Finite-time Sample Complexity Analysis of Least Square Identifying Stochastic Switched Linear System

Negin Musavi¹ and Geir E. Dullerud^{2†}

Abstract: In this paper, we examine the high-probability finite-time theoretical guarantees of the least squares method for system identification of switched linear systems with process noise and without control input. We consider two scenarios: one in which the switching is i.i.d., and the other in which the switching is according to a Markov process. We provide concentration inequalities using a martingale-type argument to bound the identification error at each mode, and we use concentration lemmas for the switching signal. Our bound is in terms of state dimension, trajectory length, finite-time gramian, and properties of the switching signal distribution. We then provide simulations to demonstrate the accuracy of the identification. Additionally, we show that the empirical convergence rate is consistent with our theoretical bound.

Keywords: Switched linear systems, Identification, Machine learning, Sample complexity

1. INTRODUCTION

The identification of linear systems is a significant task in different fields, such as control theory, time series analysis, robotics, and biology. This task involves creating a model for an unknown system by using input-output data that closely resembles the original system with minimal error. The purpose of this task can be controlling or optimizing the system, detecting faults, making predictions, or gaining knowledge of critical parameters of the unknown system through the identified model. Recently, there has been growing interest in developing algorithms with finite-time sample complexity analysis in this field. A number of works, such as those in [1-4], have addressed this problem. The least square method has gained popularity in linear system identification, and several works have provided non-asymptotic theoretical guarantees for this approach, such as those in [5, 6].

Switched linear systems are a type of hybrid system that feature multiple linear subsystems, with each of them being active for a certain period. These systems are commonly used to model practical systems like control systems, power electronics, and robotics. Identifying these systems becomes more challenging due to the presence of multiple subsystems and the switching between them. However, there has been limited research on understanding the sample complexity analysis for these systems when using the least square method.

In terms of identifying switched linear systems, the approach taken by [7] involves creating a Hankel-like matrix from multiple noisy input-output trajectories and then using ordinary least squares to obtain a good lower-order approximation of the underlying model. In contrast, [8] presents a system identification algorithm for Markov jump linear systems using a single trajectory of the system states, inputs, and modes. They also analyze non-asymptotic probabilistic sample complexity using

mixing-time arguments. In another work by [9], asymptotic but almost sure guarantees are provided for autonomous Markov jump linear systems using non-mixing arguments via martingales.

This paper presents a non-asymptotic analysis of the sample complexity of the least-squares estimator for identifying switched linear systems. We use a single trajectory of states and modes, of length T, which is excited by noise. We employ a martingale-type theory closely following the approach used in [5]. Specifically, we consider two scenarios: one in which the switching is i.i.d., and the other in which the switching is according to a Markov process. Our aim is to provide non-asymptotic bounds with high probability that take into account the trajectory length, system excitation, state dimension, and properties of the switching signal distribution. We carefully analyze the least square model with correlated data and use a concentration bound on the Markov chains to derive our finite-time bound. Finally, we present simulations to evaluate the performance of the least square estimator for identifying switched linear systems.

We start by stating the problem in Section 2, followed by the theoretical analysis in Section 3 that includes the main results of the paper. We then provide simulations in Section 4. The proof of the main results is provided in the Appendix.

2. PROBLEM STATEMENT

2.1. Notation

The set of real numbers is denoted by \mathbb{R} . For a real matrix Z, Z^T represents its transpose, |Z| denotes its maximum singular value, and $\lambda_{\min}(Z)$ represents its minimum eigenvalue. For a real symmetric matrix Z, $Z\succ 0$ and $Z\succeq 0$ indicate that Z is positive definite and positive semi-definite, respectively. Let $\ell_2(\pi)$ denote the set of functions f mapping a finite set G to \mathbb{R} that satisfy $\sum_{x\in G} f(x)^2\pi(x) < \infty$, where $\pi(x)>0$

¹Department of Mechanical Science and Engineering, University of Illinois at Urbana-Champaign, Urbana, USA nmusavi2@illinois.edu

²Department of Mechanical Science and Engineering, University of Illinois at Urbana-Champaign, Urbana, USA dullerud@illinois.edu

[†] Geir E. Dullerud is the presenter of this paper.

for all $x \in G$. The $\ell_2(\pi)$ norm of f is defined as $\|f\|_{2,\pi}^2 = \sum_{x \in G} f(x)^2 \pi(x)$. The inner product for functions f and g on $\ell_2(\pi)$ is then defined as $\langle f,g \rangle_\pi = \sum_{x \in G} f(x)g(x)\pi(x)$. For an operator $P: G \times G \to \mathbb{R}$, P denotes its conjugate on $\ell_2(\pi)$, satisfying $\langle Pf,g \rangle_\pi = \langle f,Pg \rangle_\pi$ for $f,g \in \ell_2(\pi)$. Furthermore, let $\ell_\infty(\pi)$ denote the set of functions f mapping a finite set G to \mathbb{R} , satisfying $\sup_{x \in G} |f(x)|\pi(x) < \infty$, where $\pi(x) > 0$ for all $x \in G$. The $\ell_\infty(\pi)$ norm of f is denoted as $\|f\|_{\infty,\pi} = \sup_{x \in G} |f(x)|\pi(x)$.

