# A HOMOTOPIC APPROACH TO POLICY GRADIENTS
# FOR LINEAR QUADRATIC REGULATORS WITH NONLINEAR CONTROLS

CRAIG XU CHEN AND ANDREA AGAZZI

ABSTRACT. We study the convergence of deterministic policy gradient algorithms in continuous state and action space for the prototypical Linear Quadratic Regulator (LQR) problem when the search space is not limited to the family of linear policies. We first provide a counterexample showing that extending the policy class to piecewise linear functions results in local minima of the policy gradient algorithm. To solve this problem, we develop a new approach that involves sequentially increasing a discount factor between iterations of the original policy gradient algorithm. We finally prove that this homotopic variant of policy gradient methods converges to the global optimum of the undiscounted Linear Quadratic Regulator problem for a large class of Lipschitz, non-linear policies.

## 1. INTRODUCTION

Recent advances in machine learning have strongly depended on techniques from reinforcement learning to generate their superhuman results on tasks that were previously deemed "unachievable" by artificial intelligence models. Some of the most spectacular examples are the recent breakthroughs in the games of Go and Starcraft 2 [24, 25, 31] and the robotics tasks such as learning how to walk or controlling a hand [12, 21].

A shared aspect of these recent advances is the use of *policy gradient* algorithms (and their variants) [28]. Given the simplicity of these algorithms, their success in solving non-convex optimization problems is rather surprising. In the past few years, there has been significant progress made in the development of our understanding of the convergence of these algorithms - especially in the tabular setting [13, 19]. However, many of the games and control problems mentioned in the first paragraph have incredibly large (or infinite in the case of robotics) state and action spaces, so that tabular representation is computationally prohibitive. Due to this complexity, none of the above breakthroughs would have been possible without the use of *nonlinear function approximation* – most commonly, (deep) neural networks – in combination with policy gradient algorithms: With parametric function approximation, only the parameters of the approximator need to be stored rather than a table of all possible game-states.

Nevertheless, the current body of theoretical work still does not explain the success policy gradient methods have experienced in Go and Starcraft 2 due to an insufficient understanding of how such methods behave when neural networks are used as the function approximator of choice. One obstruction in this sense stems from the nonlinear behavior of neural networks, resulting in optimization landscapes that are highly non-convex [10], preventing us from using many of the tools developed in classical optimization theory. This large gap between empirical success and theory of nonlinear function approximation in reinforcement learning algorithms is also related to our limited theoretical understanding of their interplay. In this paper, we contribute to filling this gap by proving convergence results for nonlinear deterministic policies in a prototypical control problem with continuous state and action space: the Linear Quadratic Regulator (LQR) problem. In previous work by [4, 6, 8, 33], global convergence results have been obtained in this context by restricting to the linear policy class. This choice of search space favorably aligns with the policy gradient update since the optimal policy of the LQR is known to be a linear function of the state. Our work extends these results to a larger class of non-linear policies. To do so, we first construct an example to demonstrate the potential for negative interplay between the expressiveness of the policy class and the local nature of the policy gradient algorithm. We then consider a new perspective in solving reinforcement learning problems that involves iteratively updating a discount factor. With this perspective, we reduce the global convergence problem to a local one, allowing us to work with a wider policy class and to avoid potential local minima. These stronger convergence results do not come without

a cost, though; the process of iteratively updating a discount factor necessarily results in more training iterations.

## 1.1. **Related Works.**

In light of the advantages provided by policy-based approaches to reinforcement learning, there has been a recent wave of research seeking to provide theoretical guarantees for these algorithms that have been shown to work in practice.

Most of the recent work on the convergence of policy gradient methods focuses on the finite action space setting. In this context, [18] has shown that, in finite action spaces, algorithms like TRPO [22] and PPO [23] converge to the globally optimal policy. More results on the global convergence of policy gradient dynamics and extensions to function approximation are obtained in [1, 17, 19, 32]. However, in the aforementioned publications, the assumption of a finite action space is essential in obtaining the results; when considering continuous action spaces, we use the LQR as a proxy for more general environments since the LQR setting is well understood and allows for more straightforward analysis.

In the continuous action-space setting, [8] show that deterministic policy gradient methods for the policy class of *linear policies* converge to the globally optimal policy of the LQR problem. Building off of [8], [6] provides strong positive results in an explicit characterization of the LQR when learning with linear policies showing that the cost is smooth and coercive and that the discretization of various gradient flows lead to algorithms with linear convergence rates. Similarly, [33] provides a convergence guarantee for Actor-Critic methods in the LQR setting by considering the class of linear-Gaussian policies with fixed variance. Other recent work [14] also considers the class of linear policies but works in the *finite time-horizon* setting. Rather than working with a more complex learning task such as the finite-horizon LQR, our work instead tackles the problem of nonlinear function approximation and thus we prove our results in the simpler setting of the infinite-horizon, noiseless LQR; however, unlike previous work in this setting, we consider a wider (deterministic) policy class that can be understood as linear combinations of nonlinear functions. Furthermore, most works only consider the evolution of the policy with respect to time/iterations; we also consider the evolution of the optimal policy with respect to the discount factor. This new perspective allows us to use some new tools in proving the convergence of policy gradient methods in the LQR.

Most relevant to our work, [9] discusses the effect of changing the dscount factor during training of reinforcement learning algorithms. This study, although strictly empirical, shows that our homotopic variant of policy gradient methods can reduce the number of training iterations needed to learn an optimal policy, *i.e.*, that this modified reinforcement learning procedure not only improves the optimality properties of reinforcement learning algorithms but also their convergence speed. On the theoretical side, a team at DeepMind [29] has studied how the *value* function behaves as function of the discount factor and develop a set of interpolating objectives to go from a myopic value function $V_\gamma$ to the true objective $V_{\gamma^*}$; they also show that their family of interpolating objectives yields faster convergence and better asymptotic performance on MuJoCo tasks.

## 1.2. **Contributions.**

Our work continues along a line of work using the LQR as a proxy for more general learning environments [3, 5, 7, 8, 11, 14–16, 30]. In this paper, we present a homotopy-based approach inspired by the human learning process to iteratively update the policy until reaching the globally optimal policy, and in particular, we directly extend the result of [8] to nonlinear policy classes. We show that by continuously updating the discount factor after policy improvements, one can reach the global optimum regardless of the quality of the initial policy, subject to a few assumptions on the parametrization and expressiveness of the policy (*e.g.*, our policy parametrization is linear in the parameters). Although this approach can be understood as a new type of policy-gradient algorithm, our focus is theoretical and we mainly consider the homotopy approach as a tool to prove convergence results.

The intuition behind the homotopy-based algorithm comes from the observation that when teaching children, one does not expect the child to consider the long term consequences of their actions. Instead, children are first taught to learn the basics (i.e., learn with a discount factor close to zero) before moving

on to developing their long-term planning ability (increasing to a larger discount factor). Heuristically, one can make sense of this by reasoning that it's likely easier to learn the optimal policy having learned the basics than it is learning it from scratch. Mathematically, the homotopy approach relies on the fact that the optimal policy of the LQR is continuous with respect to the discount factor. In more challenging environments such as Chess or Go, this may not be the case.

Summarizing, the main contributions of this work are as follows:

- We show that one can not always expect policy gradient methods to converge even with relatively well-behaved function approximators. We provide an example policy that is linear in the parameters, contains the globally optimal policy, but does *not* converge to such globally optimal policy when trained with a fixed discount factor. Furthermore, this example satisfies all assumptions needed for our convergence result.
- We formulate the policy gradient learning process as a homotopy between various optimal policies along the discount factor. By leveraging this idea, we provide a convergence result for an expanded policy class consisting of linear combinations of non-linear functions (in state space) under some mild assumptions on the expressiveness and parametrization of the policy class (see Assumptions 1, 2).

As a corollary, our result also shows that the assumption of a stabilizing initial policy in [8] is not necessary. By using our homotopy algorithm, *any* linear policy will converge to the optimal policy of the undiscounted LQR. This is ultimately a trade-off between initial information and number of training iterations; the homotopy approach allows us to start the learning process with less background information at the cost of more training iterations. This type of trade-off has been considered before, most notably in [24] and [25] where AlphaGo was initially trained on human games via supervised learning before engaging in self-play but AlphaZero was trained entirely via self-play and no human input.

## 1.3. **Organization.**

The paper is organized as follows. In the coming section, we provide the necessary background information to understand the context of our problem (Section 2). In Section 3 we provide our result regarding the convergence of policy gradient methods with a non-linear function approximator in the LQR using the aforementioned homotopy approach inspired by the human learning process. In this section, we also provide the counterexample where policy gradient methods do not converge to the global optimum. The discussion is contained in Section 4. The full proofs of the results, as well as some numerical examples corroborating our claims, obtained are housed in Appendices A, B and C.

## 2. PRELIMINARIES AND BACKGROUND

In this paper, we study a reinforcement learning problem where an agent interacts in discrete time with the environment with the goal of minimizing some cost function that depends on the agent's states and actions. The agent does this by sequentially choosing actions over time which influence the agent's state in the environment. Throughout, for given $n, m \in \mathbb{N}$ we respectively define $X = \mathbb{R}^n$ and $U = \mathbb{R}^m$ as the set of possible states and actions of the agent. Correspondingly, for $t \in \mathbb{Z}_{\geq 0}$ we denote by $x_t \in X$ and $u_t \in U$ the state and action of the agent at time $t$. The function that maps a state to action of the agent in that state is called the *policy* $\pi : X \to U$, and the motivation behind this paper is to investigate the conditions under which an agent will learn the optimal policy of a specific environment: the *Linear-Quadratic Regulator*.

## 2.1. **The Linear-Quadratic Regulator.**

The Linear-Quadratic Regulator (LQR) is a classic optimal control problem. In this paper, we are interested in the special case where the dynamics are *linear*, time-invariant with no disturbance or added noise and the cost function is *quadratic* in the state and the control action. We consider the discounted

infinite time-horizon problem,

$$
\text{(2.1)} \qquad\qquad \text{minimize} \quad \mathbb{E}_{x_0 \sim \varrho_0} \left[ \sum_{t=0}^{\infty} \gamma^t \left( x_t^{\mathsf{T}} Q x_t + u_t^{\mathsf{T}} R u_t \right) \right]
$$

$$
\text{(2.2)} \qquad\qquad\qquad \text{with} \quad x_{t+1} = A x_t + B u_t,
$$

where $A : X \to X$ and $B : U \to X$ are the system (or transition) matrices, $Q : X \to X$ and $R : U \to U$ are positive definite cost matrices, $x_0 \in X$ is randomly distributed according to probability distribution $\varrho_0$ on $X$, and $\gamma \in [0, 1]$ is the discount factor. For a given policy $\pi$, the expression to be minimized in (2.1) is referred to as the *cost* associated to $\pi$, which we denote by $C(\pi)$. In the entirety of this work, the pair $(A, B)$ is assumed to be controllable and $(A, D)$ is assumed to be observable (where $Q = D^{\mathsf{T}} D$).

It is known [2] that the optimal control for the LQR problem is a linear function of the state,

$$
u_t = K^* x_t.
$$

If $A, B, Q, R, \gamma$ are known, the optimal control $K^*$ can be calculated explicitly. Let $P_\gamma$ denote the unique positive definite solution to the discounted discrete-time algebraic Riccati equation (DARE)

$$
P_\gamma = \gamma A^{\mathsf{T}} P_\gamma A - \gamma^2 A^{\mathsf{T}} P_\gamma B \left( R + \gamma B^{\mathsf{T}} P_\gamma B \right)^{-1} B^{\mathsf{T}} P_\gamma A + Q \,,
$$

for which we know the solution exists and since $(A, B)$ is assumed controllable and $(A, D)$ is assumed observable [2]. We can then write the optimal control as

$$
\text{(2.3)} \qquad\qquad K^*_\gamma = -\gamma \left( 2R + \gamma B^{\mathsf{T}} P_\gamma B \right)^{-1} B^{\mathsf{T}} P_\gamma A.
$$

Later, we will sometimes omit the subscript $\gamma$ when the dependencies are clear from context.

Note that, in general, the optimal policy for the discounted LQR is different from the optimal policy for the undiscounted LQR. For instance, the optimal policy with respect to the discounted LQR may not even be stabilizing [20].

Throughout the remainder of the paper, we will use $\gamma$ to denote the discount rate which we choose given an initial policy. We assume that $\gamma \in [0, 1]$ is small enough such that the cost function is still bounded from above for all initial states $x_0 \in \text{supp} \left( \varrho_0 \right)$, where $\varrho_0$ is the distribution of initial states.

*Why LQR?.* Although the assumptions for the linear-quadratic system are rather restrictive, the LQR is still an important problem for multiple reasons. Firstly, it is one of the few problems where the closed-form of the optimal control exists. This fact by itself makes the LQR a pleasant environment to work with. Secondly, the LQR has a continuous state and action space. Thirdly, despite its simplicity, there is still a lack of theoretical guarantees regarding convergence of reinforcement learning algorithms in this setting. All of these factors have certainly contributed to its re-emergence as a benchmark environment in reinforcement learning theory; it is both nice to work with mathematically and difficult enough to be applicable to the real world.

## 2.2. Policy Gradient Methods.

Policy gradient methods are a general class of reinforcement learning algorithms where the agent directly learns a policy – which is assumed to be in a certain parametric family of functions – rather than first learning a value function. During training, the parameters of the policy are updated via gradient descent on the cost function. As mentioned in the introduction, policy gradient methods have increasingly become the approach of choice for solving difficult reinforcement learning problems. One reason for this popularity is that policy-based methods, as opposed to value-based, can easily learn stochastic policies, whereas action-value based methods have no natural *and* flexible way to do so.

The fundamental result underlying the successful application of these methods is the policy gradient theorem [28] which provides an explicit formula for the gradient of the performance in terms of the gradient of the policy w.r.t its parameters – importantly, the gradient of the state distribution is *not* needed, greatly simplifying the application of these methods.

The analogous theorem in the deterministic setting (first proven by [26]) allows one to conveniently compute the gradient of the cost function when considering a class of deterministic policies, as is the case

in this paper. Furthermore, as shown in [8], it is possible to explicitly compute the gradient of the LQR cost function for linear policies. They show that for a linear policy $K$,

$$\nabla C(K) = 2\left(\left(R + \gamma B^{\mathsf{T}} P_K B\right)K + \gamma B^{\mathsf{T}} P_K A\right)\Sigma_K$$

where $P_K$ denotes the "cost to-go" matrix (*i.e.*, $C(K) = \mathbb{E}_{x_0 \sim \varrho_0}\left[x_0^{\mathsf{T}} P_K x_0\right]$) and $\Sigma_K$ is the state covariance matrix.

In our work, we consider the same LQR setup as in [8]. That is, we also consider a deterministic policy rather than a stochastic policy, a randomly distributed initial state, and noiseless dynamics. Note that it is known [2] that the inclusion of additive zero-mean white noise to the LQR dynamics (2.2) does not change the optimal control.

## 3. Main Results

*Notation.* We use $\|\cdot\|$ as the operator norm of a matrix or the Euclidean norm of a vector. We will use $\lambda_{\min}(\cdot)$ to refer to the smallest eigenvalue of a matrix.

In the first subsection, we construct the aforementioned example where vanilla (fixed discount factor) policy gradients only yield convergence to a *local* optimum. Importantly, this policy class of the example satisfies the assumptions needed for our convergence result. In the following subsection, we provide the main theorems which show convergence of the homotopy algorithm. The lemmas and additional details for both subsections have been omitted for clarity, the full details can be found in Appendices C and B. We denote by $\{x_t\}$ the sequence generated by the dynamics $x_{t+1} = Ax_t + B\pi_\vartheta(x_t)$. Throughout, $\pi_\vartheta$ denotes the parametric policy of choice evaluated with parameters $\vartheta \in \mathbb{R}^d$. In this paper, we consider policies of the form

$$\pi_\vartheta(x_t) = \sum_{k=1}^{d} \vartheta_k f_k(x_t) \tag{3.1}$$

where each $f_k : \mathbb{R}^n \to \mathbb{R}^m$ is a non-linear, globally Lipschitz function and $\vartheta_k \in \mathbb{R}$ is a parameter to be learned.

