

Identification and Adaptive Control of Markov Jump Systems: Sample Complexity and Regret Bounds

Yahya Sattar^{*†} Zhe Du^{*‡} Davoud Ataee Tarzanagh[‡] Necmiye Ozay[‡]
Laura Balzano[‡] Samet Oymak[†]

Abstract

Learning how to effectively control unknown dynamical systems is crucial for intelligent autonomous systems. This task becomes a significant challenge when the underlying dynamics are changing with time. Motivated by this challenge, this paper considers the problem of controlling an unknown Markov jump linear system (MJS) to optimize a quadratic objective. By taking a model-based perspective, we consider identification-based adaptive control for MJSs. We first provide a system identification algorithm for MJS to learn the dynamics in each mode as well as the Markov transition matrix, underlying the evolution of the mode switches, from a single trajectory of the system states, inputs, and modes. Through mixing-time arguments, sample complexity of this algorithm is shown to be $\mathcal{O}(1/\sqrt{T})$. We then propose an adaptive control scheme that performs system identification together with certainty equivalent control to adapt the controllers in an episodic fashion. Combining our sample complexity results with recent perturbation results for certainty equivalent control, we prove that when the episode lengths are appropriately chosen, the proposed adaptive control scheme achieves $\mathcal{O}(\sqrt{T})$ regret, which can be improved to $\mathcal{O}(\text{polylog}(T))$ with partial knowledge of the system. Our proof strategy introduces innovations to handle Markovian jumps and a weaker notion of stability common in MJSs. Our analysis provides insights into system theoretic quantities that affect learning accuracy and control performance. Numerical simulations are presented to further reinforce these insights.

1 Introduction

A canonical problem at the intersection of machine learning and control is that of adaptive control of an unknown dynamical system. An intelligent autonomous system is likely to encounter such a task; from an observation of the inputs and outputs, it needs to both learn and effectively control the dynamics. A commonly used control paradigm is the Linear Quadratic Regulator (LQR), which is theoretically well understood when system dynamics are linear and known. LQR also provides an interesting benchmark, when system dynamics are unknown, for reinforcement learning (RL) with continuous state and action spaces and for adaptive control [2, 4, 9, 16, 36, 47].

A generalization of linear dynamical systems called Markov jump linear systems (MJSs) models dynamics that switch between multiple linear systems, called modes, according to an underlying finite Markov chain. MJS allows for modeling a richer set of problems where the underlying dynamics can abruptly change over time. One can, similarly, generalize the LQR paradigm to MJS by using mode-dependent cost matrices, which allow different control goals under different modes. While the MJS-LQR problem is also well understood when one has perfect knowledge of the system dynamics [12, 14], in practice, it is not always possible to have a perfect knowledge of the system dynamics and the Markov transition matrix. For instance, a Mars rover optimally exploring an unknown heterogeneous terrain, optimal solar power generation on a cloudy day, or

^{*}Equal contribution.

[†]Dept. of Electrical and Computer Engineering, Univ. of California, Riverside. Email: {ysatt001, soymak}@ucr.edu.

[‡]Dept. of Electrical Engineering and Computer Science, Univ. of Michigan Ann Arbor. Email: {zhedu, tarzanaq, necmiye, girasole}@umich.edu.

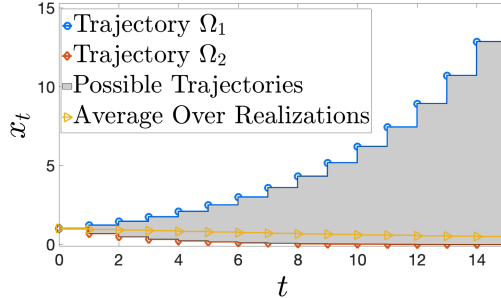


Figure 1: State trajectories for a two-modes MJS $\begin{cases} x_{t+1} = 1.2x_t \\ x_{t+1} = 0.7x_t \end{cases}$ with Markov matrix $\begin{bmatrix} 0.6 & 0.4 \\ 0.3 & 0.7 \end{bmatrix}$ and $\mathbf{x}_0 = 1$. Red and blue curves: mode switching sequences $\Omega_1 = \{1, 1, \dots\}$ and $\Omega_2 = \{2, 2, \dots\}$. Yellow curve: average over all realizations. Gray area: region for all possible trajectories.

controlling investments in financial markets may be modeled as MJS-LQR problems with unknown system dynamics [6, 8, 46, 60, 64]. Earlier works have aimed at analyzing the asymptotic properties (i.e., stability) of adaptive controllers for unknown MJSs both in continuous-time [7] and discrete-time [66] settings, however, despite the practical importance of MJSs, non-asymptotic sample complexity results and regret analysis for MJSs are lacking. The high-level challenge here is the hybrid nature of the problem that requires consideration of both the system dynamics and the underlying Markov transition matrix. A related challenge is that, typically, the stability of MJS is understood only in the *mean-square sense*. This is in stark contrast to the deterministic stability (e.g., as in LQR), where the system is guaranteed to converge towards an equilibrium point in the absence of noise. In contrast, the convergence of MJS trajectories towards an equilibrium depends heavily on how the switching between modes occurs. Figure 1 shows an example (adapted from [14]) of an MJS that is stable in the mean-square sense despite having an unstable mode. Clearly, under an unfavorable mode switching sequence, the system trajectory can still blow up. High-probability light tail bounds are therefore not applicable without very strong assumptions on the joint spectral radius of different modes (cf. [56]). Perhaps more surprisingly, there are examples of MJS with all modes individually stable, however due to switching, the system exhibits an unstable behavior on average, and the MJS is not mean-square stable (see Example 3.17 of [14]). Therefore, finding controllers to individually stabilize the mode dynamics does not guarantee that overall system will be stable when mode switches over time. This more relaxed notion of *mean-square stability* presents major challenges in learning, controlling, and statistical analysis.

Contributions: In this paper, we provide the first comprehensive system identification and regret guarantees for learning and controlling Markov jump linear systems using a single trajectory while assuming only mean-square stability (see Def. 3.1). Importantly, our guarantees are optimal in the trajectory length T . Specifically, our contributions are as follows¹:

- **System identification:** For an MJS with s modes, the system dynamics involve a Markov transition matrix $\mathbf{T} \in \mathbb{R}^{s \times s}$ and s state-input matrix pairs $(\mathbf{A}_i, \mathbf{B}_i)_{i=1}^s$. We provide an algorithm (Alg. 1) to estimate these dynamics with an error rate of $\mathcal{O}((n+p) \log(T) \sqrt{s/T})$, where n and p are the state and input dimensions respectively, and the $\mathcal{O}(1/\sqrt{T})$ dependence on the trajectory length T is optimal.
- **$\mathcal{O}(\sqrt{T})$ -regret bound:** We employ our system identification guarantees for the MJS-LQR. When the system dynamics are unknown, we show that the certainty-equivalent adaptive MJS-LQR Algorithm (Alg. 2) achieves a regret bound of $\mathcal{O}(\sqrt{T})$. Remarkably, this coincides with the optimal regret bound for the standard LQR problem obtained via certainty equivalence [47].
- **$\mathcal{O}(\text{polylog}(T))$ -regret with partial knowledge:** We also consider the practically relevant setting where the state matrices are unknown but the input matrices are known. We show that the regret bound can be significantly improved to $\mathcal{O}(\text{polylog}(T))$. This bound also coincides with the polylogarithmic regret bound for the standard LQR with the knowledge of the input matrix \mathbf{B} [10].

¹orders of magnitude here are up to polylogarithmic factors

Technical tools: Besides these key contributions to MJS control, our proof strategy introduces multiple innovations. To address Markovian mode transitions, we introduce a mixing-time argument to jointly track the approximate-dependence across the states and the modes. This in turn helps ensure each mode has sufficient samples and these samples are sufficiently informative. Secondly, as clarified further below, due to mean-square stability and mode transitions, it becomes non-trivial to determine whether the states have a light-tailed distribution (e.g., sub-gaussian or sub-exponential). To circumvent this, we develop intricate system identification arguments that allow for heavy-tailed states. Such arguments can potentially benefit other RL problems with heavy-tailed data.

2 Related Work

Our work is related to several topics in model-based reinforcement learning, system identification, and adaptive control. A comparison with the related works, in the LQR setting, is provided in Table 1.

- **System Identification:** Learning dynamical models has a long history in the control community, with major theoretical results being related to asymptotic properties under strong assumptions on persistence of excitation [45]. The problem becomes harder for hybrid and switched systems where the initial focus was on computational complexity as opposed to sample complexity of learning [40, 53]. There are some recent results on asymptotic consistency [30] in the switched system setting, a special case of MJS where the modes change in an independently and identically distributed manner. Identification of MJS has also attracted attention from different communities in the case when mode sequence is hidden [25, 63].

- **Sample Complexity of System Identification:** There is a recent surge of interest towards understanding the sample complexity of learning linear dynamical systems from a single trajectory under mild assumptions [52], using statistical tools like martingales [55, 59, 61] or mixing-time arguments [35, 50]. Recently, [34] provides precise rates for the finite-time identification of LTI systems using a single trajectory. The literature gets scarcer for switched systems. In [38], a novel approach based on Lyapunov equation is proposed for systems with stochastic switches, but no theory is built. [56] is one of the early works – and it seems to be the only work not assuming persistence of excitation – to provide finite sample analysis for learning systems with stochastic switches, yet with additional strong assumptions like independent switches and small joint spectral radius. The proof techniques developed within our work aim to obviate such assumptions. Our paper tackles the open problem of learning MJS from finite samples, obtained from a single trajectory, with theoretical guarantees under mild assumptions.

- **Learning-based Control and Regret Analysis:** As a direct application of single-trajectory system identification results, one can provide more sophisticated adaptive control guarantees from regret perspective [1, 2, 16, 23, 29, 47]. Specifically, [58] achieves $\mathcal{O}(\sqrt{T})$ regret lower bound for adaptive LQR control with unknown system dynamics, while with partial knowledge of the system [10] or persistence of excitation assumptions [37], one can achieve logarithmic regret [10, 37], as no additional excitation noise is needed to guarantee learnability of the system. However, in the MJS setting, due to the lack of well established identification analysis, prior works provide guarantees [7, 66] from the stability aspect. The case of input design without system state dynamics is considered in [5], which can be thought of as a generalization of linear bandits to have a Markovian structure in the reward function without any continuous dynamic structure. However, only a regret lower bound is provided in [5]. Finally, we refer the reader to the survey papers [28, 48, 54] for a broad overview of the recent developments on non-asymptotic system identification, adaptive control and reinforcement learning from the perspective of optimization and control.

- **Model-free Approaches:** Somehow orthogonal to the above developments, but still highly relevant, are approaches that sidestep system identification and try to learn an optimal controller (policy) directly (among many others, see e.g., [24, 49, 68, 69]). These works analyze the optimization landscape of LQR and related optimal control problems and provide polynomial-time algorithms that lead to a globally convergent search in the space of controllers. Importantly, these optimization algorithms do not require the knowledge of the system parameters as long as relevant quantities like gradients can be approximated from simulated system trajectories. More recently, this line of work is extended to MJSs in [33], significantly expanding their utility. However, these works require multiple trajectories to estimate the gradients as opposed to a controller that

adapts at run-time, therefore, they provide a complementary perspective to the single trajectory adaptation and regret analysis in our work.

A preliminary version of this work has been submitted to the American Control Conference 2022 [19], where we provide preliminary guarantees for the data-driven control of MJS. In contrast to this paper, Algorithm 1 in [19] performs a fairly-sophisticated double-subsampling to estimate the unknown MJS dynamics $(\mathbf{A}_i, \mathbf{B}_i)_{i=1}^S$ and \mathbf{T} with guarantees. On the other hand, Algorithm 1 in this paper uses all of the bounded samples from an MJS trajectory to estimate the unknown MJS dynamics. In this paper, we provide new and substantially improved results for the adaptive control of MJS with unknown state matrices but the input matrices are known. We also analyze the impact of “mean-square stability” (see Def. 3.1) on our regret bounds by replacing it with “uniform stability” (see Sec. 5.2). Hence, this paper also provides regret bounds for the adaptive control of MJS under uniform stability assumption. Furthermore, it also provides the necessary technical framework and the associated proofs. Lastly, compared to [19], the exposition of this paper is significantly improved by adding detailed discussions on the design of initial stabilizing controllers for the adaptive control of MJS and analyzing the sub-optimality gap for the offline control of MJS.

Table 1: Comparison with prior works in the LQR setting.

Model	Reference	Regret	Computational Complexity	Cost	Stabilizability/Controllability
LTI	[2]	\sqrt{T}	Exponential	Strongly Convex	Controllable
	[32]	\sqrt{T}	Exponential	Convex	Controllable
	[3] (one dim. systems)	\sqrt{T}	Polynomial	Strongly Convex	Stabilizable
	[15]	T^2	Polynomial	Convex	Stabilizable
	[47]	\sqrt{T}	Polynomial	Strongly Convex	Controllable
	[13]	\sqrt{T}	Polynomial	Strongly Convex	Strongly Stabilizable
	[22, 58]	\sqrt{T}	Polynomial	Strongly Convex	Stabilizable
MJS	[10] (known \mathbf{A} or \mathbf{B})	polylog(T)	Polynomial	Strongly Convex	Strongly Stabilizable
	Ours	$s^2\sqrt{T}$	Polynomial	Strongly Convex	MSS
	Ours (known $\mathbf{B}_{1:s}$)	$s^2\text{polylog}(T)$	Polynomial	Strongly Convex	MSS

3 Preliminaries and Problem Setup

Notations: We use boldface uppercase (lowercase) letters to denote matrices (vectors). For a matrix \mathbf{V} , $\rho(\mathbf{V})$ denotes its spectral radius. We use $\|\cdot\|$ to denote the Euclidean norm of vectors as well as the spectral norm of matrices. Similarly, we use $\|\cdot\|_1$ to denote the 1 -norm of a matrix/vector. The Kronecker product of two matrices \mathbf{M} and \mathbf{N} is denoted as $\mathbf{M} \otimes \mathbf{N}$. $\mathbf{V}_{1:s}$ denotes a set of s matrices $\{\mathbf{V}_i\}_{i=1}^s$ of same dimensions. We define $[s] := \{1, 2, \dots, s\}$ and $\|\mathbf{V}_{1:s}\| := \max_{i \in [s]} \|\mathbf{V}_i\|$. The i -th row or column of a matrix \mathbf{M} is denoted by $[\mathbf{M}]_i$, or $[\mathbf{M}]_{\cdot i}$ respectively. Orders of magnitude notation $\tilde{\mathcal{O}}(\cdot)$ hides $\log(1)$ or $\log^2(1)$ terms.

3.1 Markov Jump Linear Systems

In this paper we consider the identification and adaptive control of MJSs which are governed by the following state equation,

$$\mathbf{x}_{t+1} = \mathbf{A}_{(t)} \mathbf{x}_t + \mathbf{B}_{(t)} \mathbf{u}_t + \mathbf{w}_t \quad \text{s.t.} \quad (t) \sim \text{Markov Chain}(\mathbf{T}), \quad (3.1)$$

where $\mathbf{x}_t \in \mathbb{R}^n$, $\mathbf{u}_t \in \mathbb{R}^p$ and $\mathbf{w}_t \in \mathbb{R}^n$ are the state, input, and process noise of the MJS at time t with $\{\mathbf{w}_t\}_{t=0}^{\infty} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \frac{2}{\mathbf{w}} I_n)$. There are s modes in total, and the dynamics of mode i is given by the state matrix \mathbf{A}_i and input matrix \mathbf{B}_i . The active mode at time t is indexed by $(t) \in [s]$. Throughout, we assume the state \mathbf{x}_t and the mode (t) can be observed at time $t \geq 0$. The mode switching sequence $\{(t)\}_{t=0}^{\infty}$ follows a Markov chain with transition matrix $\mathbf{T} \in \mathbb{R}_+^{s \times s}$ such that for all $t \geq 0$, the ij -th element of \mathbf{T} denotes the

conditional probability $[\mathbf{T}]_{ij} := \mathbb{P}(\tau(t+1) = j \mid \tau(t) = i), \forall i, j \in [S]$. Throughout, we assume the initial state \mathbf{x}_0 , the mode switching sequence $\{\tau(t)\}_{t=0}$, and the noise $\{\mathbf{w}_t\}_{t=0}$ are mutually independent. We use $\text{MJS}(\mathbf{A}_{1:S}, \mathbf{B}_{1:S}, \mathbf{T})$ to refer to an MJS with state equation (3.1), parameterized by $(\mathbf{A}_{1:S}, \mathbf{B}_{1:S}, \mathbf{T})$.

For mode-dependent state-feedback controller $\mathbf{K}_{1:S}$ that yields the input $\mathbf{u}_t = \mathbf{K}_{\tau(t)} \mathbf{x}_t$, we use $\mathbf{L}_i := \mathbf{A}_i + \mathbf{B}_i \mathbf{K}_i$ to denote the closed-loop state matrix for mode i . We use $\mathbf{x}_{t+1} = \mathbf{L}_{\tau(t)} \mathbf{x}_t$ to denote the noise-free autonomous MJS, either open-loop ($\mathbf{L}_i = \mathbf{A}_i$) or closed-loop ($\mathbf{L}_i = \mathbf{A}_i + \mathbf{B}_i \mathbf{K}_i$). Due to the randomness in $\{\tau(t)\}_{t=0}$, it is common to consider the stability of MJS in the mean-square sense which is defined as follows.

Definition 3.1 (Mean-square stability [14]) We say MJS in (3.1) with $\mathbf{u}_t = 0$ is mean-square stable (MSS) if there exists $\epsilon, \delta > 0$ such that for any initial state \mathbf{x}_0 and mode $\tau(0)$, as $t \rightarrow \infty$, we have

$$\|\mathbb{E}[\mathbf{x}_t] - \mathbf{x}\| \rightarrow 0, \quad \|\mathbb{E}[\mathbf{x}_t \mathbf{x}_t^T] - \mathbf{P}\| \rightarrow 0, \quad (3.2)$$

where the expectation is over the Markovian mode switching sequence $\{\tau(t)\}_{t=0}$, the noise $\{\mathbf{w}_t\}_{t=0}$ and the initial state \mathbf{x}_0 . In the noise-free case ($\mathbf{w}_t = 0$), we have $\mathbf{x} = 0, \mathbf{P} = 0$. We say MJS in (3.1) with $\mathbf{w}_t = 0$ is (mean-square) stabilizable if there exists mode-dependent controller $\mathbf{K}_{1:S}$ such that the closed-loop MJS $\mathbf{x}_{t+1} = (\mathbf{A}_{\tau(t)} + \mathbf{B}_{\tau(t)} \mathbf{K}_{\tau(t)}) \mathbf{x}_t$ is MSS. We call such $\mathbf{K}_{1:S}$ a stabilizing controller.

The MSS of a noise-free autonomous MJS is related to the spectral radius of an augmented state matrix $\tilde{\mathbf{L}} \in \mathbb{R}^{sn^2 \times sn^2}$ with ij -th $n^2 \times n^2$ block given by $[\tilde{\mathbf{L}}]_{ij} := [\mathbf{T}]_{ji} \mathbf{L}_j \otimes \mathbf{L}_j$. As discussed in [14, Theorem 3.9], $\tilde{\mathbf{L}}$ can be viewed as the mapping from $\mathbb{E}[\mathbf{x}_t \mathbf{x}_t^T]$ to $\mathbb{E}[\mathbf{x}_{t+1} \mathbf{x}_{t+1}^T]$, thus a noise-free autonomous MJS is MSS if and only if $\rho(\tilde{\mathbf{L}}) < 1$. The analysis of this work highly depends on certain ‘‘mixing’’ of the MJS – the distributions of both state \mathbf{x}_t and mode $\tau(t)$ can converge close enough to some stationary distributions within finite time, which is guaranteed by the following assumption.

Assumption 1 The MJS in (3.1) has ergodic Markov chain and is stabilizable.

Ergodicity guarantees that the distribution of $\tau(t)$ converges to a unique strictly positive stationary distribution [27, Theorem 4.3.5]. Throughout, we let π denote the stationary distribution of \mathbf{T} and $\pi_{\min} := \min_i \pi(i)$. We further define the mixing time [43] of \mathbf{T} as $t_{\text{MC}} := \inf \{t \in \mathbb{N} : \max_{i \in [S]} \|([\mathbf{T}^t]_{i, \cdot}) - \pi\|_1 \leq 0.5\}$, to quantify its convergence rate. In our analysis, ergodicity and t_{MC} ensures that the MJS trajectory could have enough ‘‘visits’’ to every mode $i \in [S]$ thus providing us enough data to learn $[\mathbf{T}]_{i, \cdot}, \mathbf{A}_i$ and \mathbf{B}_i . On the other hand, stability (or stabilizability) characterized by the spectral radius of $\tilde{\mathbf{L}}$ guarantees the convergence/mixing of \mathbf{x}_t , which allows us to obtain weakly dependent sub-trajectories from a single trajectory of MJS, upon which the sample complexity of learning the matrices $\mathbf{A}_{1:S}$ and $\mathbf{B}_{1:S}$ can be established.

3.2 Problem Formulation

In this work we consider two major problems under the MJS setting: System identification and adaptive control, with identification being the core part of adaptive control.

(A) System Identification. This problem seeks to estimate unknown system dynamics from data, i.e. from input-output trajectory(ies), when one has the flexibility to design the inputs so that the collected data has nice statistical properties. In the MJS setting, one needs to estimate both the state/input matrices $\mathbf{A}_{1:S}, \mathbf{B}_{1:S}$ for every mode $i \in [S]$ as well as the Markov transition matrix \mathbf{T} . In this work, we seek to estimate the MJS dynamics from a single trajectory of states, inputs and mode observations $\{\mathbf{x}_t, \mathbf{u}_t, \tau(t)\}_{t=0}^T$ and provide finite sample guarantees. As mentioned earlier, MJS presents unique statistical analysis challenges due to Markovian jumps and a weaker notion of stability. Section 4 presents our system identification guarantees overcoming these challenges. These guarantees are further integrated into model-based control for MJS-LQR in Section 5.

(B) Online Linear Quadratic Regulator. In this paper, we consider the following finite-horizon Markov jump system linear quadratic regulator (MJS-LQR) problem:

$$\begin{aligned} \inf_{\mathbf{u}_{0:T}} \quad & J(\mathbf{u}_{0:T}) := \sum_{t=0}^T \mathbb{E}[\mathbf{x}_t^T \mathbf{Q}_{\tau(t)} \mathbf{x}_t + \mathbf{u}_t^T \mathbf{R}_{\tau(t)} \mathbf{u}_t], \\ \text{s.t.} \quad & \mathbf{x}_t, \tau(t) \sim \text{MJS}(\mathbf{A}_{1:S}, \mathbf{B}_{1:S}, \mathbf{T}). \end{aligned} \quad (3.3)$$

Algorithm 1: MJS-SYSID

Input: A mean square stabilizing controller \mathbf{K}_{1s} ; process and exploration noise variances $\frac{2}{\mathbf{w}}$ and $\frac{2}{\mathbf{z}}$; MJS trajectory $\{\mathbf{x}_t, \mathbf{z}_t, (t)\}_{t=0}^T$ generated using input $\mathbf{u}_t = \mathbf{K}_{(t)}\mathbf{x}_t + \mathbf{z}_t$ with $\mathbf{z}_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \frac{2}{\mathbf{z}}\mathbf{I}_p)$; and data clipping thresholds $c_{\mathbf{x}}, c_{\mathbf{z}}$.

- 1 **Estimate** $\mathbf{A}_{1s}, \mathbf{B}_{1s}$: **for all** modes $i \in [s]$ **do**
- 2 $S_i = \{t \mid (t) = i, \|\mathbf{x}_t\| \leq c_{\mathbf{x}} \sqrt{\log(T)}, \|\mathbf{z}_t\| \leq c_{\mathbf{z}}\}$ // sub-sample data
- 3 $\hat{1}_{1,i}, \hat{2}_{2,i} = \arg \min \sum_k S_i \|\mathbf{x}_{k+1} - \mathbf{A}_{1,i}\mathbf{x}_k - \mathbf{B}_{1,i}\mathbf{z}_k\|^2$ // regress with data S_i
- 4 $\hat{\mathbf{B}}_i = \hat{2}_{2,i} / \mathbf{z}, \hat{\mathbf{A}}_i = (\hat{1}_{1,i} - \hat{\mathbf{B}}_i \mathbf{K}_i) / \mathbf{w}$
- 5 **Estimate** \mathbf{T} : $[\hat{\mathbf{T}}]_{ji} = \frac{\sum_{t=1}^T \mathbf{1}_{\{(t)=i, (t-1)=j\}}}{\sum_{t=1}^T \mathbf{1}_{\{(t-1)=j\}}}$ //empirical frequency of transitions

Output: $\hat{\mathbf{A}}_{1s}, \hat{\mathbf{B}}_{1s}, \hat{\mathbf{T}}$

Here, the goal is to design control inputs to minimize the expected quadratic cost function composed of positive semi-definite cost matrices \mathbf{Q}_{1s} and \mathbf{R}_{1s} under the MJS dynamics (3.1). The quadratic cost incurred by the state \mathbf{x}_t represents the deviation from target values, e.g. desired velocity, position, angle, etc., whereas, the quadratic term in \mathbf{u}_t represents the control effort, e.g. energy consumption. The flexibility of having mode-dependent cost matrices allows one to design different control requirements or trade-offs under different circumstances. For the MJS-LQR problem (3.3), we assume the following.

Assumption 2 For all $i \in [s]$, $\mathbf{R}_i > 0, \mathbf{Q}_i > 0$.

Assumptions 1 and 2 together guarantee the solvability of MJS-LQR when the dynamics are known [14, Corollary A.21]. In the remaining of the paper, we use MJS-LQR($\mathbf{A}_{1s}, \mathbf{B}_{1s}, \mathbf{T}, \mathbf{Q}_{1s}, \mathbf{R}_{1s}$) to denote MJS-LQR problem (3.3) composed of MJS($\mathbf{A}_{1s}, \mathbf{B}_{1s}, \mathbf{T}$) and cost matrices $\mathbf{Q}_{1s}, \mathbf{R}_{1s}$.

Recall our assumption that the states \mathbf{x}_t and the modes (t) can be observed at time $t \geq 0$. With these observations, instead of a fixed and open-loop input sequence, one can design closed-loop policies that generate real-time control inputs based on the current observations, e.g. mode-dependent state-feedback controllers. When the dynamics $\mathbf{A}_{1s}, \mathbf{B}_{1s}, \mathbf{T}$ of the MJS are known, one can solve for the optimal controllers recursively via coupled discrete-time algebraic Riccati equations [14]. In this work, we assume the dynamics are unknown, and only the design parameters \mathbf{Q}_{1s} and \mathbf{R}_{1s} are known. Control schemes in this scenario are typically referred to as adaptive control, which usually involves procedures of learning, either the dynamics or directly the controllers. Adaptive control suffers additional costs as (i) the lack of the exact knowledge of the system and (ii) the exploration-exploitation trade-off – the necessity to sacrifice short-term input optimality to boost learning, so that overall long-term optimality can be improved.

Because of this, to evaluate the performance of an adaptive scheme, one is interested in the notion of regret – how much more cost it will incur if one could have applied the optimal controllers? In our setting, we compare the resulting cost against the optimal cost $T \cdot J^*$ where J^* is the optimal infinite-horizon average cost

$$J^* := \limsup_T \frac{1}{T} \inf_{\mathbf{u}_{0:T}} J(\mathbf{u}_{0:T}), \quad (3.4)$$

i.e., if one applies the optimal controller for infinitely long, how much cost one would get on average for each single time step. Compared to the regret analysis of standard adaptive LQR problem [15], in MJS-LQR setting, the cost analysis requires additional consideration of Markov chain mixing, which is addressed in this paper.

4 System Identification for MJS

Our MJS identification procedure is given in Algorithm 1. We assume one has access to an initial stabilizing controller \mathbf{K}_{1s} , which is a standard assumption in data-driven control [3, 13, 15, 32, 58] for LTI systems. For

MJSs, a thorough discussion on the validity of this assumption is provided in Section 6.1. Note that, if the open-loop MJS is already MSS, then one can simply set $\mathbf{K}_{1s} = 0$ and carry out MJS identification. Given an MJS trajectory $\{\mathbf{x}_t, \mathbf{z}_t, (t)\}_{t=0}^T$ generated using the input $\mathbf{u}_t = \mathbf{K}_{(t)}\mathbf{x}_t + \mathbf{z}_t$ (where $\mathbf{z}_t \sim \mathcal{N}(0, \frac{\sigma}{2}\mathbf{1})$ is the exploration noise), we sub-sample it, for each mode $i \in [s]$, to obtain s sub-trajectories with bounded states \mathbf{x}_t and excitations \mathbf{z}_t . This sub-sampling is required because of the mean-square stability, which can at most guarantee that the states are bounded in expectation. As a result of sub-sampling only bounded states/excitations, we obtain samples with manageable distributional properties. After appropriate scaling, we regress over these samples to obtain the estimates $\hat{\mathbf{A}}_i, \hat{\mathbf{B}}_i$ for each $i \in [s]$. Lastly, using the empirical frequency of observed modes, we obtain the estimate $\hat{\mathbf{T}}$.

The following theorem gives our main results on learning the dynamics of an unknown MJS from finite samples obtained from a single trajectory. One can refer to Theorems B.1 and B.17 in Appendix B for the detailed theorem statements and proofs.

Theorem 4.1 (Identification of MJS) *Suppose we run Algorithm 1 with $c_{\mathbf{x}} = \mathcal{O}(\sqrt{\bar{n}})$ and $c_{\mathbf{z}} = \mathcal{O}(\sqrt{\bar{p}})$. Let $\tilde{\mathbf{L}} = (\tilde{\mathbf{L}})$, where $\tilde{\mathbf{L}}$ is the augmented state matrix of the closed-loop MJS defined in Sec. 3.1. Suppose the trajectory length obeys $T \geq \tilde{\mathcal{O}}\left(\frac{\bar{s}t_{\text{MC}} \log^2(T)}{\min(1-\bar{\rho})}(n+p)\right)$. Then, under Assumption 1, with probability at least $1 - \bar{\rho}$, for all $i \in [s]$, we have*

$$\begin{aligned} \max\{\|\hat{\mathbf{A}}_i - \mathbf{A}_i\|, \|\hat{\mathbf{B}}_i - \mathbf{B}_i\|\} &\leq \tilde{\mathcal{O}}\left(\frac{(\bar{\mathbf{z}} + \bar{\mathbf{w}})(n+p)\log(T)}{\bar{\mathbf{z}} \min(1-\bar{\rho})} \sqrt{\frac{\bar{s}}{T}}\right), \\ \text{and } \|\hat{\mathbf{T}} - \mathbf{T}\| &\leq \tilde{\mathcal{O}}\left(\frac{1}{\min} \sqrt{\frac{\log(T)}{T}}\right). \end{aligned} \quad (4.1)$$

Corollary 4.2 *Consider the same setting of Theorem 4.1. Additionally, when \mathbf{B}_{1s} are known, setting $\bar{\mathbf{z}} = 0$ and solving only for the state matrices leads to the stronger upper bound $\|\hat{\mathbf{A}}_i - \mathbf{A}_i\| \leq \tilde{\mathcal{O}}\left(\frac{(n+p)\log(T)}{\min(1-\bar{\rho})} \sqrt{\frac{\bar{s}}{T}}\right)$ for all $i \in [s]$.*

Proof sketch [Theorem 4.1] Our proof strategy addresses the key challenges introduced by MJS and mean-square stability. We only emphasize the core technical challenges here. The idea is to think of the set S_i as a union of L subsets $S_i^{(k)}$ defined as follows:

$$S_i^{(k)} := \{t + kL \mid (t + kL) = i, \|\mathbf{x}_{t+kL}\| \leq c_{\mathbf{w}} + \sqrt{\bar{n}}, \|\mathbf{z}_{t+kL}\| \leq c_{\mathbf{z}} \sqrt{\bar{p}}\}, \quad (4.2)$$

where $0 \leq k \leq L-1$ is a fixed offset and $k = 1, 2, \dots, \lfloor \frac{T-L}{L} \rfloor$. The spacing of samples by $L \geq 1$ in each subset $S_i^{(k)}$ aims to reduce the statistical dependence across the samples belonging to that subset, to obtain weakly-dependent sub-trajectories. This weak dependence is due to the Markovian mode switching sequence $\{(t)\}_{t=0}^T$ – unique to the MJS setting – and the system’s memory (contributions from the past states). Thus L is primarily a function of the mixing-time (t_{MC}) of the Markov chain and the spectral radius ($\rho(\tilde{\mathbf{L}})$) of the MJS. At a high-level, by choosing sufficiently large L (e.g., $\mathcal{O}(T)$), we can upper/lower bound the empirical covariance matrix of the concatenated state vector $\mathbf{h}_k := [\frac{1}{\bar{\mathbf{w}}}\mathbf{x}_k \quad \frac{1}{\bar{\mathbf{z}}}\mathbf{z}_k]$ for all $k \in S_i^{(k)}$.

Unlike related works on system identification and regret analysis [15, 36, 37, 52, 59], mean-square stability does not lead to strong high-probability bounds, as one can only bound $\|\mathbf{x}_t\|$ or $\mathbf{x}_t\mathbf{x}_t$ in expectation. Therefore, in Algorithm 1, we sample only *bounded* state-excitation pairs $(\mathbf{x}_t, \mathbf{z}_t)$ on each mode $i \in [s]$. This boundedness enables us to control the covariance matrix of \mathbf{h}_k , despite MSS and potentially heavy-tailed states, via non-asymptotic tool-sets (e.g., Thm 5.44 of [65]). However, heavy-tailed empirical covariance lower bounds require independence, and our sub-sampled data are only “approximately independent” (coupled over modes and history). To make matters worse, the fact that we sub-sample only bounded states introduces further dependencies. To resolve this, we introduce a novel strategy to construct (for the purpose of analysis) an independent subset of *processed states* from this larger dependent set. The independence is ensured by conditioning on the mode-sequence and truncating the contribution of earlier states. We then use perturbation-based techniques (see e.g., [57]) to deal with actual (non-truncated) states. The final ingredient

is showing that, for each mode $i \in [S]$, with high probability, this carefully-crafted subset contains enough samples to ensure a well-conditioned covariance (with excitation provided by \mathbf{z}_t and \mathbf{w}_t). With this in place, after stitching together the estimation error from L sub-trajectories $\{(\mathbf{x}_{k'}, \mathbf{z}_{k'}, \mathbf{w}_{k'})\}_{k \in S_i}$ for $0 \leq i \leq L-1$, least-squares will accurately estimate \mathbf{A}_i and \mathbf{B}_i from the data $\{(\mathbf{x}_t, \mathbf{z}_t, \mathbf{w}_t)\}_{t \in S_i}$ with statistical error rate of $\tilde{O}(1/\sqrt{T})$. \blacksquare

Note that, our system identification result achieves near-optimal ($\mathcal{O}(1/\sqrt{T})$) dependence on the trajectory length T . However, the effective sample complexity of our system identification algorithm is $\mathcal{O}(s(n + \rho)^2 \log^2(T) / \lambda_{\min}^2)$, that is, the sample complexity grows quadratically in the state dimension n , which can potentially be improved to linear via a more refined analysis of the state-covariance (see e.g., [16, 59] for standard LTI systems). It also grows with the inverse of the minimum mode frequency as λ_{\min}^{-2} . Note that, λ_{\min} dictates the trajectory fraction of the least-frequent mode, thus, in the result λ_{\min}^{-1} multiplier is unavoidable. In Corollary 4.2, we show that, when $\mathbf{B}_{1:S}$ is assumed to be known, $\mathbf{A}_{1:S}$ can be estimated regardless of the exploration strength λ . This is because the excitation for the state matrix arises from noise \mathbf{w}_t . As we will see in Section 5, the distinct λ dependencies in Theorem 4.1 and Corollary 4.2 will lead to different regret bounds for MJS-LQR (albeit both bounds will be optimal up to polylog(T)).

5 Adaptive Control for MJS-LQR

Our adaptive MJS-LQR control scheme is given in Algorithm 2. It is performed on an epoch-by-epoch basis; a fixed controller is used for each epoch, and from epoch to epoch, the controller is updated using a newly collected MJS trajectory. Note that a new epoch is just a continuation of previous epochs instead of restarting the MJS. Similar to the discussion in Section 4, we assume, at the beginning of epoch 0, that one has access to a stabilizing controller $\mathbf{K}_{1:S}^{(0)}$. During epoch i , the controller $\mathbf{K}_{1:S}^{(i)}$ is used together with additive exploration noise $\mathbf{z}_t^{(i)} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \lambda \mathbf{I}_p)$ to boost learning. At the end of epoch i , the trajectory during that epoch is used to obtain a new MJS dynamics estimate $\mathbf{A}_{1:S}^{(i)}, \mathbf{B}_{1:S}^{(i)}, \mathbf{T}^{(i)}$ using Algorithm 1. Then, we set the controller $\mathbf{K}_{1:S}^{(i+1)}$ for epoch $i+1$ to be the optimal controller for the infinite-horizon MJS-LQR($\mathbf{A}_{1:S}^{(i)}, \mathbf{B}_{1:S}^{(i)}, \mathbf{T}^{(i)}, \mathbf{Q}_{1:S}, \mathbf{R}_{1:S}$), which can be computed as follows:

For a generic infinite-horizon MJS-LQR($\mathbf{A}_{1:S}, \mathbf{B}_{1:S}, \mathbf{T}, \mathbf{Q}_{1:S}, \mathbf{R}_{1:S}$), its optimal controller is given by $\mathbf{K}_{1:S}$ such that for all $j \in [S]$,

$$\mathbf{K}_j := -(\mathbf{R}_j + \mathbf{B}_{j-1}(\mathbf{P}_{1:S})\mathbf{B}_j)^{-1} \mathbf{B}_{j-1}(\mathbf{P}_{1:S})\mathbf{A}_j, \quad (5.1)$$

where $\mathbf{P}_{1:S} := \sum_{k=1}^S [\mathbf{T}]_k \mathbf{P}_k$ and $\mathbf{P}_{1:S}$ is the solution to the following coupled discrete-time algebraic Riccati equations (cDARE):

$$\mathbf{P}_j = \mathbf{A}_{j-1}(\mathbf{P}_{1:S})\mathbf{A}_j + \mathbf{Q}_j - \mathbf{A}_{j-1}(\mathbf{P}_{1:S})\mathbf{B}_j(\mathbf{R}_j + \mathbf{B}_{j-1}(\mathbf{P}_{1:S})\mathbf{B}_j)^{-1} \mathbf{B}_{j-1}(\mathbf{P}_{1:S})\mathbf{A}_j, \quad (5.2)$$

for all $j \in [S]$. In practice, cDARE can be solved efficiently via value iteration or LMIs [14]. Note that cDARE may not be solvable for arbitrary parameters, but our theory guarantees that when epoch lengths are appropriately chosen, cDARE parameterized by $\mathbf{A}_{1:S}^{(i)}, \mathbf{B}_{1:S}^{(i)}, \mathbf{T}^{(i)}, \mathbf{Q}_{1:S}, \mathbf{R}_{1:S}$ is solvable for every epoch i . This control design based on the estimated dynamics is also referred to as certainty equivalent control.

To achieve theoretically guaranteed performance, i.e., sub-linear regret, the key is to have a subtle scheduling of epoch lengths T_i and exploration noise variance $\lambda_{z,i}$. We choose T_i to increase exponentially with rate $\gamma > 1$, and set $\lambda_{z,i} = \lambda_w / \sqrt{T_i}$, which collectively guarantee $\tilde{O}(\sqrt{T})$ regret when combined with the system identification result from Theorem 4.1. Intuitively, this scheduling can be interpreted as follows: (i) the increase of epoch lengths guarantees we have more accurate MJS estimates thus more optimal controllers; (ii) as the controller becomes more optimal we can gradually decrease the exploration noise and deploy (exploit) the controller for a longer time. Note that the scheduling rate γ has a similar role to the discount factor in reinforcement learning: smaller γ aims to reduce short-term cost while larger γ aims to reduce long-term cost.

Algorithm 2: Adaptive MJS-LQR

Input: Initial epoch length T_0 ; initial stabilizing controller $\mathbf{K}_{1s}^{(0)}$; epoch incremental ratio $\gamma > 1$; and data clipping thresholds c_x, c_z

- 1 **for** $i = 0, 1, 2, \dots$ **do**
- 2 Set epoch length $T_i = \lfloor T_0 \gamma^i \rfloor$.
- 3 Set exploration noise variance $\frac{2}{z_i} = \frac{2}{\frac{w}{T_i}}$.
- 4 Evolve the MJS for T_i steps with $\mathbf{u}_t^{(i)} = \mathbf{K}_{(t)(i)}^{(i)} \mathbf{x}_t^{(i)} + \mathbf{z}_t^{(i)}$ with $\mathbf{z}_t^{(i)} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \frac{2}{z_i} \mathbf{I}_p)$ and record the trajectory $\{\mathbf{x}_t^{(i)}, \mathbf{z}_t^{(i)}, \mathbf{z}_t^{(i)}(t)\}_{t=0}^{T_i}$.
- 5 $\mathbf{A}_{1s}^{(i)}, \mathbf{B}_{1s}^{(i)}, \mathbf{T}^{(i)} = \text{MJS-SYSID}(\mathbf{K}_{1s}^{(i)}, \frac{2}{w}, \frac{2}{z_i}, \{\mathbf{x}_t^{(i)}, \mathbf{z}_t^{(i)}, \mathbf{z}_t^{(i)}(t)\}_{t=0}^{T_i}, c_x, c_z)$.
- 6 Set the controller $\mathbf{K}_{1s}^{(i+1)}$ for the next epoch to be the optimal controller for the infinite-horizon MJS-LQR($\mathbf{A}_{1s}^{(i)}, \mathbf{B}_{1s}^{(i)}, \mathbf{T}^{(i)}, \mathbf{Q}_{1s}, \mathbf{R}_{1s}$).
- 7 **end**

5.1 Regret Analysis

We define filtration $\mathcal{F}_{-1}, \mathcal{F}_0, \mathcal{F}_1, \dots$ such that $\mathcal{F}_{-1} := (\mathbf{x}_0, (0))$ is the sigma-algebra generated by the initial state and initial mode, and $\mathcal{F}_i := (\mathbf{x}_0, (0), \{\mathbf{z}_t^{(j)}(t)\}_{t=1}^{T_j}\}_{j=0}^i, \mathbf{w}_0, \{\mathbf{w}_t^{(j)}\}_{t=1}^{T_j}\}_{j=0}^i, \mathbf{z}_0, \{\mathbf{z}_t^{(j)}\}_{t=1}^{T_j}\}_{j=0}^i)$ is the sigma-algebra generated by the randomness up to epoch i . Note that the initial state $\mathbf{x}_0^{(i)}$ of epoch i is also the final state $\mathbf{x}_{T_{i-1}}^{(i-1)}$ of epoch $i-1$, therefore, $\mathbf{x}_0^{(i)}$ is \mathcal{F}_{i-1} -measurable, and so is $(i)(0)$. Suppose time step t belongs to epoch i , then we define the following conditional expected cost at time t as:

$$c_t = \mathbb{E}[\mathbf{x}_t \mathbf{Q}_{(t)} \mathbf{x}_t + \mathbf{u}_t \mathbf{R}_{(t)} \mathbf{u}_t | \mathcal{F}_{i-1}], \quad (5.3)$$

and cumulative cost as $J_T = \sum_{t=1}^T c_t$. We define the total regret and epoch- i regret as

$$\text{Regret}(T) = J_T - T J^*, \quad \text{Regret}_i = \left(\sum_{t=1}^{T_i} c_{T_0 + \dots + T_{i-1} + t} \right) - T_i J^*. \quad (5.4)$$

Then, we have $\text{Regret}(T) = \mathcal{O}(\sum_{i=1}^{\log(T/T_0)} \text{Regret}_i)$, where regret of epoch 0 is ignored as it does not scale with time T . Let \mathbf{K}_{1s} denote the optimal controller for the infinite-horizon MJS-LQR($\mathbf{A}_{1s}, \mathbf{B}_{1s}, \mathbf{T}, \mathbf{Q}_{1s}, \mathbf{R}_{1s}$) problem. $\tilde{\mathbf{L}}^{(0)}$ and $\tilde{\mathbf{L}}$ denote the closed-loop augmented state matrices under the initial controller $\mathbf{K}_{1s}^{(0)}$ and \mathbf{K}_{1s} respectively, and we let $\tilde{\rho} := \max\{\|\tilde{\mathbf{L}}^{(0)}\|, \|\tilde{\mathbf{L}}\|\}$. With these definitions, we have the following sub-linear regret guarantee. Please refer to Theorem C.11 in the appendix for the complete version and proof.

Theorem 5.1 (Sub-linear regret) *Assume that the initial state $\mathbf{x}_0 = 0$, and Assumptions 1 and 2 hold. In Algorithm 2, suppose hyper parameters $c_x = \mathcal{O}(\sqrt{n})$, $c_z = \mathcal{O}(\sqrt{p})$, and $T_0 \geq \tilde{\mathcal{O}}\left(\frac{\bar{s}_{\text{MC}} \log^2(T_0)}{\min(1-\gamma)}(n+p)\right)$. Then, with probability at least $1 - \delta$, Algorithm 2 achieves*

$$\text{Regret}(T) \leq \tilde{\mathcal{O}}\left(\frac{s^2 \rho (n^2 + p^2)}{2 \min} \frac{w}{\tilde{\rho}} \log^2(T) \sqrt{T}\right) + \mathcal{O}\left(\frac{\sqrt{ns} \log^3(T)}{\tilde{\rho}}\right). \quad (5.5)$$

Proof sketch [Theorem 5.1] For simplicity, we only show the dominant $\tilde{\mathcal{O}}(\cdot)$ term here and leave the complete proof to appendix. Define the estimation error after epoch i as $\mathbf{A}_{\mathbf{B}}^{(i)} := \max_j \max_{[s]} \{\|\mathbf{A}_j^{(i)} - \mathbf{A}_j\|, \|\mathbf{B}_j^{(i)} - \mathbf{B}_j\|\}$ and $\mathbf{T}^{(i)} := \|\mathbf{T}^{(i)} - \mathbf{T}\|$. Analyzing the finite-horizon cost and combining the infinite-horizon perturbation results in [18], we can bound epoch- i regret as $\text{Regret}_i \leq \mathcal{O}\left(T_i \frac{2}{z_i} + T_i \frac{2}{w} \left(\mathbf{A}_{\mathbf{B}}^{(i-1)} + \mathbf{T}^{(i-1)}\right)^2\right)$. Plugging in $\frac{2}{z_i} = \frac{2}{\frac{w}{T_i}}$ and the upper bounds on the estimation errors $\mathbf{A}_{\mathbf{B}}^{(i)} \leq \tilde{\mathcal{O}}\left(\frac{z_i + w}{z_i \min} \frac{\bar{s}(n+p) \log(T_i)}{T_i}\right)$ and $\mathbf{T}^{(i)} \leq \tilde{\mathcal{O}}\left(\sqrt{\frac{\log(T_i)}{T_i}}\right)$

from Theorem 4.1, we have $\text{Regret}_i \leq \tilde{\mathcal{O}}\left(\frac{s^2 p(n^2 + p^2)}{2 \min} \frac{2}{\mathbf{w}} \sqrt{T_i} \log^2(T_i)\right)$. Finally, since $T_i = \mathcal{O}(T_0^{-i})$ from Alg. 2, we have $\text{Regret}(T) = \sum_{i=1}^{\mathcal{O}(\log(\frac{T}{T_0}))} \text{Regret}_i \leq \tilde{\mathcal{O}}\left(\frac{s^2 p(n^2 + p^2)}{2 \min} \frac{2}{\mathbf{w}} \sqrt{T} \log\left(\frac{T}{T_0}\right) \left(\frac{-}{-1}\right)^3 \left(\log\left(\frac{T}{T_0}\right) - \sqrt{\log\left(\frac{T}{T_0}\right)}\right)\right) = \tilde{\mathcal{O}}\left(\frac{s^2 p(n^2 + p^2)}{2 \min} \frac{2}{\mathbf{w}} \text{polylog}(T) \sqrt{T}\right)$. \blacksquare

One can see the interplay between T and $\frac{1}{\gamma}$ from the term $\left(\frac{-}{-1}\right)^3 \left(\log\left(\frac{T}{T_0}\right) - \sqrt{\log\left(\frac{T}{T_0}\right)}\right)$ in the proof sketch. Specifically, when horizon T is smaller, a smaller $\frac{1}{\gamma}$ minimizes the upper bound, and vice versa. This further provides a mathematical justification for $\frac{1}{\gamma}$ being similar to the discount factor in reinforcement learning in early discussions.