2.2. Problem Statement

A switched linear system is a type of dynamical system that combines multiple linear subsystems, each with its own state and dynamics. The system switches between these subsystems according to a switching signal, which can be based on time, sensor measurements, or other criteria. In this paper, we consider a discrete-time switched linear system described as:

$$x_{t+1} = A_{*\phi_t} x_t + \eta_t, (1)$$

where $x_t \in \mathbb{R}^n$ represents the state of the system at time t, $\eta_t \sim \mathcal{N}(0, \sigma^2 I)$ denotes i.i.d. noise at time t with $\sigma > 0$, and $A_{*\phi_t} \in \mathbb{R}^{n \times n}$ is a matrix that defines the dynamics of the system in the current mode determined by the switching signal ϕ_t . We assume that the switching signal is a function of time and is independent of the system's states. We use $\mathcal{M} = \{1, 2, \cdots, L\}$ to denote the set of modes, where L > 0 (i.e., $\phi_t \in \mathcal{M}$).

Given a sequence of states $\{x_0, x_1, x_2, ..., x_{T+1}\}$ with $x_0 = 0$ and a sequence of modes $\{\phi_0, \phi_1, ..., \phi_T\}$, the least square (LS) estimator is the solution to the following optimization problem:

$$\hat{A}_{1}(T), \dots, \hat{A}_{L}(T) = \underset{A_{1}, \dots, A_{L} \in \mathbb{R}^{n \times n}}{\operatorname{argmin}} \frac{1}{2} \sum_{t=1}^{T} \|x_{t+1} - A_{\phi_{t}} x_{t}\|_{2}^{2}.$$
(2)

Our objective is to provide high-probability bounds on the performance of the LS estimator for each mode within a finite time frame. For a given mode $l \in \mathcal{M}$, trajectory length T, and a given $\delta \in (0,1)$, we seek to determine whether there exists an $\epsilon \in (0,1)$ satisfying:

$$\frac{|\hat{A}_l(T) - A_{*l}|}{|A_{*1}|} \le \epsilon \text{ with probability at least } 1 - \delta? \quad (3)$$

If such an ϵ exists, we aim to explore its dependence on the trajectory length T, system parameters, and δ . It is worth noting that we can rewrite the LS estimator problem in (2) as:

$$\hat{A}_{1}(T), \dots, \hat{A}_{L}(T)$$

$$= \underset{A_{1}, \dots, A_{L} \in \mathbb{R}^{n \times n}}{\operatorname{argmin}} \frac{1}{2} \sum_{t=1}^{T} \sum_{l=1}^{L} \|x_{t+1} - A_{\phi_{t}} x_{t}\|_{2}^{2} \mathbb{1} \{ \phi_{t} = l \}.$$

Assuming that we can observe the modes at each time t, the LS estimator $\hat{A}_l(T)$ for each mode $\ell \in \mathcal{M}$ is the

solution to the following optimization problem:

$$\hat{A}_{l}(T) = \operatorname*{argmin}_{A_{1} \in \mathbb{R}^{n \times n}} \frac{1}{2} \sum_{t=1}^{T} \|x_{t+1} - A_{\phi_{t}} x_{t}\|_{2}^{2} \mathbb{1} \{\phi_{t} = l\}.$$

3. ANALYSIS

In this section, we provide a theoretical assessment of the LS estimator's performance in addressing the problem described in (2). To do so, we start by introducing definitions, assumptions, and lemmas that are essential for conducting the analysis outlined in (3).

For the linear system in (1), the filtration created by the states, switching signal, and noise process at time t is defined as $\mathcal{F}_t := \{\eta_0,...,\eta_{t-1},\phi_0,...,\phi_{t-1},x_1,...,x_t\}$. A random process $(z_t)_{t\geq 1}$ is said to be $\{\mathcal{F}_t\}_{t\geq 0}$ -adapted if $z_t\in\mathcal{F}_t$ and $z_t\notin\mathcal{F}_{t-1}$. We define a sub-Gaussian noise with respect to this filtration as follows:

Definition 1: We say that $\eta_t | \mathcal{F}_t$ is a mean-zero and σ^2 -sub-Gaussian noise if $\mathbb{E}[\eta_t | \mathcal{F}_t] = 0$, and it satisfies the tail condition $\mathbb{E}[e^{\lambda \eta_t} | \mathcal{F}_t] \leq e^{\sigma^2 \lambda^2/2}$ for some $\lambda > 0$.

We define the stability of the linear system in (1) in the mean square sense as follows:

Definition 2: The linear system in (1) is said to be stable in the mean square sense if, for every initial condition x_0 and ϕ_0 , $\mathbb{E}[x_tx_t^T] \to 0$ as $t \to \infty$. It is said to be marginally stable in the mean square sense if, for every initial condition x_0 and ϕ_0 , $\mathbb{E}[x_tx_t^T] < \infty$ as $t \to \infty$.

