

Can Direct Latent Model Learning Solve Linear Quadratic Gaussian Control?

Yi Tian

Massachusetts Institute of Technology

yitian@mit.edu

Kaiqing Zhang

University of Maryland, College Park

kaiqing@umd.edu

Russ Tedrake

Massachusetts Institute of Technology

russt@mit.edu

Suvrit Sra

Massachusetts Institute of Technology

suvrit@mit.edu

Abstract

We study the task of learning state representations from potentially high-dimensional observations, with the goal of controlling an unknown partially observable system. We pursue a *direct latent model learning* approach, where a dynamic model in some latent state space is learned by predicting quantities directly related to planning (e.g., costs) without reconstructing the observations. In particular, we focus on an intuitive cost-driven state representation learning method for solving Linear Quadratic Gaussian (LQG) control, one of the most fundamental partially observable control problems. As our main results, we establish finite-sample guarantees of finding a near-optimal state representation function and a near-optimal controller using the directly learned latent model. To the best of our knowledge, despite various empirical successes, prior to this work it was unclear if such a cost-driven latent model learner enjoys finite-sample guarantees. Our work underscores the value of predicting multi-step costs, an idea that is key to our theory, and notably also an idea that is known to be empirically valuable for learning state representations.

1 Introduction

We consider state representation learning for control in partially observable systems, inspired by the recent successes of *control from pixels* (Hafner et al., 2019b,a). Control from pixels is an everyday task for human beings, but it remains challenging for learning agents. Methods to achieve it generally fall into two main categories: *model-free* and *model-based* ones. Model-free methods directly learn a visuomotor policy, also known as direct reinforcement learning (RL) (Sutton and Barto, 2018). On the other hand, model-based methods, also known as indirect RL (Sutton and Barto, 2018), attempt to learn a *latent model* that is a compact representation of the system, and to synthesize a policy in the latent model. Compared with model-free methods,

model-based ones facilitate generalization across tasks and enable efficient planning (Hafner et al., 2020), and are sometimes more sample efficient (Tu and Recht, 2019; Sun et al., 2019; Zhang et al., 2019) than the model-free ones.

In latent-model-based control, the state of the latent model is also referred to as a *state representation* in the deep RL literature, and the mapping from an observed history to a latent state is referred to as the (state) representation function. *Reconstructing the observation* often serves as a supervision for representation learning for control in the empirical RL literature (Hafner et al., 2019b,a, 2020; Fu et al., 2021; Wang et al., 2022). This is in sharp contrast to model-free methods, where the policy improvement step is completely *cost-driven*. Reconstructing observations provides a rich supervision signal for learning a task-agnostic world model, but they are high-dimensional and noisy, so the reconstruction requires an expressive reconstruction function; latent states learned by reconstruction contain irrelevant information for control, which can distract RL algorithms (Zhang et al., 2020; Fu et al., 2021; Wang et al., 2022). This is especially the case for practical visuomotor control tasks, e.g., robotic manipulation and self-driving cars, where the visual images contain predominately task-irrelevant objects and backgrounds.

Various empirical attempts (Schrittwieser et al., 2020; Zhang et al., 2020; Okada and Taniguchi, 2021; Deng et al., 2021; Yang et al., 2022) have been made to bypass observation reconstruction. Apart from observation, the interaction involves two other variables: actions (control inputs) and costs. Inverse model methods (Lamb et al., 2022) reconstruct actions; while other methods rely on costs. We argue that since neither the reconstruction function nor the inverse model is used for policy learning, cost-driven state representation learning is the most *direct* one, in that costs are directly relevant for control purposes. In this paper, we aim to examine the soundness of this methodology in linear quadratic Gaussian (LQG) control, one of the most fundamental partially observable control models.

Parallel to the empirical advances of learning for control from pixels, partially observable linear systems has been extensively studied in the context of learning for dynamic control (Oymak and Ozay, 2019; Simchowitz et al., 2020; Lale et al., 2020, 2021; Zheng et al., 2021; Minasyan et al., 2021; Umenberger et al., 2022). In this context, the representation function is more formally referred to as a *filter*, the optimal one being the Kalman filter. Most existing *model-based* learning approaches for LQG control focus on the linear time-invariant (LTI) case, and are based on the idea of *learning Markov parameters* (Ljung, 1998), the mapping from control inputs to observations. Hence, they need to predict observations by definition. Motivated by the empirical successes in control from pixels, we take a different, cost-driven route, in hope of avoiding reconstructing observations or control inputs.

We focus on finite-horizon time-varying LQG control and address the following question:

Can direct, cost-driven state representation learning provably solve LQG control?

This work answers the question in the affirmative, by establishing finite-sample guarantees for a cost-driven state representation learning method.

Challenges & Our techniques. Overall, to establish finite-sample guarantees, a major technical challenge is to deal with the *quadratic regression* problem in cost prediction, arising from

the inherent quadratic form of the cost function in LQG. Directly solving the problem for the representation function involves *quartic* optimization; instead, we propose to solve a quadratic regression problem, followed by low-rank approximate factorization. The quadratic regression problem also appears in identifying the cost matrices, which involves concentration for random variables that are fourth powers of Gaussians. We believe these techniques might be of independent interest.

Moreover, the first ℓ -step *latent* states may not be adequately *excited* (having full-rank covariance), which invalidates the use of most system identification techniques. We instead identify only *relevant directions* of the system parameters, and prove that this is sufficient for learning a near-optimal controller by analyzing state covariance mismatch. This fact is reflected in the separation in the statement of Theorem 1; developing finite-sample analysis in this case is technically challenging.

Implications. For practitioners, one takeaway from our work is the benefit of predicting *multi-step cumulative* costs in cost-driven state representation learning. Whereas the cost at a single time step may not be revealing enough of the latent state, the cumulative cost across multiple steps can be. This is an intuitive idea for the control community, given the multi-step nature in the classical definitions of controllability and observability. Its effectiveness has also been empirically observed in MuZero (Schrittwieser et al., 2020) in state representation learning for control, and our work can be viewed as a formal understanding of it in the LQG control setting.

Notation. We use 0 (resp. 1) to denote either the scalar or a matrix consisting of all zeros (resp. all ones); we use I to denote an identity matrix. The dimension, when emphasized, is specified in subscripts, e.g., $0_{d_x \times d_x}$, 1_{d_x} , I_{d_x} . Let \mathbb{I}_S denote the indicator function for set S and $\mathbb{I}_S(A)$ apply to matrix A elementwise. For some positive semidefinite P , we define $\|v\|_P := (v^\top P v)^{1/2}$. Semicolon “;” denotes stacking vectors or matrices vertically. For a collection of d -dimensional vectors $(v_t)_{t=i}^j$, let $v_{i:j} := [v_i; v_{i+1}; \dots; v_j] \in \mathbb{R}^{d(j-i+1)}$ denote the concatenation along the column. For random variable η , let $\|\eta\|_{\psi_\beta}$ denote its β -sub-Weibull norm, a special case of Orlicz norms (Zhang and Wei, 2022), with $\beta = 1, 2$ corresponding to subexponential and sub-Gaussian norms. $\sigma_i(A)$, $\sigma_{\min}(A)$, $\sigma_{\min}^+(A)$, $\sigma_{\max}(A)$ denote its i th largest, minimum, minimum positive, maximum singular values, respectively. $\|A\|_2$, $\|A\|_F$, $\|A\|_*$ denote the operator (induced by vector 2-norms), Frobenius, nuclear norms of matrix A , respectively. $\langle \cdot, \cdot \rangle_F$ denotes the Frobenius inner product between matrices. The Kronecker, symmetric Kronecker and Hadamard products between matrices are denoted by “ \otimes ”, “ \otimes_s ” and “ \odot ”, respectively. $\text{vec}(\cdot)$ and $\text{svec}(\cdot)$ denote flattening a matrix and a symmetric matrix by stacking their columns; $\text{svec}(\cdot)$ does not repeat the off-diagonal elements, but scales them by $\sqrt{2}$ (Schacke, 2004). Let $\text{diag}(\cdot)$ denote the block diagonal matrix formed by the matrices inside the parentheses.

2 Problem setup

We study a partially observable linear time-varying (LTV) dynamical system

$$x_{t+1} = A_t^* x_t + B_t^* u_t + w_t, \quad y_t = C_t^* x_t + v_t, \quad (2.1)$$

for $t = 0, 1, \dots, T - 1$ and $y_T = C_T^* x_T + v_T$. For all $t \geq 0$, we have the notation of state $x_t \in \mathbb{R}^{d_x}$, observation $y_t \in \mathbb{R}^{d_y}$, and control input $u_t \in \mathbb{R}^{d_u}$. Process noises $(w_t)_{t=0}^{T-1}$ and observation noises $(v_t)_{t=0}^T$ are i.i.d. sampled from $\mathcal{N}(0, \Sigma_{w_t})$ and $\mathcal{N}(0, \Sigma_{v_t})$, respectively. Let initial state x_0 be sampled from $\mathcal{N}(0, \Sigma_0)$.

Let $\Phi_{t,t_0} = A_{t-1}^* A_{t-2}^* \cdots A_{t_0}^*$ for $t > t_0$ and $\Phi_{t,t} = I$. Then $x_t = \Phi_{t,t_0} x_{t_0} + \sum_{\tau=t_0}^{t-1} \Phi_{t,\tau+1} w_\tau$ under zero control input. To ensure the state and the cumulative noise do not grow with time, we make the following uniform exponential stability assumption.

Assumption 1 (Uniform exponential stability). *The system is uniformly exponentially stable, i.e., there exists $\alpha > 0, \rho \in (0, 1)$ such that for any $0 \leq t_0 < t \leq T$, $\|\Phi_{t,t_0}\|_2 \leq \alpha \rho^{t-t_0}$.*

Assumption 1 is standard in controlling LTV systems (Zhou and Zhao, 2017; Minasyan et al., 2021), satisfied by a stable LTI system. It essentially says that zero control is a stabilizing policy, and can be potentially relaxed to the assumption of *being given a stabilizing policy* (Lale et al., 2020; Simchowit et al., 2020). Specifically, one can excite the system using the stabilizing policy plus Gaussian random noises.

Define the ℓ -step controllability matrix

$$\Phi_{t,\ell}^c := [B_t^*, A_t^* B_{t-1}^*, \dots, A_t^* A_{t-1}^* \cdots A_{t-\ell+2}^* B_{t-\ell+1}^*] \in \mathbb{R}^{d_x \times \ell d_u}$$

for $\ell - 1 \leq t \leq T - 1$, which reduces to the standard controllability matrix $[B, \dots, A^{\ell-1} B]$ in the LTI setting. We make the following controllability assumption.

Assumption 2 (Controllability). *For all $\ell - 1 \leq t \leq T - 1$, $\text{rank}(\Phi_{t,\ell}^c) = d_x$, $\sigma_{\min}(\Phi_{t,\ell}^c) \geq \nu > 0$.*

Under zero noise, we have

$$x_{t+\ell} = \Phi_{t+\ell,t} x_t + \Phi_{t+\ell-1,\ell}^c [u_{t+\ell-1}; \dots; u_t],$$

so Assumption 2 ensures that from any state x , there exist control inputs that drive the state to 0 in ℓ steps, and ν ensures that the equation leading to them is well conditioned. We do not assume controllability for $0 \leq t < \ell - 1$, since we do not want to impose the constraint that $d_u > d_x$. This turns out to present a significant challenge for state representation learning, as seen from the separation of the results before and after the ℓ -steps in Theorem 1.

The quadratic cost functions are given by

$$c_t(x, u) = \|x\|_{Q_t^*}^2 + \|u\|_{R_t^*}^2, \quad 0 \leq t \leq T - 1, \quad (2.2)$$

and $c_T(x) = \|x\|_{Q_T^*}^2$, where $(Q_t^*)_{t=0}^T$ are positive semidefinite matrices and $(R_t^*)_{t=0}^{T-1}$ are positive definite matrices. Sometimes the cost is defined as a function on observation y . Since the quadratic form $y^\top Q_t^* y = x^\top (C_t^*)^\top Q_t^* C_t^* x$, our analysis still applies if the assumptions on $(Q_t^*)_{t=0}^T$ hold for $((C_t^*)^\top Q_t^* C_t^*)_{t=0}^T$ instead.

(A, C) and $(A, Q^{1/2})$ observabilities are standard assumptions in controlling LTI systems. To differentiate from the former, we call the latter *cost observability*, since it implies the states are observable through costs. Whereas Markov-parameter-based approaches need to assume (A, C)

observability to identify the system, our cost-driven approach does not. Here we deal with the more difficult problem of having only the scalar cost as the supervision signal (instead of the concatenation of all observations, as in Markov-parameter-based ones). Nevertheless, the notion of cost observability is still important for our approach, formally defined as follows.

Assumption 3 (Cost observability). *For all $0 \leq t \leq \ell - 1$, $Q_t^* \succcurlyeq \mu^2 I$. For all $\ell \leq t \leq T$, there exists $m > 0$ such that the cost observability Gramian (Kailath, 1980)*

$$\sum_{\tau=t}^{t+k-1} \Phi_{\tau,t}^\top Q_\tau^* \Phi_{\tau,t} \succcurlyeq \mu^2 I,$$

where $k = m \wedge (T - t + 1)$.

This assumption ensures that without noises, if we start with a nonzero state, the cumulative cost becomes positive in m steps. The special requirement for $0 \leq t \leq \ell - 1$ results from the difficulty in lacking controllability in these time steps. The following is a regularity assumption on system parameters.

Assumption 4. *$(\sigma_{\min}(\Sigma_{v_t}))_{t=0}^T$ are uniformly lower bounded by some $\sigma_v > 0$. The operator norms of all matrices in the problem definition are uniformly upper bounded, including $(A_t^*, B_t^*, R_t^*, \Sigma_{w_t})_{t=0}^{T-1}$, $(C_t^*, Q_t^*, \Sigma_{v_t})_{t=0}^T$. In other words, they are all $\mathcal{O}(1)$.*

Let $h_t := [y_{0:t}; u_{0:(t-1)}] \in \mathbb{R}^{(t+1)d_y + td_u}$ denote the available history before deciding control u_t . A policy $\pi = (\pi_t : h_t \mapsto u_t)_{t=0}^{T-1}$ determines at time t a control input u_t based on history h_t . With a slight abuse of notation, let $c_t := c_t(x_t, u_t)$ for $0 \leq t \leq T - 1$ and $c_T := c_T(x_T)$ denote the cost at each time step. Then, $J^\pi := \mathbb{E}^\pi[\sum_{t=0}^T c_t]$ is the expected cumulative cost under policy π , where the expectation is taken over the randomness in the process noises, observation noises, and controls (if the policy π is stochastic). The objective of LQG control is to find a policy π such that J^π is minimized.

If the system parameters $((A_t^*, B_t^*, R_t^*)_{t=0}^{T-1}, (C_t^*, Q_t^*)_{t=0}^T)$ are known, the optimal control is obtained by combining the Kalman filter $z_0^* = L_0^* y_0$,

$$z_{t+1}^* = A_t^* z_t^* + B_t^* u_t + L_{t+1}^* (y_{t+1} - C_{t+1}^* (A_t^* z_t^* + B_t^* u_t))$$

for $0 \leq t \leq T - 1$, with the optimal feedback control gains of the linear quadratic regulator (LQR) $(K_t^*)_{t=0}^{T-1}$ such that $u_t^* = K_t^* z_t^*$, where $(L_t^*)_{t=0}^T$ are the Kalman gains; this is known as the *separation principle* (Åström, 2012). The Kalman gains and optimal feedback control gains are given by

$$\begin{aligned} L_t^* &= S_t^* (C_t^*)^\top (C_t^* S_t^* (C_t^*)^\top + \Sigma_{v_t})^{-1}, \\ K_t^* &= -((B_t^*)^\top P_{t+1}^* B_t^* + R_t^*)^{-1} (B_t^*)^\top P_{t+1}^* A_t^*, \end{aligned}$$

where S_t^* and P_t^* are determined by their corresponding Riccati difference equations (RDEs):

$$S_{t+1}^* = A_t^* (S_t^* - S_t^* (C_t^*)^\top (C_t^* S_t^* (C_t^*)^\top + \Sigma_{v_t})^{-1} C_t^* S_t^*) (A_t^*)^\top + \Sigma_{w_t}, \quad (2.3)$$

$$P_t^* = (A_t^*)^\top (P_{t+1}^* - P_{t+1}^* B_t^* ((B_t^*)^\top P_{t+1}^* B_t^* + R_t^*)^{-1} (B_t^*)^\top P_{t+1}^*) A_t^* + Q_t^*, \quad (2.4)$$

with $S_0^* = \Sigma_0$ and $P_T^* = Q_T^*$.

We consider data-driven control in a partially observable LTV system (2.1) with unknown cost matrices $(Q_t^*)_{t=0}^T$. For simplicity, we assume $(R_t^*)_{t=0}^T$ is known, though our approaches can be readily generalized to the case without knowing them; one can identify them in the quadratic regression (3.3).

2.1 Latent model of LQG

Under the Kalman filter, the observation prediction error $i_{t+1} := y_{t+1} - C_{t+1}^*(A_t^*z_t^* + B_t^*u_t)$ is called an *innovation*. It is known that i_t is independent of history h_t and $(i_t)_{t=1}^T$ are independent (Bertsekas, 2012).

Now we are ready to present the following proposition that represents the system in terms of the state estimates by the Kalman filter, which we shall refer to as the *latent model*.

Proposition 1. *Let $(z_t^*)_{t=0}^T$ be state estimates given by the Kalman filter. Then,*

$$z_{t+1}^* = A_t^*z_t^* + B_t^*u_t + L_{t+1}^*i_{t+1},$$

where $L_{t+1}^*i_{t+1}$ is independent of z_t^* and u_t , i.e., the state estimates follow the same linear dynamics as the underlying state, with noises $L_{t+1}^*i_{t+1}$. The cost at step t can then be reformulated as functions of the state estimates by

$$c_t = \|z_t^*\|_{Q_t^*}^2 + \|u_t\|_{R_t^*}^2 + b_t + \gamma_t + \eta_t,$$

where $b_t > 0$ is a problem-dependent constant, and $\gamma_t = \|z_t^* - x_t\|_{Q_t^*}^2 - b_t$, $\eta_t = \langle z_t^*, x_t - z_t^* \rangle_{Q_t^*}$ are both zero-mean subexponential random variables. Under Assumptions 1 and 4, $b_t = \mathcal{O}(1)$ and $\|\gamma_t\|_{\psi_1} = \mathcal{O}(d_x^{1/2})$; moreover, if control $u_t \sim \mathcal{N}(0, \sigma_u^2 I)$ for $0 \leq t \leq T$, then $\|\eta_t\|_{\psi_1} = \mathcal{O}(d_x^{1/2})$.

Proof. By the property of the Kalman filter, $z_t^* = \mathbb{E}[x_t \mid y_{0:t}, u_{0:(t-1)}]$ is a function of the past history $(y_{0:t}, u_{0:(t-1)})$. u_t is a function of the past history $(y_{0:t}, u_{0:(t-1)})$ and some independent random variables. Since i_{t+1} is independent of the past history $(y_{0:t}, u_{0:(t-1)})$, it is independent of z_t^* and u_t . For the cost function,

$$c_t = \|z_t^*\|_{Q_t^*}^2 + \|u_t\|_{R_t^*}^2 + \|z_t^* - x_t\|_{Q_t^*}^2 + 2\langle z_t^*, x_t - z_t^* \rangle_{Q_t^*}.$$

Let $b_t = \mathbb{E}[\|x_t - z_t^*\|_{Q_t^*}^2]$ be a constant that depends on system parameters $(A_t^*, B_t^*, \Sigma_{w_t})_{t=0}^{T-1}$, $(C_t^*, \Sigma_{v_t})_{t=0}^T$ and Σ_0 . Then, random variable $\gamma_t := \|z_t^* - x_t\|_{Q_t^*}^2 - b_t$ has zero mean. Since $(x_t - z_t^*)$ is Gaussian, its squared norm is subexponential. Since z_t^* and $(x_t - z_t^*)$ are independent zero-mean Gaussian random vectors (Bertsekas, 2012), their inner product η_t is a zero-mean subexponential random variable.

