh $\pi_c$, evaluate the approximate value $\widehat{\mathcal{C}}(\theta^n)$ and state covariance $\widehat{\Sigma}^n$ using the empirical distribution of these sample paths, and compute an approximate gradient $(\widehat{\nabla_{\theta_i^n}\mathcal{C}})_{i=0}^{127}$ using automatic differentiation. The iterate $\theta^n$ is updated by (2.26) with $\widehat{\Sigma}^n$, $\widehat{\nabla_{\theta^n}\mathcal{C}}$ and the stepsize $\tau = 0.01$. The performance of the scheme is measured by the errors $(\widehat{\mathcal{C}}(\theta^n) - \mathcal{C}^\star)_{n \in \mathbb{N}_0}$, where $\mathcal{C}^\star$ is the optimal cost of (4.1) obtained by Riccati equations. Further implementation details are given in Appendix B.

Figure 1 (left) exhibits the decay of $(\widehat{\mathcal{C}}(\theta^n) - \mathcal{C}^\star)_{n \in \mathbb{N}_0}$ with respect to the number of iterations, where the solid line and the shaded area indicate the sample mean and the spread over 10 repeated experiments, respectively. It clearly shows the linear convergence of (2.26), as indicated in Theorems 2.6 and 2.7. The seemingly higher noise for larger iteration numbers results from the small errors in this case, so that the fluctuations appear larger on the log scale. The variance could be reduced by increasing the number of samples.

**Robustness in action frequency.** We then compare the performance of (2.26) with a standard PG method for different policy discretisation timescales. The former (termed "scaled PG") scales the gradients with the discretisation mesh size, while the latter (termed "unscaled PG") updates the policy with unscaled gradients. More precisely, let $\theta^0 = (K^0, V^0)$ be a fixed initial guess given as above, and $\pi_m = \{i\frac{1}{m}\}_{i=0}^m$, $m \in \{8, 16, 32, 64, 128\}$ be a family of time meshes. For each $m \in \{8, 16, 32, 64, 128\}$, the scaled PG method generates the iterates $(\theta^{\pi_m, n})_{n \in \mathbb{N}_0} \subset \Theta^{\pi_m}$ according to (2.26) with $\tau = 0.01$ and $\Delta_i = 1/m$, where the required gradients for each iteration are computed as above. The unscaled PG method follows (2.26) with $\tau = 0.08$ and $\Delta_i = 1$ for all $m$. Here, a larger stepsize has been adopted for the unscaled PG method so that the two algorithms coincide for the coarsest mesh $\pi_8$.

Figure 1 (right) compares, for different discretisation timescales, the numbers of required iterations $N^{\pi_m}(0.01)$ for both schemes to achieve an accuracy of $\epsilon = 0.01$ (cf. (2.22)). One can observe clearly that the number of required iterations for the unscaled PG method exhibits a linear growth in the number of action time points. In constrast, the number of iterations for the scaled PG method remains constant for all meshes. This confirms the theoretical result in Theorem 2.7, and shows that the scaled PG method outperforms conventional PG methods for fine meshes.
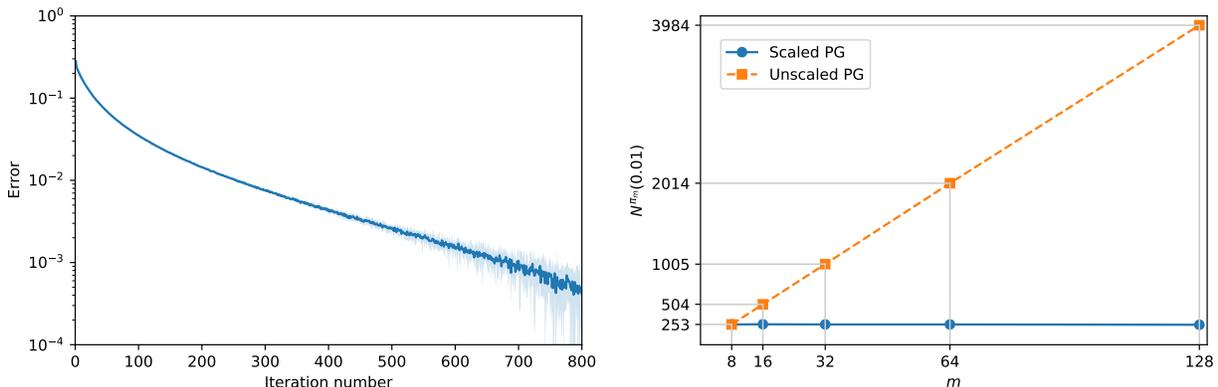


Figure 1: Convergence and robustness of the PG method (2.26).