To establish our results we make the following assumptions:

**Assumption 1.** *The functions $f_k$ are linearly independent.*

**Assumption 2.** *For the LQR problem defined by $(A, B, Q, R)$, for all $\gamma \in [0, 1]$, the $\gamma$-optimal policy $K_\gamma^*$ can be represented by $\pi_\vartheta$, i.e. $K_\gamma^* \in Span\{f_k\}$.*

Assumption 1 is a standard assumption when considering function approximation via linear combination. We need this assumption to guarantee the uniqueness of policy parametrization. Assumption 2 is necessary to guarantee that, for any choice of $\gamma$, the optimal policy is in the parametric policy class being considered.

We note that the assumption of linear independence is, however, not essential for our proofs; rather, it is made for technical convenience to simplify our proofs of Lemma C.4 and Theorem 3.5, guaranteeing local strict convexity of the landscape. For example, in the proof of Lemma C.4, if we were to relax Assumption 1, we would see that the eigenvectors that now correspond to the zero eigenvalue of the Hessian would be the directions where the policy is unaffected by changing the parameters. Therefore, we could relax the assumption to allow overparameterization and in turn to allow for wider networks in Remark 3.2.

**Lemma 3.1.** *Assumptions 1 and 2 imply that the optimal parameters are continuous in the discount factor. In other words, for $\gamma \in [0, 1]$ and for any $\varepsilon > 0$, $\exists \delta > 0$ such that for $\gamma' \in [0, 1]$ such that $|\gamma - \gamma'| < \delta$, $\|\vartheta_\gamma^* - \vartheta_{\gamma'}^*\| < \varepsilon$, where $\vartheta_\gamma^*$ denotes the parameters that correspond to $K_\gamma^*$ for the policy class of (3.1).*

*Proof.* Let $\{g_k\}$ denote the orthonormal set obtained from $\{f_k\}$ via Gram-Schmidt. We can then write

$$\pi_\vartheta = \sum_{k=1}^{d} \vartheta_k \sum_{l=1}^{d} \langle f_k, g_l \rangle g_l = \sum_{l=1}^{d} \left\langle \sum_{k=1}^{d} \vartheta_k f_k, g_l \right\rangle g_l =: \sum_{l=1}^{d} \hat{\vartheta}_l g_l$$

Thus, we see that

$$\|\pi_\vartheta - \pi_{\vartheta'}\|^2 = \sum_{l=1}^d \left( \hat{\vartheta}_l - \hat{\vartheta}'_l \right)^2$$

By Lemma C.1 we see that the optimal controls $K_\gamma^*$ are continuous in $\gamma$. Thus we conclude that

$$|\gamma - \gamma'| < \delta \implies \|\pi_{\vartheta_\gamma^*} - \pi_{\vartheta_{\gamma'}^*}\| < \varepsilon \implies \|\hat{\vartheta} - \hat{\vartheta}'\| < \varepsilon.$$

Now, notice that we can also write the hat-coefficients as

$$\begin{bmatrix} \hat{\vartheta}_1 \\ \vdots \\ \hat{\vartheta}_d \end{bmatrix} = \begin{bmatrix} \langle f_1, g_1 \rangle & \cdots & \langle f_d, g_1 \rangle \\ \vdots & \ddots & \vdots \\ \langle f_1, g_d \rangle & \cdots & \langle f_d, g_d \rangle \end{bmatrix} \begin{bmatrix} \vartheta_1 \\ \vdots \\ \vartheta_d \end{bmatrix},$$

and this change-of-basis matrix is invertible by Assumption 1. Thus, denoting the above matrix by $A$, we conclude that

$$\|\hat{\vartheta} - \hat{\vartheta}'\| < \varepsilon \implies \|\vartheta - \vartheta'\| < \|A^{-1}\|\varepsilon$$

$\square$

**Remark 3.2.** 2-layer ReLU neural networks (without bias and with a slightly modified architecture) trained in the random features regime satisfy the above assumptions. In this regime, we only train the weights of the output layer $\{\vartheta_k\} = \{a_k\} \cup \{b_k\}$. For the $i$'th component of the control, we can write

$$\pi_i(x) = \sum_{k=1}^n a_k \sigma(w_k^\mathsf{T} x) + b_k \sigma(-w_k^\mathsf{T} x),$$

where $\sigma(x) = \max(0, x)$ is the rectified linear unit (ReLU), $x \in \mathbb{R}^n$ is the state, and $w_k \sim N(0, I_n)$ are the first layer weights fixed at initialization. This type of modification with duplicated ReLU neurons has been considered before (see Appendix F of [34]).

*Proof of Remark 3.2.* Recall that $x \in \mathbb{R}^n$. For Assumption 1, we note that $\{w_k\}$ is a linearly independent set with probability 1. With this architecture, for every $x \in \mathbb{R}^n$ there will be exactly $n$ terms in the sum that are non-zero. Since $\{w_k\}$ forms a basis of $\mathbb{R}^n$, Assumption 2 follows. $\square$

**Remark 3.3.** The policy class $\{f_k\} = \{1_{i,j}\} \cup \{F\}$ for some globally Lipschitz function $F$ satisfies the above assumptions. Here $1_{i,j} \in \mathbb{R}^{m \times n}$ denotes the matrix with all 0's except for a 1 in the $i, j$'th entry. This policy can be written as

$$u_t = Kx_t + \vartheta F(x_t),$$

where $K \in \mathbb{R}^{m \times n}$, $\vartheta \in \mathbb{R}$ are the parameters to be learned. This will be the policy class we consider in Section 3.1 to construct our counterexample.

## 3.1. Counterexample with a Globally Lipschitz Policy.

We now proceed to construct the aforementioned example where vanilla (fixed discount factor) policy gradients may result in convergence to a *local* optimum. Importantly, the policy class of this example satisfies the assumptions needed for the convergence result provided in the subsequent section. The lemmas and additional details for both subsections have been omitted for clarity, the full details can be found in Appendices B and C.

We consider a simple 1-dimensional system ($X = U = \mathbb{R}, \gamma \in (0,1)$). Fix $\gamma$, we then choose $(A, B) = (0, 1)$ and $(Q, R) = (1, R)$ as our system matrices (see (2.2)) for any $R \in (0, \gamma)$. Notice that this system is controllable and observable so there is an optimal control. For $a \in \mathbb{R}$ and $\delta > 0$, we define the spike or tent function of width $2\delta$ centered at $a$:

$$\Lambda_{a,\delta}(x) := \begin{cases} \frac{x-(a-\delta)}{\delta} & \text{if } x \in [a-\delta, a] \\ \frac{(a+\delta)-x}{\delta} & \text{if } x \in [a, a+\delta] \\ 0 & \text{else} \end{cases}.$$
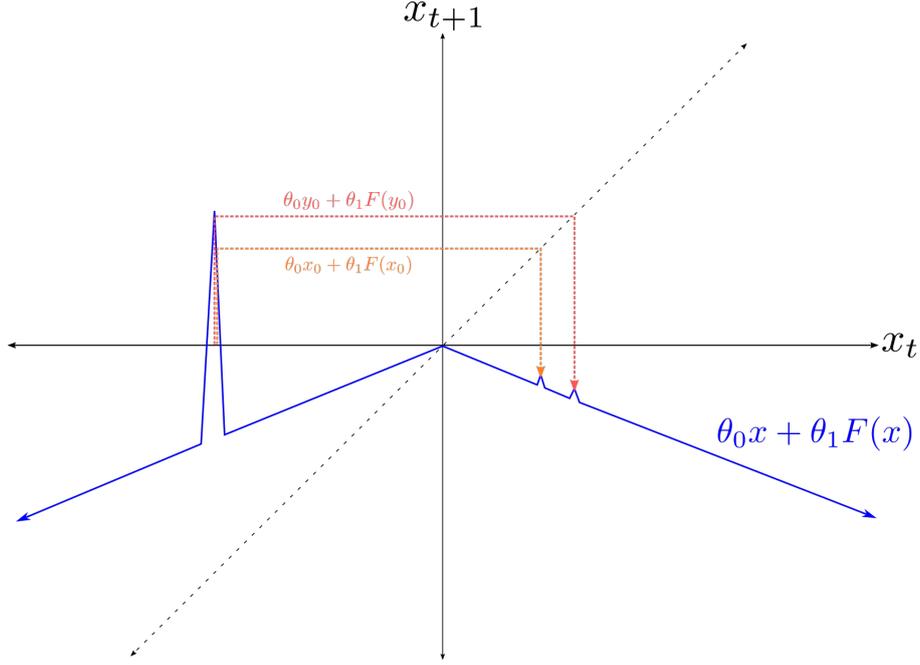
FIGURE 1. A cobweb diagram for the first step of the LQR dynamics in the counterexample. The blue line represents the dynamics of the system at initialization ($\vartheta_0 = 0, \vartheta_1 = 1$). Our initial distribution of states is heavily weighted around the two points ($x_0$ and $y_0$) that map to the two small tents on the right; this requires the points to lie somewhere in the domain of the spike on the left because all other states are mapped to negative numbers.

Next, we define the control as

$$u_t = \vartheta_0 x_t + \vartheta_1 F(x_t) \qquad \text{where } F(x_t) := \omega_0 |x_t| + \sum_{i=1}^{3} \omega_i \Lambda_{a_i, \delta_i}(x_t), \tag{3.2}$$

for a choice of parameters $\omega_0 \in (-1, 0)$ and $\{a_i, \omega_i, \delta_i\}_{i=1}^3$ where $a_2, a_3 < |a_1|$ that are *fixed* during training of the model. We set

$$\omega_0 = -0.5, \quad \delta_1 = 0.1, \quad a_1 = -2, \quad \omega_1 = 3, \quad a_2 = 1.5, \quad a_3 = 1.8, \quad \delta_2 = \delta_3 = \delta, \quad \omega_2 = \omega_3 = 0.2, \tag{3.3}$$

for a parameter $\delta$ to be chosen later. The dynamics resulting from this choice of parameters are depicted in Figure 1. We note that the optimal policy for the LQR belongs to this policy class and corresponds to the parameters $(0, 0)$.

We choose the initial condition $(\vartheta_0, \vartheta_1) = (0, 1)$ and define the distribution of the initial state as the following measure with support on $(-5, 5)$

$$\mu_0(x) = \frac{1 - \varepsilon}{2} \delta_{x_0} + \frac{1 - \varepsilon}{2} \delta_{y_0} + \frac{\varepsilon}{10}, \tag{3.4}$$

where $x_0 = \min F^{-1}(a_2)$ and $y_0 = \min F^{-1}(a_3)$, *i.e.*, we define $x_0$ and $y_0$ as the smaller of the two pre-images for $a_2$ and $a_3$, respectively.

**Proposition 3.4.** *Fix $\gamma \in (0, 1)$ and $0 < R < \gamma$. For the LQR problem with $(A, B) = (0, 1)$, $(Q, R) = (1, R)$, policy class (3.2) with parameters chosen in (3.3), and initial state distribution $\mu_0$ from (3.4), there exists $\delta > 0$ and $\varepsilon > 0$ such that the point in parameter space $\vartheta = (0, 1)$ is a local minimum of the LQR cost function.*

To prove that the point $(0, 1)$ is a local minimum of the cost function, and therefore a suboptimal fixed point of the policy gradient algorithm, we show that any infinitesimal change in $\vartheta_0$ and $\vartheta_1$ around $(\vartheta_0, \vartheta_1) = (0, 1)$ will increase such cost function.

7

To understand why this is the case, consider first the trajectory of the points $x_0$ and $y_0$. It is easy to see that, by the design of the control, any small perturbation of the parameters that change the image of $x_0$ and $y_0$ will result in a state of larger absolute value, resulting in turn in a larger cost for those trajectories. On the other hand, it is also possible that a change in parameters actually lowers the cost for some trajectories starting in $x \in [-5, 5] \setminus \{x_0, y_0\}$. Still, by our choice of $\varepsilon$ it is possible to show that the increase in cost for $x_0, y_0$ outweighs the decrease. More specifically, the proof follows the following steps (the full details are presented in Appendix B):

(1) For both $x_0, y_0$, we first prove that we can choose the parameter $\delta > 0$ (see (3.3)) such that outside of a small cone in parameter space, we can bound the change in cost from below to show that it is always positive.
(2) We then guarantee that we can pick the parameters of the model such that the cones described above corresponding to $x_0$ and $y_0$ only intersect at the origin.
(3) Next, we show that the cost difference is still positive even if the infinitesimal change in $(\vartheta_0, \vartheta_1)$ belongs inside one cone and outside of the other.
(4) Lastly, we bound the contribution of the uniform distribution term in $\mu_0$, showing that it is negligible with respect to any of the contributions to the cost function from the trajectories starting at $x_0$ or $y_0$.

Qualitatively, we can plot the cost landscapes for our example around $\vartheta = (0, 1)$ to see that this point indeed is a local minimum (see Figure 3 in the appendix). The example we have described is only one member of a family of possible counterexamples. The proofs we provide have used some of the specifics of our example; however, they can be very easily generalized to describe a set of parameters and constraints necessary for the point $\vartheta = (0, 1)$ to correspond to a local minimum of the cost function. We only present this specific example as it is sufficient to demonstrate the failure of vanilla policy gradients in the LQR.

## 3.2. Convergence Results.

In this section we present our convergence result for policies of the form

$$\pi_\vartheta(x_t) = \sum_{k=1}^{d} \vartheta_k f_k(x_t)$$

where each $f_k : \mathbb{R}^n \to \mathbb{R}^m$ is a non-linear, globally Lipschitz function and $\vartheta_k \in \mathbb{R}$ is a parameter to be learned. We prove our results under Assumptions 1 and 2, restated here for convenience.

**Assumption 1.** *The functions $f_k$ are linearly independent.*

**Assumption 2.** *For the LQR problem defined by $(A, B, Q, R)$, for all $\gamma \in [0, 1]$, the $\gamma$-optimal policy $K_\gamma^*$ can be represented by $\pi_\vartheta$, i.e. $K_\gamma^* \in Span\{f_k\}$.*

We prove the convergence of policy gradient methods for the policy class of (3.1) by using the homotopy idea discussed in the introduction. Informally, our idea is that if one updates $\gamma$ by increasing it an infinitesimally small amount to $\gamma'$, one might expect that the optimal policies for the corresponding discount factors might not be too far from each other. Then if these policies, and importantly their parameterizations, are sufficiently close, policy gradient methods will yield convergence from the $\gamma$-optimal policy to the $\gamma'$-optimal policy. It's important to remember that one should not expect this behavior in general; we have the simplicity of the LQR to thank for the continuity of the optimal policy in the discount factor, but in more challenging settings such as Chess or Go, this may not be the case.

---

**Algorithm 1:** Homotopy algorithm for LQR. Here $\text{PG}(\vartheta, \gamma)$ denotes the fixed point of the policy gradient algorithm for the given LQR problem with discount factor $\gamma$ and initial condition $\vartheta$.

---

**Result:** Optimal policy of the LQR problem
$\hat{\vartheta}^* = \vartheta(0), \{\gamma_n\}_{n=1}^N$;
**for** $n \in \{1, \ldots, N\}$ **do**
$\quad \gamma \leftarrow \gamma_n$ ;
$\quad \hat{\vartheta}^* \leftarrow \text{PG}(\hat{\vartheta}^*, \gamma)$ ;
**end**
**return** $\hat{\vartheta}^*$

---

In other words, given some initial parameters $\vartheta(0)$ and an increasing sequence of discount factors $\gamma_n \in [0, 1]$, we run the policy gradient algorithm until convergence to the optimal policy of the first discount factor. We then use this optimal policy as the initial policy for the system with the next discount factor, and let the policy gradient algorithm converge to the corresponding optimum. We continue iterating until we have converged to the optimal policy of the undiscounted LQR. To prove this algorithm converges, we need to prove that policy gradients yield convergence from the $\gamma$-optimal policy to the $\gamma'$-optimal policy if $\gamma$ and $\gamma'$ are sufficiently close, and that for any arbitrary initialization, the first iteration of the algorithm converges.