If the system is uniformly exponentially stable (Assumption 1) and the system parameters are regular (c.f. Assumption 4), then $(S_t^*)_{t=0}^T$ given by RDE (2.3) has a bounded operator norm determined by system parameters $(A_t^*, B_t^*, C_t^*, \Sigma_{w_t})_{t=0}^{T-1}$, $(\Sigma_{v_t})_{t=0}^T$ and Σ_0 (Zhang and Zhang, 2021). Since $S_t^* = \text{Cov}(x_t - z_t^*)$, $\|\gamma_t\|_{\psi_1} = \mathcal{O}(d_x^{1/2})$ by Lemma 7. By Assumption 1, if we apply zero control to the system, then $\|\text{Cov}(z_t^*)\|_2 = \mathcal{O}(1)$. By Lemma 7, $\eta_t = \langle z_t^*, x_t - z_t^* \rangle$ satisfies $\|\eta_t\|_{\psi_1} = \mathcal{O}(d_x^{1/2})$. \square

Proposition 1 states that: 1) the dynamics of the state estimates produced by the Kalman filter remains the same as the original system up to noises, determined by $(A_t^*, B_t^*)_{t=0}^{T-1}$; 2) the costs (of the latent model) are still determined by $(Q_t^*)_{t=0}^T$ and $(R_t^*)_{t=0}^{T-1}$, up to constants and noises. Hence, a latent model can be parameterized by $((A_t, B_t)_{t=0}^{T-1}, (Q_t)_{t=0}^T)$ (recall that we assume $(R_t^*)_{t=0}^T$ is known for convenience). Note that observation matrices $(C_t^*)_{t=0}^T$ are *not* involved.

Now let us take a closer look at the state representation function. The Kalman filter can be written as $z_{t+1}^* = \bar{A}_t^* z_t^* + \bar{B}_t^* u_t + L_{t+1}^* y_{t+1}$, where $\bar{A}_t^* = (I - L_{t+1}^* C_{t+1}^*) A_t^*$ and $\bar{B}_t^* = (I - L_{t+1}^* C_{t+1}^*) B_t^*$. For $0 \leq t \leq T$, unrolling the recursion gives

$$\begin{aligned} z_t^* &= \bar{A}_{t-1}^* z_{t-1}^* + \bar{B}_{t-1}^* u_{t-1} + L_t^* y_t \\ &= [\bar{A}_{t-1}^* \bar{A}_{t-2}^* \cdots \bar{A}_0^* L_0^*, \dots, L_t^*] [y_0; \dots; y_t] + [\bar{A}_{t-1}^* \bar{A}_{t-2}^* \cdots \bar{A}_1^* \bar{B}_0^*, \dots, \bar{B}_{t-1}^*] [u_0; \dots; u_{t-1}] \\ &=: M_t^* [y_{0:t}; u_{0:(t-1)}], \end{aligned}$$

where $M_t^* \in \mathbb{R}^{d_x \times ((t+1)d_y + td_u)}$. This means the optimal state representation function is *linear* in the history of observations and controls. A state representation function can then be parameterized by matrices $(M_t)_{t=0}^T$, and the latent state at step t is given by $z_t = M_t h_t$.

Overall, a policy π is a combination of state representation function $(M_t)_{t=0}^{T-1}$ (M_T is not needed) and feedback gain $(K_t)_{t=0}^{T-1}$ in the latent model; in this case, we write $\pi = (M_t, K_t)_{t=0}^{T-1}$ as the composition of the two.

3 Methodology: cost-driven state representation learning

State representation learning involves history data that contains samples of three variables: observation, control input, and cost. Each of them can potentially be used as a *supervision* signal, and be used to define a type of state representation learning algorithms. We summarize our categorization of the methods in the literature as follows.

- *Predicting observations* defines the class of *observation-reconstruction-based* methods, including methods based on Markov parameters (mapping from control actions to observations) in linear systems (Lale et al., 2021; Zheng et al., 2021) and methods that learn a mapping from states to observations in more complex systems (Ha and Schmidhuber, 2018; Hafner et al., 2019b,a). This type of method tends to recover all state components.
- *Predicting actions* defines the class of *inverse model* methods, where the control is predicted from states across different time steps (Mhammedi et al., 2020; Frandsen et al., 2022; Lamb et al., 2022). This type of method tends to recover the control-relevant state components.
- *Predicting (cumulative) costs* defines the class of *cost-driven state representation learning* methods (Zhang et al., 2020; Schrittwieser et al., 2020; Yang et al., 2022). This type of methods tend to recover the state components relevant to the cost.

Our method falls into the cost-driven category, which is more direct than the other two types, in the sense that the cost is directly relevant to planning with a dynamic model, whereas the observation reconstruction functions and inverse models are not. Another reason why we

Algorithm 1 Direct latent model learning

- 1: **Input:** sample size n , input noise magnitude σ_u , singular value threshold $\theta = \Theta(n^{-1/4})$
- 2: Collect n trajectories using $u_t \sim \mathcal{N}(0, \sigma_u^2 I)$, for $0 \leq t \leq T - 1$, to obtain data in the form of

$$\mathcal{D}_{\text{raw}} = (y_0^{(i)}, u_0^{(i)}, c_0^{(i)}, \dots, y_{T-1}^{(i)}, u_{T-1}^{(i)}, c_{T-1}^{(i)}, y_T^{(i)}, c_T^{(i)})_{i=1}^n$$

- 3: Run CoREL($\mathcal{D}_{\text{raw}}, \theta$) (Algorithm 2) to obtain state representation function estimate $(\hat{M}_t)_{t=0}^T$ and latent state estimates $(z_t^{(i)})_{t=0, i=1}^{T, n}$, so that the data are converted to

$$\mathcal{D}_{\text{state}} = (z_0^{(i)}, u_0^{(i)}, c_0^{(i)}, \dots, z_{T-1}^{(i)}, u_{T-1}^{(i)}, c_{T-1}^{(i)}, z_T^{(i)}, c_T^{(i)})_{i=1}^n$$

- 4: Run SysID($\mathcal{D}_{\text{state}}$) (Algorithm 3) to obtain system parameter estimates $((\hat{A}_t, \hat{B}_t)_{t=0}^{T-1}, (\hat{Q}_t)_{t=0}^T)$
 - 5: Find feedback gains $(\hat{K}_t)_{t=0}^{T-1}$ from $((\hat{A}_t, \hat{B}_t, R_t^*)_{t=0}^{T-1}, (\hat{Q}_t)_{t=0}^T)$ by RDE (2.4)
 - 6: **Return:** policy $\hat{\pi} = (\hat{M}_t, \hat{K}_t)_{t=0}^{T-1}$
-

call our method *direct latent model learning* is that compared with Markov parameter-based approaches for linear systems, our approach directly parameterizes the state representation function, without exploiting the structure of the Kalman filter, making our approach closer to empirical practice that was designed for general RL settings.

Reference (Subramanian et al., 2020) proposes to optimize a simple combination of cost and transition prediction errors to learn the *approximate information state* (Subramanian et al., 2020). That is, we parameterize a state representation function by matrices $(M_t)_{t=0}^T$ and a latent model by matrices $((A_t, B_t)_{t=0}^{T-1}, (Q_t)_{t=0}^T)$ and then solve

$$\min_{(M_t, Q_t, b_t)_{t=0}^T, (A_t, B_t)_{t=0}^{T-1}} \sum_{t=0}^T \sum_{i=1}^n l_t^{(i)}, \quad (3.1)$$

where $(b_t)_{t=0}^T$ are additional scalar parameters to account for noises, and the loss at step t for trajectory i is defined by

$$l_t^{(i)} = (\|M_t h_t^{(i)}\|_{Q_t}^2 + \|u_t^{(i)}\|_{R_t^*}^2 + b_t - c_t^{(i)})^2 + \|M_{t+1} h_{t+1}^{(i)} - A_t M_t h_t^{(i)} - B_t u_t^{(i)}\|^2, \quad (3.2)$$

for $0 \leq t \leq T - 1$ and $l_T^{(i)} = (\|M_T h_T^{(i)}\|_{Q_T}^2 + b_T - c_T^{(i)})^2$. The optimization problem (3.1) is nonconvex; even if we can find a global minimizer, it is unclear how to establish finite-sample guarantees for it. A main finding of this work is that for LQG, we can solve the cost and transition loss optimization problems *sequentially*, with the caveat of using *cumulative* costs.

Our method is summarized in Algorithm 1. It has three steps: cost-driven state representation function learning (CoREL, Algorithm 2), latent system identification (SysID, Algorithm 3), and planning by RDE (2.4).

This three-step approach is very similar to the World Model approach (Ha and Schmidhuber, 2018) used in empirical RL, except that in the first step, instead of using an autoencoder to learn the state representation function, we use cost values to supervise the representation learning. Most empirical state representation learning methods (Hafner et al., 2019b,a; Schrittwieser et al.,

Algorithm 2 CoREL: cost driven state representation learning

- 1: **Input:** raw data \mathcal{D}_{raw} , singular value threshold $\theta = \Theta((\ell(d_y + d_u))^{1/2} d_x^{3/4} n^{-1/4})$
- 2: Estimate the state representation function and cost constants by solving $(\hat{N}_t, \hat{b}_t)_{t=0}^T \in$

$$\underset{(N_t=N_t^\top, b_t)_{t=0}^T}{\operatorname{argmin}} \sum_{t=0}^T \sum_{i=1}^n \left(\|[y_{0:t}^{(i)}; u_{0:(t-1)}^{(i)}]\|_{N_t}^2 + \sum_{\tau=t}^{t+k-1} \|u_\tau^{(i)}\|_{R_\tau^*}^2 + b_t - \bar{c}_t^{(i)} \right)^2, \quad (3.3)$$

where $k = 1$ for $0 \leq t \leq l - 1$ and $k = m \wedge (T - t + 1)$ for $\ell \leq t \leq T$

- 3: Find $\tilde{M}_t \in \mathbb{R}^{d_x \times ((t+1)d_y + td_u)}$ such that $\tilde{M}_t^\top \tilde{M}_t$ is an approximation of \hat{N}_t
 - 4: For all $0 \leq t \leq \ell - 1$, set $\hat{M}_t = \text{TRUNCsv}(\tilde{M}_t, \theta)$; for all $\ell \leq t \leq T$, set $\hat{M}_t = \tilde{M}_t$
 - 5: Compute $\hat{z}_t^{(i)} = \hat{M}_t [y_{0:t}^{(i)}; u_{0:(t-1)}^{(i)}]$ for all $t = 0, \dots, T$ and $i = 1, \dots, n$
 - 6: **Return:** state representation estimate $(\hat{M}_t)_{t=0}^T$ and latent state estimates $(\hat{z}_t^{(i)})_{t=0, i=1}^{T, n}$
-

Algorithm 3 SysID: system identification

- 1: **Input:** data in the form of $(z_0^{(i)}, u_0^{(i)}, c_0^{(i)}, \dots, z_{T-1}^{(i)}, u_{T-1}^{(i)}, c_{T-1}^{(i)}, z_T^{(i)}, c_T^{(i)})_{i=1}^n$
- 2: Estimate the system dynamics by $(\hat{A}_t, \hat{B}_t)_{t=0}^{T-1} \in$ \triangleright pick min. Frobenius norm solution by pseudoinverse

$$\underset{(A_t, B_t)_{t=0}^{T-1}}{\operatorname{argmin}} \sum_{t=0}^{T-1} \sum_{i=1}^n \|A_t z_t^{(i)} + B_t u_t^{(i)} - z_{t+1}^{(i)}\|^2 \quad (3.4)$$

- 3: For all $0 \leq t \leq \ell - 1$ and $t = T$, set $\hat{Q}_t = I_{d_x}$
- 4: For all $\ell \leq t \leq T - 1$, obtain \tilde{Q}_t by $\tilde{Q}_t, \hat{b}_t \in$

$$\underset{Q_t=Q_t^\top, b_t}{\operatorname{argmin}} \sum_{i=1}^n (\|z_t^{(i)}\|_{Q_t}^2 + \|u_t^{(i)}\|_{R_t^*}^2 + b_t - c_t^{(i)})^2, \quad (3.5)$$

and set $\hat{Q}_t = U \max(\Lambda, 0) U^\top$, where $\tilde{Q}_t = U \Lambda U^\top$ is its eigenvalue decomposition

- 5: **Return:** system parameters $((\hat{A}_t, \hat{B}_t)_{t=0}^{T-1}, (\hat{Q}_t)_{t=0}^T)$
-

2020) use cost supervision as one loss term; the special structure of LQG allows us to use it alone and have theoretical guarantees.

CoREL (Algorithm 2) is the core of our algorithm. Once the state representation function $(\hat{M}_t)_{t=0}^T$ is obtained, SysID (Algorithm 3) identifies the latent system using linear and quadratic regression, followed by planning using RDE (2.4) to obtain controller $(\hat{K}_t)_{t=0}^{T-1}$ from $((\hat{A}_t, \hat{B}_t, R_t^*)_{t=0}^{T-1}, (\hat{Q}_t)_{t=0}^T)$. SysID consists of the standard regression procedures; the full algorithmic detail is shown in Algorithm 3. Below we explain the cost-driven state representation learning algorithm (CoREL, Algorithm 2) in detail.

3.1 Learning the state representation function

The state representation function is learned via CoREL (Algorithm 2). Given the raw data consisting of n trajectories, CoREL first solves the regression problem (3.3) to recover the symmetric matrix \hat{N}_t . The target \bar{c}_t of regression (3.3) is defined by

$$\bar{c}_t := c_t + c_{t+1} + \dots + c_{t+k-1},$$

where $k = 1$ for $0 \leq t \leq \ell - 1$ and $k = m \wedge (T - t + 1)$ for $\ell \leq t \leq T$. The superscript in $\bar{c}_t^{(i)}$ denotes the observed \bar{c}_t in the i th trajectory. The quadratic regression has a closed-form solution, by converting it to linear regression using $\|v\|_P^2 = \langle vv^\top, P \rangle_F = \langle \text{svec}(vv^\top), \text{svec}(P) \rangle$.

Why cumulative cost? The state representation function is parameterized by $(M_t)_{t=0}^T$ and the latent state at step t is given by $z_t = M_t h_t$. The single-step cost prediction (neglecting control cost $\|u_t\|_{R_t}^2$ and constant b_t) is given by $\|z_t\|_{Q_t}^2 = h_t^\top M_t^\top Q_t M_t h_t$. The regression recovers $(M_t^*)^\top Q_t^* M_t^*$ as a whole, from which we can recover $(Q_t^*)^{1/2} M_t^*$ up to an orthogonal transformation. If Q_t^* is positive definite and known, then we can further recover M_t^* from it. However, if Q_t^* does not have full rank, information about M_t^* is partially lost, and there is no way to fully recover M_t^* even if Q_t^* is known. To see why multi-step cumulative cost helps, define $\bar{Q}_t^* := \sum_{\tau=t}^{t+k-1} \Phi_{\tau,t}^\top Q_\tau^* \Phi_{\tau,t}$ for the same k above. Under zero control and zero noise, starting from x_t at step t , the k -step cumulative cost is precisely $\|x_t\|_{\bar{Q}_t^*}^2$. Under the cost observability assumption (Assumption 3), $(\bar{Q}_t^*)_{t=0}^T$ are positive definite.

The normalized parameterization. Still, since \bar{Q}_t^* is unknown, even if we recover $(M_t^*)^\top \bar{Q}_t^* M_t^*$ as a whole, it is not viable to extract M_t^* and \bar{Q}_t^* . Such ambiguity is unavoidable; in fact, for every \bar{Q}_t^* we choose, there is an equivalent parameterization of the system such that the system response is exactly the same. In partially observable LTI systems, it is well-known that the system parameters can only be recovered up to a similarity transform (Oymak and Ozay, 2019). Since every parameterization is correct, we simply choose $\bar{Q}_t^* = I$, which we refer to as the *normalized parameterization*. Concretely, let us define $x'_t = (\bar{Q}_t^*)^{1/2} x_t$. Then, the new parameterization is given by

$$x'_{t+1} = A_t^{*'} x'_t + B_t^{*'} u_t + w'_t, \quad y_t = C_t^{*'} x'_t + v_t, \quad c'_t(x', u) = \|x'\|_{Q_t^{*'}}^2 + \|u\|_{R_t^{*'}}^2,$$

and $c'_T(x') = \|x'\|_{(Q_T^*)'}$, where for all $t \geq 0$,

$$\begin{aligned} A_t^{*'} &= (\bar{Q}_{t+1}^*)^{1/2} A_t^* (\bar{Q}_t^*)^{-1/2}, & B_t^{*'} &= (\bar{Q}_{t+1}^*)^{1/2} B_t^*, & C_t^{*'} &= C_t^* (\bar{Q}_t^*)^{-1/2}, \\ w'_t &= (\bar{Q}_{t+1}^*)^{1/2} w_t, & (Q_t^*)' &= (\bar{Q}_t^*)^{-1/2} Q_t^* (\bar{Q}_t^*)^{-1/2}. \end{aligned}$$

It is easy to verify that under the normalized parameterization the system satisfies Assumptions 1, 2, 3, and 4, up to a change of some constants in the bounds. Without loss of generality, we assume system (2.1) is in the normalized parameterization; if not, the recovered state representation function and latent system are with respect to the normalized parameterization.

Low-rank approximate factorization. Regression (3.3) has a closed-form solution; solving it gives $(\hat{N}_t, \hat{b}_t)_{t=0}^T$. Constants $(\hat{b}_t)_{t=0}^T$ account for the variance of the state estimation error, and are

not part of the state representation function; $d_h \times d_h$ symmetric matrices $(\hat{N}_t)_{t=0}^T$ are estimates of $(M_t^*)^\top M_t^*$ under the normalized parameterization, where $d_h = (t+1)d_y + td_u$. M_t^* can only be recovered up to an orthogonal transformation, since for any orthogonal $S \in \mathbb{R}^{d_x \times d_x}$, $(SM_t^*)^\top SM_t^* = (M_t^*)^\top M_t^*$.

We want to recover \tilde{M}_t from \hat{N}_t such that $\hat{N}_t = \tilde{M}_t^\top \tilde{M}_t$. Let $U\Lambda U^\top = \hat{N}_t$ be its eigenvalue decomposition. Let $\Sigma := \max(\Lambda, 0)$ be the positive semidefinite diagonal matrix containing nonnegative eigenvalues, where ‘‘max’’ applies elementwise. If $d_h \leq d_x$, we can construct $\tilde{M}_t = [\Sigma^{1/2}U^\top; 0_{(d_x-d_h) \times d_h}]$ by padding zeros. If $d_h > d_x$, however, $\text{rank}(\hat{N}_t)$ may exceed d_x . Without loss of generality, assume that the diagonal elements of Σ are in descending order. Let Σ_{d_x} be the left-top $d_x \times d_x$ block of Σ and U_{d_x} be the left d_x columns of U . By the Eckart-Young-Mirsky theorem, $\tilde{M}_t = \Sigma_{d_x}^{1/2}U_{d_x}^\top$ provides the best approximation of \hat{N}_t with $\tilde{M}_t^\top \tilde{M}_t$ among $d_x \times d_h$ matrices in terms of the Frobenius norm.

