

# A Tutorial on the Non-Asymptotic Theory of System Identification

Ingvar Ziemann<sup>1</sup>, Anastasios Tsiamis<sup>2</sup>, Bruce Lee<sup>1</sup>, Yassir Jedra<sup>3</sup>, Nikolai Matni<sup>1</sup>, and George J. Pappas<sup>1</sup>

<sup>1</sup>University of Pennsylvania

<sup>2</sup>ETH Zürich

<sup>3</sup>KTH Royal Institute of Technology

**Abstract**—This tutorial serves as an introduction to recently developed non-asymptotic methods in the theory of—mainly linear—system identification. We emphasize tools we deem particularly useful for a range of problems in this domain, such as the covering technique, the Hanson-Wright Inequality and the method of self-normalized martingales. We then employ these tools to give streamlined proofs of the performance of various least-squares based estimators for identifying the parameters in autoregressive models. We conclude by sketching out how the ideas presented herein can be extended to certain nonlinear identification problems. **Note:** For reasons of space, proofs have been omitted in this version and are available in an online version: <https://arxiv.org/abs/2309.03873>.

## NOTATION

Maxima (resp. minima) of two numbers  $a, b \in \mathbb{R}$  are denoted by  $a \vee b = \max(a, b)$  ( $a \wedge b = \min(a, b)$ ). For two sequences  $\{a_t\}_{t \in \mathbb{Z}}$  and  $\{b_t\}_{t \in \mathbb{Z}}$  we introduce the shorthand  $a_t \lesssim b_t$  if there exists a universal constant  $C > 0$  and an integer  $t_0$  such that  $a_t \leq Cb_t$  for every  $t \geq t_0$ . If  $a_t \lesssim b_t$  and  $b_t \lesssim a_t$  we write  $a_t \asymp b_t$ . Let  $X \subset \mathbb{R}^d$  and let  $f, g \in X \rightarrow \mathbb{R}$ . We write  $f = O(g)$  if  $\limsup_{x \rightarrow x_0} |f(x)/g(x)| < \infty$ , where the limit point  $x_0$  is typically understood from the context. We use  $\tilde{O}$  to hide logarithmic factors and write  $f = o(g)$  if  $\limsup_{x \rightarrow x_0} |f(x)/g(x)| = 0$ . We write  $f = \Omega(g)$  if  $\limsup_{x \rightarrow x_0} |f(x)/g(x)| > 0$ . For an integer  $N$ , we also define the shorthand  $[N] \triangleq \{1, \dots, N\}$ .

Expectation (resp. probability) with respect to all the randomness of the underlying probability space is denoted by  $\mathbb{E}$  (resp.  $\mathbb{P}$ ).

The Euclidean norm on  $\mathbb{R}^d$  is denoted  $\|\cdot\|_2$ , and the unit sphere in  $\mathbb{R}^d$  is denoted  $\mathbb{S}^{d-1}$ . The standard inner product on  $\mathbb{R}^d$  is denoted  $\langle \cdot, \cdot \rangle$ . We embed matrices  $M \in \mathbb{R}^{d_1 \times d_2}$  in Euclidean space by vectorization:  $\text{vec } M \in \mathbb{R}^{d_1 d_2}$ , where  $\text{vec}$  is the operator that vertically stacks the columns of  $M$  (from left to right and from top to bottom). For a matrix  $M$  the Euclidean norm is the Frobenius norm, i.e.,  $\|M\|_F \triangleq \|\text{vec } M\|_2$ . We similarly define the inner product of two matrices  $M, N$  by  $\langle M, N \rangle \triangleq \langle \text{vec } M, \text{vec } N \rangle$ . The transpose of a matrix  $M$  is denoted by  $M^T$  and  $\text{tr } M$  denotes its trace. For a matrix  $M \in \mathbb{R}^{d_1 \times d_2}$ , we order its singular values  $\sigma_1(M), \dots, \sigma_{d_1 \wedge d_2}(M)$  in descending order by magnitude. We also write  $\|M\|_{\text{op}}$  for its largest singular value:  $\|M\|_{\text{op}} \triangleq \sigma_1(M)$ . To not carry dimensional notation, we will also use  $\sigma_{\min}(M)$  for the smallest nonzero singular

value. For square matrices  $M \in \mathbb{R}^{d \times d}$  with real eigenvalues, we similarly order the eigenvalues of  $M$  in descending order as  $\lambda_1(M), \dots, \lambda_d(M)$ . In this case,  $\lambda_{\min}(M)$  will also be used to denote the minimum (possibly zero) eigenvalue of  $M$ . For two symmetric matrices  $M, N$ , we write  $M \succ N$  ( $M \succeq N$ ) if  $M - N$  is positive (semi-)definite.

## I. INTRODUCTION

Machine learning methods are at an ever increasing pace being integrated into domains that have classically been within the purview of controls. There is a wide range of examples, including perception-based control, agile robotics, and autonomous driving and racing. As exciting as these developments may be, they have been most pronounced on the experimental and empirical sides. To deploy these systems safely, stably, and robustly into the real world, we argue that a principled and integrated theoretical understanding of a) fundamental limitations and b) statistical optimality is needed. Under the past few years, a host of new techniques have been introduced to our field. Unfortunately, existing results in this area are relatively inaccessible to a typical first or second year graduate student in control theory, as they require both sophisticated mathematical tools not typically included in a control theorist's training (e.g., high-dimensional statistics and learning theory).

This tutorial seeks to provide a streamlined exposition of some of these recent advances that are most relevant to the non-asymptotic theory of linear system identification. Our aim is not to be encyclopedic but rather to give simple proofs of the main developments and to highlight and collect the key technical tools to arrive at these results. For a broader—and less technical—overview of the literature we point the reader to our recent survey [31]. It is also worth to point out that the classical literature on system identification has done a formidable job at—often very accurately—characterizing the asymptotic performance of identification algorithms [14]. Our aim is not to supplant this literature but rather to complement the asymptotic picture with finite sample guarantees by relaying recently developed technical tools drawn from high-dimensional probability, statistics and learning theory [35, 37].