The proof is split into two main steps:

(1) Prove convergence of arbitrary initializations to the $\gamma = 0$ optimal policy.
(2) Prove local convergence of policy gradient methods. In other words, we show that for discount factors $\gamma'$ and $\gamma$ sufficiently close, policy gradient methods yield convergence from the $\gamma$-optimal policy to the $\gamma'$-optimal policy.

The first step brings policies into the "homotopy" regime, while the second step shows that we can iterate this process to go from the $\gamma = 0$ optimal policy to the $\gamma = 1$ (undiscounted) optimal policy. We remark that it is also reasonable to ignore this first step and to initialize the policy to the zero function as is occasionally done in practice.

In the statements of the theorems and their proofs, we recall that $\vartheta \in \mathbb{R}^d$ denotes the parameters of the model and $\vartheta_\gamma^* \in \mathbb{R}^d$ denotes the optimal parameters for the given value of $\gamma$, which exist and are unique by Assumptions 1, 2.

**Theorem 3.5** (Convergence at $\gamma = 0$). *Let $\{f_k\}_{k=1}^d$ be a set of globally $\text{Lip}(f_k)$-Lipschitz continuous functions, and let $\varrho_0$ (the distribution of the initial state) have full support on $\mathbb{R}^d$. Under Assumptions 1, 2, any random initialization of the parameters $\vartheta$ will converge to the optimal policy for discounted LQR with $\gamma = 0$.*

We first note that the discount factor of this first iteration does not need to be $\gamma = 0$; however, to make it easier to prove convergence of the first iteration, we can assume that $\gamma = 0$. With such high myopia, the learning task becomes much easier, since the optimal policy for the LQR with $\gamma = 0$ is simply the zero function.

The proof uses the fact that we can write a differential equation describing the evolution of the cost difference $C_\gamma(\vartheta(s)) - C_\gamma(\vartheta_\gamma^*)$ over time, where $C_\gamma(\vartheta)$ denotes the cost for the policy $\pi_\vartheta$ with discount factor $\gamma$. This results in an expression that is similar to gradient-domination bounds and from this the conclusion follows.

We now state the result about the local convergence of policy gradient methods in a neighborhood of the optimal policy for the chosen discount factor:

**Theorem 3.6** (Convergence Near $\gamma$-Optimality). *Under the same conditions as Theorem 3.5, for any $\gamma \in (0, 1)$ there exists $\delta > 0$, $\lambda > 0$ such that for all initial conditions $\vartheta(0)$ with $\Delta\vartheta(0) := \|\vartheta(0) - \vartheta_\gamma^*\| < \delta$,*

$$C_\gamma(\vartheta(s)) - C_\gamma(\vartheta_\gamma^*) \le e^{-s\lambda} \left( C_\gamma(\vartheta(0)) - C_\gamma(\vartheta_\gamma^*) \right),$$

*and that $\lim_{s \to \infty} \vartheta(s) = \vartheta_\gamma^*$.*

The main idea behind the proof of this result is to leverage a Taylor expansion of the cost function $C_\gamma(\,\cdot\,)$ around its minimum $\vartheta_\gamma^*$ for any fixed value of $\gamma \in (0,1)$. We show that by considering the second order Taylor expansion in a small neighborhood around optimality, the landscape of the policy gradient algorithm is locally convex. This allows to show that from any point within this neighborhood, policy gradient methods experience exponential convergence to the $\gamma$-optimal policy by applying Grönwall's inequality.

Lastly, we want to show that by continuously updating the discount factor, we can ensure that we are always within a small neighborhood of the optimal policy for the *next* discount factor, eventually reaching the undiscounted optimal.

**Theorem 3.7** (Convergence to Undiscounted Optimal). *For any initial parameters $\vartheta_0$, Algorithm 1 will converge to the undiscounted optimal policy $\vartheta^*$.*

**Corollary 3.8.** *There exists an increasing sequence $\{\gamma_n\}_n$ of discount factors such that Algorithm 1 converges to the global optimum of the undiscounted LQR.*

*Proof of Theorem 3.7.* Theorem 3.5 shows that any initialization of the policy will converge to the optimal policy of the $\gamma = 0$ system. Theorem 3.6 together with positive-definiteness of $H_\gamma$ then tells us that policy gradient methods experience local exponential convergence near optimality. Corollary C.2 and Lemma 3.1 tell us that for $|\gamma' - \gamma|$ sufficiently small, $\vartheta_\gamma^*$ will be in close enough to $\vartheta_{\gamma'}^*$ for Theorem 3.6 to apply. One can iterate this process until reaching the undiscounted optimal policy. □

## 4. DISCUSSION AND CONCLUSION

This paper provides a convergence guarantee for model-based policy gradient methods in the setting of the LQR with the policy class $\{\sum_{k=1}^{d} \vartheta_k f_k(x) \mid \vartheta_k \in \mathbb{R}\}$ contingent on a few assumptions on the expressiveness and parametrization of the policy class. Unlike previous works, we adopt a new approach for proving convergence, namely, the homotopy-based approach which, to the best of our knowledge, has not been used before in theoretical reinforcement learning. Using a discount factor to guarantee finite costs and well-defined updates can also be applied to solving issues involving a chaotic policy, such as the one arising naturally from piece-wise linear policies such as the tent map.

Furthermore, we provide an example illustrating a situation where expanding the policy class in the LQR case leads to local minima of vanilla policy gradient algorithms. This issue is resolved by applying the homotopic variant of policy gradients developed in this paper.

*Limitations.* It is important to note that our work considers an idealized, model-based and infinitesimal-stepsize approach. To generalize our results to the model-free perspective, one could proceed in a similar fashion to [8] by showing that when the roll-out is sufficiently long, the cost function and covariance of the state trajectory can be accurately approximated and that with enough samples, the true gradient can be approximated within a desired accuracy. Some other interesting avenues include relaxing the assumption of continuous updating, *i.e.*, to consider the stochastic approximation problem resulting from the real-life setting of finite samples and finite step-size and establishing quantitative bounds for the homotopy method.

Furthermore, the approach proposed in this work results in a significant increase in the computational cost of training when compared to the vanilla policy gradient algorithm. This results in a trade-off between computational cost and convergence guarantees for solving the control problem at hand. While this trade-off may not be of immediate interest in practice, what we propose is a new paradigm for studying the convergence of reinforcement learning algorithms.

*Future Work.* The problem of extending our proofs to establish convergence for the policy class of neural networks remains open. We expect that it should be possible to extend our result to linearized neural networks (such as the Neural Tangent Kernel regime) but we leave this result for future research.

*Societal Impact.* In the much bigger picture, our work is related to the issue of bias in machine learning. The common practice of seeding reinforcement learning agents with human training data can be problematic when the human-generated data contains implicit biases that we do not want in the agent — almost all data associated with humans will contain such biases. Our work specifically shows that we can do away

with some assumptions on the initial knowledge of reinforcement learning agents (*e.g.*, assuming the initial policy of the LQR is stabilizing), and thus, in applicable settings, do away with any dependence on human-generated data. We choose to include this reminder because it is imperative to keep this question in mind when designing machine learning models with the goal of bettering human lives. A simple example is the autonomous vehicle trolley problem. In some cultures, it may be preferable to crash into one individual over the other, and thus a model trained in one country may have disastrous consequences when deployed in another.

## REFERENCES

[1] Alekh Agarwal, Sham M. Kakade, Jason D. Lee, and Gaurav Mahajan. Optimality and approximation with policy gradient methods in markov decision processes. *CoRR*, abs/1908.00261, 2019.

[2] Brian D. O. Anderson and John B. Moore. *Optimal Control: Linear Quadratic Methods*. Prentice-Hall, Inc., USA, 1990.

[3] Sanjeev Arora, Elad Hazan, Holden Lee, Karan Singh, Cyril Zhang, and Yi Zhang. Towards provable control for unknown linear dynamical systems. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Workshop Track Proceedings*, 2018.

[4] Jalaj Bhandari and Daniel Russo. Global optimality guarantees for policy gradient methods. *CoRR*, abs/1906.01786, 2019.

[5] Steven J. Bradtke. Reinforcement learning applied to linear quadratic regulation. In *Advances in Neural Information Processing Systems 5*, pages 295–302. Morgan-Kaufmann, 1993.

[6] Jingjing Bu, Afshin Mesbahi, Maryam Fazel, and Mehran Mesbahi. Lqr through the lens of first order methods: Discrete-time case, 2019.

[7] Sarah Dean, Horia Mania, Nikolai Matni, Benjamin Recht, and Stephen Tu. On the sample complexity of the linear quadratic regulator, 2018.

[8] Maryam Fazel, Rong Ge, Sham Kakade, and Mehran Mesbahi. Global convergence of policy gradient methods for the linear quadratic regulator. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1467–1476. PMLR, 10–15 Jul 2018.

[9] Vincent François-Lavet, Raphael Fonteneau, and Damien Ernst. How to discount deep reinforcement learning: Towards new dynamic strategies, 2016.

[10] Ian Goodfellow, Oriol Vinyals, and Andrew Saxe. Qualitatively characterizing neural network optimization problems. In *International Conference on Learning Representations*, 2015.

[11] Benjamin Gravell, Peyman Mohajerin Esfahani, and Tyler Summers. Learning robust control for LQR systems with multiplicative noise via policy gradient. *CoRR*, abs/1905.13547, 2019.

[12] Tuomas Haarnoja, Sehoon Ha, Aurick Zhou, Jie Tan, George Tucker, and Sergey Levine. Learning to walk via deep reinforcement learning. In *Robotics: Science and Systems*, 2019.

[13] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1861–1870. PMLR, 2018.

[14] Ben Hambly, Renyuan Xu, and Huining Yang. Policy gradient methods for the noisy linear quadratic regulator over a finite horizon, 2020.

[15] Elad Hazan, Holden Lee, Karan Singh, Cyril Zhang, and Yi Zhang. Spectral filtering for general linear dynamical systems. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.

[16] Elad Hazan, Karan Singh, and Cyril Zhang. Learning linear dynamical systems via spectral filtering. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett,

editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.

[17] Gen Li, Yuting Wei, Yuejie Chi, Yuantao Gu, and Yuxin Chen. Softmax policy gradient methods can take exponential time to converge, 2021.

[18] Boyi Liu, Qi Cai, Zhuoran Yang, and Zhaoran Wang. Neural trust region/proximal policy optimization attains globally optimal policy. In *Advances in Neural Information Processing Systems 32*, pages 10565–10576. Curran Associates, Inc., 2019.

[19] Jincheng Mei, Chenjun Xiao, Csaba Szepesvari, and Dale Schuurmans. On the global convergence rates of softmax policy gradient methods. In *Proceedings of Machine Learning and Systems*, pages 10170–10179. PMLR, 2020.

[20] R. Postoyan, L. Buşoniu, D. Nešić, and J. Daafouz. Stability analysis of discrete-time infinite-horizon optimal control with discounted cost. *IEEE Transactions on Automatic Control*, 62:2736–2749, 2017.

[21] Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman, Emanuel Todorov, and Sergey Levine. Learning Complex Dexterous Manipulation with Deep Reinforcement Learning and Demonstrations. In *Proceedings of Robotics: Science and Systems (RSS)*, 2018.

[22] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 1889–1897. PMLR, 2015.

[23] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *CoRR*, abs/1707.06347, 2017.

[24] David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George van den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering the game of go with deep neural networks and tree search. *Nature*, 529(7587):484–489, Jan 2016.

[25] David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dharshan Kumaran, Thore Graepel, Timothy Lillicrap, Karen Simonyan, and Demis Hassabis. A general reinforcement learning algorithm that masters chess, shogi, and go through self-play. *Science*, 362(6419):1140–1144, 2018.

[26] David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller. Deterministic policy gradient algorithms. In Eric P. Xing and Tony Jebara, editors, *Proceedings of the 31st International Conference on Machine Learning*, volume 32 of *Proceedings of Machine Learning Research*, pages 387–395, Bejing, China, 22–24 Jun 2014. PMLR.

[27] A.A. Stoorvogel and A. Saberi. Continuity properties of solutions to H2 and H∞ Riccati equations. *Systems & Control Letters*, 27(4):209–222, 1996.

[28] Richard Sutton, David Mcallester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. *Adv. Neural Inf. Process. Syst*, 12:1057–1063, 2000.

[29] Yunhao Tang, Mark Rowland, Rémi Munos, and Michal Valko. Taylor expansion of discount factors, 2021.

[30] Stephen Tu and Benjamin Recht. Least-squares temporal difference learning for the linear quadratic regulator. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80, pages 5005–5014. PMLR, 2018.

[31] Oriol Vinyals, Igor Babuschkin, Wojciech M. Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H. Choi, Richard Powell, Timo Ewalds, Petko Georgiev, Junhyuk Oh, Dan Horgan, Manuel Kroiss, Ivo Danihelka, Aja Huang, Laurent Sifre, Trevor Cai, John P. Agapiou, Max Jaderberg, Alexander S. Vezhnevets, Rémi Leblond, Tobias Pohlen, Valentin Dalibard, David Budden, Yury Sulsky, James Molloy, Tom L. Paine, Caglar Gulcehre, Ziyu Wang, Tobias Pfaff, Yuhuai Wu, Roman Ring, Dani Yogatama, Dario Wünsch, Katrina McKinney, Oliver Smith, Tom Schaul, Timothy Lillicrap, Koray Kavukcuoglu, Demis Hassabis, Chris Apps, and David Silver. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575(7782):350–354, Nov 2019.

[32] Lingxiao Wang, Qi Cai, Zhuoran Yang, and Zhaoran Wang. Neural policy gradient methods: Global optimality and rates of convergence. In *International Conference on Learning Representations*, 2020.

[33] Zhuoran Yang, Yongxin Chen, Mingyi Hong, and Zhaoran Wang. Provably global convergence of actor-critic: A case for linear quadratic regulator with ergodic cost. In *Advances in Neural Information Processing Systems 32*, pages 8353–8365. Curran Associates, Inc., 2019.