# A Proofs of technical results

The following lemma establishes the well-posedness of stochastic differential equations, whose coefficients are Lipschitz continuous in state with time-dependent Lipschitz constants. The proof follows essentially the lines of Theorems 3.2.2 and 3.3.1 (Method 2) in [35], and hence is omitted.

**Lemma A.1.** *Let $T > 0$, $(\Omega, \mathcal{F}, \mathbb{F}, \mathbb{P})$ be a filtered probability space satisfying the usual condition, $b : \Omega \times [0, T] \times \mathbb{R}^d \to \mathbb{R}^d$ and $\sigma : \Omega \times [0, T] \times \mathbb{R}^d \to \mathbb{R}^{d \times d}$ be progressively measurable functions such that $b_\cdot(\cdot, 0) \in L^1(\Omega \times [0, T]; \mathbb{R}^d)$ and $\sigma_\cdot(\cdot, 0) \in L^2(\Omega \times [0, T]; \mathbb{R}^{d \times d})$. Assume that there exists $A \in L^1(0, T; \mathbb{R})$ and $C \in L^2(0, T; \mathbb{R})$ such that for all $(\omega, t) \in \Omega \times [0, T]$ and $x, x' \in \mathbb{R}^d$, $|b_t(\omega, x) - b_t(\omega, x')| \le |A_t||x - x'|$ and $|\sigma_t(\omega, x) - \sigma_t(\omega, x')| \le |C_t||x - x'|$. Then for all $\xi_0 \in L^2(\mathcal{F}_0; \mathbb{R}^d)$, there exists a unique strong solution $X \in \mathcal{S}^2(0, T; \mathbb{R}^d)$ to the following equation*

$$\mathrm{d}X_t = b_t(X_t)\,\mathrm{d}s + \sigma_t(X_t)\,\mathrm{d}W_t, \quad t \in [0, T]; \quad X_0 = \xi_0. \tag{A.1}$$

**Proposition A.2.** *Suppose (H.1(1)) holds. Then*

*(1) for all $\mathfrak{m} \in \mathcal{A}$, (2.1) admits a unique strong solution $X^{\mathfrak{m}} \in \mathcal{S}^2(0, T; \mathbb{R}^d)$.*

*(2) for all $\nu^\theta \in \mathcal{V}$, (2.8) admits a unique strong solution $X^\theta \in \mathcal{S}^2(0, T; \mathbb{R}^d)$.*

*Proof.* Let $E = [0, T] \times \mathbb{R}^k$. We verify that the coefficients of (2.1) and (2.8) satisfy the conditions of Lemma A.1.

To prove Item (1), let $\mathfrak{m} \in \mathcal{A}$ be given, and define $\Phi^{\mathfrak{m}} : \Omega \times [0, T] \times \mathbb{R}^d \to \mathbb{R}^d$ and $\Gamma^{\mathfrak{m}} : \Omega \times [0, T] \times \mathbb{R}^d \to \overline{\mathbb{S}_+^d}$ such that for all $(\omega, t, x) \in \Omega \times [0, T] \times \mathbb{R}^d$, $\Phi_t^{\mathfrak{m}}(\omega, x) = \Phi_t(x, \mathfrak{m}_t(\omega))$ and $\Gamma_t^{\mathfrak{m}}(\omega, x) = \Gamma_t(x, \mathfrak{m}_t(\omega))$, with $\Phi$ and $\Gamma$ defined in (2.2). By Fubini's theorem and Hölder's inequality,

$$\mathbb{E}\left[\int_0^T |\Phi_t^{\mathfrak{m}}(\cdot, 0)|\right]\,\mathrm{d}t \le \int_0^T \left(\mathbb{E}\left[\int_{\mathbb{R}^k} |a|\,\mathfrak{m}_t(\mathrm{d}a)\right]|B_t|\right)\,\mathrm{d}t \le \|B\|_{L^2}\left(\mathbb{E}\left[\int_E |a|^2\,\mathfrak{m}_t(\mathrm{d}t, \mathrm{d}a)\right]\right)^{\frac{1}{2}} < \infty,$$

$$\mathbb{E}\left[\int_0^T |\Gamma_t^{\mathfrak{m}}(\cdot, 0)|^2\,\mathrm{d}t\right] \le \widetilde{C}\|D\|_{L^\infty}^2 \mathbb{E}\left[\int_E |a|^2\,\mathfrak{m}_t(\mathrm{d}t, \mathrm{d}a)\right] < \infty.$$