Let the finite-time controllability gramian of the system in (1) be defined as

$$\bar{\Gamma}_t = \mathbb{E}\left[\sum_{l=0}^{t-2} \prod_{s=t-1}^{l+1} A_{\phi_s} A'_{\phi_s}\right] + I$$

for t>1 with $\bar{\Gamma}_1=I$. If the system is marginally stable in the mean square sense, then $\bar{\Gamma}_t$ is finite for any $t\geq 1$. To ensure a finite $\bar{\Gamma}_t$ for our subsequent analysis, we assume that our system is marginally stable in the mean square sense throughout the paper. Our approach is similar to that described in [5], which relies on the linear system being excitable by noise. To formalize this concept in terms of the states x_t along a trajectory of length T, the following definition is helpful.

Definition 3: Given an $\{\mathcal{F}_t\}_{t\geq 1}$ -adapted random process $(x_t)_{t\geq 1}$ in \mathbb{R}^n , we say it satisfies the $(k,\underline{\Gamma},p)$ -block martingale small ball condition for $\underline{\Gamma}\succ 0$ with some $k\geq 1,\underline{\Gamma}\succ 0$, and $p\in (0,1)$ if, for any $w\in \mathbb{R}^n$ with |w|=1 and for any $j\geq 0$, it satisfies:

$$\frac{1}{k} \sum_{i=1}^{k} \mathbb{P} \left[w^T x_{j+i} x_{j+i}^T w \ge w^T \underline{\Gamma} w | \mathcal{F}_j \right] \ge p.$$

It is worth noting that as the eigenvalues of $\underline{\Gamma}$ increase, the system becomes more susceptible to noise. The next lemma provides an estimate of the statistical performance of the LS method for identifying the system in (1) with a single mode. This result is borrowed from [5] (Theorem 2.4).

Lemma 1: For a given $\delta \in (0,1)$ and T > 0, suppose that the pair $(x_{t+1}, x_t)_{t \ge 1}$ is a random sequence satisfying the following conditions:

- It is generated by marginally stable linear system $x_{t+1} = A_* x_t + \eta_t$ with $\eta_t | \mathcal{F}_t$ being mean-zero and σ^2 -sub-Gaussian noise.
- The process $(x_t)_{t\geq 1}$ satisfies the $(k,\underline{\Gamma},p)$ -block martingale small ball condition for some $p\in(0,1)$ and $\Gamma\succ 0$.
- It holds that $\mathbb{P}[\sum_{t=1}^T x_t x_t^T \npreceq T\bar{\Gamma}] \leq \delta$ for some $\bar{\Gamma} \succeq \Gamma$

Then, if the following condition holds:

$$\frac{T}{k} \geq \frac{10}{p^2} \bigg(2n \log \frac{10}{p} + \log \frac{1}{\delta} + \log \det \left(\bar{\Gamma} \underline{\Gamma}^{-1} \right) \bigg),$$

the LS estimator in (2) achieves the following performance:

$$\mathbb{P}\left[|\hat{A}(T) - A_*|\right] > \frac{90\sigma}{p} \sqrt{\frac{n + n\log\frac{10}{p} + \log\frac{1}{\delta} + \log\det\left(\bar{\Gamma}\underline{\Gamma}^{-1}\right)}{T\lambda_{\min}\left(\underline{\Gamma}\right)}} \leq \delta.$$

The proof of this result differs from the standard LS analysis because the states $\{x_1, x_2, \cdots, x_T\}$ are correlated. The central assumption of this meta-theorem is that every length k block of states (i.e., $(x_{j+i})_{i=1}^k$ for any j) satisfies the lower bound in the block martingale small ball condition. This assumption ensures that the linear system is excitable by noise. In the probabilistic bound, this assumption appears as the term $\log \det \left(\bar{\Gamma} \underline{\bar{\Gamma}}^{-1} \right)$ and suggests that the more excitable the system, the faster the learning rate. Moreover, as T increases, the probabilistic bound improves.

To apply this result to a multi-mode case, we need additional definitions. We define a length T execution of switching signals as $\alpha_T = \{\phi_0, \phi_1, \cdots, \phi_{T-1}\}$. Within a time frame of length T, there are 2^T possible executions, denoted by α_T^i , where i ranges from 1 to 2^T . We introduce a random variable, S_T , which represents a sequence of modes within a time frame of length T. By considering each possible execution α_T^i , we can modify Theorem 2 and obtain the following results:

Lemma 2: Consider the dynamical system in (1) that is marginally stable in the mean square sense, and let $x_0=0$. Fix T>0 and $\delta\in(0,1/2)$, and let $\{x_1,x_2,\cdots,x_T\}$ be a random sequence of states generated by the dynamical system. Let α_T^i be the execution of modes that generate this sequence of states. Let T_l^i denote the number of times that mode l is active along this trajectory, i.e., $T_l^i=\sum_{t=1}^T\mathbb{1}\{\phi_t=l|S_T=\alpha_T^i\}$. Then, for any k_l satisfying

$$\frac{T_l^i}{k_l} > c_n \left(n \log \frac{n}{\delta} + \log \det \left(\bar{\Gamma}_T \bar{\Gamma}_{k_l}^{-1} \right) \right), \tag{4}$$

we have:

$$\mathbb{P}\left[|\hat{A}_l(T) - A_{*l}| > \epsilon_l^i(\delta, T, T_l^i, k_l) \middle| S_T = \alpha_T^i\right] \le \delta,$$

where

$$\epsilon_l^i(\delta, T, T_l^i, k_l) = C_n \sqrt{\frac{n \log \frac{n}{\delta} + \log \det \left(\bar{\Gamma}_T \bar{\Gamma}_{k_l}^{-1}\right)}{T_l^i \lambda_{\min} \left(\bar{\Gamma}_{k_l}\right)}}$$

for some $c_n, C_n > 0$.