### A. Problem Formulation

Let us now fix ideas. We are concerned with linear time-series models of the form:

$$Y_t = \theta^* X_t + V_t \quad t = 1, 2, \dots, T \quad (1)$$

where  $Y_{1:T}$  is a sequence of outputs (or targets) assuming values in  $\mathbb{R}^{d_y}$  and  $X_{1:T}$  is a sequence of inputs (or covariates) assuming values in  $\mathbb{R}^{d_x}$ . The goal of the user (or learner) is to recover the a priori unknown linear map  $\theta^* \in \mathbb{R}^{d_y \times d_x}$  using only the observations  $X_{1:T}$  and  $Y_{1:T}$ . The linear relationship in the regression model (1) is perturbed by a stochastic noise sequence  $V_{1:T}$  assuming values in  $\mathbb{R}^{d_y}$ . We refer to the regression model (1) as a time-series to emphasize the fact that the observations  $X_{1:T}$  and  $Y_{1:T}$  may arrive sequentially and in particular that past  $X_t$  and  $Y_t$  may influence future  $X_{t'}$  and  $Y_{t'}$  (i.e. with  $t' > t$ ).

a) *Example: Autoregressive Models.*: For instance, a model class of particular interest to us which is subsumed by (1) are the (vector) autoregressive exogenous models of order  $p$  and  $q$  (briefly ARX( $p, q$ )):

$$Y_t = \sum_{i=1}^p A_i^* Y_{t-i} + \sum_{j=1}^q B_j^* U_{t-j} + W_t \quad (2)$$

where typically  $U_{1:T-1}$  is a sequence of user specified inputs taking values in  $\mathbb{R}^{d_u}$  and  $W_{1:T}$  is an iid sequence of noise variables taking values in  $\mathbb{R}^{d_w}$ . If we are only interested in the parameters  $[A_{1:p}^* \ B_{1:q}^*]$ , we obtain the model (2) by setting

$$X_t = [Y_{t-1:t-p}^\top \ U_{t-1:t-q}^\top]^\top; \quad \theta^* = [A_{1:p}^* \ B_{1:q}^*]; \quad V_t = W_t. \quad (3)$$

We point out that that the above discussion presupposes that the order of the model, ( $p, q$ ), is known (there are ways around this).

In this tutorial we will provide the necessary tools to tackle the following problem.

**Problem I.1.** Fix  $\varepsilon > 0$ ,  $\delta \in (0, 1)$ , and a norm  $\|\cdot\|$ . Fix also a ‘reasonable’ estimator  $\hat{\theta}$  of  $\theta_*$  using a sample  $(X, Y)_{1:T}$  from (1). We seek to establish finite sample guarantees of the form

$$\|\hat{\theta} - \theta^*\| \leq \varepsilon \quad \text{with probability at least } 1 - \delta \quad (4)$$

where  $\varepsilon$  controls the accuracy (or rate) and the failure parameter  $\delta$  controls the confidence.

In the sequel, ‘reasonable’ estimator will typically mean some form of least squares estimator (7). These are introduced in Section I-B below. A bound of the form (4) is typically thought of as follows. We fix a priori the failure parameter  $\delta$  and then provide guarantees of the form  $\|\hat{\theta} - \theta^*\| \leq \varepsilon(T, \delta, P_{XY})$  where  $P_{XY}$  is the joint distribution of  $(X, Y)_{1:T}$ . Hence, the sample size  $T$ , the failure probability  $\delta$  and the distribution of the samples all

impact the performance guarantee  $\varepsilon$  we are able to establish. To be more specific,  $\varepsilon$  will typically be of the form

$$\varepsilon \propto (\text{Noise Scale}) \times \sqrt{\frac{\text{problem dimension} + \log(1/\delta)}{\text{sample size}}}. \quad (5)$$

Thus in principle, the best possible choice of  $\varepsilon^2$  can be thought of as a high probability version of the (inverse) signal-to-noise ratio of the problem at hand. The fact that the confidence parameter  $\delta$  typically affects (5) additively in  $\log(1/\delta)$  is consistent with classical asymptotic normality theory of estimators. One often expects the normalized difference  $T^{-1/2}(\hat{\theta} - \theta^*)$  to converge in law to a normal distribution [33]. In this tutorial we will provide tools that allow us to match such classical asymptotics but with a finite sample twist. Let us also remark that there often is a minimal requirement on the sample size necessary for a bound of the form (4)-(5) to hold. Such requirements are typically of the form

$$\text{sample size} \gtrsim \text{problem dimension} + \log(1/\delta). \quad (6)$$

Requirements such as (6) are called burn-in times and are related to the notion of persistence of excitation. They correspond to the rather minimal requirement that the parameter identification problem is feasible in the complete absence of observation noise.

### B. Least Squares Regression and the Path Ahead

Let us now return to the general setting of (1). Fix a subset  $M$  of  $\mathbb{R}^{d_y \times d_x}$ , called the model class. The estimator

$$\hat{\theta} \in \underset{\theta \in M}{\operatorname{argmin}} \frac{1}{T} \sum_{t=1}^T \|Y_t - \theta X_t\|_2^2 \quad (7)$$

is the least squares estimator (LSE) of  $\theta^*$  (with respect to  $M$ ). Often we simply set  $M = \mathbb{R}^{d_y \times d_x}$ . In this case, equivalently:

$$\hat{\theta} = \left( \sum_{t=1}^T Y_t X_t^\top \right) \left( \sum_{t=1}^T X_t X_t^\top \right)^\dagger \quad (8)$$

and the LSE reduces to the (minimum norm) ordinary least squares (OLS) estimator (8).

For simplicity, let us further assume that the (normalized) empirical covariance matrix:

$$\hat{\Sigma} \triangleq \frac{1}{T} \sum_{t=1}^T X_t X_t^\top; \quad (9)$$

is full rank almost surely.

a) *The Path Ahead.*: Let us now briefly sketch the path ahead to solve Problem I.1. If (9) is full rank—as required above—the estimator (8) admits the convenient error representation:

$$\begin{aligned} & \hat{\theta} - \theta^* \\ &= \left[ \left( \sum_{t=1}^T V_t X_t^\top \right) \left( \sum_{t=1}^T X_t X_t^\top \right)^{-1/2} \right] \left( \sum_{t=1}^T X_t X_t^\top \right)^{-1/2}. \end{aligned} \quad (10)$$

The leftmost term of (10) (in square brackets) can be shown to be (almost) time-scale invariant in many situations. For instance, if the noise  $V_{1:T}$  is a sub-Gaussian martingale difference sequence with respect to the filtration generated by the covariates  $X_{1:T}$ , one can invoke methods from the theory of self-normalized processes to show this [18, 1]. These methods are the topic of Section IV.