Why singular value truncation in the first ℓ steps? The latent states are used to identify the latent system dynamics, so whether they are sufficiently excited, namely having full-rank covariance, makes a big difference: if not, the system matrices can only be identified partially. Proposition 2 below confirms that the optimal latent state $z_t^* = M_t^* h_t$ indeed has full-rank covariance for $t \geq \ell$.

Proposition 2. *If system (2.1) satisfies Assumptions 2 (controllability) and 4 (regularity), then under control $(u_t)_{t=0}^{T-1}$, where $u_t \sim \mathcal{N}(0, \sigma_u^2 I)$, $\sigma_{\min}(\text{Cov}(z_t^*)) = \Omega(v^2)$, M_t^* has rank d_x and $\sigma_{\min}(M_t^*) = \Omega(vt^{-1/2})$ for all $\ell \leq t \leq T$.*

Proof. For $\ell \leq t \leq T$, unrolling the Kalman filter gives

$$\begin{aligned} z_t^* &= A_{t-1}^* z_{t-1}^* + B_{t-1}^* u_{t-1} + L_t^* i_t \\ &= A_{t-1}^* (A_{t-2}^* z_{t-2}^* + B_{t-2}^* u_{t-2} + L_{t-1}^* i_{t-1}) + B_{t-1}^* u_{t-1} + L_t^* i_t \\ &= [B_{t-1}^*, \dots, A_{t-1}^* A_{t-2}^* \dots A_{t-\ell+1}^* B_{t-\ell}^*] [u_{t-1}; \dots; u_{t-\ell}] + A_{t-1}^* A_{t-2}^* \dots A_{t-\ell}^* z_{t-\ell}^* \\ &\quad + [L_t^*, \dots, A_{t-1}^* A_{t-2}^* \dots A_{t-\ell+1}^* L_{t-\ell+1}^*] [i_t; \dots; i_{t-\ell+1}], \end{aligned}$$

where $(u_\tau)_{\tau=t-\ell}^{t-1}$, $z_{t-\ell}^*$ and $(i_\tau)_{\tau=t-\ell+1}^t$ are independent. The matrix multiplied by $[u_{t-1}; \dots; u_{t-\ell}]$ is precisely the controllability matrix $\Phi_{t-1, \ell}^c$. Then

$$\begin{aligned} \text{Cov}(z_t^*) &= \mathbb{E}[z_t^* (z_t^*)^\top] \succcurlyeq \Phi_{t-1, \ell}^c \mathbb{E}[[u_{t-1}; \dots; u_{t-\ell}][u_{t-1}; \dots; u_{t-\ell}]^\top] (\Phi_{t-1, \ell}^c)^\top \\ &= \sigma_u^2 \Phi_{t-1, \ell}^c (\Phi_{t-1, \ell}^c)^\top. \end{aligned}$$

By the controllability assumption (Assumption 2), $\text{Cov}(z_t^*)$ has full rank and

$$\sigma_{\min}(\text{Cov}(z_t^*)) \geq \sigma_u^2 v^2.$$

On the other hand, since $z_t^* = M_t^* h_t$,

$$\text{Cov}(z_t^*) = \mathbb{E}[M_t^* h_t h_t^\top (M_t^*)^\top] \preccurlyeq \sigma_{\max}(\mathbb{E}[h_t h_t^\top]) M_t^* (M_t^*)^\top.$$

Since $h_t = [y_{0:t}; u_{0:(t-1)}]$ and $(\text{Cov}(y_t))_{t=0}^T$, $(\text{Cov}(u_t))_{t=0}^{T-1}$ have $\mathcal{O}(1)$ operator norms, by Lemma 8, $\|\text{Cov}(h_t)\| = \|\mathbb{E}[h_t h_t^\top]\| = \mathcal{O}(t)$. Hence,

$$0 < \sigma_u^2 v^2 \leq \sigma_{\min}(\text{Cov}(z_t^*)) = \mathcal{O}(t) \sigma_{d_x}^2 (M_t^*).$$

This implies that $\text{rank}(M_t^*) = d_x$ and $\sigma_{\min}(M_t^*) = \Omega(vt^{-1/2})$. \square

Proposition 2 implies that for all $\ell \leq t \leq T$, N_t^* has rank d_x , so if d_x is not provided, this gives a way to discover it. For $\ell \leq t \leq T$, Proposition 2 guarantees that as long as \tilde{M}_t is close enough to M_t^* , it also has full rank, and so does $\text{Cov}(\tilde{M}_t h_t)$. Hence, we simply take the final estimate $\hat{M}_t = \tilde{M}_t$. Without further assumptions, however, there is no such a full-rank guarantee for $(\text{Cov}(z_t^*))_{t=0}^{\ell-1}$ and $(M_t^*)_{t=0}^{\ell-1}$. We make the following minimal assumption to ensure that the minimum positive singular values $(\sigma_{\min}^+(\text{Cov}(z_t^*)))_{t=0}^{\ell-1}$ are uniformly lower bounded. Note that $(\text{Cov}(z_t^*))_{t=0}^{\ell-1}$ are not required to have full rank.

Assumption 5. For $0 \leq t \leq \ell - 1$, $\sigma_{\min}^+(M_t^*) \geq \beta > 0$.

Still, for $0 \leq t \leq \ell - 1$, Assumption 5 does not guarantee the full-rankness of $\text{Cov}(\tilde{M}_t h_t)$, not even a lower bound on its minimum positive singular value; that is why we introduce TRUNCsv that truncates the singular values of \tilde{M}_t by a threshold $\theta > 0$. Concretely, we take $\hat{M}_t = (\mathbb{I}_{[\theta, +\infty)}(\Sigma_{d_x}^{1/2}) \odot \Sigma_{d_x}^{1/2}) U_{d_x}^\top$. Then, \hat{M}_t has the same singular values as \tilde{M}_t except that those below θ are zeroed. We take $\theta = \Theta(\ell(d_y + d_u)d_x^{3/4}n^{-1/4} \log^{1/4}(1/p))$ to ensure a sufficient lower bound on the minimum positive singular value of \hat{M}_t , without increasing the statistical errors.

4 Theoretical guarantees

Theorem 1 below offers a finite-sample guarantee for our approach. Overall, it confirms cost-driven latent model learning (Algorithm 1) as a viable path to solving LQG control.

Theorem 1. *Given an unknown LQG system (2.1), under Assumptions 1, 2, 3, 4 and 5, if we run Algorithm 1 with $n \geq \text{poly}(T, d_x, d_y, d_u, \log(1/p))$, then with probability at least $1 - p$, state representation function $(\hat{M}_t)_{t=0}^T$ is $\text{poly}(\ell, d_x, d_y, d_u)n^{-1/4}$ optimal in the first ℓ steps, and $\text{poly}(v^{-1}, T, d_x, d_y, d_u)n^{-1/2}$ optimal in the next $(T - \ell)$ steps. Also, the learned controller $(\hat{K}_t)_{t=0}^{T-1}$ is $\text{poly}(\ell, m, d_x, d_y, d_u)c^\ell n^{-1/4}$ optimal for some dimension-free constant $c > 0$ depending on system parameters in the first ℓ steps, and $\text{poly}(T, m, d_x, d_y, d_u, \log(1/p))n^{-1}$ optimal in the last $(T - \ell)$ steps.*

From Theorem 1, we observe a separation of the sample complexities before and after time step ℓ , resulting from the loss of the full-rankness of $(\text{Cov}(z_t^*))_{t=0}^{\ell-1}$ and $(M_t^*)_{t=0}^{\ell-1}$. In more detail, a proof sketch goes as follows. Quadratic regression guarantees that \hat{N}_t converges to N_t^* at a rate of $n^{-1/2}$ for all $0 \leq t \leq T$. Before step ℓ , \hat{M}_t suffers a square root decay of the rate to $n^{-1/4}$ because M_t^* may not have rank d_x . Since $(\hat{z}_t)_{t=0}^{\ell-1}$ may not have full-rank covariances, $(A_t^*)_{t=0}^{\ell-1}$ are only recovered partially. As a result, $(\hat{K}_t)_{t=0}^{\ell-1}$ may not stabilize $(A_t^*, B_t^*)_{t=0}^{\ell-1}$, causing the exponential dependence on ℓ . This means if n is not big enough, this controller may be inferior to zero control, since the system $(A_t^*, B_t^*)_{t=0}^{\ell-1}$ is uniformly exponential stable (Assumption 1) and zero control has suboptimality gap linear in ℓ . After step ℓ , \hat{M}_t retains the $n^{-1/2}$ sample complexity, and so do (\hat{A}_t, \hat{B}_t) ; the certainty equivalent controller then has an order of n^{-1} suboptimality gap for linear quadratic control (Mania et al., 2019).

Theorem 1 states the guarantees for the state representation function $(\hat{M}_t)_{t=0}^T$ and the controller $(\hat{K}_t)_{t=0}^{T-1}$ separately. One may wonder the suboptimality gap of $\hat{\pi} = (\hat{M}_t, \hat{K}_t)_{t=0}^{T-1}$ in combination; after all, this is the output policy. The new challenge is that a suboptimal controller is applied

to a suboptimal state estimation. An exact analysis requires more effort, but a reasonable conjecture is that $(\hat{M}_t, \hat{K}_t)_{t=0}^{T-1}$ has the same order of suboptimality gap as $(\hat{K}_t)_{t=0}^{T-1}$: before step ℓ , the extra suboptimality gap resulted from $(\hat{M}_t)_{t=0}^{\ell-1}$ can be analyzed by considering perturbation $\hat{K}_t(\hat{M}_t - M_t^*)h_t$ on controls; after step ℓ , similar to the analysis of the LQG suboptimality gap in (Mania et al., 2019), the overall suboptimality gap can be analyzed by a Taylor expansion of the value function at $(M_t^*, K_t^*)_{t=\ell}^{T-1}$, with $(\hat{K}_t\hat{M}_t - K_t^*M_t^*)_{t=\ell}^{T-1}$ being perturbations.

4.1 Proof of Theorem 1

Proof. Recall that Algorithm 1 has three main steps: state representation learning (CoREL, Algorithm 2), latent system identification (SysID, Algorithm 3), and planning by RDE (2.4). Correspondingly, the analysis below is organized around these three steps.

Recovery of the state representation function. By Proposition 3, with $u_t = \mathcal{N}(0, \sigma_u^2 I)$ for all $0 \leq t \leq T-1$ to system (2.1), the k -step cumulative cost starting from step t , where $k = 1$ for $0 \leq t \leq \ell-1$ and $k = m \wedge (T-t+1)$ for $\ell \leq t \leq T$, is given under the normalized parameterization by

$$\bar{c}_t := c_t + c_{t+1} + \dots + c_{t+k-1} = \|z_t^{*'}\|^2 + \sum_{\tau=t}^{t+k-1} \|u_\tau\|_{\mathbb{R}^r}^2 + \bar{b}_t + \bar{e}_t,$$

where $\bar{b}_t = \mathcal{O}(k)$, and \bar{e}_t is a zero-mean subexponential random variable with $\|\bar{e}_t\|_{\psi_1} = \mathcal{O}(kd_x^{1/2})$. Then, it is clear that Algorithm 2 recovers latent states $z_t^{*'} = M_t^{*'}h_t$, where $0 \leq t \leq T$, by a combination of quadratic regression and low-rank approximate factorization. Below we drop the superscript prime for notational simplicity, but keep in mind that the optimal state representation function $(M_t^*)_{t=0}^T$, the corresponding latent states $(z_t^*)_{t=0}^T$, and the true latent system parameters $((A_t^*, B_t^*)_{t=0}^{T-1}, (Q_t^*)_{t=0}^T)$ are all with respect to the normalized parameterization.

Let $N_t^* := (M_t^*)^\top M_t^*$ for all $0 \leq t \leq T$. For quadratic regression, Lemma 2 (detailed later) and the union bound over all time steps guarantees that as long as $n \geq aT^4(d_y + d_u)^4 \log(aT^3(d_y + d_u)^2/p) \log(T/p)$ for an absolute constant $a > 0$, with probability at least $1 - p$, for all $0 \leq t \leq T$,

$$\|\hat{N}_t - N_t^*\|_F = \mathcal{O}(kt^2(d_y + d_u)^2 d_x^{1/2} n^{-1/2} \log^{1/2}(1/p)).$$

Now let us bound the distance between \hat{M}_t and M_t^* . Recall that we use $d_h = (t+1)d_y + td_u$ as a shorthand. The estimate \hat{N}_t may not be positive semidefinite. Let $\hat{N}_t = U\Lambda U^\top$ be its eigenvalue decomposition, with the $d_h \times d_h$ matrix Λ having descending diagonal elements. Let $\Sigma := \max(\Lambda, 0)$. Then $\tilde{N}_t := U\Sigma U^\top$ is the projection of \hat{N}_t onto the positive semidefinite cone (Boyd and Vandenberghe, 2004, Section 8.1.1) with respect to the Frobenius norm. Since $N_t^* \succcurlyeq 0$,

$$\|\tilde{N}_t - N_t^*\|_F \leq \|\hat{N}_t - N_t^*\|_F.$$

The low-rank factorization is essentially a combination of low-rank approximation and matrix factorization. For $d_h < d_x$, $\tilde{M}_t := [\Sigma^{1/2}U^\top; 0_{(d_x-d_h) \times d_h}]$ constructed by padding zeros satisfies $\tilde{M}_t^\top \tilde{M}_t = \tilde{N}_t$. For $d_h \geq d_x$, construct $\tilde{M}_t := \Sigma_{d_x}^{1/2}U_{d_x}^\top$, where Σ_{d_x} is the left-top $d_x \times d_x$ block

in Σ and U_{d_x} consists of d_x columns of U from the left. By the Eckart-Young-Mirsky theorem, $\tilde{M}_t^\top \tilde{M}_t = U_{d_x}^\top \Sigma_{d_x} U_{d_x}$ satisfies

$$\|\tilde{M}_t^\top \tilde{M}_t - \tilde{N}_t\|_F = \min_{N \in \mathbb{R}^{d_h \times d_h}, \text{rank}(N) \leq d_x} \|N - \tilde{N}_t\|_F.$$

Hence,

$$\begin{aligned} \|\tilde{M}_t^\top \tilde{M}_t - N_t^*\|_F &\leq \|\tilde{M}_t^\top \tilde{M}_t - \tilde{N}_t\|_F + \|\tilde{N}_t - N_t^*\|_F \\ &\leq 2\|\tilde{N}_t - N_t^*\|_F \leq 2\|\hat{N}_t - N_t^*\|_F \\ &= \mathcal{O}(kt^2(d_y + d_u)^2 d_x^{1/2} n^{-1/2} \log^{1/2}(1/p)). \end{aligned}$$

From now on, we consider $0 \leq t \leq \ell - 1$ and $\ell \leq t \leq T$ separately, since, as we will show, in the latter case we have the additional condition that $\text{rank}(M_t^*) = d_x$.

For $0 \leq t \leq \ell - 1, k = 1$. By Lemma 4 to be proved later, there exists a $d_x \times d_x$ orthogonal matrix S_t , such that $\|\tilde{M}_t - S_t M_t^*\|_2 \leq \|\tilde{M}_t - S_t M_t^*\|_F = \mathcal{O}(t(d_y + d_u) d_x^{3/4} n^{-1/4} \log^{1/4}(1/p))$. Recall that $\hat{M}_t = \text{TRUNC SV}(\tilde{M}_t, \theta)$; that is, $\hat{M}_t = (\mathbb{I}_{[\theta, +\infty)}(\Sigma_{d_x}^{1/2}) \odot \Sigma_{d_x}^{1/2}) U_{d_x}^\top$. Then

$$\begin{aligned} \|\hat{M}_t - \tilde{M}_t\|_2 &= \|(\mathbb{I}_{[\theta, +\infty)}(\Sigma_{d_x}^{1/2}) \odot \Sigma_{d_x}^{1/2} - \Sigma_{d_x}^{1/2}) U_{d_x}^\top\|_2 \\ &\leq \|\mathbb{I}_{[\theta, +\infty)}(\Sigma_{d_x}^{1/2}) \odot \Sigma_{d_x}^{1/2} - \Sigma_{d_x}^{1/2}\|_2 \leq \theta. \end{aligned}$$

Hence, the distance between \hat{M}_t and M_t^* satisfies

$$\begin{aligned} \|\hat{M}_t - S_t M_t^*\|_2 &= \|\hat{M}_t - \tilde{M}_t + \tilde{M}_t - S_t M_t^*\|_2 \\ &\leq \|\hat{M}_t - \tilde{M}_t\|_2 + \|\tilde{M}_t - S_t M_t^*\|_2 \\ &\leq \theta + \mathcal{O}(t(d_y + d_u) d_x^{3/4} n^{-1/4} \log^{1/4}(1/p)). \end{aligned}$$

θ should be chosen in such a way that it keeps the error on the same order; that is, $\theta = \mathcal{O}(t(d_y + d_u) d_x^{3/4} n^{-1/4} \log^{1/4}(1/p))$. As a result, since $\hat{z}_t = \hat{M}_t h_t$ and $z_t^* = M_t^* h_t$,

$$\|\hat{z}_t - S_t z_t^*\| = \|(\hat{M}_t - S_t M_t^*) h_t\| \leq \|\hat{M}_t - S_t M_t^*\|_2 \|h_t\|.$$

Since $h_t = [y_{0:t}; u_{0:(t-1)}]$ whose ℓ_2 -norm is sub-Gaussian with its mean and sub-Gaussian norm bounded by $\mathcal{O}((t(d_y + d_u))^{1/2})$, $\|\hat{z}_t - S_t z_t^*\|$ is sub-Gaussian with its mean and sub-Gaussian norm bounded by

$$\mathcal{O}((t(d_y + d_u))^{3/2} d_x^{3/4} n^{-1/4} \log^{1/4}(1/p)).$$

On the other hand, threshold θ ensures a lower bound on $\sigma_{\min}^+(\sum_{i=1}^n z_t^{(i)} (z_t^{(i)})^\top)$. As shown in the proof of Lemma 5, this property is important for ensuring the system identification outputs \hat{A}_t and \hat{B}_t have bounded norms. Specifically,

$$\begin{aligned} \sigma_{\min}^+(\sum_{i=1}^n z_t^{(i)} (z_t^{(i)})^\top) &= \sigma_{\min}^+(\sum_{i=1}^n \hat{M}_t h_t^{(i)} (h_t^{(i)})^\top \hat{M}_t^\top) \\ &= (\sigma_{\min}^+(\hat{M}_t))^2 \sigma_{\min}^+(\sum_{i=1}^n h_t^{(i)} (h_t^{(i)})^\top) \\ &\stackrel{(i)}{=} \theta^2 \cdot \Omega(n) = \Omega(\theta^2 n), \end{aligned}$$

where (i) holds with probability at least $1 - p$, as long as $n \geq a \log(1/p)$ for some absolute constant $a > 0$. This concentration is standard in analyzing linear regression (Wainwright, 2019).