[34] Gilad Yehudai and Ohad Shamir. On the power and limitations of random features for understanding neural networks. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
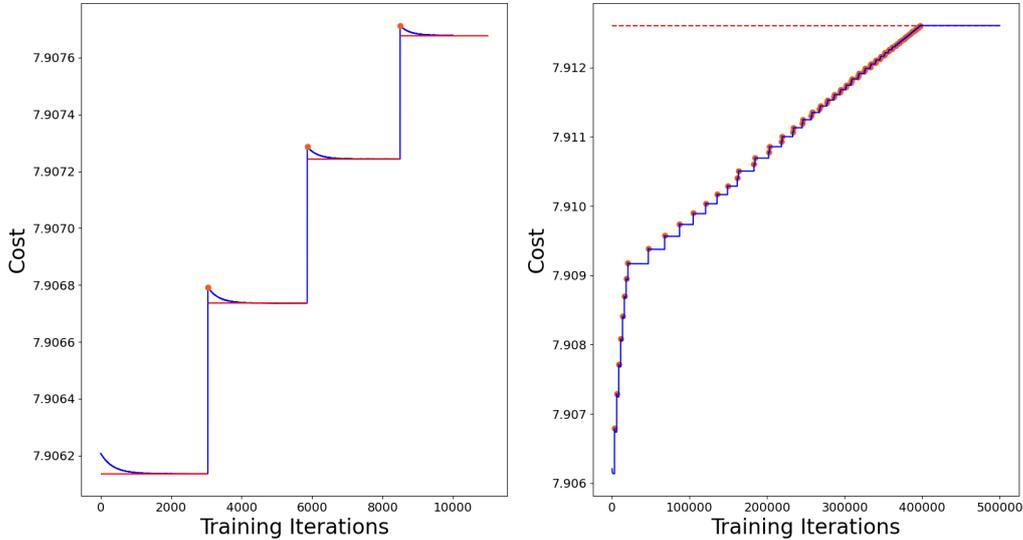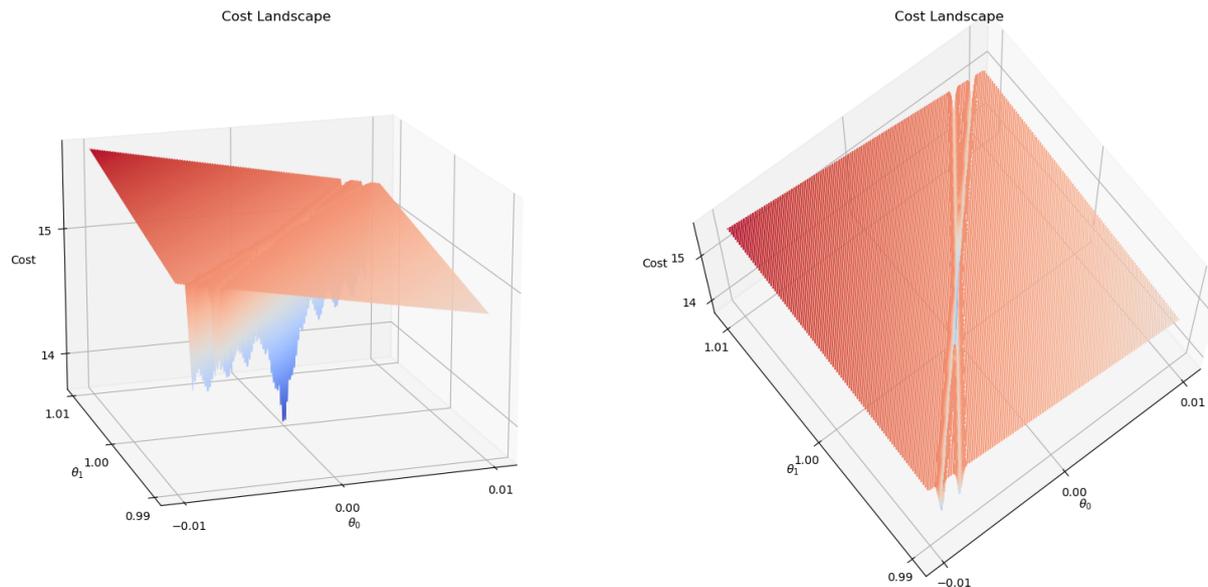
A.1. **Running Algorithm 1.**



FIGURE 2. Evolution of the cost during Algorithm 1 as a function of training iterations. The blue line is the cost during training of the homotopy algorithm in a similar setting to the one described in Section 3.1. The orange dots represent a point at which $\gamma$ was updated. The red lines are the optimal costs for the various values of $\gamma$. The graph on the left is a close-up of the cost between iterations $0$ and $10^5$, the one on the right represents the full training process. The dashed red line on the plot on the right is the optimal cost for $\gamma = 1$.
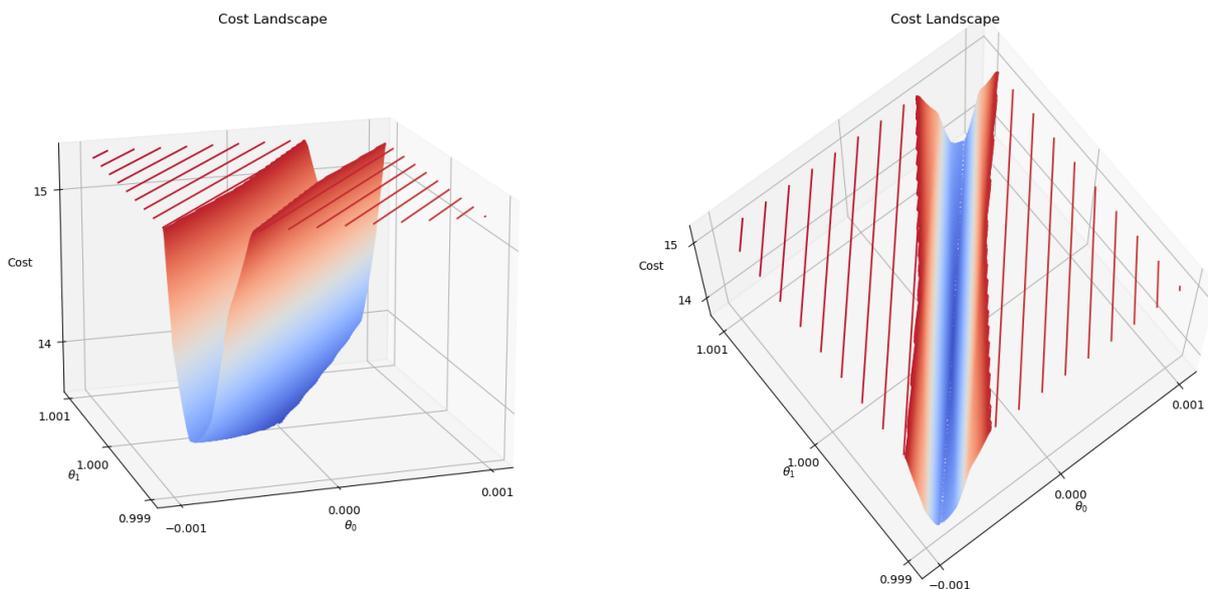
We run Algorithm 1 to solve the LQR problem (2.1), (2.2) in dimension $n = m = 1$ with $(A, B) = (0.1, 1)$ and $(Q, R) = (1, 0.1)$ with the parametric policy class described in Section 3.1 and setting $\delta = 5\text{e-}4$. We do so by writing a policy using PyTorch and updating the parameters via gradient descent with a stepsize of 1e-3. We increased $\gamma$ by 0.02 whenever the absolute values of both partial derivatives with respect to the parameters $\vartheta_0, \vartheta_1$ were less than a tolerance hyperparameter, in our case set to 1e-4. Note that the initial states $x_0$ and $y_0$ described in Section 3.1 can be explicitly calculated to be $-117/59$ and $-588/295$, respectively.

We trained the model for a total of 500000 episodes. During each episode, we interacted with the LQR system for a fixed number of steps, in our case, 5. Since the value for $\varepsilon$ determined by our model is negligibly small (*i.e.*, $\varepsilon < 1\text{e-}12$), we simultaneously simulate the roll-outs from $x_0$ and $y_0$ during training and update using the total sum of discounted costs experienced in both roll-outs. We update the parameters after every episode. To highlight the homotopy algorithm, we initialize our policy with the parameters $\vartheta = (0, 0)$ which correspond to the $\gamma = 0$ optimal policy.

14

## A.2. Counterexample.



(A) Notice, in the top-right figure, the two "good" directions in parameter space.



(B) Close-up of the above plots.

FIGURE 3. Landscape of the cost function from the example of Section 3.1 in a neighborhood around $\vartheta = (0, 1)$. The local minimum at $(0, 1)$ can be clearly seen in these plots.

We compute the landscape of the cost function to verify numerically the result of Proposition 3.4 and represent the lanscape in Figure 3. These results were generated by evaluating the cost experienced by the policy with parameters $(\vartheta_0, \vartheta_1)$ in a neighborhood of $(0, 1)$ for the choice of parameters in Section 3.1 and setting $\varepsilon = 0$ and $\delta = $ 5e-4. The cost of a policy was evaluated by rolling out the policy from both $x_0$ and $y_0$ for 5 steps and computing the discounted sum of the instantaneous costs. Plots were made using the contour plot method from matplotlib.pyplot.

Before stating and proving our results, we give a roadmap for the structure of the proof. In Lemma B.1 we show that the cost associated to trajectories starting at $x_0$ or $y_0$ is lower-bounded locally by the $\ell_1$ norm of the perturbation of the parameters $\|d\vartheta\|_1$ outside of 2 cone in parameter space (one for each initial condition). Then, in Lemma B.2 we show that a similar lower bound – with a reduced constant – holds when the perturbation $d\vartheta$ lies in one such cone (*e.g.* the one associated to $x_0$) but outside of the other (*e.g.* the one associated to $y_0$), *i.e.*, in such regions the positive contribution to the cost of one of the trajectories outweighs the one of the other. Combining this with Lemma B.3 showing that for a choice of parameters the cones do not overlap, we obtain in Corollary B.4 that the *joint* contribution to the cost function associated to trajectories starting at $x_0, y_0$ is lower-bounded locally by (a constant times) the $\ell_1$ norm of the perturbation of the parameters $\|d\vartheta\|_1$. We conclude by further combining this with Lemma B.5, showing that the contribution to the calculation of the uniform measure is negligible, so that the lower bound stated above holds for $\mu_0$.

Throughout, we define

$$C[x_0](\vartheta) := \sum_{t=0}^{\infty} \gamma^t \left(x_t(\vartheta)^\mathsf{T} Q x_t(\vartheta) + \pi_\vartheta(x_t(\vartheta))^\mathsf{T} R \pi_\vartheta(x_t(\vartheta))\right)$$

where $x_t(d\vartheta) \in \mathbb{R}^n$ denotes the trajectory of the LQR with control $\pi_{\vartheta+d\vartheta}$ for $\vartheta = (0, 1)$ and initial condition $x_0$.

*Notation.* We will use $z_t$ to denote the $t$'th state of the unperturbed trajectory and we will use $z_t(d\vartheta)$ for the corresponding state of the perturbed trajectory (*i.e.*, generated by policy $\pi_{\vartheta+d\vartheta}$).

**Lemma B.1.** *Let $\vartheta = (0, 1)$ and let $z_0 \in (x_0, y_0)$ defined in Section 3.1. For any $\alpha > 0$, there exists $\delta > 0$, $\mathcal{K} > 0$ such that for any $i \in \{2, 3\}$ and $d\vartheta = (d\vartheta_0, d\vartheta_1)$ satisfying*

$$|z_0 d\vartheta_0 + F(z_0) d\vartheta_1| \geq \alpha \|d\vartheta\|_1,$$

*the cost difference between the contribution to the cost function using the perturbed and unperturbed policies for trajectories starting from $z_0$ is strictly positive and bounded from below:*

$$C[z_0](\vartheta + d\vartheta) - C[z_0](\vartheta) \geq \frac{3\mathcal{K}}{4\delta} \|d\vartheta\|_1 \,.$$

*Proof of Lemma B.1.* Throughout, let $z, a, \omega$ be generic variables that can be replaced by either $(x, a_2, \omega_2)$ or $(y, a_3, \omega_3)$. We first consider how the cost of a trajectory changes for any infinitesimal variation of the parameters $d\vartheta = (d\vartheta_0, d\vartheta_1)$. To do so we calculate the first order corrections (in $d\vartheta$) to the first two timesteps of the trajectory $z_t(d\vartheta)$ starting at $z_0$. For the first timestep, since along the unperturbed trajectory we have $z_1 = F(z_0) = a$, we have

$$z_1(d\vartheta) = d\vartheta_0 z_0 + (1 + d\vartheta_1) F(z_0)$$
$$= z_1 + (a d\vartheta_1 + z_0 d\vartheta_0), \tag{B.1}$$

where here and throughout this section $z_t = z_t(0)$ denotes the unperturbed trajectory. For the second timestep,

$$z_2(d\vartheta) = d\vartheta_0 z_1(d\vartheta) + (1 + d\vartheta_1)\left(\omega_0 |z_1(d\vartheta_1)| + \omega \Lambda_{a,\delta}(z_1(d\vartheta_1))\right)$$

$$= d\vartheta_0 z_1 + (1 + d\vartheta_1)\bigg(\omega_0 |z_1| + \omega_0(a d\vartheta_1 + z_0 d\vartheta_0)$$

$$\qquad + \omega\left(\Lambda_{a,\delta}(z_1) + \Lambda'_{a,\delta}(z_1)(a d\vartheta_1 + z_0 d\vartheta_0)\right)\bigg)$$

$$= a d\vartheta_0 + (1 + d\vartheta_1)\left(z_2 + \omega_0(a d\vartheta_1 + z_0 d\vartheta_0) - \frac{\omega}{\delta}|a d\vartheta_1 + z_0 d\vartheta_0|\right)$$

$$= z_2 + \left(z_2 d\vartheta_1 + a d\vartheta_0 + \omega_0(a d\vartheta_1 + z_0 d\vartheta_0) - \frac{\omega}{\delta}|a d\vartheta_1 + z_0 d\vartheta_0|\right)$$

$$=: z_2 + dz_2$$

where $dz_2 := z_2 d\vartheta_1 + a d\vartheta_0 + \omega_0(a d\vartheta_1 + z_0 d\vartheta_0) - \frac{\omega}{\delta}|a d\vartheta_1 + z_0 d\vartheta_0|$ is the first order correction to the second timestep of the trajectory. We then note that since $z_2 \in (-1, 0)$ due to our choice of parameters, the future cost from $z_2$ is given by

$$\gamma^2 C[z_2](\vartheta + d\vartheta) = \gamma^2 \sum_{t=0}^{\infty} \gamma^t \left( Q(z_{t+2})^2 + R(u_{t+2})^2 \right) = \gamma^2 \sum_{t=0}^{\infty} \gamma^t \left( Q(z_{t+2})^2 + R(z_{t+3})^2 \right)$$

$$= \gamma^2 \left( Q(z_2)^2 + \sum_{t=0}^{\infty} \gamma^t \left( \gamma Q(z_{t+3})^2 + R(z_{t+3})^2 \right) \right)$$

$$= \gamma^2 \left( Q(z_2)^2 + (\gamma Q + R) \sum_{t=0}^{\infty} \gamma^t ((d\vartheta_0 + \omega_0(1 + d\vartheta_1))^{t+1} z_2)^2 \right)$$

$$= \gamma^2 (z_2)^2 \left( Q + \frac{(\gamma Q + R)(\omega_0 + (d\vartheta_0 + \omega_0 d\vartheta_1))^2}{1 - \gamma (\omega_0 + (d\vartheta_0 + \omega_0 d\vartheta_1))^2} \right)$$

Since $z_2(d\vartheta) \in (-1, 0)$ by construction as well, the above expression also holds for $z_2(d\vartheta)$. Therefore, since $Q = 1$ and $1 + R\omega_0^2 > 0$, the cost of trajectory from $z_2$ onward will vary (to first order in $d\vartheta$) by

$$\gamma^2 z_2 (d\vartheta)^2 \frac{1 + R(\omega_0 + (d\vartheta_0 + \omega_0 d\vartheta_1))^2}{1 - \gamma (\omega_0 + (d\vartheta_0 + \omega_0 d\vartheta_1))^2} = \gamma^2 \frac{1 + R(\omega_0 + (d\vartheta_0 + \omega_0 d\vartheta_1))^2}{1 - \gamma (\omega_0 + (d\vartheta_0 + \omega_0 d\vartheta_1))^2} \left( (z_2)^2 + 2 z_2 dz_2 \right)$$

$$= \frac{\gamma^2 \left( 1 + R\omega_0^2 \right)}{1 - \gamma \omega_0^2} \left( (z_2)^2 + 2 z_2 dz_2 \right)$$

$$+ \frac{2\gamma^2 \omega_0 (\gamma + R)}{\left( 1 - \gamma \omega_0^2 \right)^2} (d\vartheta_0 + \omega_0 d\vartheta_1) (z_2)^2$$

$$= \frac{\gamma^2}{1 - \gamma \omega_0^2} \left( \left( 1 + R\omega_0^2 \right) (z_2)^2 + \left( 1 + R\omega_0^2 \right) 2 z_2 dz_2 \right.$$

$$\left. + \frac{2\omega_0(\gamma + R)}{1 - \gamma \omega_0^2} (d\vartheta_0 + \omega_0 d\vartheta_1) (z_2)^2 \right)$$