This lemma is a direct consequence of the previous lemma, and its proof has been moved to the Appendix for readability. It should be noted that the bound described above applies to any k_l that satisfies (4). Therefore, it is also valid for $\bar{\epsilon}_l^i(\delta,T,T_l^i)=\inf_k \epsilon_l^i(\delta,T,T_l^i,k)$. Having established the required definitions and lemmas, we can now present the following theorem.

Theorem 1: Consider the dynamical system in (1) that is marginally stable in the mean square sense and $x_0=0$. Fix T>0 and $\delta\in(0,1/2)$. Let $\bar{T}\leq T$ and T_l^i be the number of times that mode l is active along the execution α_T^i , i.e. $T_l^i=\sum_{t=1}^T\mathbb{1}\{\phi_t=l|S_T=\alpha_T^i\}$, where $1\leq i\leq 2^T$. Let $\bar{\epsilon}_l(\delta,T,\bar{T})=\max_{T_l^i\geq \bar{T}}\bar{\epsilon}_l^i(\delta,T,T_l^i)$, then we have:

$$\mathbb{P}\left[|\hat{A}_l(T) - A_{*l}| > \bar{\epsilon}_l(\delta, T, \bar{T})\right] \le \delta + (1 - \delta)\bar{p}_l(\bar{T})$$

with
$$ar{p}_l(ar{T}) = \sum_{\substack{i = 1 \ T_1^i < ar{T}}}^{2^T} \mathbb{P}\big[S_T = lpha_T^i\big].$$

This result can be directly derived using the Law of Total Probability. For a detailed explanation of the proof, please see the Appendix. It should be noted that this theorem is valid for any $\bar{T} \leq T$, and the optimal \bar{T} is the one that minimizes both $\bar{\epsilon}_l(\delta,T,\bar{T})$ and $\bar{p}_l(\bar{T})$ simultaneously. So far, we have not made any assumptions regarding the distribution of the switching signal, other than it being independent of the system's state. However, the distribution of the switching signal affects $\bar{p}_l(\bar{T})$. We will examine two general cases for the switching signal: one with i.i.d. switching, and the other in which the switching is according to a Markov process. These two cases are studied in more detail.

3.1. i.i.d. case

In this situation, where the probability of each mode l being active at each time is p_l , we can represent the event of mode l being active as a Bernoulli random variable with parameter p_l . As a result, the following relationship holds:

$$\bar{p}_l(\bar{T}) = \sum_{\substack{i=1 : \\ T_l^i < \bar{T}}}^{2^T} \mathbb{P}\big[S_T = \alpha_T^i\big] = \sum_{k=1}^{\bar{T}} \binom{T}{k} \, p_l^k (1 - p_l)^{T-k}.$$

Lemma 3: Suppose at each time t, each mode $l \in \mathcal{M}$ is active with probability p_l , where $\sum_{l \in \mathcal{M}} p_l = 1$. Fix T > 0. Fix $\delta \in (0, 1/2)$ and let $\bar{T} < p_l T$, then we have:

$$\mathbb{P}\left[|\hat{A}_{l}(T) - A_{*l}| > \bar{\epsilon}_{l}(\delta, T, \bar{T})\right]$$

$$\leq \delta + (1 - \delta) \sum_{l=1}^{\bar{T}} {T \choose k} p_{l}^{k} (1 - p_{l})^{T-k}.$$

This is a direct result of the previous theorem. It is important to note that this applies to any $\bar{T} < p_l T$. For a larger \bar{T} , $\bar{\epsilon}_l(\delta,T,\bar{T})$ is smaller, which is desired, and $\bar{p}_l(\bar{T})$ is larger, which is not desired. Therefore, the optimal value of \bar{T} is the one that minimizes both $\bar{\epsilon}_l(\delta,T,\bar{T})$ and $\bar{p}_l(\bar{T})$. Notice for $\bar{T} < p_l T$ we can also use Hoeffding's inequality for sum of bounded independent random variables to get

$$\sum_{k=1}^{\bar{T}} {T \choose k} p_l^k (1-p_l)^{T-k} \le \exp\left(-T\left(p_l - \frac{\bar{T}}{T}\right)^2\right).$$

This implies that $\bar{p}_l(\bar{T})$ is bounded if $\bar{T} < p_l T$ and it tends to 0 as $T \to \infty$.

3.2. DTMC case

In this situation, where the modes are the states of a finite, aperiodic, and irreducible DTMC, we can utilize a concentration lemma for finite Markov chains, which is derived from Theorem 3.3 in [10].