5.2 Two Special Cases

5.2.1 Tighter probability bound under uniform stability

Note that the regret upper bound (5.5) in Theorem 5.1 has the second term depending on the failure probability $\frac{1}{\gamma}$ through $\frac{1}{\gamma}$. Though this term has a much milder dependency on the time horizon T , when setting $\frac{1}{\gamma}$ to be small, it can still easily outweigh the other $\tilde{\mathcal{O}}(\cdot)$ term in (5.5), which only has $\log(\frac{1}{\gamma})$ dependency, and can result in overly pessimistic regret bounds. The main cause of this $\frac{1}{\gamma}$ term is that in the regret analysis, one needs to factor in the cumulative impact of initial state of every epoch, i.e. $\sum_i \|\mathbf{x}_0^{(i)}\|^2$. Since MSS guarantees the stability and state convergence only in the mean-square sense, we can, at best, only bound $\mathbb{E}[\|\mathbf{x}_0^{(i)}\|^2]$ and then use the Markov inequality: with probability at least $1 - \frac{1}{\gamma}$, $\|\mathbf{x}_0^{(i)}\|^2 \leq \mathbb{E}[\|\mathbf{x}_0^{(i)}\|^2] / \frac{1}{\gamma}$. Furthermore, in Appendix C.4, we construct an MJS example that is MSS but no dependencies better than $\frac{1}{\gamma}$ can be established. Fortunately, there exists an easy workaround to get rid of this $\frac{1}{\gamma}$ dependency if the MJS is uniformly stable [42, 44], which enforces stability under arbitrary switching sequences thus is stronger than MSS. It allows us to bound $\mathbf{x}_0^{(i)}$ using tail inequalities much tighter than the Markov inequality and obtain $\|\mathbf{x}_0^{(i)}\|^2 \leq \mathcal{O}(\log(\frac{1}{\gamma}))$. In the end, in the regret bound, $\frac{1}{\gamma}$ can be improved to $\log(\frac{1}{\gamma})$.

One type of uniform stability assumption that can help us in this case is regarding the closed-loop MJS under the optimal controllers. We let $\mathbf{K}_{1:S}$ denote the optimal controller for the infinite-horizon MJS-LQR($\mathbf{A}_{1:S}, \mathbf{B}_{1:S}, \mathbf{T}, \mathbf{Q}_{1:S}, \mathbf{R}_{1:S}$) and define closed-loop state matrices $\mathbf{L}_i = \mathbf{A}_i + \mathbf{B}_i \mathbf{K}_i$ for all $i \in [S]$. We let $\rho(\mathbf{L}_{1:S})$ denote the joint spectral radius of $\mathbf{L}_{1:S}$, i.e. $\rho(\mathbf{L}_{1:S}) := \lim_{l \rightarrow \infty} \max_{i_1, \dots, i_l \in [S]^l} \|\mathbf{L}_{i_1} \cdots \mathbf{L}_{i_l}\|^{\frac{1}{l}}$, and we say $\mathbf{L}_{1:S}$ is uniformly stable if and only if $\rho(\mathbf{L}_{1:S}) < 1$. The resulting regret bound is outlined in the following theorem, with its complete version and proof provided in Theorem C.12 of Appendix C.4.1.

Theorem 5.2 (Regret Under Uniform Stability) *Assume that the initial state $\mathbf{x}_0 = 0$, Assumptions 1 and 2 hold, and $\mathbf{L}_{1:S}$ is uniformly stable. If hyper-parameters T_0 , $c_{\mathbf{x}}$, and $c_{\mathbf{z}}$ are chosen as sufficiently large, with probability at least $1 - \frac{1}{\gamma}$, Algorithm 2 achieves*

$$\text{Regret}(T) \leq \tilde{\mathcal{O}}\left(\frac{s^2 p(n^2 + p^2)}{2 \min} \frac{2}{\mathbf{w}} \log^2(T) \sqrt{T}\right). \quad (5.6)$$

5.2.2 Partial knowledge of dynamics

From Corollary 4.2, we know that when input matrices $\mathbf{B}_{1:S}$ are known, no further exploration noise is needed to identify the state matrices $\mathbf{A}_{1:S}$ or Markov matrix \mathbf{T} . This can also be applied to the adaptive MJS-LQR setting, and the resulting regret bound can improve (from $\mathcal{O}(\log^2(T) \sqrt{T})$ to $\mathcal{O}(\log^3(T))$) since exploration noise incurs additional costs. The result is given by the following corollary, and we omit the proof due to its similarity to the proofs of Theorems 5.1 and 5.2.

Corollary 5.3 (Poly-logarithmic regret) *When $\mathbf{B}_{1:S}$ are known, it suffices to set the exploration noise to be $\mathbf{z}_{:,i} = 0$ for all i in Algorithm 2. Then, the regret bound in Theorem 5.1 becomes $\text{Regret}(T) \leq$*

$$\tilde{\mathcal{O}}\left(\frac{s^2 p(n^2 + p^2)}{2} \frac{2}{\min} \log^3(T)\right) + \mathcal{O}\left(-\frac{\bar{n}s \log^3(T)}{2}\right). \text{ Additionally, the regret bound in Theorem 5.2 becomes } \text{Regret}(T) \leq \tilde{\mathcal{O}}\left(\frac{s^2 p(n^2 + p^2)}{2} \frac{2}{\min} \log^3(T)\right).$$

6 Discussion

In this section, we discuss how one may obtain the initial stabilizing controller for MJS as required in the input to Algorithms 1 and 2 and the application of our results to offline data-driven control.

6.1 Initial Stabilizing Controllers

Having access to an initial stabilizing controller has become a very common assumption in system identification (see for instance [41] and references therein) and adaptive control [3, 13, 15, 32, 58] for LTI systems. On the other hand, for work where no initial stabilizing controller is required, there is usually a separate warm-up phase at the beginning, where coarse dynamics is learned, upon which a stabilizing controller is computed. Recent non-asymptotic system identification results [21, 55] on potentially unstable LTI systems can be used to obtain coarse dynamics without stabilizing controller. One can use random linear feedback to construct a confidence set of the dynamics such that any point in this set can produce a stabilizing controller by solving Riccati equations [20]. In the model-free setting, [39] provides asymptotic results and relies on persistent excitation assumption. [11] designs subtle scaled one-hot vector input and collects the trajectory to estimate the dynamics, then a stabilizing controller can be solved via semi-definite programming. For MJS or general switched systems, to the best of our knowledge, there is no work on stabilizing unknown dynamics using single trajectory with guarantees. One challenge is, as we discussed in Section 1, the individual mode stability and overall stability does not imply each other due to mode switching. However, as outlined below, we can approach this problem leveraging what is recently done for the LTI case in the aforementioned literature (modulo some additional assumptions).

Similar to the LTI case, suppose we could obtain some coarse dynamics estimate $\hat{\mathbf{A}}_{1s}, \hat{\mathbf{B}}_{1s}, \hat{\mathbf{T}}$, then we can solve for the optimal controller $\hat{\mathbf{K}}_{1s}$ of the infinite-horizon MJS-LQR($\hat{\mathbf{A}}_{1s}, \hat{\mathbf{B}}_{1s}, \hat{\mathbf{T}}, \mathbf{Q}_{1s}, \mathbf{R}_{1s}$) via Riccati equations. To investigate when $\hat{\mathbf{K}}_{1s}$ can stabilize the MJS, the key is to obtain sample complexity guarantees for this coarse dynamics, i.e. dependence of estimation error $\|\hat{\mathbf{A}}_i - \mathbf{A}_i\|$, $\|\hat{\mathbf{B}}_i - \mathbf{B}_i\|$, and $\|\hat{\mathbf{T}} - \mathbf{T}\|$ on sample size. Fortunately [18] provides the required estimation accuracy under which $\hat{\mathbf{K}}_{1s}$ is guaranteed to be stabilizing. Thus, combining [18] with the estimation error bounds (in terms of sample size), the required accuracy can be translated to the required number of samples. Note that learning \mathbf{T} is the same as learning a Markov chain, thus using the mode transition pair frequencies in an arbitrary single MJS trajectory, we can obtain an estimate $\hat{\mathbf{T}}$ as in Algorithm 1, and its sample complexity is given in Lemma B.1 in Appendix B. The more challenging part is the identification scheme and corresponding sample complexity for $\hat{\mathbf{A}}_{1s}$ and $\hat{\mathbf{B}}_{1s}$. Here, we outline two potential schemes.

- Suppose we could generate N i.i.d. MJS rollout trajectories, each with length T (small T , e.g. $T = 1$, is preferred to avoid potential unstable behavior and for the ease of the implementation). We can obtain least squares estimates $\hat{\mathbf{A}}_{1s}, \hat{\mathbf{B}}_{1s}$ using only $\{\mathbf{x}_T, \mathbf{x}_{T-1}, \mathbf{u}_{T-1}, (T-1)\}$ from each trajectory, which is similar to the scheme in [16] for LTI systems. Since only i.i.d. data is used in the computation, one can easily obtain the sample complexity in terms of N .
- If each mode in the MJS can run in isolation (i.e. for any $i \in [S]$, $(t) = i$ for all t) so that it acts as an LTI system, we could use recent advances [21, 55] on single-trajectory open-loop LTI system identification to obtain coarse estimates together with sample complexity for $\hat{\mathbf{A}}_i$ and $\hat{\mathbf{B}}_i$ for every mode i .

We also note that while finding an initial stabilizing controller is theoretically very interesting and challenging, most results we know of are limited to simulated or numerical examples (see for instance [41] and references therein). This is because, from a practical standpoint, an initial stabilizing controller is almost

required in model-based approaches since running experiments with open-loop unstable plants can be very dangerous as the state could explode quickly.

6.2 Offline Data-Driven Control

In many scenarios, we may not be able to perform learning and control in real time due to limited onboard computing resources or measurement sensors. In this case, the dynamics is usually learned in a one-shot way at the beginning, and the resulting controller will be deployed forever without any further update. The controller suboptimality in this non-adaptive setting does not improve over time, thus the regret will increase linearly over time rather than sublinearly as in our work. The natural performance metric in this case is the time-averaged regret, which can also be viewed as the slope of the cumulative regret with respect to time. The system identification scheme and corresponding sample complexity developed in this paper can also help address this problem.

Suppose we obtain MJS estimate $\hat{\mathbf{A}}_{1s}, \hat{\mathbf{B}}_{1s}, \hat{\mathbf{T}}$ from a length- T_0 rollout trajectory using Algorithm 1 and solve for the controller $\hat{\mathbf{K}}_{1s}$ that is optimal for the infinite-horizon MJS-LQR($\hat{\mathbf{A}}_{1s}, \hat{\mathbf{B}}_{1s}, \hat{\mathbf{T}}, \mathbf{Q}_{1s}, \mathbf{R}_{1s}$) via Riccati equations. Let $\hat{J} := \limsup_T \frac{1}{T} \mathcal{J}(\hat{\mathbf{K}}_{1s}(t), \mathbf{x}_t)$ denote the infinite-horizon average cost incurred when we deploy $\hat{\mathbf{K}}_{1s}$ indefinitely. Combining our identification sample complexity result in Theorem 4.1 with the infinite-horizon MJS-LQR perturbation result in [18], we can easily obtain an upper bound on the suboptimality, $\hat{J} - J \leq \tilde{\mathcal{O}}(\log^2(T_0)/T_0)$, which provides the required rollout trajectory length T_0 if certain suboptimality is desired.

7 Numerical Experiments

We provide experiments to investigate the efficiency and verify the theory of the proposed algorithms on synthetic datasets. Throughout, we show results from a synthetic experiment where entries of the true system matrices ($\mathbf{A}_{1s}, \mathbf{B}_{1s}$) were generated randomly from a standard normal distribution. We further scale each \mathbf{A}_i to have $\|\mathbf{A}_i\| \leq 0.5$. Since this guarantees the MJS itself is MSS, as we discussed in Sec 4, we set controller $\mathbf{K}_{1s} = 0$ in system identification Algorithm 1 and initial stabilizing controller $\mathbf{K}_{1s}^{(0)} = 0$ in adaptive MJS-LQR Algorithm 2. For the cost matrices ($\mathbf{Q}_{1s}, \mathbf{R}_{1s}$), we set $\mathbf{Q}_j = \underline{\mathbf{Q}}_j \underline{\mathbf{Q}}_j$, and $\mathbf{R}_j = \underline{\mathbf{R}}_j \underline{\mathbf{R}}_j$ where $\underline{\mathbf{Q}}_j \in \mathbb{R}^{n \times n}$ and $\underline{\mathbf{R}}_j \in \mathbb{R}^{p \times p}$ were generated from a standard normal distribution. The Markov matrix $\mathbf{T} \in \mathbb{R}_+^{s \times s}$ was sampled from a Dirichlet distribution $\text{Dir}((s-1) \cdot \mathbf{I}_s + 1)$, where \mathbf{I}_s denotes the identity matrix. We assume that we had equal probability of starting in any initial mode.

Since for system identification, our main contribution is estimating \mathbf{A}_{1s} and \mathbf{B}_{1s} of the MJS, we omit the plots for estimating \mathbf{T} . Let $\hat{\Psi}_j = [\hat{\mathbf{A}}_j, \hat{\mathbf{B}}_j]$ and $\Psi_j = [\mathbf{A}_j, \mathbf{B}_j]$. We use $\|\hat{\Psi} - \Psi\|/\|\Psi\| := \max_j \|\hat{\Psi}_j - \Psi_j\|/\|\Psi_j\|$ to investigate the convergence behaviour of MJS-SYSID Algorithm 1. The clipping constants in this algorithm, i.e., $C_{\text{sub}}, c_{\mathbf{x}}$, and $c_{\mathbf{z}}$ are chosen based on their lower bounds provided in Theorem 5.1. In all the aforementioned algorithms, the depicted results are averaged over 10 independent replications.

7.1 Performance of MJS-SYSID

In this section, we investigate the performance of our MJS-SYSID method, i.e., Algorithm 1. We first empirically evaluate the effect of the noise variances \mathbf{w} and \mathbf{z} . In particular, we study how the system errors vary with (i) $\mathbf{w} = 0.01, \mathbf{z} \in \{0.01, 0.02, 0.1\}$ and (ii) $\mathbf{z} = 0.01, \mathbf{w} \in \{0.01, 0.02, 0.1\}$. The number of states, inputs, and modes are set to $n = 5, p = 3$, and $s = 5$, respectively. Fig. 2 (a) and (b) demonstrate how the relative estimation error $\|\hat{\Psi} - \Psi\|/\|\Psi\|$ changes as T increases. Each curve on the plot represents a fixed \mathbf{w} and \mathbf{z} . These empirical results are all consistent with the theoretical bound of MJS-SYSID given in (4.1). In particular, the estimation errors degrade with increasing \mathbf{w} and decreasing \mathbf{z} , respectively.

Now, we fix $\mathbf{w} = \mathbf{z} = 0.01$ and investigate the performance of the MJS-SYSID with varying number of states, inputs, and modes. Fig. 2 (c) and (d) show how the estimation error $\|\hat{\Psi} - \Psi\|/\|\Psi\|$ changes with (left) $s = 5, n \in \{5, 10, 20\}, p = n - 2$ and (right) $n = 5, p = n - 2, s \in \{5, 10, 20\}$. As we can see, the MJS-SYSID has better performance with small n, p and s which is consistent with (4.1).

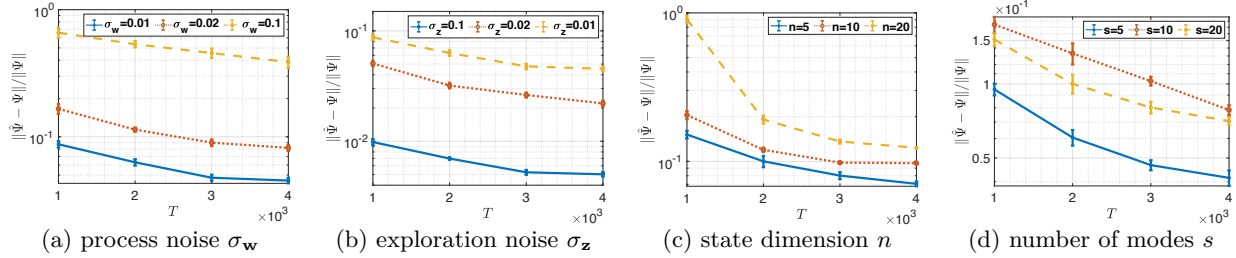


Figure 2: Performance profiles of MJS-SYSID with varying: (a) process noise \mathbf{w} , (b) exploration noise \mathbf{z} , (c) state dimension n , and (d) number of modes s .

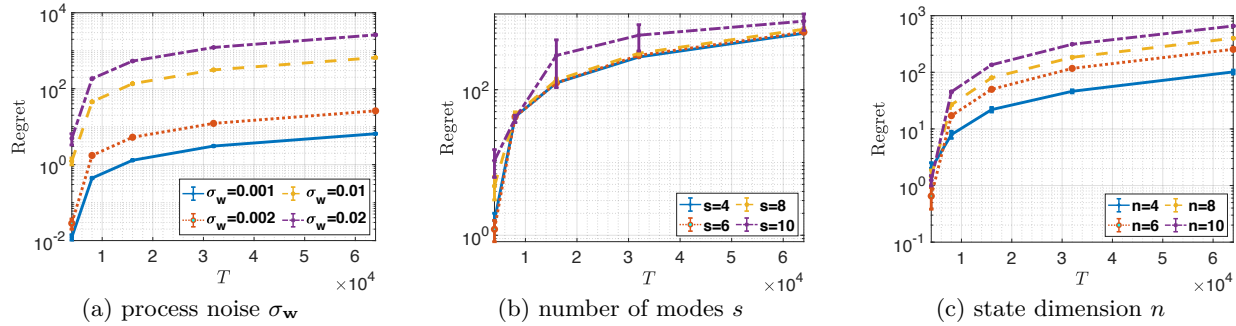


Figure 3: Performance profiles of Adaptive MJS-LQR with varying: (a) process noise \mathbf{w} , (b) number of modes s , (c) state dimension n .

7.2 Performance of Adaptive MJS-LQR

In our next series of experiments, we explore the sensitivity of the regret bounds to the system parameters. In these experiments, we set the initial epoch length $T_0 = 2000$ and incremental ratio $\rho = 2$. We select five epochs to run Algorithm 2. As an intermediate step for computing controller $\mathbf{K}_1^{(i+1)}$ in Algorithm 2, the coupled Riccati equations (5.2) are solved via value iteration, and the iteration stops when the parameter variation between two iterations falls below 10^{-6} , or iteration number reaches 10^4 .

Fig. 3 demonstrates how regret bounds vary with (a) $\mathbf{w} \in \{0.001, 0.002, 0.01, 0.02\}$, $n = 10$, $\rho = s = 5$; (b) $\mathbf{w} = 0.01$, $n = 10$, $\rho = 5$, $s \in \{4, 6, 8, 10\}$, and (c) $\mathbf{w} = 0.01$, $s = 10$, $\rho = 5$, $n \in \{4, 6, 8, 10\}$. We see that the regret degrades as \mathbf{w} , n , and s increase. We also see that when \mathbf{w} is large (T is small), the regret becomes worse quickly as n and s grow larger. These results are consistent with the theoretical bounds in Theorem 5.1.

8 Conclusions and Discussion

Markov jump systems are fundamental to a rich class of control problems where the underlying dynamics are changing with time. Despite its importance, statistical understanding (system identification and regret bounds) of MJS have been lacking due to the technicalities such as Markovian transitions and weaker notion of mean-square stability. At a high-level, this work overcomes (much of) these challenges to provide finite sample system identification and model-based adaptive control guarantees for MJS. Notably, resulting estimation error and regret bounds are optimal in the trajectory length and coincide with the standard LQR up to polylogarithmic factors. As a future work, it would be interesting and of practical importance to investigate the case when mode is not observed, which makes both system identification and adaptive quadratic control problems non-trivial.

We want to mention possible negative societal impacts. While our work is theoretical and has many

potential positive impacts in reinforcement learning, robotics, and autonomous systems, there are also potential negative applications in the military (e.g. with drone control) and for malicious actors (e.g. computer network hackers), among others. Additionally, all our work was built on stochastic noise assumptions, whereas in reality intelligent autonomous systems may instead encounter adversarial behavior. There is potential here for future work to extend our approach to non-stochastic noise or even non-Markovian / non-random switching among states.

Acknowledgements

Y. Sattar and S. Oymak were supported in part by NSF grant CNS-1932254 and S. Oymak was supported in part by NSF CAREER award CCF-2046816 and ARO MURI grant W911NF-21-1-0312. Z. Du and N. Ozay were supported in part by ONR under grant N00014-18-1-2501 and N. Ozay was supported in part by NSF under grant CNS-1931982 and ONR under grant N00014-21-1-2431. D. Ataee Tarzanagh, Z. Du, and L. Balzano were supported in part by NSF CAREER award CCF-1845076 and AFOSR YIP award FA9550-19-1-0026.

References

- [1] Yasin Abbasi-Yadkori, Nevena Lazic, and Csaba Szepesvári. Model-free linear quadratic control via reduction to expert prediction. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 3108–3117. PMLR, 2019.
- [2] Yasin Abbasi-Yadkori and Csaba Szepesvári. Regret bounds for the adaptive control of linear quadratic systems. In *Proc. of COLT*, pages 1–26. JMLR Workshop and Conference Proceedings, 2011.
- [3] Marc Abeille and Alessandro Lazaric. Improved regret bounds for thompson sampling in linear quadratic control problems. In *International Conference on Machine Learning*, pages 1–9. PMLR, 2018.
- [4] Marc Abeille and Alessandro Lazaric. Efficient optimistic exploration in linear-quadratic regulators via lagrangian relaxation. In *ICML*, pages 23–31. PMLR, 2020.
- [5] Sevi Baltaoglu, Lang Tong, and Qing Zhao. Online learning and optimization of markov jump affine models. *arXiv preprint arXiv:1605.02213*, 2016.
- [6] Lars Blackmore, Stanislav Funiak, and Brian C Williams. Combining stochastic and greedy search in hybrid estimation. In *AAAI*, pages 282–287, 2005.
- [7] Peter E Caines and Ji-Feng Zhang. On the adaptive control of jump parameter systems via nonlinear filtering. *SIAM J. Control Optim.*, 33(6):1758–1777, 1995.
- [8] DO Cajueiro. *Stochastic optimal control of jumping Markov parameter processes with applications to finance*. PhD thesis, PhD thesis, 2002, Instituto Tecnológico de Aeronáutica-ITA, Brazil, 2002.
- [9] Marco C Campi and PR Kumar. Adaptive linear quadratic gaussian control: the cost-biased approach revisited. *SIAM J. Control Optim.*, 36(6):1890–1907, 1998.
- [10] Asaf Cassel, Alon Cohen, and Tomer Koren. Logarithmic regret for learning linear quadratic regulators efficiently. In *International Conference on Machine Learning*, pages 1328–1337. PMLR, 2020.
- [11] Xinyi Chen and Elad Hazan. Black-box control for linear dynamical systems. In *Conference on Learning Theory*, pages 1114–1143. PMLR, 2021.
- [12] Howard J Chizeck, Alan S Willsky, and D Castanon. Discrete-time markovian-jump linear quadratic optimal control. *International Journal of Control*, 43(1):213–231, 1986.

- [13] Alon Cohen, Tomer Koren, and Yishay Mansour. Learning linear-quadratic regulators efficiently with only \sqrt{T} regret. In *International Conference on Machine Learning*, pages 1300–1309. PMLR, 2019.
- [14] Oswaldo Luiz Valle Costa, Marcelo Dutra Fragoso, and Ricardo Paulino Marques. *Discrete-time Markov jump linear systems*. Springer, 2006.
- [15] Sarah Dean, Horia Mania, Nikolai Matni, Benjamin Recht, and Stephen Tu. Regret bounds for robust adaptive control of the linear quadratic regulator. In *Advances in Neural Information Processing Systems*, pages 4188–4197, 2018.
- [16] Sarah Dean, Horia Mania, Nikolai Matni, Benjamin Recht, and Stephen Tu. On the sample complexity of the linear quadratic regulator. *FOCM*, pages 1–47, 2019.
- [17] Zhe Du, Necmiye Ozay, and Laura Balzano. Mode clustering for markov jump systems. *arXiv preprint arXiv:1910.02193*, 2019.
- [18] Zhe Du, Yahya Sattar, Davoud Ataee Tarzanagh, Laura Balzano, Samet Oymak, and Necmiye Ozay. Certainty equivalent quadratic control for markov jump systems. *arXiv preprint arXiv:2105.12358*, 2021.
- [19] Zhe Du, Yahya Sattar, Davoud Ataee Tarzanagh, Laura Balzano, Necmiye Ozay, and Samet Oymak. Data-driven control of markov jump systems: Sample complexity and regret bounds. https://www.dropbox.com/s/s0xp09Ixzgtsv4f/AdaptiveMJS_LQR_ACC.pdf, 2021.
- [20] Mohamad Kazem Shirani Faradonbeh, Ambuj Tewari, and George Michailidis. Finite-time adaptive stabilization of linear systems. *IEEE Transactions on Automatic Control*, 64(8):3498–3505, 2018.
- [21] Mohamad Kazem Shirani Faradonbeh, Ambuj Tewari, and George Michailidis. Finite time identification in unstable linear systems. *Automatica*, 96:342–353, 2018.
- [22] Mohamad Kazem Shirani Faradonbeh, Ambuj Tewari, and George Michailidis. On adaptive linear-quadratic regulators. *Automatica*, 117:108982, 2020.
- [23] Mohamad Kazem Shirani Faradonbeh, Ambuj Tewari, and George Michailidis. Optimism-based adaptive regulation of linear-quadratic systems. *IEEE Transactions on Automatic Control*, 2020.
- [24] 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.
- [25] Emily B Fox, Erik B Sudderth, Michael I Jordan, and Alan S Willsky. Bayesian nonparametric methods for learning markov switching processes. *IEEE Signal Processing Magazine*, 27(6):43–54, 2010.
- [26] David A Freedman. On tail probabilities for martingales. *the Annals of Probability*, pages 100–118, 1975.
- [27] Robert G Gallager. *Stochastic processes: theory for applications*. Cambridge University Press, 2013.
- [28] Joseph E Gaudio, Travis E Gibson, Anuradha M Annaswamy, Michael A Bolender, and Eugene Lavretsky. Connections between adaptive control and optimization in machine learning. In *2019 IEEE 58th Conference on Decision and Control (CDC)*, pages 4563–4568. IEEE, 2019.
- [29] Elad Hazan, Sham Kakade, and Karan Singh. The nonstochastic control problem. In *Algorithmic Learning Theory*, pages 408–421. PMLR, 2020.
- [30] Pedro Hespanhol and Anil Aswani. Statistical consistency of set-membership estimator for linear systems. *IEEE Control Systems Letters*, 4(3):668–673, 2020.
- [31] Daniel Hsu, Sham Kakade, Tong Zhang, et al. A tail inequality for quadratic forms of subgaussian random vectors. *Electronic Communications in Probability*, 17, 2012.

- [32] Morteza Ibrahimi, Adel Javanmard, and Benjamin Van Roy. Efficient reinforcement learning for high dimensional linear quadratic systems. In *NeurIPS*, pages 2645–2653, 2012.
- [33] Joao Paulo Jansch-Porto, Bin Hu, and Geir Dullerud. Policy learning of mdps with mixed continuous/discrete variables: A case study on model-free control of markovian jump systems. In *Learning for Dynamics and Control*, pages 947–957. PMLR, 2020.
- [34] Yassir Jedra and Alexandre Proutiere. Finite-time identification of stable linear systems optimality of the least-squares estimator. In *2020 59th IEEE Conference on Decision and Control (CDC)*, pages 996–1001. IEEE, 2020.
- [35] Vitaly Kuznetsov and Mehryar Mohri. Generalization bounds for non-stationary mixing processes. *Machine Learning*, 106(1):93–117, 2017.
- [36] Sahin Lale, Kamyar Azizzadenesheli, Babak Hassibi, and Anima Anandkumar. Explore more and improve regret in linear quadratic regulators. *arXiv preprint arXiv:2007.12291*, 2020.
- [37] Sahin Lale, Kamyar Azizzadenesheli, Babak Hassibi, and Anima Anandkumar. Logarithmic regret bound in partially observable linear dynamical systems. In *Advances in Neural Information Processing Systems*, 2020.
- [38] Sahin Lale, Oguzhan Teke, Babak Hassibi, and Anima Anandkumar. Stability and identification of random asynchronous linear time-invariant systems. In *Learning for Dynamics and Control*, pages 651–663. PMLR, 2021.
- [39] Andrew Lamperski. Computing stabilizing linear controllers via policy iteration. In *2020 59th IEEE Conference on Decision and Control (CDC)*, pages 1902–1907. IEEE, 2020.
- [40] F Lauer and G Bloch. Hybrid system identification: Theory and algorithms for learning switching models, vol. 478. *Cham, Switzerland: Springer*, 2018.
- [41] Bruce Lee and Andrew Lamperski. Non-asymptotic closed-loop system identification using autoregressive processes and hankel model reduction. In *2020 59th IEEE Conference on Decision and Control (CDC)*, pages 3419–3424. IEEE, 2020.
- [42] Ji-Woong Lee and Geir E. Dullerud. Uniform stabilization of discrete-time switched and markovian jump linear systems. *Automatica*, 42(2):205–218, 2006.
- [43] David A Levin and Yuval Peres. *Markov chains and mixing times*, volume 107. American Mathematical Soc., 2017.
- [44] Daniel Liberzon. *Switching in systems and control*. Springer Science & Business Media, 2003.
- [45] Lennart Ljung. System identification. *Wiley Encyclopedia of Electrical and Electronics Engineering*, pages 1–19, 1999.
- [46] KA Loparo and F Abdel-Malek. A probabilistic approach to dynamic power system security. *IEEE transactions on circuits and systems*, 37(6):787–798, 1990.
- [47] Horia Mania, Stephen Tu, and Benjamin Recht. Certainty equivalence is efficient for linear quadratic control. In *NeurIPS*, 2019.
- [48] Nikolai Matni, Alexandre Proutiere, Anders Rantzer, and Stephen Tu. From self-tuning regulators to reinforcement learning and back again. In *2019 IEEE 58th Conference on Decision and Control (CDC)*, pages 3724–3740. IEEE, 2019.
- [49] Hesameddin Mohammadi, Mahdi Soltanolkotabi, and Mihailo R Jovanović. On the linear convergence of random search for discrete-time lqr. *IEEE Control Systems Letters*, 5(3):989–994, 2020.

- [50] Mehryar Mohri and Afshin Rostamizadeh. Stability bounds for non-iid processes. In *Advances in Neural Information Processing Systems*, pages 1025–1032, 2008.
- [51] Samet Oymak. Stochastic gradient descent learns state equations with nonlinear activations. In *Conference on Learning Theory*, pages 2551–2579, 2019.
- [52] Samet Oymak and Necmiye Ozay. Non-asymptotic identification of lti systems from a single trajectory. *American Control Conference*, 2019.
- [53] Necmiye Ozay, Mario Sznaier, Constantino M Lagoa, and Octavia I Camps. A sparsification approach to set membership identification of switched affine systems. *IEEE Transactions on Automatic Control*, 57(3):634–648, 2011.
- [54] Benjamin Recht. A tour of reinforcement learning: The view from continuous control. *Annual Review of Control, Robotics, and Autonomous Systems*, 2:253–279, 2019.
- [55] Tuhin Sarkar and Alexander Rakhlin. Near optimal finite time identification of arbitrary linear dynamical systems. In *ICML*, pages 5610–5618. PMLR, 2019.
- [56] Tuhin Sarkar, Alexander Rakhlin, and Munther Dahleh. Nonparametric system identification of stochastic switched linear systems. In *2019 IEEE 58th Conference on Decision and Control (CDC)*, pages 3623–3628. IEEE, 2019.
- [57] Yahya Sattar and Samet Oymak. Non-asymptotic and accurate learning of nonlinear dynamical systems. *arXiv preprint arXiv:2002.08538*, 2020.
- [58] Max Simchowitz and Dylan Foster. Naive exploration is optimal for online lqr. In *ICML*, pages 8937–8948. PMLR, 2020.
- [59] Max Simchowitz, Horia Mania, Stephen Tu, Michael I Jordan, and Benjamin Recht. Learning without mixing: Towards a sharp analysis of linear system identification. In *Conference On Learning Theory*, pages 439–473. PMLR, 2018.
- [60] Lars EO Svensson, Noah Williams, et al. Optimal monetary policy under uncertainty: a markov jump-linear-quadratic approach. *Federal Reserve Bank of St. Louis Review*, 90(4):275–293, 2008.
- [61] Anastasios Tsiamis and George J Pappas. Finite sample analysis of stochastic system identification. In *2019 IEEE 58th Conference on Decision and Control (CDC)*, pages 3648–3654. IEEE, 2019.
- [62] Stephen Tu, Ross Boczar, Andrew Packard, and Benjamin Recht. Non-asymptotic analysis of robust control from coarse-grained identification. *arXiv preprint arXiv:1707.04791*, 2017.
- [63] Jitendra Tugnait. Adaptive estimation and identification for discrete systems with markov jump parameters. *IEEE Transactions on Automatic control*, 27(5):1054–1065, 1982.
- [64] Valery Ugrinovskii* and Hemanshu R Pota. Decentralized control of power systems via robust control of uncertain markov jump parameter systems. *International Journal of Control*, 78(9):662–677, 2005.
- [65] Roman Vershynin. *Introduction to the non-asymptotic analysis of random matrices*, page 210–268. Cambridge University Press, 2012.
- [66] F Xue and L Guo. Necessary and sufficient conditions for adaptive stabilizability of jump linear systems. *Communications in Information and Systems*, 1(2):205–224, 2001.
- [67] Anru Zhang and Mengdi Wang. Spectral state compression of markov processes. *IEEE transactions on information theory*, 66(5):3202–3231, 2019.

- [68] Kaiqing Zhang, Bin Hu, and Tamer Basar. Policy optimization for \mathcal{H}_2 linear control with \mathcal{H}_∞ robustness guarantee: Implicit regularization and global convergence. In *Learning for Dynamics and Control*, pages 179–190. PMLR, 2020.
- [69] Yang Zheng, Yujie Tang, and Na Li. Analysis of the optimization landscape of linear quadratic gaussian (lqg) control. *arXiv preprint arXiv:2102.04393*, 2021.

A Preliminaries

In addition to the notations defined in Section 1, we define a few more here to be used throughout the appendix. For a matrix \mathbf{V} , $\lambda_{\min}(\mathbf{V})$, $\|\mathbf{V}\|_1$, and $\|\mathbf{V}\|_F$ denote its smallest singular value, ℓ_1 norm and Frobenius norm, respectively. We use $\text{vec}(\mathbf{V})$ to denote the vectorization of a matrix \mathbf{V} and define $\|\mathbf{V}\|_+ := \|\mathbf{V}\|_1 + 1$. We define $\lambda_{\min}(\mathbf{V}_{1:S}) := \min_{i \in [S]} \lambda_{\min}(\mathbf{V}_i)$ and $\|\mathbf{V}_{1:S}\|_+ := \max_{i \in [S]} \|\mathbf{V}_i\|_+$. We use \mathbf{I}_n to denote the identity matrix of dimension n . $\mathbf{1}_n$ denotes the all 1 vector of dimension n and $\mathbf{1}_{\{\cdot\}}$ denotes the indicator function. Lastly, we use \lesssim and \gtrsim for inequalities that hold up to a constant factor.

To begin, we define the following quantity which will be used throughout to quantify the decay of a square matrix \mathbf{M} .

Definition A.1 For a square matrix \mathbf{M} with $\lambda_{\min}(\mathbf{M}) \leq 1$, we have

$$\tau(\mathbf{M}) := \sup_{k \in \mathbb{N}} \{\|\mathbf{M}^k\| / (\lambda_{\min}(\mathbf{M}))^k\}. \quad (\text{A.1})$$

Note that $\tau(\mathbf{M})$ is finite by Gelfand’s formula, and it is easy to see that $\tau(\mathbf{M}) \geq 1$. This quantity measures the transient response of a non-switching system with state matrix \mathbf{M} and can be upper bounded by its \mathcal{H}_∞ norm [62]. In this work, we will mainly use this quantity to evaluate the augmented state matrix for an MJS defined in Section 3.1.

For a Markov chain with transition matrix \mathbf{T} , we let $\mathbf{x}_0 \in \mathbb{R}^S$ denote the initial state distribution and \mathbf{x}_t denote the transient state distribution, i.e. $P(\mathbf{x}(t) = i) = \mathbf{x}_t(i)$. Then, it is easy to see $\mathbf{x}_t = \mathbf{x}_0 \mathbf{T}^t$. Note that \mathbf{x}_t is essentially a convex combination of rows of matrix \mathbf{T}^t , then by triangle inequality, we have $\|\mathbf{x}_t - \mathbf{x}_\infty\|_1 \leq \max_{i \in [S]} \|([\mathbf{T}^t]_i) - \mathbf{x}_\infty\|_1$. Thus, for an ergodic Markov matrix \mathbf{T} , we define the following to quantify the convergence of $\|\mathbf{x}_t - \mathbf{x}_\infty\|_1$.

Definition A.2 For an ergodic Markov matrix $\mathbf{T} \in \mathbb{R}^{S \times S}$, let $\tau_{MC} > 0$ and $\beta_{MC} \in [0, 1)$ be two constants [43, Theorem 4.9] such that

$$\max_{i \in [S]} \|([\mathbf{T}^t]_i) - \mathbf{x}_\infty\|_1 \leq \tau_{MC} \beta_{MC}^t. \quad (\text{A.2})$$

Furthermore, we define

$$t_{MC}(\cdot) := \min \left\{ t \in \mathbb{N} : \max_{i \in [S]} \frac{1}{2} \|([\mathbf{T}^t]_i) - \mathbf{x}_\infty\|_1 \leq \beta_{MC} \right\}. \quad (\text{A.3})$$

When parameter β_{MC} is omitted, it denotes $t_{MC} := t_{MC}(\frac{1}{4})$, i.e. the mixing time defined in Section 3.1.

Note that $\tau(\mathbf{M})$ and τ_{MC} have similar roles except $\tau(\mathbf{M})$ is usually used to study state matrices while τ_{MC} is for Markov matrices. For \mathbf{M} , we have $\|\mathbf{M}^k\| \leq \tau(\mathbf{M}) (\lambda_{\min}(\mathbf{M}))^k$, and for a Markov matrix $\|\mathbf{T}^t - \mathbf{1}_S\|_1 \leq \tau_{MC} \beta_{MC}^t$.

In this section, we define a few notations to ease the exposition in the appendix. Note that, for notations under parameterized form, i.e., notations which are functions of (\cdot, \cdot, \cdot) etc., one can choose these parameters freely to get different deterministic quantities.

Table 2 introduces notations and constants related to the choice of tuning parameter c_x, c_z , and the shortest trajectory (initial epoch) length such that theoretical performance guarantees can be achieved. Recall that $\mathbf{K}_{1:S}^{(0)}$ is the stabilizing controller for epoch 0 in Algorithm 2. We let $\mathbf{L}_i^{(0)} := \mathbf{A}_i + \mathbf{B}_i \mathbf{K}_i^{(0)}$, for all $i \in [S]$, denote the closed-loop state matrix, and $\tilde{\mathbf{L}}^{(0)} \in \mathbb{R}^{Sm^2 \times Sm^2}$ denotes the augmented closed-loop state matrix

Table 2: Notations — Tuning Parameters and Trajectory Length

\bar{z} (depending on context)	\mathbf{z} or \mathbf{z}_0 or $\sqrt{\ \mathbf{z}\ }$
$\bar{\mathbf{w}}$ (depending on context)	\mathbf{w} or $\sqrt{\ \mathbf{w}\ }$
C_z	$\bar{z}/\bar{\mathbf{w}}$
-2	$\ \mathbf{B}_{1s}\ ^{2-2} \frac{\bar{z}}{\bar{\mathbf{w}}} + \frac{-2}{\bar{\mathbf{w}}}$
$\underline{C}_x(\cdot, \cdot)$	$3\sqrt{\frac{18n\bar{s}^{-2}}{\min\{2, (1-\cdot)\}}}$
\underline{C}_z	$\max\left\{(\sqrt{3} + \sqrt{6})\sqrt{\bar{\rho}}, \sqrt{3\log\left(\frac{6}{\min}\right)}\right\}$
$-$	$\max\{\tilde{\mathbf{L}}^{(0)}, \tilde{\mathbf{L}}\}$
$-$	$\max\{\tilde{\mathbf{L}}^{(0)}, \frac{1+^*}{2}\}$
C_{MC}	$t_{MC} \cdot \max\{3, 3 - 3\log(\max \log(s))\}$
$\underline{I}_{MC,1}(\mathcal{C}, \cdot)$	$(68\mathcal{C} \max_{\min}^{-2} \log(\frac{\mathcal{C}}{s}))^2$
$\underline{I}_{MC}(\mathcal{C}, \cdot)$	$(612\mathcal{C} \max_{\min}^{-2} \log(\frac{2s}{\mathcal{C}}))^2$
$\underline{I}_{cl,1}(\cdot, \cdot)$	$\frac{(1-\cdot)^2}{4n^{1.5}\bar{s}^{-4}}$
$\underline{I}_N(\mathcal{C}, \cdot, \cdot, \cdot)$	$\max\{\underline{I}_{MC}(\mathcal{C}, \frac{\cdot}{2}), \underline{I}_{cl,1}(\cdot, \cdot)\}$
$\underline{I}_{id}(\mathcal{C}, \cdot, T, \cdot, \cdot)$	$\frac{\bar{s}\mathcal{C}\log(T)}{\min(1-\cdot)} \left((\sqrt{2\log(nT)} + \sqrt{2\log(2/\cdot)})^2 + C_z^2 \ \mathbf{B}_{1s}\ ^2 \log\left(\frac{36s(n+p)\mathcal{C}\log(T)}{\cdot}\right)(n+p) \right)$
$\underline{I}_{id,N}(L, \cdot, T, \cdot, \cdot)$	$\max\left\{\underline{I}_{id}\left(\frac{L}{\log(T)}, \frac{L}{2L}, T, \cdot, \cdot\right), \underline{I}_N\left(\frac{L}{\log(T)}, \frac{L}{2L}, \cdot, \cdot\right)\right\}$
$\underline{I}_{rgt,-}(\cdot, T)$	$\mathcal{O}\left(\frac{\bar{s}(n+p)}{\min}^{-4} \mathbf{A}, \mathbf{B}, \mathbf{T} \log(1) \log^4(T)\right)$
	$\mathcal{O}\left(\frac{\bar{s}(n+p)}{\min}^{-2} \mathbf{A}, \mathbf{B}, \mathbf{T} \log(1) \log^2(T)\right)$ (when \mathbf{B}_{1s} is known)
$\underline{I}_{\mathbf{x}_0}(\cdot)$	$\frac{1}{\log(1-\cdot)} \max\left\{\frac{2}{\log(\cdot)}, \log\left(\frac{2}{3} \frac{\bar{s}}{\cdot}\right)\right\}$
$\underline{I}_{rgt}(\cdot, T)$	$\max\{\underline{I}_{\mathbf{x}_0}(\cdot), \underline{I}_{rgt,-}(\cdot, T), \underline{I}_{MC,1}(\cdot), \underline{I}_{id,N}(\underline{L}, \cdot, T, \cdot, \cdot)\}$

with ij -th $n^2 \times n^2$ block given by $[\tilde{\mathbf{L}}^{(0)}]_{ij} = [\mathbf{T}]_{ji} \mathbf{L}_j^{(0)} \otimes \mathbf{L}_i^{(0)}$. (\cdot) is as in Definition A.1 and $\rho(\cdot)$ denotes the spectral radius. For the infinite-horizon MJS-LQR $(\mathbf{A}_{1s}, \mathbf{B}_{1s}, \mathbf{T}, \mathbf{Q}_{1s}, \mathbf{R}_{1s})$ problem, we let \mathbf{P}_{1s} denote the solution to cDARE given by (5.2) and \mathbf{K}_{1s} denotes the optimal controller which can be computed via (5.1) with \mathbf{P}_{1s} . Similarly, we define \mathbf{L}_{1s} and $\tilde{\mathbf{L}}$ to be the corresponding closed-loop state matrix and augmented closed-loop state matrix respectively and $\bar{\rho} := \rho(\tilde{\mathbf{L}})$. \max and \min are the largest and smallest elements in the stationary distribution of the ergodic Markov matrix \mathbf{T} . For the definition of $\underline{I}_{rgt,-}(\cdot, T)$, notation $\bar{\mathbf{A}}, \bar{\mathbf{B}}, \bar{\mathbf{T}}$ is defined in Table 3. As a slight abuse of notation, T in $\underline{I}_{rgt,-}(\cdot, T)$ (as well as Table 4) and \mathcal{C} are merely arguments to be replaced with specific quantities depending on the context.

Table 3 lists the notations related to infinite-horizon MJS perturbation results closely following the notations in [18]. It provides several sensitivity parameters, e.g., how the optimal controller \mathbf{K}_{1s} varies with perturbations in the MJS parameters \mathbf{A}_{1s} , \mathbf{B}_{1s} , and \mathbf{T} and how the MJS-LQR cost \mathcal{J} varies with the controller \mathbf{K}_{1s} . It also provides certain upper bounds on the variations in \mathbf{A}_{1s} , \mathbf{B}_{1s} , \mathbf{T} , and \mathbf{K}_{1s} such that the perturbation theory holds. In this table, $\mathbf{R}_{1s}^{-1} := \{\mathbf{R}_i^{-1}\}_{i=1}^S$ and recall $\|\cdot\|_+ := \|\cdot\| + 1$.

A.1 MJS Covariance Dynamics Under MSS

Consider MJS $(\mathbf{A}_{1s}, \mathbf{B}_{1s}, \mathbf{T})$ with process noise $\mathbf{w}_t \sim \mathcal{N}(0, \bar{\mathbf{w}})$ and input $\mathbf{u}_t = \mathbf{K}_{(t)} \mathbf{x}_t + \mathbf{z}_t$ under a stabilizing controller \mathbf{K}_{1s} and excitation for exploration $\mathbf{z}_t \sim \mathcal{N}(0, \bar{\mathbf{z}})$. Let $\mathbf{L}_i := \mathbf{A}_i + \mathbf{B}_i \mathbf{K}_i$ be the closed-loop state matrix. Let $\tilde{\mathbf{L}} \in \mathbb{R}^{sn^2 \times sn^2}$ be the augmented closed-loop state matrix with ij -th $n^2 \times n^2$ block given by $[\tilde{\mathbf{L}}]_{ij} = [\mathbf{T}]_{ji} \mathbf{L}_j \otimes \mathbf{L}_i$. Let $\bar{\zeta} > 0$ and $\underline{\zeta} \in [0, 1)$ be two constants such that $\|\tilde{\mathbf{L}}^k\| \leq \bar{\zeta} \frac{\bar{\zeta}}{\underline{\zeta}}$. By definitions of $\tilde{\mathbf{L}}$

Table 3: Notations — MJS-LQR Perturbation

	$\min\{\ \mathbf{B}_{1s}\ _+^{-2}\ \mathbf{R}_{1s}^{-1}\ _+^{-1}\ \mathbf{L}_{1s}\ _+^{-2}, _-(\mathbf{P}_{1s})\}$
	$\ \mathbf{A}_{1s}\ _+^2\ \mathbf{B}_{1s}\ _+^4\ \mathbf{P}_{1s}\ _+^3\ \mathbf{R}_{1s}^{-1}\ _+^2$
Γ	$\max\{\ \mathbf{A}_{1s}\ _+, \ \mathbf{B}_{1s}\ _+, \ \mathbf{P}_{1s}\ _+, \ \mathbf{K}_{1s}\ _+\}$
$C_{\mathbf{A},\mathbf{B},\mathbf{T}}^{\mathbf{K}}$	$28\sqrt{ns}(\tilde{\mathbf{L}})(1 -)^{-1}(_-(\mathbf{R}_{1s})^{-1} + \Gamma^3_-(\mathbf{R}_{1s})^{-2})\Gamma^3$
$C_{\mathbf{K}}^{\mathbf{J}}$	$2s^{1.5}\sqrt{n}\min\{n, p\}(\ \mathbf{R}_{1s}\ + \Gamma^3)\frac{(\tilde{\mathbf{L}}^*)}{1 - _*$
$_-\mathbf{K}$	$\min\left\{\ \mathbf{K}_{1s}\ , \frac{1 - _*$
$_-\text{LQR}$ $\mathbf{A}, \mathbf{B}, \mathbf{T}$	$\frac{(1 - _*)\min\{_*, _-(\mathbf{R}_{1:s})^2 _-\mathbf{K}\}}{28\sqrt{ns}(\tilde{\mathbf{L}}^*)^3(_-(\mathbf{R}_{1:s}) + _*)} - 1$
$_-\mathbf{A}, \mathbf{B}, \mathbf{T}$	$\min\left\{\frac{(1 - _*)^2}{204ns(\tilde{\mathbf{L}}^*)^2}, \ \mathbf{B}_{1s}\ , _-(\mathbf{Q}_{1s}), \frac{_-\text{LQR}}{\mathbf{A}, \mathbf{B}, \mathbf{T}}\right\}$

and $(\tilde{\mathbf{L}})$, one can choose them for $\tilde{\mathbf{L}}$ and $\tilde{\mathbf{L}}$ respectively. Let $_i(t) := \mathbb{E}[\mathbf{x}_t\mathbf{x}_t^T \mathbf{1}_{\{(t)=i\}}]$, $_t(t) := \mathbb{E}[\mathbf{x}_t\mathbf{x}_t^T]$,

$$\mathbf{s}_t := \begin{bmatrix} \text{vec}(_1(t)) \\ \vdots \\ \text{vec}(_s(t)) \end{bmatrix}, \tilde{\mathbf{B}}_t := \begin{bmatrix} \sum_{j=1}^S _t(j)\mathbf{T}_{j1}(\mathbf{B}_j \otimes \mathbf{B}_j) \\ \vdots \\ \sum_{j=1}^S _t(j)\mathbf{T}_{js}(\mathbf{B}_j \otimes \mathbf{B}_j) \end{bmatrix}, \text{ and } \tilde{\mathbf{L}}_t := _t \otimes \mathbf{I}_{n^2}. \quad (\text{A.4})$$

The following lemma shows how \mathbf{s}_t depends on \mathbf{s}_0 , \mathbf{z} , and \mathbf{w} , which will be used to upper bound $\mathbb{E}[\|\mathbf{x}_t\|^2]$ in Lemma A.4.