For all $(\omega, t) \in \Omega \times [0, T]$, using $\mathfrak{m}_t(\omega) \in \mathcal{P}(\mathbb{R}^k)$, $|\Phi_t^{\mathfrak{m}}(\omega, x) - \Phi_t^{\mathfrak{m}}(\omega, x')| \le |A_t||x - x'|$ for all $x, x' \in \mathbb{R}^d$. To prove the Lipschitz continuity of $\Gamma^{\mathfrak{m}}$, observe that for all $(t, x, m) \in [0, T] \times \mathbb{R}^d \times \mathcal{P}_2(\mathbb{R}^k)$, $\Gamma_t(x, m) = (M_{m,t}(x)M_{m,t}(x)^\top + N_{m,t}N_{m,t}^\top)^{1/2}$, where $M_{m,t}(x) := C_t x + D_t \int_{\mathbb{R}^k} a\, m(\mathrm{d}a)$ and $N_{m,t} := D_t \left(\int_{\mathbb{R}^k} aa^\top m(\mathrm{d}a)\right)^{1/2}$. This implies that

$$\begin{pmatrix} \Gamma_t(x, m) & 0_{d \times d} \\ 0_{d \times d} & 0_{d \times d} \end{pmatrix} = \begin{pmatrix} M_{m,t}(x)M_{m,t}(x)^\top + N_{m,t}N_{m,t}^\top & 0_{d \times d} \\ 0_{d \times d} & 0_{d \times d} \end{pmatrix}^{\frac{1}{2}} = \left|\begin{pmatrix} M_{m,t}(x) & N_{m,t} \\ 0_{d \times 1} & 0_{d \times 1} \end{pmatrix}\right|_{\mathrm{mat}},$$

where $0_{m \times n}$ is $m \times n$ zero matrix, and $|\cdot|_{\mathrm{mat}}$ is the matrix absolute value defined by $|M|_{\mathrm{mat}} = (MM^\top)^{1/2}$ for any matrix $M$. Then, for all $(t, m) \in [0, T] \times \mathcal{P}_2(\mathbb{R}^k)$ and $x, x' \in \mathbb{R}^d$,

$$\begin{aligned}
&|\Gamma_t(x, m) - \Gamma_t(x', m)| \\
&= \left|\begin{pmatrix} \Gamma_t(x, m) - \Gamma_t(x', m) & 0_{d \times d} \\ 0_{d \times d} & 0_{d \times d} \end{pmatrix}\right| = \left|\begin{pmatrix} \Gamma_t(x, m) & 0_{d \times d} \\ 0_{d \times d} & 0_{d \times d} \end{pmatrix} - \begin{pmatrix} \Gamma_t(x', m) & 0_{d \times d} \\ 0_{d \times d} & 0_{d \times d} \end{pmatrix}\right| \\
&= \left|\left|\begin{pmatrix} M_{m,t}(x) & N_{m,t} \\ 0_{d \times 1} & 0_{d \times 1} \end{pmatrix}\right|_{\mathrm{mat}} - \left|\begin{pmatrix} M_{m,t}(x') & N_{m,t} \\ 0_{d \times 1} & 0_{d \times 1} \end{pmatrix}\right|_{\mathrm{mat}}\right| \\
&\le \sqrt{2}\left|\begin{pmatrix} M_{m,t}(x) & N_{m,t} \\ 0_{d \times 1} & 0_{d \times 1} \end{pmatrix} - \begin{pmatrix} M_{m,t}(x') & N_{m,t} \\ 0_{d \times 1} & 0_{d \times 1} \end{pmatrix}\right| = \sqrt{2}|M_{m,t}(x) - M_{m,t}(x')|,
\end{aligned} \tag{A.2}$$