Whenever this is the case, the dominant term in the rate of convergence of the least squares estimator is  $(\sum_{t=1}^T X_t X_t^\top)^{-1/2}$ . In other words, providing control of the smallest eigenvalue of (9) effectively yields control of the rate of convergence of the least squares estimator in many situations. Thus, to analyze the rate of convergence of (7) when  $M = \mathbb{R}^{d_y \times d_x}$  it suffices to:

- Analyze the smallest eigenvalue (or lower tail) of (9). We provide such analyses in Section III
- Analyze the scale invariant term (in square brackets) of (10). This can in many situations be handled for instance by the self-normalized martingale method described in Section IV.

### C. Overview

Before covering these more technical topics in Section III and Section IV, we also briefly review some preliminaries from probability theory in Section II. We then demonstrate how to apply these ideas in the setting of identifying the parameters of an ARX( $p, q$ ) model of the form (2) in Section V. An alternative perspective not based on the decomposition (10) for more general least squares algorithms is given in Section VI. We conclude with a brief discussion on how the tools in Section VI can be extended to study more general nonlinear phenomena in Section VII.

## II. PRELIMINARIES: CONCENTRATION INEQUALITIES, PACKING AND COVERING

Before we proceed to tackle the more advanced question of analyzing the LSE (7), let us discuss a few preliminary inequalities that control the tail of a random variable. Our first inequality is Markov's.

**Lemma II.1.** *Let  $X$  be a nonnegative random variable. For every  $s > 0$  we have that*

$$\mathbf{P}(X \geq s) \leq s^{-1} \mathbf{E}[X]. \quad (11)$$

*Proof.* We have that  $\mathbf{E}[X] \geq \mathbf{E}[\mathbf{1}_{X \geq s} X] \geq s \mathbf{E}[\mathbf{1}_{X \geq s}]$ . Since  $\mathbf{E}[\mathbf{1}_{X \geq s}] = \mathbf{P}(X \geq s)$  the result follows by rearranging. ■

Typically, Markov's inequality itself is insufficient for our goals: we seek deviation inequalities that taper off exponentially fast in  $s$  and not as  $s^{-1}$ . Such scaling is for instance predicted asymptotically by the central limit theorem by the asymptotic normality of renormalized sums of square integrable iid random variables; that is, sums of the form  $S_n/\sqrt{n} = (X_1 + X_2 + \dots + X_n)/\sqrt{n}$  where the  $X_i, i \in [n]$  are independent and square integrable. For random variables possessing a moment generating function, Markov's inequality

can be "boosted" by the so-called "Chernoff trick". Namely, we apply Markov's inequality to the moment generating function of the random variable instead of applying it directly to the random variable itself.

**Corollary II.1** (Chernoff). *Fix  $s > 0$  and suppose that  $\mathbf{E} \exp(\lambda X)$  exists. Then*

$$\mathbf{P}(X \geq s) \leq \min_{\lambda \geq 0} e^{-\lambda s} \mathbf{E} \exp(\lambda X). \quad (12)$$

*Proof.* Fix  $\lambda \geq 0$ . We have:

$$\begin{aligned} \mathbf{P}(X \geq s) &= \mathbf{P}(\exp(\lambda X) \geq \exp(\lambda s)) \quad (\text{mono: } x \mapsto e^{\lambda x}) \\ &\leq e^{-\lambda s} \mathbf{E} \exp(\lambda X) \quad (\text{Markov}). \end{aligned}$$

The result follows by optimizing. ■

Recall that the function  $\psi_X(\lambda) \triangleq \mathbf{E} \exp(\lambda X)$  is the moment generating function of  $X$ . For instance, if  $X$  has univariate Gaussian distribution with mean zero and variance  $\sigma^2$ , the moment generating function appearing in (12) is just  $\mathbf{E} \exp(\lambda X) = \exp(\lambda^2 \sigma^2 / 2)$ . Hence the probability that said Gaussian exceeds  $s$  is upper-bounded:

$$\mathbf{P}(X > s) \leq \min_{\lambda \geq 0} e^{-\lambda s} \exp(\lambda^2 \sigma^2 / 2) = \exp\left(\frac{-s^2}{2\sigma^2}\right) \quad (13)$$

which (almost) exhibits the correct Gaussian tails as compared to (11).<sup>1</sup> It should be pointed out that assumptions stronger than those of the Central Limit Theorem (finite variance) are indeed needed for a non-asymptotic theory with sub-Gaussian tails as in (13). An assumption of this kind which is relatively standard in the literature is introduced next.

### A. Sub-Gaussian Concentration and the Hanson-Wright Inequality

In the sequel, we will not want to impose the Gaussian assumption. Instead, we define a class of random variables that admit reasoning analogous to (13).

**Definition II.1.** *We say that a random vector  $W$  taking values in  $\mathbb{R}^d$  is  $\sigma^2$ -sub-Gaussian ( $\sigma^2$ -subG) if for every  $v \in \mathbb{R}^d$  we have that:*

$$\mathbf{E} \exp(\langle v, W \rangle) \leq \exp\left(\frac{\sigma^2 \|v\|^2}{2}\right). \quad (14)$$

*Similarly, we say that  $W$  is  $\sigma^2$ -conditionally sub-Gaussian with respect to a  $\sigma$ -field  $\mathcal{F}$  if (14) holds with  $\mathbf{E}[\cdot]$  replaced by  $\mathbf{E}[\cdot | \mathcal{F}]$ .*

The term  $\sigma^2$  appearing in (14) is called the variance proxy of a sub-Gaussian random variable. The significance of this definition is that the one-dimensional projections  $X = \langle v, W \rangle$  (with  $\|v\| = 1$ ) satisfy the tail inequality (13). While obviously Gaussian random variables are sub-Gaussian with their variance as variance-proxy, there are many examples beyond Gaussians that fit into this framework. It is for instance straightforward to show that bounded random