Lemma A.3 *The vectorized covariance \mathbf{s}_t has the following dynamics,*

$$\mathbf{s}_t = \tilde{\mathbf{L}}^t \mathbf{s}_0 + (\tilde{\mathbf{B}}_t + \tilde{\mathbf{L}}\tilde{\mathbf{B}}_{t-1} + \dots + \tilde{\mathbf{L}}^{t-1}\tilde{\mathbf{B}}_1)\text{vec}(\mathbf{z}) + (\tilde{\mathbf{L}}_t + \tilde{\mathbf{L}}\tilde{\mathbf{L}}_{t-1} + \dots + \tilde{\mathbf{L}}^{t-1}\tilde{\mathbf{L}}_1)\text{vec}(\mathbf{w}).$$

Proof To begin, we evaluate $_i(t)$, from the equivalent MJS dynamics $\mathbf{x}_{t+1} = \mathbf{L}(_t)\mathbf{x}_t + \mathbf{B}(_t)\mathbf{z}_t + \mathbf{w}_t$, as follows,

$$\begin{aligned} \mathbb{E}[\mathbf{x}_{t+1}\mathbf{x}_{t+1}^T \mathbf{1}_{\{(t+1)=i\}}] &= \sum_{j=1}^S \mathbb{E}[\mathbf{L}_j\mathbf{x}_t\mathbf{x}_t^T\mathbf{L}_j^T \mathbf{1}_{\{(t+1)=i, (t)=j\}}] \\ &+ \sum_{j=1}^S \mathbb{E}[\mathbf{B}_j\mathbf{z}_t\mathbf{z}_t^T\mathbf{B}_j^T \mathbf{1}_{\{(t+1)=i, (t)=j\}}] + \mathbb{E}[\mathbf{w}_t\mathbf{w}_t^T \mathbf{1}_{\{(t+1)=i\}}]. \end{aligned} \quad (\text{A.5})$$

Since $\mathbf{w}_t \sim \mathcal{N}(0, \mathbf{w})$ and $\mathbf{z}_t \sim \mathcal{N}(0, \mathbf{z})$, we get

$$_i(t+1) = \sum_{j=1}^S \mathbf{T}_{ji}\mathbf{L}_j _j(t)\mathbf{L}_j + \sum_{j=1}^S _t(j)\mathbf{T}_{ji}\mathbf{B}_j \mathbf{z}\mathbf{B}_j + _t+1(i)\mathbf{w}. \quad (\text{A.6})$$

Vectorizing both sides of the above equation, we have

$$\begin{aligned} \text{vec}(_i(t+1)) &= \sum_{j=1}^S \mathbf{T}_{ji}(\mathbf{L}_j \otimes \mathbf{L}_j)\text{vec}(_j(t)) \\ &+ \sum_{j=1}^S _t(j)\mathbf{T}_{ji}(\mathbf{B}_j \otimes \mathbf{B}_j)\text{vec}(\mathbf{z}) + _t+1(i)\text{vec}(\mathbf{w}). \end{aligned}$$

Stacking this for every $i \in [s]$, we obtain

$$\begin{bmatrix} \text{vec}(_1(t+1)) \\ \vdots \\ \text{vec}(_s(t+1)) \end{bmatrix} = \tilde{\mathbf{L}} \begin{bmatrix} \text{vec}(_1(t)) \\ \vdots \\ \text{vec}(_s(t)) \end{bmatrix} + \tilde{\mathbf{B}}_{t+1}\text{vec}(\mathbf{z}) + \tilde{\mathbf{L}}_{t+1}\text{vec}(\mathbf{w}). \quad (\text{A.7})$$

Propagating this dynamics from t to 0 gives the desired result. \blacksquare

We next provide a key lemma that upper bounds $E[\|\mathbf{x}_t\|^2]$ and $\| \cdot (t) \|_F$, which are later used extensively in system identification analysis.

Lemma A.4 For $E[\|\mathbf{x}_t\|^2]$ and $\| \cdot (t) \|_F$, under MSS given in Definition 3.1, we have

$$E[\|\mathbf{x}_t\|^2] \leq \sqrt{ns} \cdot \tilde{\mathbf{L}}^t \cdot E[\|\mathbf{x}_0\|^2] + n\sqrt{s}(\|\mathbf{B}_{1s}\|^2 \| \mathbf{z} \| + \| \mathbf{w} \|) \frac{\tilde{\mathbf{L}}}{1 - \tilde{\mathbf{L}}}, \quad (\text{A.8a})$$

$$\| \cdot (t) \|_F \leq \sqrt{s} \cdot \tilde{\mathbf{L}}^t \cdot E[\|\mathbf{x}_0\|^2] + \sqrt{ns}(\|\mathbf{B}_{1s}\|^2 \| \mathbf{z} \| + \| \mathbf{w} \|) \frac{\tilde{\mathbf{L}}}{1 - \tilde{\mathbf{L}}}. \quad (\text{A.8b})$$

Proof First we derive an upper bound for $E[\|\mathbf{x}_t\|^2]$. The upper bound for $\| \cdot (t) \|_F$ follows similarly. For state \mathbf{x}_t , we have

$$\begin{aligned} E[\|\mathbf{x}_t\|^2] &= \sum_{i=1}^S E[\|\mathbf{x}_t\|^2 \mathbf{1}_{\{(t)=i\}}] = \sum_{i=1}^S \text{tr} (E[\mathbf{x}_t \mathbf{x}_t^T \mathbf{1}_{\{(t)=i\}}]) = \sum_{i=1}^S \text{tr} (\cdot (t)) \\ &= \sum_{i=1}^S \sum_{j=1}^n \lambda_j (\cdot (t)) \leq \sqrt{ns} \sqrt{\sum_{i=1}^S \sum_{j=1}^n \lambda_j^2 (\cdot (t))} \\ &\leq \sqrt{ns} \sqrt{\sum_{i=1}^S \| \cdot (t) \|_F^2}. \end{aligned}$$

Then, by definition of \mathbf{s}_t in (A.4), we have

$$E[\|\mathbf{x}_t\|^2] \leq \sqrt{ns} \|\mathbf{s}_t\|. \quad (\text{A.9})$$

Now, applying the dynamics of \mathbf{s}_t from Lemma A.3, we have

$$\begin{aligned} E[\|\mathbf{x}_t\|^2] &\leq \sqrt{ns} \left(\|\tilde{\mathbf{L}}^t\| \|\mathbf{s}_0\| + \sum_{t'=1}^t \|\tilde{\mathbf{L}}^{t-t'}\| \|\tilde{\mathbf{B}}_{t'} \text{vec}(\mathbf{z})\| + \sum_{t'=1}^t \|\tilde{\mathbf{L}}^{t-t'}\| \|\tilde{\mathbf{w}}_{t'} \text{vec}(\mathbf{w})\| \right) \\ &\leq \sqrt{ns} \tilde{\mathbf{L}}^t \left(\|\mathbf{s}_0\| + \sum_{t'=1}^t \tilde{\mathbf{L}}^{t-t'} \|\tilde{\mathbf{B}}_{t'} \text{vec}(\mathbf{z})\| + \sum_{t'=1}^t \tilde{\mathbf{L}}^{t-t'} \|\tilde{\mathbf{w}}_{t'} \text{vec}(\mathbf{w})\| \right), \end{aligned} \quad (\text{A.10})$$

where the second line follows from $\|\tilde{\mathbf{L}}^t\| \leq \tilde{\mathbf{L}}^t$.

Now, we evaluate $\|\mathbf{s}_0\|$, $\|\tilde{\mathbf{B}}_{t'} \text{vec}(\mathbf{z})\|$, and $\|\tilde{\mathbf{w}}_{t'} \text{vec}(\mathbf{w})\|$ separately. For the first term, we have

$$\|\mathbf{s}_0\| = \sqrt{\sum_{i=1}^S \| \cdot (0) \|_F^2} = \sqrt{\sum_{i=1}^S \lambda_i(0)^2 E[\|\mathbf{x}_0 \mathbf{x}_0^T\|_F^2]} \leq E[\|\mathbf{x}_0 \mathbf{x}_0^T\|_F] \leq E[\|\mathbf{x}_0\|^2]. \quad (\text{A.11})$$

Let $[\tilde{\mathbf{B}}_{t'}]_i$ denote the i th block of $\tilde{\mathbf{B}}_{t'}$, i.e., $[\tilde{\mathbf{B}}_{t'}]_i = \sum_{j=1}^S \mathbf{1}_{t-1}(j) \mathbf{T}_{ji} (\mathbf{B}_j \otimes \mathbf{B}_j)$, then

$$\begin{aligned} \|\tilde{\mathbf{B}}_{t'} \text{vec}(\mathbf{z})\| &= \sqrt{\sum_{i=1}^S \| [\tilde{\mathbf{B}}_{t'}]_i \text{vec}(\mathbf{z}) \|^2} \leq \sum_{i=1}^S \| [\tilde{\mathbf{B}}_{t'}]_i \text{vec}(\mathbf{z}) \| \\ &= \sum_{i=1}^S \left\| \sum_{j=1}^S \mathbf{1}_{t-1}(j) \mathbf{T}_{ji} (\mathbf{B}_j \otimes \mathbf{B}_j) \text{vec}(\mathbf{z}) \right\| \\ &= \sum_{i=1}^S \left\| \sum_{j=1}^S \mathbf{1}_{t-1}(j) \mathbf{T}_{ji} (\mathbf{B}_j \mathbf{z} \mathbf{B}_j) \right\|_F \\ &\leq \|\mathbf{B}_{1s}\|^2 \| \mathbf{z} \| \cdot \sum_{i=1}^S \left\| \sum_{j=1}^S \mathbf{1}_{t-1}(j) \mathbf{T}_{ji} \mathbf{I}_n \right\|_F \\ &= \|\mathbf{B}_{1s}\|^2 \| \mathbf{z} \| \cdot \sum_{i=1}^S \| \cdot (i) \mathbf{I}_n \|_F \\ &\leq \sqrt{n} \|\mathbf{B}_{1s}\|^2 \| \mathbf{z} \|. \end{aligned} \quad (\text{A.12})$$

Lastly, we have

$$\|\tilde{\mathbf{r}} \text{vec}(\mathbf{w})\| = \sqrt{\sum_{i=1}^S \|\tilde{\mathbf{r}}(i) \text{vec}(\mathbf{w})\|^2} \leq \|\text{vec}(\mathbf{w})\| = \|\mathbf{w}\|_F = \sqrt{\bar{n}} \|\mathbf{w}\|. \quad (\text{A.13})$$

Plugging (A.11)–(A.13) into (A.10), we obtain

$$\begin{aligned} \mathbb{E}[\|\mathbf{x}_t\|^2] &\leq \sqrt{\bar{n}S} \bar{\zeta} \left(\bar{\zeta} \mathbb{E}[\|\mathbf{x}_0\|^2] + \sqrt{\bar{n}} \|\mathbf{B}_{1s}\|^2 \|\mathbf{z}\| \sum_{t'=1}^t \bar{\zeta}^{t-t'} + \sqrt{\bar{n}} \|\mathbf{w}\| \sum_{t'=1}^t \bar{\zeta}^{t-t'} \right), \\ &\leq \sqrt{\bar{n}S} \cdot \bar{\zeta} \cdot \bar{\zeta} \cdot \mathbb{E}[\|\mathbf{x}_0\|^2] + n\sqrt{S}(\|\mathbf{B}_{1s}\|^2 \|\mathbf{z}\| + \|\mathbf{w}\|) \frac{\bar{\zeta}}{1 - \bar{\zeta}}, \end{aligned} \quad (\text{A.14})$$

which gives the bound for $\mathbb{E}[\|\mathbf{x}_t\|^2]$ in (A.8a).

To obtain the bound for $\|\mathbf{x}(t)\|_F$ in (A.8b), note that $\|\mathbf{x}(t)\|_F = \|\sum_{i=1}^S \mathbf{x}_i(t)\|_F \leq \sqrt{S} \sqrt{\sum_{i=1}^S \|\mathbf{x}_i(t)\|_F^2} \leq \sqrt{S} \|\mathbf{s}_t\|$. We then follow a similar line of reasoning as above to get the statement of the lemma. This completes the proof. \blacksquare

A.2 Supporting Lemmas

In this section, we provide a list of lemmas that will be useful for the subsequent proofs.

Lemma A.5 Suppose $\mathbf{z} \sim \mathcal{N}(0, \Sigma_{\mathbf{z}})$ with $\Sigma_{\mathbf{z}} \in \mathbb{R}^{p \times p}$. For any $t \geq (3 + 2\sqrt{2})p$, we have

$$\mathbb{P}(\|\mathbf{z}\|^2 \geq 3\|\Sigma_{\mathbf{z}}\|t) \leq e^{-t}.$$

Proof From [31, Proposition 1], we have for any $t > 0$,

$$\mathbb{P}(\|\mathbf{z}\|^2 \geq \text{tr}(\Sigma_{\mathbf{z}}) + 2\sqrt{\text{tr}(\Sigma_{\mathbf{z}})}t + 2\|\Sigma_{\mathbf{z}}\|t) \leq e^{-t},$$

which implies

$$\mathbb{P}(\|\mathbf{z}\|^2 \geq p\|\Sigma_{\mathbf{z}}\| + 2\sqrt{p}\|\Sigma_{\mathbf{z}}\|\sqrt{t} + 2\|\Sigma_{\mathbf{z}}\|t) \leq e^{-t}.$$

We can see that when $t \geq (3 + 2\sqrt{2})p$, we have $p + 2\sqrt{p}\sqrt{t} \leq t$, which implies $p\|\Sigma_{\mathbf{z}}\| + 2\sqrt{p}\|\Sigma_{\mathbf{z}}\|\sqrt{t} \leq \|\Sigma_{\mathbf{z}}\|t$. Therefore, we have

$$\mathbb{P}(\|\mathbf{z}\|^2 \geq 3\|\Sigma_{\mathbf{z}}\|t) \leq e^{-t}. \quad \blacksquare$$

Lemma A.6 Let \mathbf{x}_t be the MJS state and define the noise-removed state $\tilde{\mathbf{x}}_t = \mathbf{x}_t - \mathbf{w}_{t-1}$ which is independent of \mathbf{w}_{t-1} . Let $\tilde{\mathbf{x}}_t$ be zero mean with $\mathbb{E}[\|\tilde{\mathbf{x}}_t\|] \leq B$ and \mathbf{w}_t has i.i.d. entries with variance $\frac{2}{c_{\mathbf{w}}}$ bounded in absolute value by $c_{\mathbf{w}}$ for some $c_{\mathbf{w}} > 0$. Consider the conditional random vector

$$\mathbf{y}_t \sim \{\mathbf{x}_t \mid \|\mathbf{x}_t\| \leq 3B\}.$$

If $c_{\mathbf{w}} \sqrt{\bar{n}} \leq B$, then $\text{Cov}[\mathbf{y}_t \mathbf{y}_t] \geq \frac{2}{c_{\mathbf{w}}} \mathbf{I}_n / 2$.

Proof Observe that $\|\mathbf{w}_t\| \leq c_{\mathbf{w}} \sqrt{\bar{n}} \leq B$. Define the events

$$E_1 = \{\|\mathbf{x}_t\| \leq 3B\}, \quad E_2 = \{\|\tilde{\mathbf{x}}_t\| \leq 2B\}.$$

Clearly $E_2 \subset E_1$ as $\|\mathbf{w}_t\| \leq B$. Now, observe that

$$\begin{aligned} \text{Cov}[\mathbf{y}_t \mathbf{y}_t] &= \text{Cov}[\mathbf{y}_t \mathbf{y}_t \mid E_2] \mathbb{P}(E_2 \mid E_1) \\ &\geq \text{Cov}[\mathbf{y}_t \mathbf{y}_t \mid E_2] \mathbb{P}(E_2). \end{aligned}$$

Note that $P(E_2) \geq 0.5$ from Markov bound as $E[\|\tilde{\mathbf{x}}_t\|] \leq B$. Additionally, on the event E_2 , $\tilde{\mathbf{x}}_t$ and \mathbf{w}_{t-1} are independent. Thus, we further have

$$\begin{aligned} \text{Cov}[\mathbf{y}_t \mathbf{y}_t] &\geq \text{Cov}[\mathbf{y}_t \mathbf{y}_t \mid E_2] P(E_2) \\ &\geq \text{Cov}[\mathbf{w}_{t-1} \mathbf{w}_{t-1} \mid E_2] P(E_2) \\ &\geq 0.5 \cdot \text{Cov}[\mathbf{w}_{t-1} \mathbf{w}_{t-1}] \\ &\geq \frac{2}{\mathbf{w}} \mathbf{I}_n / 2. \end{aligned}$$

■

Lemma A.7 Let $\mathbf{z} \sim \mathcal{N}(0, \frac{2}{\mathbf{z}} \mathbf{I}_p)$. Consider the conditional random vector $\mathbf{y} \sim \{\mathbf{z} \mid \|\mathbf{z}\| \leq c \sqrt{\mathbf{z}} \sqrt{\rho}\}$, where $c \geq 6$ is a fixed constant. Then $\text{Cov}[\mathbf{y} \mathbf{y}] \geq \frac{2}{\mathbf{z}} \mathbf{I}_p / 2$.

Proof This proof gives a lower bound on the covariance of truncated Gaussian vector $\mathbf{z} \mid \|\mathbf{z}\| \leq c \sqrt{\mathbf{z}} \sqrt{\rho}$. Note that, $\mathbf{z} = \mathbf{z} / \sqrt{\mathbf{z}}$ is $\mathcal{N}(0, \mathbf{I}_p)$. Set variable $X = \|\mathbf{z}\|^2$. We have the following Lipschitz Gaussian tail bound (we use Lipschitzness of the $\sqrt{\cdot}$ norm and use minor calculus and relaxations)

$$P(\|\mathbf{z}\|^2 \geq 4tp) \leq \begin{cases} 1 & \text{if } t \leq 1 \\ e^{-tp/2} & \text{if } t \geq 1. \end{cases}$$

This implies the following tail bound for X

$$Q(t) = P(X \geq t) \leq \begin{cases} 1 & \text{if } t \leq 4p \\ e^{-t/8} & \text{if } t \geq 4p. \end{cases}$$

Fix $\rho \geq 4$. Using integration-by-parts, this implies that

$$\begin{aligned} E[X \mid X \geq \rho] P(X \geq \rho) &= - \int_{\rho}^{\infty} x dQ(x) = -[xQ(x)]_{\rho}^{\infty} + \int_{\rho}^{\infty} Q(x) dx, \\ &\leq (\rho + 8) e^{-\rho/8}. \end{aligned} \tag{A.15}$$

The final line as a function of ρ is decreasing when $\rho \geq 36$. Specifically it is upper bounded by $1/2$ when $\rho \geq 36$ (as $\rho \geq 1$). Now define the event

$$E_z = \{\|\mathbf{z}\| \leq \sqrt{\mathbf{z}} \sqrt{\rho}\}.$$

For $\sqrt{\rho} \geq 6$, $\sqrt{\cdot}$ will map to the c in the statement of the lemma. Observe that this is also the event $X \leq \rho$. Following (A.15), this implies

$$\begin{aligned} E[\|\mathbf{z}\|^2 \mid E_z^c] P(E_z^c) &\leq E[\frac{2}{\mathbf{z}} X \mid E_z^c] P(E_z^c) \\ &\leq \frac{2}{\mathbf{z}} E[X \mid E_z^c] P(E_z^c) \\ &\leq \frac{2}{\mathbf{z}} / 2. \end{aligned}$$

This also yields the covariance bound of the tail event

$$E[\mathbf{z} \mathbf{z} \mid E_z^c] P(E_z^c) \leq E[\|\mathbf{z}\|^2 \mathbf{I}_p \mid E_z^c] P(E_z^c) \leq \frac{2}{\mathbf{z}} \mathbf{I}_p / 2.$$

Finally, from the conditional decomposition, observe that

$$E[\mathbf{z} \mathbf{z}] = E[\mathbf{z} \mathbf{z} \mid E_z^c] P(E_z^c) + E[\mathbf{z} \mathbf{z} \mid E_z] P(E_z) \implies E[\mathbf{z} \mathbf{z} \mid E_z] P(E_z) \geq \frac{2}{\mathbf{z}} \mathbf{I}_p / 2.$$

To conclude, observe that $E[\mathbf{z} \mathbf{z} \mid E_z] = E[\mathbf{y} \mathbf{y}]$, where \mathbf{y} is the conditional vector defined by truncating \mathbf{z} . Thus, we found

$$E[\mathbf{y} \mathbf{y}] P(E_z) \geq \frac{2}{\mathbf{z}} \mathbf{I}_p / 2 \implies E[\mathbf{y} \mathbf{y}] \geq \frac{2}{\mathbf{z}} \mathbf{I}_p / 2.$$

■

Theorem A.8 [65, Theorem 5.41 (Isotropic)] Let \mathbf{X} be an $N \times d$ matrix whose rows $\mathbf{x}_i \in \mathbb{R}^d$ are independent isotropic. Let m be such that $\|\mathbf{x}_i\| \leq \sqrt{m}$ almost surely for all $i \in [N]$. Then, for every $t \geq 0$, with probability $1 - 2d \cdot e^{-ct^2}$, we have

$$\sqrt{N} - t\sqrt{m} \leq \lambda_{\min}(\mathbf{X}) \leq \|\mathbf{X}\| \leq \sqrt{N} + t\sqrt{m}.$$

Corollary A.9 (Non-isotropic) Let \mathbf{X} be an $N \times d$ matrix whose rows $\mathbf{x}_i \in \mathbb{R}^d$ are independent with covariance Σ_i . Suppose each covariance obeys

$$\lambda_{\min}^2(\Sigma_i) \leq \|\Sigma_i\| \leq \lambda_{\max}^2.$$

Let m be such that $\|\mathbf{x}_i\| \leq \sqrt{m}$ almost surely for all $i \in [N]$. Then, for every $t \geq 0$, with probability $1 - 2d \cdot e^{-ct^2}$, we have

$$\lambda_{\min} \sqrt{N} - t\sqrt{m} \leq \lambda_{\min}(\mathbf{X}) \leq \|\mathbf{X}\| \leq \lambda_{\max} \sqrt{N} + t \frac{\lambda_{\max}}{\lambda_{\min}} \sqrt{m}.$$

Proof Let $\mathbf{x}_i = \lambda_{\min}^{-1/2} \mathbf{z}_i$. Observe that \mathbf{z}_i are independent isotropic. Define the matrix \mathbf{Z} with rows \mathbf{z}_i . Note that $\|\mathbf{x}_i\| \leq \|\mathbf{z}_i\| / \lambda_{\min} \leq \lambda_{\min}^{-1/2} \sqrt{m}$. Thus, applying Theorem A.8 on \mathbf{Z} , for every $t \geq 0$, with probability $1 - 2d \cdot e^{-ct^2}$, we have

$$\sqrt{N} - t \lambda_{\min}^{-1/2} \sqrt{m} \leq \lambda_{\min}(\mathbf{Z}) \leq \|\mathbf{Z}\| \leq \sqrt{N} + t \lambda_{\min}^{-1/2} \sqrt{m}. \quad (\text{A.16})$$

Next, observing that $\mathbf{X} \mathbf{X}^T = \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T = \sum_{i=1}^N \sqrt{\lambda_{\min}} \mathbf{z}_i \mathbf{z}_i^T \sqrt{\lambda_{\min}}$, we find that

$$\begin{aligned} \lambda_{\min}^2 \mathbf{X} \mathbf{X}^T &= \lambda_{\min}^2 \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T \leq \mathbf{X} \mathbf{X}^T = \sum_{i=1}^N \sqrt{\lambda_{\min}} \mathbf{z}_i \mathbf{z}_i^T \sqrt{\lambda_{\min}} \\ &\leq \lambda_{\max}^2 \sum_{i=1}^N \mathbf{z}_i \mathbf{z}_i^T = \lambda_{\max}^2 \mathbf{Z}^T \mathbf{Z}, \end{aligned}$$

which implies that

$$\lambda_{\min}(\mathbf{X}) \leq \lambda_{\min}(\mathbf{Z}) \leq \|\mathbf{Z}\| \leq \lambda_{\max} \|\mathbf{X}\|.$$

Plugging this into (A.16) completes the proof. \blacksquare

B Sys ID Analysis

We first list in Table 4 a few shorthand notations to be used in this appendix. They are mainly used in the fictional sub-trajectories analysis in Appendices B.2 and B.4. Notations on the inside the parentheses are arguments to be replaced with context-depending variables.

B.1 Estimating \mathbf{T}

The following theorem adapted from [67, Lemma 7] provides the sample complexity result for estimating Markov matrix \mathbf{T} , which is a corresponds to the sample complexity on $\|\hat{\mathbf{T}} - \mathbf{T}\|$ in Theorem 4.1.

Theorem B.1 Suppose we have an ergodic Markov chain $\mathbf{T} \in \mathbb{R}^{S \times S}$ with mixing time t_{MC} and stationary distribution $\pi \in \mathbb{R}^S$. Let $\lambda_{\max} := \max_i \pi(s_i)$ and $\lambda_{\min} := \min_i \pi(s_i)$. Given a state sequence $(0), (1), \dots, (T)$ of the Markov chain, define the empirical estimator $\hat{\mathbf{T}}$ of the Markov matrix as follows,

$$[\hat{\mathbf{T}}]_{ij} = \frac{\sum_{t=1}^{T-1} \mathbf{1}_{\{(t)=i, (t+1)=j\}}}{\sum_{t=1}^{T-1} \mathbf{1}_{\{(t)=i\}}},$$

Table 4: Notations — Sampling Periods

$\underline{C}_{\mathbf{w}}(T, \cdot)$	$\sqrt{2\log(nT)} + \sqrt{2\log(2/\cdot)}$
$\cdot + (\cdot, \cdot, \underline{C}_{\mathbf{w}})$	$\sqrt{\frac{2 \cdot \bar{s}(c_{\mathbf{w}}^2 + C_{\mathbf{z}}^2 \mathbf{B}_{1:s}^{-2})}{1-\cdot}}$
$\cdot + (\cdot, \cdot, \underline{C}_{\mathbf{w}}, \mathbf{K}_{1:s})$	$\underline{C}_{\mathbf{w}} + \cdot + (\cdot, \cdot, \underline{C}_{\mathbf{w}})(\ \mathbf{A}_{1:s}\ + \ \mathbf{B}_{1:s}\ \ \mathbf{K}_{1:s}\) + C_{\mathbf{z}}\sqrt{\rho/n}\ \mathbf{B}_{1:s}\ $
$\underline{C}_{\text{Sub},\mathbf{x}}(\bar{X}_0, \cdot, T, \cdot, \cdot)$	$\frac{2}{\log(\cdot^{-1})} + \frac{2}{\log(\cdot^{-1})\log(T)} \log\left(\frac{24n \cdot \bar{s} \max\{\bar{x}_0^2, \cdot^{-2}\}}{(1-\cdot)}\right)\}$
$\underline{C}_{\text{Sub},\bar{\mathbf{x}}}(\cdot, T, \cdot, \cdot)$	$\frac{1}{\log(\cdot^{-1})} + \frac{1}{\log(\cdot^{-1})\log(T)} \log\left(\frac{72 \cdot \bar{n}s^{1.5}}{\cdot}\right)$
$\underline{C}_{\text{Sub},N}(\bar{X}_0, \cdot, T, \cdot, \cdot)$	$\max\{C_{MC}, \underline{C}_{\text{Sub},\mathbf{x}}(\bar{X}_0, \bar{2}, T, \cdot, \cdot), \underline{C}_{\text{Sub},\bar{\mathbf{x}}}(\cdot, T, \cdot, \cdot)\}$
$\underline{L}_{id,t_0}(\bar{X}_0, \cdot, \cdot, \underline{C}_{\mathbf{w}})$	$\frac{\log((1-\cdot)\bar{x}_0^2 (c_{\mathbf{w}}^{-2} \cdot^{-2} + \frac{2}{\mathbf{z}} \mathbf{B}_{1:s}^{-2}))}{1-\cdot}$
$\underline{L}_{id,cov}(\cdot, \cdot, \underline{C}_{\mathbf{w}}, \mathbf{K}_{1:s})$	$1 + \frac{2\log(8c^2 \cdot + (\cdot, \cdot, \underline{C}_{\mathbf{w}}) \cdot + (\cdot, \cdot, \underline{C}_{\mathbf{w}}, \mathbf{K}_{1:s})n \cdot \bar{n}s)}{1-\cdot}$
$\underline{L}_{id,tr1}(\cdot, T, \cdot, \cdot, \underline{C}_{\mathbf{w}}, \mathbf{K}_{1:s})$	$1 + \frac{2\log(2 \cdot \bar{n}s \cdot T \cdot + (\cdot, \cdot, \underline{C}_{\mathbf{w}}, \mathbf{K}_{1:s}) \cdot + (\cdot, \cdot, \underline{C}_{\mathbf{w}}) \cdot)}{1-\cdot}$
$\underline{L}_{id,tr2}(\cdot, T, \cdot, \cdot, \underline{C}_{\mathbf{w}}, \mathbf{K}_{1:s})$	$1 + \frac{2}{(1-\cdot)} \log\left(\frac{192c^2 \cdot + (\cdot, \cdot, \underline{C}_{\mathbf{w}}, \mathbf{K}_{1:s})(1 + (\cdot, \cdot, \underline{C}_{\mathbf{w}})n \cdot \bar{s}(n+p)T)}{\cdot}\right)$
$\underline{L}_{id,tr3}(\cdot, T, \cdot, \cdot, \underline{C}_{\mathbf{w}}, \mathbf{K}_{1:s})$	$1 + \frac{2}{(1-\cdot)} \log\left(\frac{c_{\mathbf{w}} \cdot \cdot + (\cdot, \cdot, \underline{C}_{\mathbf{w}}, \mathbf{K}_{1:s}) \cdot n \cdot \bar{n}s^2}{(1 + (\cdot, \cdot, \underline{C}_{\mathbf{w}}) \cdot) \cdot (n+p) (C_{\mathbf{w}} \cdot \bar{n} + \bar{p} + C_0 \log(2s \cdot))}\right)$
$\underline{L}_{id}(\bar{X}_0, \cdot, T, \cdot, \cdot, \underline{C}_{\mathbf{w}}, \mathbf{K}_{1:s}, L)$	$\max\{\underline{L}_{id,t_0}(\bar{X}_0, \cdot, \cdot, \underline{C}_{\mathbf{w}}), \underline{L}_{id,cov}(\cdot, \cdot, \underline{C}_{\mathbf{w}}, \mathbf{K}_{1:s}),$ $\underline{L}_{id,tr1}(\frac{\cdot}{36L}, T, \cdot, \cdot, \underline{C}_{\mathbf{w}}, \mathbf{K}_{1:s}), \underline{L}_{id,tr2}(\frac{\cdot}{36L}, T, \cdot, \cdot, \underline{C}_{\mathbf{w}}, \mathbf{K}_{1:s}),$ $\underline{L}_{id,tr3}(\frac{\cdot}{36L}, T, \cdot, \cdot, \underline{C}_{\mathbf{w}}, \mathbf{K}_{1:s}), \underline{C}_{\text{Sub},N}(\bar{X}_0, \frac{\cdot}{2L}, T, \cdot, \cdot) \log(T)\}$

Assume for some $\cdot > 0$, $T \geq \underline{I}_{MC,1}(C_{MC}, \bar{4}) := (68C_{MC} \max_{\min}^{-2} \log(4s))^2$, where C_{MC} is defined in Table 2. Then, we have with probability at least $1 - \cdot$,

$$\|\hat{\mathbf{T}} - \mathbf{T}\| \leq 4 \frac{\cdot^{-1}}{\min} \|\mathbf{T}\| \sqrt{\frac{17 \max C_{MC} \log(T)}{T} \log\left(\frac{4sC_{MC} \log(T)}{\cdot}\right)}. \quad (\text{B.1})$$

Proof We first consider the estimators computed using a sub-trajectory of $(0), (1), \dots, (T)$, then combine them together to show the error bound for $\hat{\mathbf{T}}$ in the claim. For C_{MC} defined in Table 2, let $L = C_{MC} \log(T)$. Then, for $l = 0, 1, \dots, L-1$, define $\hat{\mathbf{T}}^{(l)} \in \mathbb{R}^{s \times s}$ such that $[\hat{\mathbf{T}}^{(l)}]_{ij} = \frac{\mathbf{1}_{\{k=1}^{\lfloor T/L \rfloor} \mathbf{1}_{\{(kL+l)=i, (kL+1+l)=j\}}}}{\mathbf{1}_{\{k=1}^{\lfloor T/L \rfloor} \mathbf{1}_{\{(kL+l)=i\}}}}$. In other words, $\hat{\mathbf{T}}^{(l)}$ is the estimator computed using data with subsampling period L . Following the proof of [67, Lemma 7], we know for any $\cdot < \min/2$, suppose $L \geq 6t_{MC} \log(\cdot^{-1})$.

$$\mathbb{P}(\|\hat{\mathbf{T}}^{(l)} - \mathbf{T}\| \leq 4 \frac{\cdot^{-1}}{\min} \|\mathbf{T}\|) \geq 1 - 4s \exp\left(-\frac{T^2}{17 \max L}\right). \quad (\text{B.2})$$

By setting $\cdot = 4s \exp\left(-\frac{T^2}{17 \max L}\right)$, one can also interpret the above result as: for all $\cdot > 0$, suppose

$$L \geq 3t_{MC} \log\left(\frac{T}{17 \max L \log(4s)}\right), \quad (\text{B.3})$$

then when

$$T \geq 68L \max_{\min}^{-2} \log\left(\frac{4s}{\cdot}\right), \quad (\text{B.4})$$

we have with probability at least $1 - \cdot$

$$\|\hat{\mathbf{T}}^{(l)} - \mathbf{T}\| \leq 4 \frac{\cdot^{-1}}{\min} \|\mathbf{T}\| \sqrt{\frac{17 \max C_{MC} \log(T) \log(4s)}{T}}. \quad (\text{B.5})$$

One can verify (B.3) holds by plugging in $L = C_{MC} \log(T)$ and using definition $C_{MC} := t_{MC} \cdot \max\{3, 3 - 3 \log(\max \log(s))\}$; (B.4) holds under the premise condition $T \geq \underline{L}_{MC,1}(C_{MC}, \frac{1}{4}) := (68 C_{MC} \max \frac{1}{\min} \log(\frac{4s}{s}))^2$.

Note that by definition, $\hat{\mathbf{T}}$ can be viewed as a convex combination of $\hat{\mathbf{T}}^{(l)}$ for all $l = 0, 1, \dots, L$, thus by triangle inequality and union bound, we have with probability $1 - L^{-1}$,

$$\|\hat{\mathbf{T}} - \mathbf{T}\| \leq 4 \frac{1}{\min} \|\mathbf{T}\| \sqrt{\frac{17 \max C_{MC} \log(T) \log(\frac{4s}{s})}{T}}. \quad (\text{B.6})$$

Finally, by replacing L with $\frac{1}{L^{-1}}$, we could show (B.1) and conclude the proof. \blacksquare

B.2 Estimation of \mathbf{A}_{1s} and \mathbf{B}_{1s} from a Single Trajectory (Main SYSID Analysis)

In this section, we estimate the MJS dynamics \mathbf{A}_{1s} and \mathbf{B}_{1s} from finite samples obtained from a single trajectory of (3.1). To estimate a coarse model of the unknown system dynamics, we use the method of linear least squares. By running experiments in which the system starts at state \mathbf{x}_0 and the dynamics evolve with a given input, we can record the resulting state, excitation and mode observations. Let \mathbf{K}_{1s} stabilizes the system (3.1) in the mean square sense according to Def. 3.1. Then, choosing the input to be $\mathbf{u}_t = \mathbf{K}_{1s} \mathbf{x}_t + \mathbf{z}_t$, the state updates as follows,

$$\mathbf{x}_{t+1} = (\mathbf{A}_{1s}(t) + \mathbf{B}_{1s}(t) \mathbf{K}_{1s}(t)) \mathbf{x}_t + \mathbf{B}_{1s}(t) \mathbf{z}_t + \mathbf{w}_t = \mathbf{L}_{1s}(t) \mathbf{x}_t + \mathbf{B}_{1s}(t) \mathbf{z}_t + \mathbf{w}_t, \quad (\text{B.7})$$

where $\{\mathbf{z}_t\}_{t=0}^T \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \frac{2}{\rho} \mathbf{I}_p)$ is the i.i.d. excitation for exploration and we let $\mathbf{L}_{1s}(t) := \mathbf{A}_{1s}(t) + \mathbf{B}_{1s}(t) \mathbf{K}_{1s}(t)$. Observe that the closed-loop state update (B.7) can be expanded as follows,

$$\mathbf{x}_t = \begin{cases} \mathbf{x}_0 & \text{if } t = 0, \\ \mathbf{L}_{1s}(0) \mathbf{x}_0 + \mathbf{B}_{1s}(0) \mathbf{z}_0 + \mathbf{w}_0 & \text{if } t = 1, \\ \prod_{j=0}^{t-1} \mathbf{L}_{1s}(j) \mathbf{x}_0 + \sum_{\ell=0}^{t-2} \prod_{j=\ell+1}^{t-1} \mathbf{L}_{1s}(j) \mathbf{B}_{1s}(\ell) \mathbf{z}_\ell + \mathbf{B}_{1s}(t-1) \mathbf{z}_{t-1} \\ \quad + \sum_{\ell=0}^{t-2} \prod_{j=\ell+1}^{t-1} \mathbf{L}_{1s}(j) \mathbf{w}_\ell + \mathbf{w}_{t-1} & \text{if } t \geq 2. \end{cases} \quad (\text{B.8})$$

B.2.1 Estimation from Bounded States

To estimate the unknown system dynamics, we run the system for T time-steps and collect the samples $(\mathbf{x}_t, \mathbf{z}_t, \mathbf{x}_t, \mathbf{x}_{t+1})_{t=0}^{T-1}$. Then, we run Alg. 1 to get the estimates $(\hat{\mathbf{A}}_i, \hat{\mathbf{B}}_i)$ for all $i \in [S]$. Our learning method is described in Alg. 1. For the ease of analysis, we first derive the estimation error bounds with the following assumption on the noise.

Assumption 3 (subGaussian noise) Let $\{\mathbf{w}_t\}_{t=0}^{T-1} \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}_w$. There exists $c_w > 0$ and $c_w \geq 1$ such that, each entry of \mathbf{w}_t is i.i.d. zero-mean subGaussian with variance $\frac{2}{c_w}$ and we have $\|\mathbf{w}_t\| \leq c_w \frac{2}{c_w}$.

Later on, we will relax this assumption to get the estimation error bounds with the Gaussian noise. To proceed, we first show that the Euclidean norm of the states \mathbf{x}_t in (B.7) can be upper bounded in expectation. The following result, which is a corollary of Lemma A.4, accomplishes this.

Corollary B.2 (Bounded states) Let \mathbf{x}_t be the state at time t of the MJS (B.7), with initial state $\mathbf{x}_0 \sim \mathcal{D}_x$ such that $\mathbb{E}[\mathbf{x}_0] = 0$, $\mathbb{E}[\|\mathbf{x}_0\|^2] \leq \frac{2}{\rho} n$ for some $\rho > 0$. Suppose Assumption 1 on the system and the Markov chain and Assumption 3 on the process noise hold. Suppose $\{\mathbf{z}_t\}_{t=0}^T \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \frac{2}{\rho} \mathbf{I}_p)$. Let $C_z := \frac{2}{\rho}$ be a constant, $\mathbf{h}_t := [\frac{1}{c_w} \mathbf{x}_t \quad \frac{1}{c_w} \mathbf{z}_t]$ be the concatenated state and define

$$t_0 := \frac{\log((1 - \tilde{\epsilon}) \frac{2}{\rho} n / (c_w^2 \frac{2}{c_w} + \frac{2}{\rho} \|\mathbf{B}_{1s}\|))}{1 - \tilde{\epsilon}}, \quad (\text{B.9})$$

$$t_+ := \frac{2\sqrt{s}(c_w^2 + C_z^2 \|\mathbf{B}_{1s}\|^2) \tilde{\epsilon}}{1 - \tilde{\epsilon}}. \quad (\text{B.10})$$

Then, for all $t \geq t_0$, we have

$$\mathbb{E}[\|\mathbf{x}_t\|^2] \leq \frac{2}{\mathbf{w}} \frac{2}{\mathbf{z}} n \quad \text{and} \quad \mathbb{E}[\|\mathbf{h}_t\|^2] \leq (1 + \frac{2}{\mathbf{z}})(n + \rho). \quad (\text{B.11})$$

Proof Recall from Lemma A.4 that the states \mathbf{x}_t can be bounded in expectation as follows,

$$\mathbb{E}[\|\mathbf{x}_t\|^2] \leq \frac{2}{\mathbf{z}} \sqrt{s} \left(\frac{t}{\mathbf{z}} \frac{2}{\mathbf{z}} \sqrt{n} + \frac{c_{\mathbf{w}}^2 \frac{2}{\mathbf{w}} + \frac{2}{\mathbf{z}} \|\mathbf{B}_{1s}\|^2}{1 - \mathbf{z}} \right) n \leq \frac{2}{\mathbf{w}} \frac{2}{\mathbf{z}} \sqrt{s} (c_{\mathbf{w}}^2 + C_{\mathbf{z}}^2 \|\mathbf{B}_{1s}\|^2) \frac{\mathbf{z}}{1 - \mathbf{z}} n, \quad (\text{B.12})$$

where we get the last inequality by choosing the timestep t to satisfy the following lower bound,

$$\frac{t}{\mathbf{z}} \leq \frac{c_{\mathbf{w}}^2 \frac{2}{\mathbf{w}} + \frac{2}{\mathbf{z}} \|\mathbf{B}_{1s}\|^2}{(1 - \mathbf{z}) \frac{2}{\mathbf{z}} \sqrt{n}} \iff t \geq t_0 := \frac{\log((1 - \mathbf{z}) \frac{2}{\mathbf{z}} n / (c_{\mathbf{w}}^2 \frac{2}{\mathbf{w}} + \frac{2}{\mathbf{z}} \|\mathbf{B}_{1s}\|))}{1 - \mathbf{z}}. \quad (\text{B.13})$$

This gives the advertised upper bound on $\mathbb{E}[\|\mathbf{x}_t\|^2]$ for $t \geq t_0$. Using Jensen's inequality, this further implies

$$\mathbb{E}[\|\mathbf{x}_t\|] \leq \frac{2}{\mathbf{w}} \sqrt{\frac{2s^{1/2} (c_{\mathbf{w}}^2 + C_{\mathbf{z}}^2 \|\mathbf{B}_{1s}\|^2) \frac{\mathbf{z}}{1 - \mathbf{z}} n}{1 - \mathbf{z}}} \quad \text{for } t \geq t_0. \quad (\text{B.14})$$

Next, using standard results on the distribution of squared Euclidean norm of a Gaussian vector, we have $\mathbb{E}[\|\mathbf{z}_t\|^2] = \frac{2}{\mathbf{z}} \rho$ for all $t \geq 0$. Combining this with (B.12), we get the following upper bound on the expected squared norm of $\mathbf{h}_t := [\frac{1}{\mathbf{w}} \mathbf{x}_t \quad \frac{1}{\mathbf{z}} \mathbf{z}_t]$, that is, for all $t \geq t_0$, we have

$$\begin{aligned} \mathbb{E}[\|\mathbf{h}_t\|^2] &= \frac{1}{\mathbf{w}} \mathbb{E}[\|\mathbf{x}_t\|^2] + \frac{1}{\mathbf{z}} \mathbb{E}[\|\mathbf{z}_t\|^2] \leq \frac{2\sqrt{s}(c_{\mathbf{w}}^2 + C_{\mathbf{z}}^2 \|\mathbf{B}_{1s}\|^2) \frac{\mathbf{z}}{1 - \mathbf{z}} n}{1 - \mathbf{z}} + \rho, \\ &\leq \left(1 + \frac{2\sqrt{s}(c_{\mathbf{w}}^2 + C_{\mathbf{z}}^2 \|\mathbf{B}_{1s}\|^2) \frac{\mathbf{z}}{1 - \mathbf{z}}}{1 - \mathbf{z}} \right) (n + \rho). \end{aligned} \quad (\text{B.15})$$

This gives the advertised upper bound on $\mathbb{E}[\|\mathbf{h}_t\|^2]$ for $t \geq t_0$. Using Jensen's inequality, this further implies

$$\mathbb{E}[\|\mathbf{h}_t\|] \leq \sqrt{\left(1 + \frac{2s^{1/2} (c_{\mathbf{w}}^2 + C_{\mathbf{z}}^2 \|\mathbf{B}_{1s}\|^2) \frac{\mathbf{z}}{1 - \mathbf{z}}}{1 - \mathbf{z}} \right) (n + \rho)} \quad \text{for } t \geq t_0. \quad (\text{B.16})$$

This completes the proof. \blacksquare

Suppose the MJS in (B.7) is run for T timesteps and we have access to the trajectory $(\mathbf{x}_t, \mathbf{z}_t)_{t=0}^T$. Our proof strategy is based on the observation that a single MJS trajectory can be split into multiple weakly dependent sub-trajectories, defined as follows.

Definition B.3 (Sub-trajectories of bounded states) Let sampling period $L \geq 1$ be an integer. Let $\mathcal{K} = \{k = 0 + kL \mid 0 \leq k \leq L - 1\}$ be a fixed set and $k = 1, 2, \dots, \lfloor \frac{T-L}{L} \rfloor$. We sub-sample the trajectory $(\mathbf{x}_t, \mathbf{z}_t)_{t=0}^T$ at time indices $\mathcal{K} \in S_i^{(k)}$, where

$$S_i^{(k)} := \{k \mid (\mathbf{x}_k) = i, \|\mathbf{x}_k\| \leq c_{\mathbf{w}} \sqrt{n}, \|\mathbf{z}_k\| \leq c_{\mathbf{z}} \sqrt{\rho}\}, \quad (\text{B.17})$$

to obtain the i th sub-trajectory $\{(\mathbf{x}_{k+1}, \mathbf{x}_k, \mathbf{z}_{k+1}, \mathbf{z}_k)_{k \in S_i^{(k)}}\}$.

Note that, $S_i = \bigcup_{k=0}^{L-1} S_i^{(k)}$, where S_i is as defined in Alg. 1. This shows that a single trajectory with bounded states and excitations $\{(\mathbf{x}_{t+1}, \mathbf{x}_t, \mathbf{z}_{t+1}, \mathbf{z}_t)\}_{t \in S_i}$ can be split into L weakly dependent sub-trajectories $\{(\mathbf{x}_{k+1}, \mathbf{x}_k, \mathbf{z}_{k+1}, \mathbf{z}_k)_{k \in S_i^{(k)}}\}_{k \in S_i^{(k)}}$ for $0 \leq k \leq L - 1$. To proceed, we first show that the covariance of the bounded states has the following properties.