30

where the last inequality used the Lipschitz continuity of the matrix absolute value $|\cdot|_{\mathrm{mat}}$ (see [1]). Therefore, by the definition of $M_{m,t}(x)$, for all $(\omega, t) \in \Omega \times [0, T]$ and $x, x' \in \mathbb{R}^d$,

$$|\Gamma_t^{\mathfrak{m}}(\omega, x) - \Gamma_t^{\mathfrak{m}}(\omega, x')| \leq \sqrt{2}|M_{\mathfrak{m}_t(\omega),t}(x) - M_{\mathfrak{m}_t(\omega),t}(x')| \leq \sqrt{2}|C_t||x - x'|.$$

As $A \in L^1(0, T; \mathbb{R}^{d \times d})$ and $C \in L^2(0, T; \mathbb{R}^{d \times d})$, the coefficients of (2.1) satisfy the conditions of Lemma A.1, which subsequently implies the well-posedness of (2.1).

To prove Item (2), let $\nu^\theta \in \mathcal{V}$ be given, and define $\Phi^\theta : [0, T] \times \mathbb{R}^d \to \mathbb{R}^d$ and $\Gamma^\theta : [0, T] \times \mathbb{R}^d \to \overline{\mathbb{S}_+^d}$ such that for all $(t, x) \in [0, T] \times \mathbb{R}^d$, $\Phi_t^\theta(x) = \Phi_t(x, \nu_t^\theta(x))$ and $\Gamma_t^\theta(x) = \Gamma_t(x, \nu_t^\theta(x))$, with $\Phi$ and $\Gamma$ defined in (2.2). Then by Lemma 3.1, $\Phi_{\cdot}^\theta(0) = 0$ and $\Gamma_{\cdot}^\theta(0) = DV^{\frac{1}{2}} \in L^2(0, T; \mathbb{R}^{d \times d})$. Moreover, for all $(t, x) \in [0, T] \times \mathbb{R}^d$, $|\Phi_t^\theta(x) - \Phi_t^\theta(x')| \leq (|A_t| + |B_t K_t|)|x - x'|$ and by (A.2), $|\Gamma_t^\theta(x) - \Gamma_t^\theta(x')| \leq (|C_t| + |D_t K_t|)|x - x'|$. By (H.1(1)) and $K \in L^2(0, T; \mathbb{R}^{d \times k})$, $|A| + |BK| \in L^1(0, T; \mathbb{R})$ and $|C| + |DK| \in L^2(0, T; \mathbb{R})$. This proves that the coefficients of (2.8) satisfy the conditions of Lemma A.1, and hence the well-posedness of (2.8). $\qquad \square$

*Proof of Proposition 2.4.* For each $\varepsilon > 0$, let $X^\varepsilon = X^{K^\varepsilon}$ be such that $X_t^\varepsilon = \exp(-\int_0^t (1 + \varepsilon - s)^{-1} \, \mathrm{d}s) = \frac{1 + \varepsilon - t}{1 + \varepsilon}$ for all $t \in [0, 1]$. Thus for all $\varepsilon > 0$, $\mathcal{C}(K^\varepsilon) = \frac{1}{(1+\varepsilon)^2}$ but $\|K^\varepsilon\|_{L^1} = \log\left(\frac{1+\varepsilon}{\varepsilon}\right)$.

Now let $\tilde{K}^\varepsilon = 0.5 K^\varepsilon$, and $\tilde{X}^\varepsilon$ be such that $\tilde{X}_t^\varepsilon = \exp(-0.5 \int_0^t (1 + \varepsilon - s)^{-1} \, \mathrm{d}s) = \sqrt{\frac{1+\varepsilon-t}{1+\varepsilon}}$ for all $t \in [0, 1]$. Thus for all $\varepsilon > 0$,

$$\mathcal{C}(\tilde{K}^\varepsilon) = 0.5^2 \int_0^1 (K_t^\varepsilon \tilde{X}_t^\varepsilon)^2 \, \mathrm{d}t = \frac{0.25}{1+\varepsilon} \int_0^1 (1 + \varepsilon - t)^{-1} \, \mathrm{d}t = \frac{0.25}{1+\varepsilon} \log\left(\frac{1+\varepsilon}{\varepsilon}\right) > 0.$$

Hence $\mathcal{C}(\tilde{K}^\varepsilon) > \mathcal{C}(\mathbf{0})$ and $\lim_{\varepsilon \to 0} \frac{\mathcal{C}(\tilde{K}^\varepsilon)}{\mathcal{C}(K)} = \lim_{\varepsilon \to 0} 0.25(1 + \varepsilon) \log\left(\frac{1+\varepsilon}{\varepsilon}\right) = \infty.$ $\qquad \square$

# B Experiment details

This section presents additional details for the numerical experiments in Section 4.