<sup>1</sup>We write almost because  $\exp(-s^2/2\sigma^2) \approx \mathbf{P}(V > s)$  where  $V \sim N(0, \sigma^2)$  but the expression is not exact.

variables have variance proxy proportional to the square of their width [see eg. 37, Examples 2.3 and 2.4]. Moreover, it is readily verified that the normalized sum mentioned above— $S_n/\sqrt{n} = (X_1 + \dots + X_n)/\sqrt{n}$ —satisfies the same bound (13) provided that the entries of  $X_{1:n}$  are independent, mean zero and  $\sigma^2$ -sub-Gaussian. To see this, notice that the moment generating function “tensorizes” across products. Namely, for every  $\lambda \in \mathbb{R}$ :

$$\begin{aligned} \mathbf{E} \exp \left( \frac{\lambda}{\sqrt{n}} \sum_{i=1}^n X_i \right) &= \prod_{i=1}^n \mathbf{E} \exp \left( \frac{\lambda}{\sqrt{n}} X_i \right) \\ &\leq \prod_{i=1}^n \exp \left( \frac{\lambda^2 \sigma^2}{2n} \right) = \exp \left( \frac{\lambda^2 \sigma^2}{2} \right). \end{aligned}$$

Hence, by the exact same reasoning leading up to (13) such normalized sub-Gaussian sums satisfy the same tail bound (13).

When analyzing linear regression models, most variables of interest are typically either linear or quadratic in the variables of interest (cf. (10)). Hence, we also need to understand how squares of sub-Gaussian random variables behave. The next result shows that sub-Gaussian quadratic forms exhibit similar tail behavior to the Chi-squared distribution (often in the literature referred to as sub-exponential tails). It is known as the Hanson-Wright Inequality.

**Theorem II.1** ([6, 22]). *Let  $M \in \mathbb{R}^{d \times d}$ . Fix a random variable  $W = W_{1:d}$  where each  $W_i, i \in [d]$  is a scalar, mean zero and independent  $\sigma^2$ -sub-Gaussian random variable. Then for every  $s \in [0, \infty)$ :*

$$\begin{aligned} \mathbf{P} (|W^\top M W - \mathbf{E} W^\top M W| > s) \\ \leq 2 \exp \left( - \min \left( \frac{s^2}{144\sigma^4 \|M\|_F^2}, \frac{s}{16\sqrt{2}\sigma^2 \|M\|_{\text{op}}} \right) \right). \end{aligned}$$

The proof of Theorem II.1 is rather long and technical and thus relegated to the appendix. There, the reader may also find further useful concentration inequalities for quadratic forms in sub-Gaussian variables. In fact, there are plethora of useful concentration inequalities not covered here and the interested reader is urged to consult the first few chapters of [35].

### B. Covering and Discretization Arguments

We will often find ourselves in a situation where it is possible to obtain a scalar concentration bound but need this to hold uniformly for many random variables at once. The  $\varepsilon$ -net argument, which proceeds via the notion of covering numbers, is a relatively straightforward way of converting concentration inequalities for scalars into their counterparts for vectors, matrices and functions more generally.

The reader will for instance notice that the quantity being controlled by Theorem II.1 is a scalar quadratic form in sub-Gaussian random variables. By contrast, the empirical covariance matrix (9) is a matrix and so a conversion step is needed. This idea will be used frequently and in various

forms throughout the manuscript, so we review it briefly here for the particular case of controlling the operator norm of a random matrix. To this end, we notice that for any matrix  $M \in \mathbb{R}^{m \times d}$ :

$$\|M\|_{\text{op}}^2 = \max_{v \in \mathbb{S}^{d-1}} \langle Mv, Mv \rangle. \quad (15)$$

Hence, the operator norm of a random matrix is a maximum of scalar random variables indexed by the unit sphere  $\mathbb{S}^{d-1}$ .

Recall now that the union bound states that the probability that the maximum of a *finite collection* ( $|S| < \infty$ )  $\{X_i\}_{i \in S}$  of random variables exceeds a certain threshold can be bounded by the sum of their probabilities:

$$\mathbf{P} \left( \max_{i \in S} X_i > t \right) \leq \sum_{i \in S} \mathbf{P} (X_i > t). \quad (16)$$

Unfortunately, the unit sphere appearing (15) is not a finite set and so the union bound (16) cannot be directly applied. However, when the domain of optimization has geometric structure, one can often exploit this to leverage the union bound not directly but rather in combination with a discretization argument. Returning to our example of the operator norm of a matrix, the set  $S$  appearing in (16) will be a discretized version of the unit sphere  $\mathbb{S}^{d-1}$ .

The following notion is key.

**Definition II.2.** *Let  $(X, d)$  be a compact metric space and fix  $\varepsilon > 0$ . A subset  $\mathcal{N}$  of  $X$  is called an  $\varepsilon$ -net of  $X$  if every point of  $X$  is within radius  $\varepsilon$  of a point of  $\mathcal{N}$ :*

$$\sup_{x \in X} \inf_{x' \in \mathcal{N}} d(x, x') \leq \varepsilon. \quad (17)$$

Moreover, the minimal cardinality of  $\mathcal{N}$  necessary such that (17) holds is called the covering number at resolution  $\varepsilon$  of  $(X, d)$  and is denoted  $\mathcal{N}(\varepsilon, X, d)$ .

We will not explore this notion in full, but simply content ourselves to note that it plays very nicely with the notion of operator norm.

**Lemma II.2** (Lemma 4.4.1 in [35]). *Let  $M \in \mathbb{R}^{m \times d}$  and let  $\varepsilon \in (0, 1)$ . Then for any  $\varepsilon$ -net  $\mathcal{N}$  of  $(\mathbb{S}^{d-1}, \|\cdot\|_2)$  we have that:*

$$\|M\|_{\text{op}} \leq \frac{1}{1 - \varepsilon} \sup_{v \in \mathcal{N}} \|Mv\|_2. \quad (18)$$

Hence at a small multiplicative cost, the computation of the operator norm can be restricted to the discretized sphere  $\mathcal{N}$ . Our intention is now to apply the union bound (16) to the right hand side of (18). To do so, we also need control of the size (cardinality) of the  $\varepsilon$ -net.