Lemma B.4 (Covariance of bounded states) Consider the same setup of Corollary B.2. Let t_0 and ϵ be as in (B.9) and (B.10) respectively. Let $c \geq 6$ be a fixed constant. For all $t \geq t_0$, we have

$$\left(\frac{2}{\mathbf{w}}/2\right)I_n \leq \mathbb{E}[\mathbf{x}_t | \|\mathbf{x}_t\| \leq c \sqrt{\mathbf{w}} + \sqrt{n}, \|\mathbf{z}_t\| \leq c \sqrt{\mathbf{z}}\sqrt{\rho}] \leq c^2 \frac{2}{\mathbf{w}} \frac{2}{\mathbf{z}} n I_n, \quad (\text{B.18})$$

$$(1/2)I_{n+p} \leq \mathbb{E}[\mathbf{h}_t | \|\mathbf{x}_t\| \leq c \sqrt{\mathbf{w}} + \sqrt{n}, \|\mathbf{z}_t\| \leq c \sqrt{\mathbf{z}}\sqrt{\rho}] \leq c^2(1 + \frac{2}{\mathbf{z}})(n+p)I_{n+p}. \quad (\text{B.19})$$

Proof Let \mathbf{x}_t be the state at time t of the MJS given by (B.7), with initial state $\mathbf{x}_0 \sim \mathcal{D}_x$ such that $\mathbb{E}[\mathbf{x}_0] = 0$ and define the noise-removed state $\tilde{\mathbf{x}}_t = \mathbf{x}_t - \mathbf{w}_{t-1}$ which is independent of \mathbf{w}_{t-1} . Observe that $\mathbb{E}[\tilde{\mathbf{x}}_t] = 0$ because both \mathbf{x}_t and \mathbf{w}_{t-1} are zero-mean. Next, from Corollary B.2, we know that, for all $t \geq t_0$, we have $\mathbb{E}[\|\mathbf{x}_t\|] \leq \sqrt{\mathbf{w}} + \sqrt{n}$. Combining this with $\mathbb{E}[\|\mathbf{w}_{t-1}\|] \leq c \sqrt{\mathbf{w}} \sqrt{n}$, we have

$$\mathbb{E}[\|\tilde{\mathbf{x}}_t\|] \leq \mathbb{E}[\|\mathbf{x}_t\|] + \mathbb{E}[\|\mathbf{w}_{t-1}\|] \leq 2 \sqrt{\mathbf{w}} + \sqrt{n}. \quad (\text{B.20})$$

To proceed, consider the conditional random variable $\mathbf{y}_t \sim \{\mathbf{x}_t | \|\mathbf{x}_t\| \leq c \sqrt{\mathbf{w}} + \sqrt{n}\}$. To lower bound the covariance matrix $\mathbb{E}[\mathbf{y}_t \mathbf{y}_t^T]$, observe that $\|\mathbf{w}_{t-1}\| \leq c \sqrt{\mathbf{w}} \sqrt{n} \leq \sqrt{\mathbf{w}} + \sqrt{n}$ and $\mathbb{E}[\|\tilde{\mathbf{x}}_t\|] \leq 2 \sqrt{\mathbf{w}} + \sqrt{n}$, where ϵ is given by (B.10). Therefore, using Lemma A.6 from Section A.2 with $B = 2 \sqrt{\mathbf{w}} + \sqrt{n}$, we can lower bound $\mathbb{E}[\mathbf{y}_t \mathbf{y}_t^T]$ as follows,

$$\mathbb{E}[\mathbf{y}_t \mathbf{y}_t^T] = \mathbb{E}[\mathbf{x}_t | \|\mathbf{x}_t\| \leq c \sqrt{\mathbf{w}} + \sqrt{n}] \geq \left(\frac{2}{\mathbf{w}}/2\right)I_n, \quad (\text{B.21})$$

where we use $c \geq 6$ to get the last inequality. Next, we upper bound the covariance matrix $\mathbb{E}[\mathbf{y}_t \mathbf{y}_t^T]$ as follows,

$$\|\mathbb{E}[\mathbf{y}_t \mathbf{y}_t^T]\| = \mathbb{E}[\|\mathbf{y}_t \mathbf{y}_t^T\|] \leq \mathbb{E}[\|\mathbf{y}_t\|^2] = \mathbb{E}[\|\mathbf{x}_t\|^2 | \|\mathbf{x}_t\| \leq c \sqrt{\mathbf{w}} + \sqrt{n}] \leq c^2 \frac{2}{\mathbf{w}} \frac{2}{\mathbf{z}} n, \quad (\text{B.22})$$

Combining (B.21) and (B.22), we get the first statement of the lemma. Next, using a similar argument with Lemma A.7, we can show that, when $c \geq 6$, we also have

$$\left(\frac{2}{\mathbf{z}}/2\right)I_p \leq \mathbb{E}[\mathbf{z}_t | \|\mathbf{z}_t\| \leq c \sqrt{\mathbf{z}}\sqrt{\rho}] \leq (2 \frac{2}{\mathbf{z}}\rho)I_p. \quad (\text{B.23})$$

Finally, combining the derived bounds on $\mathbb{E}[\mathbf{x}_t | \|\mathbf{x}_t\| \leq c \sqrt{\mathbf{w}} + \sqrt{n}]$ and $\mathbb{E}[\mathbf{z}_t | \|\mathbf{z}_t\| \leq c \sqrt{\mathbf{z}}\sqrt{\rho}]$, we get the second statement of the lemma,

$$(1/2)I_{n+p} \leq \mathbb{E}[\mathbf{h}_t | \|\mathbf{x}_t\| \leq c \sqrt{\mathbf{w}} + \sqrt{n}, \|\mathbf{z}_t\| \leq c \sqrt{\mathbf{z}}\sqrt{\rho}] \leq c^2(1 + \frac{2}{\mathbf{z}})(n+p)I_{n+p}. \quad (\text{B.24})$$

This completes the proof. \blacksquare

To proceed, let $\mathbf{h}_t := [\frac{1}{\mathbf{w}}\mathbf{x}_t \quad \frac{1}{\mathbf{z}}\mathbf{z}_t]$ be the concatenated state and $\mathbf{h}_i := [\mathbf{w}_i \mathbf{L}_i \quad \mathbf{z}_i \mathbf{B}_i]$ for all $i \in [S]$. Then the output of each sample in $\{(\mathbf{x}_{t+1}, \mathbf{x}_t, \mathbf{z}_t, \mathbf{y}_t) | t \in S_i\}$ can be related to the inputs as follows,

$$\mathbf{x}_{t+1} = \mathbf{A}_t \mathbf{h}_t + \mathbf{w}_t \quad \text{for all } t \in S_i. \quad (\text{B.25})$$

Next, to carry out finite sample identification of \mathbf{h}_i using the method of linear least squares, we define the following concatenated matrices,

$$\mathbf{Y}_i = \begin{bmatrix} \mathbf{x}_{t_1+1} \\ \mathbf{x}_{t_2+1} \\ \vdots \\ \mathbf{x}_{t_{|S_i|}+1} \end{bmatrix}, \quad \mathbf{H}_i = \begin{bmatrix} \mathbf{h}_{t_1} \\ \mathbf{h}_{t_2} \\ \vdots \\ \mathbf{h}_{t_{|S_i|}} \end{bmatrix}, \quad \mathbf{W}_i = \begin{bmatrix} \mathbf{w}_{t_1} \\ \mathbf{w}_{t_2} \\ \vdots \\ \mathbf{w}_{t_{|S_i|}} \end{bmatrix}, \quad (\text{B.26})$$

that is, \mathbf{Y}_i has $\{\mathbf{x}_{t+1}\}_{t \in S_i}$ on its rows, \mathbf{H}_i has $\{\mathbf{h}_t\}_{t \in S_i}$ on its rows and \mathbf{W}_i has $\{\mathbf{w}_t\}_{t \in S_i}$ on its rows. Similarly, we also construct $\mathbf{Y}_i^{(\cdot)}$, $\mathbf{H}_i^{(\cdot)}$ and $\mathbf{W}_i^{(\cdot)}$ by (row-wise) stacking $\{\mathbf{x}_{k+1}\}_{k \in S_i^{(\cdot)}}$, $\{\mathbf{h}_k\}_{k \in S_i^{(\cdot)}}$ and $\{\mathbf{w}_k\}_{k \in S_i^{(\cdot)}}$ respectively. Then, we have $\mathbf{Y}_i = \mathbf{H}_i \mathbf{h}_i + \mathbf{W}_i$. Our goal in this paper is to solve the following least squares problems,

$$\hat{\mathbf{h}}_i = \arg \min_{\mathbf{h}_i} \frac{1}{2|S_i|} \|\mathbf{Y}_i - \mathbf{H}_i \mathbf{h}_i\|_F^2. \quad (\text{B.27})$$

The least squares estimator of \mathbf{y}_i is $\hat{\mathbf{y}}_i = \mathbf{H}_i^\dagger \mathbf{Y}_i = (\mathbf{H}_i \mathbf{H}_i)^{-1} \mathbf{H}_i \mathbf{Y}_i$ and its estimation error is given by

$$\begin{aligned} \|\hat{\mathbf{y}}_i - \mathbf{y}_i\| &= \|(\mathbf{H}_i \mathbf{H}_i)^{-1} \mathbf{H}_i \mathbf{W}_i\| \leq \|(\mathbf{H}_i \mathbf{H}_i)^{-1}\| \|\mathbf{H}_i \mathbf{W}_i\| = \frac{\|\mathbf{H}_i \mathbf{W}_i\|}{\min(\mathbf{H}_i \mathbf{H}_i)}, \\ &\stackrel{(a)}{=} \frac{\|\sum_{i=0}^{L-1} \mathbf{H}_i^{(\cdot)} \mathbf{W}_i^{(\cdot)}\|}{\min(\sum_{i=0}^{L-1} \mathbf{H}_i^{(\cdot)} \mathbf{H}_i^{(\cdot)})}, \\ &\stackrel{(b)}{=} \frac{\sum_{i=0}^{L-1} \|\mathbf{H}_i^{(\cdot)} \mathbf{W}_i^{(\cdot)}\|}{\sum_{i=0}^{L-1} \min(\mathbf{H}_i^{(\cdot)} \mathbf{H}_i^{(\cdot)})}, \end{aligned} \quad (\text{B.28})$$

where we obtain (a) from the fact that $S_i = \cup_{i=0}^{L-1} S_i^{(\cdot)}$ and (b) follows from using triangular inequality and Weyl's inequality for Hermitian matrices. If we upper bound the terms $\|\mathbf{H}_i^{(\cdot)} \mathbf{W}_i^{(\cdot)}\|$ and lower bound the terms $\min(\mathbf{H}_i^{(\cdot)} \mathbf{H}_i^{(\cdot)})$, for all $0 \leq i \leq L-1$, we can use (B.28) to upper bound the estimation error $\|\hat{\mathbf{y}}_i - \mathbf{y}_i\|$. However, because $\mathbf{H}_i^{(\cdot)}$ has non-i.i.d. rows, it is not straightforward to bound the terms $\|\mathbf{H}_i^{(\cdot)} \mathbf{W}_i^{(\cdot)}\|$ and $\min(\mathbf{H}_i^{(\cdot)} \mathbf{H}_i^{(\cdot)})$ directly. To resolve this issue, we rely on the notion of stability and use perturbation based techniques to indirectly bound these terms in the following sub-sections.

B.2.2 Estimation from Truncated States

Definition B.5 (Truncated state vector [51]) Consider the state equation (B.7). Given $t \geq L > 0$, the L -truncation of \mathbf{x}_t is denoted by $\mathbf{x}_{t,L}$ and is obtained by driving the system with excitation \mathbf{z} and additive noise \mathbf{w} until time t , where

$$\mathbf{v} = \begin{cases} 0 & \text{if } t < t-L \\ \mathbf{v} & \text{else} \end{cases}. \quad (\text{B.29})$$

In words, the L -truncated state vector $\mathbf{x}_{t,L}$ is obtained by unrolling \mathbf{x}_t until time $t-L$ and setting \mathbf{x}_{t-L} to 0.

Using a truncation argument, we can obtain independent samples from a single trajectory which will be used to capture the effect of learning from a single trajectory. With high probability over the mode observation, truncated states can be made very close to the original states with sufficiently large truncation length. From (B.8), we have

$$\mathbf{x}_t - \mathbf{x}_{t,L} = \prod_{j=t-L}^{t-1} \mathbf{L}^{(j)} \mathbf{x}_{t-L}. \quad (\text{B.30})$$

As a corollary of Lemma A.4, observe that for a closed loop autonomous system $\mathbf{x}_{t+1} = \mathbf{L}^{(t)} \mathbf{x}_t$, mean-square stability implies that, for any initial conditions \mathbf{x}_0 and $\mathbf{0}$, we have $E[\|\mathbf{x}_t\|^2] \leq \sqrt{ns} \prod_{i=0}^{t-1} \|\mathbf{x}_0\|^2$. Combining this argument with (B.30), we have

$$\begin{aligned} E[\|\mathbf{x}_t - \mathbf{x}_{t,L}\|^2] &= E[\|\prod_{j=t-L}^{t-1} \mathbf{L}^{(j)} \mathbf{x}_{t-L}\|^2] \leq \sqrt{ns} \prod_{i=0}^{t-1} \|\mathbf{x}_{t-L}\|^2, \\ \implies E[\|\mathbf{x}_t - \mathbf{x}_{t,L}\|] &\leq \sqrt{(ns)^{1/2}} \prod_{i=0}^{t-1} \|\mathbf{x}_{t-L}\| \leq \sqrt{ns} \prod_{i=0}^{t-1} \|\mathbf{x}_{t-L}\|, \end{aligned} \quad (\text{B.31})$$

where the expectation is over the Markov modes $\{\mathbf{L}^{(j)}\}_{j=t-L}^{t-1}$ and we get the last relation by using Jensen's inequality. Moreover, if we also have $\|\mathbf{x}_{t-L}\| \leq c \|\mathbf{w}\| + \sqrt{n}$, then we can make $E[\|\mathbf{x}_t - \mathbf{x}_{t,L}\|]$ arbitrarily small by picking a sufficiently large truncation length $L \geq 1$,

$$E[\|\mathbf{x}_t - \mathbf{x}_{t,L}\| \mid \|\mathbf{x}_{t-L}\| \leq c \|\mathbf{w}\| + \sqrt{n}] \leq \prod_{i=0}^{t-1} \|\mathbf{x}_{t-L}\| \leq c \|\mathbf{w}\| + n\sqrt{s} \prod_{i=0}^{t-1} \|\mathbf{x}_{t-L}\|. \quad (\text{B.32})$$

To proceed, we carry out the truncation of the sub-trajectories introduced in Def. B.3 to get the truncated sub-trajectories defined as follows.

Definition B.6 (Truncated sub-trajectories) Let sampling period $L \geq 1$ be an integer. Let $k = i + kL$ be the sub-sampling indices, where $0 \leq i \leq L-1$ is a fixed offset and $k = 1, 2, \dots, \lfloor \frac{T-L}{L} \rfloor$. Let $S_i^{(s)}$ be as in Def. B.3. For each $k \in S_i^{(s)}$, let $k_{-k} \in S_i^{(s)}$ denotes the largest time index smaller than k . Given the i -th sub-trajectory $\{(\mathbf{x}_{k+1}, \mathbf{x}_k, \mathbf{z}_k, \mathbf{L}_{(k)}, \mathbf{B}_{(k)})\}_{k \in S_i^{(s)}}$ from Def. B.3, we truncate each state \mathbf{x}_k by $k - k_{-k} - 1$ to get the i -th truncated sub-trajectory $\{(\bar{\mathbf{x}}_{k+1}, \bar{\mathbf{x}}_k, \mathbf{z}_k, \mathbf{L}_{(k)}, \mathbf{B}_{(k)})\}_{k \in S_i^{(s)}}$, where

$$\bar{\mathbf{x}}_k := \mathbf{x}_{k - k_{-k} - 1} \quad \text{and} \quad \bar{\mathbf{x}}_{k+1} := \mathbf{L}_{(k)} \bar{\mathbf{x}}_k + \mathbf{B}_{(k)} \mathbf{z}_k + \mathbf{w}_k. \quad (\text{B.33})$$

If k_{-k} is the smallest time index in $S_i^{(s)}$, we set $k_{-k} = 0$.

Note that $k - k_{-k} \geq L$ by definition. Hence, the truncation lengths used to obtain $\{\bar{\mathbf{x}}_k\}_{k \in S_i^{(s)}}$ are always larger than $L-1$. Next, we show that when $L \geq 1$ is sufficiently large enough, then the truncated states $\{\bar{\mathbf{x}}_k\}_{k \in S_i^{(s)}}$ as well as the Euclidean distance between the truncated and non-truncated states can be bounded with high probability over the modes.

Lemma B.7 (Bounded states (truncated)) Consider the same setup of Corollary B.2. Let $\{\mathbf{x}_k\}_{k \in S_i^{(s)}}$ be the bounded states and $\{\bar{\mathbf{x}}_k\}_{k \in S_i^{(s)}}$ be the truncated states from Def. B.3 and B.6 respectively. Let

$$c_{\mathbf{w}} := c_{\mathbf{w}} + \|\mathbf{L}_{1s}\| + C_{\mathbf{z}} \sqrt{\rho/n} \|\mathbf{B}_{1s}\|, \quad (\text{B.34})$$

$$L_{tr1}(\bar{\epsilon}, \epsilon) := 1 + \frac{2 \log(2\sqrt{n} \bar{\epsilon} T + 1/\epsilon)}{1 - \bar{\epsilon}}, \quad (\text{B.35})$$

$$\text{and } L \geq \max\{t_0, L_{tr1}(\bar{\epsilon}, \epsilon)\}. \quad (\text{B.36})$$

Then, with probability at least $1 - \epsilon$ over the modes, for all $k \in S_i^{(s)}$ and all $i \in [s]$, we have

$$\|\mathbf{x}_k - \bar{\mathbf{x}}_k\| \leq (1/2) c_{\mathbf{w}} + \sqrt{n} \quad \text{and} \quad \|\bar{\mathbf{x}}_k\| \leq (3/2) c_{\mathbf{w}} + \sqrt{n}. \quad (\text{B.37})$$

Proof To begin, using Assumption 1 and (B.31), the impact of truncation can be bounded in expectation over the modes as follows,

$$\begin{aligned} \mathbb{E}[\|\mathbf{x}_k - \bar{\mathbf{x}}_k\|] &= \mathbb{E}[\|\mathbf{x}_{k+L} - \mathbf{x}_{k+L, (k-k')L-1}\|] \leq \sqrt{n} \bar{\epsilon} \bar{\epsilon}^{((k-k')L-1)/L} \|\mathbf{x}_{k'+L+1}\|, \\ &\leq \sqrt{n} \bar{\epsilon} \bar{\epsilon}^{(L-1)/L} \|\mathbf{x}_{k'+1}\|, \end{aligned} \quad (\text{B.38})$$

where we get the last inequality from the fact that $k - k_{-k} \geq 1$, and the expectation is over the Markov modes at timesteps $k + kL + 1, k + kL + 2, \dots, k + kL - 1$. To proceed, observe that, for all $k \in S_i^{(s)}$, we have

$$\begin{aligned} \|\mathbf{x}_{k+1}\| &= \|\mathbf{L}_{(k)} \mathbf{x}_k + \mathbf{B}_{(k)} \mathbf{z}_k + \mathbf{w}_k\| \leq \max_{i \in [s]} \|\mathbf{L}_i\| \|\mathbf{x}_k\| + \max_{i \in [s]} \|\mathbf{B}_i\| \|\mathbf{z}_k\| + \|\mathbf{w}_k\|, \\ &\leq c_{\mathbf{w}} + \sqrt{n} \|\mathbf{L}_{1s}\| + C_{\mathbf{z}} \sqrt{\rho/n} \|\mathbf{B}_{1s}\| + c_{\mathbf{w}} \sqrt{n}, \\ &\leq c_{\mathbf{w}} (c_{\mathbf{w}} + \|\mathbf{L}_{1s}\| + C_{\mathbf{z}} \sqrt{\rho/n} \|\mathbf{B}_{1s}\|) \sqrt{n}, \\ &= c_{\mathbf{w}} + \sqrt{n}, \end{aligned} \quad (\text{B.39})$$

where we set $c_{\mathbf{w}} := c_{\mathbf{w}} + \|\mathbf{L}_{1s}\| + C_{\mathbf{z}} \sqrt{\rho/n} \|\mathbf{B}_{1s}\|$ and $C_{\mathbf{z}} := C_{\mathbf{z}} \sqrt{\rho/n}$. Combining (B.39) with (B.38), for all $k \in S_i^{(s)}$ and all $i \in [s]$, we have

$$\mathbb{E}[\|\mathbf{x}_k - \bar{\mathbf{x}}_k\|] \leq c_{\mathbf{w}} + n\sqrt{\bar{\epsilon}} \bar{\epsilon}^{(L-1)/L}, \quad (\text{B.40})$$

$$\implies \mathbb{P}(\|\mathbf{x}_k - \bar{\mathbf{x}}_k\| \leq \frac{c_{\mathbf{w}} + n\sqrt{\bar{\epsilon}} \bar{\epsilon}^{(L-1)/L} T}{\epsilon}) \geq 1 - \epsilon, \quad (\text{B.41})$$

where we get the high probability bound by using Markov inequality and union bounding over all bounded states. This further implies that, with probability at least $1 - \epsilon$ over the modes, for all $k \in S_i^{(\cdot)}$ and all $i \in [S]$, we have

$$\|\bar{\mathbf{x}}_k\| \leq \|\mathbf{x}_k\| + \|\mathbf{x}_k - \bar{\mathbf{x}}_k\| \leq c_{\mathbf{w}} + \sqrt{\bar{n}} + \frac{c_{\mathbf{w}} + n\sqrt{s} \bar{\epsilon} \bar{\epsilon}^{(L-1)2} T}{1 - \bar{\epsilon}} \leq (3/2)c_{\mathbf{w}} + \sqrt{\bar{n}}, \quad (\text{B.42})$$

where we get the last inequality by choosing $L \geq 1$ via

$$\begin{aligned} \frac{c_{\mathbf{w}} + n\sqrt{s} \bar{\epsilon} \bar{\epsilon}^{(L-1)2} T}{1 - \bar{\epsilon}} \leq c_{\mathbf{w}} + \sqrt{\bar{n}}/2 &\iff \frac{(L-1)^2}{\bar{\epsilon}} \leq \frac{c_{\mathbf{w}} + \sqrt{\bar{n}}/2}{2\sqrt{ns} \bar{\epsilon} T}, \\ &\iff L \geq 1 + \frac{2 \log(2\sqrt{ns} \bar{\epsilon} T) + c_{\mathbf{w}} + \sqrt{\bar{n}}/2}{1 - \bar{\epsilon}}. \end{aligned} \quad (\text{B.43})$$

This also implies that, with probability at least $1 - \epsilon$ over the modes, for all $k \in S_i^{(\cdot)}$ and all $i \in [S]$, we have $\|\mathbf{x}_k - \bar{\mathbf{x}}_k\| \leq (1/2)c_{\mathbf{w}} + \sqrt{\bar{n}}$. This completes the proof. \blacksquare

By construction, conditioned on the modes, $\bar{\mathbf{x}}_k = \mathbf{x}_{k, k-L+1}$ only depends on the excitation and noise $\{\mathbf{z}_t, \mathbf{w}_t\}_{t=k-L+1}^{+kL-1}$. Note that the dependence ranges $[k-L+1, kL-1]$ are disjoint intervals for each (k, i) pairs. Hence, $\{\bar{\mathbf{x}}_k\}_{k \in S_i^{(\cdot)}}$ should all be independent of each other. However, this is not the case because $\{\bar{\mathbf{x}}_k\}_{k \in S_i^{(\cdot)}}$ are obtained by truncating only bounded states $\{\mathbf{x}_k\}_{k \in S_i^{(\cdot)}}$. Therefore, we will look for a subset of independent truncated states within $\{\bar{\mathbf{x}}_k\}_{k \in S_i^{(\cdot)}}$, as follows.

Definition B.8 (Subset of bounded states) Let sampling period $L \geq 1$ be an integer. Let $k = +kL$ be the sub-sampling indices, where $0 \leq k \leq L-1$ is a fixed offset and $k = 1, 2, \dots, \lfloor \frac{T-L}{L} \rfloor$. We sub-sample the trajectory $(\mathbf{x}_t, \mathbf{z}_t, \mathbf{w}_t)_{t=0}^T$ at time indices $k \in \bar{S}_i^{(\cdot)}$, where

$$\bar{S}_i^{(\cdot)} := \{k \mid (k) = i, \|\mathbf{x}_k\| \leq c_{\mathbf{w}} + \sqrt{\bar{n}}, \|\bar{\mathbf{x}}_k\| \leq (1/2)c_{\mathbf{w}} + \sqrt{\bar{n}}, \|\mathbf{z}_k\| \leq c_{\mathbf{z}}\sqrt{\bar{\rho}}\}, \quad (\text{B.44})$$

to obtain a subset of the i th sub-trajectory, denoted by $\{(\mathbf{x}_{k+1}, \bar{\mathbf{x}}_k, \mathbf{z}_k, \mathbf{w}_k)_{k \in \bar{S}_i^{(\cdot)}}\}$.

Next, we show that, conditioned on the modes, the samples in $\{(\bar{\mathbf{x}}_{k+1}, \bar{\mathbf{x}}_k, \mathbf{z}_k, \mathbf{w}_k)_{k \in \bar{S}_i^{(\cdot)}}\}$ are independent.

Lemma B.9 (Conditional independence) Consider the MJS (B.7). Suppose $\{\mathbf{z}_t\}_{t=0}^{i.i.d.} \mathcal{N}(0, \frac{2}{L}\mathbf{I}_p)$ and $\{\mathbf{w}_t\}_{t=0}^{i.i.d.} \mathcal{D}_{\mathbf{w}}$ satisfies Assumption 3. Suppose the sampling period $L \geq 1$ satisfies (B.36). Let $S_i^{(\cdot)}$ and $\bar{S}_i^{(\cdot)}$ be as in Def. B.3 and B.8 respectively. Then, with probability at least $1 - \epsilon$ over the mode, we have, (a) $\{\bar{\mathbf{x}}_k\}_{k \in \bar{S}_i^{(\cdot)}}$ is a subset of $\{\bar{\mathbf{x}}_k\}_{k \in S_i^{(\cdot)}}$, (b) conditioned on the modes, $\{\bar{\mathbf{x}}_k\}_{k \in \bar{S}_i^{(\cdot)}}$ are all independent, (c) conditioned on the modes, $\{\bar{\mathbf{x}}_k\}_{k \in \bar{S}_i^{(\cdot)}}$, $\{\mathbf{z}_k\}_{k \in \bar{S}_i^{(\cdot)}}$ and $\{\mathbf{w}_k\}_{k \in \bar{S}_i^{(\cdot)}}$ are all independent of each other.

Proof The first statement is a direct implication of Def. B.3 and B.8. To prove the second statement, consider $\{\mathbf{x}_k\}_{k \in S_i^{(\cdot)}}$ which contains states bounded by $c_{\mathbf{w}} + \sqrt{\bar{n}}$. From (B.8), observe that, each state can be decomposed into $\mathbf{x}_k = \bar{\mathbf{x}}_k + \tilde{\mathbf{x}}_k$, where $\bar{\mathbf{x}}_k$ is the truncated state and $\tilde{\mathbf{x}}_k$ captures the impact of the past states at time index $k-L+1$. When the sampling period L satisfies (B.36), then from Lemma B.7, with probability at least $1 - \epsilon$ over the modes, for all $k \in S_i^{(\cdot)}$ and all $i \in [S]$, we have $\|\tilde{\mathbf{x}}_k\| \leq (1/2)c_{\mathbf{w}} + \sqrt{\bar{n}}$. Furthermore, from Def. B.8, for all $k \in \bar{S}_i^{(\cdot)}$ and all $i \in [S]$, we have $\|\bar{\mathbf{x}}_k\| \leq (1/2)c_{\mathbf{w}} + \sqrt{\bar{n}}$. Combining these results, with probability at least $1 - \epsilon$ over the modes, for all $k \in \bar{S}_i^{(\cdot)}$ and all $i \in [S]$, we have,

$$\|\mathbf{x}_k\| \leq \|\bar{\mathbf{x}}_k\| + \|\tilde{\mathbf{x}}_k\| \leq c_{\mathbf{w}} + \sqrt{\bar{n}}. \quad (\text{B.45})$$

Secondly, by construction, conditioned on the modes, $\bar{\mathbf{x}}_k = \mathbf{x}_{k, k-L+1}$ only depends on the excitation and noise $\{\mathbf{z}_t, \mathbf{w}_t\}_{t=k-L+1}^{+kL-1}$. Note that the dependence ranges $[k-L+1, kL-1]$ are disjoint intervals for

each (k, k) pairs. Hence, conditioned on the modes, the samples in the set $\{\bar{\mathbf{x}}_k\}_{k \in \bar{S}_i^{(j)}}$ are all independent of each other.

To show the independence of $\{\bar{\mathbf{x}}_k\}_{k \in \bar{S}_i^{(j)}}$ and $\{\mathbf{z}_k\}_{k \in \bar{S}_i^{(j)}}$, observe that $\mathbf{z}_k = \mathbf{z}_{+kL}$ have timestamps $+kL$; which is not covered by $[+kL+1, +kL-1]$ – the dependence ranges of $\{\bar{\mathbf{x}}_k\}_{k \in \bar{S}_i^{(j)}}$. Identical argument shows the independence of $\{\bar{\mathbf{x}}_k\}_{k \in \bar{S}_i^{(j)}}$ and $\{\mathbf{w}_k\}_{k \in \bar{S}_i^{(j)}}$. Lastly, $\{\mathbf{z}_k\}_{k \in \bar{S}_i^{(j)}}$ and $\{\mathbf{w}_k\}_{k \in \bar{S}_i^{(j)}}$ are independent of each other by definition. Hence, $\{\bar{\mathbf{x}}_k\}_{k \in \bar{S}_i^{(j)}}$, $\{\mathbf{z}_k\}_{k \in \bar{S}_i^{(j)}}$ and $\{\mathbf{w}_k\}_{k \in \bar{S}_i^{(j)}}$ are all independent of each other. This completes the proof. \blacksquare

Next, we state a lemma similar to Lemma B.4 to show that the truncated states have nice covariance properties.

Lemma B.10 (Covariance of truncated states $\{\bar{\mathbf{x}}_k\}_{k \in \bar{S}_i^{(j)}}$) *Consider the setup of Corollary B.2. Let t_0 , σ_w and σ_z be as in (B.9), (B.10) and (B.34) respectively. Let $c \geq 6$ be a fixed constant and $\{\bar{\mathbf{x}}_k\}_{k \in \bar{S}_i^{(j)}}$ be as in Lemma B.9. Define*

$$L_{cov}(\bar{\zeta}) := 1 + \frac{2 \log(8c^2 \sigma_w^2 + n\sqrt{ns} \bar{\zeta})}{1 - \bar{\zeta}}, \quad (\text{B.46})$$

and suppose, the sampling period $L \geq 1$ obeys,

$$L \geq \max\{t_0, L_{cov}(\bar{\zeta})\}. \quad (\text{B.47})$$

Then, for all $k \in \bar{S}_i^{(j)}$, we have

$$(\sigma_w^2/4)I_n \leq [\bar{\mathbf{x}}_k] \leq 2c^2 \frac{\sigma_w^2}{\sigma_w^2} n I_n. \quad (\text{B.48})$$

$$(1/4)I_{n+p} \leq [\bar{\mathbf{h}}_k] \leq 2c^2(1 + \frac{\sigma_z^2}{\sigma_w^2})(n+p)I_{n+p}. \quad (\text{B.49})$$

Proof To begin, for all $k \in \bar{S}_i^{(j)}$, we upper bound the difference between the covariance of truncated and non-truncated states as follows,

$$\begin{aligned} \|\mathbb{E}[\bar{\mathbf{x}}_k \bar{\mathbf{x}}_k - \mathbf{x}_k \mathbf{x}_k]\| &= \|\mathbb{E}[\bar{\mathbf{x}}_k \bar{\mathbf{x}}_k - \mathbf{x}_k \bar{\mathbf{x}}_k + \mathbf{x}_k \bar{\mathbf{x}}_k - \mathbf{x}_k \mathbf{x}_k]\|, \\ &\leq \mathbb{E}[\|\bar{\mathbf{x}}_k\| \|\mathbf{x}_k - \bar{\mathbf{x}}_k\|] + \mathbb{E}[\|\mathbf{x}_k\| \|\mathbf{x}_k - \bar{\mathbf{x}}_k\|], \\ &\leq (1/2)c \frac{\sigma_w^2}{\sigma_w^2} + \sqrt{n} \mathbb{E}[\|\mathbf{x}_k - \bar{\mathbf{x}}_k\|] + c \frac{\sigma_w^2}{\sigma_w^2} + \sqrt{n} \mathbb{E}[\|\mathbf{x}_k - \bar{\mathbf{x}}_k\|], \\ &\leq 2c^2 \frac{\sigma_w^2}{\sigma_w^2} + n\sqrt{ns} \bar{\zeta} \bar{\zeta}^{(L-1)2}, \end{aligned} \quad (\text{B.50})$$

where we used Def. B.8 to obtain the second last inequality and (B.40) to obtain the last inequality. Combining this with Lemma B.4, for all $k \in \bar{S}_i^{(j)}$, assuming $c \geq 6$, we have

$$\begin{aligned} \min([\bar{\mathbf{x}}_k]) &\geq \min([\mathbf{x}_k]) - \|\mathbb{E}[\bar{\mathbf{x}}_k \bar{\mathbf{x}}_k - \mathbf{x}_k \mathbf{x}_k]\|, \\ &\geq \frac{\sigma_w^2}{\sigma_w^2}/2 - 2c^2 \frac{\sigma_w^2}{\sigma_w^2} + n\sqrt{ns} \bar{\zeta} \bar{\zeta}^{(L-1)2} \geq \frac{\sigma_w^2}{\sigma_w^2}/4, \end{aligned} \quad (\text{B.51})$$

where we get the last inequality by choosing $L \geq 1$ via

$$\begin{aligned} \frac{\sigma_w^2}{\sigma_w^2}/4 \geq 2c^2 \frac{\sigma_w^2}{\sigma_w^2} + n\sqrt{ns} \bar{\zeta} \bar{\zeta}^{(L-1)2} &\iff \bar{\zeta} \bar{\zeta}^{(L-1)2} \leq \frac{1}{8c^2 \frac{\sigma_w^2}{\sigma_w^2} + n\sqrt{ns} \bar{\zeta}}, \\ &\iff L \geq 1 + \frac{2 \log(8c^2 \frac{\sigma_w^2}{\sigma_w^2} + n\sqrt{ns} \bar{\zeta})}{1 - \bar{\zeta}}. \end{aligned} \quad (\text{B.52})$$

This also implies that we have the following upper bound on the covariance spectral norm, that is, for all $k \in \bar{S}_i^{(j)}$, assuming $c \geq 6$, we have

$$\|\bar{\mathbf{x}}_k\| \leq \|\mathbf{x}_k\| + \|\mathbb{E}[\bar{\mathbf{x}}_k \bar{\mathbf{x}}_k - \mathbf{x}_k \mathbf{x}_k]\| \leq c^2 \frac{\sigma_w^2}{\sigma_w^2} n + \frac{\sigma_w^2}{\sigma_w^2}/4 \leq 2c^2 \frac{\sigma_w^2}{\sigma_w^2} n. \quad (\text{B.53})$$

Combining (B.51) and (B.53), we get the first statement of the lemma, which is combined with (B.23) to obtain the second statement of the lemma. This completes the proof. \blacksquare

To proceed, consider the t th truncated sub-trajectory $\{(\bar{\mathbf{x}}_{k+1}, \bar{\mathbf{x}}_k, \mathbf{z}_k, (k))\}_{k \in \bar{S}_i^{(t)}}$ given by Def. B.6.

Let $\bar{\mathbf{h}}_k := [\frac{1}{\mathbf{w}} \bar{\mathbf{x}}_k, \frac{1}{\mathbf{z}} \mathbf{z}_k]$. Similar to (B.26), we construct $\bar{\mathbf{Y}}_i^{(t)}$, $\bar{\mathbf{H}}_i^{(t)}$ and $\mathbf{W}_i^{(t)}$ by (row-wise) stacking $\{\bar{\mathbf{x}}_{k+1}\}_{k \in \bar{S}_i^{(t)}}$, $\{\bar{\mathbf{h}}_k\}_{k \in \bar{S}_i^{(t)}}$ and $\{\mathbf{w}_k\}_{k \in \bar{S}_i^{(t)}}$ respectively. As an intermediate step, we lower bound $\min(\bar{\mathbf{H}}_i^{(t)} \bar{\mathbf{H}}_i^{(t)})$ and upper bound $\|\bar{\mathbf{H}}_i^{(t)} \mathbf{W}_i^{(t)}\|$. This will, in turn, allow us to lower and upper bound the non-truncated terms $\min(\mathbf{H}_i^{(t)} \mathbf{H}_i^{(t)})$ and $\|\mathbf{H}_i^{(t)} \mathbf{W}_i^{(t)}\|$ respectively via,

$$\min(\mathbf{H}_i^{(t)} \mathbf{H}_i^{(t)}) \geq \min(\bar{\mathbf{H}}_i^{(t)} \bar{\mathbf{H}}_i^{(t)}) - \|\mathbf{H}_i^{(t)} \mathbf{H}_i^{(t)} - \bar{\mathbf{H}}_i^{(t)} \bar{\mathbf{H}}_i^{(t)}\|, \quad (\text{B.54})$$

$$\|\mathbf{H}_i^{(t)} \mathbf{W}_i^{(t)}\| \leq \|\bar{\mathbf{H}}_i^{(t)} \mathbf{W}_i^{(t)}\| + \|\mathbf{H}_i^{(t)} \mathbf{W}_i^{(t)} - \bar{\mathbf{H}}_i^{(t)} \mathbf{W}_i^{(t)}\|. \quad (\text{B.55})$$

For this purpose, our next lemma lower bounds the eigenvalues of the matrix $\bar{\mathbf{H}}_i^{(t)} \bar{\mathbf{H}}_i^{(t)}$ and upper bounds the error term $\|\bar{\mathbf{H}}_i^{(t)} \mathbf{W}_i^{(t)}\|$.

Theorem B.11 (Bounding $\min(\bar{\mathbf{H}}_i^{(t)} \bar{\mathbf{H}}_i^{(t)})$ and $\|\bar{\mathbf{H}}_i^{(t)} \mathbf{W}_i^{(t)}\|$) Consider the setup of Corollary B.2. Let $t_0, \alpha, \beta, L_{\text{tr1}}(\underline{\mathbf{c}})$ and $L_{\text{cov}}(\underline{\mathbf{c}})$ be as in (B.9), (B.10), (B.34), (B.35) and (B.46), respectively. Let $C, C_0 > 0$ and $c \geq 6$ be fixed constants. Let $\bar{\mathbf{H}}_i^{(t)}$, $\tilde{\mathbf{H}}_i^{(t)}$ and $\mathbf{W}_i^{(t)}$ be constructed by (row-wise) stacking $\{\bar{\mathbf{h}}_k\}_{k \in \bar{S}_i^{(t)}}$, $\{\tilde{\mathbf{h}}_k\}_{k \in \tilde{S}_i^{(t)}}$ and $\{\mathbf{w}_k\}_{k \in \bar{S}_i^{(t)}}$ respectively. Suppose the sampling period $L \geq 1$ and the number of independent samples $|\bar{S}_i^{(t)}|$ satisfy the following lower bounds,

$$L \geq \max\{t_0, L_{\text{tr1}}(\underline{\mathbf{c}}), L_{\text{cov}}(\underline{\mathbf{c}})\}, \quad (\text{B.56})$$

$$|\bar{S}_i^{(t)}| \geq 16c^2(1 + \alpha) \log\left(\frac{2s(n+p)}{\beta}\right)(n+p). \quad (\text{B.57})$$

Then, with probability at least $1 - 3^{-s}$, for all $i \in [s]$, we have

$$\min(\bar{\mathbf{H}}_i^{(t)} \bar{\mathbf{H}}_i^{(t)}) \geq \min(\tilde{\mathbf{H}}_i^{(t)} \tilde{\mathbf{H}}_i^{(t)}) \geq \frac{|\bar{S}_i^{(t)}|}{16}, \quad (\text{B.58})$$

$$\|\bar{\mathbf{H}}_i^{(t)} \mathbf{W}_i^{(t)}\| \leq 2c(1 + \alpha) \sqrt{|\bar{S}_i^{(t)}|(n+p)} \left(C_{\mathbf{w}} \sqrt{n+p} + C_0 \sqrt{\log\left(\frac{2s}{\beta}\right)} \right). \quad (\text{B.59})$$

Proof To begin, recall that not all the rows of $\bar{\mathbf{H}}_i^{(t)}$ are independent. Therefore, to lower bound $\min(\bar{\mathbf{H}}_i^{(t)} \bar{\mathbf{H}}_i^{(t)})$, we first consider the matrix $\tilde{\mathbf{H}}_i^{(t)}$ which is constructed by (row-wise) stacking $\{\tilde{\mathbf{h}}_k\}_{k \in \tilde{S}_i^{(t)}}$.

Observe that, conditioned on the modes, the matrix $\tilde{\mathbf{H}}_i^{(t)}$, which is a sub-matrix of $\bar{\mathbf{H}}_i^{(t)}$, has independent rows from Lemma B.9.

• **Lower bounding $\min(\tilde{\mathbf{H}}_i^{(t)})$:** Using Lemma B.9, we observe that, conditioned on the modes, the rows of $\tilde{\mathbf{H}}_i^{(t)}$ are all independent. Secondly, by definition, each row of $\tilde{\mathbf{H}}_i^{(t)}$ can be deterministically bounded as follows: for all $k \in \tilde{S}_i^{(t)}$, we have

$$\|\tilde{\mathbf{h}}_k\|^2 \leq \frac{1}{\mathbf{w}} \|\bar{\mathbf{x}}_k\|^2 + \frac{1}{\mathbf{z}} \|\mathbf{z}_k\|^2 \leq (1/4)c^2 \alpha n + c^2 p \leq c^2(1 + \alpha)(n+p). \quad (\text{B.60})$$

Thirdly, from Lemma B.10, when $c \geq 6$ and $L \geq \max\{t_0, L_{\text{cov}}(\underline{\mathbf{c}})\}$, then for all $k \in \tilde{S}_i^{(t)}$, we have

$$(1/4)I_{n+p} \leq [\tilde{\mathbf{h}}_k] \leq 2c^2(1 + \alpha)(n+p)I_{n+p}. \quad (\text{B.61})$$

Therefore, we can use Corollary A.9 to lower bound $\min(\tilde{\mathbf{H}}_i^{(t)})$. Specifically, using Corollary A.9 with $\min = 1/2$ and $m = c^2(1 + \alpha)(n+p)$, with probability at least $1 - 3^{-s}$, for all $i \in [s]$, we have

$$\min(\tilde{\mathbf{H}}_i^{(t)}) \geq \frac{\sqrt{|\tilde{S}_i^{(t)}|}}{2} - c \sqrt{(1 + \alpha)(n+p) \log\left(\frac{2s(n+p)}{\beta}\right)} \geq \frac{\sqrt{|\tilde{S}_i^{(t)}|}}{4}, \quad (\text{B.62})$$

as long as $|\bar{S}_i^{(\cdot)}|$ satisfies the following lower bound,

$$\begin{aligned} \frac{\sqrt{|\bar{S}_i^{(\cdot)}|}}{4} &\geq c\sqrt{(1+\frac{2}{\epsilon})(n+p)\log\left(\frac{2s(n+p)}{\epsilon}\right)} \\ \iff |\bar{S}_i^{(\cdot)}| &\geq 16c^2(1+\frac{2}{\epsilon})(n+p)\log\left(\frac{2s(n+p)}{\epsilon}\right). \end{aligned} \quad (\text{B.63})$$

• **Lower bounding** $\min(\bar{\mathbf{H}}_i^{(\cdot)} \bar{\mathbf{H}}_i^{(\cdot)})$: Using Lemma B.9, we have, $\{\bar{\mathbf{h}}_\kappa\}_{\kappa \in \bar{S}_i^{(\cdot)}}$ is a subset of $\{\bar{\mathbf{h}}_\kappa\}_{\kappa \in S_i^{(\cdot)}}$. As a result, (B.62) also implies that, with probability at least $1 - \frac{\epsilon}{2}$, for all $i \in [S]$, we have

$$\lambda_{\min}(\bar{\mathbf{H}}_i^{(\cdot)}) \geq \lambda_{\min}(\tilde{\mathbf{H}}_i^{(\cdot)}) \geq \frac{\sqrt{|\bar{S}_i^{(\cdot)}|}}{4} \implies \min(\bar{\mathbf{H}}_i^{(\cdot)} \bar{\mathbf{H}}_i^{(\cdot)}) \geq \frac{|\bar{S}_i^{(\cdot)}|}{16}, \quad (\text{B.64})$$

as long as $|\bar{S}_i^{(\cdot)}|$ satisfies the lower bound in (B.63).

• **Upper bounding** $\|\bar{\mathbf{H}}_i^{(\cdot)} \mathbf{W}_i^{(\cdot)}\|$: Using Lemma B.7, when $L \geq \max\{t_0, L_{tr1}(\frac{2s}{\epsilon})\}$, with probability at least $1 - \frac{\epsilon}{2}$ over the modes, for all $\kappa \in S_i^{(\cdot)}$ and all $i \in [S]$, we have, $\|\bar{\mathbf{h}}_\kappa\|^2 \leq c^2(1 + (9/4)\frac{2}{\epsilon})(n+p)$. This implies that, with probability at least $1 - \frac{\epsilon}{2}$ over the modes, for all $i \in [S]$, we have

$$\|\bar{\mathbf{H}}_i^{(\cdot)}\| \leq \|\bar{\mathbf{H}}_i^{(\cdot)}\|_F \leq c(1 + 2\frac{2}{\epsilon})\sqrt{|S_i^{(\cdot)}|(n+p)} \leq 2c(1 + \frac{2}{\epsilon})\sqrt{|S_i^{(\cdot)}|(n+p)}.$$

To proceed, let $\bar{\mathbf{H}}_i^{(\cdot)}$ have singular value decomposition $\mathbf{U} \mathbf{V}$ with $\|\mathbf{V}\| \leq 2c(1 + \frac{2}{\epsilon})\sqrt{|S_i^{(\cdot)}|(n+p)}$. Since $\mathbf{W}_i^{(\cdot)}$ has i.i.d. \mathbf{w} -subGaussian entries (Assumption 3), $\mathbf{U} \mathbf{W}_i^{(\cdot)} \in \mathbb{R}^{(n+p) \times n}$ has i.i.d. \mathbf{w} -subGaussian columns. As a result, applying Theorem 5.39 of [65], with probability at least $1 - \frac{\epsilon}{2}$, for all $i \in [S]$, we have

$$\|\mathbf{U} \mathbf{W}_i^{(\cdot)}\| \leq C \mathbf{w}\sqrt{n+p} + C_0\sqrt{\log\left(\frac{2s}{\epsilon}\right)}. \quad (\text{B.65})$$

This implies that, with probability at least $1 - 2\frac{\epsilon}{2}$, for all $i \in [S]$, we have

$$\|\bar{\mathbf{H}}_i^{(\cdot)} \mathbf{W}_i^{(\cdot)}\| \leq \|\mathbf{U} \mathbf{W}_i^{(\cdot)}\| \leq 2c(1 + \frac{2}{\epsilon})\sqrt{|S_i^{(\cdot)}|(n+p)} \left(C \mathbf{w}\sqrt{n+p} + C_0\sqrt{\log\left(\frac{2s}{\epsilon}\right)} \right). \quad (\text{B.66})$$

This completes the proof. ■

B.2.3 Estimation from Non-truncated States

Coming back to the original problem of estimating the unknown dynamics from dependent samples, observe that the estimation error (B.28) in the case of single trajectory can be upper bounded as follows,

$$\begin{aligned} \|\hat{\mathbf{A}}_i - \mathbf{A}_i\| &\leq \frac{\sum_{l=0}^{L-1} \|\mathbf{H}_i^{(l)} \mathbf{W}_i^{(l)}\|}{\sum_{l=0}^{L-1} \min(\mathbf{H}_i^{(l)} \mathbf{H}_i^{(l)})'} \\ &\leq \frac{\sum_{l=0}^{L-1} (\|\bar{\mathbf{H}}_i^{(l)} \mathbf{W}_i^{(l)}\| + \|\mathbf{H}_i^{(l)} \mathbf{W}_i^{(l)} - \bar{\mathbf{H}}_i^{(l)} \mathbf{W}_i^{(l)}\|)}{\sum_{l=0}^{L-1} (\min(\bar{\mathbf{H}}_i^{(l)} \bar{\mathbf{H}}_i^{(l)}) - \|\mathbf{H}_i^{(l)} \mathbf{H}_i^{(l)} - \bar{\mathbf{H}}_i^{(l)} \bar{\mathbf{H}}_i^{(l)}\|)}. \end{aligned} \quad (\text{B.67})$$

Therefore, to upper bound the estimation error $\|\hat{\mathbf{A}}_i - \mathbf{A}_i\|$ with dependent samples, we need to upper bound the impact of truncation, captured by $\|\mathbf{H}_i^{(l)} \mathbf{W}_i^{(l)} - \bar{\mathbf{H}}_i^{(l)} \mathbf{W}_i^{(l)}\|$ and $\|\mathbf{H}_i^{(l)} \mathbf{H}_i^{(l)} - \bar{\mathbf{H}}_i^{(l)} \bar{\mathbf{H}}_i^{(l)}\|$ for all $i \in [S]$ and all $0 \leq l \leq L-1$. This is done in the following theorem.