**Optimal cost.** Let $P^\star \in C([0, T]; \mathbb{R})$ solve the following Riccati equation: for all $t \in [0, T]$,

$$\left(\tfrac{\mathrm{d}P}{\mathrm{d}t}\right)_t - B_t \left(P_t \sum_{j=1}^3 (D^{(j)})^\top D^{(j)} + \rho \bar{V}^{-1}\right)^{-1} B_t^\top P_t^2 = 0; \quad P_T = \tfrac{\mu}{2}. \tag{B.1}$$

Then the optimal policy of (4.1)-(4.2) satisfies $\nu_t^\star(x) = \mathcal{N}(K_t^\star x, V_t^\star)$ for all $(t, x) \in [0, T] \times \mathbb{R}$, where

$$K_t^\star = -\left(P_t^\star \sum_{j=1}^3 (D^{(j)})^\top D^{(j)} + \rho \bar{V}^{-1}\right)^{-1} B_t^\top P_t^\star, \quad V_t^\star = \rho \left(P_t^\star \sum_{j=1}^3 (D^{(j)})^\top D^{(j)} + \rho \bar{V}^{-1}\right)^{-1}.$$

Moreover, let $\varphi^\star \in C([0, T]; \mathbb{R})$ satisfy for all $t \in [0, T]$,

$$\left(\tfrac{\mathrm{d}}{\mathrm{d}t}\varphi\right)_t + \tfrac{1}{2}\mathrm{tr}\left(\left(P_t^\star \sum_{j=1}^3 (D^{(j)})^\top D^{(j)} + \rho \bar{V}^{-1}\right) V_t^\star\right) + \tfrac{\rho}{2}\left(-3 + \ln\left(\tfrac{\det(\bar{V})}{\det(V_t)}\right)\right) = 0; \quad \varphi_T = 0, \tag{B.2}$$

Then the optimal cost of (4.1)-(4.2) is given by $\mathcal{C}^\star = \tfrac{1}{2}\mathbb{E}[\xi_0^\top \xi_0]P_0^\star + \varphi_0^\star.$

**Implementation details.** The numerical experiments are coded by using Tensorflow. To examine the linear convergence, the scheme (2.26) is implemented on the uniform time grid $\pi_c$ with mesh size $\Delta t = 1/128$. Indeed, let $K^0 \equiv (1/3, 1/3, 1/3)$ and $V^0 \equiv 0.1 D^\top D$ be the initial guess. For each $n \in \mathbb{N}_0$, given $\theta^n = (K_i^n, V_i^n)_{i=0}^{127}$, consider the Euler–Maruyama discretisation of (4.3): $X_0 = \xi_0$ and for all $i = 0, \ldots, 127$,

$$X_{i+1} = X_i + B_{i\Delta t}\left(K_i^n X_i + (V_i^n)^{\frac{1}{2}}\zeta_i\right)\Delta t + \sum_{j=1}^3 D^{(j)}\left(K_i^n X_i + (V_i^n)^{\frac{1}{2}}\zeta_i\right)\Delta W_i^{(j)}, \qquad \text{(B.3)}$$

where $(\Delta W_i^{(j)})_{i=0,\ldots,127,j=1,\ldots3}$ are independent normal random variables with mean zero and variance $1/128$, and $(\zeta_i)_{i=0}^{127}$ are independent standard normal random vectors in $\mathbb{R}^3$. We simulate $N_{\mathrm{MC}} = 10^5$ independent trajectories of (B.3) and approximate $\mathcal{C}(\theta^n)$ as follows (cf. (3.3)):

$$\widehat{\mathcal{C}}(\theta^n) := \frac{1}{N_{\mathrm{MC}}}\sum_{l=1}^{N_{\mathrm{MC}}}\frac{1}{2}\left(\mu X_{128,l}^2 + \rho\sum_{i=0}^{127}\left((K_i^n)^\top \bar{V}^{-1}K_i^n X_{i,l}^2 + \mathrm{tr}(\bar{V}^{-1}V_i^n) - 3 + \ln\left(\frac{\det(\bar{V})}{\det(V_i^n)}\right)\right)\Delta t\right),$$