For $\ell \leq t \leq T$, $k \leq m$. By Proposition 2, $\text{Cov}(z_t^*)$ has full rank, $\sigma_{\min}(\text{Cov}(z_t^*)) = \Omega(v^2)$ and $\sigma_{\min}(M_t^*) = \Omega(vt^{-1/2})$. Recall that for $\ell \leq t \leq T$, we simply set $\hat{M}_t = \tilde{M}_t$. Then, by Lemma 3, there exists a $d_x \times d_x$ orthogonal matrix S_t , such that

$$\begin{aligned} \|\hat{M}_t - S_t M_t^*\|_F &= \mathcal{O}(\sigma_{\min}^{-1}(M_t^*)) \|\tilde{M}_t^\top \tilde{M}_t - N_t^*\|_F \\ &= \mathcal{O}(v^{-1} m t^{5/2} (d_y + d_u)^2 d_x^{1/2} n^{-1/2} \log^{1/2}(1/p)), \end{aligned}$$

which is also an upper bound on $\|\hat{M}_t - S_t M_t^*\|_2$. As a result,

$$\|\hat{z}_t - S_t z_t^*\|_2 \leq \|\hat{M}_t - S_t M_t^*\|_2 \|h_t\|,$$

from which we have that $\|\hat{z}_t - S_t z_t^*\|_2$ is sub-Gaussian with its mean and sub-Gaussian norm bounded by

$$\mathcal{O}(v^{-1} m t^3 (d_y + d_u)^{5/2} d_x^{1/2} n^{-1/2} \log^{1/2}(1/p)).$$

Consider

$$\begin{aligned} \left\| \sum_{i=1}^n \hat{z}_t^{(i)} (\hat{z}_t^{(i)})^\top - S_t (z_t^*)^{(i)} ((z_t^*)^{(i)})^\top S_t^\top \right\|_2 &= \sum_{i=1}^n (\|\hat{z}_t^{(i)}\| + \|(z_t^*)^{(i)}\|) \cdot \|\hat{z}_t^{(i)} - S_t (z_t^*)^{(i)}\| \\ &\stackrel{(i)}{=} n d_x^{1/2} \log^{1/2}(n/p) \cdot \|\hat{z}_t^{(i)} - S_t (z_t^*)^{(i)}\|, \end{aligned}$$

where (i) holds with probability $1 - p$. Hence, there exists an absolute constant $c > 0$, such that if $n \geq c v^{-6} m^2 T^6 (d_y + d_u)^5 d_x^2 \log^2(n/p)$,

$$\begin{aligned} \sigma_{\min} \left(\sum_{i=1}^n \hat{z}_t^{(i)} (\hat{z}_t^{(i)})^\top \right) &\geq \sigma_{\min} \left(\sum_{i=1}^n (z_t^*)^{(i)} ((z_t^*)^{(i)})^\top \right) \\ &\quad - \left\| \sum_{i=1}^n \hat{z}_t^{(i)} (\hat{z}_t^{(i)})^\top - S_t (z_t^*)^{(i)} ((z_t^*)^{(i)})^\top S_t^\top \right\| \\ &\geq \sigma_{\min} \left(\sum_{i=1}^n (z_t^*)^{(i)} ((z_t^*)^{(i)})^\top \right) / 2 = \Omega(v^2 n). \end{aligned}$$

This is needed for the analysis of latent system identification in the next step.

Recovery of the latent dynamics. The latent system $(A_t^*, B_t^*)_{t=0}^{T-1}$ is identified in Algorithm 3, using $(\hat{z}_t^{(i)})_{i=1, t=0}^{N, T}$ produced by Algorithm 2, by ordinary least squares. Recall from Proposition 1 that $z_{t+1}^* = A_t^* z_t^* + B_t^* u_t + L_{t+1} i_{t+1}$. With the transforms on z_t^* and z_{t+1}^* , we have

$$S_{t+1} z_{t+1}^* = (S_{t+1} A_t^* S_t^\top) S_t z_t^* + S_{t+1} B_t^* u_t + S_{t+1} L_{t+1} i_{t+1},$$

and $(z_t^*)^\top Q_t^* z_t^* = (S_t z_t^*)^\top S_t Q_t^* S_t^\top S_t z_t^*$. Under control $u_t \sim \mathcal{N}(0, \sigma_u^2 I_{d_u})$ for $0 \leq t \leq T-1$, we know that z_t^* is a zero-mean Gaussian random vector; so is $S_t z_t^*$. Let $\Sigma_t^* = \mathbb{E}[S_t z_t^* (z_t^*)^\top S_t^\top]$ be its covariance.

For $0 \leq t \leq \ell - 1$, we need a bound for the estimation error of rank-deficient and noisy linear regression. By Lemma 5, there exists an absolute constant $c > 0$, such that as long as $n \geq c(d_x + d_u + \log(1/p))$, with probability at least $1 - p$,

$$\begin{aligned} &\|([\hat{A}_t, \hat{B}_t] - [S_{t+1} A_t^* S_t^\top, S_{t+1} B_t^*]) \text{diag}((\Sigma_t^*)^{1/2}, \sigma_u I_{d_u})\|_2 \\ &= \mathcal{O}((1 + \beta^{-1}) \|z_t - S_t z_t^*\| \log^{1/2}(n/p) + n^{-1/2} (d_x + d_u + \log(1/p))^{1/2}), \end{aligned}$$

which implies

$$\|(\hat{A}_t - S_{t+1}A_t^*S_t^\top)(\Sigma_t^*)^{1/2}\|_2 = \mathcal{O}((1 + \beta^{-1})(t(d_y + d_u))^{3/2}d_x^{3/4}n^{-1/4}\log^{3/4}(n/p)),$$

which is also a bound on $\|\hat{B}_t - S_{t+1}B_t^*\|_2$.

For $\ell \leq t \leq T - 1$, by Lemma 6, with probability at least $1 - p$,

$$\begin{aligned} & \|[\hat{A}_t, \hat{B}_t] - [S_{t+1}A_t^*S_t^\top, S_{t+1}B_t^*]\|_2 \\ &= \mathcal{O}((\nu^{-1} + \sigma_u^{-1})(\|\hat{z}_t - S_t z_t^*\| \log^{1/2}(n/p) + n^{-1/2}(d_x + d_u + \log(1/p))^{1/2})), \end{aligned}$$

which implies

$$\|\hat{A}_t - S_{t+1}A_t^*S_t^\top\|_2 = \mathcal{O}\left((1 + \nu^{-2})(mt^3(d_y + d_u)^{5/2}d_x^{1/2}\log^{1/2}(1/p))n^{-1/2}\right),$$

which is also a bound on $\|\hat{B}_t - S_{t+1}B_t^*\|_2$. We note that since the observation of u_t is exact, a tighter bound on $\|\hat{B}_t - S_{t+1}B_t^*\|_2$ is possible; we keep the current bound since it does not affect the final bounds in Theorem 1.

By Assumption 3, $(Q_t^*)_{t=0}^{\ell-1}$ and Q_T^* are positive definite; they are identity matrices under the normalized parameterization. For $(Q_t^*)_{t=\ell}^{T-1}$, which may not be positive definite, we identify them in SysID (Algorithm 3) by (3.5). By Lemma 2,

$$\begin{aligned} \|\tilde{Q}_t - S_t Q_t^* S_t^\top\|_F &= \mathcal{O}(md_x^2 \log(n/p) \cdot mt^3(d_y + d_u)^{5/2}d_x^{1/2}n^{-1/2}\log^{1/2}(1/p) \\ &\quad + d_x^2 md_x^{1/2}n^{-1/2}\log^{1/2}(1/p)) \\ &= \mathcal{O}(m^2 t^3 (d_y + d_u)^{5/2} d_x^{5/2} n^{-1/2} \log^{1/2}(1/p)), \end{aligned}$$

where we have used $\nu = \Omega(1)$. By construction, \hat{Q}_t is the projection of \tilde{Q}_t onto the positive semidefinite cone with respect to the Frobenius norm (Boyd and Vandenberghe, 2004, Section 8.1.1). Since $S_t Q_t^* S_t^\top \succcurlyeq 0$,

$$\|\hat{Q}_t - S_t Q_t^* S_t^\top\|_F \leq \|\tilde{Q}_t - S_t Q_t^* S_t^\top\|_F.$$

Certainty equivalent linear quadratic control. The last step of Algorithm 1 is to compute the optimal controller in the estimated system $((\hat{A}_t, \hat{B}_t, R_t^*)_{t=0}^{T-1}, (\hat{Q}_t)_{t=0}^T)$ by RDE (2.4).

RDE (2.4) proceeds backward. For $\ell \leq t \leq T - 1$, $\|\hat{A}_t - S_{t+1}A_t^*S_t^\top\|_2$, $\|\hat{B}_t - S_{t+1}B_t^*\|_2$ and $\|\hat{Q}_t - S_t Q_t^* S_t^\top\|_2$ are all bounded by

$$\mathcal{O}(m^2 T^3 (d_y + d_u)^{5/2} d_x^{5/2} n^{-1/2} \log(1/p)^{1/2}).$$

By Lemma 10 on certainty equivalent linear quadratic control, for $n \geq \text{poly}(T, d_x, d_y, d_u, \log(1/p))$, where hidden constants are dimension-free and depend polynomially on system parameters, the controller $(\hat{K}_t)_{t=\ell}^{T-1}$ is ϵ -optimal in system $((S_{t+1}A_t^*S_t^{-1}, S_{t+1}B_t^*, R_t^*)_{t=\ell}^{T-1}, (S_t Q_t^* S_t^\top)_{t=\ell}^T)$, for

$$\epsilon = \mathcal{O}((d_x \wedge d_u)(d_y + d_u)^5 d_x^5 m^4 T^7 \log(1/p) n^{-1}), \quad (4.1)$$

that is, if

$$n \geq cm^4(d_y + d_u)^5 d_x^6 T^7 \log(1/p) \epsilon^{-1},$$

for a dimension-free constant $c > 0$ that depends on system parameters.

For $0 \leq t \leq \ell - 1$, $(\hat{K}_t)_{t=0}^{\ell-1}$ and $(K_t^*)_{t=0}^{\ell-1}$ are optimal controllers in the ℓ -step systems with terminal costs given by \hat{P}_ℓ and P_ℓ^* , respectively, where \hat{P}_ℓ and P_ℓ^* are the solutions to RDE (2.4) in systems $((\hat{A}_t, \hat{B}_t, R_t^*)_{t=0}^{T-1}, (\hat{Q}_t)_{t=0}^T)$ and $((S_{t+1}A_t^*S_t^{-1}, S_{t+1}B_t^*, R_t^*)_{t=0}^{T-1}, (S_tQ_t^*S_t^\top)_{t=0}^T)$ at step ℓ , respectively. By Lemma 10,

$$\|\hat{P}_\ell - P_\ell^*\|_2 = \mathcal{O}(m^2 t^3 (d_y + d_u)^{5/2} d_x^{5/2} n^{-1/2} \log^{1/2}(1/p)),$$

which is dominated by the $n^{-1/4}$ rate of $\|(\hat{A}_t - S_{t+1}A_t^*S_t^\top)(\Sigma_t^*)^{1/2}\|_2$ and $\|\hat{B}_t - S_{t+1}B_t^*\|_2$ for $0 \leq t \leq \ell - 1$. Then, by Lemma 9, $(\hat{K}_t)_{t=0}^{\ell-1}$ is ϵ -optimal in system $((S_{t+1}A_t^*S_t^{-1}, S_{t+1}B_t^*, R_t^*)_{t=0}^{\ell-1}, (S_tQ_t^*S_t^\top)_{t=0}^{\ell-1})$ with terminal cost matrix P_t^* , for

$$\begin{aligned} \epsilon &= \mathcal{O}(a^\ell d_x (\ell(d_y + d_u))^{3/2} d_x^{3/4} n^{-1/4} \log^{3/4}(n/p) + d_x (d_x + d_u + \log(1/p))^{1/2} n^{-1/2}) \\ &= \mathcal{O}(d_x^{7/4} (\ell(d_y + d_u))^{3/2} a^\ell n^{-1/4} \log^{3/4}(n/p)), \end{aligned} \quad (4.2)$$

where dimension-free constant $a > 0$ depends on the system parameters; that is, if

$$n \geq a_0 d_x^7 (d_y + d_u)^6 \ell^6 a_1^\ell \log^3(n/p) \epsilon^{-4},$$

for some dimension-free constants $a_0, a_1 > 0$ that depend on system parameters. The bound (4.2) is worse than (4.1), because for $0 \leq t \leq \ell - 1$, z_t^* does not have full-rank covariance, and $S_{t+1}A_t^*S_t^\top$ is only recovered partially. Even with large enough data, linear regression has no guarantee for $\|\hat{A}_t - S_{t+1}A_t^*S_t^\top\|_2$ to be small; we do not know the controllability of $(\hat{A}_t, \hat{B}_t)_{t=0}^{\ell-1}$, not even its stabilizability. \square

Next, we provide several key technical lemmas regarding the regression and factorization procedures used in our algorithm.

4.2 Quadratic regression bound

As noted in Section 3.1, the quadratic regression can be converted to linear regression using $\|h\|_P^2 = \langle hh^\top, P \rangle_F = \langle \text{svec}(hh^\top), \text{svec}(P) \rangle$. To analyze this linear regression with an intercept, we need the following lemma. We note that a similar lemma without considering the intercept has been proved in (Jadbabaie et al., 2021, Proposition).

Lemma 1. *Let $(h_0^{(i)})_{i=1}^n$ be n independent observations of the d -dimensional random vector $h_0 \sim \mathcal{N}(0, I_d)$. Let $f_0^{(i)} := \text{svec}(h_0^{(i)}(h_0^{(i)})^\top)$ and $\bar{f}_0^{(i)} := [f_0^{(i)}; 1]$. There exists an absolute constant $a > 0$, such that as long as $n \geq ar^4 \log(ar^2/p)$, with probability at least $1 - p$,*

$$\sigma_{\min} \left(\sum_{i=1}^n \bar{f}_0^{(i)} (\bar{f}_0^{(i)})^\top \right) = \Omega(d^{-1}n).$$

Proof. Let $f_0 = \text{svec}(h_0 h_0^\top)$ and $\bar{f}_0 = [f_0; 1]$. We first show that $\lambda_{\min}(\mathbb{E}[\bar{f} \bar{f}^\top])$ is lower bounded. Consider

$$\bar{\Sigma} = \mathbb{E}[\bar{f} \bar{f}^\top] = \begin{bmatrix} 2I_{d(d-1)/2} & 0 & 0 \\ 0 & 2I_d + 1_d 1_d^\top & 1_d \\ 0 & 1_d^\top & 1 \end{bmatrix}.$$

To lower bound its smallest eigenvalue, let us compute its inverse. By the Sherman-Morrison formula,

$$(2I_d + 1_d 1_d^\top)^{-1} = (2I_d)^{-1} - \frac{(2I_d)^{-1} 1_d 1_d^\top (2I_d)^{-1}}{1 + 1^\top (2I_d)^{-1} 1} = \frac{1}{2} I_d - \frac{1}{2d+4} 1_d 1_d^\top.$$

Then, by the inverse of a block matrix,

$$\begin{bmatrix} 2I_d + 1_d 1_d^\top & 1_d \\ 1_d^\top & 1 \end{bmatrix}^{-1} = \frac{1}{2} \begin{bmatrix} I_d & -1_d \\ -1_d^\top & d+2 \end{bmatrix}.$$

Therefore,

$$\begin{aligned} \|\bar{\Sigma}^{-1}\|_2 &\leq \frac{1}{2} (\|I_{d(d-1)/2}\|_2 + \|I_d\|_2 + 2\| -1_d\|_2 + \|d+2\|_2) \\ &= \frac{1}{2} d + d^{1/2} + 2 \leq 4d. \end{aligned}$$

Hence,

$$\lambda_{\min}(\bar{\Sigma}) = \|\bar{\Sigma}^{-1}\|_2^{-1} \geq (4d)^{-1} = \Omega(d^{-1}).$$

Then, with the similar concentration arguments to those in (Jadbabaie et al., 2021, Appendix A.1), we can show that with probability at least $1 - p$,

$$\lambda_{\min}(F_0^\top F_0) = \Omega(d^{-1}n).$$

□

Lemma 1 lower bounds the minimum singular value of a matrix that contains the fourth powers of elements in standard Gaussian random vectors. The following lemma is the main result in Section 4.2.

Lemma 2 (Quadratic regression). *Define random variable $c := (h^*)^\top N^* h^* + b^* + e$, where $h^* \sim \mathcal{N}(0, \Sigma_*)$ is a d -dimensional Gaussian random vector, $N^* \in \mathbb{R}^{d \times d}$ is a positive semidefinite matrix, $b^* \in \mathbb{R}$ is a constant and e is a zero-mean subexponential random variable with $\|e\|_{\psi_1} \leq E$. Assume that $\|N^*\|_2$ and $\|\Sigma_*\|_2$ are $\mathcal{O}(1)$. Let $\sigma_{\min}(\Sigma_*^{1/2}) \geq \beta > 0$. Assume that $\beta = \Omega(1)$. Define $h := h^* + \delta$ where the noise vector δ can be correlated with h^* and its ℓ_2 norm is sub-Gaussian with $\mathbb{E}[\|\delta\|] \leq \epsilon$, $\|\|\delta\|\|_{\psi_2} \leq \epsilon$. Assume that $\epsilon \leq \min((d\|\Sigma_*\|^{1/2}), a(\beta \wedge 1)d^{-3/2}\|\Sigma_*\|_2^{-1/2}(\log(n/p))^{-1})$ for some absolute constant $a > 0$. Suppose we get n observations $h^{(i)}$ and $c^{(i)}$ of h and c , where $((h^*)^{(i)})_{i=1}^n$ are independent and $(\delta^{(i)})_{i=1}^n$ can be correlated. Consider the regression problem*

$$\hat{N}, \hat{b} \in \underset{N=N^\top, b}{\operatorname{argmin}} \sum_{i=1}^n (c^{(i)} - \|h^{(i)}\|_N^2 - b)^2. \quad (4.3)$$

There exists an absolute constant $a_0 > 0$, such that as long as $n \geq a_0 d^4 \log(a_0 d^2/p) \log(1/p)$, with probability at least $1 - p$, $\|\hat{N} - N^*\|_F$ and $|\hat{b} - b^*|$ are bounded by

$$\mathcal{O}((d(1 + Ed^{1/2}) \log(n/p) + \|b^*\|)\epsilon + d^2 En^{-1/2} \log^{1/2}(1/p)).$$

Proof. Regression (4.3) can be written as

$$\min_{\text{svec}(N), b} \sum_{i=1}^n (c^{(i)} - (\text{svec}(h^{(i)}(h^{(i)})^\top)^\top \text{svec}(N) - b)^2.$$

Let $f^{(i)} := \text{svec}(h^{(i)}(h^{(i)})^\top)$ denote the covariates and $\bar{f}^{(i)} = [f^{(i)}; 1]$ denote the extended covariates. Define $(f^*)^{(i)}$ and $(\bar{f}^*)^{(i)}$ similarly by replacing $h^{(i)}$ with $(h^*)^{(i)}$. Then, regression (4.3) can be written as

$$\min_{\text{svec}(N), b} \sum_{i=1}^n (c^{(i)} - (f^{(i)})^\top \text{svec}(N) - b)^2. \quad (4.4)$$

Let $\bar{F} := [\bar{f}^{(1)}, \dots, \bar{f}^{(n)}]^\top$ be an $n \times \frac{d(d+3)}{2}$ matrix whose i th row is $(\bar{f}^{(i)})^\top$. Define \bar{F}^* similarly by replacing $\bar{f}^{(i)}$ with $(\bar{f}^*)^{(i)}$. Solving linear regression (4.4) gives

$$\bar{F}^\top \bar{F} [\text{svec}(\hat{N}); \hat{b}] = \sum_{i=1}^n \bar{f}^{(i)} c^{(i)}.$$