Inserting the expression of $dz_2$,

$$\text{RHS} = \frac{\gamma^2 \left( 1 + R\omega_0^2 \right)}{1 - \gamma \omega_0^2} \left( (z_2)^2 + 2 z_2 \left( - \frac{\omega}{\delta} |a d\vartheta_1 + z_0 d\vartheta_0| + z_2 d\vartheta_1 + a d\vartheta_0 \right. \right.$$

$$\left. \left. + \omega_0(a d\vartheta_1 + z_0 d\vartheta_0) + \frac{\omega_0(\gamma + R)}{\left( 1 + R\omega_0^2 \right) \left( 1 - \gamma \omega_0^2 \right)} z_2(d\vartheta_0 + \omega_0 d\vartheta_1) \right) \right) \quad \text{(B.2)}$$

$$\geq \frac{\gamma^2 \left( 1 + R\omega_0^2 \right)}{1 - \gamma \omega_0^2} \left( (z_2)^2 - 2 z_2 \frac{3\omega\alpha}{4\delta} \|d\vartheta\|_1 \right)$$

$$= \frac{\gamma^2 \left( 1 + R\omega_0^2 \right)}{1 - \gamma \omega_0^2} \left( (z_2)^2 + 2|z_2| \frac{3\omega\alpha}{4\delta} \|d\vartheta\|_1 \right),$$

since $z_2$ is negative. The inequality comes from the fact that whenever $|z_0 d\vartheta_0 + a d\vartheta_1| \geq \alpha \|d\vartheta\|_1$ we can choose $\delta > 0$ sufficiently small such that

$$\left| \frac{\omega |a d\vartheta_1 + z_0 d\vartheta_0|}{4\delta} \right| \geq \left| \frac{\omega\alpha \|d\vartheta\|_1}{4\delta} \right|$$

$$\geq 4 \max \left( \left| z_2 d\vartheta_1 + a d\vartheta_0 \right|, \left| \omega_0(a d\vartheta_1 + z_0 d\vartheta_0) \right|, \quad \text{(B.3)} \right.$$

$$\left. \left| \frac{\omega_0(\gamma + R)}{\left( 1 + R\omega_0^2 \right) \left( 1 - \gamma \omega_0^2 \right)} z_2(d\vartheta_0 + \omega_0 d\vartheta_1) \right|,$$

$$\left|\frac{(\gamma + R)a|ad\vartheta_1 + z_0 d\vartheta_0|}{z_2\gamma^2\left(1 + R\omega_0^2\right)/\left(1 - \gamma\omega_0^2\right)}\right|\right)$$

Bounding the right-hand-side by triangle inequality, the above is satisfied when

$$\delta \leq \frac{|\omega\alpha|}{4\max\{M_0, M_1\}},$$

where

$$M_0 = 4\max\left\{|z_2|, \ |\omega_0 a|, \ \left|\frac{|z_2|\omega_0^2(\gamma + R)}{\left(1 + R\omega_0^2\right)\left(1 - \gamma\omega_0^2\right)}\right|, \ \left|\frac{a^2(\gamma + R)(1 - \gamma\omega_0^2)}{|z_2|\gamma^2(1 + R\omega_0^2)}\right|\right\}$$

$$M_1 = 4\max\left\{|a|, \ |\omega_0 z_0|, \ \left|\frac{z_2\omega_0(\gamma + R)}{\left(1 + R\omega_0^2\right)\left(1 - \gamma\omega_0^2\right)}\right|, \ \left|\frac{az_0(\gamma + R)(1 - \gamma\omega_0^2)}{z_2\gamma^2(1 + R\omega_0^2)}\right|\right\}.$$

In other words, we have shown that for all $d\vartheta \in \left\{(d\vartheta_0, d\vartheta_1) \in \mathbb{R}^2 \mid |z_0 d\vartheta_0 + F(z_0)d\vartheta_1| \geq \alpha\|d\vartheta\|_1\right\}$,

$$\begin{aligned}
C[z_0](\vartheta + d\vartheta) - C[z_0](\vartheta) &= \left(z_0^2 + (\gamma + R)z_1(d\vartheta)^2 + \gamma R z_2(d\vartheta)^2 + \gamma^2 C[z_2(d\vartheta)](\vartheta + d\vartheta)\right) \\
&\quad - \left(z_0^2 + (\gamma + R)z_1^2 + \gamma R z_2^2 + \gamma^2 C[z_2](\vartheta)\right) \\
&\geq \gamma^2\left(C[z_2(d\vartheta)](\vartheta + d\vartheta) - C[z_2](\vartheta)\right) - 2(\gamma + R)|z_1||z_1(d\vartheta) - z_1| \\
&= \gamma^2\left(C[z_2(d\vartheta)](\vartheta + d\vartheta) - C[z_2](\vartheta)\right) - 2(\gamma + R)a\left(ad\vartheta_1 + z_0 d\vartheta_0\right) \\
&\geq 3\frac{2|z_2|\gamma^2\left(1 + R\omega_0^2\right)\omega\alpha}{\left(1 - \gamma\omega_0^2\right)4\delta}\|d\vartheta\|_1 \quad\quad\quad\quad\quad\text{(B.4)} \\
&=: \frac{3\mathcal{K}}{4\delta}\|d\vartheta\|_1.
\end{aligned}$$

for

$$\mathcal{K} := \frac{2|z_2|\gamma^2\left(1 + R\omega_0^2\right)\omega\alpha}{1 - \gamma\omega_0^2} \quad\quad\quad\quad\quad\text{(B.5)}$$

where the first inequality comes from our assumption that $d\vartheta$ lies outside of the cone in parameter space ($|z_0 d\vartheta_0 + F(z_0)d\vartheta_1| \geq \alpha\|d\vartheta\|_1$) and the last one results from our choice of $\delta$ small enough from (B.3). $\qquad\square$

Based on the above, the only variation of the parameters that does not increase the cost of the policy is when $d\vartheta$ is inside the cone

$$|z_0 d\vartheta_0 + F(z_0)d\vartheta_1| \leq \alpha\|d\vartheta\|_1,$$

for $z_0 \in \{x_0, y_0\}$. For the statement of the next lemma, for a choice of $z_0 \in \{x_0, y_0\}$ we define $\bar{z}_0 \in \{x_0, y_0\} \setminus z_0$ as the *other* initial condition and by $\bar{z}_t(d\vartheta)$ the corresponding perturbed trajectory.

**Lemma B.2.** *Suppose $d\vartheta \in \{|z_0 d\vartheta_0 + F(z_0)d\vartheta_1| \leq \alpha\|d\vartheta\|_1\} \cap \{|\bar{z}_0 d\vartheta_0 + F(\bar{z}_0)d\vartheta_1| \geq \alpha\|d\vartheta\|_1\}$. For the same $\delta, \mathcal{K} > 0$ as determined by Lemma B.1, the cost difference is strictly positive and can be bounded by*

$$\int_{\mathbb{R}}(C[x](\vartheta + d\vartheta) - C[x](\vartheta))\mu_{0,pp}(\mathrm{d}x) \geq \frac{1 - \varepsilon}{2}\frac{\mathcal{K}}{12\delta}\|d\vartheta\|_1.$$

*Here $\mu_{0,pp} = \frac{1-\varepsilon}{2}(\delta_{x_0} + \delta_{y_0})$ denotes the pure-point, or discrete, part of $\mu_0$.*

*Proof of Lemma B.2.* The proof proceeds similarly to the proof of Lemma B.1. We start by bounding from below the variation of the cost due to the trajectory $z_t$ starting at $z_0$ with $|z_0 d\vartheta_0 + F(z_0)d\vartheta_1| \leq \alpha\|d\vartheta\|_1$. Recalling the expression for the first-order variation of the cost (B.2) from the proof of Lemma B.1 and that $z_2 < 0$ by our choice of parameters, we bound the negative part of the contribution to the cost function for $t \geq 2$ as:

$$\begin{aligned}
(\gamma^2 C[z_2(d\vartheta)](\vartheta + d\vartheta) &- \gamma^2 C[z_2](\vartheta))_- \\
&= \frac{\gamma^2\left(1 + R\omega_0^2\right)}{1 - \gamma\omega_0^2}2z_2\Big(z_2 d\vartheta_1 + ad\vartheta_0 + \omega_0(ad\vartheta_1 + z_0 d\vartheta_0)
\end{aligned}$$

$$
\left. - \frac{\omega}{\delta}|ad\vartheta_1 + z_0 d\vartheta_0| + \frac{\omega_0(\gamma + R)}{\left(1 + R\omega_0^2\right)\left(1 - \gamma\omega_0^2\right)} z_2(d\vartheta_0 + \omega_0 d\vartheta_1) \right)_+
$$

$$
= \frac{\gamma^2\left(1 + R\omega_0^2\right)}{1 - \gamma\omega_0^2} \frac{\omega_2}{\delta} 2|z_2||ad\vartheta_1 + z_0 d\vartheta_0|
$$

$$
+ \frac{\gamma^2\left(1 + R\omega_0^2\right)}{1 - \gamma\omega_0^2} 2z_2\Big( z_2 d\vartheta_1 + ad\vartheta_0 + \omega_0(ad\vartheta_1 + z_0 d\vartheta_0)
$$

$$
\left. + \frac{\omega_0(\gamma + R)}{\left(1 + R\omega_0^2\right)\left(1 - \gamma\omega_0^2\right)} z_2(d\vartheta_0 + \omega_0 d\vartheta_1) \right)_+
$$

$$
\geq \frac{\gamma^2\left(1 + R\omega_0^2\right)}{1 - \gamma\omega_0^2} 2z_2\Big( z_2 d\vartheta_1 + ad\vartheta_0 + \omega_0(ad\vartheta_1 + z_0 d\vartheta_0)
$$

$$
\left. + \frac{\omega_0(\gamma + R)}{\left(1 + R\omega_0^2\right)\left(1 - \gamma\omega_0^2\right)} z_2(d\vartheta_0 + \omega_0 d\vartheta_1) \right)_+
$$

where $(b)_- := \min(0, b), (b)_+ := \max(0, b)$ denote the positive and negative part of $b \in \mathbb{R}$ respectively. We can assume that the quantity on the LHS above is negative; if it were positive, then the proof is done since we are aiming to lower bound the total difference, so we could ignore this term if it were positive. Since we are using the same $\delta$ as chosen in Lemma B.1, recalling the definition of $\mathcal{K} > 0$ and that $z_2 < 0$ we obtain

$$
\gamma^2\left(C[z_2(d\vartheta)](\vartheta + d\vartheta)| - C[z_2](\vartheta)\right)_- \geq -\frac{2\gamma^2|z_2|\left(1 + R\omega_0^2\right)}{1 - \gamma\omega_0^2}\left|\frac{\omega_2\alpha\|d\vartheta_1\|}{4\delta}\right| = -\frac{\mathcal{K}}{4\delta}\|d\vartheta\|_1 .
$$

Then, recalling that $z_1 = F(z_0) = a$ and $0 < R < \gamma$, the total cost difference from the initial point $z_0$ (to first order) can then be bounded as follows for $\|d\vartheta\|$ sufficiently small:

$$
\begin{aligned}
C[z_0](\vartheta + d\vartheta) - C[z_0](\vartheta) &\geq (C[z_0](\vartheta + d\vartheta) - C[z_0](\vartheta))_- \\
&= \Big( \big(z_0^2 + (\gamma + R)z_1(d\vartheta)^2 + \gamma R z_2(d\vartheta)^2 + \gamma^2 C[z_2(d\vartheta)](\vartheta + d\vartheta)\big) \\
&\quad - \big(z_0^2 + (\gamma + R)z_1^2 + \gamma R z_2^2 + \gamma^2 C[z_2](\vartheta)\big) \Big)_- \\
&\geq -(\gamma + R)\big|z_1(d\vartheta)^2 - z_1^2\big| + \gamma R \left(z_2(d\vartheta)^2 - z_2^2\right)_- \\
&\quad + \gamma^2\left(C[z_2(d\vartheta)](\vartheta + d\vartheta) - C[z_2](\vartheta)\right)_- \\
&\geq -2(\gamma + R)|z_1||z_1(d\vartheta) - z_1| \\
&\quad + 2\gamma^2\left(C[z_2(d\vartheta)](\vartheta + d\vartheta) - C[z_2](\vartheta)\right)_- \\
&\geq -2(\gamma + R)|z_1|\alpha\|d\vartheta\|_1 - \frac{2\mathcal{K}}{4\delta}\|d\vartheta\|_1 \\
&\geq -\frac{2\mathcal{K}}{3\delta}\|d\vartheta\|_1.
\end{aligned}
$$

where in the second inequality we have used that $z_2(d\vartheta)^2 - z_2^2 > 0$. For any $d\vartheta \in \{|z_0 d\vartheta_0 + F(z_0)d\vartheta_1| \leq \alpha\|d\vartheta\|_1\} \cap \{|\bar{z}_0 d\vartheta_0 + F(\bar{z}_0)d\vartheta_1| \geq \alpha\|d\vartheta\|_1\}$, we can then compute that

$$
\int_{\mathbb{R}} \left(C[x](\vartheta + d\vartheta) - C[x](\vartheta)\right) \mu_{0,\mathrm{pp}}(dx)
$$

$$
= \frac{1 - \varepsilon}{2}\left(\int_{\mathbb{R}} C[x](\vartheta + d\vartheta) - C[x](\vartheta)\delta(z_0)\mathrm{d}x + \int_{\mathbb{R}} C[x](\vartheta + d\vartheta) - C[x](\vartheta)\delta(\bar{z}_0)\mathrm{d}x\right)
$$

$$
\geq \frac{1 - \varepsilon}{2}\left(\frac{3\mathcal{K}}{4\delta}\|d\vartheta\|_1 - \frac{2\mathcal{K}}{3\delta}\|d\vartheta\|_1\right)
$$

$$
\geq \frac{1 - \varepsilon}{2} \frac{\mathcal{K}}{12\delta}\|d\vartheta\|_1
$$

where $\mu_{0,\mathrm{pp}}$ denotes the discrete part of the measure $\mu_0$. $\qquad\square$

**Lemma B.3.** *For any $x_0 \neq y_0$ such that $x_0, y_0 < 0$ and $F(x_0), F(y_0) > 0$, there exists an $\alpha > 0$ such that*

$$\{|x_0 d\vartheta_0 + F(x_0)d\vartheta_1| \leq \alpha\|d\vartheta\|_1\} \cap \{|y_0 d\vartheta_0 + F(y_0)d\vartheta_1| \leq \alpha\|d\vartheta\|_1\} = \{(0,0)\}.$$

The above directly implies that

**Corollary B.4.** *For the same $\delta, \mathcal{K} > 0$ as determined by Lemma B.1, and for every $d\vartheta \in \mathbb{R}^d$ with $\|d\vartheta\| < \delta$ we have that*

$$\int_{\mathbb{R}} \left(C[x](\vartheta + d\vartheta) - C[x](\vartheta)\right) \mu_{0,pp}(\mathrm{d}x) \geq \frac{1-\varepsilon}{2} \frac{\mathcal{K}}{12\delta}\|d\vartheta\|_1 \, .$$

*Here $\mu_{0,pp}$ denotes the pure-point, or discrete, part of $\mu_0$.*

*Proof of Corollary B.4.* For our choice of parameters (3.3), Lemma B.3 gives us an $\alpha > 0$ such that the cones determined by $x_0$ and $y_0$ in parameter space do *not* overlap except at $d\vartheta = (0,0)$. By Lemma B.1, for any such $\alpha$ there exists a $\delta > 0$ such that the claim holds. $\square$

*Proof of Lemma B.3.* By assumption, the lines

$$\frac{d\vartheta_1}{d\vartheta_0} = \frac{-x_0}{F(x_0)} \quad \text{and} \quad \frac{d\vartheta_1}{d\vartheta_0} = \frac{-y_0}{F(y_0)}$$

only intersect at the origin. Then, the bisector of the central angle between these two lines can be written as

$$\frac{d\vartheta_1}{d\vartheta_0} = \frac{-\frac{\sqrt{x_0^2 + F(x_0)^2}}{\sqrt{y_0^2 + F(y_0)^2}}y_0 - x_0}{\frac{\sqrt{x_0^2 + F(x_0)^2}}{\sqrt{y_0^2 + F(y_0)^2}}F(y_0) + F(x_0)} =: M > 0 \, .$$

We can assume without loss of generality that the line corresponding to $x_0$ lies *above* the bisector (*i.e.*, $-x_0/F(x_0) > M$). Then, it suffices to choose $\alpha$ such that neither of the sets

$$\{|x_0 d\vartheta_0 + F(x_0)d\vartheta_1| \leq \alpha\|d\vartheta\|_1\} \quad , \qquad \{|y_0 d\vartheta_0 + F(y_0)d\vartheta_1| \leq \alpha\|d\vartheta\|_1\} \qquad \text{(B.6)}$$

intersects the line $\ell_M = \{(d\vartheta_0, d\vartheta_1) \in \mathbb{R}^2 \ : \ d\vartheta_0 = M d\vartheta_1\}$ outside of $(0,0)$.