where $(X_{i,l})_{i=0}^{128}$, $l = 1, \ldots, N_{\mathrm{MC}}$, represents the $l$-th trajectory of (B.3). The required gradients $(\widehat{\nabla_{K_i^n}\mathcal{C}}, \widehat{\nabla_{V_i^n}\mathcal{C}})_{i=1}^{127}$ are computed using automatic differentiation along these paths, and for each $i = 0, \ldots, 127$, the state covariance $\Sigma_{i\Delta t}^{\theta^n}$ is estimated by $\widehat{\Sigma}_i^n := \frac{1}{N_{\mathrm{MC}}}\sum_{l=1}^{N_{\mathrm{MC}}}X_{i,l}^2$. The policy is then updated as follows (cf. (2.26)): for all $i = 0, \cdots, 127$,

$$K_i^{n+1} = K_i^n - \frac{\tau}{\Delta t\widehat{\Sigma}_i^n}\widehat{\nabla_{K_i^n}\mathcal{C}}, \quad V_i^{n+1} = V_i^n - \frac{\tau}{\Delta t}\left(\widehat{\nabla_{V_i^n}\mathcal{C}}\,V_i^n + V_i^n\,\widehat{\nabla_{V_i^n}\mathcal{C}}\right).$$

The optimal cost of (4.1)-(4.2) is computed by solving (B.1) and (B.2) with the explicit Euler scheme on $\pi_c$, which leads to the value $\mathcal{C}^\star = 0.0402$.

To examine the robustness of (2.26) in time discretisation, a family of coarser time grids $\pi_m = \{i\frac{1}{m}\}_{i=0}^m \subset \pi_c$, $m \in \{8, 16, 32, 64, 128\}$, have been introduced. The PG scheme only updates policy parameters at the grid points of these coarser grids. However, to mimic a continuous-time environment, the performance of each policy iterate is still evaluated by simulating (B.3) on the fine grid $\pi_c$ (with mesh size $\Delta t = 1/128$). In particular, let $(K^0, V^0)$ be given as above. For each $m \in \{8, 16, 32, 64, 128\}$ and $n \in \mathbb{N}_0$, given $\theta^n = (K_j^n, V_j^n)_{j=0}^{m-1}$, consider the following Euler-Maruyama discretisation of (4.3): $X_0 = \xi_0$ and for all $j = 0, \ldots, m-1$, and all $i = 0, \ldots, 127$ such that $\frac{j}{m} \leq i\Delta t < \frac{j+1}{m}$,

$$X_{i+1} = X_i + B_{i\Delta t}\left(K_j^n X_i + (V_j^n)^{\frac{1}{2}}\zeta_i\right)\Delta t + \sum_{j=1}^3 D^{(j)}\left(K_j^n X_i + (V_j^n)^{\frac{1}{2}}\zeta_i\right)\Delta W_i^{(j)}, \qquad \text{(B.4)}$$

where $(\Delta W_i^{(j)})_{i=0,\ldots,127,j=1,\ldots3}$ and $(\zeta_i)_{i=0}^{127}$ are independent random variables as in (B.3). We shall sample $10^5$ independent trajectories of (B.4), and use them to approximate the gradients in $(K_j^n, V_j^n)_{j=0}^{m-1}$ and the state covariance $(\Sigma_{j/m}^{\theta^n})_{j=0}^{m-1}$ with similar methods as above. The scaled PG method (2.26) then updates the parameters by: for all $j = 0, \ldots, m-1$,

$$K_j^{n+1} = K_j^n - \frac{m\tau}{\widehat{\Sigma}_j^n}\widehat{\nabla_{K_j^n}\mathcal{C}}, \quad V_j^{n+1} = V_j^n - m\tau\left(\widehat{\nabla_{V_j^n}\mathcal{C}}\,V_j^n + V_j^n\,\widehat{\nabla_{V_j^n}\mathcal{C}}\right), \quad \text{with } \tau = 0.01, \quad \text{(B.5)}$$

while the unscaled PG method updates the parameters by: for all $j = 0, \ldots, m-1$,

$$K_j^{n+1} = K_j^n - \frac{\tau}{\widehat{\Sigma}_j^n}\widehat{\nabla_{K_j^n}\mathcal{C}}, \quad V_j^{n+1} = V_j^n - \tau\left(\widehat{\nabla_{V_j^n}\mathcal{C}}\,V_j^n + V_j^n\,\widehat{\nabla_{V_j^n}\mathcal{C}}\right), \quad \text{with } \tau = 0.08. \quad \text{(B.6)}$$