**Lemma II.3** (Corollary 4.2.13 in [35]). *For any  $\varepsilon > 0$  the covering numbers of  $\mathbb{S}^{d-1}$  satisfy*

$$\mathcal{N}(\varepsilon, \mathbb{S}^{d-1}, \|\cdot\|) \leq \left(1 + \frac{1}{2\varepsilon}\right)^d. \quad (19)$$

We now provide two instances of this covering argument combined with the union bound. The second of these uses an

alternative variational characterization of the operator norm but otherwise similar ideas.

**Lemma II.4.** *Let  $M$  be an  $m \times d$  random matrix, and  $\epsilon \in (0, 1)$ . Furthermore, let  $\mathcal{N}$  be an  $\epsilon$ -net of  $\mathbb{S}^{d-1}$  of minimal cardinality. Then for all  $\rho > 0$ , we have*

$$\mathbf{P}(\|M\|_{\text{op}} > \rho) \leq \left(\frac{2}{\epsilon} + 1\right)^d \max_{v \in \mathcal{N}} \mathbf{P}(\|Mv\|_2 > (1 - \epsilon)\rho).$$

**Lemma II.5.** *Let  $M$  be an  $d \times d$  symmetric random matrix, and let  $\epsilon \in (0, 1/2)$ . Furthermore, let  $\mathcal{N}$  be an  $\epsilon$ -net of  $\mathbb{S}^{d-1}$  with minimal cardinality. Then for all  $\rho > 0$ , we have*

$$\mathbf{P}(\|M\|_{\text{op}} > \rho) \leq \left(\frac{2}{\epsilon} + 1\right)^d \max_{v \in \mathcal{N}} \mathbf{P}(|v^\top Mv| > (1 - 2\epsilon)\rho).$$

Lemma II.4 and Lemma II.5 exploit two different variational forms of the operator norm. Namely for any  $M$  we have that  $\|M\|_{\text{op}}^2 = \sup_{v \in \mathbb{S}^{d-1}} \|Mv\|_2^2$  and in addition, when  $M$  is symmetric we also have,  $\|M\|_{\text{op}} = \sup_{v \in \mathbb{S}^{d-1}} |v^\top Mv|$ . The proof of these last two lemmas are standard and can be found for example in [35, Chapter 4].

### C. Concentration of the Covariance Matrix of Linear Systems

To not get lost in the weeds, let us provide an example showcasing the use of Theorem II.1 due to [8]. Recall that the matrix  $\widehat{\Sigma}$  appearing in (9) is crucial to the performance of the least squares estimator. We will now see that this matrix is well-conditioned when we consider stable first order auto-regressions of the form:

$$X_{t+1} = A^* X_t + W_t \quad t = 1, \dots, T \quad (20)$$

taking values in  $\mathbb{R}^{d_X}$  and with  $W_{1:T}$  iid isotropic and  $K^2$ -subG. By stable we mean that the spectrum of  $A^*$  is contained in the unit disc.

The following result is a consequence of the Hanson-Wright inequality together with the discretization strategy outlined in Section II-B.

**Theorem II.2.** *Let  $\epsilon > 0$  and set  $M \triangleq \left(\sum_{t=1}^T \sum_{k=0}^{t-1} (A^*)^k (A^{*,\top})^k\right)^{-\frac{1}{2}}$ . Let also  $\mathbf{L}$  be the linear operator such that  $X_{1:T} = \mathbf{L}W_{1:T}$ . Then simultaneously for every  $i \in [d_X]$ :*

$$\begin{aligned} (1-\epsilon)^2 \lambda_{\min} \left( \sum_{t=1}^T \sum_{k=0}^{t-1} (A^*)^k (A^{*,\top})^k \right) &\leq \lambda_i \left( \sum_{t=1}^T X_t X_t^\top \right) \\ &\leq (1+\epsilon)^2 \lambda_{\max} \left( \sum_{t=1}^T \sum_{k=0}^{t-1} (A^*)^k (A^{*,\top})^k \right) \end{aligned}$$

holds with probability at least

$$1 - \exp\left(-\frac{\epsilon^2}{576 K^2 \|M\|_{\text{op}}^2 \|\mathbf{L}\|_{\text{op}}^2} + d_X \log(18)\right). \quad (21)$$

Put differently, on the same event as in Theorem II.2, the spectrum of

$$\widehat{\Sigma} = \frac{1}{T} \sum_{t=1}^T X_t X_t^\top \quad (22)$$

is sandwiched by its population counterpart within a small multiplicative factor. The result holds with high probability for strictly stable systems.

The quantity  $\|L\|_{\text{op}}$  in (21) grows very quickly as the spectral radius of  $A^*$  tends to 1; Theorem II.2 becomes vacuous in the marginally stable regime. It turns out that requirement of two-sided concentration—the sandwiching of the entire spectrum—is too stringent a requirement to obtain bounds that degrade gracefully with the stability of the system. Fortunately, we only need sharp control of the lower half of the spectrum to control the error (10). This motivates Section III below, in which we will see how to relax the stability assumption and analyze more general linear systems.

### D. Notes

The basic program carried out in Section II-C can be summarized as follows: (1) introduce a discretization of the problem considered—for matrices this is typically a discretization of the unit sphere; (2) prove an exponential inequality for a family of scalar random variables corresponding to one-dimensional projection of the discretization—in our case: prove bounds on the moment generating function of quadratic forms in random matrices; and (3) conclude to obtain a uniform bound by using the union bound across the discretization. This roughly summarizes the proof of Theorem II.2. These tools are thematic throughout this manuscript.

## III. THE LOWER SPECTRUM OF THE EMPIRICAL COVARIANCE

Recall that our outline of the analysis of the least squares estimator in Section I-B consists of two main components, one of which being the lower tail of the empirical covariance matrix (9). In this section we provide a self-contained analysis of this random matrix for a class of "causal" systems. Moreover, we will emphasize only the lower tail of this random matrix as to sidestep issues with bounds degrading with the stability of the system considered. This allows us to quantitatively separate the notions of persistence of excitation and stability.