Theorem B.12 (Small impact of truncation) Consider the same setup of Corollary B.2. Let $t_0, \epsilon, \delta, \eta$ and $L_{tr1}(\underline{L}, \epsilon)$ be as in (B.9), (B.10), (B.34) and (B.35), respectively. Suppose the sampling period L obeys $L \geq \max\{t_0, L_{tr1}(\underline{L}, \epsilon)\}$. Let $\mathbf{H}_i^{(\cdot)}$ and $\mathbf{W}_i^{(\cdot)}$ be constructed by (row-wise) stacking $\{\mathbf{h}_k\}_{k \in S_i^{(\cdot)}}$ and $\{\mathbf{w}_k\}_{k \in S_i^{(\cdot)}}$ respectively. Then, with probability at least $1 - \delta$ over the modes, for all $i \in [S]$, we have

$$\|\mathbf{H}_i^{(\cdot)} \mathbf{H}_i^{(\cdot)} - \bar{\mathbf{H}}_i^{(\cdot)} \bar{\mathbf{H}}_i^{(\cdot)}\| \leq \frac{3c^2 \epsilon (1 + \epsilon) \underline{L} \underline{L}^{(L-1)2} n \sqrt{s(n+p)} |S_i^{(\cdot)}| T}{\epsilon}, \quad (\text{B.68})$$

$$\|\mathbf{H}_i^{(\cdot)} \mathbf{W}_i^{(\cdot)} - \bar{\mathbf{H}}_i^{(\cdot)} \mathbf{W}_i^{(\cdot)}\| \leq \frac{cc_{\mathbf{w}} \epsilon \epsilon + \underline{L} \underline{L}^{(L-1)2} n \sqrt{ns} |S_i^{(\cdot)}| T}{\epsilon}. \quad (\text{B.69})$$

Proof We begin by simplifying the the term $\|\mathbf{H}_i^{(\cdot)} \mathbf{H}_i^{(\cdot)} - \bar{\mathbf{H}}_i^{(\cdot)} \bar{\mathbf{H}}_i^{(\cdot)}\|$ as follows,

$$\begin{aligned} \|\mathbf{H}_i^{(\cdot)} \mathbf{H}_i^{(\cdot)} - \bar{\mathbf{H}}_i^{(\cdot)} \bar{\mathbf{H}}_i^{(\cdot)}\| &= \left\| \sum_{k \in S_i^{(\cdot)}} (\mathbf{h}_k \mathbf{h}_k - \bar{\mathbf{h}}_k \bar{\mathbf{h}}_k) \right\|, \\ &\leq |S_i^{(\cdot)}| \max_{k \in S_i^{(\cdot)}} \|\mathbf{h}_k \mathbf{h}_k - \bar{\mathbf{h}}_k \bar{\mathbf{h}}_k\|, \\ &= |S_i^{(\cdot)}| \max_{k \in S_i^{(\cdot)}} \|\mathbf{h}_k \mathbf{h}_k - \mathbf{h}_k \bar{\mathbf{h}}_k + \mathbf{h}_k \bar{\mathbf{h}}_k - \bar{\mathbf{h}}_k \bar{\mathbf{h}}_k\|, \\ &\leq |S_i^{(\cdot)}| \max_{k \in S_i^{(\cdot)}} (\|\mathbf{h}_k\| \|\mathbf{h}_k - \bar{\mathbf{h}}_k\| + \|\bar{\mathbf{h}}_k\| \|\mathbf{h}_k - \bar{\mathbf{h}}_k\|). \end{aligned} \quad (\text{B.70})$$

We will upper bound each of these terms separately and combine them together in (B.70) to get the desired upper bound. First of all, observe that, for all $i \in [S]$, each row of $\mathbf{H}_i^{(\cdot)}$ is deterministically bounded as follows,

$$\|\mathbf{h}_k\|^2 \leq \frac{1}{2} \|\mathbf{x}_k\|^2 + \frac{1}{2} \|\mathbf{z}_k\|^2 \leq c^2 \epsilon^2 n + c^2 p \leq c^2 (1 + \epsilon^2) (n + p). \quad (\text{B.71})$$

Similarly, from Lemma B.7, when $L \geq \max\{t_0, L_{tr1}(\underline{L}, \epsilon)\}$, we observe that, with probability at least $1 - \delta$ over the modes, for all $i \in [S]$, each row of $\bar{\mathbf{H}}_i^{(\cdot)}$ can be bounded as follows,

$$\|\bar{\mathbf{h}}_k\|^2 \leq \frac{1}{2} \|\bar{\mathbf{x}}_k\|^2 + \frac{1}{2} \|\mathbf{z}_k\|^2 \leq c^2 (9/4) \epsilon^2 n + c^2 p \leq 4c^2 (1 + \epsilon^2) (n + p). \quad (\text{B.72})$$

To proceed, recall from (B.41) that, with probability at least $1 - \delta$ over the modes, for all $k \in S_i^{(\cdot)}$ and all $i \in [S]$, we have

$$\|\mathbf{h}_k - \bar{\mathbf{h}}_k\| = \left\| \begin{bmatrix} \frac{1}{\epsilon} \mathbf{x}_k \\ \frac{1}{\epsilon} \mathbf{z}_k \end{bmatrix} - \begin{bmatrix} \frac{1}{\epsilon} \bar{\mathbf{x}}_k \\ \frac{1}{\epsilon} \mathbf{z}_k \end{bmatrix} \right\|_2 = \frac{1}{\epsilon} \|\mathbf{x}_k - \bar{\mathbf{x}}_k\| \leq \frac{c \epsilon + n \sqrt{s} \epsilon \underline{L} \underline{L}^{(L-1)2} T}{\epsilon}. \quad (\text{B.73})$$

Combining (B.71), (B.72) and (B.73) into (B.70), with probability at least $1 - \delta$ over the modes, for all $i \in [S]$, we have

$$\|\mathbf{H}_i^{(\cdot)} \mathbf{H}_i^{(\cdot)} - \bar{\mathbf{H}}_i^{(\cdot)} \bar{\mathbf{H}}_i^{(\cdot)}\| \leq \frac{3c^2 \epsilon (1 + \epsilon) \underline{L} \underline{L}^{(L-1)2} n \sqrt{s(n+p)} |S_i^{(\cdot)}| T}{\epsilon}. \quad (\text{B.74})$$

Using a similar line of reasoning, with probability at least $1 - \delta$ over the modes, for all $i \in [S]$, we also have

$$\begin{aligned} \|\mathbf{H}_i^{(\cdot)} \mathbf{W}_i^{(\cdot)} - \bar{\mathbf{H}}_i^{(\cdot)} \mathbf{W}_i^{(\cdot)}\| &= \left\| \sum_{k \in S_i^{(\cdot)}} (\mathbf{h}_k \mathbf{w}_k - \bar{\mathbf{h}}_k \mathbf{w}_k) \right\|, \\ &\leq |S_i^{(\cdot)}| \max_{k \in S_i^{(\cdot)}} \|\mathbf{h}_k - \bar{\mathbf{h}}_k\| \|\mathbf{w}_{(j,k)}\|, \\ &\leq \frac{cc_{\mathbf{w}} \epsilon \epsilon + \underline{L} \underline{L}^{(L-1)2} n \sqrt{ns} |S_i^{(\cdot)}| T}{\epsilon}. \end{aligned} \quad (\text{B.75})$$

This completes the proof. \blacksquare

Combining Theorems B.11 and B.12, we obtain our result on the estimation of MJS in (B.7) from finite samples obtained from a single trajectory.

Theorem B.13 (Learning with bounded noise) *Consider the same setup of Corollary B.2. Let $t_0, \dots, L_{tr1}(\bar{\zeta}, \bar{\nu})$ and $L_{cov}(\bar{\zeta})$ be as in (B.9), (B.10), (B.34), (B.35) and (B.46), respectively. Let $S_i^{(\cdot)}$ and $\bar{S}_i^{(\cdot)}$ be as in Def. B.3 and B.8 respectively and assume $|\bar{S}_i^{(\cdot)}| \geq \frac{\min T}{2L}$, for all $i \in [s]$, with probability at least $1 - \epsilon$. Suppose $\|\mathbf{K}_{1s}\| \leq C_K$ for some constant $C_K > 0$. Let $C, C_0 > 0$, and $c \geq 6$ be fixed constants. Define*

$$L_{tr2}(\bar{\zeta}, \bar{\nu}) := 1 + \frac{2}{(1 - \bar{\zeta})} \log \left(\frac{192c^2 \bar{\zeta} + (1 + \bar{\nu})n\sqrt{s(n+p)}T}{\min} \right), \quad (\text{B.76})$$

$$L_{tr3}(\bar{\zeta}, \bar{\nu}) := 1 + \frac{2}{(1 - \bar{\zeta})} \log \left(\frac{c_{\mathbf{w}} \mathbf{w} + \bar{\zeta} n\sqrt{n\bar{s}}T^2}{(1 + \bar{\nu})\sqrt{(n+p)}(C_{\mathbf{w}}\sqrt{n+p} + C_0\sqrt{\log(2s/\epsilon)})} \right). \quad (\text{B.77})$$

Suppose the sampling period L and the trajectory length T satisfy

$$L \geq \max\{t_0, L_{cov}(\bar{\zeta}), L_{tr1}(\bar{\zeta}, \frac{\bar{\nu}}{18L}), L_{tr2}(\bar{\zeta}, \frac{\bar{\nu}}{18L}), L_{tr3}(\bar{\zeta}, \frac{\bar{\nu}}{18L})\} \quad (\text{B.78})$$

$$T \geq \frac{32L}{\min} c^2 (1 + \bar{\nu}) \log \left(\frac{36sL(n+p)}{\min} \right) (n+p). \quad (\text{B.79})$$

Then, solving the least-squares problem (B.27), with probability at least $1 - \epsilon/2$, for all $i \in [s]$, we have

$$\begin{aligned} \|\hat{\mathbf{A}}_i - \mathbf{A}_i\| &\leq \frac{192c(1 + \bar{\nu})}{\min} \frac{(z + C_K \mathbf{w})}{z} \sqrt{\frac{L(n+p)}{T}} \left(C\sqrt{n+p} + (C_0/\mathbf{w})\sqrt{\log\left(\frac{36sL}{\min}\right)} \right), \\ \|\hat{\mathbf{B}}_i - \mathbf{B}_i\| &\leq \frac{192c(1 + \bar{\nu})}{\min} \frac{\mathbf{w}}{z} \sqrt{\frac{L(n+p)}{T}} \left(C\sqrt{n+p} + (C_0/\mathbf{w})\sqrt{\log\left(\frac{36sL}{\min}\right)} \right). \end{aligned}$$

Proof To begin, using Theorem B.11 along-with the assumption made in the statement of the theorem regarding $|\bar{S}_i^{(\cdot)}|$, with probability at least $1 - 4\epsilon$, for all $i \in [s]$, we have

$$\min(\bar{\mathbf{H}}_i^{(\cdot)} \bar{\mathbf{H}}_i^{(\cdot)}) \geq \frac{|\bar{S}_i^{(\cdot)}|}{16} \geq \frac{\min T}{32L}, \quad (\text{B.80})$$

as long as the trajectory length T satisfies the following lower bound,

$$T \geq \frac{32L}{\min} c^2 (1 + \bar{\nu}) \log \left(\frac{2s(n+p)}{\min} \right) (n+p). \quad (\text{B.81})$$

Combining this with Theorem B.12, with probability at least $1 - 5\epsilon$, for all $i \in [s]$, we have

$$\begin{aligned} \min(\mathbf{H}_i^{(\cdot)} \mathbf{H}_i^{(\cdot)}) &\geq \min(\bar{\mathbf{H}}_i^{(\cdot)} \bar{\mathbf{H}}_i^{(\cdot)}) - \|\mathbf{H}_i^{(\cdot)} \mathbf{H}_i^{(\cdot)} - \bar{\mathbf{H}}_i^{(\cdot)} \bar{\mathbf{H}}_i^{(\cdot)}\|, \\ &\geq \frac{\min T}{32L} - \frac{3c^2 + (1 + \bar{\nu}) \bar{\zeta} \bar{\zeta}^{(L-1)2} n\sqrt{s(n+p)}T^2}{L}, \\ &\geq \frac{\min T}{64L}, \end{aligned} \quad (\text{B.82})$$

where we get the last inequality by choosing $L \geq 1$ via,

$$\begin{aligned} \frac{\min T}{64L} &\geq \frac{3c^2 + (1 + \bar{\nu}) \bar{\zeta} \bar{\zeta}^{(L-1)2} n\sqrt{s(n+p)}T^2}{L}, \\ \Leftrightarrow \bar{\zeta}^{(L-1)2} &\leq \frac{\min}{192c^2 \bar{\zeta} + (1 + \bar{\nu})n\sqrt{s(n+p)}T}, \\ \Leftrightarrow L &\geq 1 + \frac{2 \log(192c^2 \bar{\zeta} + (1 + \bar{\nu})n\sqrt{s(n+p)}T) / (\min)}{(1 - \bar{\zeta})}. \end{aligned} \quad (\text{B.83})$$

Similarly, combing Theorems B.11 and B.12, with probability at least $1 - 4^{-\ell}$, for all $i \in [S]$, we also have

$$\begin{aligned} \|\mathbf{H}_i^{(\ell)} \mathbf{W}_i^{(\ell)}\| &\leq \|\bar{\mathbf{H}}_i^{(\ell)} \mathbf{W}_i^{(\ell)}\| + \|\mathbf{H}_i^{(\ell)} \mathbf{W}_i^{(\ell)} - \bar{\mathbf{H}}_i^{(\ell)} \mathbf{W}_i^{(\ell)}\|, \\ &\leq 2c(1 + \epsilon) \sqrt{\frac{T(n+p)}{L}} \left(C_{\mathbf{w}} \sqrt{n+p} + C_0 \sqrt{\log\left(\frac{2S}{\epsilon}\right)} \right) + \frac{c c_{\mathbf{w}} \mathbf{w} + \bar{\epsilon} \bar{L}^{(L-1)^2} n \sqrt{nS} T^2}{L}, \\ &\leq 3c(1 + \epsilon) \sqrt{\frac{T(n+p)}{L}} \left(C_{\mathbf{w}} \sqrt{n+p} + C_0 \sqrt{\log\left(\frac{2S}{\epsilon}\right)} \right), \end{aligned} \quad (\text{B.84})$$

where we get the last inequality by choosing $L \geq 1$ via,

$$\begin{aligned} \frac{c c_{\mathbf{w}} \mathbf{w} + \bar{\epsilon} \bar{L}^{(L-1)^2} n \sqrt{nS} T^2}{L} &\leq c(1 + \epsilon) \sqrt{\frac{T(n+p)}{L}} \left(C_{\mathbf{w}} \sqrt{n+p} + C_0 \sqrt{\log\left(\frac{2S}{\epsilon}\right)} \right) \\ \Leftrightarrow \bar{L}^{(L-1)^2} &\leq \frac{(1 + \epsilon) \sqrt{TL(n+p)} (C_{\mathbf{w}} \sqrt{n+p} + C_0 \sqrt{\log(2S/\epsilon)})}{c_{\mathbf{w}} \mathbf{w} + \bar{\epsilon} n \sqrt{nS} T^2}, \\ \Leftarrow L &\geq 1 + \frac{2}{(1 - \bar{\epsilon})} \log\left(\frac{c_{\mathbf{w}} \mathbf{w} + \bar{\epsilon} n \sqrt{nS} T^2}{(1 + \epsilon) \sqrt{(n+p)} (C_{\mathbf{w}} \sqrt{n+p} + C_0 \sqrt{\log(2S/\epsilon)})} \right). \end{aligned} \quad (\text{B.85})$$

Finally combining (B.82) and (B.84) into (B.28) and union bounding over all $0 \leq \ell \leq L-1$, with probability at least $1 - 9L^{-\ell}$, for all $i \in [S]$, we have

$$\begin{aligned} \|\hat{\mathbf{L}}_i - \mathbf{L}_i\| &\leq \frac{\sum_{\ell=0}^{L-1} \|\mathbf{H}_i^{(\ell)} \mathbf{W}_i^{(\ell)}\|}{\sum_{\ell=0}^{L-1} \min(\mathbf{H}_i^{(\ell)}, \mathbf{H}_i^{(\ell)})}, \\ &\leq \frac{192c(1 + \epsilon)}{\min} \sqrt{\frac{L(n+p)}{T}} \left(C_{\mathbf{w}} \sqrt{n+p} + C_0 \sqrt{\log\left(\frac{2S}{\epsilon}\right)} \right). \end{aligned} \quad (\text{B.86})$$

To proceed, using standard result from linear algebra that the spectral norm of a sub-matrix is upper bounded by the norm of the original matrix, with probability at least $1 - 9L^{-\ell}$, for all $i \in [S]$, we have

$$\begin{aligned} \|\hat{\mathbf{L}}_i - \mathbf{L}_i\| &\leq \frac{192c(1 + \epsilon)}{\min} \sqrt{\frac{L(n+p)}{T}} \left(C_{\mathbf{w}} \sqrt{n+p} + (C_0/\mathbf{w}) \sqrt{\log\left(\frac{2S}{\epsilon}\right)} \right), \\ \|\hat{\mathbf{B}}_i - \mathbf{B}_i\| &\leq \frac{192c(1 + \epsilon)}{\min} \frac{\mathbf{w}}{\mathbf{z}} \sqrt{\frac{L(n+p)}{T}} \left(C_{\mathbf{w}} \sqrt{n+p} + (C_0/\mathbf{w}) \sqrt{\log\left(\frac{2S}{\epsilon}\right)} \right), \\ \Rightarrow \|\hat{\mathbf{A}}_i - \mathbf{A}_i\| &\leq \|\hat{\mathbf{L}}_i - \mathbf{L}_i\| + \|\mathbf{K}_i\| \|\hat{\mathbf{B}}_i - \mathbf{B}_i\|, \\ &\leq \frac{192c(1 + \epsilon)}{\min} \frac{(\mathbf{z} + C_K \mathbf{w})}{\mathbf{z}} \sqrt{\frac{L(n+p)}{T}} \left(C_{\mathbf{w}} \sqrt{n+p} + (C_0/\mathbf{w}) \sqrt{\log\left(\frac{2S}{\epsilon}\right)} \right). \end{aligned} \quad (\text{B.87})$$

Finally replacing ϵ with $\epsilon/(18L)$ we get the statement of the theorem. This completes the proof. \blacksquare

Next, we use the following lemma to relax the Assumption 3 on the noise.

Lemma B.14 (From Bounded to Unbounded Noise) *Let $\mathbf{g} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \frac{2}{\mathbf{w}} \mathbf{I}_{nT})$ and \mathbf{a} be two independent vectors. Let \mathbf{g} be the truncated Gaussian distribution $\mathbf{g} \sim \{\mathbf{g} \mid \|\mathbf{g}\|_{\infty} \leq c_{\mathbf{w}} \mathbf{w}\}$. Let $S_{\mathbf{g}, \mathbf{a}}$ be the indicator function of an event defined on vectors \mathbf{g}, \mathbf{a} s.t.*

$$\mathbb{E}[S_{\mathbf{g}, \mathbf{a}}] \geq 1 - \epsilon/2.$$

That is, the event holds, on the bounded variable \mathbf{g} , with probability at least $1 - \epsilon/2$. Then, if the bound above holds for $c_{\mathbf{w}} > C := \sqrt{2 \log(nT)} + \sqrt{2 \log(2/\epsilon)}$, we also have that

$$\mathbb{E}[S_{\mathbf{g}, \mathbf{a}}] \geq 1 - \epsilon.$$

That is, the probability that event holds on the unbounded variable \mathbf{g} is at least $1 - \epsilon$.

Proof Let E be the event $\{\mathbf{g} \mid \|\mathbf{g}\|_\infty \leq c_{\mathbf{w}} \mathbf{w}\}$. If $c_{\mathbf{w}} \geq \sqrt{2\log(nT)} + \sqrt{2\log(2/\delta)}$, using Gaussian tail bound and the fact that $E[\|\mathbf{g}\|_\infty] \leq \mathbf{w}\sqrt{2\log(nT)}$, observe that $P(E) \geq 1 - e^{-\frac{(c_{\mathbf{w}} \mathbf{w} - E[\|\mathbf{g}\|_\infty])^2}{2}} \geq 1 - \delta/2$. Therefore, we have

$$E[S_{\mathbf{g},\mathbf{a}}] \geq E[S_{\mathbf{g},\mathbf{a}}|E]P(E) = E[S_{\mathbf{g},\mathbf{a}}]P(E) \geq (1 - \delta/2)^2 \geq 1 - \delta.$$

This completes the proof. \blacksquare

Combining Theorem B.13 and Lemma B.14, we get the following result on learning the MJS dynamics when the process noise is Gaussian.

Corollary B.15 (Learning with un-bounded noise) *Consider the same setup of Theorem B.13 except that Assumption 3 is replaced with $\{\mathbf{w}_t\}_{t=0}^{i.i.d.} \mathcal{N}(0, \frac{2}{\mathbf{w}} \mathbf{I}_n)$ and the threshold for bounding the noise satisfies,*

$$c_{\mathbf{w}} \geq \sqrt{2\log(nT)} + \sqrt{2\log(2/\delta)}. \quad (\text{B.88})$$

Suppose $\|\mathbf{B}_{1:s}\| \leq C_B$ for some $C_B > 0$ and the trajectory length T satisfies,

$$T \gtrsim \frac{\underline{\mathbf{z}}\sqrt{sL}}{\min(1 - \underline{\mathbf{z}})} \left((\sqrt{2\log(nT)} + \sqrt{2\log(2/\delta)})^2 + C_{\mathbf{z}}^2 C_B^2 \right) \log\left(\frac{36sL(n+p)}{\delta}\right) (n+p). \quad (\text{B.89})$$

Then, solving the least-squares problem (B.27), with probability at least $1 - \delta$, for all $i \in [s]$, we have

$$\begin{aligned} \|\hat{\mathbf{A}}_i - \mathbf{A}_i\| &\leq \frac{(\mathbf{z} + C_K \mathbf{w})}{\mathbf{z}} \frac{\underline{\mathbf{z}}\sqrt{s}(\sqrt{2\log(nT)} + \sqrt{2\log(2/\delta)} + C_{\mathbf{z}} C_B)}{\min(1 - \underline{\mathbf{z}})} \sqrt{\frac{L(n+p)}{T}} \\ &\quad \times \left(C\sqrt{n+p} + (C_0/\mathbf{w})\sqrt{\log\left(\frac{36sL}{\delta}\right)} \right), \end{aligned} \quad (\text{B.90})$$

$$\begin{aligned} \|\hat{\mathbf{B}}_i - \mathbf{B}_i\| &\leq \frac{\mathbf{w}}{\mathbf{z}} \frac{\underline{\mathbf{z}}\sqrt{s}(\sqrt{2\log(nT)} + \sqrt{2\log(2/\delta)} + C_{\mathbf{z}} C_B)}{\min(1 - \underline{\mathbf{z}})} \sqrt{\frac{L(n+p)}{T}} \\ &\quad \times \left(C\sqrt{n+p} + (C_0/\mathbf{w})\sqrt{\log\left(\frac{36sL}{\delta}\right)} \right). \end{aligned} \quad (\text{B.91})$$

At this point, we are only left with verifying the assumption that, for all $i \in [s]$, with probability at least $1 - \delta$, we have $|\bar{S}_i^{(\cdot)}| \geq \frac{\min T}{2L}$ for some choice of L and T . In the following, we will state a lemma to show that the above assumption indeed holds for certain choice of L and T . The detailed analysis for obtaining a lower bound on $|\bar{S}_i^{(\cdot)}|$ is given in Sec. B.4. Specifically, the following result can be obtained by applying union bound to Lemma B.21 over $i = 0, 1, \dots, L-1$.

Lemma B.16 *Let $\bar{S}_i^{(\cdot)}$ be as in Def. B.8 and consider the setup of Alg. 1. Assume $c_{\mathbf{x}} \geq \underline{c}_{\mathbf{x}}(\underline{\mathbf{z}}, \underline{\mathbf{z}})$, $c_{\mathbf{z}} \geq \underline{c}_{\mathbf{z}}$, $L \geq \underline{C}_{\text{sub},N}(0\sqrt{n}, \underline{\mathbf{z}}, T, \underline{\mathbf{z}}, \underline{\mathbf{z}})\log(T)$, and $T \geq \underline{T}_N(\frac{L}{\log(T)}, \underline{\mathbf{z}}, \underline{\mathbf{z}})$, where $\underline{c}_{\mathbf{x}}(\cdot, \cdot)$, $\underline{c}_{\mathbf{z}}$, and $\underline{T}_N(\cdot, \cdot, \cdot)$ are defined in Table 2, and $\underline{C}_{\text{sub},N}(x_0, \cdot, T, \cdot, \cdot)$ is defined in Table 4. Then with probability at least $1 - \delta$, for all $i \in [s]$ and all $i = 0, 1, \dots, L-1$ we have*

$$|\bar{S}_i^{(\cdot)}| \geq \frac{\min T}{2L}. \quad (\text{B.92})$$

B.2.4 Finalizing the SYSID: Proof of Theorem 4.1

To finalize, we combine Corollary B.15 and Lemma B.16 to get our main result on learning the unknown MJS dynamics. The following theorem is a more refined and precise version of our main system identification result in Theorem 4.1.

Theorem B.17 (Main result) Consider the MJS (B.7), with initial state $\mathbf{x}_0 \sim \mathcal{D}_x$ such that $E[\mathbf{x}_0] = 0$, $E[\|\mathbf{x}_0\|^2] \leq \frac{2}{\delta} n$ for some $\delta > 0$. Suppose Assumption 1 on the system and the Markov chain holds. Suppose $\{\mathbf{z}_t\}_{t=0}^{i.i.d.} \sim \mathcal{N}(0, \frac{2}{z} I_p)$, $\{\mathbf{w}_t\}_{t=0}^{i.i.d.} \sim \mathcal{N}(0, \frac{2}{w} I_n)$ and the threshold for bounding the noise satisfies,

$$c_w \geq \sqrt{2 \log(nT)} + \sqrt{2 \log(2/\delta)}. \quad (\text{B.93})$$

Suppose $\|\mathbf{B}_{1,s}\| \leq C_B$ and $\|\mathbf{K}_{1,s}\| \leq C_K$ for some $C_B, C_K > 0$. Let $t_0, \dots, L_{tr1}(\bar{c}, \bar{c}), L_{cov}(\bar{c}), L_{tr2}(\bar{c}, \bar{c})$ and $L_{tr3}(\bar{c}, \bar{c})$ be as in (B.9), (B.10), (B.34), (B.35), (B.46), (B.76) and (B.77), respectively. Suppose $c_x \geq \underline{c}_x(\bar{c}, \bar{c}), c_z \geq \underline{c}_z$, where $\underline{c}_x(\bar{c}, \bar{c})$ and \underline{c}_z are defined in Table 2. Let $C, C_0 > 0$, and $c \geq 6$ be fixed constants. Suppose the trajectory length T satisfies

$$T \gtrsim \max \left\{ \frac{\bar{c} \bar{s} L}{\min(1-\bar{c})} \left((\sqrt{2 \log(nT)} + \sqrt{2 \log(2/\delta)})^2 + C_z^2 C_B^2 \right) \log \left(\frac{36sL(n+p)}{\bar{c}} \right) (n+p), \right. \\ \left. \underline{I}_N \left(\frac{L}{\log(T)}, \bar{c}, \bar{c}, \bar{c} \right) \right\} \quad (\text{B.94})$$

$$\text{where } L \geq \max \left\{ t_0, L_{cov}(\bar{c}), L_{tr1}(\bar{c}, \frac{1}{18L}), L_{tr2}(\bar{c}, \frac{1}{18L}), L_{tr3}(\bar{c}, \frac{1}{18L}), \right. \\ \left. \underline{C}_{sub,N} \left(\frac{1}{\delta} \sqrt{n}, \bar{c}, T, \bar{c}, \bar{c} \right) \log(T) \right\}, \quad (\text{B.95})$$

where $\underline{I}_N(C, \bar{c}, \bar{c}, \bar{c})$ and $\underline{C}_{sub,N}(\bar{x}_0, \bar{c}, T, \bar{c}, \bar{c})$ are defined in Table 2 and 4 respectively. Then, solving the least-squares problem (B.27), with probability at least $1 - \delta$, for all $i \in [s]$, we have

$$\|\hat{\mathbf{A}}_i - \mathbf{A}_i\| \leq \frac{(z + C_K w) \bar{c} \sqrt{s} (\sqrt{2 \log(nT)} + \sqrt{2 \log(2/\delta)} + C_z C_B)}{z \min(1-\bar{c})} \sqrt{\frac{L(n+p)}{T}} \quad (\text{B.96})$$

$$\times \left(C \sqrt{n+p} + (C_0/w) \sqrt{\log \left(\frac{36sL}{\bar{c}} \right)} \right),$$

$$\|\hat{\mathbf{B}}_i - \mathbf{B}_i\| \leq \frac{w \bar{c} \sqrt{s} (\sqrt{2 \log(nT)} + \sqrt{2 \log(2/\delta)} + C_z C_B)}{z \min(1-\bar{c})} \sqrt{\frac{L(n+p)}{T}} \quad (\text{B.97})$$

$$\times \left(C \sqrt{n+p} + (C_0/w) \sqrt{\log \left(\frac{36sL}{\bar{c}} \right)} \right).$$

Remark B.3 Note that in Theorem B.17, with the shorthand notations defined in Tables 2 and 4, the premise conditions (B.93), (B.94), and (B.95) can also be interpreted as the following.

$$c_w = \underline{c}_w(T, \bar{c}) \quad (\text{B.98})$$

$$T \geq \underline{I}_{id,N}(L, \bar{c}, T, \bar{c}, \bar{c}) \quad (\text{B.99})$$

$$L \geq \underline{L}_{id} \left(\frac{1}{\delta} \sqrt{n}, \bar{c}, T, \bar{c}, \bar{c}, \underline{c}_w(T, \bar{c}), \mathbf{K}_{1,s}, L \right). \quad (\text{B.100})$$

From the definition of \underline{L}_{id} , one can see there exists $\underline{L} = \tilde{O}(\log(T))$ such that (B.100) holds by choosing $L = \underline{L}$. Define shorthand notation $\underline{I}_{id,N,L}(\bar{c}, T, \bar{c}, \bar{c}) := \underline{I}_{id,N}(\underline{L}, \bar{c}, T, \bar{c}, \bar{c})$, then the premise conditions (B.93), (B.94), and (B.95) can be implied by the single condition $T \geq \underline{I}_{id,N,L}(\bar{c}, T, \bar{c}, \bar{c})$, under which the main results in Theorem B.17 still hold.

Discussion

• **Sample complexity:** Here, a few remarks are in place. First, the result appears to be convoluted however most of the dependencies are logarithmic (specifically dependency on the failure probability δ and $\log(T)$ terms). Besides these, the dominant term (when estimating \mathbf{A}) reduces to

$$\frac{(z + C_K w) \bar{c} \sqrt{s}}{z \min(1-\bar{c})} \frac{(n+p)}{\sqrt{T}}.$$

which is identical to our statement in Theorem 4.1. Note that the overall sample complexity grows as $T \gtrsim s(n+p)^2 / \frac{2}{\min}$. We remark that, this quadratic growth is somewhat undesirable. A degrees-of-freedom

counting argument would lead to an ideal dependency of $T \gtrsim s(n+p)/\min$. The reason is that, each vector state equation we fit has n scalar equations. The total degrees of freedom for each dynamics pair $(\mathbf{A}_i, \mathbf{B}_i)$ is $n \times (n+p)$. Additionally, for the least-frequent mode, in steady-state, we should observe $\min T$ equations. Putting these together, we would minimally need $n \times \min T \geq n \times (n+p)$, which means we need $T \geq s(n+p)/\min$ samples to estimate s dynamic pairs $(\mathbf{A}_{1s}, \mathbf{B}_{1s})$. Our analysis indicates that this sub-optimality (at least the quadratic growth in n) can be addressed to achieve optimal dependence by establishing a stronger control on the state covariance (e.g. refining (B.18)) as well as a better control on the degree of independence across sampled states (this issue arises during the proof of Theorem B.11).

• **To what extent subsampling is necessary?** We recall that our argument is based on mixing-time arguments which are well-studied in the literature. In Alg. 1, we sub-sample the trajectory for bounded samples and use them to estimate the unknown MJS dynamics. Unfortunately, such a sub-sampling seems unavoidable as long as we don't have a good tail control on the distribution of the state vectors. Specifically, as long as the feature vectors (in our case state vectors) are allowed to be heavy-tailed, existing – to the best of our knowledge – minimum singular value concentration guarantees for the empirical covariance apply under the assumption of boundedness [65]. More recently, self-normalized martingale arguments are employed to address temporal dependencies [55, 59]. We remark that, using martingale-based arguments, it is possible to mitigate the spectral radius dependency by shaving a factor of $1/(1 - \bar{\rho})$ (e.g. martingale based arguments have milder $\bar{\rho}$ dependence [55, 59]).

B.4 Lower Bounding $|S_i^{(k)}|$

To begin, we define sub-sampling period $L = C_{sub} \log(T)$, sub-sampling indices $k = \lfloor i + kL \rfloor$ for $k = 1, 2, \dots, \lfloor T/L \rfloor$, and the time index set

$$S_i^{(k)} = \{k \mid (k) = i, \|\mathbf{x}_k\| \leq c_x \sqrt{\|\mathbf{w}\| \log(T)}, \|\mathbf{z}_k\| \leq c_z \sqrt{\|\mathbf{z}\|}\}$$

by bounding $\|\mathbf{x}_t\|$ and $\|\mathbf{z}_t\|$, which is used to estimate \mathbf{A}_{1s} and \mathbf{B}_{1s} through least squares (Here we generalize isotropic noise $\mathbf{w}_t \sim \mathcal{N}(0, \frac{2}{\mathbf{w}} I_n)$ and $\mathbf{z}_t \sim \mathcal{N}(0, \frac{2}{\mathbf{z}} I_p)$ to $\mathcal{N}(0, \mathbf{w})$ and $\mathcal{N}(0, \mathbf{z})$, respectively.). A fundamental question is: Is $|S_i^{(k)}|$ big enough such that there will be enough data available when applying least squares? We provide answer to this question in this section. Lemma B.18 acts as a building block for the later result; Lemma B.19 provides the lower bound on $|S_i^{(k)}|$; Corollary B.20 gives a more interpretable lower bound on $|S_i^{(k)}|$ when c_x and c_z are large enough; and finally, Lemma B.21 shows how many samples in $S_i^{(k)}$ are “weakly” independent, which is the quantity that essentially determines the sample complexity of estimating \mathbf{A}_{1s} and \mathbf{B}_{1s} .

For clarity, we reiterate some definitions and define a few new ones here. We are given an MJS $(\mathbf{A}_{1s}, \mathbf{B}_{1s}, \mathbf{T})$ with process noise $\mathbf{w}_t \sim \mathcal{N}(0, \mathbf{w})$ and ergodic Markov matrix \mathbf{T} . With some stabilizing controller \mathbf{K}_{1s} , the input is given by $\mathbf{u}_t = \mathbf{K}_{1s} \mathbf{x}_t + \mathbf{z}_t$ where $\mathbf{z}_t \sim \mathcal{N}(0, \mathbf{z})$. Let $\bar{\rho} := \|\mathbf{B}_{1s}\|^2 \|\mathbf{z}\| + \|\mathbf{w}\|$. Let $\mathbf{L}_i := \mathbf{A}_i + \mathbf{B}_i \mathbf{K}_{1s}$. Let $\tilde{\mathbf{L}} \in \mathbb{R}^{sn^2 \times sn^2}$ denote the augmented closed-loop state matrix with ij -th $n^2 \times n^2$ block given by $[\tilde{\mathbf{L}}]_{ij} := [\mathbf{T}]_{ji} \mathbf{L}_j \otimes \mathbf{L}_i$. Let $\bar{\rho} \in [0, 1)$ and $\bar{\rho} > 0$ be two constants such that $\|\tilde{\mathbf{L}}^k\| \leq \bar{\rho}^k$. By definition, one available choice for $\bar{\rho}$ and $\bar{\rho}$ are $\bar{\rho}$ and $(\tilde{\mathbf{L}})$, respectively. Let $t_{MC}(\cdot)$ and t_{MC} denote the mixing time of \mathbf{T} as in Definition A.2. Let π denote the stationary distribution of \mathbf{T} , $\pi_{\min} = \min_i \pi(i)$, and $\pi_{\max} = \max_i \pi(i)$. Assume the initial state \mathbf{x}_0 satisfies $\mathbb{E}[\|\mathbf{x}_0\|^2] \leq \bar{\rho}_0^2$ for some $\bar{\rho}_0 \geq 0$. Lastly, without loss of generality, we consider the sub-trajectory with zero shift, that is, $\tau = 0$, which is identical to any $0 \leq \tau \leq L - 1$.

Lemma B.18 *Suppose the Markov chain trajectory $\{(0), (1), \dots\}$ and a sequence of events $\{A_0, A_1, \dots\}$ are both adapted to filtration $\{\mathcal{F}_0, \mathcal{F}_1, \dots\}$, i.e. (t) and $\mathbf{1}_{\{A_t\}}$ are both \mathcal{F}_t -measurable. We assume $\mathbb{E}[\mathbf{1}_{\{(t)=j\}} \mid \mathcal{F}_{t-r}] = \mathbb{P}((t)=j \mid (t-r))$ for all $j \in [s], t$, and $r < t$. For all $i \in [s]$, let*

$$N_i = \sum_{k=1}^{\lfloor T/L \rfloor} \mathbf{1}_{\{(kL)=i\}} \mathbf{1}_{\{A_{kL}\}} \quad (\text{B.101})$$

and suppose

$$\mathbb{E}[\mathbf{1}_{\{A_t\}} \mid \mathcal{F}_{t-L}] \geq 1 - p_t \quad (\text{B.102})$$

for some $p_t \in [0, 1)$ and $L = C_{\text{sub}} \log(T)$. Assume $C_{\text{sub}} \geq C_{MC}$, and for some $\epsilon > 0$, $T \geq \underline{I}_{MC,1}(C_{\text{sub}}, \epsilon)$, where C_{MC} and $\underline{I}_{MC,1}(C, \epsilon)$ are defined in Tables 2 and 4, respectively. Then we have

$$\mathbb{P}\left(\bigcap_{i=1}^S \left\{ N_i \geq \frac{T}{C_{\text{sub}} \log(T)} \left(1 - \frac{1}{(i)} \sqrt{\log\left(\frac{S}{i}\right) \frac{17 C_{\text{sub}} \max \log(T)}{T}}\right) - \sum_{k=1}^{T/L} \rho_{kL} \right\}\right) \geq 1 - \epsilon. \quad (\text{B.103})$$

Proof For some $\epsilon < \min/2$, we temporarily let $L \geq 6 t_{MC} \log(\epsilon^{-1})$. From the proof for [17, Lemma 13 (47)], we know this guarantees $L \geq t_{MC}(\epsilon/2)$. By definition of $t_{MC}(\cdot)$, we know $\max_i \|([\mathbf{T}^L]_i) - \mathbf{1}\|_1 \leq \epsilon/2$, and since $([\mathbf{T}^L]_i) \mathbf{1} = \mathbf{1}$, we further have

$$\max_i \|([\mathbf{T}^L]_i) - \mathbf{1}\| \leq \frac{\epsilon}{2} \leq \frac{\min}{4}. \quad (\text{B.104})$$

For simplicity, we assume $\lfloor T/L \rfloor = T/L =: \tilde{T}$. To ease the notation, we let $\tilde{\sim}(k) := (kL)$, $\tilde{A}_k := A_{kL}$, and $\tilde{\mathcal{F}}_k := \mathcal{F}_{kL}$. Then, one can see $\tilde{\sim}(k)$ and \tilde{A}_k are both $\tilde{\mathcal{F}}_k$ -measurable. Define $\tilde{\rho}_k, \tilde{\rho}_k \in \mathbb{R}^S$ such that

$$j(i) := \mathbf{1}_{\{\tilde{\sim}(j)=i\}} \mathbf{1}_{\{\tilde{A}_j\}} - \mathbb{E}[\mathbf{1}_{\{\tilde{\sim}(j)=i\}} \mathbf{1}_{\{\tilde{A}_j\}} \mid \tilde{\mathcal{F}}_{j-1}], \quad (\text{B.105})$$

$$\tilde{\rho}_k(i) := \sum_{j=1}^k j(i). \quad (\text{B.106})$$

Note that for all $i \in [S]$, $\{\tilde{\rho}_k(i), \tilde{\mathcal{F}}_k\}$ forms a martingale as

$$\begin{aligned} \mathbb{E}[\tilde{\rho}_{k+1}(i) \mid \tilde{\mathcal{F}}_k] &= \mathbb{E}\left[\sum_{j=1}^{k+1} j(i) \mid \tilde{\mathcal{F}}_k\right] \\ &= \sum_{j=1}^k j(i) + \mathbb{E}[\mathbf{1}_{\{\tilde{\sim}(k+1)=i\}} \mathbf{1}_{\{\tilde{A}_{k+1}\}} - \mathbb{E}[\mathbf{1}_{\{\tilde{\sim}(k+1)=i\}} \mathbf{1}_{\{\tilde{A}_{k+1}\}} \mid \tilde{\mathcal{F}}_k] \mid \tilde{\mathcal{F}}_k] \\ &= \sum_{j=1}^k j(i) = \tilde{\rho}_k(i), \end{aligned} \quad (\text{B.107})$$

thus $\tilde{\rho}_k(i) - \tilde{\rho}_{k-1}(i)$ can be viewed as the martingale difference sequence. Since $\mathbb{E}[\tilde{\rho}_k(i) \mid \tilde{\mathcal{F}}_{k-1}] = 0$, we have $\mathbb{E}[\tilde{\rho}_k(i)^2 \mid \tilde{\mathcal{F}}_{k-1}] = \text{Var}(\tilde{\rho}_k(i) \mid \tilde{\mathcal{F}}_{k-1}) = \text{Var}(\mathbf{1}_{\{\tilde{\sim}(j)=i\}} \mathbf{1}_{\{\tilde{A}_j\}} \mid \tilde{\mathcal{F}}_{k-1}) \leq \mathbb{E}[\mathbf{1}_{\{\tilde{\sim}(k)=i\}}^2 \mathbf{1}_{\{\tilde{A}_k\}}^2 \mid \tilde{\mathcal{F}}_{k-1}] \leq \mathbb{E}[\mathbf{1}_{\{\tilde{\sim}(k)=i\}} \mid \tilde{\mathcal{F}}_{k-1}] = \mathbb{P}(\tilde{\sim}(k) = i \mid \tilde{\sim}(k-1)) = [\mathbf{T}^L]_{((k-1)L), i}$. By the choice of L , using (B.104), we know $[\mathbf{T}^L]_{((k-1)L), i} \leq (i) + \max_j \|([\mathbf{T}^L]_j) - \mathbf{1}\| \leq 2 \max$. Thus,

$$\sum_{k=1}^{\tilde{T}} \mathbb{E}[\tilde{\rho}_k(i)^2 \mid \tilde{\mathcal{F}}_{k-1}] \leq 2 \max \tilde{T}. \quad (\text{B.108})$$

With this, and the fact that $|\tilde{\rho}_k(i)| < 1$, we have

$$\begin{aligned} \mathbb{P}\left(N_i - \sum_{k=1}^{\tilde{T}} \mathbb{E}[\mathbf{1}_{\{\tilde{\sim}(k)=i\}} \mathbf{1}_{\{\tilde{A}_k\}} \mid \tilde{\mathcal{F}}_{k-1}] \geq \tilde{T}/2\right) &\stackrel{(i)}{=} \mathbb{P}\left(\tilde{\tau}(i) \geq \tilde{T}/2\right) \\ &\stackrel{(ii)}{\leq} \exp\left(-\frac{\tilde{T}^2/8}{2 \max + \tilde{T}/6}\right) \\ &\stackrel{(iii)}{\leq} \exp\left(-\frac{\tilde{T}^2}{17 \max L}\right), \end{aligned} \quad (\text{B.109})$$

where (i) follows from the definition of N_i and $\tilde{\tau}(i)$; (ii) follows from Freedman's inequality [26, Theorem 1.6], and (iii) follows since $\leq \min/2$. Note that

$$\begin{aligned}
& \left| \sum_{k=1}^{\tilde{T}} \mathbb{E}[\mathbf{1}_{\{\tilde{\tau}(k)=i\}} \mathbf{1}_{\{\tilde{A}_k\}} \mid \tilde{\mathcal{F}}_{k-1}] - \tilde{T}(i) \right| \\
& \leq \left| \sum_{k=1}^{\tilde{T}} \mathbb{E}[\mathbf{1}_{\{\tilde{\tau}(k)=i\}} \mid \tilde{\mathcal{F}}_{k-1}] - \tilde{T}(i) \right| + \left| \sum_{k=1}^{\tilde{T}} \mathbb{E}[\mathbf{1}_{\{\tilde{\tau}(k)=i\}} \mid \tilde{\mathcal{F}}_{k-1}] - \mathbb{E}[\mathbf{1}_{\{\tilde{\tau}(k)=i\}} \mathbf{1}_{\{\tilde{A}_k\}} \mid \tilde{\mathcal{F}}_{k-1}] \right| \\
& \leq \tilde{T} \max_j |[\mathbf{T}^L]_{j,i} - \tilde{T}(i)| + \left| \sum_{k=1}^{\tilde{T}} \mathbb{E}[\mathbf{1}_{\{\tilde{A}_k\}} \mid \tilde{\mathcal{F}}_{k-1}] \right| \\
& \leq \tilde{T} \frac{1}{2} + \sum_{k=1}^{\tilde{T}} \rho_{kL}.
\end{aligned} \tag{B.110}$$

Then, combining this with (B.109) and applying union bound, we have with probability at least $1 - \text{Sexp}(-\frac{T^2}{17 \max L})$,

$$\bigcap_{i=1}^s \left\{ N_i \geq \frac{T}{L} \tilde{T}(i) - \frac{T}{L} - \sum_{k=1}^{\tilde{T}} \rho_{kL} \right\} \tag{B.111}$$

when $\leq \min/2$ and $L \geq 6t_{MC} \log(\frac{1}{\epsilon})$. Then, similar to the proof of Lemma B.1, we know if we pick $L = C_{sub} \log(T)$ with $C_{sub} \geq t_{MC} \cdot \max\{3, 3 - 3 \log(\frac{1}{\epsilon} \log(s))\}$, and for some $\epsilon > 0$, we pick the trajectory length $T \geq (68 C_{sub} \max \frac{1}{\min} \log(\frac{2s}{\epsilon}))^2$, with probability at least $1 - \epsilon$, we have

$$\bigcap_{i=1}^s \left\{ N_i \geq \frac{T}{C_{sub} \log(T)} \tilde{T}(i) \left(1 - \frac{1}{\tilde{T}(i)} \sqrt{\log(\frac{2s}{\epsilon}) \frac{17 C_{sub} \max \log(T)}{T}} \right) - \sum_{k=1}^{\tilde{T}} \rho_{kL} \right\}. \tag{B.112}$$

Lemma B.19 For some $\epsilon > 0$, we assume $c_z \geq (\sqrt{3} + \sqrt{6})\sqrt{\rho}$, $C_{sub} \geq \max\{C_{MC}, \underline{C}_{sub, \mathbf{x}}(\bar{\mathbf{x}}_0, \bar{\mathbf{z}}, T, \bar{\mathbf{c}}, \bar{\mathbf{c}})\}$ and $T \geq \max\{\underline{I}_{MC,1}(C_{sub}, \frac{1}{2}), \underline{I}_{cl,1}(\bar{\mathbf{c}}, \bar{\mathbf{c}})\}$, where C_{MC} , $\underline{I}_{MC,1}(C, \cdot)$ and $\underline{I}_{cl,1}(\cdot, \cdot)$ are defined in Table 2, and $\underline{C}_{sub, \mathbf{x}}(\bar{\mathbf{x}}_0, \bar{\mathbf{z}}, T, \bar{\mathbf{c}}, \bar{\mathbf{c}})$ is defined in Table 4. Then, with probability at least $1 - \epsilon$, the following intersected events occur

$$\bigcap_{i=1}^s \left\{ |S_i(\cdot)| \geq \frac{T}{C_{sub} \log(T)} \tilde{T}(i) \left(1 - \frac{1}{\tilde{T}(i)} \sqrt{\log(\frac{2s}{\epsilon}) \frac{17 C_{sub} \max \log(T)}{T}} \right) - \frac{2n\sqrt{s} \bar{\mathbf{c}}^{-2}}{(i)c_{\mathbf{x}}^2 \|\mathbf{w}\| \log(T)(1 - \bar{\mathbf{c}})} - \frac{1}{(i)} e^{-\frac{c_z^2}{3}} \right\}. \tag{B.113}$$