32

Let $(\theta^{\pi_m,n})_{n\in\mathbb{N}_0}$ be the policy iterate generated by (B.5), define $N^{\pi_m}(0.01)$ by

$$N^{\pi_m}(0.01) := \min\left\{n\in\mathbb{N}_0 \,|\, \widehat{\mathcal{C}}(\theta^{\pi_m,n}) - \mathcal{C}^{\star}_{\pi_m}) < 0.01\right\},$$

where $\mathcal{C}^{\star}_{\pi_m} := \frac{1}{50}\sum_{n=951}^{1000}\widehat{\mathcal{C}}(\theta^{\pi_m,n})$ approximates the optimal cost among all piecewise constant polices on $\pi_m$. The quantity $N^{\pi_m}(0.01)$ is defined similarly for the iterates generated by (B.6).

# References

[1] H. ARAKI AND S. YAMAGAMI, *An inequality for Hilbert-Schmidt norm*, Communications in Mathematical Physics, 81 (1981), pp. 89–96.

[2] A. S. BERAHAS, L. CAO, K. CHOROMANSKI, AND K. SCHEINBERG, *A theoretical and empirical comparison of gradient approximations in derivative-free optimization*, Foundations of Computational Mathematics, 22 (2022), pp. 507–560.

[3] J. BU, A. MESBAHI, AND M. MESBAHI, *Policy gradient-based algorithms for continuous-time linear quadratic control*, arXiv preprint arXiv:2006.09178, (2020).

[4] R. CARMONA, *Lectures on BSDEs, stochastic control, and stochastic differential games with financial applications*, SIAM, 2016.

[5] A. CARTEA, S. JAIMUNGAL, AND J. RICCI, *Algorithmic trading, stochastic control, and mutually exciting processes*, SIAM Review, 60 (2018), pp. 673–703.

[6] M. FAZEL, R. GE, S. KAKADE, AND M. MESBAHI, *Global convergence of policy gradient methods for the linear quadratic regulator*, in International Conference on Machine Learning, PMLR, 2018, pp. 1467–1476.

[7] D. FIROOZI AND S. JAIMUNGAL, *Exploratory LQG mean field games with entropy regularization*, Automatica, 139 (2022), p. 110177.

[8] M. GEIST, B. SCHERRER, AND O. PIETQUIN, *A theory of regularized Markov decision processes*, in International Conference on Machine Learning, PMLR, 2019, pp. 2160–2169.

[9] P. J. GRABER, *Linear quadratic mean field type control and mean field games with common noise, with application to production of an exhaustible resource*, Applied Mathematics & Optimization, 74 (2016), pp. 459–486.

[10] B. GRAVELL, P. M. ESFAHANI, AND T. SUMMERS, *Learning optimal controllers for linear systems with multiplicative noise via policy gradient*, IEEE Transactions on Automatic Control, 66 (2020), pp. 5283–5298.

[11] B. HAMBLY, R. XU, AND H. YANG, *Policy gradient methods find the Nash equilibrium in n-player general-sum linear-quadratic games*, arXiv preprint arXiv:2107.13090, (2021).

[12] B. M. HAMBLY, R. XU, AND H. YANG, *Policy gradient methods for the noisy linear quadratic regulator over a finite horizon*, Available at SSRN, (2020).

[13] A. HAN, B. MISHRA, P. K. JAWANPURIA, AND J. GAO, *On Riemannian optimization over positive definite matrices with the Bures-Wasserstein geometry*, Advances in Neural Information Processing Systems, 34 (2021), pp. 8940–8953.

[14] Y. Jia and X. Y. Zhou, *Policy gradient and actor-critic learning in continuous time and space: Theory and algorithms*, The Journal of Machine Learning Research, 23 (2022), pp. 12603–12652.

[15] ——, *q-learning in continuous time.*, The Journal of Machine Learning Research, 24 (2023), pp. 161–1.