Lemma 4: Let's consider an aperiodic and irreducible discrete-time Markov chain (DTMC) on a finite set \mathcal{M} , with a probability transition matrix P and a stationary distribution π . Suppose the multiplicative symmetrization $K=P^*P$ is also irreducible. We have a function $f:\mathcal{M}\to\mathbb{R}$ such that $\pi f=0, \|f\|_{\pi,\infty}<1$, and $0<\|f\|_{\pi,2}^2\leq b^2$ for some b>0. Now, for an initial distribution q, a positive integer T, and $0<\gamma\leq 1$, we have the following inequality:

$$\mathbb{P}\left[\frac{1}{T}\sum_{t=1}^{T} f(s_t) \ge \gamma\right] \le N_q \exp\left(\frac{-T\gamma^2 \left(1 - \lambda_1(K)\right)}{8b^2 \left(1 + h\left(\frac{5\gamma}{b^2}\right)\right)}\right),$$

where $\{s_1,\cdots,s_T\}$ are sequence of states generated by this DTMC up to time $T,\ N_q=\|\frac{q}{\pi}\|_{\pi,2}$, $\lambda_1(P^*P)$ is the second largest eigenvalue of P^*P and $h(s)=\frac{1}{2}\left(\sqrt{1+s}-(1-\frac{s}{2})\right)$.

This concentration lemma is distinct from the analysis of sums of independent random variables, as the states of the Markov chain are correlated. However, the rate of convergence is still exponentially decreasing as T grows, with the additional dependence of the spectral gap of the probability transition matrix. A larger spectral gap results in a faster rate of convergence. With this concentration lemma, we can state the following lemma:

Lemma 5: Suppose the switching signal is according to a aperiodic and irreducible DTMC with probability transition matrix P and stationary distribution $\pi = (\pi_1, \cdots, \pi_L)$. Suppose P^*P is also irreducible. Fix T>0 and $\delta\in(0,1/2)$. For each mode $l\in\mathcal{M}$ if $\bar{T}<\pi_l T$, the following holds with $\gamma(\bar{T})=\pi_l-\frac{\bar{T}}{T}$:

$$\begin{split} &\mathbb{P}\bigg[|\hat{A}_l(T) - A_{*l}| > \bar{\epsilon}_l(\delta, T, \bar{T})\bigg] \\ &\leq \delta + (1-\delta)N_0 \exp\bigg(\frac{-T\gamma(\bar{T})^2\big(1-\lambda_1(P^*P)\big)}{8(1-\pi_l)^2\big(1+h\big(\frac{5\gamma(\bar{T})}{(1-\pi_l)^2}\big)\big)}\bigg), \end{split}$$
 To keep the readability of the main text, the proof is

To keep the readability of the main text, the proof is provided in the Appendix.

3.3. Asymptotic Bound

Here we aim to derive the asymptotic bound for the DTMC case. Since $\bar{T} < \pi_l T$, it also holds that $\bar{T} = (\pi_l - \theta)T$ for any $\theta \in (0, \pi_l)$. Thus we have:

$$\lim_{T \to \infty} \bar{\epsilon}_l(\delta, T, \bar{T})$$

$$= \lim_{T \to \infty} \inf_k C_n \sqrt{\frac{n \log \frac{n}{\delta} + \log \det \left(\bar{\Gamma}_T \bar{\Gamma}_k^{-1}\right)}{\bar{T} \lambda_{\min}(\bar{\Gamma}_k)}}$$

$$= \lim_{T \to \infty} \inf_k C_n \sqrt{\frac{n \log \frac{n}{\delta} + \log \det \left(\bar{\Gamma}_T \bar{\Gamma}_k^{-1}\right)}{(\pi_l - \theta) T \lambda_{\min}(\bar{\Gamma}_k)}}$$

$$= \lim_{T \to \infty} \sqrt{\frac{n \log \frac{n}{\delta}}{(\pi_l - \theta) T}}.$$

Moreover,

$$\lim_{T\to\infty} N_q \exp\left(\frac{-T\gamma(\bar{T})^2\left(1-\lambda_1(P^*P)\right)}{8(1-\pi_l)^2\left(1+h\left(\frac{5\gamma(\bar{T})}{(1-\pi_l)^2}\right)\right)}\right) = 0.$$

With these considerations, we can conclude that

$$\lim_{T \to \infty} |\hat{A}_l(T) - A_{*l}| \le \lim_{T \to \infty} \sqrt{\frac{n \log \frac{n}{\delta}}{(\pi_l - \theta)T}},$$

with a probability of at least $1-\delta$. This implies that the asymptotic rate of convergence is $\mathcal{O}\left(\sqrt{\frac{\log \frac{1}{\delta}}{T}}\right)$ with a probability of at least $1-\delta$. Similar asymptotic rates can be concluded for the i.i.d. case if $\bar{T} < p_l T$.