Let us now carry out this program. Fix two integers  $T$  and  $k$  such that  $T/k \in \mathbb{N}$ . We consider causal processes of the form  $X_{1:T} = (X_1^\top, \dots, X_T^\top)^\top$  evolving on  $\mathbb{R}^d$ . More precisely, we assume the existence of an isotropic sub-Gaussian process evolving on  $\mathbb{R}^p$ ,  $W_{1:T}$  with  $\mathbf{E}W_{1:T}W_{1:T}^\top = I_{pT}$  and a (block-) lower-triangular matrix  $\mathbf{L} \in \mathbb{R}^{dT \times pT}$  such that

$$X_{1:T} = \mathbf{L}W_{1:T}. \quad (23)$$

We will assume that all the  $pT$ -many entries of  $W_{1:T}$  are independent  $K^2$ -sub-Gaussian for some positive  $K \in \mathbb{R}$ .

We say that  $X_{1:T}$  is  $k$ -causal if the matrix  $\mathbf{L}$  has the block lower-triangular form:

$$\mathbf{L} = \begin{bmatrix} \mathbf{L}_{1,1} & 0 & 0 & 0 & 0 \\ \mathbf{L}_{2,1} & \mathbf{L}_{2,2} & 0 & 0 & 0 \\ \mathbf{L}_{3,1} & \mathbf{L}_{3,2} & \mathbf{L}_{3,3} & 0 & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \mathbf{L}_{T/k,1} & \dots & \dots & \dots & \dots \mathbf{L}_{T/k,T/k} \end{bmatrix} = \begin{bmatrix} \mathbf{L}_1 \\ \mathbf{L}_2 \\ \mathbf{L}_3 \\ \vdots \\ \mathbf{L}_{T/k} \end{bmatrix} \quad (24)$$

where each  $\mathbf{L}_{ij} \in \mathbb{R}^{dk \times pk}$ ,  $i, j \in [T/k] \triangleq \{1, 2, \dots, T/k\}$ . In brief, we say that  $X_{0:T-1}$  satisfying the above construction is  $k$ -causal with independent  $K^2$ -sub-Gaussian increments.

Obviously, every 1-causal process is  $k$ -causal for every  $k \in \mathbb{N}$  as long as the divisibility condition holds. To analyze the lower tail of the empirical covariance of  $X_{0:T-1}$  we will also associate a decoupled random process

$$\tilde{X}_{1:T} = \text{blkdiag}(\mathbf{L}_{11}, \dots, \mathbf{L}_{T/k, T/k}) W_{1:T}.$$

Hence, the process  $\tilde{X}_{1:T}$  is generated in much the same way as  $X_{1:T}$  but by removing the sub-diagonal entries of  $\mathbf{L}$ :

$$\tilde{\mathbf{L}} \triangleq \begin{bmatrix} \mathbf{L}_{1,1} & 0 & 0 & 0 \\ 0 & \mathbf{L}_{2,2} & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & \mathbf{L}_{T/k, T/k} \end{bmatrix} \implies \tilde{X}_{1:T} = \tilde{\mathbf{L}} W_{1:T}.$$

We emphasize that by our assumptions on  $W_{1:T}$  and the block-diagonal structure of  $\tilde{\mathbf{L}}$  the variables  $\tilde{X}_{1:k}, \tilde{X}_{k+1:2k}, \dots, \tilde{X}_{T-k+1:T}$  are all independent of each other; they have been decoupled. This decoupled process will effectively dictate our lower bound, and we will show under relatively mild assumptions that

$$\lambda_{\min} \left( \frac{1}{T} \sum_{t=1}^T X_t X_t^\top \right) \gtrsim \lambda_{\min} \left( \frac{1}{T} \sum_{t=1}^T \mathbf{E} \tilde{X}_t \tilde{X}_t^\top \right) \quad (25)$$

with probability that approaches 1 at an exponential rate in the sample size  $T$ . More precisely, the following statement is the main result of this section.

**Theorem III.1.** *Fix an integer  $k \in \mathbb{N}$ , let  $T \in \mathbb{N}$  be divisible by  $k$  and suppose  $X_{1:T}$  is a  $k$ -causal process taking values in  $\mathbb{R}^d$  with  $K^2$ -sub-Gaussian increments. Suppose further that the diagonal blocks are all equal:  $\mathbf{L}_{j,j} = \mathbf{L}_{1,1}$  for all  $j \in [T/k]$ . Suppose  $\lambda_{\min} \left( \sum_{t=1}^T \mathbf{E} \tilde{X}_t \tilde{X}_t^\top \right) > 0$ . We have that:*

$$\mathbf{P} \left( \frac{1}{T} \sum_{t=1}^T X_t X_t^\top \not\geq \frac{1}{8T} \sum_{t=1}^T \mathbf{E} \tilde{X}_t \tilde{X}_t^\top \right) \leq (C_{\text{sys}})^d \exp \left( -\frac{T}{576K^2k} \right) \quad (26)$$

where

$$C_{\text{sys}} \triangleq 1 + \frac{\left( \frac{T \|\mathbf{L}\mathbf{L}^\top\|_{\text{op}}}{18k \lambda_{\min} \left( \sum_{t=1}^T \mathbf{E} X_t X_t^\top \right)} + 9 \right) \lambda_{\max} \left( \sum_{t=1}^T \mathbf{E} X_t X_t^\top \right)}{2\sqrt{2} \lambda_{\min} \left( \sum_{t=1}^T \mathbf{E} \tilde{X}_t \tilde{X}_t^\top \right)}. \quad (27)$$

To parse Theorem III.1, note that it simply informs us that there exist a system-dependent constant  $C_{\text{sys}}$ —which itself has no more than polynomial dependence on relevant quantities—such that if

$$T/k \geq 576K^2(d \log C_{\text{sys}} + \log(1/\delta)) \quad (28)$$

then on an event with probability mass at least  $1 - \delta$ :

$$\frac{1}{T} \sum_{t=1}^T X_t X_t^\top \succeq \frac{1}{8T} \sum_{t=1}^T \mathbf{E} \tilde{X}_t \tilde{X}_t^\top.$$

**Remark III.1.** *Since the blocks of  $\mathbf{L}$  can be regarded to specify the noise-to-output map, the assumption that the diagonal blocks are constant is for instance satisfied by linear time-invariant (LTI) systems. The assumption can be removed at the cost of a more complicated expression.*

The next example serves as the archetype for the reduction from  $\mathbf{L}$  to  $\tilde{\mathbf{L}}$ .