[16] Z. Jin, J. M. Schmitt, and Z. Wen, *On the analysis of model-free methods for the linear quadratic regulator*, arXiv preprint arXiv:2007.03861, (2020).

[17] S. M. Kakade, *A natural policy gradient*, Advances in neural information processing systems, 14 (2001).

[18] B. Kerimkulov, J.-M. Leahy, D. Šiška, and Ł. Szpruch, *Convergence of policy gradient for entropy regularized MDPs with neural network approximation in the mean-field regime*, arXiv preprint arXiv:2201.07296, (2022).

[19] V. Konda and J. Tsitsiklis, *Actor-critic algorithms*, Advances in neural information processing systems, 12 (1999).

[20] J. Mei, C. Xiao, C. Szepesvari, and D. Schuurmans, *On the global convergence rates of softmax policy gradient methods*, in International Conference on Machine Learning, PMLR, 2020, pp. 6820–6829.

[21] R. Munos, *Policy gradient in continuous time*, Journal of Machine Learning Research, 7 (2006), pp. 771–791.

[22] S. Park, J. Kim, and G. Kim, *Time discretization-invariant safe action repetition for policy gradient methods*, Advances in Neural Information Processing Systems, 34 (2021), pp. 267–279.

[23] C. Reisinger, W. Stockinger, and Y. Zhang, *Linear convergence of a policy gradient method for finite horizon continuous time stochastic control problems*, arXiv preprint arXiv:2203.11758, (2022).

[24] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz, *Trust region policy optimization*, in International conference on machine learning, PMLR, 2015, pp. 1889–1897.

[25] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, *Proximal policy optimization algorithms*, arXiv preprint arXiv:1707.06347, (2017).

[26] D. Šiška and Ł. Szpruch, *Gradient flows for regularized stochastic control problems*, arXiv preprint arXiv:2006.05956, (2020).

[27] J. Sun, X. Li, and J. Yong, *Open-loop and closed-loop solvabilities for stochastic linear quadratic optimal control problems*, SIAM Journal on Control and Optimization, 54 (2016), pp. 2274–2308.

[28] R. S. Sutton, D. McAllester, S. Singh, and Y. Mansour, *Policy gradient methods for reinforcement learning with function approximation*, Advances in neural information processing systems, 12 (1999).

[29] L. Szpruch, T. Treetanthiploet, and Y. Zhang, *Optimal scheduling of entropy regulariser for continuous-time linear-quadratic reinforcement learning*, arXiv preprint arXiv:2208.04466, (2022).

[30] C. Tallec, L. Blier, and Y. Ollivier, *Making Deep Q-learning methods robust to time discretization*, arXiv preprint arXiv:1901.09732, (2019).

[31] H. Wang, T. Zariphopoulou, and X. Y. Zhou, *Reinforcement learning in continuous time and space: A stochastic control approach*, Journal of Machine Learning Research, 21 (2020), pp. 1–34.

[32] H. Wang and X. Y. Zhou, *Continuous-time mean–variance portfolio selection: A reinforcement learning framework*, Mathematical Finance, 30 (2020), pp. 1273–1308.

[33] W. Wang, J. Han, Z. Yang, and Z. Wang, *Global convergence of policy gradient for linear-quadratic mean-field control/game in continuous time*, in International Conference on Machine Learning, PMLR, 2021, pp. 10772–10782.

[34] J. Yong and X. Y. Zhou, *Stochastic Controls: Hamiltonian Systems and HJB Equations*, vol. 43, Springer Science & Business Media, 1999.

[35] J. Zhang, *Backward Stochastic Differential Equations*, Springer, 2017.

[36] K. Zhang, B. Hu, and T. Basar, *Policy optimization for $\mathcal{H}_2$ linear control with $\mathcal{H}_\infty$ robustness guarantee: Implicit regularization and global convergence*, SIAM Journal on Control and Optimization, 59 (2021), pp. 4081–4109.

[37] K. Zhang, X. Zhang, B. Hu, and T. Basar, *Derivative-free policy optimization for linear risk-sensitive and robust control design: Implicit regularization and sample complexity*, Advances in Neural Information Processing Systems, 34 (2021), pp. 2949–2964.

[38] X. Y. Zhou and D. Li, *Continuous-time mean-variance portfolio selection: A stochastic LQ framework*, Applied Mathematics and Optimization, 42 (2000), pp. 19–33.