Notice that the boundaries of the sets in (B.6) that are closest to $\ell_M$ can be explicitly written as the lines

$$\frac{d\vartheta_1}{d\vartheta_0} = \frac{-\alpha - x_0}{\alpha + F(x_0)} \quad \text{and} \quad \frac{d\vartheta_1}{d\vartheta_0} = \frac{\alpha - y_0}{-\alpha + F(y_0)} \, .$$

Then, noting that we can replace $-x_0$ with $|x_0|$ and $-y_0$ with $|y_0|$ since $x_0, y_0 < 0$, we have that the only intersection of the sets in (B.6) with $\ell_M$ is $(0,0)$ if

$$M < \frac{-\alpha - x_0}{\alpha + F(x_0)} \quad \text{and} \quad M > \frac{\alpha - y_0}{-\alpha + F(y_0)} \, ,$$

which is equivalent to the conditions

$$\alpha < \frac{|x_0| - MF(x_0)}{1 + M} \quad \text{and} \quad \alpha < \frac{MF(y_0) - |y_0|}{1 + M} \, .$$

Since we have assumed that

$$\frac{|x_0|}{F(x_0)} > M > \frac{|y_0|}{F(y_0)}$$

and that $F(x_0), F(y_0) > 0$, we know that the bounds in both constraints are strictly positive. Thus, we can choose any

$$\alpha < \frac{1}{2} \min\left(\frac{|x_0| - MF(x_0)}{1 + M}, \frac{MF(y_0) - |y_0|}{1 + M}\right)$$

and this proves the claim. $\square$

**Lemma B.5.** *There exists $\mathcal{K}' > 0$ such that for any initial point $u_0 \sim U(-5,5)$, the average cost difference can be bounded from above by*

$$\left| \int_{-5}^{5} (C[u](\vartheta + d\vartheta) - C[u](\vartheta)) \frac{du}{10} \right| \leq \mathcal{K}' \|d\vartheta\|_1 .$$

*Proof of Lemma B.5.* For any $u_0 \in [-5, 5]$, we can write the cost difference as

$$
\begin{aligned}
\Delta C[u_0](d\vartheta) &:= C[u_0](\vartheta + d\vartheta) - C[u_0](\vartheta) \\
&= \left( u_0^2 + Ru_1(d\vartheta)^2 + \gamma C[u_1(d\vartheta)](\vartheta + d\vartheta) \right) - \left( u_0^2 + R(u_1)^2 + \gamma C[u_1](\vartheta) \right) \\
&= \gamma \Delta C[u_1(d\vartheta)](d\vartheta) + R(u_1(d\vartheta)^2 - (u_1)^2) + \gamma \left( C[u_1(d\vartheta)](\vartheta) - C[u_1](\vartheta) \right) \\
&\leq \gamma \Delta C[u_0](d\vartheta) + R(u_1(d\vartheta)^2 - (u_1)^2) + \gamma \left( C[u_1(d\vartheta)](\vartheta) - C[u_1](\vartheta) \right) \\
&\leq \gamma \Delta C[u_0](d\vartheta) + 2Ru_1|F(u_0)d\vartheta_1 + u_0 d\vartheta_0| + \gamma \left( C[u_1(d\vartheta)](\vartheta) - C[u_1](\vartheta) \right) \\
&\leq \gamma \Delta C[u_0](d\vartheta) + 2R(2.5)(5)\|d\vartheta\|_1 + \gamma \left( C[u_1(d\vartheta)](\vartheta) - C[u_1](\vartheta) \right),
\end{aligned}
$$

since we assume, in the construction of the example, that $0 < R < \gamma$. Iterating the bound established above we obtain

$$
\begin{aligned}
\Delta C[u_0](d\vartheta) &\leq \sum_{t=0}^{\infty} \gamma^t \left( 2R(2.5)(5)\|d\vartheta\|_1 + \gamma \left( C[u_1(d\vartheta)](\vartheta) - C[u_1](\vartheta) \right) \right) \\
&= \frac{2R(2.5)(5)\|d\vartheta\|_1 + \gamma \left( C[u_1(d\vartheta)](\vartheta) - C[u_1](\vartheta) \right)}{1 - \gamma}
\end{aligned}
$$

(B.7)

We now proceed to estimate the numerator of the above expression. To do so, let $L = \frac{|\omega_0| + |\omega_1|}{\delta_1}$ denote the Lipschitz constant of $F$ on $(-5, 5)$. Notice that in the interval $(-1, 1)$, the Lipschitz constant of $F$ restricted to this interval is $|\omega_0|$ and that on $(-5, 5)$ $F$ is non-expansive ($|F(x)| \leq |x|$) and $F^2$ is contractive ($\exists k \in [0, 1)$ s.t. $|F^2(x)| \leq k|x|$) by construction. Since $5|\omega_0|^2 > 1$ but $5|\omega_0|^3 < 1$, any initial point will take at most three steps before it is mapped into the region with Lipschitz constant $|\omega_0|$.

$$
\begin{aligned}
C[u_1(d\vartheta)](\vartheta) - C[u_1](\vartheta) &\leq (1 + R) \sum_{t=0}^{\gamma} \gamma^t \left( \underbrace{F \circ \cdots \circ F}_{t \text{ times}}(u_1(d\vartheta))^2 - \underbrace{F \circ \cdots \circ F}_{t \text{ times}}(u_1)^2 \right) \\
&\leq (1 + R) \left( \sum_{t=0}^{3} \gamma^t L^{2t} \left( u_1(d\vartheta)^2 - (u_1)^2 \right) + \gamma^4 \sum_{t=0}^{\infty} \gamma^t |\omega_0|^{2t} \left( u_1(d\vartheta)^2 - (u_1)^2 \right) \right) \\
&\leq (1 + R) \left( L^6 + \frac{\gamma^4}{1 - \gamma\omega_0^2} \right) \left( u_1(d\vartheta)^2 - (u_1)^2 \right) \\
&= (1 + R) \left( L^6 + \frac{\gamma^4}{1 - \gamma\omega_0^2} \right) \left( 2u_1 \left( F(u_0)d\vartheta_1 + u_0 d\vartheta_0 \right) + (F(u_0)d\vartheta_1 + u_0 d\vartheta_0)^2 \right) \\
&\leq L'' \|d\vartheta\|_1,
\end{aligned}
$$

for $\|d\vartheta\|$ small enough, taking only the first order term. Combining the above with (B.7) and defining $L' = \gamma L'' + 25R$ finally gives

$$\Delta C[u_0](d\vartheta) \leq \mathcal{K}' \|d\vartheta\|_1, \qquad \text{for} \qquad \mathcal{K}' := \frac{L'}{1 - \gamma} \|d\vartheta\|_1, \tag{B.8}$$

as claimed. $\qquad \square$

**Proposition 3.4.** *Fix $\gamma \in (0, 1)$ and $0 < R < \gamma$. For the LQR problem with $(A, B) = (0, 1)$, $(Q, R) = (1, R)$, policy class (3.2) with parameters chosen in (3.3), and initial state distribution $\mu_0$ from (3.4), there exists $\delta > 0$ and $\varepsilon > 0$ such that the point in parameter space $\vartheta = (0, 1)$ is a local minimum of the LQR cost function.*

*Proof of Proposition 3.4.* Recalling the definitions of $\mathcal{K}, \mathcal{K}'$ from (B.5) and (B.8) respectively, for any

$$0 < \varepsilon < \frac{\mathcal{K}}{24\delta\mathcal{K}' + \mathcal{K}}, \tag{B.9}$$

and for $d\vartheta \neq 0$ sufficiently small, since $\mu_0 = \mu_{0,\text{pp}} + \frac{\varepsilon}{10}$ for $\mu_{0,\text{pp}} = \frac{1-\varepsilon}{2}\delta_{x_0} + \frac{1-\varepsilon}{2}\delta_{y_0}$ we have that

$$\int_{-5}^{5} (C(\vartheta + d\vartheta) - C(\vartheta)) \mathrm{d}\mu_0(x) =$$

$$= \frac{1-\varepsilon}{2} \int_{-5}^{5} (C(\vartheta + d\vartheta) - C(\vartheta))\, \mathrm{d}\mu_{0,\text{pp}} + \varepsilon \int_{-5}^{5} (C(\vartheta + d\vartheta) - C(\vartheta))\, \frac{\mathrm{d}x}{10}$$

$$\geq \frac{1-\varepsilon}{2} \frac{\mathcal{K}}{12\delta} \|d\vartheta\|_1 - \varepsilon \mathcal{K}' \|d\vartheta\|_1$$

$$= \frac{\mathcal{K}}{24\delta} \|d\vartheta\|_1 - \varepsilon \left( \mathcal{K}' + \frac{\mathcal{K}}{24\delta} \right) \|d\vartheta\|_1 > 0,$$

where in the second line we have combined Corollary B.4 and Lemma B.5. Thus, we see that the point in parameter space corresponding to $\vartheta = (0, 1)$ is a local minimum of the cost function.

$\square$

## APPENDIX C. PROOFS OF CONVERGENCE OF THE HOMOTOPY ALGORITHM

For notational convenience, we define

$$\mathcal{A} := A + BK$$
$$\overline{Q} := Q + K^\mathsf{T} RK$$

We use $\Delta\vartheta = (\Delta\vartheta_1, \ldots \Delta\vartheta_d)$ to denote a perturbation of the parameters. We study the behavior of the parameters in the setting where $\vartheta = \vartheta^* + \Delta\vartheta$.

**Lemma C.1.** *For a fixed quadruple $(A, B, Q, R)$ that is controllable and observable, let $P_\gamma$ denote the unique positive semi-definite solution to the discounted discrete algebraic Riccati equation (discounted DARE) with discount factor $\gamma \in [0, 1]$*

$$P_\gamma = \gamma A^\mathsf{T} P_\gamma A - \gamma^2 A^\mathsf{T} P_\gamma B \left( R + \gamma B^\mathsf{T} P_\gamma B \right)^{-1} B^\mathsf{T} P_\gamma A + Q$$

*Then, $f : (0, 1) \to \mathbb{R}^{n \times n}$, $f(\gamma) = P_\gamma$ is a continuous function.*

*Proof of Lemma C.1.* For any $\gamma \in [0, 1]$, since $(A, B)$ is controllable, we know that $(\sqrt{\gamma}A, B)$ also controllable and thus stabilizable. Similarly, $(\sqrt{\gamma}A, D)$ is also observable, and by the positive definiteness of $R$ and positive semidefiniteness of $P_\gamma$, $(R + \gamma B^\mathsf{T} P_\gamma B)^{-1}$ exists and is continuous. By [27, Theorem 2.4], we conclude that $P_\gamma$ is a continuous function of $\gamma \in [0, 1]$. $\square$

Since the optimal policy is a continuous function of $P_\gamma$, Lemma C.1 immediately gives us

**Corollary C.2.** *The optimal policy $K_\gamma^*$ is continuous in $\gamma$.*

We preface the next two lemmas by reminding the reader that the expansion of the cost function that we consider is around the *optimal* policy $K_\gamma^*$. However, since this notation is a bit tedious, we simply use $K$ to refer to this optimal policy in our coming computations. Here, we are using Assumption 2 when we stipulate that $\pi_{\vartheta^*} = K_\gamma^*$.

**Lemma C.3.** *The first-order term in the Taylor expansion of the cost (with respect to the parameters) evaluates to zero at optimality.*

Let $\pi$ denote the policy and let $\tilde{\pi}$ denote the perturbation of $\pi$ (*i.e.*, $\pi_{\vartheta^* + \Delta\vartheta} = \pi_{\vartheta^*} + \tilde{\pi}$). Furthermore, recall that our policy representation is linear in the parameters. We first compute the zero'th, first, and second order terms of the expansion of the trajectory.

$$x_t^{(0)} = \mathcal{A}^t x_0$$

$$x_t^{(1)} = \sum_{l=0}^{t-1} \mathcal{A}^l B \tilde{\pi}(x_{t-l-1}^{(0)})$$

$$x_t^{(2)} = \sum_{l=1}^{t-1} \mathcal{A}^l B \left( \tilde{\pi}(x_{t-l-1}^{(0)} + x_{t-l-1}^{(1)}) - \tilde{\pi}(x_{t-l-1}^{(0)}) \right)$$

where we ignore all non-second-order terms to reach the last equality. We note that, by the Lipschitz continuity of policy $\pi$, all remaining terms in the first order term of the Volterra expansion for the given system are of order higher than 1.