4. SIMULATION RESULTS

This section presents simulation outcomes that demonstrate the practical effectiveness of LS for a switched linear system in (1) consisting of two modes. The nominal system matrices associated with each mode

are as
$$A_{*1} = T_1 \begin{pmatrix} \Lambda_1 & 0 & 0 \\ 0 & \Lambda_2 & 0 \\ 0 & 0 & \Lambda_3 \end{pmatrix} T_1^{-1} \text{ and } A_{*2} = T_2 \begin{pmatrix} \Lambda_4 & 0 \\ 0 & \Lambda_5 \end{pmatrix} T_2^{-1}$$
, with $\Lambda_1 = \begin{pmatrix} 0.198 & 0.969 \\ -0.969 & 0.198 \end{pmatrix}$, $\Lambda_2 = \begin{pmatrix} 0.4 & 0.693 \\ -0.693 & 0.4 \end{pmatrix}$, $\Lambda_3 = 0.01I$, $\Lambda_4 = \begin{pmatrix} 0.49 & 0.499 \\ -0.499 & 0.49 \end{pmatrix}$, $\Lambda_5 = 0.1I$, and non-singular T_1 and T_2 . The spectral radius of the system matrix linked to the first mode is higher than that of the second mode. This means that the first mode is more susceptible to noise and, with the same activation rate, is anticipated to be identi-

The upcoming Figure 1 illustrates the performance of LS in a scenario where the switching signal is an i.i.d. random variable with probability p. The plot displays the outcomes for p=0.1,0.3,0.5 and each graph corresponds to the average of 20 distinct trajectories with a duration of T. Upon analyzing the individual plots for each mode, it is apparent that the LS estimation error improves

fied more effortlessly.

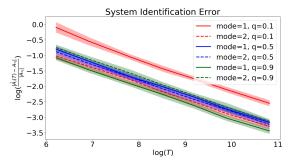


Fig. 1 Empirical performance of LS: the switching signal is a binomial random process with parameter p.

as the probability of activation increases. When considering the plots jointly for both modes, it is observed that when the probability of activation is the same for both modes (p = 0.5), the estimation error for the first mode is smaller than that of the second mode. On the extreme end, when the probability of activation for the first mode is 0.1 and for the second mode is 0.9, the performance of LS is better for the second mode, although identifying it is a more challenging task. Another scenario arises when the probability of activation for the first mode is 0.3, which is still less than that of the second mode, which is 0.7, but the LS performance for both modes is comparable. In addition, notice that by inspecting it more closely we see that $\log \frac{|\hat{A}_l(T) - A_{*1}|}{|A_{*l}|}$ can be approximated with $-\frac{1}{2} \log T$. This indicate that the rate of convergence empirically is approximately $\frac{1}{\sqrt{T}}$. This is consistent with our asymptotic convergence rate.

Similar results can be seen in the upcoming Figure 2 which illustrates the performance of LS in a scenario where the switching is according to a Markov process with probability transition matrix and stationary distribution as $P = \begin{pmatrix} 0.1 & 0.9 \\ 0.5 & 0.5 \end{pmatrix}, \ \pi = \begin{pmatrix} 0.36, & 0.64 \end{pmatrix}$, respectively. The plot displays the outcomes for different

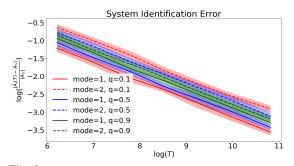


Fig. 2 Empirical performance of LS: the switching is according to a Markov process with initial distribution (q, 1-q).

q=0.1,0.5,0.9 where q is the parameter in initial distribution $\left(q,-1-q\right)$. Each graph corresponds to the average of 20 distinct trajectories with a duration of T. Similar to the i.i.d. case the plot reveals that $\log \frac{|\hat{A}_l(T) - A_{*1}|}{|A_{*l}|}$ can be approximated with $-\frac{1}{2}\log T$ which again implies the empirical rate of convergence approximately as $\frac{1}{\sqrt{T}}$ which is consistent with our theoretical bound.

5. ACKNOWLEDGEMENT

The authors would like to thank Sanjay Shakkottai, Dawei Sun, and Sayan Mitra for the discussions that helped lead to this work.

REFERENCES

- [1] Y. Jedra and A. Proutiere, "Finite-time identification of stable linear systems optimality of the least-squares estimator", In *IEEE 59th Conference on Decision and Control (CDC)*, pp. 996–1001, IEEE, 2020.
- [2] A. Tsiamis and G. J. Pappas, "Finite sample analysis of stochastic system identification", In *IEEE* 58th Conference on Decision and Control (CDC), pp. 3648–3654, IEEE, 2019.
- [3] S. Oymak and N. Ozay, "Non-asymptotic identification of lti systems from a single trajectory", In *IEEE American control conference (ACC)*, pp. 5655–5661, IEEE, 2019.
- [4] T. Sarkar and A. Rakhlin, "Near optimal finite time identification of arbitrary linear dynamical systems", In *International Conference on Machine Learning*, pp. 5610–5618, PMLR, 2019.
- [5] M. Simchowitz and H. Mania, "Learning without mixing: Towards a sharp analysis of linear system identification", In *Conference on Learning Theory*, pp. 439–473, PMLR, 2018.
- [6] M. K. S. Faradonbeh, A. Tewari, and G. Michailidis, "Finite time identification in unstable linear systems", *Automatica*, Vol. 96, pp. 342–353, Elsevier, 2018.
- [7] T. Sarkar, A. Rakhlin and M. A. Dahleh, "Data driven estimation of stochastic switched linear systems of unknown order", *arXiv preprint arXiv:1909.04617*, 2019.
- [8] Y. Sattar, Z. Du, D. A. Tarzanagh, L. Balzano, N. Ozay, and S. Oymak, "Identification and adaptive control of markov jump systems: Sample complexity and regret bounds", arXiv preprint arXiv:2111.07018, 2021.
- [9] B. Sayedana, M. Afshari, P. E. Caines, and A. Mahajan, "Consistency and rate of convergence of switched least squares system identification for autonomous Markov jump linear systems", In *IEEE 61st Conference on Decision and Control (CDC)*, pp. 6678–6685, IEEE, 2022.
- [10] P. Lezaud, "Chernoff-type bound for finite Markov chains", *Annals of Applied Probability*, pp. 849–867, 1998.