Substituting $c^{(i)} = ((\bar{f}^*)^{(i)})^\top [\text{svec}(N^*); b^*] + e^{(i)}$ into the above equation yields

$$\bar{F}^\top \bar{F} [\text{svec}(\hat{N}); \hat{b}] = \bar{F}^\top \bar{F}^* [\text{svec}(N^*); b^*] + \bar{F}^\top \zeta, \quad (4.5)$$

where ζ denotes the vector whose i th element is $e^{(i)}$. Rearranging the terms, we have

$$\bar{F}^\top \bar{F} [\text{svec}(\hat{N} - N^*); \hat{b} - b^*] = \bar{F}^\top (\bar{F}^* - \bar{F}) [\text{svec}(N^*); b^*] + \bar{F}^\top \zeta. \quad (4.6)$$

Next, we show that $\bar{F}^\top \bar{F}$ is invertible with high probability. We can represent h^* by $\Sigma_*^{1/2} h_0$, where $h_0 \sim \mathcal{N}(0, I_d)$ is an d -dimensional standard Gaussian random vector. Correspondingly, an independent observation $(h^*)^{(i)}$ can be expressed as $\Sigma_*^{1/2} h_0^{(i)}$, where $h_0^{(i)}$ is an independent observation of h_0 . It follows that $h^*(h^*)^\top = \Sigma_*^{1/2} h_0 h_0^\top \Sigma_*^{1/2}$ and $(h^*)^{(i)} ((h^*)^{(i)})^\top = \Sigma_*^{1/2} h_0^{(i)} (h_0^{(i)})^\top \Sigma_*^{1/2}$. Define $f_0 := \text{svec}(h_0 h_0^\top)$, $f_0^{(i)} := \text{svec}(h_0^{(i)} (h_0^{(i)})^\top)$, and $F_0 := [f_0^{(1)}, \dots, f_0^{(n)}]^\top$ be an $n \times \frac{r(r+1)}{2}$ matrix whose i th row is $(f_0^{(i)})^\top$. Define $\bar{f}_0, \bar{f}_0^{(i)}$ and \bar{F}_0 as the extended counterparts. Then,

$$f = \text{svec}(hh^\top) = \text{svec}(\Sigma_*^{1/2} h_0 h_0^\top \Sigma_*^{1/2}) = (\Sigma_*^{1/2} \otimes_s \Sigma_*^{1/2}) \text{svec}(h_0 h_0^\top) = \Phi^* f_0,$$

where $\Phi^* := \Sigma_*^{1/2} \otimes_s \Sigma_*^{1/2}$ is a $\frac{d(d+1)}{2} \times \frac{d(d+1)}{2}$ matrix. Then, $F^* = F_0 (\Phi^*)^\top$. By the properties of the symmetric Kronecker product (Schacke, 2004),

$$\sigma_{\min}(\Phi^*) = \sigma_{\min}(\Sigma_*^{1/2})^2 = \sigma_{\min}(\Sigma_*) = \beta > 0.$$

By Lemma 1, there exist absolute constants $a_0, a_1 > 0$, such that if $n \geq a_0 d^4 \log(a_0 d^2/p)$, with probability at least $1 - p$, $\lambda_{\min}(\bar{F}_0^\top \bar{F}_0) \geq a_1 d^{-1} n$. Since $\bar{f} = \text{diag}(\Phi^*, 1) \bar{f}_0$ and $\bar{F}^* = \bar{F}_0 \text{diag}((\Phi^*)^\top, 1)$,

$$\lambda_{\min}((\bar{F}^*)^\top \bar{F}^*) \geq \sigma_{\min}(\text{diag}(\Phi^*, 1))^2 a_1 d^{-1} n = (\beta^2 \wedge 1) a_1 d^{-1} n.$$

By Weyl's inequality,

$$|\sigma_{\min}(\bar{F}) - \sigma_{\min}(\bar{F}^*)| \leq \|\bar{F} - \bar{F}^*\|_2 = \|F - F^*\|_2.$$

Hence, we want to bound $\|F^* - F\|_2$, which satisfies

$$\|F^* - F\|_2^2 \leq \|F^* - F\|_F^2 = \sum_{i=1}^n \|(h^*)^{(i)} ((h^*)^{(i)})^\top - h^{(i)} (h^{(i)})^\top\|_F^2.$$

Since $h^* (h^*)^\top - h h^\top$ has at most rank two, we have

$$\begin{aligned} \|h^* (h^*)^\top - h h^\top\|_F &\leq \sqrt{2} \|h^* (h^*)^\top - h h^\top\|_2 \\ &= \|h^* (h^* - h)^\top + (h^* - h) h^\top\|_2 \\ &\leq (\|h^*\| + \|h\|) \|\delta\|. \end{aligned}$$

Since $h^* \sim \mathcal{N}(0, \Sigma_*)$, $\|h^*\|$ is sub-Gaussian with its mean and sub-Gaussian norm bounded by $\mathcal{O}((d \|\Sigma_*\|)^{1/2})$. Since $\|\delta\|$ is sub-Gaussian with its mean and sub-Gaussian norm bounded by $\epsilon \leq (d \|\Sigma_*\|)^{1/2}$, we conclude that $\|h^* (h^*)^\top - h h^\top\|_F$ is subexponential with its mean and subexponential norm bounded by $\mathcal{O}(\epsilon (d \|\Sigma_*\|)^{1/2})$. Hence, with probability at least $1 - p$,

$$\|h^* (h^*)^\top - h h^\top\|_F^2 = \mathcal{O}(\epsilon^2 d \|\Sigma_*\| \log^2(n/p)).$$

Therefore,

$$\|F^* - F\|_F^2 = \sum_{i=1}^n \|(h^*)^{(i)} ((h^*)^{(i)})^\top - h^{(i)} (h^{(i)})^\top\|_F^2 = \mathcal{O}(\epsilon^2 d \|\Sigma_*\|_2 n \log^2(n/p)).$$

It follows that

$$\|F^* - F\|_2 = \mathcal{O}(\epsilon (d \|\Sigma_*\|_2 n)^{1/2} \log(n/p)).$$

Hence, there exists some absolute constant $a > 0$, such that as long as

$$\epsilon \leq a(\beta \wedge 1) d^{-3/2} \|\Sigma_*\|_2^{-1/2} (\log(n/p))^{-1},$$

we have

$$|\sigma_{\min}(F) - \sigma_{\min}(F^*)| \leq (\beta \wedge 1) a_1 d^{-1/2} n^{1/2} / 2.$$

It follows that

$$\lambda_{\min}(\bar{F}^\top \bar{F}) = \Omega((\beta^2 \wedge 1) d^{-1} n) = \Omega(d^{-1} n).$$

Now, we return to (4.6). By inverting $\bar{F}^\top \bar{F}$, we obtain

$$\begin{aligned} \|\text{svec}(\hat{N} - N^*); \hat{b} - b^*\| &= \|\bar{F}^\top (\bar{F}^* - \bar{F}) [\text{svec}(N^*); b^*] + \bar{F}^\top \xi\| \\ &\leq \underbrace{\|\bar{F}^\top (\bar{F}^* - \bar{F}) [\text{svec}(N^*); b^*]\|}_{(a)} + \underbrace{\|\bar{F}^\top \xi\|}_{(b)}. \end{aligned} \quad (4.7)$$

Term (a) is upper bounded by

$$\begin{aligned} \sigma_{\min}(\bar{F})^{-1} \|\bar{F}^* - \bar{F}\| [\text{svec}(N^*); b^*] &= \mathcal{O}(\sigma_{\min}(\bar{F})^{-1}) \|(F^* - F) \text{svec}(N^*) + \mathbf{1}_n b^*\| \\ &= \mathcal{O}(d^{1/2} n^{-1/2}) (\|(F^* - F) \text{svec}(N^*)\| + \|\mathbf{1}_n b^*\|). \end{aligned}$$

Using arguments similar to those in (Mhammedi et al., 2020, Section B.2.13), we have

$$\begin{aligned} \|(F^* - F) \text{svec}(N^*)\|^2 &= \sum_{i=1}^n \left\langle \text{svec}((h^*)^{(i)} ((h^*)^{(i)})^\top) - \text{svec}(h^{(i)} (h^{(i)})^\top), \text{svec}(N^*) \right\rangle \\ &\stackrel{(i)}{=} \sum_{i=1}^n \left\langle (h^*)^{(i)} ((h^*)^{(i)})^\top - h^{(i)} (h^{(i)})^\top, N^* \right\rangle_F \\ &\stackrel{(ii)}{\leq} \|N^*\|_2^2 \sum_{i=1}^n \|(h^*)^{(i)} ((h^*)^{(i)})^\top - h^{(i)} (h^{(i)})^\top\|_*^2 \\ &\stackrel{(iii)}{\leq} 2 \|N^*\|_2^2 \sum_{i=1}^n \|(h^*)^{(i)} ((h^*)^{(i)})^\top - h^{(i)} (h^{(i)})^\top\|_F^2 \\ &= 2 \|N^*\|_2^2 \|F^* - F\|_F^2, \end{aligned}$$

where $\langle \cdot, \cdot \rangle_F$ denotes the Frobenius product between matrices in (i), $\|\cdot\|_*$ denotes the nuclear norm in (ii), and (iii) follows from the fact that the matrix $(h^*)^{(i)} ((h^*)^{(i)})^\top - h^{(i)} (h^{(i)})^\top$ has at most rank two. Combining with $\|\mathbf{1}_n b^*\| = n^{1/2} \|b^*\|$, we obtain that the term (a) in (4.7) is bounded by

$$\begin{aligned} &\mathcal{O}(d^{1/2} n^{-1/2} \cdot (\|N^*\|_2 \cdot \epsilon (d \|\Sigma_*\|_2)^{1/2} n^{1/2} \log(n/p) + n^{1/2} \|b^*\|)) \\ &= \mathcal{O}(d \log(n/p) + \|b^*\|) \epsilon. \end{aligned}$$

Now we consider term (b) in (4.7):

$$(b) = \|\bar{F}^\top \xi\| \leq \sigma_{\min}(\bar{F}^\top \bar{F})^{-1} \|\bar{F}^\top \xi\| = \mathcal{O}(dn^{-1}) (\|(\bar{F} - \bar{F}^*)^\top \xi\| + \|(\bar{F}^*)^\top \xi\|).$$

Since ξ is a vector of zero-mean subexponential variables with subexponential norms bounded by E , $\|\xi\| = \mathcal{O}(En^{1/2} \log(n/p))$. Hence,

$$\begin{aligned} \|(\bar{F} - \bar{F}^*)^\top \xi\| &\leq \|\bar{F} - \bar{F}^*\|_2 \|\xi\| \\ &= \mathcal{O}(\epsilon (d \|\Sigma_*\|_2 n)^{1/2} \log(n/p)) \cdot \mathcal{O}(En^{1/2} \log(n/p)) \\ &= \mathcal{O}(d^{1/2} En \log^2(n/p) \epsilon). \end{aligned}$$

To bound $\|(\bar{F}^*)^\top \xi\| = \|\text{diag}(\Phi^*, \mathbf{1}) \bar{F}_0^\top \xi\|$, note that $\|\bar{F}_0^\top \xi\| = \|\sum_{i=1}^n \bar{f}_0^{(i)} e^{(i)}\|$. Consider the j th component in the summation $\sum_{i=1}^n [\bar{f}_0^{(i)}]_j (e^{(i)})^{(j)}$. Recall that $f_0 = \text{svec}(h_0 h_0^\top)$, so $[\bar{f}_0]_j$ is either the square of a standard Gaussian random variable, $\sqrt{2}$ times the product of two independent

standard Gaussian random variables, or one. Hence, $[f_0]_j$ is subexponential with mean and $\|[f_0]_j\|_{\psi_1}$ both bounded by $\mathcal{O}(1)$. As a result, the product $[f_0]_j e$ is $\frac{1}{2}$ -sub-Weibull, with the sub-Weibull norm being $\mathcal{O}(E)$. By (Hao et al., 2019, Theorem 3.1),

$$\sum_{i=1}^n [f_0^{(i)}]_j e^{(i)} = \mathcal{O}(En^{1/2} \log^{1/2}(1/p)).$$

Hence, the norm of the $d(d+1)/2$ -dimensional vector $F_0^\top \xi$ is $\mathcal{O}(dEn^{1/2} \log^{1/2}(1/p))$. By the properties of the symmetric Kronecker product (Schacke, 2004), $\|\Phi^*\|_2 = \|\Sigma_*^{1/2}\|_2^2 = \|\Sigma_*\|_2 = \mathcal{O}(1)$. Then,

$$\|(\bar{F}^*)^\top \xi\| = \mathcal{O}(dEn^{1/2} \log^{1/2}(1/p)).$$

Eventually, term (b) is bounded by

$$\begin{aligned} & \mathcal{O}(dn^{-1}) \cdot \mathcal{O}(d^{1/2}En \log(n/p)\epsilon + dEn^{1/2} \log^{1/2}(1/p)) \\ &= \mathcal{O}(d^{3/2}E \log(n/p)\epsilon + d^2En^{-1/2} \log^{1/2}(1/p)). \end{aligned}$$

Combining the bounds on (a) and (b), we have

$$\begin{aligned} & \|[\text{svec}(\hat{N} - N^*); \hat{b} - b^*]\| \\ &= \mathcal{O}((d \log(n/p) + \|b^*\| + d^{3/2}E \log(n/p))\epsilon + d^2En^{-1/2} \log^{1/2}(1/p)) \\ &= \mathcal{O}((d(1 + Ed^{1/2}) \log(n/p) + \|b^*\|)\epsilon + d^2En^{-1/2} \log^{1/2}(1/p)), \end{aligned}$$

which concludes the proof. \square

4.3 Matrix factorization bound

Given two $m \times n$ matrices A, B , we are interested in bounding $\min_{S^\top S=I} \|SA - B\|_F$ using $\|A^\top A - B^\top B\|_F$. The minimum problem is known as the orthogonal Procrustes problem, solved in (Schoenemann, 1964; Schönemann, 1966). Specifically, the minimum is attained at $S = UV^\top$, where $U\Sigma V^\top = BA^\top$ is its singular value decomposition.

If $m \leq n$ and $\text{rank}(A) = m$, then the following lemma from (Tu et al., 2016) establishes that the distance between A and B is of the same order of $\|A^\top A - B^\top B\|_F$.

Lemma 3 ((Tu et al., 2016, Lemma 5.4)). *For $m \times n$ matrices A, B , let $\sigma_m(A)$ denote its m th largest singular value. Then*

$$\min_{S^\top S=I} \|SA - B\|_F \leq (2(\sqrt{2} - 1))^{-1/2} \sigma_m(A)^{-1} \|A^\top A - B^\top B\|_F.$$

If $\sigma_{\min}(A)$ equals zero, the above bound becomes vacuous. In general, the following lemma shows that the distance is of the order of the square root of the $\|A^\top A - B^\top B\|_F$, with a multiplicative \sqrt{d} factor, where $d = \min(2m, n)$.

Lemma 4. *For $m \times n$ matrices A, B , $\min_{S^\top S=I} \|SA - B\|_F^2 \leq \sqrt{d} \|A^\top A - B^\top B\|_F$, where $d = \min(2m, n)$.*

Proof. Let $U\Sigma V^\top = BA^\top$ be its singular value decomposition. By substituting the solution UV^\top of the orthogonal Procrustes problem, the square of the attained minimum equals

$$\begin{aligned}\|UV^\top A - B\|_F^2 &= \left\langle UV^\top A - B, UV^\top A - B \right\rangle_F \\ &= \|A\|_F^2 + \|B\|_F^2 - 2 \left\langle UV^\top A, B \right\rangle_F \\ &\stackrel{(i)}{=} \|A\|_F^2 + \|B\|_F^2 - 2\text{tr}(\Sigma) \\ &= \|A^\top A\|_* + \|B^\top B\|_* - 2\|BA^\top\|_*,\end{aligned}$$

where (i) is due to the property of U, V .

To establish the relationship between $\|A^\top A\|_* + \|B^\top B\|_* - 2\|BA^\top\|_*$ and $\|A^\top A - B^\top B\|_F$, we need to operate in the space of singular values. For $m \times n$ matrix M , let $(\sigma_1(M), \dots, \sigma_{d'}(M))$ be its singular values in descending order, where $d' = m \wedge n$.

In terms of singular values,

$$\begin{aligned}\|A^\top A\|_* + \|B^\top B\|_* - 2\|BA^\top\|_* &\stackrel{(i)}{=} \|A^\top A + B^\top B\|_* - 2\|BA^\top\|_* \\ &\stackrel{(ii)}{=} \sum_{i=1}^d \sigma_i(A^\top A + B^\top B) - 2 \sum_{i=1}^{d'} \sigma_i(BA^\top),\end{aligned}$$

where (i) holds since $A^\top A$ and $B^\top B$ are positive semidefinite matrices, and in (ii) $d := \min(2m, n)$ since $\text{rank}(A^\top A + B^\top B) \leq n$ and $\text{rank}(A^\top A + B^\top B) \leq \text{rank}(A^\top A) + \text{rank}(B^\top B) \leq 2m$. If $x \geq y > 0$, then $x - y \leq \sqrt{x^2 - y^2}$. For all $1 \leq i \leq d$, $2\sigma_i(BA^\top) \leq \sigma_i(A^\top A + B^\top B)$ (Bhatia and Kittaneh, 1990). Take $\sigma_i(A^\top A + B^\top B)$ as x and $2\sigma_i(BA^\top)$ as y ; it follows that

$$\sigma_i(A^\top A + B^\top B) - 2\sigma_i(BA^\top) \leq \sqrt{\sigma_i^2(A^\top A + B^\top B) - 4\sigma_i^2(BA^\top)}.$$

Let $\sigma_i(BA^\top) := 0$ for $d' < i \leq d$. Combining the above yields

$$\begin{aligned}\|A^\top A\|_* + \|B^\top B\|_* - 2\|BA^\top\|_* &\leq \sum_{i=1}^d \sqrt{\sigma_i^2(A^\top A + B^\top B) - 4\sigma_i^2(BA^\top)} \\ &\stackrel{(i)}{\leq} \sqrt{d} \sqrt{\sum_{i=1}^d \sigma_i^2(A^\top A + B^\top B) - 4 \sum_{i=1}^{d'} \sigma_i^2(BA^\top)} \\ &\stackrel{(ii)}{\leq} \sqrt{d} \sqrt{\|A^\top A + B^\top B\|_F^2 - 4 \langle A^\top A, B^\top B \rangle_F} \\ &\leq \sqrt{d} \|A^\top A - B^\top B\|_F,\end{aligned}$$

where (i) is due to the Cauchy-Schwarz inequality, and (ii) uses

$$\sum_{i=1}^{d'} \sigma_i^2(BA^\top) = \|BA^\top\|_F^2 = \text{tr}(AB^\top BA^\top) = \text{tr}(A^\top AB^\top B) = \left\langle A^\top A, B^\top B \right\rangle_F.$$

This completes the proof. □

4.4 Linear regression bound

A standard assumption in analyzing linear regression $y = A^*x + e$ is that $\text{Cov}(x)$ has *full rank*. However, as discussed in Section 3.1, we need to handle rank deficient $\text{Cov}(x)$ in the first ℓ steps of system identification.

The following lemma is the main result in Section 4.4.