We can bound the Euclidean norm of the second order term in the expansion using the assumed Lipschitz continuity of $\pi$:

$$
\begin{aligned}
\|x_t^{(2)}\| &= \left\| \sum_{l=1}^{t-1} \mathcal{A}^l B \left( \tilde{\pi}(x_{t-l-1}^{(0)} + x_{t-l-1}^{(1)}) - \tilde{\pi}(x_{t-l-1}^{(0)}) \right) \right\| \\
&\le \sum_{l=1}^{t-1} \|\mathcal{A}^l B\| \|\tilde{\pi}(x_{t-l-1}^{(0)} + x_{t-l-1}^{(1)}) - \tilde{\pi}(x_{t-l-1}^{(0)})\| \\
&\le \mathrm{Lip}(\tilde{\pi}) \sum_{l=1}^{t-1} \|\mathcal{A}^l B\| \|x_{t-l-1}^{(1)}\|
\end{aligned}
$$

*Proof of Lemma C.3.* We consider a perturbation of the parameters of the policy. For any $x_0 \in \mathrm{supp}\,(\varrho_0)$, the first order term of the Taylor approximation is

$$
\begin{aligned}
C^{(1)}[x_0] &= \sum_{t=1}^{\infty} \gamma^t \left( \left( x_t^{(1)} \right)^{\mathsf{T}} \overline{Q} x_t^{(0)} + \left( x_t^{(0)} \right)^{\mathsf{T}} \overline{Q} \left( x_t^{(1)} \right) \right) \\
&\quad + \sum_{t=0}^{\infty} \gamma^t \left( \tilde{\pi}(x_t^{(0)})^{\mathsf{T}} R K x_t^{(0)} + \left( x_t^{(0)} \right)^{\mathsf{T}} K^{\mathsf{T}} R \tilde{\pi}(x_t^{(0)}) \right) \\
&= \sum_{t=1}^{\infty} \gamma^t \sum_{l=0}^{t-1} \left( \mathcal{A}^{t-l-1} B \tilde{\pi}(x_l^{(0)}) \right)^{\mathsf{T}} \overline{Q} \mathcal{A}^t x_0 + \sum_{t=1}^{\infty} \gamma^t \left( \mathcal{A}^t x_0 \right)^{\mathsf{T}} \overline{Q} \sum_{l=0}^{t-1} \left( \mathcal{A}^{t-l-1} B \tilde{\pi}(x_l^{(0)}) \right) \\
&\quad + \sum_{t=0}^{\infty} \gamma^t \left( \tilde{\pi}(x_t^{(0)})^{\mathsf{T}} R K \mathcal{A}^t x_0 + \left( \mathcal{A}^t x_0 \right)^{\mathsf{T}} K^{\mathsf{T}} R \tilde{\pi}(x_t^{(0)}) \right) \\
&= \sum_{t=1}^{\infty} \gamma^t \left( \sum_{l=0}^{t-1} \left( \mathcal{A}^{t-l-1} B \tilde{\pi}(x_l^{(0)}) \right)^{\mathsf{T}} \overline{Q} \mathcal{A}^t x_0 + \left( \mathcal{A}^t x_0 \right)^{\mathsf{T}} \overline{Q} \sum_{l=0}^{t-1} \omega \left( \mathcal{A}^{t-l-1} B \tilde{\pi}(x_l^{(0)}) \right) \right) \\
&\quad + \sum_{t=0}^{\infty} \gamma^t \left( \tilde{\pi}(x_t^{(0)})^{\mathsf{T}} R K \mathcal{A}^t x_0 + \left( \mathcal{A}^t x_0 \right)^{\mathsf{T}} K^{\mathsf{T}} R \tilde{\pi}(x_t^{(0)}) \right)
\end{aligned}
$$

Let $s = t - l$, changing the order of summation, we get

$$
\begin{aligned}
&= \sum_{l=0}^{\infty} \gamma^l \tilde{\pi}(x_l^{(0)})^{\mathsf{T}} \left[ B^{\mathsf{T}} \left( \sum_{s=1}^{\infty} \gamma^s \left( \mathcal{A}^{s-1} \right)^{\mathsf{T}} \overline{Q} \mathcal{A}^{s-1} \right) \mathcal{A} \right] \mathcal{A}^l x_0 \\
&\quad + \sum_{t=0}^{\infty} \gamma^t \tilde{\pi}(x_t^{(0)})^{\mathsf{T}} [RK] \mathcal{A}^t x_0 \\
&\quad + \sum_{l=0}^{\infty} \gamma^l \left( \mathcal{A}^l x_0 \right)^{\mathsf{T}} \left[ \mathcal{A}^{\mathsf{T}} \left( \sum_{s=1}^{\infty} \gamma^s \left( \mathcal{A}^{s-1} \right)^{\mathsf{T}} \overline{Q} \mathcal{A}^{s-1} \right) B \right] \tilde{\pi}(x_l^{(0)}) \\
&\quad + \sum_{t=0}^{\infty} \gamma^t \left( \mathcal{A}^t x_0 \right)^{\mathsf{T}} \left[ K^{\mathsf{T}} R \right] \tilde{\pi}(x_t^{(0)}) \\
&= \sum_{l=0}^{\infty} \gamma^l \tilde{\pi}(x_l^{(0)})^{\mathsf{T}} \left( \gamma B^{\mathsf{T}} P_{\gamma} \left( A + BK \right) + RK \right) \mathcal{A}^l x_0 \\
&\quad + \sum_{l=0}^{\infty} \gamma^l \left( \mathcal{A}^l x_l \right)^{\mathsf{T}} \left( \gamma \left( A + BK \right)^{\mathsf{T}} P_{\gamma} B + K^{\mathsf{T}} R \right) \tilde{\pi}(x_l^{(0)})
\end{aligned}
$$

which evaluates to zero at $K^*$ since $K^* = -\gamma \left(R + \gamma B^\mathsf{T} P_\gamma B\right)^{-1} B^\mathsf{T} P_\gamma A$. This is consistent with the fact that the first order term is the gradient multiplied by $\Delta\vartheta$. This calculation gives us an explicit form of the term that will simplify the proof of the next lemma. $\qquad\square$

**Lemma C.4.** *Let $\varrho_0$ have full support, then the Hessian is symmetric and strictly positive definite at optimality.*

*Proof of Lemma C.4.* For any $x_0 \in \mathrm{supp}\left(\varrho_0\right)$, the second order term in the approximation of the cost is as follows.

$$
C^{(2)}[x_0] = \sum_{t=2}^\infty \gamma^t \left(x_t^{(2)}\right)^\mathsf{T} \overline{Q} x_t^{(0)} + \sum_{t=2}^\infty \gamma^t \left(x_t^{(0)}\right)^\mathsf{T} \overline{Q} x_t^{(2)}
$$

$$
+ \sum_{t=1}^\infty \gamma^t \left(x_t^{(1)}\right)^\mathsf{T} \overline{Q} x_t^{(1)} + \underline{\sum_{t=0}^\infty \gamma^t \left(\tilde{\pi}(x_t^{(0)})\right)^\mathsf{T} R \left(\tilde{\pi}(x_t^{(0)})\right)}
$$

$$
+ \sum_{t=1}^\infty \gamma^t \left(\left(x_t^{(1)}\right)^\mathsf{T} K^\mathsf{T} R \tilde{\pi}(x_t^{(0)}) + \left(\tilde{\pi}(x_t^{(0)} + x_t^{(1)}) - \tilde{\pi}(x_t^{(0)})\right)^\mathsf{T} RK x_t^{(0)}\right)
$$

$$
+ \sum_{t=1}^\infty \gamma^t \left(\tilde{\pi}(x_t^{(0)})^\mathsf{T} RK x_t^{(1)} + \left(x_t^{(0)}\right)^\mathsf{T} K^\mathsf{T} R \left(\tilde{\pi}(x_t^{(0)} + x_t^{(1)}) - \tilde{\pi}(x_t^{(0)})\right)\right)
$$

Notice we can ignore the underlined sum since it is already symmetric. In the remainder of this calculation, we will use a red underline to highlight the symmetric terms that we will omit for clarity. The remaining terms of interest are

$$
= \sum_{t=1}^\infty \gamma^t \left(x_t^{(1)}\right)^\mathsf{T} \overline{Q} x_t^{(1)} + \sum_{t=2}^\infty \gamma^t \sum_{l=1}^{t-1} \left(\tilde{\pi}(x_l^{(0)} + x_l^{(1)}) - \tilde{\pi}(x_l^{(0)})\right)^\mathsf{T} B^\mathsf{T} \left(\mathcal{A}^{t-l-1}\right)^\mathsf{T} \overline{Q} \mathcal{A}^t x_0
$$

$$
+ \sum_{t=2}^\infty \gamma^t \sum_{l=1}^{t-1} \left(\mathcal{A}^t x_0\right)^\mathsf{T} \overline{Q} \mathcal{A}^{t-l-1} B \left(\tilde{\pi}(x_l^{(0)} + x_l^{(1)}) - \tilde{\pi}(x_l^{(0)})\right)
$$

$$
+ \sum_{t=1}^\infty \gamma^t \left(\left(x_t^{(1)}\right)^\mathsf{T} K^\mathsf{T} R \tilde{\pi}(x_t^{(0)}) + \left(\tilde{\pi}(x_t^{(0)} + x_t^{(1)}) - \tilde{\pi}(x_t^{(0)})\right)^\mathsf{T} RK \mathcal{A}^t x_0\right)
$$

$$
+ \sum_{t=1}^\infty \gamma^t \left(\tilde{\pi}(x_t^{(0)})^\mathsf{T} RK x_t^{(1)} + \left(\mathcal{A}^t x_0\right)^\mathsf{T} K^\mathsf{T} R \left(\tilde{\pi}(x_t^{(0)} + x_t^{(1)}) - \tilde{\pi}(x_t^{(0)})\right)\right)
$$

Like the proof of Lemma C.3, we define $s = t - l$ and change the order of summation to get

$$
= \sum_{l=1}^\infty \gamma^l \left(\tilde{\pi}(x_l^{(0)} + x_l^{(1)}) - \tilde{\pi}(x_l^{(0)})\right)^\mathsf{T} \left[B^\mathsf{T} \sum_{s=1}^\infty \gamma^s \left(\mathcal{A}^{s-1}\right)^\mathsf{T} \overline{Q} \mathcal{A}^{s-1} (A + BK)\right] x_l^{(0)}
$$

$$
+ \sum_{l=1}^\infty \gamma^l \left(x_l^{(0)}\right)^\mathsf{T} \left[(A + BK)^\mathsf{T} \sum_{s=1}^\infty \gamma^s \left(\mathcal{A}^{s-1}\right)^\mathsf{T} \overline{Q} \mathcal{A}^{s-1} B\right] \left(\tilde{\pi}(x_l^{(0)} + x_l^{(1)}) - \tilde{\pi}(x_l^{(0)})\right)
$$

$$
+ \sum_{t=1}^\infty \gamma^t \left(\tilde{\pi}(x_t^{(0)} + x_t^{(1)}) - \tilde{\pi}(x_t^{(0)})\right)^\mathsf{T} \left[RK\right] x_t^{(0)} + \sum_{t=1}^\infty \gamma^t \left(x_t^{(1)}\right)^\mathsf{T} K^\mathsf{T} R \tilde{\pi}(x_t^{(0)})
$$

$$
+ \sum_{t=1}^\infty \gamma^t \left(x_t^{(0)}\right)^\mathsf{T} \left[RK\right] \left(\tilde{\pi}(x_t^{(0)} + x_t^{(1)}) - \tilde{\pi}(x_t^{(0)})\right) + \sum_{t=1}^\infty \gamma^t \tilde{\pi}(x_t^{(0)})^\mathsf{T} RK x_t^{(1)}
$$

$$
+ \sum_{t=2}^\infty \gamma^t \left(x_t^{(1)}\right)^\mathsf{T} \overline{Q} \sum_{m=1}^{t-1} \left(\mathcal{A}^{t-m-1} B \tilde{\pi}(x_m^{(0)})\right) + \sum_{t=2}^\infty \gamma^t \sum_{m=1}^{t-1} \left(\mathcal{A}^{t-m-1} B \tilde{\pi}(x_m^{(0)})\right)^\mathsf{T} \overline{Q} x_t^{(1)}
$$

$$
+ \underline{\sum_{t=1}^\infty \gamma^t \left(\tilde{\pi}(x_0)\right)^\mathsf{T} B^\mathsf{T} \left(\mathcal{A}^{t-1}\right)^\mathsf{T} \overline{Q} \mathcal{A}^{t-1} B \tilde{\pi}(x_0)}
$$

We have dealt with the term containing two copies of $x_t^{(1)}$ by splitting it into two pieces: the sum of all terms where both indices of the inner sums are 0 and everything else. The former is symmetric, so we can omit it. Noticing that $P_\gamma = \sum_{t=0}^{\infty} \gamma^t (\mathcal{A}^t)^\top \bar{Q} \mathcal{A}^t$ we continue from above

$$
\begin{aligned}
&= \sum_{l=1}^{\infty} \gamma^l \left( \tilde{\pi}(x_l^{(0)} + x_l^{(1)}) - \tilde{\pi}(x_l^{(0)}) \right)^\mathsf{T} \left[ \gamma(R + \gamma B^\mathsf{T} P_\gamma B)K + \gamma B^\mathsf{T} P_\gamma A \right] x_l^{(0)} \\
&\quad + \sum_{l=1}^{\infty} \gamma^l \left( x_l^{(0)} \right)^\mathsf{T} \left[ \gamma K^\mathsf{T}(R + \gamma B^\mathsf{T} P_\gamma B) + \gamma A^\mathsf{T} P_\gamma B \right] \left( \tilde{\pi}(x_l^{(0)} + x_l^{(1)}) - \tilde{\pi}(x_l^{(0)}) \right) \\
&\quad + \sum_{t=2}^{\infty} \gamma^t \sum_{l=0}^{t-1} \tilde{\pi}(x_l^{(0)})^\mathsf{T} B^\mathsf{T} \left( \mathcal{A}^{t-l-1} \right)^\mathsf{T} \overline{Q} \sum_{m=1}^{t-1} \left( \mathcal{A}^{t-m-1} B \tilde{\pi}(x_m^{(0)}) \right) \\
&\quad + \sum_{t=2}^{\infty} \gamma^t \sum_{m=1}^{t-1} \left( \mathcal{A}^{t-m-1} B \tilde{\pi}(x_m^{(0)}) \right)^\mathsf{T} \overline{Q} \sum_{l=0}^{t-1} \mathcal{A}^{t-l-1} B \tilde{\pi}(x_l^{(0)}) \\
&\quad + \sum_{t=1}^{\infty} \gamma^t \left( x_t^{(1)} \right)^\mathsf{T} K^\mathsf{T} R \tilde{\pi}(x_t^{(0)}) + \sum_{t=1}^{\infty} \gamma^t \tilde{\pi}(x_t^{(0)})^\mathsf{T} R K x_t^{(1)} \\[2mm]
&= \sum_{l=1}^{\infty} \gamma^l \left( \tilde{\pi}(x_l^{(0)} + x_l^{(1)}) - \tilde{\pi}(x_l^{(0)}) \right)^\mathsf{T} \left[ \gamma(R + \gamma B^\mathsf{T} P_\gamma B)K + \gamma B^\mathsf{T} P_\gamma A \right] x_l^{(0)} \\
&\quad + \sum_{l=1}^{\infty} \gamma^l \left( x_l^{(0)} \right)^\mathsf{T} \left[ \gamma K^\mathsf{T}(R + \gamma B^\mathsf{T} P_\gamma B) + \gamma A^\mathsf{T} P_\gamma B \right] \left( \tilde{\pi}(x_l^{(0)} + x_l^{(1)}) - \tilde{\pi}(x_l^{(0)}) \right) \\
&\quad + \sum_{t=2}^{\infty} \gamma^t \sum_{l=0}^{t-1} \tilde{\pi}(x_l^{(0)})^\mathsf{T} B^\mathsf{T} \left( \mathcal{A}^{t-l-1} \right)^\mathsf{T} \overline{Q} \sum_{m=1}^{t-1} \left( \mathcal{A}^{t-m-1} B \tilde{\pi}(x_m^{(0)}) \right) \\
&\quad + \sum_{t=2}^{\infty} \gamma^t \sum_{m=1}^{t-1} \left( \mathcal{A}^{t-m-1} B \tilde{\pi}(x_m^{(0)}) \right)^\mathsf{T} \overline{Q} \sum_{l=0}^{t-1} \mathcal{A}^{t-l-1} B \tilde{\pi}(x_l^{(0)}) \\
&\quad + \sum_{t=1}^{\infty} \gamma^t \sum_{l=0}^{t-1} \left( \tilde{\pi}(x_l^{(0)}) \right)^\mathsf{T} B^\mathsf{T} \left( \mathcal{A}^{t-l-1} \right)^\mathsf{T} K^\mathsf{T} R \tilde{\pi}(x_t^{(0)}) \\
&\quad + \sum_{t=1}^{\infty} \gamma^t \sum_{l=0}^{t-1} \tilde{\pi}(x_t^{(0)})^\mathsf{T} R K \mathcal{A}^{t-l-1} B \tilde{\pi}(x_l^{(0)}) \\[2mm]
&= \sum_{l=1}^{\infty} \gamma^l \left( \tilde{\pi}(x_l^{(0)} + x_l^{(1)}) - \tilde{\pi}(x_l^{(0)}) \right)^\mathsf{T} \left[ \gamma(R + \gamma B^\mathsf{T} P_\gamma B)K + \gamma B^\mathsf{T} P_\gamma A \right] x_l^{(0)} \\
&\quad + \sum_{l=1}^{\infty} \gamma^l \left( x_l^{(0)} \right)^\mathsf{T} \left[ \gamma K^\mathsf{T}(R + \gamma B^\mathsf{T} P_\gamma B) + \gamma A^\mathsf{T} P_\gamma B \right] \left( \tilde{\pi}(x_l^{(0)} + x_l^{(1)}) - \tilde{\pi}(x_l^{(0)}) \right) \\
&\quad + \sum_{l=0}^{\infty} \gamma^l \left( \tilde{\pi}(x_l^{(0)}) \right)^\mathsf{T} B^\mathsf{T} \left[ \sum_{s=1}^{\infty} \gamma^s \left( \mathcal{A}^{s-1} \right)^\mathsf{T} K^\mathsf{T} R \tilde{\pi}(x_{l+s}^{(0)}) \right. \\
&\qquad\qquad\qquad\qquad \left. + \sum_{s=2}^{\infty} \gamma^s \left( \mathcal{A}^{s-1} \right)^\mathsf{T} \overline{Q} \sum_{r=1}^{s-1} \left( \mathcal{A}^{s-r-1} B \tilde{\pi}(x_{l+r}^{(0)}) \right) \right] \\
&\quad + \sum_{l=0}^{\infty} \gamma^l \left[ \sum_{s=1}^{\infty} \gamma^s \tilde{\pi}(x_{l+s}^{(0)})^\mathsf{T} R K \mathcal{A}^{s-1} \right.
\end{aligned}
$$

$$+ \sum_{s=2}^{\infty} \gamma^s \sum_{r=1}^{s-1} \left( \mathcal{A}^{s-r-1} B \tilde{\pi}(x_{l+r}^{(0)}) \right)^{\mathsf{T}} \overline{Q} \mathcal{A}^{s-1} \Bigg] B \tilde{\pi}(x_l^{(0)})$$