APPENDIX

Proof: (proof of Lemma 2) To prove the lemma, it is sufficient to demonstrate the existence of $0 \prec \underline{\Gamma} \preceq \bar{\Gamma}$ satisfying the following conditions for any $l \in \mathcal{M}$:

• The $\{\mathcal{F}_t\}_{t\geq 1}$ -adapted process $(x_t)_{t\geq 1:\phi_t=l}$ for a given execution α_T^i satisfies the $(k,\underline{\Gamma},p)$ block martingale small ball condition.

- It holds that $\mathbb{P}\left[\sum_{t=1}^{T} x_t x_t^T \mathbb{1}\{\phi_t = l | \mathcal{S}_T = \alpha_T^i\} \not\preceq T\bar{\Gamma}\right] \leq \delta$.
- \bullet For any $w \in \mathbb{R}^n$ with |w|=1, it holds that $\mathbb{P}\big[\sum_{t=1}^T w^T x_t x_t^T w \mathbb{1}\{\phi_t=l|\mathcal{S}_T=\alpha_T^i\} \leq \frac{w^T \underline{\Gamma} w p^2}{8} k \lfloor \frac{T_l^i}{k} \rfloor \big] \leq e^{-\frac{\lfloor \frac{T_l^i}{k} \rfloor p^2}{8}}$, and \bullet A martingale-type argument, similar to Lemma
- A martingale-type argument, similar to Lemma 4.2 in [5], holds for the $\{\mathcal{F}_t\}_{t\geq 1}$ -adapted process $(x_t)_{t\geq 1: \phi_t=l}$ for a given execution α_T^i .

Due to limitations of space, we leave the statements mentioned above without proof. However, it can be observed that all of these statements remain valid with the following choices of parameters: $\bar{\Gamma} = \frac{\sigma^2 d}{\delta} \bar{\Gamma}_T$, $p = \frac{3}{20}$, and $\underline{\Gamma} = \sigma^2 \bar{\Gamma}_r$. Here, r represents the index in the execution α_T^i where mode l is active for the $\lfloor \frac{k}{2} \rfloor$ -th time, meaning that $\sum_{t=1}^{r+1} \mathbb{1}\{\phi_t = l | \mathcal{S}_T = \alpha_T^i\} = \lfloor \frac{k}{2} \rfloor$.

Proof: (proof of Theorem 1) Using the law of total probability, we can derive an expression for the probability of the identification error of A_{*l} being greater than $\epsilon>0$ as:

$$\begin{split} & \mathbb{P}\left[|\hat{A}_{l}(T) - A_{*l}| > \epsilon\right] \\ & = \mathbb{P}\left[\left. \bigcup_{i=1}^{2^{T}} \left\{\left\{|\hat{A}_{l}(T) - A_{*l}| > \epsilon\right\} \cap \left\{S_{T} = \alpha_{T}^{i}\right\}\right\}\right] \\ & = \sum_{i=1}^{2^{T}} \mathbb{P}\left[\left\{\left\{|\hat{A}_{l}(T) - A_{*l}| > \epsilon\right\} \cap \left\{S_{T} = \alpha_{T}^{i}\right\}\right\}\right] \\ & = \sum_{i=1}^{2^{T}} \mathbb{P}\left[\left|\hat{A}_{l}(T) - A_{*l}| > \epsilon\right| S_{T} = \alpha_{T}^{i}\right] \mathbb{P}\left[S_{T} = \alpha_{T}^{i}\right]. \end{split}$$

This expression tells us that the probability of the identification error depends on both the error probability for each execution and the probability of each execution occurring. However, it's important to note that not all executions have the same probability of occurring, and some executions that are less likely to occur might contribute less in the identification error. To account for this, we split the sum into two parts:

$$\begin{split} & \mathbb{P}\bigg[|\hat{A}_l(T) - A_{*l}| > \epsilon \bigg] \\ & = \sum_{\substack{i \, = \, 1 \, : \\ T_l^i \, < \, \bar{T}}}^{2^T} \mathbb{P}\bigg[|\hat{A}_l(T) - A_{*l}| > \epsilon \big| S_T = \alpha_T^i \bigg] \mathbb{P}\big[S_T = \alpha_T^i \big] \\ & + \sum_{\substack{i \, = \, 1 \, : \\ T_l^i \, > \, \bar{T}}}^{2^T} \mathbb{P}\bigg[|\hat{A}_l(T) - A_{*l}| > \epsilon \big| S_T = \alpha_T^i \bigg] \mathbb{P}\big[S_T = \alpha_T^i \big]. \end{split}$$

In the first sum, where $T_l^i < \bar{T}$, we consider executions where the subsystem with A_{*l} appears less frequently. Since these trajectories are less likely to contribute to a better estimation error, we set an upper bound of 1 on their probability.