Lemma 5 (Noisy rank deficient linear regression). *Define random vector $y^* := A^*x^* + e$, where $x^* \sim \mathcal{N}(0, \Sigma_*)$ and $e \sim \mathcal{N}(0, \Sigma_e)$ are d_1 and d_2 dimensional Gaussian random vectors, respectively. Define $x := x^* + \delta_x$ and $y := y^* + \delta_y$ where the noise vectors δ_x and δ_y can be correlated with x^* and y^* , and their ℓ_2 -norms are sub-Gaussian with $\mathbb{E}[\|\delta^x\|] \leq \epsilon_x$, $\|\|\delta^x\|\|_{\psi_2} \leq \epsilon_x$, and $\mathbb{E}[\|\delta^y\|] \leq \epsilon_y$, $\|\|\delta^y\|\|_{\psi_2} \leq \epsilon_y$. Assume that $\|A^*\|_2$, $\|\Sigma_*\|_2$ and $\|\Sigma_e\|_2$ are $\mathcal{O}(1)$. Let $\sigma_{\min}^+(\Sigma_*^{1/2}) \geq \beta > 0$. Suppose we get n observations $x^{(i)}$ and $y^{(i)}$ of x and y , where $((x^*)^{(i)})_{i=1}^n$ are independent and $(\delta_x^{(i)}, \delta_y^{(i)})_{i=1}^n$ can be correlated. Assume that $\sigma_{\min}^+(\sum_{i=1}^n x^{(i)}(x^{(i)})^\top) \geq \theta^2 n$, for some $\theta > 0$ that satisfies $\theta = \Omega(\epsilon_x)$, $\theta = \Omega(\epsilon_y)$ for absolute constants and θ has at least $n^{-1/4}$ dependence on n . Consider the minimum Frobenius norm solution*

$$\hat{A} \in \underset{A}{\operatorname{argmin}} \sum_{i=1}^n \|y^{(i)} - Ax^{(i)}\|^2.$$

Then, there exists an absolute constant $c > 0$, such that if $n \geq c(d_1 + d_2 + \log(1/p))$, with probability at least $1 - p$,

$$\|(\hat{A} - A^*)\Sigma_*^{1/2}\|_2 = \mathcal{O}((\beta^{-1}\epsilon_x + \epsilon_y) \log^{1/2}(n/p) + n^{-1/2}(d_1 + d_2 + \log(1/p))^{1/2}).$$

Proof. Let $r = \text{rank}(\Sigma_*)$ and $\Sigma_* = DD^\top$ where $D \in \mathbb{R}^{d_1 \times r}$. We can view x^* as generated from an r -dimensional standard Gaussian $g \sim \mathcal{N}(0, I_r)$, by $x^* = Dg$; $x^{(i)}$ can then be viewed as $Dg^{(i)} + \delta_x^{(i)}$, where $(g^{(i)})_{i=1}^n$ are independent observations of g . Let X denote the matrix whose i th row is $(x^{(i)})^\top$; $X^*, Y, G, E, \Delta_x, \Delta_y$ are defined similarly.

To solve the regression problem, we set its gradient to be zero and substitute $Y = X^*(A^*)^\top + E + \Delta_y$ to obtain

$$\hat{A}X^\top X = A^*(X^*)^\top X + E^\top X + \Delta_y^\top X.$$

Substituting X by $GD^\top + \Delta_x$ gives

$$\begin{aligned} & \hat{A}DG^\top GD^\top + \hat{A}(DG^\top \Delta_x + \Delta_x^\top GD^\top + \Delta_x^\top \Delta_x) \\ &= A^*DG^\top GD^\top + A^*DG^\top \Delta_x + (E^\top + \Delta_y^\top)(GD^\top + \Delta_x). \end{aligned}$$

By rearranging the terms, we have

$$\begin{aligned} & (\hat{A} - A^*)DG^\top GD^\top \\ &= A^*DG^\top \Delta_x - \hat{A}(DG^\top \Delta_x + \Delta_x^\top GD^\top + \Delta_x^\top \Delta_x) + (E^\top + \Delta_y^\top)(GD^\top + \Delta_x). \end{aligned}$$

Hence,

$$\begin{aligned} & \|(\hat{A} - A^*)D\|_2 \\ &= \|(A^*DG^\top \Delta_x - \hat{A}(DG^\top \Delta_x + \Delta_x^\top GD^\top + \Delta_x^\top \Delta_x) + (E^\top + \Delta_y^\top)(GD^\top + \Delta_x))(D^\top)^\dagger (G^\top G)^{-1}\|_2. \end{aligned}$$

We claim that as long as $n \geq 16(d_1 + d_2 + \log(1/p))$, with probability at least $1 - 4p$, $\|\hat{A}\|_2 = \mathcal{O}(\|A^*\|_2)$ (Claim 1). Then, since $\|D^\dagger\|_2 = (\sigma_{\min}^+(\Sigma_*)^{-1}) \leq \beta^{-1}$,

$$\begin{aligned} & \|(\hat{A} - A^*)D\|_2 \\ &= \mathcal{O}(\beta^{-1}\|A^*\|_2)(\|D\|_2\|G\|_2\|\Delta_x\|_2 + \|\Delta_x\|_2^2)\|(G^\top G)^{-1}\|_2 \\ & \quad + \mathcal{O}(1)(\|G^\dagger E\|_2 + \|G^\dagger \Delta_y\|_2 + \beta^{-1}(\|E\|_2 + \|\Delta_y\|_2)\|\Delta_x\|_2)\|(G^\top G)^{-1}\|_2. \end{aligned}$$

By (Wainwright, 2019, Theorem 6.1), the Gaussian ensemble G satisfies that with probability at least $1 - p$,

$$\begin{aligned} \|G\|_2 &\leq (n\|I_r\|_2)^{1/2} + (\text{tr}(I_r))^{1/2} + (2\|I_r\|_2 \log(1/p))^{1/2} \\ &\leq n^{1/2} + d_1^{1/2} + (2\log(1/p))^{1/2}, \\ \sigma_{\min}(G) &\geq (n\lambda_{\min}(I_r))^{1/2} - (\text{tr}(I_r))^{1/2} - (2\lambda_{\min}(I_r) \log(1/p))^{1/2} \\ &\geq n^{1/2} - d_1^{1/2} - (2\log(1/p))^{1/2}. \end{aligned}$$

Since $n \geq 8d_1 + 16\log(1/p)$, we have $\|G\|_2 = \mathcal{O}(n^{1/2})$ and $\sigma_{\min}(G) = \Omega(n^{1/2})$. It follows that $\|(G^\top G)^{-1}\|_2 = \mathcal{O}(n^{-1})$ and $\|G^\dagger\|_2 = \mathcal{O}(n^{-1/2})$. Similarly, $\|E\|_2 = \mathcal{O}((\|\Sigma_e\|_2 n)^{1/2})$. Note that $\Delta_x^\top \Delta_x = \sum_{i=1}^n \delta_x^{(i)} (\delta_x^{(i)})^\top$. Since $\|\delta_x^{(i)}\|$ is sub-Gaussian with $\mathbb{E}[\|\delta_x^{(i)}\|] \leq \epsilon_x$ and $\|\|\delta_x^{(i)}\|\|_{\psi_2} \leq \epsilon_x$, with probability at least $1 - p$, $\|\delta_x^{(i)}\| = \mathcal{O}(\epsilon_x \log^{1/2}(n/p))$. Hence,

$$\|\Delta_x\|_2 \leq \|\Delta_x\|_F = \left(\sum_{i=1}^n \|\delta_x^{(i)}\|^2\right)^{1/2} = \mathcal{O}(\epsilon_x n^{1/2} \log^{1/2}(n/p)).$$

Similarly, $\|\Delta_y\|_2 = \mathcal{O}(\epsilon_y n^{1/2} \log^{1/2}(n/p))$. Hence,

$$\begin{aligned} \|(\hat{A} - A^*)D\|_2 &= \mathcal{O}(\beta^{-1}\|A^*\|_2)(\|D\|_2\epsilon_x \log^{1/2}(n/p) + \epsilon_x^2 \log(n/p)) \\ & \quad + \mathcal{O}(1)(\|G^\dagger E\|_2 + \epsilon_y \log^{1/2}(n/p) \\ & \quad + \beta^{-1}(\|\Sigma_e^{1/2}\|_2 + \epsilon_y \log^{1/2}(n/p))\epsilon_x \log^{1/2}(n/p)) \\ &= \mathcal{O}(\beta^{-1}\|A^*\|_2\epsilon_x \log^{1/2}(n/p) + \epsilon_y \log^{1/2}(n/p) + \|G^\dagger E\|_2), \end{aligned}$$

where we consider ϵ_x and ϵ_y as quantities much smaller than one such that terms like $\epsilon_x^2, \epsilon_x \epsilon_y$ are absorbed into ϵ_x, ϵ_y . It remains to control $\|G^\dagger E\|_2$. (Mhammedi et al., 2020, Section B.2.11) proves via a covering number argument that with probability at least $1 - p$,

$$\begin{aligned} \|G^\dagger E\|_2 &= \mathcal{O}(1)\sigma_{\min}(G)^{-1}\|\Sigma_e^{1/2}\|_2(d_1 + d_2 + \log(1/p))^{1/2} \\ &= \mathcal{O}(n^{-1/2}(d_1 + d_2 + \log(1/p))^{1/2}). \end{aligned}$$

Overall, we obtain that

$$\begin{aligned} \|(\hat{A} - A^*)\Sigma_*^{1/2}\|_2 &= \|(\hat{A} - A^*)D\|_2 \\ &= \mathcal{O}(\beta^{-1}\|A^*\|_2\epsilon_x \log^{1/2}(n/p) + \epsilon_y \log^{1/2}(n/p) \\ & \quad + n^{-1/2}(d_1 + d_2 + \log(1/p))^{1/2}), \end{aligned}$$

which completes the proof. \square

Claim 1. Under the conditions in Lemma 5, as long as $n \geq 16(d_1 + d_2 + \log(1/p))$, with probability at least $1 - 4p$, $\|\hat{A}\|_2 = \mathcal{O}(\|A^*\|_2)$.

Proof. A minimum norm solution is given by the following closed-form solution of \hat{A} using pseudoinverse (Moore-Penrose inverse):

$$\begin{aligned}\hat{A} &= (A^*(X^*)^\top + E^\top + \Delta_y^\top)X(X^\top X)^\dagger \\ &= A^*(X^\dagger X^*)^\top + (X^\dagger E)^\top + (X^\dagger \Delta_y)^\top \\ &= A^*(X^\dagger X)^\top - A^*(X^\dagger \Delta_x)^\top + (X^\dagger E)^\top + (X^\dagger \Delta_y)^\top.\end{aligned}$$

Then

$$\|\hat{A}\|_2 \leq \|A^*\|_2 + \|X^\dagger E\|_2 + (\|A^*\|_2 \|\Delta_x\|_2 + \|\Delta_y\|_2)(\sigma_{\min}^+(X))^{-1},$$

where we note that $\|X^\dagger\|_2 = \sigma_{\min}^+(X)^{-1}$ when $X \neq 0$. Since $\sigma_{\min}^+(X) = (\sigma_{\min}^+(X^\top X))^{1/2} \geq \theta n^{1/2}$,

$$(\sigma_{\min}^+(X))^{-1} \leq \theta^{-1} n^{-1/2}.$$

We have shown in the proof of Lemma 5 that with probability at least $1 - 2p$,

$$\|\Delta_x\|_2 = \mathcal{O}(\epsilon_x n^{1/2} \log^{1/2}(n/p)), \quad \|\Delta_y\|_2 = \mathcal{O}(\epsilon_y n^{1/2} \log^{1/2}(n/p)).$$

Similar to the proof of Lemma 5, by (Mhammedi et al., 2020, Section B.2.11), with probability at least $1 - p$,

$$\begin{aligned}\|X^\dagger E\|_2 &= \mathcal{O}(1) \sigma_{\min}^+(X)^{-1} \|\Sigma_e^{1/2}\|_2 (d_1 + d_2 + \log(1/p))^{1/2} \\ &= \mathcal{O}(\theta^{-1} n^{-1/2} (d_1 + d_2 + \log(1/p))^{1/2}).\end{aligned}$$

Combining the bounds above, we obtain

$$\begin{aligned}\|\hat{A}\|_2 &\leq \|A^*\|_2 + \mathcal{O}(\theta^{-1} n^{-1/2} (d_1 + d_2 + \log(1/p))^{1/2}) \\ &\quad + \mathcal{O}((\|A^*\|_2 \epsilon_x + \epsilon_y) n^{1/2} \log^{1/2}(n/p) \theta^{-1} n^{-1/2}) \\ &= \mathcal{O}(\|A^*\|_2 (1 + \theta^{-1} \epsilon_x \log^{1/2}(n/p)) + \theta^{-1} \epsilon_y \log^{1/2}(n/p)) \\ &\quad + \mathcal{O}(\theta^{-1} n^{-1/2} (d_1 + d_2 + \log(1/p))^{1/2}).\end{aligned}$$

Hence, as long as $\theta = \Omega(\epsilon_x)$, $\theta = \Omega(\epsilon_y)$ for absolute constants and θ has at least $n^{-1/4}$ dependence on n , $\|\hat{A}\|_2$ is bounded by $c\|A^*\|_2$ for an absolute constant $c > 0$. \square

In Lemma 5, if Σ_* has full rank and $\sigma_{\min}(\Sigma_*) \geq \beta > 0$ and $\sigma_{\min}(\sum_{i=1}^n x^{(i)}(x^{(i)})^\top) = \Omega(\beta^2 n)$, then

$$\|\hat{A} - A^*\|_2 = \mathcal{O}(\beta^{-1} ((\beta^{-1} \epsilon_x + \epsilon_y) \log^{1/2}(n/p) + n^{-1/2} (d_1 + d_2 + \log(1/p))^{1/2})).$$

The following lemma shows that we can strengthen the result by removing the β^{-1} factor before ϵ_x .

Lemma 6 (Noisy linear regression). Define random variable $y^* = A^*x^* + e$, where $x^* \sim \mathcal{N}(0, \Sigma_*)$ and $e \sim \mathcal{N}(0, \Sigma_e)$ are d_1 and d_2 dimensional random vectors. Assume that $\|A^*\|_2$, $\|\Sigma_*\|_2$ and $\|\Sigma_e\|_2$ are $\mathcal{O}(1)$, and $\sigma_{\min}(\Sigma_*^{1/2}) \geq \beta > 0$. Define $x := x^* + \delta_x$ and $y := y^* + \delta_y$ where the noise vectors δ_x and δ_y can be correlated with x^* and y^* , and their ℓ_2 norms are sub-Gaussian with $\mathbb{E}[\|\delta^x\|] \leq \epsilon_x$, $\|\|\delta^x\|\|_{\psi_2} \leq \epsilon_x$ and $\mathbb{E}[\|\delta^y\|] \leq \epsilon_y$, $\|\|\delta^y\|\|_{\psi_2} \leq \epsilon_y$. Suppose we get n independent observations $x^{(i)}, y^{(i)}$ of x and y . Assume that $\sigma_{\min}(\sum_{i=1}^n x^{(i)}(x^{(i)})^\top) = \Omega(\beta^2 n)$. Consider the minimum Frobenius norm solution

$$\hat{A} \in \underset{A}{\operatorname{argmin}} \sum_{i=1}^n (y^{(i)} - Ax^{(i)})^2.$$

Then, there exists an absolute constant $c > 0$, such that if $n \geq c(d_1 + d_2 + \log(1/p))$, with probability at least $1 - p$,

$$\|\hat{A} - A^*\|_2 = \mathcal{O}(\beta^{-1}((\epsilon_x + \epsilon_y) \log^{1/2}(n/p) + n^{-1/2}(d_1 + d_2 + \log(1/p))^{1/2})).$$

Proof. Following the proof of Claim 1, we have

$$\|\hat{A} - A^*\|_2 \leq \|A^*\| \|X^\top \Delta_x\|_2 + \|X^\top E\|_2 + \|X^\top \Delta_y\|_2.$$

Combining with the bounds on $\|X^\top \Delta_x\|_2$, $\|X^\top \Delta_y\|_2$ and $\|X^\top E\|_2$ concludes the proof. \square

5 Concluding remarks

We examined cost-driven state representation learning methods in time-varying LQG control. With a finite-sample analysis, we showed that a direct, cost-driven state representation learning algorithm effectively solves LQG. In the analysis, we revealed the importance of using multi-step cumulative costs as the supervision signal, and a separation of the convergence rates before and after step ℓ , the controllability index, due to early-stage insufficient excitement of the system. A major limitation of our method is the use of history-based state representation functions; recovering the recursive Kalman filter would be ideal.

This work has opened up many opportunities for future research. An immediate question is how our cost-driven state representation learning approach performs in the infinite-horizon LTI setting. Moreover, one may wonder about the extent to which cost-driven state representation learning generalizes to nonlinear observations or systems. Investigating the connection with reduced-order control is also an interesting question, which may reveal the unique advantage of cost-driven state representation learning. Finally, one argument for favoring latent-model-based over model-free methods is their ability to generalize across different tasks; cost-driven state representation learning may offer a perspective to formalize this intuition.

Acknowledgement

YT, SS acknowledge partial support from the NSF BIGDATA grant (number 1741341). KZ's work was mainly done while at MIT, and acknowledges partial support from Simons-Berkeley

Research Fellowship. The authors also thank Xiang Fu, Horia Mania, and Alexandre Megretski for helpful discussions.

References

- Karl J Åström. *Introduction to Stochastic Control Theory*. Courier Corporation, 2012.
- Dimitri Bertsekas. *Dynamic Programming and Optimal Control: Volume I*, volume 1. Athena Scientific, 2012.
- Rajendra Bhatia and Fuad Kittaneh. On the singular values of a product of operators. *SIAM Journal on Matrix Analysis and Applications*, 11(2):272–277, 1990.
- Stephen P Boyd and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- Fei Deng, Ingook Jang, and Sungjin Ahn. Dreamerpro: Reconstruction-free model-based reinforcement learning with prototypical representations. *arXiv preprint arXiv:2110.14565*, 2021.
- Maryam Fazel, Rong Ge, Sham Kakade, and Mehran Mesbahi. Global convergence of policy gradient methods for the linear quadratic regulator. In *International Conference on Machine Learning*, pages 1467–1476. PMLR, 2018.
- Abraham Frandsen, Rong Ge, and Holden Lee. Extracting latent state representations with linear dynamics from rich observations. In *International Conference on Machine Learning*, pages 6705–6725. PMLR, 2022.
- Xiang Fu, Ge Yang, Pulkit Agrawal, and Tommi Jaakkola. Learning task informed abstractions. In *International Conference on Machine Learning*, pages 3480–3491. PMLR, 2021.
- David Ha and Jürgen Schmidhuber. World models. *arXiv preprint arXiv:1803.10122*, 2018.
- Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. *arXiv preprint arXiv:1912.01603*, 2019a.
- Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James Davidson. Learning latent dynamics for planning from pixels. In *International conference on machine learning*, pages 2555–2565. PMLR, 2019b.
- Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari with discrete world models. *arXiv preprint arXiv:2010.02193*, 2020.
- Botao Hao, Yasin Abbasi Yadkori, Zheng Wen, and Guang Cheng. Bootstrapping upper confidence bound. *Advances in Neural Information Processing Systems*, 32, 2019.
- Ali Jadbabaie, Horia Mania, Devavrat Shah, and Suvrit Sra. Time varying regression with hidden linear dynamics. *arXiv preprint arXiv:2112.14862*, 2021.