Proof For simplicity, we assume $\lfloor T/L \rfloor = T/L$. We let \mathcal{F}_t denote the sigma algebra generated by $\{\{(r)\}_{r=0}^t, \mathbf{w}_0, \mathbf{z}_0, \mathbf{x}_0\}$, and let $A_t = \{\|\mathbf{x}_t\| \leq c_x \sqrt{\|\mathbf{w}\| \log(T)}, \|\mathbf{z}_t\| \leq c_z \sqrt{\|\mathbf{z}\|}\}$, then by Lemma B.18, when $C_{sub} \geq C_{MC} = t_{MC} \cdot \max\{3, 3 - 3 \log(\frac{1}{\epsilon} \log(s))\}$ and $T \geq \underline{I}_{MC,1}(C_{sub}, \frac{1}{2}) = (68 C_{sub} \max \frac{1}{\min} \log(\frac{2s}{\epsilon}))^2$, with probability at least $1 - \frac{\epsilon}{2}$, we have

$$\bigcap_{i=1}^s \left\{ |S_i(\cdot)| \geq \frac{T}{C_{sub} \log(T)} \tilde{T}(i) \cdot \left(1 - \frac{1}{\tilde{T}(i)} \sqrt{\log(\frac{2s}{\epsilon}) \frac{17 C_{sub} \max \log(T)}{T}} \right) - \underbrace{\sum_{k=1}^{\tilde{T}} \mathbb{P}(A_k^c \mid \mathbf{x}_{(k-1)L+1}, \dots, ((k-1)L))}_{=P} \right\}. \tag{B.114}$$

For term P , we have

$$\begin{aligned}
P &= \sum_{k=1}^{T-L} \mathbb{P} \left(\|\mathbf{x}_{kL}\| \geq c_{\mathbf{x}} \sqrt{\|\mathbf{w}\| \log(T)} \cup \|\mathbf{z}_{kL}\| \geq c_{\mathbf{z}} \sqrt{\|\mathbf{z}\|} \mid \mathbf{x}_{(k-1)L+1}, ((k-1)L) \right) \\
&\leq \underbrace{\sum_{k=1}^{T-L} \mathbb{P} \left(\|\mathbf{z}_{kL}\| \geq c_{\mathbf{z}} \sqrt{\|\mathbf{z}\|} \mid \mathbf{x}_{(k-1)L+1}, ((k-1)L) \right)}_{= P_1} \\
&\quad + \underbrace{\sum_{k=1}^{T-L} \mathbb{P} \left(\|\mathbf{x}_{kL}\| \geq c_{\mathbf{x}} \sqrt{\|\mathbf{w}\| \log(T)} \mid \mathbf{x}_{(k-1)L+1}, ((k-1)L) \right)}_{= P_2}.
\end{aligned} \tag{B.115}$$

For term P_1 , we know from Lemma A.5 that when $c_{\mathbf{z}} \geq (\sqrt{3} + \sqrt{6})\sqrt{\rho}$, we have $P_1 = \sum_{k=1}^{T-L} \mathbb{P} \left(\|\mathbf{z}_{kL}\| \geq c_{\mathbf{z}} \sqrt{\|\mathbf{z}\|} \right) \leq \frac{T}{L} e^{-\frac{c_{\mathbf{z}}^2}{3}}$. Now we consider term P_2 . From Lemma A.4, we know

$$\mathbb{E}[\|\mathbf{x}_{kL}\|^2 \mid \mathbf{x}_{(k-1)L+1}, ((k-1)L)] \leq \sqrt{ns} \bar{\zeta} \bar{\zeta}^{L-1} \|\mathbf{x}_{(k-1)L+1}\|^2 + \frac{n\sqrt{s} \bar{\zeta}^{-2}}{1 - \bar{\zeta}}, \tag{B.116}$$

thus by Markov inequality, we have

$$\begin{aligned}
P_2 &\leq \sum_{k=1}^{T-L} \frac{1}{c_{\mathbf{x}}^2 \|\mathbf{w}\| \log(T)} \left(\sqrt{ns} \bar{\zeta} \bar{\zeta}^{L-1} \|\mathbf{x}_{(k-1)L+1}\|^2 + \frac{n\sqrt{s} \bar{\zeta}^{-2}}{1 - \bar{\zeta}} \right) \\
&\leq \frac{1}{c_{\mathbf{x}}^2 \|\mathbf{w}\| \log(T)} \left(\frac{T}{L} \frac{n\sqrt{s} \bar{\zeta}^{-2}}{1 - \bar{\zeta}} + \sqrt{ns} \bar{\zeta} \bar{\zeta}^{L-1} \sum_{k=1}^{T-L} \|\mathbf{x}_{(k-1)L+1}\|^2 \right).
\end{aligned} \tag{B.117}$$

Now, we seek to upper bound $\bar{\zeta}^{L-1} \sum_{k=1}^{T-L} \|\mathbf{x}_{(k-1)L+1}\|^2$ with high probability. Note that the assumption $C_{sub} \geq \underline{C}_{sub, \mathbf{x}}(\bar{\chi}_0, T, \bar{\zeta}, \bar{\zeta})$ implies the following

$$L = C_{sub} \log(T) \geq \frac{1}{\log(\bar{\zeta}^{-1})} \max \left\{ \log(2), 2 \log\left(\frac{8\sqrt{ns} \bar{\zeta} \bar{\chi}_0^2}{1 - \bar{\zeta}}\right), 2 \log\left(4T \frac{n\sqrt{s} \bar{\zeta}^{-2}}{(1 - \bar{\zeta})} + 2\right) \right\}. \tag{B.118}$$

Then, we have

$$\begin{aligned}
\mathbb{P} \left(\bar{\zeta}^{L-1} \sum_{k=1}^{T-L} \|\mathbf{x}_{(k-1)L+1}\|^2 \leq \frac{T}{L} \frac{(1 - \bar{\zeta})}{4Tn\sqrt{s} \bar{\zeta}^{-2}} \right) &\stackrel{(i)}{\geq} \mathbb{P} \left(\bar{\zeta}^{L-1} \sum_{k=1}^{T-L} \|\mathbf{x}_{(k-1)L+1}\|^2 \leq \frac{T}{L} \bar{\zeta}^{\frac{L}{2}-1} \right) \\
&\geq \mathbb{P} \left(\bigcap_{k=1}^{T-L} \left\{ \|\mathbf{x}_{(k-1)L+1}\|^2 \leq \frac{\bar{\zeta}^{-\frac{L}{2}}}{L} \right\} \right) \\
&\geq 1 - \sum_{k=1}^{T-L} \mathbb{P} \left(\|\mathbf{x}_{(k-1)L+1}\|^2 \geq \frac{\bar{\zeta}^{-\frac{L}{2}}}{L} \right) \\
&\stackrel{(ii)}{\geq} 1 - \sum_{k=1}^{T-L} \frac{\bar{\zeta}^{\frac{L}{2}}}{L} \left(\sqrt{ns} \bar{\zeta} \bar{\zeta}^{(k-1)L+1} \bar{\chi}_0^2 + \frac{n\sqrt{s} \bar{\zeta}^{-2}}{1 - \bar{\zeta}} \right) \\
&\geq 1 - \frac{\bar{\zeta}^{\frac{L}{2}+1} \sqrt{ns} \bar{\zeta} \bar{\chi}_0^2}{1 - \bar{\zeta}} - \frac{T}{L} \frac{\bar{\zeta}^{\frac{L}{2}} n\sqrt{s} \bar{\zeta}^{-2}}{1 - \bar{\zeta}} \\
&\stackrel{(iii)}{\geq} 1 - 2 \frac{\bar{\zeta}^{\frac{L}{2}} \sqrt{ns} \bar{\zeta} \bar{\chi}_0^2}{4L} \\
&\stackrel{(iv)}{\geq} 1 - \frac{\bar{\zeta}^{\frac{L}{2}}}{4} - \frac{\bar{\zeta}^{\frac{L}{2}}}{4} = 1 - \frac{\bar{\zeta}^{\frac{L}{2}}}{2},
\end{aligned} \tag{B.119}$$

where (i) follows from (B.118) which gives $\frac{\underline{\zeta}}{\underline{\zeta}}^{-1} \leq \frac{(1-\underline{\zeta})}{4Tn\frac{\underline{\zeta}}{\underline{\zeta}}^{-2}}$; (ii) follows from Lemma A.4 and Markov inequality; (iii) follows from (B.118) which gives $\frac{\underline{\zeta}}{\underline{\zeta}} \leq \frac{1}{2}$ and $\frac{\underline{\zeta}}{\underline{\zeta}} \leq \frac{(1-\underline{\zeta})}{4Tn\frac{\underline{\zeta}}{\underline{\zeta}}^{-2}}$ and (iv) follows from (B.118) which gives $\frac{\underline{\zeta}}{\underline{\zeta}} \leq \frac{1}{8\frac{\underline{\zeta}}{\underline{\zeta}}}$. Therefore, we have with probability at least $1 - \frac{1}{2}$

$$P_2 \leq \frac{1}{c_{\mathbf{x}}^2 \|\mathbf{w}\| \log(T)} \left(\frac{T n \sqrt{s} \underline{\zeta}^{-2}}{L(1-\underline{\zeta})} + \frac{1-\underline{\zeta}}{4L\sqrt{n}^{-2}} \right), \quad (\text{B.120})$$

and thus,

$$\begin{aligned} P \leq P_1 + P_2 &\leq \frac{1}{c_{\mathbf{x}}^2 \|\mathbf{w}\| \log(T)} \left(\frac{T n \sqrt{s} \underline{\zeta}^{-2}}{L(1-\underline{\zeta})} + \frac{1-\underline{\zeta}}{4L\sqrt{n}^{-2}} \right) + \frac{T}{L} e^{-\frac{c_{\mathbf{z}}^2}{3}} \\ &\leq \frac{1}{c_{\mathbf{x}}^2 \|\mathbf{w}\| \log(T)} \left(\frac{T 2n \sqrt{s} \underline{\zeta}^{-2}}{L(1-\underline{\zeta})} \right) + \frac{T}{L} e^{-\frac{c_{\mathbf{z}}^2}{3}}, \end{aligned} \quad (\text{B.121})$$

where the second inequality follows from $T \geq \underline{I}_{cl,1}(\underline{\zeta}, \underline{\zeta})$. Plugging this into (B.114), we have with probability at least $1 - \frac{1}{2}$,

$$\bigcap_{i=1}^S \left\{ |S_i^{(\cdot)}| \geq \frac{T}{C_{sub} \log(T)} \cdot \left(1 - \frac{1}{(i)} \sqrt{\log\left(\frac{2S}{(i)}\right)} \frac{17C_{sub} \max \log(T)}{T} - \frac{2n\sqrt{s} \underline{\zeta}^{-2}}{(i)c_{\mathbf{x}}^2 \|\mathbf{w}\| \log(T)(1-\underline{\zeta})} - \frac{1}{(i)} e^{-\frac{c_{\mathbf{z}}^2}{3}} \right) \right\}, \quad (\text{B.122})$$

which concludes the proof. \blacksquare

Note that when T , $c_{\mathbf{x}}$, and $c_{\mathbf{z}}$ are sufficiently large enough, we could obtain a more interpretable version of Lemma B.19 which is presented as follows.

Corollary B.20 *For some $\delta > 0$, assume $c_{\mathbf{x}} \geq \frac{1}{3} \underline{c}_{\mathbf{x}}(\underline{\zeta}, \underline{\zeta})$, $c_{\mathbf{z}} \geq \underline{c}_{\mathbf{z}}$, $C_{sub} \geq \max\{C_{MC}, \underline{C}_{sub,\mathbf{x}}(\bar{x}_0, \cdot, T, \underline{\zeta}, \underline{\zeta})\}$, and $T \geq \underline{I}_N(C_{sub}, 2, \cdot, \underline{\zeta}, \underline{\zeta}) := \max\{\underline{I}_{MC}(C_{sub}, \cdot), \underline{I}_{cl,1}(\underline{\zeta}, \underline{\zeta})\}$, where $\underline{c}_{\mathbf{x}}(\cdot, \cdot)$, $\underline{c}_{\mathbf{z}}$, C_{MC} , $\underline{I}_{MC}(C, \cdot)$, $\underline{I}_{cl,1}(\cdot, \cdot)$ are defined in Table 2, and $\underline{C}_{sub,\mathbf{x}}(\bar{x}_0, \cdot, T, \cdot, \cdot)$ is defined in Table 4. Then, with probability at least $1 - \delta/2$, the following intersected events occur*

$$\bigcap_{i=1}^S \left\{ |S_i^{(\cdot)}| \geq \frac{T \min}{2C_{sub} \log(T)} \right\}. \quad (\text{B.123})$$

Now we provide a result on how many data in $S_i^{(\cdot)}$ are ‘‘weakly’’ independent, which is the quantity that essentially determines the sample complexity of estimating \mathbf{A}_{1S} and \mathbf{B}_{1S} in Algorithm 1. We first define a few notations. Let $i_{,1}, \dots, i_{,S_i^{(\cdot)}}$ denote the elements in $S_i^{(\cdot)}$, and let $i_{,0} = 0$. Define $\bar{\mathbf{x}}_{i,k}$ such that

$$\bar{\mathbf{x}}_{i,k} = \sum_{j=1}^{i_{,k-1}} \left(\prod_{k=1}^{j-1} \mathbf{L}_{(t-k)} \right) (\mathbf{B}_{(t-j)} \mathbf{z}_{t-j} + \mathbf{w}_{t-j}) + \mathbf{B}_{(i_{,k-1})} \mathbf{z}_{i_{,k-1}} + \mathbf{w}_{i_{,k-1}}. \quad (\text{B.124})$$

One can view $\bar{\mathbf{x}}_{i,k}$ as follows: set $\mathbf{x}_{i,k-1} = 0$, then propagate the dynamics to time $i_{,k}$ following the same noise and mode switching sequences, $\mathbf{w}_{i_{,k-1} \dots i_{,k-1}}$, $\mathbf{z}_{i_{,k-1} \dots i_{,k-1}}$, $\{\mathbf{L}_{(t)}\}_{t=i_{,k-1}}^{i_{,k}-1}$. Or, one can also view $\bar{\mathbf{x}}_{i,k}$ as the contribution of noise \mathbf{x}_t and \mathbf{z}_t that propagate $\mathbf{x}_{i_{,k-1}}$ to $\mathbf{x}_{i,k}$. And it is easy to see that

$$\mathbf{x}_{i,k} - \bar{\mathbf{x}}_{i,k} = \left(\prod_{k=1}^{i_{,k-1}} \mathbf{L}_{(t-k)} \right) \mathbf{x}_{i_{,k-1}}. \quad (\text{B.125})$$

Define $\bar{S}_i^{(\cdot)} \subseteq S_i^{(\cdot)}$ such that

$$\bar{S}_i^{(\cdot)} := \left\{ k \mid (k) = i, \|\mathbf{x}_k\| \leq c_{\mathbf{x}} \sqrt{\|\mathbf{w}\| \log(T)}, \|\mathbf{z}_k\| \leq c_{\mathbf{z}} \sqrt{\|\mathbf{z}\|}, \|\bar{\mathbf{x}}_k\| \leq \frac{c_{\mathbf{x}} \sqrt{\|\mathbf{w}\| \log(T)}}{2} \right\}.$$

The next lemma provides a lower bound on $|\bar{S}_i^{(\cdot)}|$.

Lemma B.21 Assume $c_{\mathbf{x}} \geq \underline{c}_{\mathbf{x}}(\underline{\bar{L}}, \underline{\bar{L}})$, $c_{\mathbf{z}} \geq \underline{c}_{\mathbf{z}}$, $C_{\text{sub}} \geq \underline{C}_{\text{sub}, N}(\bar{X}_0, T, \underline{\bar{L}}, \underline{\bar{L}}) := \max\{C_{MC}, \underline{C}_{\text{sub}, \mathbf{x}}(\bar{X}_0, \bar{z}, T, \underline{\bar{L}}, \underline{\bar{L}}), \underline{C}_{\text{sub}, \bar{\mathbf{x}}}(\bar{z}, T, \underline{\bar{L}}, \underline{\bar{L}})\}$, and $T \geq \underline{T}_N(\underline{C}_{\text{sub}}, \underline{\bar{L}}, \underline{\bar{L}})$, where $\underline{c}_{\mathbf{x}}(\cdot, \cdot)$, $\underline{c}_{\mathbf{z}}$, $\underline{C}_{\text{sub}, \mathbf{x}}(\bar{X}_0, T, \cdot, \cdot)$, and $\underline{C}_{\text{sub}, \bar{\mathbf{x}}}(\bar{z}, T, \cdot, \cdot)$ are defined in Table 4, and C_{MC} and $\underline{T}_N(\cdot, \cdot, \cdot)$ are defined in Table 2. Then with probability at least $1 - \bar{\epsilon}$, the following intersected events occur

$$\bigcap_{i=1}^s \left\{ |\bar{S}_i^{(\cdot)}| \geq \frac{T \min}{2 C_{\text{sub}} \log(T)} \right\}. \quad (\text{B.126})$$

Proof We define sets $R_i^{(\cdot)} \subseteq S_i^{(\cdot)}$ and $\bar{R}_i^{(\cdot)} \subseteq \bar{S}_i^{(\cdot)}$ such that

$$R_i^{(\cdot)} := \left\{ k \mid (k) = i, \|\mathbf{x}_k\| \leq \frac{c_{\mathbf{x}} \sqrt{\|\mathbf{w}\| \log(T)}}{3}, \|\mathbf{z}_k\| \leq c_{\mathbf{z}} \sqrt{\|\mathbf{z}\|} \right\},$$

$$\bar{R}_i^{(\cdot)} := \left\{ k \mid (k) = i, \|\mathbf{x}_k\| \leq \frac{c_{\mathbf{x}} \sqrt{\|\mathbf{w}\| \log(T)}}{3}, \|\mathbf{z}_k\| \leq c_{\mathbf{z}} \sqrt{\|\mathbf{z}\|}, \|\bar{\mathbf{x}}_k\| \leq \frac{c_{\mathbf{x}} \sqrt{\|\mathbf{w}\| \log(T)}}{2} \right\}.$$

Note that $\bar{R}_i^{(\cdot)} \subseteq R_i^{(\cdot)}$. We will first (i) lower bound $|R_i^{(\cdot)}|$ and (ii) show $|\bar{R}_i^{(\cdot)}| = |R_i^{(\cdot)}|$, then we could lower bound $|\bar{S}_i^{(\cdot)}|$ since $|\bar{S}_i^{(\cdot)}| \geq |\bar{R}_i^{(\cdot)}|$ and conclude the proof.

Using Corollary B.20, we see under given assumptions, with probability at least $1 - \bar{\epsilon}$,

$$\bigcap_{i=1}^s \left\{ |R_i^{(\cdot)}| \geq \frac{T \min}{2 C_{\text{sub}} \log(T)} \right\}. \quad (\text{B.127})$$

Let $i_{1,1}, \dots, i_{i, R_i^{(\cdot)}}$ denote the elements in $R_i^{(\cdot)}$. It is easy to see $\{i_{1,1}, \dots, i_{i, R_i^{(\cdot)}}\} \subseteq \{i_{1,1}, \dots, i_{i, S_i^{(\cdot)}}\}$. Consider an arbitrary $i_{ij} \in R_i^{(\cdot)}$ and $i_{ij'} \in S_i^{(\cdot)}$ denote the counterpart of i_{ij} such that $i_{ij'} = i_{ij}$. By definition of $R_i^{(\cdot)}$, we have

$$\|\mathbf{x}_{i_{j'}}\| \leq \frac{c_{\mathbf{x}} \sqrt{\|\mathbf{w}\| \log(T)}}{3}. \quad (\text{B.128})$$

From (B.125), together with Lemma A.4, we have

$$\begin{aligned} \mathbb{E}[\|\mathbf{x}_{i_{j'}} - \bar{\mathbf{x}}_{i_{j'}}\|^2] &\leq \sqrt{ns} \underline{L} \underline{L}^{i_{j'} - i_{j'-1}} \mathbb{E}[\|\mathbf{x}_{i_{j'-1}}\|^2] \\ &\leq \sqrt{ns} \underline{L} \underline{L}^{\frac{L}{L}} (c_{\mathbf{x}}^2 \|\mathbf{w}\| \log(T)), \end{aligned} \quad (\text{B.129})$$

where the second inequality follows from $i_{j'} - i_{j'-1} \geq L$ and $\|\mathbf{x}_{i_{j'-1}}\| \leq c_{\mathbf{x}} \sqrt{\|\mathbf{w}\| \log(T)}$ by definition of $S_i^{(\cdot)}$. Then, by Markov inequality, we have

$$\mathbb{P}\left(\|\mathbf{x}_{i_{j'}} - \bar{\mathbf{x}}_{i_{j'}}\| \leq \frac{c_{\mathbf{x}} \sqrt{\|\mathbf{w}\| \log(T)}}{6}\right) \geq 1 - 36\sqrt{ns} \underline{L} \underline{L}^{\frac{L}{L}}. \quad (\text{B.130})$$

Then, using union bound, we have

$$\begin{aligned}
& \mathbb{P} \left(\bigcap_{i \in [s]} \bigcap_{j'} \left\{ \|\mathbf{x}_{i,j'} - \bar{\mathbf{x}}_{i,j'}\| \leq \frac{c_{\mathbf{x}} \sqrt{\|\mathbf{w}\| \log(T)}}{6} \right\} \right) \\
& \geq 1 - 36\sqrt{\bar{n}} s^{1.5} |\mathcal{R}_i^{(\cdot)}| \bar{\zeta} \frac{L}{\bar{\zeta}} \\
& \geq 1 - 36\sqrt{\bar{n}} s^{1.5} T \bar{\zeta} \frac{L}{\bar{\zeta}} \\
& \geq 1 - \bar{\zeta},
\end{aligned} \tag{B.131}$$

where the last line follows from $L = C_{sub} \log(T)$ and $C_{sub} \geq \underline{C}_{sub, \bar{\mathbf{x}}}(T, \bar{\zeta}, \bar{\zeta})$ in the assumption. Note that $\|\bar{\mathbf{x}}_{i,j}\| = \|\bar{\mathbf{x}}_{i,j'}\| \leq \|\mathbf{x}_{i,j'}\| + \|\mathbf{x}_{i,j'} - \bar{\mathbf{x}}_{i,j'}\|$. This together with (B.128) and (B.131) gives, with probability at least $1 - \bar{\zeta}$,

$$\bigcap_{i \in [s]} \bigcap_{j \in [\mathcal{R}_i^{(\cdot)}]} \left\{ \|\bar{\mathbf{x}}_{i,j}\| \leq \frac{c_{\mathbf{x}} \sqrt{\|\mathbf{w}\| \log(T)}}{2} \right\}. \tag{B.132}$$

This implies for any i , for any $k \in \mathcal{R}_i^{(\cdot)}$, we have $k \in \bar{\mathcal{R}}_i^{(\cdot)}$, i.e. $\mathcal{R}_i^{(\cdot)} \subseteq \bar{\mathcal{R}}_i^{(\cdot)}$. Thus, we have $\mathcal{R}_i^{(\cdot)} = \bar{\mathcal{R}}_i^{(\cdot)}$ and $|\mathcal{R}_i^{(\cdot)}| = |\bar{\mathcal{R}}_i^{(\cdot)}|$. Combining this with (B.127), we have with probability at least $1 - \bar{\zeta}$,

$$\bigcap_{i=1}^s \left\{ |\bar{\mathcal{R}}_i^{(\cdot)}| \geq \frac{T \min}{2C_{sub} \log(T)} \right\}. \tag{B.133}$$

Finally, we could conclude the proof by noticing $|\bar{\mathcal{S}}_i^{(\cdot)}| \geq |\bar{\mathcal{R}}_i^{(\cdot)}|$. ■

C MJS Regret Analysis

Consider MJS-LQR($\mathbf{A}_{1s}, \mathbf{B}_{1s}, \mathbf{T}, \mathbf{Q}_{1s}, \mathbf{R}_{1s}$) with dynamics noise $\mathbf{w}_t \sim \mathcal{N}(0, \mathbf{w})$, some arbitrary initial state \mathbf{x}_0 and stabilizing controller \mathbf{K}_{1s} . The input is $\mathbf{u}_t = \mathbf{K}_{(t)} \mathbf{x}_t + \mathbf{z}_t$ where exploration noise $\mathbf{z}_t \sim \mathcal{N}(0, \mathbf{z})$. Let $\mathbf{L}_j := \mathbf{A}_j + \mathbf{B}_j \mathbf{K}_j$. Let $\tilde{\mathbf{L}} \in \mathbb{R}^{sn^2 \times sn^2}$ denote the augmented closed-loop state matrix with ij -th $n^2 \times n^2$ block given by $[\tilde{\mathbf{L}}]_{ij} := [\mathbf{T}]_{ji} \mathbf{L}_j \otimes \mathbf{L}_j$. Let $\bar{\zeta} > 0$ and $\underline{\zeta} \in [0, 1)$ be two constants such that $\|\tilde{\mathbf{L}}^k\| \leq \bar{\zeta} \frac{k}{\underline{\zeta}}$. By definition, one available choice for $\bar{\zeta}$ and $\underline{\zeta}$ are $(\tilde{\mathbf{L}})$ and $(\tilde{\mathbf{L}})$.

We define the following cumulative cost conditioned on the initial state \mathbf{x}_0 , initial mode (0) , and controller \mathbf{K}_{1s} .

$$J_T(\mathbf{x}_0, (0), \{\mathbf{K}_{1s}, \mathbf{z}\}) := \sum_{t=1}^T \mathbb{E}[\mathbf{x}_t \mathbf{Q}_{(t)} \mathbf{x}_t + \mathbf{u}_t \mathbf{R}_{(t)} \mathbf{u}_t \mid \mathbf{x}_0, (0), \mathbf{K}_{1s}]. \tag{C.1}$$

The definition of this cumulative cost coincides with the cost $\sum_{t=1}^{T_i} c_{T_0 + T_{i-1} + t}$ in the definition of Regret $_i$ in (5.4) with $\mathbf{x}_0, (0), \mathbf{K}_{1s}$ setting to $\mathbf{x}_0^{(i)}, (i)(0), \mathbf{K}_{1s}^{(i)}$ since Regret $_i$ depends on randomness in \mathcal{F}_{i-1} only through $\mathbf{x}_0^{(i)}, (i)(0), \mathbf{K}_{1s}^{(i)}$. In the remainder of this appendix, for simplicity, we will drop the conditions $\mathbf{x}_0, (0), \mathbf{K}_{1s}$ in the expectation and simply write $\mathbb{E}[\cdot \mid \mathbf{x}_0, (0), \mathbf{K}_{1s}]$ as $\mathbb{E}[\cdot]$. So, for any measurable function f , $\mathbb{E}[f(\mathbf{x}_0, (0), \mathbf{K}_{1s})] = f(\mathbf{x}_0, (0), \mathbf{K}_{1s})$. Note that even though the results in this appendix are derived for conditional expectation $\mathbb{E}[\cdot \mid \mathbf{x}_0, (0), \mathbf{K}_{1s}]$, most of them also hold for the total expectation $\mathbb{E}[\cdot]$.

For the infinite-horizon case, we define the following infinite-horizon average cost without exploration noise \mathbf{z}_t and starting from $\mathbf{x}_0 = 0$.

$$J(0, (0), \{\mathbf{K}_{1s}\}) := \limsup_T \frac{1}{T} J_T(0, (0), \{\mathbf{K}_{1s}, 0\}). \tag{C.2}$$

Let \mathbf{P}_{1s} denote the solution to $\text{cDARE}(\mathbf{A}_{1s}, \mathbf{B}_{1s}, \mathbf{T}, \mathbf{Q}_{1s}, \mathbf{R}_{1s})$ defined in (5.2). Let \mathbf{K}_{1s} denote the resulting infinite-horizon optimal controller computed using \mathbf{P}_{1s} and following (5.1). Note that the infinite-horizon optimal average cost J in (3.4) is achieved if the optimal controller \mathbf{K}_{1s} is used, i.e.

$$J = J(0, (0), \{\mathbf{K}_{1s}\}). \quad (\text{C.3})$$

Note that if the underlying Markov chain \mathbf{T} is ergodic, for any initial state \mathbf{x}_0 and mode (0) , $J = J(\mathbf{x}_0, (0), \{\mathbf{K}_{1s}\})$. Let $\mathbf{L}_i = \mathbf{A}_i + \mathbf{B}_i \mathbf{K}_i$ for all $i \in [s]$ denote the closed-loop state matrix when the optimal controller \mathbf{K}_{1s} is used. Define the augmented state matrix $\tilde{\mathbf{L}}$ such that its ij -th block is given by $[\tilde{\mathbf{L}}]_{ij} := [\mathbf{T}]_{ji} \mathbf{L}_j \otimes \mathbf{L}_j$. From [14], we know \mathbf{K}_{1s} stabilizes the MJS, thus $\rho(\tilde{\mathbf{L}}) < 1$.

Since Regret_i defined in (5.4) can be written as

$$\text{Regret}_i = J_T(\mathbf{x}_0^{(i)}, (i)(0), \{\mathbf{K}_{1s}^{(i)}, \mathbf{z}_i\}) - TJ, \quad (\text{C.4})$$

to evaluate $\text{Regret}(T)$, it suffices to evaluate $J_T(\mathbf{x}_0, (0), \{\mathbf{K}_{1s}, \mathbf{z}\}) - TJ$ for generic $\mathbf{x}_0, (0), \mathbf{K}_{1s}$, and \mathbf{z} . The outline of this Appendix C is as follows.

- In Appendix C.1, we restate perturbation results [18] on $J(0, (0), \{\mathbf{K}_{1s}\}) - J$.
- In Appendix C.2, we evaluate $J_T(\mathbf{x}_0, (0), \{\mathbf{K}_{1s}, \mathbf{z}\}) - TJ(0, (0), \{\mathbf{K}_{1s}\})$. Then, applying the results in Appendix C.1, we can bound the single epoch regret $J_T(\mathbf{x}_0, (0), \{\mathbf{K}_{1s}, \mathbf{z}\}) - TJ$.
- In Appendix C.3, we stitch regrets for all epochs together, and combine them with identification results in Appendix B to bound $\text{Regret}(T)$.

C.1 MJS-LQR Perturbation Results

We first present a lemma on the perturbation of augmented closed-loop state matrix if we use a controller \mathbf{K}_{1s} that is close to the optimal \mathbf{K}_{1s} .

Lemma C.1 (Lemma 9 in [18]) *For an arbitrary controller \mathbf{K}_{1s} , let $\mathbf{L}_i = \mathbf{A}_i + \mathbf{B}_i \mathbf{K}_i$ for all $i \in [s]$, and let $\tilde{\mathbf{L}}$ be the augmented state matrix such that its ij -th block is given by $[\tilde{\mathbf{L}}]_{ij} := [\mathbf{T}]_{ji} \mathbf{L}_j \otimes \mathbf{L}_j$. Assume $\|\mathbf{K}_{1s} - \mathbf{K}_{1s}\| \leq \bar{\kappa}$, where $\bar{\kappa}$ is defined in Table 3. Then, we have*

$$\|\tilde{\mathbf{L}}^k\| \leq (\tilde{\mathbf{L}}) \left(\frac{1+\bar{\kappa}}{2}\right)^k, \quad \forall k \in \mathbb{N}, \quad (\text{C.5})$$

$$(\tilde{\mathbf{L}}) \leq \frac{1+\bar{\kappa}}{2}. \quad (\text{C.6})$$

Thus controller \mathbf{K}_{1s} is stabilizing.

The following perturbation results show how much the infinite-horizon average cost deviates depending on the deviations from the optimal controller and how much the optimal controller deviates depending on the model accuracy for the MJS-LQR problem.

Lemma C.2 (Perturbation of Infinite-horizon MJS-LQR, Corollary 11 and Theorem 6 in [18]) *Infinite-horizon MJS-LQR($\mathbf{A}_{1s}, \mathbf{B}_{1s}, \mathbf{T}, \mathbf{Q}_{1s}, \mathbf{R}_{1s}$) problems have the following perturbation results. Note that notations $\bar{\kappa}, \bar{\mathbf{A}}, \bar{\mathbf{B}}, \bar{\mathbf{T}}$, and $C_{\mathbf{A}, \mathbf{B}, \mathbf{T}}^{\mathbf{K}}$ are defined in Table 3.*

1. Suppose we have an arbitrary controller \mathbf{K}_{1s} such that $\|\mathbf{K}_{1s} - \mathbf{K}_{1s}\| \leq \bar{\kappa}$. Then, we have

$$J(0, (0), \{\mathbf{K}_{1s}\}) - J \leq C_{\mathbf{K}}^J \|\mathbf{w}\| \|\mathbf{K}_{1s} - \mathbf{K}_{1s}\|^2. \quad (\text{C.7})$$

2. Suppose there is an arbitrary MJS($\hat{\mathbf{A}}_{1s}, \hat{\mathbf{B}}_{1s}, \hat{\mathbf{T}}$) with $\bar{\mathbf{A}}, \bar{\mathbf{B}} := \max\{\|\hat{\mathbf{A}}_{1s} - \mathbf{A}_{1s}\|, \|\hat{\mathbf{B}}_{1s} - \mathbf{B}_{1s}\|\} \leq \bar{\mathbf{A}}, \bar{\mathbf{B}}, \bar{\mathbf{T}}$, and $\bar{\mathbf{T}} := \|\hat{\mathbf{T}} - \mathbf{T}\| \leq \bar{\mathbf{A}}, \bar{\mathbf{B}}, \bar{\mathbf{T}}$. Then, there exists an optimal controller \mathbf{K}_{1s} to the infinite-horizon MJS-LQR($\hat{\mathbf{A}}_{1s}, \hat{\mathbf{B}}_{1s}, \hat{\mathbf{T}}, \mathbf{Q}_{1s}, \mathbf{R}_{1s}$) and it can be computed using (5.1) and (5.2), and we have

$$\|\mathbf{K}_{1s} - \mathbf{K}_{1s}\| \leq C_{\mathbf{A}, \mathbf{B}, \mathbf{T}}^{\mathbf{K}} (\bar{\mathbf{A}}, \bar{\mathbf{B}} + \bar{\mathbf{T}}). \quad (\text{C.8})$$

By definition of $\bar{\mathbf{A}}, \bar{\mathbf{B}}, \bar{\mathbf{T}}$, we see $\|\mathbf{K}_{1s} - \mathbf{K}_{1s}\| \leq \bar{\kappa}$, thus Lemma C.1 is applicable.

C.2 Single Epoch Regret Analysis

Recall the definitions of $\tilde{\mathbf{B}}_t$ and $\tilde{\mathbf{z}}_t$ in (A.4) of Appendix A.1. Furthermore, we define

$$\tilde{\mathbf{R}}_t = \mathbf{I}_{n^2} \otimes \mathbf{I}_{n^2}, \quad \tilde{\mathbf{R}}_t = \sum_{i=1}^S \tilde{\mathbf{z}}_t(i) \mathbf{R}_i. \quad (\text{C.9})$$

For a set of matrices $\mathbf{V}_{1:S}$, define the following reshaping mapping

$$\mathcal{H}\left(\begin{bmatrix} \mathbf{V}_1 \\ \vdots \\ \mathbf{V}_S \end{bmatrix}\right) = \begin{bmatrix} \text{vec}(\mathbf{V}_1) \\ \vdots \\ \text{vec}(\mathbf{V}_S) \end{bmatrix}, \quad (\text{C.10})$$

and let \mathcal{H}^{-1} denote the inverse mapping of \mathcal{H} . Let

$$\mathbf{M}_i := \mathbf{Q}_i + \mathbf{K}_i \mathbf{R}_i \mathbf{K}_i, \quad \mathbf{M} := [\mathbf{M}_1, \dots, \mathbf{M}_S]. \quad (\text{C.11})$$

We define

$$\begin{aligned} N_{0,t} &= \text{tr}\left(\mathbf{M} \mathcal{H}^{-1}\left(\tilde{\mathbf{L}}^t \begin{bmatrix} \text{vec}(\mathbf{z}_1(0)) \\ \vdots \\ \text{vec}(\mathbf{z}_S(0)) \end{bmatrix}\right)\right), \\ N_{\mathbf{z},1,t} &= \text{tr}\left(\mathbf{M} \mathcal{H}^{-1}\left(\left(\tilde{\mathbf{B}}_t + \tilde{\mathbf{L}} \tilde{\mathbf{B}}_{t-1} + \dots + \tilde{\mathbf{L}}^{t-1} \tilde{\mathbf{B}}_1\right) \text{vec}(\mathbf{z})\right)\right), \\ N_{\mathbf{w},t} &= \text{tr}\left(\mathbf{M} \mathcal{H}^{-1}\left(\left(\tilde{\mathbf{z}}_t + \tilde{\mathbf{L}} \tilde{\mathbf{z}}_{t-1} + \dots + \tilde{\mathbf{L}}^{t-1} \tilde{\mathbf{z}}_1\right) \text{vec}(\mathbf{w})\right)\right), \\ N_{\mathbf{z},2,t} &= \text{tr}(\tilde{\mathbf{R}}_t \mathbf{z}), \end{aligned} \quad (\text{C.12})$$

and

$$S_{0,T} = \sum_{t=1}^T N_{0,t}, \quad S_{\mathbf{z},1,T} = \sum_{t=1}^T N_{\mathbf{z},1,t}, \quad S_{\mathbf{w},T} = \sum_{t=1}^T N_{\mathbf{w},t}, \quad S_{\mathbf{z},2,T} = \sum_{t=1}^T N_{\mathbf{z},2,t}. \quad (\text{C.13})$$

First, we provide the exact expression for the cumulative cost. It will be used later to analyze the regret.

Lemma C.3 (Cumulative Cost Expression) *For the cost $J_T(\mathbf{x}_0, \mathbf{z}(0), \{\mathbf{K}_{1:S}, \mathbf{z}\})$ defined in (C.1), we have*

$$J_T(\mathbf{x}_0, \mathbf{z}(0), \{\mathbf{K}_{1:S}, \mathbf{z}\}) = S_{0,T} + S_{\mathbf{z},1,T} + S_{\mathbf{z},2,T} + S_{\mathbf{w},T}. \quad (\text{C.14})$$

Proof For the expected cost at time t , we have

$$\begin{aligned} \mathbb{E}[\mathbf{x}_t \mathbf{Q}_t \mathbf{x}_t + \mathbf{u}_t \mathbf{R}_t \mathbf{u}_t] &= \sum_{i=1}^S \text{tr}\left(\mathbb{E}[\mathbf{Q}_t \mathbf{x}_t \mathbf{x}_t^T \mathbf{1}_{\{(t)=i\}}] + \mathbb{E}[\mathbf{R}_t \mathbf{u}_t \mathbf{u}_t^T \mathbf{1}_{\{(t)=i\}}]\right) \\ &= \sum_{i=1}^S \text{tr}\left(\left(\mathbf{Q}_i + \mathbf{K}_i \mathbf{R}_i \mathbf{K}_i\right) \mathbf{z}_i(t) + \tilde{\mathbf{z}}_t(i) \mathbf{R}_i \mathbf{z}\right) \\ &= \sum_{i=1}^S \text{tr}\left(\mathbf{M}_i \mathbf{z}_i(t)\right) + N_{\mathbf{z},2,t}, \end{aligned} \quad (\text{C.15})$$

where the second equality follows since $\mathbf{u}_t = \mathbf{K}_t \mathbf{x}_t + \mathbf{z}_t$. Now plugging in the dynamics of $\mathbf{z}_i(t)$ in Lemma A.3, we can conclude the proof. \blacksquare

Next, before proceeding, we provide several properties regarding the operator $\text{tr}(\mathbf{M} \mathcal{H}(\cdot))$ that shows up in (C.12) and (C.13), which will be used later to evaluate $J_T(\mathbf{x}_0, \mathbf{z}(0), \{\mathbf{K}_{1:S}, \mathbf{z}\}) - TJ(0, \mathbf{z}(0), \{\mathbf{K}_{1:S}\})$.

Lemma C.4 (Properties of Cost Building Bricks) *For any $t, t' \in \mathbb{N}$, we have*

$$(L1) \quad \text{tr}(\mathbf{M} \mathcal{H}^{-1}(\tilde{\mathbf{L}}^t \mathbf{v})) \leq \sqrt{ns} \|\mathbf{M}_{1:S}\| \|\tilde{\mathbf{L}}^t\| \|\mathbf{v}\|, \quad \text{where } \mathbf{v} := [\text{vec}(\mathbf{V}_1), \dots, \text{vec}(\mathbf{V}_S)] \text{ for some } \mathbf{V}_{1:S} \text{ such that } \mathbf{V}_i \geq 0 \text{ for all } i \in [S];$$

$$(L2) \quad \text{tr}(\mathbf{M}\mathcal{H}^{-1}(\tilde{\mathbf{L}}^t \tilde{\mathbf{B}}_{\mathcal{H}} \text{vec}(\mathbf{z}))) \leq n\sqrt{s} \|\mathbf{M}_{1s}\| \|\tilde{\mathbf{L}}^t\| \|\mathbf{B}_{1s}\|^2 \|\mathbf{z}\|;$$

$$(L3) \quad \text{tr}(\mathbf{M}\mathcal{H}^{-1}(\tilde{\mathbf{L}}^t \tilde{\mathbf{v}}_{\mathcal{H}} \text{vec}(\mathbf{w}))) \leq n\sqrt{s} \|\mathbf{M}_{1s}\| \|\tilde{\mathbf{L}}^t\| \|\mathbf{w}\|;$$

$$(L4) \quad |\text{tr}(\mathbf{M}\mathcal{H}^{-1}(\tilde{\mathbf{L}}^t (\tilde{\mathbf{v}}_{\mathcal{H}} - \tilde{\mathbf{v}}_{\mathcal{H}}) \text{vec}(\mathbf{w})))| \leq M_C n\sqrt{s} \|\mathbf{M}_{1s}\| \|\tilde{\mathbf{L}}^t\| \|\mathbf{w}\| \frac{t}{M_C}, \text{ where } M_C \text{ are } M_C \text{ are given in Definition A.2, and } t \text{ is given in (C.9)}$$

Proof Let $[\cdot]_i$ denote the i th sub-block of an $s \times 1$ block matrix. Let vec^{-1} denote the inverse mapping of vec , i.e., $\text{vec}^{-1}([\mathbf{v}_1, \dots, \mathbf{v}_r]) = [\mathbf{v}_1, \dots, \mathbf{v}_r]$ for a set of vectors $\{\mathbf{v}_i\}_{i=1}^r$. It can be easily seen that for any set of matrices $\mathbf{A}, \mathbf{B}, \mathbf{C}$ and \mathbf{X} , we have $\mathbf{A}\mathbf{X}\mathbf{B} = \mathbf{C}$ if and only if $(\mathbf{B} \otimes \mathbf{A})\text{vec}(\mathbf{X}) = \text{vec}(\mathbf{C})$. This together with the definitions of $\tilde{\mathbf{B}}_t, \tilde{\mathbf{v}}_t$ in (A.4), $\tilde{\mathbf{L}}^t, \tilde{\mathbf{R}}_t$ in (C.9), and $\mathcal{H}(\cdot)$ in (C.10) yields the following preliminary results

$$[\mathcal{H}^{-1}(\tilde{\mathbf{L}}^t \mathbf{v})]_i \geq 0, \quad (\text{C.16a})$$

$$\text{vec}^{-1}([\tilde{\mathbf{B}}_{\mathcal{H}} \text{vec}(\mathbf{z})]_i) \geq 0, \quad (\text{C.16b})$$

$$\text{vec}^{-1}([\tilde{\mathbf{v}}_{\mathcal{H}} \text{vec}(\mathbf{w})]_i) \geq 0, \quad (\text{C.16c})$$

$$\text{vec}^{-1}([\tilde{\mathbf{v}}_{\mathcal{H}} - \tilde{\mathbf{v}}_{\mathcal{H}} | \text{vec}(\mathbf{w})]_i) \geq 0, \quad (\text{C.16d})$$

$$|\text{tr}(\mathbf{M}\mathcal{H}^{-1}(\tilde{\mathbf{L}}^t (\tilde{\mathbf{v}}_{\mathcal{H}} - \tilde{\mathbf{v}}_{\mathcal{H}}) \text{vec}(\mathbf{w})))| \leq \text{tr}(\mathbf{M}\mathcal{H}^{-1}(\tilde{\mathbf{L}}^t |\tilde{\mathbf{v}}_{\mathcal{H}} - \tilde{\mathbf{v}}_{\mathcal{H}}| \text{vec}(\mathbf{w}))), \quad (\text{C.16e})$$

where $|\cdot|$ here denotes the element-wise absolute value of a matrix. Now, let us consider (L1). We observe that

$$\begin{aligned} \text{tr}(\mathbf{M}\mathcal{H}^{-1}(\tilde{\mathbf{L}}^t \mathbf{v})) &= \text{tr}\left(\sum_{i=1}^s \mathbf{M}_i [\mathcal{H}^{-1}(\tilde{\mathbf{L}}^t \mathbf{v})]_i\right) \leq \|\mathbf{M}_{1s}\| \cdot \text{tr}\left(\sum_{i=1}^s [\mathcal{H}^{-1}(\tilde{\mathbf{L}}^t \mathbf{v})]_i\right) \\ &\leq \sqrt{n} \|\mathbf{M}_{1s}\| \sum_{i=1}^s \|\mathcal{H}^{-1}(\tilde{\mathbf{L}}^t \mathbf{v})\|_i, \end{aligned} \quad (\text{C.17})$$

where the first inequality uses (C.16a) and the definition that $\|\mathbf{M}_{1s}\| = \max_i \|\mathbf{M}_i\|$; and the last inequality follows from Cauchy-Schwarz inequality and the fact that $[\mathcal{H}^{-1}(\tilde{\mathbf{L}}^t \mathbf{v})]_i \in \mathbb{R}^{n \times n}$. Now, for the last term on the R.H.S. of (C.17), we have

$$\begin{aligned} \sum_{i=1}^s \|\mathcal{H}^{-1}(\tilde{\mathbf{L}}^t \mathbf{v})\|_i &\leq \sum_{i=1}^s \|\mathcal{H}^{-1}(\tilde{\mathbf{L}}^t \mathbf{v})\|_F \leq \sqrt{s} \sqrt{\sum_{i=1}^s \|\mathcal{H}^{-1}(\tilde{\mathbf{L}}^t \mathbf{v})\|_F^2} \\ &= \sqrt{s} \|\mathcal{H}^{-1}(\tilde{\mathbf{L}}^t \mathbf{v})\|_F \\ &= \sqrt{s} \|\tilde{\mathbf{L}}^t \mathbf{v}\| \\ &\leq \sqrt{s} \|\tilde{\mathbf{L}}^t\| \|\mathbf{v}\|, \end{aligned} \quad (\text{C.18})$$

where the second equality holds since \mathcal{H}^{-1} is a reshaping operator, and $\tilde{\mathbf{L}}^t \mathbf{v}$ is a vector. Substituting (C.18) into (C.17) gives (L1).

To show (L2), we combine (C.16b) with (L1) to get $\text{tr}(\mathbf{M}\mathcal{H}^{-1}(\tilde{\mathbf{L}}^t \tilde{\mathbf{B}}_{\mathcal{H}} \text{vec}(\mathbf{z}))) \leq \sqrt{ns} \|\mathbf{M}_{1s}\| \|\tilde{\mathbf{L}}^t\| \|\tilde{\mathbf{B}}_{\mathcal{H}} \text{vec}(\mathbf{z})\|$. Then, using the upper bound for $\|\tilde{\mathbf{B}}_{\mathcal{H}} \text{vec}(\mathbf{z})\|$ derived in (A.12) completes the proof of (L2).

To establish (L3), we combine (C.16c) with (L1) to obtain

$$\text{tr}(\mathbf{M}\mathcal{H}^{-1}(\tilde{\mathbf{L}}^t \tilde{\mathbf{v}}_{\mathcal{H}} \text{vec}(\mathbf{w}))) \leq \sqrt{ns} \|\mathbf{M}_{1s}\| \|\tilde{\mathbf{L}}^t\| \|\tilde{\mathbf{v}}_{\mathcal{H}} \text{vec}(\mathbf{w})\|. \quad (\text{C.19})$$

Then, using the upper bound for $\|\tilde{\mathbf{v}}_{\mathcal{H}} \text{vec}(\mathbf{w})\|$ derived in (A.13) gives (L2).

Finally, let us consider (L4). It follows from (C.16e) and (C.16d) in conjunction with (L1) that

$$|\text{tr}(\mathbf{M}\mathcal{H}^{-1}(\tilde{\mathbf{L}}^t |\tilde{\mathbf{v}}_{\mathcal{H}} - \tilde{\mathbf{v}}_{\mathcal{H}}| \text{vec}(\mathbf{w})))| \leq \sqrt{ns} \|\mathbf{M}_{1s}\| \|\tilde{\mathbf{L}}^t\| \|\tilde{\mathbf{v}}_{\mathcal{H}} - \tilde{\mathbf{v}}_{\mathcal{H}}\| \|\text{vec}(\mathbf{w})\|. \quad (\text{C.20})$$

Now, using (C.9), we obtain

$$\begin{aligned}
\| \tilde{r} - \tilde{r} - \text{vec}(\mathbf{w}) \| &= \sqrt{\sum_{i=1}^s \| [\tilde{r}]_i - [\tilde{r}]_i + \text{vec}(\mathbf{w}) \|^2} \\
&= \sqrt{\sum_{i=1}^s \| r(i) - (i) + \text{vec}(\mathbf{w}) \|^2} \\
&= \| r - \| \text{vec}(\mathbf{w}) \| \\
&\leq \| r - \|_1 \| \mathbf{w} \|_F \\
&\leq MC\sqrt{n} \| \mathbf{w} \|_{MC}^t,
\end{aligned}$$

where the last line follows from Definition A.2. Substituting the above inequality in (C.20) completes the proof of (L4). \blacksquare

The following lemma provides a bound for the difference $J_T(\mathbf{x}_0, (0), \{\mathbf{K}_{1s}, \mathbf{z}\}) - TJ(0, (0), \{\mathbf{K}_{1s}\})$ using an arbitrary stabilizing controller \mathbf{K}_{1s} . Based on this result, we will provide in Proposition C.6 a uniform upper bound for this difference when using any controllers \mathbf{K}_{1s} that are close to \mathbf{K}_{1s} .