To finish the proof, notice that the bracketed terms are almost identical the first-order term from Lemma C.3. We also note that the terms multiplied to the bracketed terms are guaranteed to be finite due to the assumed Lipschitz continuity of $F$.

$$\begin{bmatrix} \cdots \end{bmatrix} = 0 \qquad \text{since } K^* = -\gamma \left( R + \gamma B^{\mathsf{T}} P_\gamma B \right)^{-1} B^{\mathsf{T}} P_\gamma A$$

$$\begin{bmatrix} \cdots \end{bmatrix} = \sum_{s=2}^{\infty} \gamma^s \sum_{r=0}^{s-2} \left( \mathcal{A}^{r+1} \mathcal{A}^{s-r-2} \right)^{\mathsf{T}} \overline{Q} \mathcal{A}^{s-r-2} B \left( \tilde{\pi}(x_{l+r+1}^{(0)}) \right)$$

$$+ \sum_{s=1}^{\infty} \gamma^s \left( \mathcal{A}^{s-1} \right)^{\mathsf{T}} K^{\mathsf{T}} R \left( \tilde{\pi}(x_{l+s}^{(0)}) \right)$$

$$= \sum_{r=0}^{\infty} \gamma^{r+1} \left( \mathcal{A}^r \right)^{\mathsf{T}} \left[ \mathcal{A}^{\mathsf{T}} \sum_{m=2}^{\infty} \gamma^{m-1} \left( \mathcal{A}^{m-2} \right)^{\mathsf{T}} \overline{Q} \mathcal{A}^{m-2} B \right] \left( \tilde{\pi}(x_{l+r+1}^{(0)}) \right)$$

$$+ \sum_{s=0}^{\infty} \gamma^{s+1} \left( \mathcal{A}^s \right)^{\mathsf{T}} K^{\mathsf{T}} R \left( \tilde{\pi}(x_{l+s+1}^{(0)}) \right)$$

$$= \sum_{r=0}^{\infty} \gamma^r \left( \mathcal{A}^r \right)^{\mathsf{T}} \left[ \gamma \left( A + BK \right)^{\mathsf{T}} P_\gamma B + K^{\mathsf{T}} R \right] \left( \tilde{\pi}(x_{l+r+1}^{(0)}) \right)$$

$$= 0$$

Thus, the second-order term of the expansion of the cost function around optimality for the trajectory with initial condition $x_0$ reads

$$C^{(2)}[x_0] = \sum_{t=0}^{\infty} \gamma^t \left( \tilde{\pi}(x_t^{(0)})^{\mathsf{T}} R \tilde{\pi}(x_t^{(0)}) + \tilde{\pi}(x_0)^{\mathsf{T}} B^{\mathsf{T}} \left( \mathcal{A}^t \right)^{\mathsf{T}} \overline{Q} \mathcal{A}^t B \tilde{\pi}(x_0) \right)$$

$$= \tilde{\pi}(x_0)^{\mathsf{T}} B^{\mathsf{T}} P_\gamma B \tilde{\pi}(x_0) + \sum_{t=0}^{\infty} \gamma^t \tilde{\pi}(x_t^{(0)})^{\mathsf{T}} R \tilde{\pi}(x_t^{(0)})$$

Since the policy is linear in the parameters, we can formally write this as an inner product where we define multiplication between an $m \times m$ matrix and a vector of $m \times 1$ vectors to be carried out component-wise.

$$= \Delta \vartheta^{\mathsf{T}} \left( \begin{bmatrix} f_1(x_0) \\ \vdots \\ f_d(x_0) \end{bmatrix}^{\mathsf{T}} B^{\mathsf{T}} P B \begin{bmatrix} f_1(x_0) \\ \vdots \\ f_d(x_0) \end{bmatrix} + \sum_{t=0}^{\infty} \gamma^t \begin{bmatrix} f_1(x_t) \\ \vdots \\ f_d(x_t) \end{bmatrix}^{\mathsf{T}} R \begin{bmatrix} f_1(x_t) \\ \vdots \\ f_d(x_t) \end{bmatrix} \right) \Delta \vartheta$$

$$= \Delta \vartheta^{\mathsf{T}} \left( \begin{bmatrix} f_1(x_0)^{\mathsf{T}} B^{\mathsf{T}} P_\gamma B f_1(x_0) & \cdots & f_1(x_0)^{\mathsf{T}} B^{\mathsf{T}} P_\gamma B f_d(x_0) \\ \vdots & \ddots & \vdots \\ f_d(x_0)^{\mathsf{T}} B^{\mathsf{T}} P_\gamma B f_1(x_0) & \cdots & f_d(x_0)^{\mathsf{T}} B^{\mathsf{T}} P_\gamma B f_d(x_0) \end{bmatrix} \right.$$

$$\left. + \sum_{t=0}^{\infty} \gamma^t \begin{bmatrix} f_1(x_t)^{\mathsf{T}} R f_1(x_t) & \cdots & f_1(x_t)^{\mathsf{T}} R f_d(x_t) \\ \vdots & \ddots & \vdots \\ f_d(x_t)^{\mathsf{T}} R f_1(x_t) & \cdots & f_d(x_t)^{\mathsf{T}} R f_d(x_t) \end{bmatrix} \right) \Delta \vartheta$$

To avoid writing this large expression again, we denote the matrix in the parentheses by $H_\gamma$. We use this matrix to define the quadratic form used as the Lyapunov function in the proof of Theorem 3.6. Given this expression for $C^{(2)}$, we notice that it is always non-negative and is zero if and only if the policy is optimal for the given initial condition, $x$ (*i.e.*, $\Delta u_t = 0 \, \forall t$). Consequently, integrating this term against the initial distribution $\varrho_0$ will result in a nonzero outcome if and only if the trajectory of almost all initial

condition is optimal. Since $\varrho_0$ has full support, and by Assumption 1

$$C^{(2)} = \int_{\mathbb{R}} C^{(2)}[x] \, d\varrho_0(x) = 0 \iff C^{(2)}[x] = 0 \; [\varrho_0]\text{-a.e.} \iff \tilde{\pi}(x_t^{(0)}) = \vec{0} \; \forall t \iff \Delta\vartheta = \vec{0}$$

$\square$

**Theorem 3.5** (Convergence at $\gamma = 0$). *Let $\{f_k\}_{k=1}^d$ be a set of globally $Lip(f_k)$-Lipschitz continuous functions, and let $\varrho_0$ (the distribution of the initial state) have full support on $\mathbb{R}^d$. Under Assumptions 1, 2, any random initialization of the parameters $\vartheta$ will converge to the optimal policy for discounted LQR with $\gamma = 0$.*

*Proof of Theorem 3.5.* Let $\vartheta$ be the randomly initialized parameters. When $\gamma = 0$, the cost can be written as

$$C_0(\vartheta) = \mathbb{E}_{x_0 \sim \varrho_0} \left[ x_0^\mathsf{T} Q x_0 + u_0^\mathsf{T} R u_0 \right]$$

For convenience, we will omit the subscripts. Notice that the optimal parameters are $\vartheta^* = \vec{0}$ which correspond to the optimal action of $u = 0$. We can calculate

$$
\begin{aligned}
\frac{d}{ds}\left(C_0(\vartheta(s)) - C_0(\vartheta_0^*)\right) &= \mathbb{E}\left[\frac{d}{ds}\left(u^\mathsf{T} R u\right)\right] = \mathbb{E}\left[2u^\mathsf{T} R \frac{d}{ds} u\right] \\
&= \mathbb{E}\left[2u^\mathsf{T} R \left(\frac{d}{ds}\pi_{\vartheta(s)}(x)\right)\right] \\
&= \mathbb{E}\left[2u^\mathsf{T} R \langle \nabla_\vartheta \pi_\vartheta(x), \frac{d}{ds}\vartheta(s)\rangle\right] \\
&= \mathbb{E}\left[2u^\mathsf{T} R \langle \nabla_\vartheta \pi_\vartheta(x), -\nabla_\vartheta C_0(\vartheta)\rangle\right] \\
&= \mathbb{E}\left[2u^\mathsf{T} R \langle \nabla_\vartheta \pi_\vartheta(x), -\frac{d}{du}C_0(\vartheta) \cdot \nabla_\vartheta \pi_\vartheta(x)\rangle\right] \\
&= \mathbb{E}\left[-4\langle u^\mathsf{T} R \nabla_\vartheta \pi_\vartheta(x), u^\mathsf{T} R \nabla_\vartheta \pi_\vartheta(x)\rangle\right] \\
&= -\operatorname{Tr}\left(\nabla_\vartheta C_0^\mathsf{T} \nabla_\vartheta C_0\right)
\end{aligned}
$$

where $\langle \cdot, \cdot \rangle$ denotes the scalar product in $\mathbb{R}^d$. To prove convergence, what remains is to show that the only extremum is the global minimum. In other words, we want to show

$$\nabla_\vartheta C_0(\vartheta) = 0 \iff \vartheta = \vartheta_0^*$$

($\Leftarrow$) This direction is straightforward. At $\vartheta = \vartheta_\gamma^*$, we see that $u = 0$ for any $x$, thus, $\nabla_C = 0$.

($\Rightarrow$) Recall that for $f \geq 0$, $\int f d\mu = 0 \iff f = 0 \; [\mu]$-a.e. and if $f$ is continuous, then $f \equiv 0$. This tells us

$$-4\langle \pi_\vartheta(x)^\mathsf{T} R \nabla_\vartheta \pi_\vartheta(x), \pi_\vartheta(x)^\mathsf{T} R \nabla_\vartheta \pi_\vartheta(x)\rangle = -\|2\pi_\vartheta(x)^\mathsf{T} R \nabla_\vartheta \pi_\vartheta(x)\|^2 = 0 \quad \forall x$$

This condition is equivalent to

$$\left(\sum_{k=0}^d \vartheta_k f_k(x) R f_i(x)\right)^2 = 0 \quad \forall x \in \mathbb{R}^n, i \in \{1, \dots, d\}$$

and

$$\sum_{k=0}^d \vartheta_k f_k(x) R f_i(x) = 0 \qquad \forall x \quad \forall x \in \mathbb{R}^n, i \in \{1, \dots, d\}$$

Since we assumed that $\{f_k\}$ is linearly independent (Assumption 1), by positive definiteness of $R$ we get that $\vartheta = \vartheta_\gamma^*$.

$\square$

For the next theorem, we recall the definition of the positive definite Hessian $H_\gamma$ appearing in the proof of Lemma C.4:

$$\Delta\vartheta^\mathsf{T} H_\gamma \Delta\vartheta^\mathsf{T} = \Delta\vartheta^\mathsf{T} \left( \begin{bmatrix} f_1(x_0) \\ \vdots \\ f_d(x_0) \end{bmatrix}^\mathsf{T} B^\mathsf{T} P B \begin{bmatrix} f_1(x_0) \\ \vdots \\ f_d(x_0) \end{bmatrix} + \sum_{t=0}^\infty \gamma^t \begin{bmatrix} f_1(x_t) \\ \vdots \\ f_d(x_t) \end{bmatrix}^\mathsf{T} R \begin{bmatrix} f_1(x_t) \\ \vdots \\ f_d(x_t) \end{bmatrix} \right) \Delta\vartheta$$

**Theorem 3.6** (Convergence Near $\gamma$-Optimality)**.** *Under the same conditions as Theorem 3.5, for any $\gamma \in (0,1)$ there exists $\delta > 0$, $\lambda > 0$ such that for all initial conditions $\vartheta(0)$ with $\Delta\vartheta(0) := \|\vartheta(0) - \vartheta_\gamma^*\| < \delta$,*

$$C_\gamma(\vartheta(s)) - C_\gamma(\vartheta_\gamma^*) \le e^{-s\lambda} \left( C_\gamma(\vartheta(0)) - C_\gamma(\vartheta_\gamma^*) \right),$$

*and that $\lim_{s\to\infty} \vartheta(s) = \vartheta_\gamma^*$.*

*Proof of Theorem 3.6.* Consider the Lyapunov function

$$U(s) = \frac{1}{2} \left( \vartheta(s) - \vartheta_\gamma^* \right)^\mathsf{T} H_\gamma \left( \vartheta(s) - \vartheta_\gamma^* \right)$$

where $H_\gamma$ is a symmetric positive-definite matrix defined at the end of Lemma C.4.

In this regime we can expand the cost function $C_\gamma(\,\cdot\,)$ around its minimum $\vartheta_\gamma^*$. By Lemmas C.3 and C.4, the second-order Taylor expansion of the cost is

$$C_\gamma(\vartheta(s)) = C_\gamma^{(0)}|_{\vartheta_\gamma^*} + C_\gamma^{(1)}|_{\vartheta_\gamma^*} + C_\gamma^{(2)}|_{\vartheta_\gamma^*} + o(\|\Delta\vartheta(s)\|^2)$$

$$= C_\gamma(\vartheta_\gamma^*) + \frac{1}{2}\Delta\vartheta(s)^\mathsf{T} H_\gamma \Delta\vartheta(s) + o(\|\Delta\vartheta(s)\|^2)$$

$$C_\gamma(\vartheta(s)) - C_\gamma(\vartheta_\gamma^*) = U(s) + o(\|\Delta\vartheta(s)\|^2)$$

By chain rule,

$$\frac{d}{ds}U(s) = \langle \nabla_{\Delta\vartheta} U(s), \frac{d}{ds}\Delta\vartheta(s) \rangle$$

$$= \langle H_\gamma \Delta\vartheta(s), -\nabla_\vartheta C(\vartheta(s)) \rangle$$

$$= -\Delta\vartheta(s)^\mathsf{T} H_\gamma^2 \Delta\vartheta(s) + o(\|\Delta\vartheta(s)\|^2)$$

$$\le -\lambda_{\min}(H_\gamma) U(s) + o(\|\Delta\vartheta(s)\|^2)$$

where $\lambda_{\min}(H_\gamma)$ is the smallest of $H_\gamma$. By Grönwall's Inequality, this immediately gives, for $\|\Delta\vartheta(0)\| < \delta$ for $\delta > 0$ sufficiently small, that

$$U(s) \le e^{-s\lambda_{\min}(H_\gamma)/2} U(0),$$

recovering the claim for $\lambda = \lambda_{\min}(H_\gamma)/2$. $\qquad\square$

DEPARTMENT OF MATHEMATICS, DUKE UNIVERSITY, DURHAM, NC 27708
*Email address*: craig.chen@duke.edu and agazzi@math.duke.edu