**Example III.1.** *Suppose that (23) is specified via*

$$X_t = A_* X_{t-1} + B_* W_t \quad (29)$$

for  $t \in [T]$  and where  $(A_*, B_*) \in \mathbb{R}^{d_x \times d_x + d_x \times d_w}$ . We set  $d = d_x$  and  $p = d_w$  in the theorem above. The reduction from  $X_{1:T} = \mathbf{L} W_{1:T}$  to  $\tilde{X}_{1:T} = \text{blkdiag}(\mathbf{L}_{11}, \dots, \mathbf{L}_{T/k, T/k}) W_{1:T}$  corresponds to replacing a single trajectory from the linear system (29) of length  $T$  by  $T/k$  trajectories of length  $k$  each and sampled independently of each other. The price we pay for decoupling these systems is that our lower bound is dictated by the gramians up to range  $k$ :

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \mathbf{E} \tilde{X}_t \tilde{X}_t^\top &= \frac{1}{k} \sum_{t=1}^k \mathbf{E} \tilde{X}_t \tilde{X}_t^\top \\ &= \frac{1}{k} \sum_{t=1}^k \sum_{j=0}^{t-1} (A^*)^j B^* B^{*\top} (A^{*\top})^j \end{aligned}$$

instead of the gramians up to range  $T$ :

$$\frac{1}{T} \sum_{t=1}^T \mathbf{E} X_t X_t^\top = \frac{1}{T} \sum_{t=1}^T \sum_{j=0}^{t-1} (A^*)^j B^* B^{*\top} (A^{*\top})^j.$$

Put differently, the reduction from  $\mathbf{L}$  to  $\tilde{\mathbf{L}}$  can be thought of as restarting the system every  $k$  steps.

Comparing with Theorem II.2, the advantage of Theorem III.1 is that it allows us to provide persistence-of-excitation type guarantees that do not rely strongly on the stability of the underlying system. While Theorem II.2 gives in

principle stronger two-sided concentration results, it comes at the cost of the guarantees becoming vacuous as the spectral radius of  $A^*$  in Example III.1 tends to marginal stability (tends to 1). By contrast, Theorem III.1 does not exhibit such a blow-up since the dependence on  $C_{\text{sys}}$  in (28) is logarithmic (instead of polynomial). The distinction might seem small, but it is qualitatively important as it (almost) decouples the phenomena of stability and persistence of excitation.

#### A. A Decoupling Inequality for sub-Gaussian Quadratic Forms

Our proof of Theorem III.1 will make heavy use of Proposition III.1 below. This is the crucial probabilistic inequality that allows us to decouple—or restart as discussed in Example III.1.

**Proposition III.1.** *Fix  $K \geq 1$ ,  $x \in \mathbb{R}^n$  and a symmetric positive semidefinite  $Q \in \mathbb{R}^{(n+m) \times (n+m)}$  of the form  $Q = \begin{bmatrix} Q_{11} & Q_{12} \\ Q_{21} & Q_{22} \end{bmatrix}$  with  $Q_{22} \succ 0$ . Let  $W$  be an  $m$ -dimensional mean zero, isotropic and  $K^2$ -sub-Gaussian random vector with independent entries. Then for every  $\lambda \in \left[0, \frac{1}{8\sqrt{2}K^2\|Q_{22}\|_{\text{op}}}\right]$  it holds true that:*

$$\mathbf{E} \exp \left( -\lambda \begin{bmatrix} x \\ W \end{bmatrix}^\top \begin{bmatrix} Q_{11} & Q_{12} \\ Q_{21} & Q_{22} \end{bmatrix} \begin{bmatrix} x \\ W \end{bmatrix} \right) \leq \exp \left( -\lambda \text{tr} Q_{22} + 36K^4\lambda^2 \text{tr} Q_{22}^2 \right). \quad (30)$$

By combining Lemma III.1 below with the exponential form of Hanson-Wright we obtain the exponential inequality (30), which in the sequel will allow us to control the lower tail of the conditionally random quadratic form

$$\begin{bmatrix} x \\ W \end{bmatrix}^\top \begin{bmatrix} Q_{11} & Q_{12} \\ Q_{21} & Q_{22} \end{bmatrix} \begin{bmatrix} x \\ W \end{bmatrix}.$$

We point out that (30) is not the best possible if the entries of  $W$  are independent and Gaussian as opposed to just isotropic and sub-Gaussian. In this case, the factor  $36K^4\lambda^2(\text{tr} Q_{22})^2$  in (30) can be improved to  $\frac{\lambda^2}{2} \text{tr} Q_{22}^2$  and the inequality can be shown to hold for the entire range of non-negative  $\lambda$  [38, Lemma 2.1]. Irrespectively, we will see in the sequel that it captures the correct qualitative behavior.

**Lemma III.1** (sub-Gaussian Decoupling). *Fix  $K \geq 1$ ,  $x \in \mathbb{R}^n$  and a symmetric positive semidefinite  $Q \in \mathbb{R}^{(n+m) \times (n+m)}$  of the form  $Q = \begin{bmatrix} Q_{11} & Q_{12} \\ Q_{21} & Q_{22} \end{bmatrix}$ . Let  $W$  be an  $m$ -dimensional mean zero and  $K^2$ -sub-Gaussian random vector. Then for every  $\lambda \in \left[0, \frac{1}{4K^2\|Q_{22}\|_{\text{op}}}\right]$  it holds true that:*

$$\mathbf{E} \exp \left( -\lambda \begin{bmatrix} x \\ W \end{bmatrix}^\top \begin{bmatrix} Q_{11} & Q_{12} \\ Q_{21} & Q_{22} \end{bmatrix} \begin{bmatrix} x \\ W \end{bmatrix} \right) \leq \sqrt{\mathbf{E} \exp \left( -2\lambda W^\top Q_{22} W \right)}. \quad (31)$$

Once equipped with (31), Proposition III.1 follows immediately.

#### B. The Lower Tail of the Empirical Covariance of Causal sub-Gaussian Processes

Repeated application of Proposition III.1 to the process  $X_{1:T} = \mathbf{L}W_{1:T}$  in combination with the tower property of conditional expectation yields the following exponential inequality that controls the lower tail of (9) in any fixed direction.