- Thomas Kailath. *Linear Systems*, volume 156. Prentice-Hall Englewood Cliffs, NJ, 1980.
- Nicholas Komaroff. On bounds for the solution of the Riccati equation for discrete-time control systems. In *Control and Dynamic Systems*, volume 78, pages 275–311. Elsevier, 1996.
- Sahin Lale, Kamyar Azizzadenesheli, Babak Hassibi, and Anima Anandkumar. Logarithmic regret bound in partially observable linear dynamical systems. *Advances in Neural Information Processing Systems*, 33:20876–20888, 2020.
- Sahin Lale, Kamyar Azizzadenesheli, Babak Hassibi, and Anima Anandkumar. Adaptive control and regret minimization in linear quadratic Gaussian (LQG) setting. In *2021 American Control Conference (ACC)*, pages 2517–2522. IEEE, 2021.
- Alex Lamb, Riashat Islam, Yonathan Efroni, Aniket Didolkar, Dipendra Misra, Dylan Foster, Lekan Molu, Rajan Chari, Akshay Krishnamurthy, and John Langford. Guaranteed discovery of controllable latent states with multi-step inverse models. *arXiv preprint arXiv:2207.08229*, 2022.
- Lennart Ljung. System identification. In *Signal Analysis and Prediction*, pages 163–173. Springer, 1998.
- Horia Mania, Stephen Tu, and Benjamin Recht. Certainty equivalence is efficient for linear quadratic control. *Advances in Neural Information Processing Systems*, 32, 2019.
- Zakaria Mhammedi, Dylan J Foster, Max Simchowitz, Dipendra Misra, Wen Sun, Akshay Krishnamurthy, Alexander Rakhlin, and John Langford. Learning the linear quadratic regulator from nonlinear observations. *Advances in Neural Information Processing Systems*, 33:14532–14543, 2020.
- Edgar Minasyan, Paula Gradu, Max Simchowitz, and Elad Hazan. Online control of unknown time-varying dynamical systems. *Advances in Neural Information Processing Systems*, 34, 2021.
- Masashi Okada and Tadahiro Taniguchi. Dreaming: Model-based reinforcement learning by latent imagination without reconstruction. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 4209–4215. IEEE, 2021.
- Samet Oymak and Necmiye Ozay. Non-asymptotic identification of LTI systems from a single trajectory. In *2019 American control conference (ACC)*, pages 5655–5661. IEEE, 2019.
- Kathrin Schacke. On the Kronecker product. *Master’s Thesis, University of Waterloo*, 2004.
- Peter Hans Schoenemann. *A solution of the orthogonal Procrustes problem with applications to orthogonal and oblique rotation*. University of Illinois at Urbana-Champaign, 1964.
- Peter H Schönemann. A generalized solution of the orthogonal procrustes problem. *Psychometrika*, 31(1):1–10, 1966.
- Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. Mastering Atari, Go, Chess and Shogi by planning with a learned model. *Nature*, 588(7839):604–609, 2020.

- Max Simchowitz, Karan Singh, and Elad Hazan. Improper learning for non-stochastic control. In *Conference on Learning Theory*, pages 3320–3436. PMLR, 2020.
- Jayakumar Subramanian, Amit Sinha, Raihan Seraj, and Aditya Mahajan. Approximate information state for approximate planning and reinforcement learning in partially observed systems. *arXiv preprint arXiv:2010.08843*, 2020.
- Wen Sun, Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, and John Langford. Model-based RL in contextual decision processes: PAC bounds and exponential improvements over model-free approaches. In *Conference on Learning Theory*, pages 2898–2933. PMLR, 2019.
- Richard S Sutton and Andrew G Barto. *Reinforcement Learning: An Introduction*. MIT Press, 2018.
- Stephen Tu and Benjamin Recht. The gap between model-based and model-free methods on the linear quadratic regulator: An asymptotic viewpoint. In *Conference on Learning Theory*, pages 3036–3083. PMLR, 2019.
- Stephen Tu, Ross Boczar, Max Simchowitz, Mahdi Soltanolkotabi, and Ben Recht. Low-rank solutions of linear matrix equations via procrustes flow. In *International Conference on Machine Learning*, pages 964–973. PMLR, 2016.
- Jack Umenberger, Max Simchowitz, Juan Carlos Perdomo, Kaiqing Zhang, and Russ Tedrake. Globally convergent policy search for output estimation. In *Advances in Neural Information Processing Systems*, 2022.
- Roman Vershynin. *High-dimensional Probability: An Introduction with Applications in Data Science*, volume 47. Cambridge University press, 2018.
- Martin J Wainwright. *High-dimensional Statistics: A Non-asymptotic Viewpoint*, volume 48. Cambridge University Press, 2019.
- Tongzhou Wang, Simon S Du, Antonio Torralba, Phillip Isola, Amy Zhang, and Yuandong Tian. Denoised MDPs: Learning world models better than the world itself. *arXiv preprint arXiv:2206.15477*, 2022.
- Lujie Yang, Kaiqing Zhang, Alexandre Amice, Yunzhu Li, and Russ Tedrake. Discrete approximate information states in partially observable environments. In *2022 American Control Conference (ACC)*, pages 1406–1413. IEEE, 2022.
- Amy Zhang, Rowan McAllister, Roberto Calandra, Yarin Gal, and Sergey Levine. Learning invariant representations for reinforcement learning without reconstruction. *arXiv preprint arXiv:2006.10742*, 2020.
- Huiming Zhang and Haoyu Wei. Sharper sub-weibull concentrations. *Mathematics*, 10(13):2252, 2022.
- Kaiqing Zhang, Xiangyuan Zhang, Bin Hu, and Tamer Basar. Derivative-free policy optimization for linear risk-sensitive and robust control design: Implicit regularization and sample complexity. *Advances in Neural Information Processing Systems*, 34:2949–2964, 2021.

Marvin Zhang, Sharad Vikram, Laura Smith, Pieter Abbeel, Matthew Johnson, and Sergey Levine. Solar: Deep structured representations for model-based reinforcement learning. In *International Conference on Machine Learning*, pages 7444–7453. PMLR, 2019.

Qinghua Zhang and Lianguan Zhang. Boundedness of the Kalman filter revisited. *IFAC-PapersOnLine*, 54(7):334–338, 2021.

Yang Zheng, Luca Furieri, Maryam Kamgarpour, and Na Li. Sample complexity of linear quadratic Gaussian (LQG) control for output feedback systems. In *Learning for Dynamics and Control*, pages 559–570. PMLR, 2021.

Bin Zhou and Tianrui Zhao. On asymptotic stability of discrete-time linear time-varying systems. *IEEE Transactions on Automatic Control*, 62(8):4274–4281, 2017.

A Proposition on multi-step cumulative costs

The following proposition does not appear in the main body, but is important for analyzing CoREL (Algorithm 2) in Section 4.1.

Proposition 3. *Let $(z_t^{*'})_{t=0}^T$ be the state estimates by the Kalman filter under the normalized parameterization. If we apply $u_t \sim \mathcal{N}(0, \sigma_u^2 I)$ for all $0 \leq t \leq T-1$, then for $0 \leq t \leq T$,*

$$\bar{c}_t := c_t + c_{t+1} + \dots + c_{t+k-1} = \|z_t^{*'}\|^2 + \sum_{\tau=t}^{t+k-1} \|u_\tau\|_{R_\tau^*}^2 + b'_t + e'_t,$$

where $k = 1$ for $0 \leq t \leq \ell - 1$ and $k = m \wedge (T - t + 1)$ for $\ell \leq t \leq T$, $b'_t = \mathcal{O}(k)$, and e'_t is a zero-mean subexponential random variable with $\|e'_t\|_{\psi_1} = \mathcal{O}(kd_x^{1/2})$.

Proof. By Proposition 1, $z_{t+1}^{*'} = A_t^{*'} z_t^{*'} + B_t^{*'} u_t + L_{t+1}^{*'} i'_{t+1}$, where $L_{t+1}^{*'}, i'_{t+1}$ are the Kalman gain and the innovation under the normalized parameterization, respectively. Under Assumptions 1 and 4, $(i'_t)_{t=0}^T$ are Gaussian random vectors whose covariances have $\mathcal{O}(1)$ operator norms, and $(L_t^{*'})_{t=0}^T$ have $\mathcal{O}(1)$ operator norms (Zhang and Zhang, 2021). Hence, The covariance of $L_{t+1}^{*'} i'_{t+1}$ has $\mathcal{O}(1)$ operator norm. Since $u_t \sim \mathcal{N}(0, \sigma_u^2 I)$, $j_t := B_t^{*'} u_t + i'_t$ can be viewed as a Gaussian noise vector whose covariance has $\mathcal{O}(1)$ operator norm. By Proposition 1,

$$c_t = \|z_t^{*'}\|_{Q_t^{*'}}^2 + \|u_t\|_{R_t^*}^2 + b'_t + e'_t,$$

where $e'_t := \gamma'_t + \eta'_t$ is subexponential with $\|e'_t\|_{\psi_1} = \mathcal{O}(d_x^{1/2})$. Let $\Phi'_{t,t_0} = A_{t-1}^{*'} A_{t-2}^{*'} \dots A_{t_0}^{*'}$ for $t > t_0$ and $\Phi'_{t,t} = I$. Then, for $\tau \geq t$,

$$z_\tau^{*'} = \Phi'_{\tau,t} z_t^{*'} + \sum_{s=t}^{\tau-1} \Phi'_{\tau,s} j_s := \Phi'_{\tau,t} z_t^{*'} + j'_{\tau,t},$$

where $j'_{t,t} = 0$ and for $\tau > t$, $j'_{\tau,t}$ is a Gaussian random vector with bounded covariance due to uniform exponential stability (Assumption 1). Therefore,

$$\begin{aligned}
\bar{c}_t &= \sum_{\tau=t}^{t+k-1} c_\tau \\
&= \sum_{\tau=t}^{t+k-1} (\|\Phi'_{\tau,t} z_t^{*'} + j'_{\tau,t}\|_{Q_\tau^{*'}}^2 + \|u_\tau\|_{R_\tau^*}^2 + b'_\tau + e'_\tau) \\
&= (z_t^{*'})^\top \left(\sum_{\tau=t}^{t+k-1} (\Phi'_{\tau,t})^\top Q_t^{*'} \Phi'_{\tau,t} \right) z_t^{*'} + \sum_{\tau=t}^{t+k-1} \|u_\tau\|_{R_\tau^*}^2 \\
&\quad + \sum_{\tau=t}^{t+k-1} (\|j'_{\tau,t}\|_{Q_\tau^{*'}}^2 + (j'_{\tau,t})^\top Q_\tau^{*'} \Phi'_{\tau,t} z_t^{*'} + b'_\tau + e'_\tau) \\
&= \|z_t^{*'}\|^2 + \sum_{\tau=t}^{t+k-1} \|u_\tau\|_{R_\tau^*}^2 + \bar{b}_t + \bar{e}_t,
\end{aligned}$$

where $\sum_{\tau=t}^{t+k-1} (\Phi'_{\tau,t})^\top Q_t^{*'} \Phi'_{\tau,t} = I$ is due to the normalized parameterization, $\bar{b}_t := \sum_{\tau=t}^{t+k-1} (b_\tau + \mathbb{E}[\|j'_\tau\|_{Q_\tau^{*'}}^2]) = \mathcal{O}(k)$, and

$$\bar{e}_t := \sum_{\tau=t}^{t+k-1} (\|j'_\tau\|_{Q_\tau^{*'}}^2 - \mathbb{E}[\|j'_\tau\|_{Q_\tau^{*'}}^2] + (j'_\tau)^\top Q_\tau^{*'} \Phi'_{\tau,t} z_t^{*'} + e'_\tau)$$

has zero mean and is subexponential with $\|\bar{e}_t\|_{\psi_1} = \mathcal{O}(kd_x^{1/2})$. \square

B Auxiliary results

Lemma 7. Let $x \sim \mathcal{N}(0, \Sigma_x)$ and $y \sim \mathcal{N}(0, \Sigma_y)$ be d -dimensional Gaussian random vectors. Let Q be a $d \times d$ positive semidefinite matrix. Then, there exists an absolute constant $c > 0$ such that

$$\|\langle x, y \rangle_Q\|_{\psi_1} \leq c\sqrt{d} \|Q\|_2 \sqrt{\|\Sigma_x\|_2 \|\Sigma_y\|_2}.$$

Proof. Since $|\langle x, y \rangle_Q| = |x^\top Q y| \leq \|x\| \|Q\|_2 \|y\|$,

$$\|\langle x, y \rangle_Q\|_{\psi_1} = \|\|\langle x, y \rangle_Q\|\|_{\psi_1} \leq \|\|x\| \|Q\|_2 \|y\|\|_{\psi_1} = \|Q\|_2 \cdot \|\|x\| \|y\|\|_{\psi_1}.$$

For $x \sim \mathcal{N}(0, \Sigma_x)$, we know that $\|x\|$ is sub-Gaussian. Actually, by writing $x = \Sigma_x^{1/2} g$ for $g \sim \mathcal{N}(0, I)$, we have

$$\|\|x\|\|_{\psi_2} = \|\|\Sigma_x^{1/2} g\|\|_{\psi_2} \leq \|\|\Sigma_x^{1/2}\|_2 \|g\|\|_{\psi_2} = \|\Sigma_x^{1/2}\|_2 \|\|g\|\|_{\psi_2}.$$

The distribution of $\|g\|_2$ is known as χ distribution, and we know that $\|\|g\|\|_{\psi_2} = c'd^{1/4}$ for an absolute constant $c' > 0$ (see, e.g., (Wainwright, 2019)). Hence, $\|\|x\|\|_{\psi_2} \leq c'd^{1/4} \|\Sigma_x\|_2^{1/2}$. Similarly, $\|\|y\|\|_{\psi_2} \leq c'd^{1/4} \|\Sigma_y\|_2^{1/2}$. Since $\|\|x\| \|y\|\|_{\psi_1} \leq \|\|x\|\|_{\psi_2} \|\|y\|\|_{\psi_2}$ (see, e.g., (Vershynin, 2018, Lemma 2.7.7)), we have

$$\|\langle x, y \rangle_Q\|_{\psi_1} \leq (c')^2 \sqrt{d} \|Q\|_2 \sqrt{\|\Sigma_x\|_2 \|\Sigma_y\|_2}.$$

Taking $c = (c')^2$ concludes the proof. \square

Lemma 8. Let x, y be random vectors of dimensions d_x, d_y , respectively, defined on the same probability space. Then, $\|\text{Cov}([x; y])\|_2 \leq \|\text{Cov}(x)\|_2 + \|\text{Cov}(y)\|_2$.

Proof. Let $\text{Cov}([x; y]) = DD^\top$ be a factorization of the positive semidefinite matrix $\text{Cov}([x; y])$, where $D \in \mathbb{R}^{(d_x+d_y) \times (d_x+d_y)}$. Let D_x and D_y be the matrices consisting of the first d_x rows and the last d_y rows of D , respectively. Then,

$$\begin{aligned} \text{Cov}([x; y]) &= DD^\top = [D_x; D_y][D_x^\top, D_y^\top] \\ &= \begin{bmatrix} D_x D_x^\top & D_x D_y^\top \\ D_y D_x^\top & D_y D_y^\top \end{bmatrix}. \end{aligned}$$

Hence, $\text{Cov}(x) = D_x D_x^\top$ and $\text{Cov}(y) = D_y D_y^\top$. The proof is completed by noticing that

$$\begin{aligned} \|\text{Cov}([x; y])\|_2 &= \|D^\top D\|_2 = \|[D_x^\top, D_y^\top][D_x; D_y]\|_2 \\ &= \|D_x^\top D_x + D_y^\top D_y\|_2 \\ &\leq \|D_x^\top D_x\|_2 + \|D_y^\top D_y\|_2 \\ &= \|\text{Cov}(x)\|_2 + \|\text{Cov}(y)\|_2. \end{aligned}$$

□

C Certainty equivalent linear quadratic control

As shown in Lemma 5, if the input of linear regression does not have full-rank covariance, then the parameters can only be identified in certain directions. The following lemma studies the performance of the certainty equivalent optimal controller in this case.

Lemma 9 (Rank deficient linear quadratic control). *Consider the finite-horizon time-varying linear dynamical system $x_{t+1} = A_t^* x_t + B_t^* u_t + w_t$ with unknown $(A_t^*, B_t^*)_{t=0}^{T-1}$.*

- Assume that $(A_t^*)_{t=0}^{T-1}$ is uniformly exponentially stable (Assumption 1).
- Assume that $x_0 \sim \mathcal{N}(0, \Sigma_0)$, $w_t \sim \mathcal{N}(0, \Sigma_{w_t})$ for all $0 \leq t \leq T-1$, are independent, where Σ_0 and $(\Sigma_{w_t})_{t=0}^{T-1}$ do not necessarily have full rank.
- Instead of x_t , we observe $x'_t = x_t + \delta_{x_t}$, where the Gaussian noise vector δ_{x_t} can be correlated with x_t and satisfy $\|\text{Cov}(\delta_{x_t})^{1/2}\| \leq \varepsilon$ for all $0 \leq t \leq T$.
- Let $\Sigma_t := \text{Cov}(x_t)$ and $\Sigma'_t := \text{Cov}(x'_t)$ under control $u_t \sim \mathcal{N}(0, \sigma_u^2 I)$ for all $0 \leq t \leq T-1$. Assume $\sigma_{\min}^+(\Sigma_t)^{1/2} \geq \beta > 0$ and $\sigma_{\min}^+(\Sigma'_t)^{1/2} \geq \theta > 0$ for all $0 \leq t \leq T$, and that $\theta = \Omega(\varepsilon)$ with at least $n^{-1/4}$ dependence on n .
- Assume $(Q_t^*)_{t=0}^T, (\hat{Q}_t)_{t=0}^T$ are positive definite with $\mathcal{O}(1)$ operator norms, and $\|\hat{Q}_t - Q_t^*\|_2 \leq \varepsilon$ for all $0 \leq t \leq T$.

Collect n trajectories using $u_t \sim \mathcal{N}(0, \sigma_u^2 I)$ for some $\sigma_u > 0$ and $0 \leq t \leq T-1$. Identify $(\hat{A}_t, \hat{B}_t)_{t=0}^{T-1}$ by SysID (Algorithm 3), which uses ordinary least squares and takes a minimum Frobenius norm solution.

Let $(\hat{K}_t)_{t=0}^{T-1}$ be the optimal controller in system $((\hat{A}_t, \hat{B}_t, R_t^*)_{t=0}^{T-1}, (\hat{Q}_t)_{t=0}^T)$. Then, there exists an absolute constant $a > 0$, such that if $n \geq a(d_x + d_u + \log(1/p))$, with probability at least $1 - p$, $(\hat{K}_t)_{t=0}^{T-1}$ is ϵ -optimal for solving the LQR problem $((A_t^*, B_t^*, R_t^*)_{t=0}^{T-1}, (Q_t^*)_{t=0}^T)$ for $\epsilon =$

$$\mathcal{O}(a^T((1 + \beta^{-1})d_x\epsilon + d_x(d_x + d_u + \log(1/p))^{1/2}n^{-1/2})),$$

where dimension-free constant $a > 0$ depends on the system parameters.