Lemma C.5 *For an arbitrary stabilizing controller \mathbf{K}_{1s} , we have*

$$\begin{aligned}
&J_T(\mathbf{x}_0, (0), \{\mathbf{K}_{1s}, \mathbf{z}\}) - TJ(0, (0), \{\mathbf{K}_{1s}\}) \\
&\leq \sqrt{nS} \|\mathbf{M}_{1s}\| \cdot \|\mathbf{x}_0\|^2 + \frac{n\sqrt{S} \tilde{\mathbf{L}}}{1 - \tilde{\mathbf{L}}} \|\mathbf{M}_{1s}\| \|\mathbf{B}_{1s}\|^2 \| \mathbf{z} \| T \\
&\quad + n \|\mathbf{R}_{1s}\| \| \mathbf{z} \| T + n\sqrt{S} MC \tilde{\mathbf{L}} \|\mathbf{M}_{1s}\| \| \mathbf{w} \| \frac{MC}{MC - \tilde{\mathbf{L}}} \left(\frac{MC}{1 - MC} - \frac{\tilde{\mathbf{L}}}{1 - \tilde{\mathbf{L}}} \right),
\end{aligned} \tag{C.21}$$

where MC and MC are given in Definition A.2, $\tilde{\mathbf{L}}$ and $\tilde{\mathbf{L}}$ are constants defined in the beginning of Appendix C, and $\mathbf{M} = [\mathbf{M}_1, \dots, \mathbf{M}_s]$ with $\mathbf{M}_i = \mathbf{Q}_i + \mathbf{K}_i \mathbf{R}_i \mathbf{K}_i$.

Proof From Lemma C.3, we know

$$\begin{aligned}
J_T(\mathbf{x}_0, (0), \{\mathbf{K}_{1s}, \mathbf{z}\}) &= S_{0,T} + S_{z,1,T} + S_{z,2,T} + S_{w,T}, \\
J(0, (0), \{\mathbf{K}_{1s}\}) &= \limsup_T \frac{1}{T} (S_{0,T} + S_{w,T}) =: S_0 + S_w.
\end{aligned}$$

where $S_0 := \limsup_T \frac{1}{T} S_{0,T}$ and $S_w := \limsup_T \frac{1}{T} S_{w,T}$. Next, we will evaluate each term on the RHSs separately.

For $S_{0,T}$, letting $\mathbf{s}_0 = \begin{bmatrix} \text{vec}(\mathbf{x}_1(0)) \\ \vdots \\ \text{vec}(\mathbf{x}_s(0)) \end{bmatrix}$, we have

$$\begin{aligned}
S_{0,T} &= \sum_{t=1}^T \text{tr}(\mathbf{M} \mathcal{H}^{-1}(\tilde{\mathbf{L}}^t \mathbf{s}_0)) \leq \sqrt{nS} \|\mathbf{M}_{1s}\| \|\tilde{\mathbf{L}}^t\| \|\mathbf{s}_0\| \\
&\leq \sqrt{nS} \|\mathbf{M}_{1s}\| \cdot \mathbb{E}[\|\mathbf{x}_0\|^2] \\
&= \sqrt{nS} \|\mathbf{M}_{1s}\| \cdot \|\mathbf{x}_0\|^2,
\end{aligned}$$

where the second line follows from Item (L1) in Lemma C.4; the third line follows from (A.11) in Lemma A.4. And from the discussion at the beginning of Appendix C, we can get rid of $\mathbb{E}[\cdot]$. Then it is easy to see $S_0 = 0$, as long as $\|\mathbf{x}_0\|^2$ is bounded.

For $S_{\mathbf{z},1,T}$, we have

$$\begin{aligned}
S_{\mathbf{z},1,T} &= \sum_{t=1}^T \sum_{t'=0}^{t-1} \text{tr}(\mathbf{M}\mathcal{H}^{-1}(\tilde{\mathbf{L}}^{t'} \tilde{\mathbf{B}}_{t-t'} \text{vec}(\mathbf{z}))) \\
&\leq n\sqrt{s} \|\mathbf{M}_{1s}\| \|\mathbf{B}_{1s}\|^2 \|\mathbf{z}\| \left(\sum_{t=1}^T \sum_{t'=0}^{t-1} \|\tilde{\mathbf{L}}^{t'}\| \right) \\
&\leq \frac{n\sqrt{s} \bar{\rho}}{1 - \bar{\rho}} \|\mathbf{M}_{1s}\| \|\mathbf{B}_{1s}\|^2 \|\mathbf{z}\| T,
\end{aligned} \tag{C.22}$$

where the first inequality follows from Item (L2) in Lemma C.4, and the second inequality follows from the fact $\|\tilde{\mathbf{L}}^{t'}\| \leq \bar{\rho} \frac{t'}{\bar{\rho}}$

For $S_{\mathbf{z},2,T}$, we have

$$S_{\mathbf{z},2,T} = \sum_{t=1}^T \text{tr} \left(\sum_{i=1}^s \mathbf{r}_i(i) \mathbf{R}_i \mathbf{z} \right) \leq n \|\mathbf{R}_{1s}\| \|\mathbf{z}\| T. \tag{C.23}$$

For $S_{\mathbf{w},T}$, we have

$$S_{\mathbf{w},T} = \sum_{t=1}^T \sum_{t'=0}^{t-1} \text{tr}(\mathbf{M}\mathcal{H}^{-1}(\tilde{\mathbf{L}}^{t'} \tilde{\mathbf{w}}_{t-t'} \text{vec}(\mathbf{w}))). \tag{C.24}$$

To evaluate it, we first define the following terms:

$$S_{\mathbf{w},T}^{(\cdot)} := \sum_{t=1}^T \sum_{t'=0}^{t-1} \text{tr}(\mathbf{M}\mathcal{H}^{-1}(\tilde{\mathbf{L}}^{t'} \tilde{\mathbf{w}}_{t-t'} \text{vec}(\mathbf{w}))), \tag{C.25}$$

$$S_{\mathbf{w}}^{(\cdot)} := \limsup_T \frac{1}{T} S_{\mathbf{w},T}^{(\cdot)}, \tag{C.26}$$

where $\tilde{\mathbf{w}}$ is defined in (C.9). Note that $S_{\mathbf{w},T}^{(\cdot)}$ and $S_{\mathbf{w}}^{(\cdot)}$ are the counterparts of $S_{\mathbf{w},T}$ and $S_{\mathbf{w}}$ except that the initial mode distribution $\tilde{\mathbf{w}}_0$ is the stationary distribution $\tilde{\mathbf{w}}$.

Then, we have

$$\begin{aligned}
|S_{\mathbf{w},T} - S_{\mathbf{w},T}^{(\cdot)}| &= \left| \sum_{t=1}^T \sum_{t'=0}^{t-1} \text{tr}(\mathbf{M}\mathcal{H}^{-1}(\tilde{\mathbf{L}}^{t'} (\tilde{\mathbf{w}}_{t-t'} - \tilde{\mathbf{w}}) \text{vec}(\mathbf{w}))) \right| \\
&\leq MC n \sqrt{s} \|\mathbf{M}_{1s}\| \|\mathbf{w}\| \left(\sum_{t=1}^T \sum_{t'=0}^{t-1} \|\tilde{\mathbf{L}}^{t'}\| \frac{t-t'}{MC} \right) \\
&\leq MC n \sqrt{s} \|\mathbf{M}_{1s}\| \|\mathbf{w}\| \left(\sum_{t=1}^T \sum_{t'=0}^{t-1} \bar{\rho} \frac{t'}{\bar{\rho}} \frac{t-t'}{MC} \right) \\
&\leq n \sqrt{s} MC \bar{\rho} \|\mathbf{M}_{1s}\| \|\mathbf{w}\| \frac{MC}{MC - \bar{\rho}} \left(\frac{MC}{1 - MC} - \frac{\bar{\rho}}{1 - \bar{\rho}} \right)
\end{aligned} \tag{C.27}$$

where the first inequality follows from Item (L4) in Lemma C.4. Thus,

$$S_{\mathbf{w}} = \limsup_T \frac{1}{T} S_{\mathbf{w},T} = \limsup_T \frac{1}{T} (S_{\mathbf{w},T} - S_{\mathbf{w},T}^{(\cdot)}) + \limsup_T \frac{1}{T} S_{\mathbf{w},T}^{(\cdot)} = S_{\mathbf{w}}^{(\cdot)}. \tag{C.28}$$

Since $\sum_{t=1}^T \sum_{t'=0}^{t-1} \tilde{\mathbf{L}}^{t'} = (\mathbf{I} - \tilde{\mathbf{L}})^{-1} T - (\mathbf{I} - \tilde{\mathbf{L}})^{-2} \tilde{\mathbf{L}} (\mathbf{I} - \tilde{\mathbf{L}}^T)$ and $\sum_{t'=0}^{t-1} \tilde{\mathbf{L}}^{t'} = (\mathbf{I} - \tilde{\mathbf{L}})^{-1}$ we have

$$\begin{aligned}
S_{\mathbf{w}} = S_{\mathbf{w}}^{(\cdot)} &= \text{tr}(\mathbf{M}\mathcal{H}^{-1}(\limsup_T \frac{1}{T} \sum_{t=1}^T \sum_{t'=0}^{t-1} \tilde{\mathbf{L}}^{t'} \tilde{\mathbf{w}}_{t-t'} \text{vec}(\mathbf{w}))) \\
&= \text{tr}(\mathbf{M}\mathcal{H}^{-1}((\mathbf{I} - \tilde{\mathbf{L}})^{-1} \tilde{\mathbf{w}} \text{vec}(\mathbf{w}))) \\
&= \sum_{t'=0}^{\infty} \text{tr}(\mathbf{M}\mathcal{H}^{-1}(\tilde{\mathbf{L}}^{t'} \tilde{\mathbf{w}}_{t-t'} \text{vec}(\mathbf{w}))).
\end{aligned} \tag{C.29}$$

Thus,

$$\begin{aligned}
TS_{\mathbf{w}} &= TS_{\mathbf{w}}^{(\cdot)} = \sum_{t=1}^T \sum_{t'=0}^{t-1} \text{tr}(\mathbf{M}\mathcal{H}^{-1}(\tilde{\mathbf{L}}^{t'} \sim \text{vec}(\mathbf{w}))) \\
&\geq \sum_{t=1}^T \sum_{t'=0}^{t-1} \text{tr}(\mathbf{M}\mathcal{H}^{-1}(\tilde{\mathbf{L}}^{t'} \sim \text{vec}(\mathbf{w}))) \\
&= S_{\mathbf{w},T}^{(\cdot)}
\end{aligned} \tag{C.30}$$

where the inequality holds since each trace summand is non-negative. Therefore,

$$\begin{aligned}
S_{\mathbf{w},T} &\leq S_{\mathbf{w},T}^{(\cdot)} + |S_{\mathbf{w},T} - S_{\mathbf{w},T}^{(\cdot)}| \\
&\stackrel{\text{(C.28)}}{\leq} TS_{\mathbf{w}} + |S_{\mathbf{w},T} - S_{\mathbf{w},T}^{(\cdot)}| \\
&\stackrel{\text{(C.27)}}{\leq} TS_{\mathbf{w}} + n\sqrt{s} \frac{MC}{\tilde{\mathbf{L}}} \|\mathbf{M}_{1s}\| \|\mathbf{w}\| \frac{MC}{MC - \tilde{\mathbf{L}}} \left(\frac{MC}{1 - MC} - \frac{\tilde{\mathbf{L}}}{1 - \tilde{\mathbf{L}}} \right).
\end{aligned} \tag{C.31}$$

Finally, combining all the results we have so far, we have

$$\begin{aligned}
&J_T(\mathbf{x}_0, (0), \{\mathbf{K}_{1s}, \mathbf{z}\}) - TJ(0, (0), \{\mathbf{K}_{1s}\}) \\
&= S_{0,T} + S_{z,1,T} + S_{z,2,T} + S_{\mathbf{w},T} - T(S_0 + S_{\mathbf{w}}) \\
&\leq \sqrt{n\tilde{s}} \|\mathbf{M}_{1s}\| \cdot \|\mathbf{x}_0\|^2 \\
&\quad + \frac{n\sqrt{s} \tilde{\mathbf{L}}}{1 - \tilde{\mathbf{L}}} \|\mathbf{M}_{1s}\| \|\mathbf{B}_{1s}\|^2 \|\mathbf{z}\| T \\
&\quad + n \|\mathbf{R}_{1s}\| \|\mathbf{z}\| T \\
&\quad + n\sqrt{s} \frac{MC}{\tilde{\mathbf{L}}} \|\mathbf{M}_{1s}\| \|\mathbf{w}\| \frac{MC}{MC - \tilde{\mathbf{L}}} \left(\frac{MC}{1 - MC} - \frac{\tilde{\mathbf{L}}}{1 - \tilde{\mathbf{L}}} \right)
\end{aligned} \tag{C.32}$$

which concludes the proof. \blacksquare

We now provide a uniform upper bound on the regret $J_T(\mathbf{x}_0, (0), \{\mathbf{K}_{1s}, \mathbf{z}\}) - TJ$ for any stabilizing controller \mathbf{K}_{1s} that is close enough to the optimal controller \mathbf{K}_{1s} .

Proposition C.6 *For every \mathbf{K}_{1s} such that $\|\mathbf{K}_{1s} - \mathbf{K}_{1s}\| \leq \bar{\kappa}$, we have*

$$\begin{aligned}
J_T(\mathbf{x}_0, (0), \{\mathbf{K}_{1s}, \mathbf{z}\}) - TJ &\leq C_{\mathbf{K}}^J \|\mathbf{K}_{1s} - \mathbf{K}_{1s}\|^2 \|\mathbf{w}\| T \\
&\quad + \sqrt{n\tilde{s}M} \|\mathbf{x}_0\|^2 \\
&\quad + n\sqrt{s} \frac{2(\tilde{\mathbf{L}})}{1 - \tilde{\mathbf{L}}} \|\mathbf{B}_{1s}\|^2 M \|\mathbf{z}\| T \\
&\quad + n \|\mathbf{R}_{1s}\| \|\mathbf{z}\| T \\
&\quad + n\sqrt{s} \frac{2(\tilde{\mathbf{L}}) MC M}{2 MC - 1 - \tilde{\mathbf{L}}} \frac{MC}{1 - MC} \left(\frac{MC}{1 - MC} - \frac{1 + \tilde{\mathbf{L}}}{1 - \tilde{\mathbf{L}}} \right) \|\mathbf{w}\|,
\end{aligned} \tag{C.33}$$

where $M := \|\mathbf{Q}_{1s}\| + 4\|\mathbf{R}_{1s}\| \|\mathbf{K}_{1s}\|^2$, and $\bar{\kappa}$ and $C_{\mathbf{K}}^J$ are defined in Table 3.

Proof When $\|\mathbf{K}_{1s} - \mathbf{K}_{1s}\| \leq \bar{\kappa}$, from Lemma C.1, we know $\|\tilde{\mathbf{L}}^k\| \leq (\tilde{\mathbf{L}})(\frac{1+\tilde{\mathbf{L}}}{2})^k$, thus we could set $\tilde{\mathbf{L}}$ and $\tilde{\mathbf{L}}$ to be $(\tilde{\mathbf{L}})$ and $\frac{1+\tilde{\mathbf{L}}}{2}$. By definition, we know $\bar{\kappa} \leq \|\mathbf{K}_{1s}\|$, thus $\|\mathbf{M}_{1s}\| \leq \|\mathbf{Q}_{1s}\| + \|\mathbf{R}_{1s}\| \|\mathbf{K}_{1s}\|^2 \leq$

$\|\mathbf{Q}_{1s}\| + \|\mathbf{R}_{1s}\|(\|\mathbf{K}_{1s}\| + \bar{\kappa})^2 \leq \|\mathbf{Q}_{1s}\| + 4\|\mathbf{R}_{1s}\|\|\mathbf{K}_{1s}\|^2 = M$. Then applying Lemma C.5, we have

$$\begin{aligned}
& J_T(\mathbf{x}_0, (0), \{\mathbf{K}_{1s}, \mathbf{z}\}) - TJ(0, (0), \{\mathbf{K}_{1s}\}) \\
& \leq \sqrt{n\bar{s}M}\|\mathbf{x}_0\|^2 \\
& \quad + n\sqrt{\bar{s}} \frac{2(\tilde{\mathbf{L}})\|\mathbf{B}_{1s}\|^2 M}{1-} \|\mathbf{z}\|_T \\
& \quad + n\|\mathbf{R}_{1s}\| \|\mathbf{z}\|_T \\
& \quad + n\sqrt{\bar{s}} \frac{2(\tilde{\mathbf{L}})_{MC} M_{MC}}{2_{MC} - 1 -} \left(\frac{MC}{1 - MC} - \frac{1}{1 -} \right) \|\mathbf{w}\|
\end{aligned} \tag{C.34}$$

Now note that when $\|\mathbf{K}_{1s} - \mathbf{K}_{1s}\| \leq \bar{\kappa}$, we have $J(0, (0), \{\mathbf{K}_{1s}\}) - J \leq C_{\mathbf{K}}^J \|\mathbf{w}\| \|\mathbf{K}_{1s} - \mathbf{K}_{1s}\|^2$ using Lemma C.2. Combining this with (C.34), we could conclude the proof. \blacksquare

C.3 Stitching Every Epoch

In this section, we stitch the upper bounds on Regret_i for every epoch i and build a bound on the overall regret $\text{Regret}(T)$.

We define the estimation error after epoch i as $\bar{\mathbf{A}}, \bar{\mathbf{B}} = \max\{\|\mathbf{A}_{1s}^{(i)} - \mathbf{A}_{1s}\|, \|\mathbf{B}_{1s}^{(i)} - \mathbf{B}_{1s}\|\}$, $\bar{\mathbf{T}} = \|\mathbf{T}^{(i)} - \mathbf{T}\|$. Furthermore, we also define $\bar{\mathbf{K}} := \|\mathbf{K}_{1s}^{(i)} - \mathbf{K}_{1s}\|$ where \mathbf{K}_{1s} is the optimal controller for the infinite-horizon MJS-LQR($\mathbf{A}_{1s}, \mathbf{B}_{1s}, \mathbf{T}, \mathbf{Q}_{1s}, \mathbf{R}_{1s}$). We define the following events for every epoch i .

$$\begin{aligned}
\mathcal{A}_i &= \left\{ \text{Regret}_i \leq \mathcal{O} \left(sp \left(\bar{\mathbf{A}}, \bar{\mathbf{B}} + \bar{\mathbf{T}} \right)^2 \frac{2}{\mathbf{w}} T_i + \sqrt{n\bar{s}} \|\mathbf{x}_0^{(i)}\|^2 + n\sqrt{\bar{s}} \frac{2}{\mathbf{z}_i} T_i + c_A \right) \right\} \\
\mathcal{B}_i &= \left\{ \bar{\mathbf{A}}, \bar{\mathbf{B}} \leq \bar{\mathbf{A}}, \bar{\mathbf{B}}, \bar{\mathbf{T}}, \bar{\mathbf{T}} \leq \bar{\mathbf{A}}, \bar{\mathbf{B}}, \bar{\mathbf{T}}, \bar{\mathbf{K}} \leq \bar{\kappa} \right\} \\
\mathcal{C}_i &= \left\{ \bar{\mathbf{A}}, \bar{\mathbf{B}} \leq \mathcal{O} \left(\log \left(\frac{1}{id,i} \right) \frac{\mathbf{z}_i + \mathbf{w}}{\mathbf{z}_i \min} \frac{\sqrt{\bar{s}}(n+p) \log(T_i)}{\sqrt{T_i}} \right) \right\}, \\
\bar{\mathbf{T}} &\leq \mathcal{O} \left(\log \left(\frac{1}{id,i} \right) \frac{1}{\min} \sqrt{\frac{\log(T_i)}{T_i}} \right) \\
\mathcal{D}_i &= \left\{ \|\mathbf{x}_0^{(i+1)}\|^2 = \|\mathbf{x}_{T_i}^{(i)}\|^2 \leq \frac{\bar{\mathbf{x}}_0^2}{\mathbf{x}_{0,i}} \right\}.
\end{aligned} \tag{C.35}$$

where $c_A, \bar{\mathbf{x}}_0$ are constants, $\bar{\mathbf{A}}, \bar{\mathbf{B}}, \bar{\mathbf{T}}$ and $\bar{\kappa}$ are defined in Table 3, and id,i and $\mathbf{x}_{0,i}$ within $[0, 1]$ denotes the failure probability for event \mathcal{C}_i and \mathcal{D}_i . Note that $\mathcal{O}(\cdot)$ hides terms that are invariant to epochs such as $\|\mathbf{A}_{1s}\|, \|\mathbf{B}_{1s}\|$, etc.

Event \mathcal{A}_i describes how epoch i regret depends on initial state $\|\mathbf{x}_0^{(i)}\|^2$, exploration noise variance $\frac{2}{\mathbf{z}_i}$, and the accuracy of the estimated MJS dynamics $\mathbf{A}_{1s}^{(i-1)}, \mathbf{B}_{1s}^{(i-1)}, \hat{\mathbf{T}}$ after epoch $i-1$, which is used to compute epoch i controller $\mathbf{K}_{1s}^{(i)}$. Event \mathcal{B}_i indicates whether the estimated dynamics and resulting controllers are good enough. \mathcal{C}_i describes the dynamics estimation error after epoch i , and when epoch T_i is chosen appropriately, \mathcal{B}_i can be implied. Lastly, event \mathcal{D}_i bounds the initial state of each epoch, as the initial state plays a vital role in regret upper bound \mathcal{A}_i . We see events $\mathcal{A}_{i+1}, \mathcal{B}_i, \mathcal{C}_i, \mathcal{D}_i$ are \mathcal{F}_i -measurable, i.e. these events can be determined using random variables $\mathbf{x}_0, \mathbf{w}_t, \mathbf{z}_t, \hat{\mathbf{T}}(t)$ up to epoch i . Let $\mathcal{E}_i := \mathcal{A}_{i+1} \cap \mathcal{B}_i \cap \mathcal{C}_i \cap \mathcal{D}_i$. Note that even though \mathcal{A}_{i+1} is for the conditional expected regret of the epoch $i+1$ with randomness coming from $\mathbf{x}_0^{(i+1)} = \mathbf{x}_{T_i}^{(i)}, \hat{\mathbf{T}}^{(i+1)}(0) = \hat{\mathbf{T}}(T_i)$, and controller $\mathbf{K}_{1s}^{(i+1)}$ computed from $\mathbf{A}_{1s}^{(i)}, \mathbf{B}_{1s}^{(i)}, \hat{\mathbf{T}}^{(i)}$, thus \mathcal{A}_{i+1} is \mathcal{F}_i -measurable.

Then, we have the following results regarding the conditional probabilities of these events. First, Proposition C.7 says given the event \mathcal{B}_{i-1} (a good controller is applied during epoch i) and event \mathcal{D}_i (the initial state of epoch i , $\mathbf{x}_0^{(i)}$ is bounded), then \mathcal{D}_i could occur, i.e. $\mathbf{x}_{T_i}^{(i)}$, the final state of epoch i , a.k.a. $\mathbf{x}_0^{(i+1)}$ the initial state of epoch $i+1$, is also bounded.

Proposition C.7 Suppose $\frac{\bar{n}\bar{s}^{-T_i}}{\mathbf{x}_{0,i-1}} < 1$ and $\bar{x}_0^2 \geq \frac{n\bar{s}(\mathbf{B}_{1,s}^{2+1})\frac{2}{\mathbf{w}}}{(1-\bar{\kappa})(1-\frac{\bar{n}\bar{s}^{-T_i}}{\mathbf{x}_{0,i-1}})}$ for $i \geq 1$. Then,

$$\mathbb{P}(\mathcal{D}_i | \cap_{j=0}^{i-1} \mathcal{E}_j) = \mathbb{P}(\mathcal{D}_i | \mathcal{B}_{i-1} \mathcal{D}_{i-1}) > 1 - \mathbf{x}_{0,i},$$

and $\mathbb{P}(\mathcal{D}_0) \geq 1 - \mathbf{x}_{0,0}$.

Proof For epoch $i = 1, 2, \dots$, given event \mathcal{B}_{i-1} , we know $\bar{\kappa}^{(i+1)} \leq \bar{\kappa}$. Let $\tilde{\mathbf{L}}^{(i)}$ denote the augmented closed-loop state matrix. By Lemma C.1, we know $\|(\tilde{\mathbf{L}}^{(i)})^k\| \leq (\tilde{\mathbf{L}}^{(i)})^{(1+\frac{\bar{\kappa}}{2})^k}$. Thus, if we pick $\bar{\kappa} := \max\{(\tilde{\mathbf{L}}^{(0)}), (\tilde{\mathbf{L}}^{(1)})\}$, $\bar{\kappa} := \max\{(\tilde{\mathbf{L}}^{(0)}), \frac{1+\bar{\kappa}}{2}\}$, this can be generalized to $i = 0$ case, i.e. for every epoch $i = 0, 1, 2, \dots$, we have $\|(\tilde{\mathbf{L}}^{(i)})^k\| \leq \bar{\kappa}^k$.

For $i = 1, 2, \dots$, event \mathcal{D}_{i-1} implies $\|\mathbf{x}_0^{(i)}\|^2 \leq \frac{\bar{x}_0^2}{\mathbf{x}_{0,i-1}}$. Then, according to Lemma A.4, we know

$$\begin{aligned} \mathbb{E}[\|\mathbf{x}_{T_i}^{(i)}\|^2 | \mathcal{B}_{i-1}, \mathcal{D}_{i-1}] &\leq \sqrt{n\bar{s}} \cdot \bar{s}^{-T_i} \frac{\bar{x}_0^2}{\mathbf{x}_{0,i-1}} + n\sqrt{\bar{s}}(\|\mathbf{B}_{1,s}\|^2 \frac{2}{\sqrt{T_i}} + \frac{2}{\mathbf{w}}) \frac{1}{1-\bar{\kappa}} \\ &\leq \frac{\sqrt{n\bar{s}} \cdot \bar{s}^{-T_i}}{\mathbf{x}_{0,i-1}} \bar{x}_0^2 + (1 - \frac{\sqrt{n\bar{s}} \cdot \bar{s}^{-T_i}}{\mathbf{x}_{0,i-1}}) \bar{x}_0^2 \\ &\leq \bar{x}_0^2, \end{aligned} \quad (\text{C.36})$$

where the second line follows from the assumptions in the proposition statement. Using Markov inequality, we have

$$\mathbb{P}(\|\mathbf{x}_{T_i}^{(i)}\|^2 \leq \frac{\bar{x}_0^2}{\mathbf{x}_{0,i}} | \mathcal{B}_{i-1}, \mathcal{D}_{i-1}) \geq 1 - \mathbf{x}_{0,i},$$

which implies $\mathbb{P}(\mathcal{D}_i | \mathcal{B}_{i-1}, \mathcal{D}_{i-1}) \geq 1 - \mathbf{x}_{0,i}$. For $i = 0$, similarly, we have $\mathbb{E}[\|\mathbf{x}_{T_0}^{(0)}\|^2] \leq n\sqrt{\bar{s}}(\|\mathbf{B}_{1,s}\|^2 \frac{2}{T_0} + \frac{2}{\mathbf{w}}) \frac{1}{1-\bar{\kappa}} \leq \bar{x}_0^2$, thus $\mathbb{P}(\mathcal{D}_0) \geq 1 - \mathbf{x}_{0,0}$.

Finally, note that given a good stabilizing controller (event \mathcal{B}_{i-1}) and a bounded initial state (event \mathcal{D}_{i-1}) for epoch i , the final state of epoch i only depends on randomness in epoch i , thus $\mathbb{P}(\mathcal{D}_i | \cap_{j=0}^{i-1} \mathcal{E}_j) = \mathbb{P}(\mathcal{D}_i | \mathcal{B}_{i-1} \mathcal{D}_{i-1})$. ■

Proposition C.8 describes that given the event \mathcal{C}_i (the estimated MJS dynamics after epoch i has estimation errors decays with T_i), when epoch i has length T_i large enough, then the event \mathcal{B}_i (the estimated dynamics and controllers computed with it will be good enough) occurs.

Proposition C.8 Suppose every epoch i has length $T_i \geq \underline{T}_{rgt, (id,i, T_i)}$. Then,

$$\mathbb{P}(\mathcal{B}_i | \mathcal{C}_i, \cap_{j=0}^{i-1} \mathcal{E}_j) = \mathbb{P}(\mathcal{B}_i | \mathcal{C}_i) = 1$$

Proof When \mathcal{C}_i occurs, since $\frac{2}{\mathbf{z}_i} = \frac{2}{\mathbf{w}_i T_i}$, we have

$$\bar{\mathbf{A}}_{\mathbf{B}}^{(i)} \leq \mathcal{O}(\log(\frac{1}{id,i}) \frac{\sqrt{\bar{s}}(n+p) \log(T_i)}{\min T_i^{0.25}}), \quad \bar{\mathbf{T}}^{(i)} \leq \mathcal{O}(\log(\frac{1}{id,i}) \frac{1}{\min T_i} \sqrt{\log(T_i)}).$$

We know when $T_i \geq \mathcal{O}(\frac{\bar{s}(n+p)^{-4}}{\min \bar{\mathbf{A}}_{\mathbf{B}, \mathbf{T}} \log(\frac{1}{id,i})} \log^4(T_i)) =: \underline{T}_{rgt, (id,i, T_i)}$, we have $\bar{\mathbf{A}}_{\mathbf{B}}^{(i)} \leq \bar{\mathbf{A}}_{\mathbf{B}, \mathbf{T}}^{(i)}$, $\bar{\mathbf{T}}^{(i)} \leq \bar{\mathbf{A}}_{\mathbf{B}, \mathbf{T}}^{(i)}$. Then according to Lemma C.2, we have $\bar{\kappa}^{(i+1)} \leq \bar{\kappa}$. Thus $\mathbb{P}(\mathcal{B}_i | \mathcal{C}_i) = 1$. Finally, note that given the estimation error sample complexity in \mathcal{C}_i for epoch i , events happen before epoch i does not influence \mathcal{B}_i , so $\mathbb{P}(\mathcal{B}_i | \mathcal{C}_i, \cap_{j=0}^{i-1} \mathcal{E}_j) = \mathbb{P}(\mathcal{B}_i | \mathcal{C}_i) = 1$. ■

Next, Proposition C.9 says given the \mathcal{B}_{i-1} (a good controller is used in epoch i), then the event \mathcal{C}_i could occur, i.e. dynamics learned using the trajectory of epoch i , will be accurate enough.

Proposition C.9 For $c_{\mathbf{x}} \geq \underline{c}_{\mathbf{x}}(\bar{\cdot}, \bar{\cdot})$, $c_{\mathbf{z}} \geq \underline{c}_{\mathbf{z}}$, $T_i \geq \max\{\underline{I}_{MC,1}(\frac{id,i}{8}), \underline{I}_{id,N}(\frac{id,i}{2}, \bar{\cdot}, \bar{\cdot})\}$, we have for $i = 1, 2, \dots$,

$$\mathbb{P}(\mathcal{C}_i | \cap_{j=0}^{i-1} \mathcal{E}_j) = \mathbb{P}(\mathcal{C}_i | \mathcal{B}_{i-1}) \geq 1 - id,i. \quad (\text{C.37})$$

And $\mathbb{P}(\mathcal{C}_0) \geq 1 - id,0$.

Proof By Lemma B.1, we know for every epoch $i = 0, 1, \dots$, when $T_i \geq \underline{I}_{MC,1}(\frac{id,i}{8})$, we have with probability at least $1 - \frac{id,i}{2}$, $\frac{(i)}{\mathbf{T}} \leq \mathcal{O}\left(\log\left(\frac{1}{id,i}\right) \frac{1}{\min} \sqrt{\frac{\log(T_i)}{T_i}}\right)$.

For epoch $i = 1, 2, \dots$, given event \mathcal{B}_{i-1} , we know $\frac{(i)}{\mathbf{K}} \leq \bar{\mathbf{K}}$. Let $\tilde{\mathbf{L}}^{(i)}$ denote the augmented closed-loop state matrix. By Lemma C.1, we know $\|(\tilde{\mathbf{L}}^{(i)})^k\| \leq (\tilde{\mathbf{L}})(\frac{1+}{2})^k$. Thus, if we pick $\bar{\cdot} := \max\{(\tilde{\mathbf{L}}^{(0)}), (\tilde{\mathbf{L}})\}$, $\bar{\cdot} := \max\{(\tilde{\mathbf{L}}^{(0)}), \frac{1+}{2}\}$, this can be generalized to $i = 0$ case, i.e. for every epoch $i = 0, 1, 2, \dots$, we have $\|(\tilde{\mathbf{L}}^{(i)})^k\| \leq \bar{\cdot}^k$.

Suppose $c_{\mathbf{x}} \geq \underline{c}_{\mathbf{x}}(\bar{\cdot}, \bar{\cdot})$, $c_{\mathbf{z}} \geq \underline{c}_{\mathbf{z}}$, and $T_i \geq \underline{I}_{id,N}(\frac{id,i}{2}, \bar{\cdot}, \bar{\cdot})$ hold for $i = 0, 1, \dots$. Then, from Theorem B.17, we know for every $i = 0, 1, \dots$, with probability at least $1 - \frac{id,i}{2}$, $\frac{(i)}{\mathbf{A}, \mathbf{B}} \leq \mathcal{O}\left(\log\left(\frac{1}{id,i}\right) \frac{z,i+ w}{z,i \min} \frac{\bar{s}(n+p)\log(T_i)}{T_i}\right)$.

Applying union bound to $\frac{(i)}{\mathbf{T}}$ and $\frac{(i)}{\mathbf{A}, \mathbf{B}}$, we could show $\mathbb{P}(\mathcal{C}_0) \geq 1 - id,i$ and $\mathbb{P}(\mathcal{C}_i | \mathcal{B}_{i-1}, \mathcal{D}_{i-1}) \geq 1 - id,i$. Finally, note that given a good stabilizing controller (event \mathcal{B}_{i-1}) and bounded initial state (event \mathcal{D}_{i-1}) for epoch i , the estimation error sample complexity (event \mathcal{C}_i) does not depend on events happen before epoch i , so $\mathbb{P}(\mathcal{C}_i | \cap_{j=0}^{i-1} \mathcal{E}_j) = \mathbb{P}(\mathcal{C}_i | \mathcal{B}_{i-1}, \mathcal{D}_{i-1})$. ■

Finally, Proposition C.10 simply describes how the regret of epoch i depends on the accuracy of the estimated dynamics after epoch $i-1$.

Proposition C.10 For $\mathcal{A}_i - \mathcal{C}_i$ given in (C.35), we have

$$\mathbb{P}(\mathcal{A}_i | \mathcal{B}_{i-1}, \mathcal{C}_{i-1}, \mathcal{D}_{i-1}, \cap_{j=0}^{i-1} \mathcal{E}_j) = \mathbb{P}(\mathcal{A}_i | \mathcal{B}_{i-1}) = 1.$$

Proof From Proposition C.6, we know that for every epoch $i = 1, 2, \dots$, given $\|\mathbf{K}_{1s}^{(i)} - \mathbf{K}_{1s}\| \leq \bar{\mathbf{K}}$ in \mathcal{B}_{i-1} , we have with probability 1

$$\begin{aligned} \text{Regret}_i &\leq C_{\mathbf{K}}^J \|\mathbf{K}_{1s}^{(i)} - \mathbf{K}_{1s}\|^2 \frac{2}{w} T_i \\ &\quad + \sqrt{nsM} \|\mathbf{x}_0^{(i)}\|^2 \\ &\quad + n\sqrt{s} \frac{2}{1-} (\tilde{\mathbf{L}}) \|\mathbf{B}_{1s}\|^2 M \frac{2}{z,i} T_i \\ &\quad + n \|\mathbf{R}_{1s}\| \frac{2}{z,i} T_i \\ &\quad + n\sqrt{s} \frac{2}{2} \frac{(\tilde{\mathbf{L}})_{MC} M}{MC-1-} \frac{MC}{1-} \left(\frac{MC}{1-} - \frac{1+}{MC} - \frac{1+}{1-}\right) \frac{2}{w}. \end{aligned} \quad (\text{C.38})$$

Let c_A denote the last term in (C.38), which is a constant over epochs. Note that from $\frac{(i-1)}{\mathbf{A}, \mathbf{B}} \leq \bar{\mathbf{A}, \mathbf{B}, \mathbf{T}}$, $\frac{(i-1)}{\mathbf{T}} \leq \bar{\mathbf{A}, \mathbf{B}, \mathbf{T}}$ in event \mathcal{B}_{i-1} , we know $\|\mathbf{K}_{1s}^{(i)} - \mathbf{K}_{1s}\| \leq C_{\mathbf{A}, \mathbf{B}, \mathbf{T}}^{\mathbf{K}} \left(\frac{(i-1)}{\mathbf{A}, \mathbf{B}} + \frac{(i-1)}{\mathbf{T}}\right)$ by Lemma C.2. Plugging this into (C.38), we have

$$\text{Regret}_i \leq O\left(s \cdot \rho \left(\frac{(i-1)}{\mathbf{A}, \mathbf{B}} + \frac{(i-1)}{\mathbf{T}}\right)^2 \frac{2}{w} T_i + \sqrt{ns} \|\mathbf{x}_0^{(i)}\|^2 + n\sqrt{s} \frac{2}{z,i} T_i + c_A\right) \quad (\text{C.39})$$

where term $s \cdot \rho$ comes from term $s \min\{n, p\}$ in the definition of $C_{\mathbf{K}}^J$ in Appendix C.1. This shows $\mathbb{P}(\mathcal{A}_i | \mathcal{B}_{i-1}) = 1$. Finally, note that given a good controller (event \mathcal{B}_{i-1}) for epoch i , the regret for epoch i can be upper bounded (event \mathcal{A}_i) without dependence on other events, thus $\mathbb{P}(\mathcal{A}_i | \mathcal{B}_{i-1}, \mathcal{C}_{i-1}, \mathcal{D}_{i-1}, \cap_{j=0}^{i-1} \mathcal{E}_j) = \mathbb{P}(\mathcal{A}_i | \mathcal{B}_{i-1})$. ■

C.3.1 Proof for Theorem 5.1

Theorem C.11 (Complete version of Theorem 5.1) *Assume that the initial state $\mathbf{x}_0 = 0$, and Assumption 1 and 2 hold. Suppose $c_{\mathbf{x}} \geq c_{\mathbf{x}}(\cdot, \cdot)$, $c_{\mathbf{z}} \geq c_{\mathbf{z}}$, $T_0 \geq \mathcal{O}(\mathcal{I}_{\text{rgt}}(\cdot, T_0))$, and $\bar{x}_0^2 = \frac{n \bar{s}(\mathbf{B}_{1:s}^{2+1}) \frac{2}{\mathbf{w}}}{(1-\cdot)(1-\frac{\bar{s}}{n\bar{s}} - T_0 \frac{2}{3})}$. Then, with probability at least $1 - \cdot$, Algorithm 2 achieves*

$$\text{Regret}(T) \leq \mathcal{O}\left(\frac{s^2 \rho(n^2 + \rho^2) \frac{2}{\mathbf{w}}}{2_{\min}} \log^2\left(\frac{\log^2(T)}{\cdot}\right) \log^2(T) \sqrt{T} + \frac{\sqrt{n\bar{s}} \log^3(T)}{\cdot}\right). \quad (\text{C.40})$$

Proof In this proof, we will first show the intersected event $\cap_i \mathcal{E}_i = \cap_i \{\mathcal{A}_{i+1} \cap \mathcal{B}_i \cap \mathcal{C}_i \cap \mathcal{D}_i\}$ implies the desired regret bound, then we evaluate the occurrence probability of $\cap_i \mathcal{E}_i$ using Proposition C.8 to C.10. In the following, we set $i_{d,i} = \mathbf{x}_{0,i} = \frac{3}{2} \cdot \frac{2}{(i+1)^2}$.

With the choices $T_i = T_{i-1}$, $\frac{2}{\mathbf{z}_i} = -\frac{2}{T_i}$, and $i_{d,i} = \mathbf{x}_{0,i} = \frac{3}{2} \cdot \frac{2}{(i+1)^2}$, event $\mathcal{E}_i = \mathcal{A}_{i+1} \cap \mathcal{B}_i \cap \mathcal{C}_i \cap \mathcal{D}_i$ implies the following.

$$\begin{aligned} & \text{Regret}_{i+1} \\ & \leq \mathcal{O}(1) \log^2\left(\frac{(i+1)^2}{\cdot}\right) sp\left(\frac{\mathbf{z}_i + \mathbf{w}}{\mathbf{z}_i \min} \cdot \frac{\sqrt{s}(n+\rho) \log(T_i)}{\sqrt{T_i}} + \frac{\sqrt{\log(T_i)}}{\min \sqrt{T_i}}\right)^2 \frac{2}{\mathbf{w}} T_{i+1} \\ & \quad + \mathcal{O}\left(\frac{(i+1)^2}{\cdot}\right) \sqrt{n\bar{s}} \bar{x}_0^2 + \mathcal{O}(n\sqrt{s} \frac{2}{\mathbf{z}_{i+1}} T_{i+1}) + \mathcal{O}(1) \\ & \leq \mathcal{O}(1) \log^2\left(\frac{(i+1)^2}{\cdot}\right) \frac{s^2 \rho(n^2 + \rho^2)}{2_{\min}} \frac{(\mathbf{z}_i + \mathbf{w})^2}{\mathbf{z}_i} \frac{2}{\mathbf{w}} \log^2(T_i) \\ & \quad + \mathcal{O}\left(\frac{(i+1)^2}{\cdot}\right) \sqrt{n\bar{s}} \bar{x}_0^2 + \mathcal{O}(n\sqrt{s} \frac{2}{\mathbf{z}_{i+1}} T_{i+1}) \\ & \leq \mathcal{O}(1) \log^2\left(\frac{(i+1)^2}{\cdot}\right) \frac{s^2 \rho(n^2 + \rho^2)}{2_{\min}} \left(\frac{4}{\mathbf{z}_i} \log^2(T_i) + \frac{2}{\mathbf{z}_i} T_i\right) \\ & \quad + \mathcal{O}\left(\frac{(i+1)^2}{\cdot}\right) \sqrt{n\bar{s}} \bar{x}_0^2 \\ & \leq \mathcal{O}(1) \log^2\left(\frac{(i+1)^2}{\cdot}\right) \frac{s^2 \rho(n^2 + \rho^2)}{2_{\min}} \frac{2}{\mathbf{w}} \sqrt{T_i} \log^2(T_i) + \mathcal{O}\left(\frac{(i+1)^2}{\cdot}\right) \sqrt{n\bar{s}} \bar{x}_0^2 \end{aligned} \quad (\text{C.41})$$

We have $M := \mathcal{O}(\log(\frac{T}{T_0}))$ epochs at time T . Using the fact $T_i = \mathcal{O}(T_0^{-i})$, event $\cap_{i=0}^{M-1} \mathcal{E}_i$ implies

$$\begin{aligned} & \text{Regret}(T) \\ & = \mathcal{O}\left(\sum_{i=1}^M \text{Regret}_i\right) \\ & \leq \mathcal{O}(1) \log^2\left(\frac{\log^2(T)}{\cdot}\right) \frac{s^2 \rho(n^2 + \rho^2) \frac{2}{\mathbf{w}}}{2_{\min}} \left(\sum_{i=1}^M \sqrt{T_i} \log^2(T_i)\right) + \mathcal{O}\left(\frac{\sqrt{n\bar{s}} \log^3(T)}{\cdot}\right) \end{aligned} \quad (\text{C.42})$$

For the term $\sum_{i=1}^M \sqrt{T_i} \log^2(T_i)$, we have

$$\begin{aligned}
& \sum_{i=1}^M \sqrt{T_i} \log^2(T_i) \\
& \leq \mathcal{O}(1) \sqrt{T_0} \left(\log^2(T_0) \sum_{i=1}^M \sqrt{-i} + \log^2(\cdot) \sum_{i=1}^M \sqrt{-i} \rho^2 \right) \\
& \leq \mathcal{O}(1) \sqrt{T_0} \log^2(\cdot) \sum_{i=1}^M \sqrt{-i} \rho^2 \\
& \leq \mathcal{O}(1) \sqrt{T_0} \log^2(\cdot) M \sqrt{-M} \left(\frac{\sqrt{-}}{\sqrt{-}-1} \right)^3 \left(M - \frac{1}{\sqrt{-}} \right) \\
& \leq \mathcal{O}(1) \sqrt{T} \log^2(\cdot) \frac{\log(\frac{T}{T_0})}{\log(\cdot)} \left(\frac{\sqrt{-}}{\sqrt{-}-1} \right)^3 \left(\frac{\log(\frac{T}{T_0})}{\log(\cdot)} - \frac{1}{\sqrt{-}} \right) \\
& \leq \mathcal{O}(1) \sqrt{T} \log(\frac{T}{T_0}) \left(\frac{\sqrt{-}}{\sqrt{-}-1} \right)^3 \left(\log(\frac{T}{T_0}) - \sqrt{-} \log(\cdot) \right) \\
& \leq \mathcal{O}(\log^2(T) \sqrt{T}).
\end{aligned} \tag{C.43}$$

Plugging this back into (C.42), we have

$$\text{Regret}(T) \leq \mathcal{O} \left(\frac{s^2 \rho (\eta^2 + \rho^2)}{2 \min} \frac{2}{\mathbf{w}} \log^2 \left(\frac{\log^2(T)}{\log(\cdot)} \right) \log^2(T) \sqrt{T} + \frac{\sqrt{nS} \log^3(T)}{\log(\cdot)} \right) \tag{C.44}$$

which shows the regret bound in (C.40).

Now we are only left to show the occurrence probability of regret bound (C.40) is larger than $1 - \gamma$. To do this, we will combine Proposition C.7, C.8, C.9, and C.10 over all $i = 0, 1, \dots, M-1$. Note that for each individual i , these propositions hold only when certain prerequisite conditions on hyper-parameters $c_{\mathbf{x}}, c_{\mathbf{z}}, T_0$, and \bar{x}_0 are satisfied. We first show that under the choices $T_i = T_{i-1}$, $\frac{2}{\mathbf{z}_i} = \frac{2}{T_i}$, and $\bar{x}_{0,i} = \frac{3}{2} \cdot \frac{1}{(i+1)^2}$ these hyper-parameter conditions can be satisfied for all $i = 0, 1, \dots, M-1$.

- Proposition C.7 requires that for $i = 1, 2, \dots$, conditions $\frac{\bar{nS}^{-i} T_0^{-i} i^2}{3} < 1$ and $\bar{x}_0^2 \geq \frac{n \bar{s} (\mathbf{B}_{1:s}^{2+1}) \frac{2}{\mathbf{w}}}{(1-\gamma)(1-\frac{\bar{nS}^{-i} T_0^{-i} i^2}{3})}$ need to be satisfied. One can check when $T_0 \geq \frac{1}{\log(1-\gamma)} \max\{\frac{2}{\log(\cdot)}, \log(\frac{2}{3} \frac{\bar{nS}^{-i}}{\cdot})\} =: \underline{I}_{\mathbf{x}_0}(\cdot)$, and picking $\bar{x}_0^2 \geq \frac{n \bar{s} (\mathbf{B}_{1:s}^{2+1}) \frac{2}{\mathbf{w}}}{(1-\gamma)(1-\frac{\bar{nS}^{-i} T_0^{-i} i^2}{3})}$ would suffice.
- Proposition C.8 requires that for $i = 0, 1, \dots$, condition $T_0^{-i} \geq \underline{I}_{\text{rgt},-}(\frac{3}{2(i+1)^2}, T_0^{-i})$ holds, which can be satisfied when one chooses $T_0 \geq \mathcal{O}(\underline{I}_{\text{rgt},-}(\cdot, T_0))$.
- Proposition C.9 requires the following to hold: $c_{\mathbf{x}} \geq \underline{c}_{\mathbf{x}}(\cdot, \cdot)$, $c_{\mathbf{z}} \geq \underline{c}_{\mathbf{z}}$, and $T_0^{-i} \geq \max\{\underline{I}_{MC,1}(\frac{3}{8^{2i^2}}), \underline{I}_{id,N}(\frac{3}{2^{2(i+1)^2}}, \cdot, \cdot)\}$. The last one can be satisfied when $T_0 \geq \mathcal{O}(\max\{\underline{I}_{MC,1}(\cdot), \underline{I}_{id,N}(\cdot, \cdot, \cdot)\})$.
- Proposition C.10 requires no conditions on hyper-parameters.