For any $\delta \in (0, 1/2)$, we define $\epsilon = \bar{\epsilon}_l(\delta, T, \bar{T})$, which represents the maximum value of $\bar{\epsilon}_l^i(\delta, T, T_l^i)$ among all

 $T_l^i \geq \bar{T}$. It is worth noting that in this case, $\bar{\epsilon}_l^i(\delta, T, T_l^i) \leq \bar{\epsilon}_l(\delta, T, \bar{T})$ for all executions α_T^i . Based on this, can further simplify the terms in the second sum as:

$$\mathbb{P}\left[|\hat{A}_{l}(T) - A_{*l}| > \epsilon \middle| S_{T} = \alpha_{T}^{i}\right]$$

$$\leq \mathbb{P}\left[|\hat{A}_{l}(T) - A_{*l}| > \bar{\epsilon}_{l}^{i}(\delta, T, T_{l}^{i})\middle| S_{T} = \alpha_{T}^{i}\right]$$

$$< \delta.$$

This allows us to conclude that:

$$\mathbb{P}\left[|\hat{A}_{l}(T) - A_{*l}| > \bar{\epsilon}_{l}(\delta, T, \bar{T})\right]$$

$$\leq \sum_{\substack{i=1 : \\ T_{l}^{i} < \bar{T}}}^{2^{T}} \mathbb{P}\left[S_{T} = \alpha_{T}^{i}\right] + \delta \sum_{\substack{i=1 : \\ T_{l}^{i} \geq \bar{T}}}^{2^{T}} \mathbb{P}\left[S_{T} = \alpha_{T}^{i}\right]$$

$$= \bar{p}_{l}(\bar{T}) + \delta(1 - \bar{p}_{l}(\bar{T})),$$

where $\bar{p}_l(\bar{T})$ denotes the probability of the subsystem with A_{*l} appearing less frequently. This concludes the proof.

Proof: (proof of Lemma 5) We first need to determine an upper bound for $\bar{p}_l(\bar{T})$. We define \mathcal{M}_{-l} as the set of all modes except mode l. We can express $\bar{p}_l(\bar{T})$ as follows:

$$\bar{p}_{l}(\bar{T}) = \sum_{\substack{i=1 : \\ T_{l}^{i} < \bar{T}}}^{2^{T}} \mathbb{P}\left[S_{T} = \alpha_{T}^{i}\right] = \mathbb{P}\left[\sum_{t=1}^{T} \mathbb{1}\{\phi_{t} = l\} < \bar{T}\right]$$
$$= \mathbb{P}\left[\frac{1}{T}\sum_{t=1}^{T} \mathbb{1}\{\phi_{t} \in \mathcal{M}_{-l}\} \ge \frac{T - \bar{T}}{T}\right].$$

By defining $f(s) = \mathbb{1}\{s \in \mathcal{M}_{-1}\}$, we have:

$$\bar{p}_l(\bar{T}) = \mathbb{P}\left[\frac{1}{T}\sum_{t=1}^T f(\phi_t) \ge \frac{T - \bar{T}}{T}\right].$$

It is straightforward to verify that $\pi f = 1 - \pi_l$, $||f||_{\pi,\infty} = \pi_l$, and $||f||_{\pi,2}^2 = (1 - \pi_l)^2$. Therefore, we can utilize concentration Lemma 4 with any $0 < \gamma \le 1$:

$$\mathbb{P}_{q}\left[\frac{1}{T}\sum_{t=1}^{T}f(\phi_{t})-\pi f \geq \gamma\right]$$

$$\leq N_{q}\exp\left(-\frac{T\gamma^{2}\epsilon(K)}{8b^{2}\left(1+h\left(\frac{5\gamma}{b^{2}}\right)\right)}\right),$$

Now, let $\frac{T-\bar{T}}{T}=\gamma+\pi f$. Since the above result holds for any $0<\gamma<1$, it also holds for $\bar{T}<\pi_l T$. Therefore, we have:

$$\bar{p}_l(\bar{T}) \le N_q \exp\bigg(-\frac{T\gamma^2 \epsilon(K)}{8b^2 \Big(1+h\Big(\frac{5\gamma}{b^2}\Big)\Big)}\bigg),$$

for any $\gamma = \pi_1 - \frac{\bar{T}}{T}$ with $\bar{T} < \pi_l T$. Finally, we can complete the proof by substituting this result into Theorem 1.