Proof. First of all, we establish that $(K_t^*)_{t=0}^{T-1}, (\hat{K}_t)_{t=0}^{T-1}$ are bounded. Since $(Q_t^*)_{t=0}^T, (\hat{Q}_t)_{t=0}^T$ are positive definite, the system is fully cost observable. By (Zhang and Zhang, 2021), the optimal feedback gains $(K_t^*)_{t=0}^{T-1}$ have operator norms bounded by dimension-free constants depending on the system parameters. Note that (Zhang and Zhang, 2021) studies the boundedness of the Kalman filter, which is dual to our optimal control problem, but their results carry over. One can also see the boundedness from the compact formulation (Zhang et al., 2021), as described in the proof of Lemma 10, using known bounds for RDE solutions (Komaroff, 1996). The same argument applies to $(\hat{K}_t)_{t=0}^{T-1}$ if we can show $(\hat{A}_t, \hat{B}_t)_{t=0}^{T-1}$ have $\mathcal{O}(1)$ operator norms. SysID (Algorithm 3) identifies $[\hat{A}_t, \hat{B}_t]$ by ordinary least squares for all $0 \leq t \leq T - 1$. By Claim 1, with a union bound for all $0 \leq t \leq T - 1$, as long as $n \geq 16(d_x + d_u + \log(T/p))$, with probability at least $1 - 4p$, $\|[\hat{A}_t, \hat{B}_t]\|_2 = \mathcal{O}(\|[A_t^*, B_t^*]\|_2) = \mathcal{O}(1)$. Hence, with the same probability, $(\hat{K}_t)_{t=0}^{T-1}$ have operator norms bounded by dimension-free constants depending on the system parameters. In the following, we assume $\|K_t^*\|_2, \|\hat{K}_t\|_2 \leq \kappa$ for all $0 \leq t \leq T - 1$, where κ is a dimension-free constant that depends on system parameters.

Since Σ_t can be rank deficient, we cannot guarantee $\|\hat{A}_t - A_t^*\|_2$ is small. Instead, for all $0 \leq t \leq T - 1$, by Lemma 5,

$$([\hat{A}_t, \hat{B}_t] - [A_t^*, B_t^*])\text{diag}(\Sigma_t^{1/2}, \sigma_u I) = \mathcal{O}(\delta),$$

where $\delta := (1 + \beta^{-1})\epsilon + (d_x + d_u + \log(1/p))^{1/2}n^{-1/2}$, and $\epsilon_x, \epsilon_y, d_1, d_2$ in Lemma 5 correspond to $\epsilon, \epsilon, d_x + d_u, d_x$ here, respectively. This implies that $\|(\hat{A}_t - A_t^*)(\Sigma_t)^{1/2}\|_2 = \mathcal{O}(\delta)$ and $\|\hat{B}_t - B_t^*\|_2 = \mathcal{O}(\delta)$ for all $0 \leq t \leq T - 1$.

Let Ξ_t denote the covariance of x_t under state feedback controller $(K_t)_{t=0}^{T-1}$, where $\|K_t\| \leq \kappa$ for all $0 \leq t \leq T - 1$. We claim that there exists $b_t > 0$, such that $\Xi_t \preceq b_t \Sigma_t$. We prove it by induction. At step 0, clearly, $\Xi_0 = \Sigma_0$. Suppose $\Xi_t \preceq b_t \Sigma_t$. For step $t + 1$,

$$\begin{aligned} \Sigma_{t+1} &= A_t^* \Sigma_t (A_t^*)^\top + \sigma_u^2 B_t^* (B_t^*)^\top + \Sigma_{w_t}, \\ \Xi_{t+1} &= (A_t^* + B_t^* K_t) \Xi_t (A_t^* + B_t^* K_t)^\top + \Sigma_{w_t} \\ &\preceq 2A_t^* \Xi_t (A_t^*)^\top + 2B_t^* K_t \Xi_t K_t^\top (B_t^*)^\top + \Sigma_{w_t}. \end{aligned}$$

Hence, for $b \geq 1$,

$$b\Sigma_{t+1} - \Xi_{t+1} \succeq (b - 2b_t)A_t^* \Sigma_t (A_t^*)^\top + b\sigma_u^2 B_t^* (B_t^*)^\top - 2B_t^* K_t \Xi_t K_t^\top (B_t^*)^\top.$$

To ensure $b_{t+1}\Sigma_{t+1} \succeq \Xi_{t+1}$, it suffices to take $b_{t+1} = \max\{2, 2\sigma_u^{-2}\kappa^2\|\Sigma_t\|_2\}b_t$. Hence, $b_t \leq a_0^t$, where constant $a_0 > 0$ is dimension-free and depends on system parameters. Note that if $(K_t)_{t=0}^{T-1}$ stabilizes $(A_t^*, B_t^*)_{t=0}^{T-1}$, then $\|\Xi_t\|_2 = \mathcal{O}(1)$, and the bound on b_t can be improved to

$\sigma_{\min}^+(\Sigma_t)^{-1} \|\Xi_t\|_2$; so the exponential dependence on t results from not knowing whether $(K_t)_{t=0}^{T-1}$ stabilizes the system. By the definition of the operator norm, we have $\|(\hat{A}_t - A_t^*)\Xi_t^{1/2}\|_2 \leq \|(\hat{A}_t - A_t^*)(b_t \Sigma_t)^{1/2}\|_2 = \mathcal{O}(b_t^{1/2} \delta)$ for all $0 \leq t \leq T-1$.

Let $\hat{\Xi}_t$ denote the covariance of x_t under state feedback controller $(K_t)_{t=0}^{T-1}$ in the system $(\hat{A}_t, \hat{B}_t)_{t=0}^{T-1}$. Then, by definition,

$$\begin{aligned} & \|\Xi_{t+1} - \hat{\Xi}_{t+1}\|_2 \\ &= \left\| (A_t^* + B_t^* K_t) \Xi_t (A_t^* + B_t^* K_t)^\top + \Sigma_{w_t} - \left((\hat{A}_t + \hat{B}_t K_t) \hat{\Xi}_t (\hat{A}_t + \hat{B}_t K_t)^\top + \Sigma_{w_t} \right) \right\|_2 \\ &= \left\| (A_t^* + B_t^* K_t) \Xi_t (A_t^* + B_t^* K_t)^\top - (\hat{A}_t + \hat{B}_t K_t) \Xi_t (\hat{A}_t + \hat{B}_t K_t)^\top \right. \\ & \quad \left. + (\hat{A}_t + \hat{B}_t K_t) (\Xi_t - \hat{\Xi}_t) (\hat{A}_t + \hat{B}_t K_t)^\top \right\|_2 \\ &= \mathcal{O} \left(\kappa b_t^{1/2} \left\| (A_t^* + B_t^* K_t) \Xi_t^{1/2} - (\hat{A}_t + \hat{B}_t K_t) \Xi_t^{1/2} \right\|_2 \right) + \mathcal{O}(1 + \kappa^2) \|\hat{\Xi}_t - \Xi_t\|_2, \end{aligned}$$

where the operator norms of A_t^*, B_t^* are $\mathcal{O}(1)$ and the operator norms of \hat{A}_t, \hat{B}_t are $\mathcal{O}(1)$ with probability at least $1 - 4p$. Since

$$\begin{aligned} \left\| (A_t^* + B_t^* K_t) \Xi_t^{1/2} - (\hat{A}_t + \hat{B}_t K_t) \Xi_t^{1/2} \right\|_2 &= \left\| (A_t^* - \hat{A}_t) \Xi_t^{1/2} + (B_t^* - \hat{B}_t) K_t \Xi_t^{1/2} \right\|_2 \\ &= \mathcal{O}(\delta \kappa b_t^{1/2}), \end{aligned}$$

we have $\|\hat{\Xi}_{t+1} - \Xi_{t+1}\|_2 = \mathcal{O}(\delta \kappa^2 b_t + (1 + \kappa^2) \|\hat{\Xi}_t - \Xi_t\|_2)$. Combining with $\hat{\Xi}_0 = \Xi_0 = \Sigma_0$ gives

$$\|\hat{\Xi}_t - \Xi_t\|_2 = \mathcal{O}(b_t(a_1 + a_1 \kappa^2)^t \delta) = \mathcal{O}((a_2 + a_2 \kappa^2)^t \delta)$$

for some dimension-free constants $a_1, a_2 > 0$ that depend on the system parameters.

Let $J((K_t)_{t=0}^{T-1})$ denote the expected cumulative cost under state feedback controller $(K_t)_{t=0}^{T-1}$ in system $((A_t^*, B_t^*, R_t^*)_{t=0}^{T-1}, (Q_t^*)_{t=0}^T)$ and $\hat{J}((K_t)_{t=0}^{T-1})$ the corresponding expected cumulative cost in system $((\hat{A}_t, \hat{B}_t, R_t^*)_{t=0}^{T-1}, (\hat{Q}_t)_{t=0}^T)$. Notice that

$$\begin{aligned} J((K_t)_{t=0}^{T-1}) &= \mathbb{E} \left[\sum_{t=0}^T c_t \right] = \mathbb{E} \left[\sum_{t=0}^T x_t^\top (Q_t^*)' x_t \right] \\ &= \mathbb{E} \left[\sum_{t=0}^T \langle (Q_t^*)', x_t x_t^\top \rangle_F \right] \\ &= \sum_{t=0}^T \langle (Q_t^*)', \Xi_t \rangle_F, \end{aligned}$$

where $(Q_t^*)' := (Q_t^* + K_t^\top R_t^* K_t)$ for $0 \leq t \leq T-1$ and $(Q_T^*)' := Q_T^*$. Similarly, $\hat{J}((K_t)_{t=0}^{T-1}) = \sum_{t=0}^T \langle \hat{Q}_t', \hat{\Xi}_t \rangle_F$, where $\hat{Q}_t' = \hat{Q}_t + K_t^\top R_t^* K_t$ for $0 \leq t \leq T-1$ and $\hat{Q}_T' = \hat{Q}_T$. Hence,

$$\begin{aligned} |J((K_t)_{t=0}^{T-1}) - \hat{J}((K_t)_{t=0}^{T-1})| &= \sum_{t=0}^T \langle (Q_t^*)' - \hat{Q}_t', \Xi_t \rangle_F + \sum_{t=0}^T \langle \hat{Q}_t', \Xi_t - \hat{\Xi}_t \rangle_F \\ &= \varepsilon d_x \sum_{t=0}^T \mathcal{O}(b_t) + \kappa^2 d_x \sum_{t=0}^T \mathcal{O}((a_2 + a_2 \kappa^2)^t \delta) \\ &= \varepsilon d_x \sum_{t=0}^T \mathcal{O}(a_0^t) + \delta \kappa^2 d_x \sum_{t=0}^T \mathcal{O}((a_2 + a_2 \kappa^2)^t) \\ &= \mathcal{O}(\delta d_x a^T), \end{aligned}$$

where $\|\hat{Q}_t'\|_2 = \mathcal{O}(\kappa^2)$, ε is absorbed into $\delta = (1 + \beta^{-1})\varepsilon + (d_x + d_u + \log(1/p))^{1/2}n^{-1/2}$, and dimension-free constant $a > 0$ depends on the system parameters.

Finally, let $(K_t^*)_{t=0}^{T-1}$ be the optimal feedback gains in system $((A_t^*, B_t^*, R_t^*)_{t=0}^{T-1}, (Q_t^*)_{t=0}^T)$. By the union bound, with probability at least $1 - 8p$,

$$|J((\hat{K}_t)_{t=0}^{T-1}) - \hat{J}((\hat{K}_t)_{t=0}^{T-1})| = \mathcal{O}(\delta d_x a^T), \quad |J((K_t^*)_{t=0}^{T-1}) - \hat{J}((K_t^*)_{t=0}^{T-1})| = \mathcal{O}(\delta d_x a^T).$$

Therefore,

$$\begin{aligned} & J((\hat{K}_t)_{t=0}^{T-1}) - J((K_t^*)_{t=0}^{T-1}) \\ &= J((\hat{K}_t)_{t=0}^{T-1}) - \hat{J}((\hat{K}_t)_{t=0}^{T-1}) + \hat{J}((\hat{K}_t)_{t=0}^{T-1}) - \hat{J}((K_t^*)_{t=0}^{T-1}) + \hat{J}((K_t^*)_{t=0}^{T-1}) - J((K_t^*)_{t=0}^{T-1}) \\ &= \mathcal{O}(d_x a^T \delta) \\ &= \mathcal{O}(a^T((1 + \beta^{-1})d_x \varepsilon + d_x(d_x + d_u + \log(1/p))^{1/2}n^{-1/2})). \end{aligned}$$

The proof is completed by rescaling $8p$ to p . \square

If the input of linear regression has full-rank covariance, then by Lemma 6, the system parameters can be fully identified. The certainty equivalent optimal controller in this case has a much better guarantee compared to the rank-deficient case, as shown in the following lemma.

Lemma 10 (Linear quadratic control). *Assume the finite-horizon linear time-varying system $((A_t^*, B_t^*, R_t^*)_{t=0}^{T-1}, (Q_t^*)_{t=0}^T)$ is stabilizable. Let $(K_t^*)_{t=0}^{T-1}$ be the optimal controller and $(P_t^*)_{t=0}^T$ be the solution to RDE (2.4) of system $((A_t^*, B_t^*, R_t^*)_{t=0}^{T-1}, (Q_t^*)_{t=0}^T)$. Let $(\hat{K}_t)_{t=0}^{T-1}$ be the optimal controller and $(\hat{P}_t)_{t=0}^T$ be the solution to RDE (2.4) of system $((\hat{A}_t, \hat{B}_t, R_t^*)_{t=0}^{T-1}, (\hat{Q}_t)_{t=0}^T)$, where $\|\hat{A}_t - A_t^*\|_2 \leq \varepsilon$, $\|\hat{B}_t - B_t^*\|_2 \leq \varepsilon$ and $\|\hat{Q}_t - Q_t^*\|_2 \leq \varepsilon$. Then there exists dimension-free constant $\varepsilon_0 > 0$ with ε_0^{-1} depending polynomially on system parameters, such that as long as $\varepsilon \leq \varepsilon_0$, $\|\hat{P}_t - P_t^*\|_2 = \mathcal{O}(\varepsilon)$, $\|\hat{K}_t - K_t^*\|_2 = \mathcal{O}(\varepsilon)$ for all $t \geq 0$, and $(\hat{K}_t)_{t=0}^{T-1}$ is $\mathcal{O}((d_x \wedge d_u)T\varepsilon^2)$ -optimal in system $((A_t^*, B_t^*, R_t^*)_{t=0}^{T-1}, (Q_t^*)_{t=0}^T)$.*

Proof. (Mania et al., 2019) has studied this problem in the infinite-horizon LTI setting; here we extend their result to the finite-horizon LTV setting.

We adopt the following compact formulation of a finite-horizon LTV system, introduced in (Zhang et al., 2021, Section 3), to reduce our setting to the infinite-horizon LTI setting; since noisy and noiseless LQR has the same associated Riccati equation and optimal controller, we consider the noiseless case here:

$$\begin{aligned} x &= [x_0; \dots; x_T], \quad u = [u_0; \dots; u_{T-1}], \quad w = [x_0; w_0; \dots; w_{T-1}], \\ A^* &= \begin{bmatrix} 0_{d_x \times d_x T} & 0_{d_x \times d_x} \\ \text{diag}(A_0^*, \dots, A_{T-1}^*) & 0_{d_x T \times d_x} \end{bmatrix}, \quad B^* = \begin{bmatrix} 0_{d_x \times d_u T} \\ \text{diag}(B_0^*, \dots, B_{T-1}^*) \end{bmatrix}, \\ Q^* &= \text{diag}(Q_0^*, \dots, Q_T^*), \quad R^* = \text{diag}(R_0^*, \dots, R_{T-1}^*), \quad K = [\text{diag}(K_0, \dots, K_{T-1}), 0_{d_u T \times d_x}]. \end{aligned}$$

The control inputs using state feedback controller $(K_t)_{t=0}^{T-1}$ can be characterized by $u = Kx$. Let $(P_t^K)_{t=0}^T$ be the associated cumulative cost matrix starting from step t . Then

$$P_t^K = (A_t^* + B_t^* K_t)^\top P_{t+1}^K (A_t^* + B_t^* K_t) + Q_t^* + K_t^\top R_t^* K_t,$$

and $P^K := \text{diag}(P_0^K, \dots, P_T^K)$ is the solution to

$$P^K = (A^* + B^*K)^\top P^K (A^* + B^*K) + Q^* + K^\top R^* K.$$

Similarly, the optimal cumulative cost matrix

$$P^* := \text{diag}(P_0^*, \dots, P_T^*),$$

produced by the RDE (2.4) in system (A^*, B^*, R^*, Q^*) (that is, system $((A_t^*, B_t^*, R_t^*)_{t=0}^{T-1}, (Q_t^*)_{t=0}^T)$), satisfies

$$P^* = (A^*)^\top (P^* - P^* B^* ((B^*)^\top P^* B^* + R^*)^{-1} (B^*)^\top P^*) A^* + Q^*.$$

Let $\hat{P} = \text{diag}(\hat{P}_0, \dots, \hat{P}_T)$ be the optimal cumulative cost matrices in system $(\hat{A}, \hat{B}, \hat{Q}, R^*)$ (that is, system $((\hat{A}_t, \hat{B}_t, R_t^*)_{t=0}^{T-1}, (\hat{Q}_t)_{t=0}^T)$) by the RDE (2.4). Define $K^* := [\text{diag}(K_0^*, \dots, K_T^*), 0_{d_u T \times d_x}]$, where $(K_t^*)_{t=0}^{T-1}$ is the optimal controller in system (A^*, B^*, Q^*, R^*) , and define \hat{K} similarly for system $(\hat{A}, \hat{B}, \hat{Q}, R^*)$. By definition, (A^*, B^*, Q^*) is stabilizable and observable in the sense of LTI systems. Therefore, by (Mania et al., 2019, Propositions 1 and 2), there exists dimension-free constant $\epsilon_0 > 0$ with ϵ_0^{-1} depending polynomially on system parameters such that as long as $\epsilon \leq \epsilon_0$,

$$\|\hat{P} - P^*\|_2 = \mathcal{O}(\epsilon), \quad \|\hat{K} - K^*\|_2 = \mathcal{O}(\epsilon),$$

and that \hat{K} stabilizes system (A^*, B^*) . By (Fazel et al., 2018, Lemma 12),

$$\begin{aligned} J((\hat{K}_t)_{t=0}^{T-1}) - J((K^*)_{t=0}^{T-1}) &= \sum_{t=0}^T \text{tr}(\Sigma_t (\hat{K}_t - K_t^*)^\top (R_t^* + (B_t^*)^\top P_{t+1}^* B_t^*) (\hat{K}_t - K_t^*)) \\ &= \text{tr}(\Sigma (\hat{K} - K^*)^\top (R^* + (B^*)^\top P^* B^*) (\hat{K} - K^*)), \end{aligned}$$

where $\Sigma = \text{diag}(\Sigma_0, \dots, \Sigma_T)$ and Σ_t is $\mathbb{E}[x_t x_t^\top]$ in system (A^*, B^*) under state feedback controller \hat{K} . As a result,

$$J((\hat{K}_t)_{t=0}^{T-1}) - J((K^*)_{t=0}^{T-1}) \leq \|\Sigma\|_2 \|R^* + (B^*)^\top P^* B^*\|_2 \|\hat{K} - K^*\|_F^2.$$

Since \hat{K} stabilizes system (A^*, B^*) , $\|\Sigma\|_2 = \mathcal{O}(1)$. Since

$$\|\hat{K} - K^*\|_F \leq ((d_x \wedge d_u)T)^{1/2} \|\hat{K} - K^*\|_2,$$

we conclude that

$$J((\hat{K}_t)_{t=0}^{T-1}) - J((K^*)_{t=0}^{T-1}) = \mathcal{O}((d_x \wedge d_u)T\epsilon^2).$$

□