Therefore, when $c_{\mathbf{x}} \geq \underline{c}_{\mathbf{x}}(\cdot, \cdot)$, $c_{\mathbf{z}} \geq \underline{c}_{\mathbf{z}}$,

$$T_0 \geq \mathcal{O}(\max\{\underline{I}_{\mathbf{x}_0}(\cdot), \underline{I}_{\text{rgt},-}(\cdot, T_0), \underline{I}_{MC,1}(\cdot), \underline{I}_{id,N}(\cdot, \cdot, \cdot)\}) =: \mathcal{O}(\underline{I}_{\text{rgt}}(\cdot, T_0)),$$

we can apply Propositions C.7, C.8, C.9, and C.10 to every epoch $i = 0, 1, \dots, M-1$. First note that Propositions C.7 and C.9 give the following

$$\begin{aligned}
\mathbb{P}(\mathcal{D}_i | \cap_{j=0}^{i-1} \mathcal{E}_j) &= \mathbb{P}(\mathcal{D}_i | \mathcal{B}_{i-1} \mathcal{D}_{i-1}) > 1 - \frac{3}{2(i+1)^2}, \quad \mathbb{P}(\mathcal{D}_0) \geq 1 - \frac{3}{2} \\
\mathbb{P}(\mathcal{C}_i | \cap_{j=0}^{i-1} \mathcal{E}_j) &= \mathbb{P}(\mathcal{C}_i | \mathcal{B}_{i-1}) \geq 1 - \frac{3}{2(i+1)^2}, \quad \mathbb{P}(\mathcal{C}_0) \geq 1 - \frac{3}{2}.
\end{aligned}$$

Then combining the probability bounds in Propositions C.7, C.8, C.9, and C.10, we have

$$\begin{aligned}
& \text{P}(\text{Regret bounds in (C.40) holds}) \\
& \geq \text{P}(\cap_{i=0}^{M-1} \mathcal{E}_i) \\
& = \text{P}(\mathcal{A}_M, \mathcal{B}_{M-1}, \mathcal{C}_{M-1}, \mathcal{D}_{M-1} \mid \cap_{i=0}^{M-2} \mathcal{E}_i) \cdot \text{P}(\cap_{i=0}^{M-2} \mathcal{E}_i) \\
& = \text{P}(\mathcal{C}_{M-1}, \mathcal{D}_{M-1} \mid \cap_{i=0}^{M-2} \mathcal{E}_i) \cdot \text{P}(\cap_{i=0}^{M-2} \mathcal{E}_i) \\
& \geq (1 - \text{id}, M-1 - \mathbf{x}_{0, M-1}) \cdot \text{P}(\cap_{i=0}^{M-2} \mathcal{E}_i) \\
& \geq \prod_{i=0}^{M-1} (1 - \text{id}, i - \mathbf{x}_{0, i}) \\
& \geq 1 - \sum_{i=0}^{M-1} (\text{id}, i + \mathbf{x}_{0, i}) \\
& \geq 1 - \dots
\end{aligned} \tag{C.45}$$

where the last line holds since $\sum_{i=0}^{M-1} \frac{1}{(i+1)^2} \leq \frac{2}{6}$. ■

C.4 Regret Under Uniform Stability — Proof for Theorem 5.2

As we discussed in Section 5.2, under MSS, the regret upper bound in Theorem 5.1 (or the complete version Theorem C.11) involves $\frac{1}{2}$ dependency on failure probability δ . By checking the proof for Theorem C.11, we can see the only source for $\frac{1}{2}$ is event \mathcal{D}_i in (C.35) and the corresponding Proposition C.7, which provides $1 - \delta$ probability bound for event \mathcal{D}_i – the initial state $\mathbf{x}_0^{(i+1)}$ of epoch $i+1$, a.k.a. the final state $\mathbf{x}_{T_i}^{(i)}$ of epoch i , is bounded by $\|\mathbf{x}_0^{(i+1)}\|^2 = \|\mathbf{x}_{T_i}^{(i)}\|^2 \leq \mathcal{O}(\delta)$. In Proposition C.7, we get this bound using Markov inequality $\|\mathbf{x}_{T_i}^{(i)}\|^2 \leq \mathbb{E}[\|\mathbf{x}_{T_i}^{(i)}\|^2] / (1 - \delta)$ and Lemma A.4 which provides an upper bound on the numerator $\mathbb{E}[\|\mathbf{x}_{T_i}^{(i)}\|^2]$ under MSS. From event \mathcal{A}_i in (C.35) we see the regret of epoch i directly depends on its epoch initial state $\|\mathbf{x}_0^{(i)}\|^2$, thus in the final cumulative regret, the cumulative impact of initial states from all epochs, $\sum_i \|\mathbf{x}_0^{(i)}\|^2$ with order $\frac{1}{2}$, will show up, as given in (C.42). Therefore, whether $\frac{1}{2}$ terms can be relaxed directly hinges on whether one could refine Proposition C.7 to get a tighter dependency on δ .

This refinement, however, is not possible under the MSS assumption only, and we can easily construct a toy example to show that the $\frac{1}{2}$ dependency resulting from the Markov inequality cannot be improved. Consider a two-mode, one-dimensional, autonomous MJS:

$$\begin{cases} x_{t+1} = 2x_t \\ x_{t+1} = 0.5x_t \end{cases} \text{ with Markov matrix } \mathbf{T} = \begin{bmatrix} 0.1 & 0.9 \\ 0.1 & 0.9 \end{bmatrix}$$

with $x_0 \sim \mathcal{N}(0, 1)$, and $\text{P}(x_0 = 0) = 0.1$. It is easy to check this MJS is MSS by the spectral radius criterion discussed below Definition 3.1. Also note that with probability 0.1^t , $(0 : t-1) = 1$ and $x_t = 2^t x_0$. Therefore, for any $a > 0$,

$$\text{P}(x_t \geq a) = \sum_{(0 : t-1)} \text{P}(x_t \geq a \mid (0 : t-1)) \text{P}((0 : t-1)) \tag{C.46}$$

$$\geq \text{P}(x_t \geq a \mid (0 : t-1) = 1) \text{P}((0 : t-1) = 1) = 0.1^t \cdot \text{P}(x_0 \geq 2^{-t} a) . \tag{C.47}$$

where the inequality in (C.47) is extremely loose since we condition only on the most improbable event. For standard Gaussian x_0 , $\text{P}(x_0 \geq a) \geq \frac{C}{a} \exp(-\frac{a^2}{2})$ for some absolute constant C . Thus $\text{P}(x_t \geq a) \geq \frac{C \cdot 0.2^t}{a} \exp(-\frac{2^{-2t} a^2}{2})$. From this, we see that for any $a > 0$, any $t \geq \log(a) / \log(2)$, we have $\text{P}(x_t \geq a) \geq C \frac{0.2^t}{ea}$. We can observe that though when t grows slower than $\log(a)$, the tail of x_t has exponential decay, the Markov

inequality decay, i.e. $\frac{1}{a}$, will eventually show up when t gets larger. Interpretation from failure probability perspective is the following: letting $\delta = C\frac{0.2^t}{ea}$, we have $\mathbb{P}(X_t \leq C\frac{0.2^t}{e}) \leq 1 - \delta$, which means any dependency lighter than $\frac{1}{a}$ must have probability less than $1 - \delta$. This further implies that in the regret analysis of adaptive control, in order to obtain better probability dependency, the time horizon has to be limited, which greatly impairs its value in practice.

Intuitively, MSS assumption only provides us with stable behavior of $\|\mathbf{x}_t\|^2$ in the expectation (w.r.t. mode switchings) sense, and having only this first-order moment information is of little use compared with the deterministic Lyapunov stability typically used for LTI systems, which allows one to bound $\|\mathbf{x}_t\|^2$ with only $\log(\frac{1}{\delta})$ dependence ([15, Lemma C.5]). Then, one may wonder naturally: Does there exist a deterministic version of stability for switched systems? Can this stability (if exists) help build similar dependence for switched systems? The answers to both questions are yes and will be discussed in this appendix. In short, if there exists uniform stability for the MJS, we can adapt Proposition C.7 such that $\|\mathbf{x}_0^{(i)}\|^2$ can instead be bounded much more tightly by $\|\mathbf{x}_0^{(i)}\|^2 \leq \mathcal{O}(\log(\frac{1}{\delta}))$, thus the $\frac{1}{a}$ dependency can improve to $\log(\frac{1}{\delta})$ in the regret bound (5.5) (or (C.40)). The final improved regret bound is presented in Theorem 5.2. In order to show it, we will need to adapt Proposition C.7 together with several related results (Lemma A.4, Lemma C.1, Lemma C.2) to the uniform stability case, and we append suffix ‘‘a’’ in the result label to denote the adapted versions.

To begin with, recall \mathbf{K}_{1s} is the optimal controller for the infinite-horizon MJS-LQR($\mathbf{A}_{1s}, \mathbf{B}_{1s}, \mathbf{T}, \mathbf{Q}_{1s}, \mathbf{R}_{1s}$) and define the closed-loop state matrix $\mathbf{L}_i = \mathbf{A}_i + \mathbf{B}_i \mathbf{K}_i$ for all i . We let $\rho_{\mathbf{L}_{1s}}$ denote the joint spectral radius of \mathbf{L}_{1s} , i.e. $\rho_{\mathbf{L}_{1s}} := \lim_{l \rightarrow \infty} \max_{1:l} \max_{[s]^l} \|\mathbf{L}_{i_1} \cdots \mathbf{L}_{i_l}\|^{1/l}$. We say \mathbf{L}_{1s} is uniformly stable if and only if $\rho_{\mathbf{L}_{1s}} < 1$. Similar to in Definition A.1, define $\bar{\rho}_{\mathbf{L}_{1s}} := \sup_{l \in \mathbb{N}} \max_{1:l} \max_{[s]^l} \|\mathbf{L}_{i_1} \cdots \mathbf{L}_{i_l}\| / (\rho_{\mathbf{L}_{1s}})^l$. Note that the pair $\{\rho_{\mathbf{L}_{1s}}, \bar{\rho}_{\mathbf{L}_{1s}}\}$ for uniform stability is just the counterpart of $\{\rho_{\mathbf{L}}, \bar{\rho}_{\mathbf{L}}\}$ for MSS defined in Appendix C. Similar as before, Table 5 lists all the shorthand notations to be used in this appendix for quick reference.

Table 5: Notations — Uniform Stability

$\bar{\rho}_{\mathbf{L}_{1s}}$	$\ \mathbf{B}_{1s}\ ^2 \ \mathbf{z}\ + \ \mathbf{w}\ $ or $\ \mathbf{B}_{1s}\ ^2 \frac{\bar{\rho}_{\mathbf{L}_{1s}}}{2} + \frac{\bar{\rho}_{\mathbf{L}_{1s}}}{2}$ $(1 + \bar{\rho}_{\mathbf{L}_{1s}})/2$
$\bar{\rho}_{\mathbf{L}_{1s}}^{US}$	$\frac{1 - \rho_{\mathbf{L}_{1s}}}{2 \rho_{\mathbf{L}_{1s}}}$
$\bar{\mathbf{K}}$	$\min\{\bar{\mathbf{K}}^{US}, \bar{\mathbf{K}}\}$
$\bar{\mathbf{A}}, \bar{\mathbf{B}}, \bar{\mathbf{T}}$	$\min\{\bar{\mathbf{A}}, \bar{\mathbf{B}}, \bar{\mathbf{T}}, \frac{\bar{\mathbf{K}}}{2C_{\mathbf{A}, \mathbf{B}, \mathbf{T}}}\}$
$\bar{\chi}^{US}$	$2^{-2-2} (6 \max\{\sqrt{\bar{n}}e^{3\bar{n}}, \sqrt{\bar{\rho}}e^{3\bar{\rho}}\} + \frac{5}{(1-\bar{\rho})^2})^2$
$\bar{I}_{\mathbf{x}_0}^{US}(\cdot)$	$\max\{\frac{54^{-4-2}}{(1-\bar{\rho})^{\bar{\chi}^{US}} \log(1-\bar{\rho}) \log(\bar{\rho})}, \frac{1}{\log(1-\bar{\rho})} \log(6^{-2} + \frac{54\bar{n} \bar{\rho}^{-4-2} \log(\bar{\rho}^2/3)}{(1-\bar{\rho})(1-\bar{\rho})^{\bar{\chi}^{US}}})\}$
$\bar{I}_{rgt}^{US}(\cdot, T)$	$\mathcal{O}(\frac{\bar{\rho}^{(n+p)}}{\min_{\mathbf{A}, \mathbf{B}, \mathbf{T}} \bar{\rho}^{-4}} \log(\frac{1}{\bar{\rho}}) \log^4(T))$ $\mathcal{O}(\frac{\bar{\rho}^{(n+p)}}{\min_{\mathbf{A}, \mathbf{B}, \mathbf{T}} \bar{\rho}^{-2}} \log(\frac{1}{\bar{\rho}}) \log^2(T))$ (when \mathbf{B}_{1s} is known)
$\bar{I}_{rgt}^{US}(\cdot, T)$	$\max\{\bar{I}_{\mathbf{x}_0}^{US}(\cdot), \bar{I}_{rgt}^{US}(\cdot, T), \bar{I}_{MC,1}(\cdot), \bar{I}_{id,N}(\underline{\mathbf{L}}, T, \bar{\rho}, \bar{\rho})\}$

The following Lemma A.4a bounds the state \mathbf{x}_t under the designed input in this work. Compared with its counterpart Lemma A.4 which is only able to bound $\mathbb{E}[\|\mathbf{x}_t\|^2]$, Lemma A.4a provides high-probability bound for $\|\mathbf{x}_t\|^2$.

Lemma A.4a Consider an MJS($\mathbf{A}_{1s}, \mathbf{B}_{1s}, \mathbf{T}$) with noise $\mathbf{w}_t \sim \mathcal{N}(0, \mathbf{w})$. Consider controller \mathbf{K}_{1s} , and let \mathbf{L}_{1s} denote the closed-loop state matrices with $\mathbf{L}_i = \mathbf{A}_i + \mathbf{B}_i \mathbf{K}_i$. Assume there exist constants $\bar{\rho}$ and $\bar{\rho} \in [0, 1)$ such that, for any sequence $1:l \in [s]^l$ with any length l , $\|\mathbf{L}_{i_1} \cdots \mathbf{L}_{i_l}\| \leq \bar{\rho}^l$. Let the input be $\mathbf{u}_t = \mathbf{K}_{1s} \mathbf{x}_t + \mathbf{z}_t$

with $\mathbf{z}_t \sim \mathcal{N}(0, \Sigma_{\mathbf{z}})$. Then, for any $t \geq e^{6 \max\{n, p\}}$, with probability at least $1 - \frac{1}{2}$, we have

$$\|\mathbf{x}_t\|^2 \leq 3^{-2} 2^t \|\mathbf{x}_0\|^2 + \frac{18^{-2-2}}{(1-\frac{1}{2})^2} \log\left(\frac{1}{2}\right) + c \quad (\text{C.48})$$

where $\frac{1}{2} := \|\mathbf{w}\| + \|\mathbf{B}_{1_s}\|^2 \|\mathbf{z}\|$ and $c := 2^{-2-2} (6 \max\{\sqrt{n}e^{3n}, \sqrt{p}e^{3p}\} + \frac{5}{(1-\frac{1}{2})^2})^2$.

Proof From the MJS dynamics (3.1) and plugging in the input $\mathbf{u}_t = \mathbf{K}_{(t)} \mathbf{x}_t + \mathbf{z}_t$, we have the following.

$$\begin{aligned} \mathbf{x}_t = & \left(\prod_{h=0}^{t-1} \mathbf{L}_{(h)} \right) \mathbf{x}_0 + \sum_{i=0}^{t-2} \left(\prod_{h=i+1}^{t-1} \mathbf{L}_{(h)} \right) \mathbf{B}_{(i)} \mathbf{z}_i + \mathbf{B}_{(t-1)} \mathbf{z}_{t-1} \\ & \sum_{i=0}^{t-2} \left(\prod_{h=i+1}^{t-1} \mathbf{L}_{(h)} \right) \mathbf{w}_i + \mathbf{w}_{t-1}. \end{aligned} \quad (\text{C.49})$$

Then, by triangle inequality and the assumption that $\|\mathbf{L}_{(1)} \cdots \mathbf{L}_{(i)}\| \leq \frac{1}{i}$, we have

$$\begin{aligned} \|\mathbf{x}_t\| & \leq \frac{1}{t} \|\mathbf{x}_0\| + \|\mathbf{B}_{1_s}\| \sum_{i=0}^{t-1} \frac{1}{t-i-1} \|\mathbf{z}_i\| + \sum_{i=0}^{t-1} \frac{1}{t-i-1} \|\mathbf{w}_i\| \\ & = \frac{1}{t} \|\mathbf{x}_0\| + \|\mathbf{B}_{1_s}\| \sum_{i=0}^{t-1} \frac{1}{i} \|\mathbf{z}_{t-i-1}\| + \sum_{i=0}^{t-1} \frac{1}{i} \|\mathbf{w}_{t-i-1}\|. \end{aligned} \quad (\text{C.50})$$

For each \mathbf{w}_{t-i-1} , using Lemma A.5 (replacing e^{-t} with $\frac{1}{i}$), we have with probability $1 - \frac{1}{i}$,

$$\|\mathbf{w}_{t-i-1}\| \leq \sqrt{3 \|\mathbf{w}\|} \log^{0.5} \left(\frac{1}{\min\{i, \frac{1}{n}\}} \right), \quad (\text{C.51})$$

where $\frac{1}{n} := e^{-(3+2\frac{1}{2})n}$, and n is the dimension of vector \mathbf{w}_{t-i-1} . In the following, for all $i = 0, 1, \dots, t-1$, we set $\frac{1}{i} = \frac{3}{2} \frac{1}{(\bar{i}+1)^2}$. First note that when $i \geq \bar{i} := \sqrt{\frac{3}{2}n} - 1$, we have $\min\{i, \frac{1}{n}\} = \frac{1}{n}$, i.e. $\frac{1}{i} \leq \frac{1}{n}$, and $\min\{i, \frac{1}{n}\} = \frac{1}{i}$ otherwise. Then, applying union bound for all i , we know with probability at least $1 - \frac{1}{2}$,

$$\begin{aligned} \sum_{i=0}^{t-1} \frac{1}{i} \|\mathbf{w}_{t-i-1}\| & \leq \sqrt{3 \|\mathbf{w}\|} \sum_{i=0}^{t-1} \frac{1}{i} \log^{0.5} \left(\frac{1}{\min\{i, \frac{1}{n}\}} \right) \\ & \leq \sqrt{3 \|\mathbf{w}\|} \left(\sum_{i=0}^{t-1} \frac{1}{i} \log^{0.5} \left(\frac{1}{i} \right) + (\bar{i}+1) \log^{0.5} \left(\frac{1}{n} \right) \right). \end{aligned} \quad (\text{C.52})$$

In the above equation for the term $\sum_{i=0}^{t-1} \frac{1}{i} \log^{0.5} \left(\frac{1}{i} \right)$, we have $\sum_{i=0}^{t-1} \frac{1}{i} \log^{0.5} \left(\frac{1}{i} \right) = \sum_i \frac{1}{i} \log^{0.5} \left(\frac{2(i+1)^2}{3} \right) \leq \sum_i \frac{1}{i} \left(\log^{0.5} \left(\frac{1}{i} \right) + \sqrt{2} \log^{0.5} \left(\frac{(i+1)}{3} \right) \right) \leq \frac{1}{1-\frac{1}{2}} \log^{0.5} \left(\frac{1}{1-\frac{1}{2}} \right) + \sqrt{2} \sum_i \frac{1}{i} \frac{(i+1)}{3} \leq \frac{1}{1-\frac{1}{2}} \log^{0.5} \left(\frac{1}{1-\frac{1}{2}} \right) + \frac{\sqrt{2}}{3} \frac{1}{(1-\frac{1}{2})^2}$. And for the term $(\bar{i}+1) \log^{0.5} \left(\frac{1}{n} \right)$ in (C.52), by the definitions of \bar{i} and $\frac{1}{n}$, we have $(\bar{i}+1) \log^{0.5} \left(\frac{1}{n} \right) \leq \sqrt{2n} e^{3n}$. Plugging these two results back into (C.52), we have, with probability at least $1 - \frac{1}{2}$,

$$\sum_{i=0}^{t-1} \frac{1}{i} \|\mathbf{w}_{t-i-1}\| \leq \frac{\sqrt{3 \|\mathbf{w}\|}}{1-\frac{1}{2}} \log^{0.5} \left(\frac{1}{1-\frac{1}{2}} \right) + \frac{5\sqrt{\|\mathbf{w}\|}}{(1-\frac{1}{2})^2} + 3\sqrt{n} e^{3n} \sqrt{\|\mathbf{w}\|}. \quad (\text{C.53})$$

Similarly, with probability at least $1 - \frac{1}{2}$,

$$\sum_{i=0}^{t-1} \frac{1}{i} \|\mathbf{z}_{t-i-1}\| \leq \frac{\sqrt{3 \|\mathbf{z}\|}}{1-\frac{1}{2}} \log^{0.5} \left(\frac{1}{1-\frac{1}{2}} \right) + \frac{5\sqrt{\|\mathbf{z}\|}}{(1-\frac{1}{2})^2} + 3\sqrt{p} e^{3p} \sqrt{\|\mathbf{z}\|}. \quad (\text{C.54})$$

Plugging (C.53) and (C.54) back into (C.50) and applying union bound, we have, with probability $1 - \bar{\epsilon}$,

$$\begin{aligned} \|\mathbf{x}_t\| \leq & \bar{\epsilon} \|\mathbf{x}_0\| + \frac{\sqrt{3} (\sqrt{\|\mathbf{w}\|} + \|\mathbf{B}_{1s}\| \sqrt{\|\mathbf{z}\|})}{(1 - \bar{\epsilon})^2} \log^{0.5} \left(\frac{1}{\bar{\epsilon}} \right) \\ & + (\sqrt{\|\mathbf{w}\|} + \|\mathbf{B}_{1s}\| \sqrt{\|\mathbf{z}\|}) \left(3 \max\{\sqrt{\bar{n}} e^{3\bar{n}}, \sqrt{\bar{p}} e^{3\bar{p}}\} + \frac{5}{(1 - \bar{\epsilon})^2} \right). \end{aligned} \quad (\text{C.55})$$

Taking squares of both sides and using Cauchy-Schwartz inequality, we have

$$\|\mathbf{x}_t\|^2 \leq 3 \bar{\epsilon}^2 \|\mathbf{x}_0\|^2 + \frac{18}{1 - \bar{\epsilon}} \log \left(\frac{1}{\bar{\epsilon}} \right) + c \quad (\text{C.56})$$

where $\bar{\epsilon}^{-2} := \|\mathbf{w}\| + \|\mathbf{B}_{1s}\|^2 \|\mathbf{z}\|$ and $c := 6 \bar{\epsilon}^{-2} (3 \max\{\sqrt{\bar{n}} e^{3\bar{n}}, \sqrt{\bar{p}} e^{3\bar{p}}\} + \frac{5}{(1 - \bar{\epsilon})^2})^2$. \blacksquare

The following Lemma C.1a describes that given a set of matrices that have joint spectral radius smaller than 1, i.e. uniformly stable, moderate perturbation can preserve the uniform stability. On the other hand, its counterpart, Lemma C.1, considers perturbation results for MSS.

Lemma C.1a (Joint Spectral Radius Perturbation) *Assume $\bar{\epsilon} < 1$. For an arbitrary controller \mathbf{K}_{1s} and resulting closed-loop state matrices \mathbf{L}_{1s} with $\mathbf{L}_i = \mathbf{A}_i + \mathbf{B}_i \mathbf{K}_i$, let $\bar{\rho}(\mathbf{L}_{1s})$ denote the joint spectral radius of \mathbf{L}_{1s} . Assume $\|\mathbf{K}_{1s} - \mathbf{K}_{1s}\| \leq \frac{\bar{\epsilon}^{-us}}{2 \bar{\epsilon}^* \mathbf{B}_{1s}}$, then for any sequence $\mathbf{l}_{1:l} \in [s]^l$ with any length l ,*

$$\left\| \prod_{j=1}^l \mathbf{L}_{j, \mathbf{l}_j} \right\| \leq \bar{\epsilon}^{-l} \quad (\text{C.57})$$

$$\bar{\rho}(\mathbf{L}_{1s}) \leq \bar{\epsilon}. \quad (\text{C.58})$$

where $\bar{\epsilon} = \bar{\epsilon}$ and $\bar{\epsilon}^* = \frac{1 + \bar{\epsilon}}{2}$.

Proof Let $\mathbf{E}_j := \mathbf{L}_j - \mathbf{L}_{j, \mathbf{l}_j}$, then we see $\|\mathbf{E}_j\| \leq \|\mathbf{B}_{1s}\| \frac{\bar{\epsilon}^{-us}}{\mathbf{K}}$ and $\prod_{j=1}^l \mathbf{L}_{j, \mathbf{l}_j} = \prod_{j=1}^l (\mathbf{L}_j + \mathbf{E}_j)$. In the expansion of $\prod_{j=1}^l (\mathbf{L}_j + \mathbf{E}_j)$, for each $h = 0, 1, \dots, l$, there are $\binom{l}{h}$ terms, each of which is a product where \mathbf{E} has degree h and \mathbf{L} has degree $l - h$. We let $\mathbf{F}_{h,l}$ with $h = 0, 1, \dots, l$ and $l \in \left[\binom{l}{h} \right]$ to index such terms. Note that $\|\mathbf{F}_{h,l}\| \leq \binom{l}{h} \bar{\epsilon}^{h+1} \bar{\epsilon}^{l-h} (\|\mathbf{B}_{1s}\| \frac{\bar{\epsilon}^{-us}}{\mathbf{K}})^h$. Then, we have

$$\begin{aligned} \left\| \prod_{j=1}^l \mathbf{L}_{j, \mathbf{l}_j} \right\| & \leq \sum_{h=0}^l \sum_{l \in \left[\binom{l}{h} \right]} \|\mathbf{F}_{h,l}\| \\ & \leq \sum_{h=0}^l \binom{l}{h} \bar{\epsilon}^{h+1} \bar{\epsilon}^{l-h} (\|\mathbf{B}_{1s}\| \frac{\bar{\epsilon}^{-us}}{\mathbf{K}})^h \\ & \leq \left(\|\mathbf{B}_{1s}\| \frac{\bar{\epsilon}^{-us}}{\mathbf{K}} + \bar{\epsilon} \right)^l. \end{aligned} \quad (\text{C.59})$$

Then (C.57) follows from the fact that $\frac{\bar{\epsilon}^{-us}}{\mathbf{K}} \leq \frac{1 - \bar{\epsilon}^*}{2 \bar{\epsilon}^* \mathbf{B}_{1s}}$ and $\bar{\epsilon}^* := \frac{1 + \bar{\epsilon}}{2}$. To proceed, noticing that $\bar{\rho}(\mathbf{L}_{1s}) = \lim_{l \rightarrow \infty} \max_{\mathbf{l}_{1:l} \in [s]^l} \left\| \prod_{j=1}^l \mathbf{L}_{j, \mathbf{l}_j} \right\|^{\frac{1}{l}}$ and using the result in (C.57), we can show (C.58). \blacksquare

In Lemma C.1a, if \mathbf{K}_{1s} is computed by solving the infinite-horizon MJS-LQR $(\hat{\mathbf{A}}_{1s}, \hat{\mathbf{B}}_{1s}, \hat{\mathbf{T}}, \mathbf{Q}_{1s}, \mathbf{R}_{1s})$ for some estimated MJS $(\hat{\mathbf{A}}_{1s}, \hat{\mathbf{B}}_{1s}, \hat{\mathbf{T}})$, the following result provides the required estimation accuracy such that the resulting \mathbf{K}_{1s} is uniformly stabilizing.

Lemma C.2a *Under the setup of Lemma C.2, if $\max\{\bar{\epsilon}_{\mathbf{A}, \mathbf{B}}, \bar{\epsilon}_{\mathbf{T}}\} \leq \bar{\epsilon}_{\mathbf{A}, \mathbf{B}, \mathbf{T}}$, then we have $\|\mathbf{K}_{1s} - \mathbf{K}_{1s}\| \leq \bar{\epsilon}_{\mathbf{K}}$, and Lemma C.1a is applicable.*

Recall we defined events $\mathcal{A}_i, \mathcal{B}_i, \mathcal{C}_i, \mathcal{D}_i$ in (C.60) to analyze the events happen in each epoch of the regret. To adapt to the uniform stability assumption, we redefine event \mathcal{B}_i and \mathcal{D}_i while keep \mathcal{A}_i and \mathcal{C}_i as before. For easier reference, We list all of them below.

$$\begin{aligned}
\mathcal{A}_i &= \left\{ \text{Regret}_i \leq \mathcal{O} \left(sp \left(\frac{(i-1)}{\mathbf{A}, \mathbf{B}} + \frac{(i-1)}{\mathbf{T}} \right)^2 \frac{2}{\mathbf{w}} T_i + \sqrt{nS} \|\mathbf{x}_0^{(i)}\|^2 + n\sqrt{S} \frac{2}{\mathbf{z}, i} T_i + c_A \right) \right\} \\
\mathcal{B}_i &= \left\{ \frac{(i)}{\mathbf{A}, \mathbf{B}} \leq \frac{-}{\mathbf{A}, \mathbf{B}, \mathbf{T}}, \frac{(i)}{\mathbf{T}} \leq \frac{-}{\mathbf{A}, \mathbf{B}, \mathbf{T}}, \frac{(i+1)}{\mathbf{K}} \leq \frac{-}{\mathbf{K}} \right\}, \forall i = 0, 1, \dots \\
\mathcal{C}_i &= \left\{ \frac{(i)}{\mathbf{A}, \mathbf{B}} \leq \mathcal{O} \left(\log \left(\frac{1}{id, i} \right) \frac{\mathbf{z}, i + \mathbf{w}}{\mathbf{z}, i} \frac{\sqrt{S}(n+\rho) \log(T_i)}{\sqrt{T_i}} \right), \right. \\
&\quad \left. \frac{(i)}{\mathbf{T}} \leq \mathcal{O} \left(\log \left(\frac{1}{id, i} \right) \frac{1}{\min} \sqrt{\frac{\log(T_i)}{T_i}} \right) \right\} \\
\mathcal{D}_i &= \left\{ \|\mathbf{x}_0^{(i+1)}\|^2 = \|\mathbf{x}_{T_i}^{(i)}\|^2 \leq \frac{18^{-2-2}}{(1-)^2} \log \left(\frac{1}{\mathbf{x}_{0, i}} \right) + 2\bar{X}^{US} \right\}, \forall i = 1, 2, \dots, \\
\mathcal{D}_0 &= \left\{ \|\mathbf{x}_0^{(1)}\|^2 = \|\mathbf{x}_{T_0}^{(0)}\|^2 \leq \frac{n\sqrt{S}^{-2}/(1-)}{\mathbf{x}_{0, 0}} \right\},
\end{aligned} \tag{C.60}$$

where $\bar{X}^{US} := 2^{-2-2} (6 \max\{\sqrt{n}e^{3n}, \sqrt{\rho}e^{3\rho}\} + \frac{5}{(1-)^2})^2$, $\frac{-}{\mathbf{K}} := \min\{\frac{-US}{\mathbf{K}}, \frac{-}{\mathbf{K}}\}$, $\frac{-}{\mathbf{A}, \mathbf{B}, \mathbf{T}} := \min\{\frac{-}{\mathbf{A}, \mathbf{B}, \mathbf{T}}, \frac{\frac{-}{\mathbf{K}}}{2C_{\mathbf{A}, \mathbf{B}, \mathbf{T}}}\}$ and $\|\mathbf{B}_{1s}\|^2 \frac{2}{\mathbf{z}, 0} + \frac{2}{\mathbf{w}}$. Event \mathcal{D}_i describes the initial state magnitude of epoch $i+1$. Since Algorithm 2 requires initial MSS stabilizing controller $\mathbf{K}_{1s}^{(0)}$ for epoch 0, and as in the proof for the following Proposition C.7a, epoch 1, 2, ... have uniformly stabilizing controller, thus we define \mathcal{D}_0 and $\mathcal{D}_1, \mathcal{D}_2, \dots$ separately.

Proposition C.7a Assuming $T_i \geq \frac{1}{2 \log(1-)} \log \left(6^{-2} + \frac{54^{-4-2}}{(1-)^{\bar{X}^{US}}} \log \left(\frac{1}{\mathbf{x}_{0, i-1}} \right) \right)$; $T_1 \geq \frac{1}{2 \log(1-)} \log \left(\frac{3n \bar{S}^{-2-2}}{(1-)^{\bar{X}^{US}} \mathbf{x}_{0, 0}} \right)$, we have

$$P(\mathcal{D}_i | \mathcal{B}_{i-1}, \mathcal{D}_{i-1}) \geq 1 - \mathbf{x}_{0, i} \tag{C.61}$$

and $P(\mathcal{D}_0) \geq 1 - \mathbf{x}_{0, 0}$.

Proof For the initial epoch 0, i.e. $i = 0$, since we assume in Algorithm 2 that the initial controller $\mathbf{K}_{1s}^{(0)}$ stabilizes the MJS in the mean-squared sense, similar to the proof for Proposition C.7, we have $E[\|\mathbf{x}_{T_0}^{(0)}\|^2] \leq n\sqrt{S}(\|\mathbf{B}_{1s}\|^2 \frac{2}{\mathbf{z}, 0} + \frac{2}{\mathbf{w}}) \frac{-}{1-}$. Then by Markov inequality, with probability $1 - \mathbf{x}_{0, 0}$, $\|\mathbf{x}_{T_0}^{(0)}\|^2 \leq \frac{n \bar{S}^{-2} (1-)}{\mathbf{x}_{0, 0}}$ where $\frac{-}{2} := \|\mathbf{B}_{1s}\|^2 \frac{2}{\mathbf{z}, 0} + \frac{2}{\mathbf{w}}$. This shows $P(\mathcal{D}_0) \geq 1 - \mathbf{x}_{0, 0}$.

For epoch $i = 1, 2, \dots$, given event \mathcal{B}_{i-1} , we know $\frac{(i)}{\mathbf{K}} \leq \frac{-}{\mathbf{K}} \leq \frac{-US}{\mathbf{K}}$. Let $\mathbf{L}_{1s}^{(i)}$ denote the closed-loop state matrices for epoch i , then by Lemma C.1a, $\frac{(i)}{\mathbf{K}} \leq \frac{-US}{\mathbf{K}}$ implies that for any l and any sequence $\mathbf{s}_l \in [S]^l$, $\|\prod_{j=1}^l \mathbf{L}_{1s}^{(i)}\| \leq \frac{-}{1-}$. Then using the bound on $\|\mathbf{x}_t\|$ in Lemma A.4a, we have, with probability $1 - \mathbf{x}_{0, i}$,

$$\|\mathbf{x}_{T_i}^{(i)}\|^2 \leq \frac{18^{-2-2}}{(1-)^2} \log \left(\frac{1}{\mathbf{x}_{0, i}} \right) + 3^{-2-2T_i} \|\mathbf{x}_0^{(i)}\|^2 + \bar{X}^{US} \tag{C.62}$$

where $\bar{X}^{US} := 2^{-2-2} (6 \max\{\sqrt{n}e^{3n}, \sqrt{\rho}e^{3\rho}\} + \frac{5}{(1-)^2})^2$.

- When $i = 1$, given \mathcal{D}_0 , i.e. $\|\mathbf{x}_0^{(1)}\|^2 \leq \frac{n \bar{S}^{-2} (1-)}{\mathbf{x}_{0, 0}}$, the above (C.62) gives $\|\mathbf{x}_{T_1}^{(1)}\|^2 \leq \frac{18^{-2-2}}{(1-)^2} \log \left(\frac{1}{\mathbf{x}_{0, 1}} \right) + 3^{-2-2T_1} \frac{n \bar{S}^{-2} (1-)}{\mathbf{x}_{0, 0}} + \bar{X}^{US}$. One can check that when $T_1 \geq \frac{1}{2 \log(1-)} \log \left(\frac{3n \bar{S}^{-2-2}}{(1-)^{\bar{X}^{US}} \mathbf{x}_{0, 0}} \right)$, we have that $3^{-2-2T_1} \frac{n \bar{S}^{-2} (1-)}{\mathbf{x}_{0, 0}} \leq \bar{X}^{US}$, which gives

$$\|\mathbf{x}_{T_1}^{(1)}\|^2 \leq \frac{18^{-2-2}}{(1-)^2} \log \left(\frac{1}{\mathbf{x}_{0, 1}} \right) + 2\bar{X}^{US}. \tag{C.63}$$

- When $i = 2, 3, \dots$, given event \mathcal{D}_{i-1} , i.e. $\|\mathbf{x}_0^{(i)}\|^2 \leq \frac{18^{-2-2}}{(1-\bar{\gamma})^2} \log\left(\frac{1}{\mathbf{x}_{0,i-1}}\right) + 2\bar{X}^{US}$, the above (C.62) gives $\|\mathbf{x}_{T_i}^{(i)}\|^2 \leq \frac{18^{-2-2}}{(1-\bar{\gamma})^2} \log\left(\frac{1}{\mathbf{x}_{0,i}}\right) + 3^{-2-2T_i} \left(\frac{18^{-2-2}}{(1-\bar{\gamma})^2} \log\left(\frac{1}{\mathbf{x}_{0,i-1}}\right) + 2\bar{X}^{US} \right) + \bar{X}^{US}$. Similarly, when the trajectory length during the i th epoch, $T_i \geq \frac{1}{2\log(1-\bar{\gamma})} \log\left(6^{-2} + \frac{54^{-4-2}}{(1-\bar{\gamma})\bar{X}^{US}} \log\left(\frac{1}{\mathbf{x}_{0,i-1}}\right)\right)$, we further have

$$\|\mathbf{x}_{T_i}^{(i)}\|^2 \leq \frac{18^{-2-2}}{(1-\bar{\gamma})^2} \log\left(\frac{1}{\mathbf{x}_{0,i}}\right) + 2\bar{X}^{US}. \quad (\text{C.64})$$

Combining (C.63) and (C.64), we can claim: for epoch $i = 1, 2, \dots$, when $T_1 \geq \frac{1}{2\log(1-\bar{\gamma})} \log\left(\frac{3n\bar{s}^{-2-2}}{(1-\bar{\gamma})\bar{X}^{US}\mathbf{x}_{0,0}}\right)$ and $T_i \geq \frac{1}{2\log(1-\bar{\gamma})} \log\left(6^{-2} + \frac{54^{-4-2}}{(1-\bar{\gamma})\bar{X}^{US}} \log\left(\frac{1}{\mathbf{x}_{0,i-1}}\right)\right)$, we have $\mathbb{P}(\mathcal{D}_i | \mathcal{B}_{i-1}, \mathcal{D}_{i-1}) \geq 1 - \bar{\gamma}_{\mathbf{x}_{0,i}}$. ■

The following Proposition C.8a says that if a good controller is used in epoch i , then the final state $\mathbf{x}_{T_i}^{(i)}$ of epoch i (the initial state of epoch $i+1$) can be bounded.

Proposition C.8a *Suppose every epoch i has length $T_i \geq \underline{T}_{\text{rgt}}^{US}(\bar{\gamma}_{\mathbf{x}_{0,i}}, T_i)$. Then,*

$$\mathbb{P}(\mathcal{B}_i | \mathcal{C}_i, \cap_{j=0}^{i-1} \mathcal{E}_j) = \mathbb{P}(\mathcal{B}_i | \mathcal{C}_i) = 1 \quad (\text{C.65})$$

C.4.1 Proof for Theorem 5.2

Theorem C.12 (Complete version of Theorem 5.2) *Assume that the initial state $\mathbf{x}_0 = 0$, Assumptions 1 and 2 hold, and \mathbf{L}_{1s} is uniformly stable. Suppose $c_{\mathbf{x}} \geq \underline{c}_{\mathbf{x}}(\bar{\gamma}, \bar{\gamma})$, $c_{\mathbf{z}} \geq \underline{c}_{\mathbf{z}}$, $T_0 \geq \mathcal{O}(\underline{T}_{\text{rgt}}^{US}(\bar{\gamma}, T_0))$. Then, with probability at least $1 - \bar{\gamma}$, Algorithm 2 achieves*

$$\text{Regret}(T) \leq \mathcal{O}\left(\frac{s^2 p(n^2 + p^2)}{2 \min} \frac{2}{\mathbf{w}} \log\left(\frac{\log^2(T)}{\log^2(T)}\right) \log^2(T) \sqrt{T}\right). \quad (\text{C.66})$$

Proof The proof is almost the same as the proof for the MSS regret upper bound in Theorem C.11 in Appendix C.3.1, thus we only present the key steps and omit certain details of intermediate steps.

In the following, we set $\bar{\gamma}_{\mathbf{x}_{0,i}} = \bar{\gamma}_{\mathbf{x}_{0,i}} = \frac{3}{2} \cdot \frac{1}{(i+1)^2}$. Similar to the counterpart (C.41), event $\mathcal{E}_i = \mathcal{A}_{i+1} \cap \mathcal{B}_i \cap \mathcal{C}_i \cap \mathcal{D}_i$ implies the following: for $i = 1, 2, \dots$,

$$\begin{aligned} & \text{Regret}_{i+1} \\ & \leq \mathcal{O}(1) \log\left(\frac{(i+1)^2}{\bar{\gamma}_{\mathbf{x}_{0,i}}}\right) s p \left(\frac{\mathbf{z}_{z,i} + \mathbf{w}}{2 \min} \cdot \frac{\sqrt{s}(n+p) \log(T_i)}{\sqrt{T_i}} + \frac{\sqrt{\log(T_i)}}{\min \sqrt{T_i}} \right)^2 \frac{2}{\mathbf{w}} T_{i+1} \\ & \quad + \mathcal{O}(1) \log\left(\frac{i+1}{(1-\bar{\gamma})^2}\right) \frac{18\sqrt{n}s^{-2-2}}{(1-\bar{\gamma})^2} + \mathcal{O}(n\sqrt{s} \frac{2}{\mathbf{z}_{z,i+1}} T_{i+1}) + \mathcal{O}(1) \\ & \leq \mathcal{O}(1) \log\left(\frac{(i+1)^2}{\bar{\gamma}_{\mathbf{x}_{0,i}}}\right) \frac{s^2 p(n^2 + p^2)}{2 \min} \frac{2}{\mathbf{w}} \sqrt{T_i} \log^2(T_i) + \mathcal{O}(1) \log\left(\frac{i+1}{(1-\bar{\gamma})^2}\right) \frac{18\sqrt{n}s^{-2-2}}{(1-\bar{\gamma})^2}; \end{aligned} \quad (\text{C.67})$$

and for $i = 0$,

$$\text{Regret}_1 \leq \mathcal{O}(1) \log\left(\frac{1}{\bar{\gamma}_{\mathbf{x}_{0,0}}}\right) \frac{s^2 p(n^2 + p^2)}{2 \min} \frac{2}{\mathbf{w}} \sqrt{T_0} \log^2(T_0) + \mathcal{O}(1) \left(\frac{1}{1-\bar{\gamma}}\right) \frac{n^{1.5} s^{-2-2}}{1-\bar{\gamma}}. \quad (\text{C.68})$$

Note that the difference between (C.67) ($i = 1, 2, \dots$) and (C.68) ($i = 0$) is due to the difference between the event \mathcal{D}_i for $i = 1, 2, \dots$ and event \mathcal{D}_0 . Compared with the MSS counterpart (C.41), we see the $\frac{(i+1)^2}{\bar{\gamma}_{\mathbf{x}_{0,i}}}$ dependence in (C.41) is now replaced with $\log\left(\frac{i+1}{\bar{\gamma}_{\mathbf{x}_{0,i}}}\right)$. For all $M := \mathcal{O}(\log\left(\frac{T}{T_0}\right))$ epochs, similar to the

counterpart (C.42), event $\cap_{i=0}^{M-1} \mathcal{E}_i$ implies

$$\begin{aligned}
& \text{Regret}(T) \\
& = \mathcal{O}\left(\sum_{i=1}^M \text{Regret}_i\right) \\
& \leq \mathcal{O}\left(\frac{s^2 \rho(n^2 + p^2)}{2 \min} \frac{2}{\mathbf{w}} \log\left(\frac{\log^2(T)}{\log(1-\gamma)}\right) \sqrt{T} \log^2(T) + \frac{18\sqrt{n} \bar{s}^{-2-2}}{(1-\gamma)^2} \log\left(\frac{\log(T)}{\log(1-\gamma)}\right) \log(T)\right)
\end{aligned} \tag{C.69}$$

which shows the main result (C.66). Note that in the above summation, we have omit $\frac{1}{2}$ term in Regret_1 since it does not scale with time and can be dominated by the rest.

Now we are only left to show the occurrence probability of regret bound (C.66) is larger than $1 - \delta$. To do this, we will combine Proposition C.7a, C.8a, C.9, and C.10 over all $i = 0, 1, \dots, M-1$. Note that for each individual i , these propositions hold only when certain prerequisite conditions on hyper-parameters $c_{\mathbf{x}}$, $c_{\mathbf{z}}$, and T_0 are satisfied. We first show that under the choices $T_i = T_{i-1}$, $\frac{2}{\mathbf{z}, i} = \frac{2}{\mathbf{w}_i}$, and $\frac{1}{id, i} = \frac{3}{2} \cdot \frac{1}{(i+1)^2}$ these hyper-parameter conditions can be satisfied for all $i = 0, 1, \dots, M-1$.

- Proposition C.7a requires the following to hold: $T_0 \geq \frac{1}{2 \log(1-\gamma)} \log\left(6^{-2} + \frac{54^{-4-2}}{(1-\gamma)^{2\bar{x}^{us}}} \log\left(\frac{i^2}{3}\right)\right)$ and $T_0 \geq \frac{1}{2 \log(1-\gamma)} \log\left(\frac{2n \bar{s}^{-2-2}}{(1-\gamma)^{2\bar{x}^{us}}}\right)$. One can check $T_0 \geq \max\left\{\frac{54^{-4-2}}{(1-\gamma)^{2\bar{x}^{us}} \log(1-\gamma) \log(\gamma)}, \frac{1}{\log(1-\gamma)} \log(6^{-2} + \frac{54n \bar{s}^{-4-2} \log(\frac{2}{3})}{(1-\gamma)(1-\gamma)^{2\bar{x}^{us}}})\right\} =: \underline{T}_{\mathbf{x}_0}^{US}(\gamma)$ would suffice.
- Proposition C.8a requires that for $i = 0, 1, \dots$, condition $T_0 \geq \underline{T}_{rgt, -}^{US}(\frac{3}{2(i+1)^2}, T_0)$ holds, which can be satisfied when one chooses $T_0 \geq \mathcal{O}(\underline{T}_{rgt, -}^{US}(\frac{3}{2}, T_0))$.
- Proposition C.9 requires the following to hold: $c_{\mathbf{x}} \geq \underline{c}_{\mathbf{x}}(\gamma, \bar{s})$, $c_{\mathbf{z}} \geq \underline{c}_{\mathbf{z}}$, and $T_0 \geq \max\left\{\underline{T}_{MC, 1}(\frac{3}{8 \cdot 2^{i/2}}), \underline{T}_{id, N}(\frac{3}{2 \cdot 2^{i/2}}, \gamma, \bar{s})\right\}$. The last one can be satisfied when we have $T_0 \geq \mathcal{O}(\max\{\underline{T}_{MC, 1}(\gamma), \underline{T}_{id, N}(\gamma, \bar{s})\})$.
- Proposition C.10 requires no conditions on hyper-parameters.

Therefore, when $c_{\mathbf{x}} \geq \underline{c}_{\mathbf{x}}(\gamma, \bar{s})$, $c_{\mathbf{z}} \geq \underline{c}_{\mathbf{z}}$, $T_0 \geq \mathcal{O}(\max\{\underline{T}_{\mathbf{x}_0}^{US}(\gamma), \underline{T}_{rgt, -}^{US}(\frac{3}{2}, T_0), \underline{T}_{MC, 1}(\gamma), \underline{T}_{id, N}(\gamma, \bar{s})\}) =: \mathcal{O}(\underline{T}_{rgt, -}^{US}(\gamma, T_0))$, we can apply Proposition C.7a, C.8a, C.9, and C.10 to every epoch $i = 0, 1, \dots, M-1$. Similar to (C.45), this gives $\mathbb{P}(\text{Regret bounds in (C.66) holds}) \geq \mathbb{P}(\cap_{i=0}^{M-1} \mathcal{E}_i) \geq 1 - \delta$. ■

Remark C.5 Note that though system identification result Theorem 4.1 might also benefit from the newly added uniform stability in this appendix, but since the dependencies \sqrt{T} and $\log(\frac{1}{\delta})$ in Theorem 4.1 are close to the optimal ones for LTI systems, so we only focus on refining the regret upper bound under uniform stability and leave adapting the entire framework to uniform stability to potential future work.