**Theorem III.2.** *Fix an integer  $k \in \mathbb{N}$ , let  $T \in \mathbb{N}$  be divisible by  $k$  and suppose  $X_{1:T}$  is a  $k$ -causal process driven by independent  $K^2$ -sub-Gaussian increments as described in Section III. Fix also a matrix  $\Delta \in \mathbb{R}^{d' \times d}$ . Let  $Q_{\text{max}} \triangleq \max_{j \in [T/k]} \|\mathbf{L}_{j,j}^\top \text{blkdiag}(\Delta^\top \Delta) \mathbf{L}_{j,j}\|_{\text{op}}$ . Then for every  $\lambda \in \left[0, \frac{1}{8\sqrt{2}K^2 Q_{\text{max}}}\right]$ :*

$$\begin{aligned} & \mathbf{E} \exp \left( -\lambda \sum_{t=1}^T \|\Delta X_t\|_2^2 \right) \\ & \leq \exp \left( -\lambda \sum_{j=1}^{T/k} \text{tr} \left( \mathbf{L}_{j,j}^\top \text{blkdiag}(\Delta^\top \Delta) \mathbf{L}_{j,j} \right) \right. \\ & \quad \left. + 36K^4\lambda^2 \sum_{j=1}^{T/k} \text{tr} \left( \mathbf{L}_{j,j}^\top \text{blkdiag}(\Delta^\top \Delta) \mathbf{L}_{j,j} \right)^2 \right). \end{aligned}$$

To appreciate the terms appearing in Theorem III.2, it is worth to point out that

$$\sum_{j=1}^{T/k} \text{tr} \left( \mathbf{L}_{j,j}^\top \text{blkdiag}(\Delta^\top \Delta) \mathbf{L}_{j,j} \right) = \sum_{t=1}^T \mathbf{E} \|\Delta \tilde{X}_t\|_2^2.$$

Hence Theorem III.2 effectively passes the expectation inside the exponential at the cost of working with the possibly less excited process  $\tilde{X}_{1:T}$  and a quadratic correction term. Note also that the assumption that  $T$  is divisible by  $k$  is not particularly important. If not, let  $T'$  be the largest integer such that  $T'/k \in \mathbb{N}$  and  $T' \leq T$  and apply the result with  $T'$  in place of  $T$ .

The significance of Theorem III.2 is demonstrated by the following simple observation, which is just the Chernoff approach applied to the exponential inequality in Theorem III.2.

**Lemma III.2.** *Fix an integer  $k \in \mathbb{N}$ , let  $T \in \mathbb{N}$  be divisible by  $k$  and suppose  $X_{1:T}$  is a  $k$ -causal process with independent  $K^2$ -sub-Gaussian increments. Suppose further that the diagonal blocks are all equal:  $\mathbf{L}_{j,j} = \mathbf{L}_{1,1}$  for all  $j \in [T/k]$ . For every size-conforming matrix  $\Delta$  we have that:*

$$\mathbf{P} \left( \sum_{t=1}^T \|\Delta X_t\|_2^2 \leq \frac{1}{2} \sum_{t=1}^T \mathbf{E} \|\Delta \tilde{X}_t\|_2^2 \right) \leq \exp \left( -\frac{T}{576K^2k} \right). \quad (32)$$

Note that Lemma III.2 only yields *pointwise* control of the empirical covariance—i.e. pointwise on the sphere  $\mathbb{S}^{d-1}$ . By setting  $\Delta = v \in \mathbb{S}^{d-1}$ , the result holds for a fixed vector on the sphere, but not uniformly for all such vectors at once. Thus, returning to our over-arching goal of providing control

of the smallest eigenvalue of the empirical covariance matrix (9), we now combine (32) (using  $d' = 1$ ) with a union bound. This approach yields Theorem III.1.<sup>2</sup>

### C. Notes

In this manuscript we have chosen a perhaps less well-known but conceptually simpler approach to establishing lower bounds on the empirical covariance matrix Equation (25). The first proof of a statement similar to Theorem III.1 is due to [27] which in turn relies on a more advanced notion from probability theory known as the small-ball method, due to [15]. The emphasis therein is on anti-concentration—which can hold under milder moment assumptions—rather than concentration. However, the introduction of this tool is not necessary for Gaussian (or sub-Gaussian) system identification. For instance, Sarkar and Rakhlin [24] leverage the method of self-normalized martingales introduced in Section IV below.

Our motivation for providing a different proof is to streamline the exposition as to fit control of the lower tail into the “standard machinery”, which roughly consists of: (1) prove a family of scalar exponential inequalities, (2) invoke the Chernoff method, and (3) conclude by a discretization argument and a union bound to port the result from scalars to matrices. Our proof here follows this outline and emphasizes the exponential inequality in Theorem III.2. We finally remark that the proof presented here is new to the literature and extends a result in [38] from the Gaussian setting to the sub-Gaussian setting.

## IV. SELF-NORMALIZED MARTINGALE BOUNDS

The objective in this section is to bound the operator and Frobenius norms of the self-normalized term of (10):

$$\left( \sum_{t=1}^T V_t X_t^\top \right) \left( \sum_{t=1}^T X_t X_t^\top \right)^{-1/2}. \quad (33)$$

This object has special structure. Firstly, in many cases of interest, e.g. the autoregressive model in (2), the noise  $V_t$  is independent of  $X_k$  for all  $k \leq t$ . This is what provides martingale structure, as will be made precise shortly. Secondly, it is self-normalized: if the covariates  $X_t$  are large for some  $t$ , then any increase in the left sum will be compensated by an increase in the sum in the term on the right. Together, these properties make the object above a *self-normalized martingale* term.

To express results generally and compactly, several definitions are in order.

**Definition IV.1. (Filtration and Adapted Process)** A sequence of sub- $\sigma$ -algebras  $\{\mathcal{F}_t\}_{t=1}^T$  is said to be a filtration if  $\mathcal{F}_t \subseteq \mathcal{F}_k$  for  $t \leq k$ . A stochastic process  $\{W_t\}_{t=1}^T$  is said to be adapted to the filtration  $\{\mathcal{F}_t\}_{t=1}^T$  if for all  $t \geq 1$ ,  $W_t$  is  $\mathcal{F}_t$ -measurable.

Conditioning on a sub- $\sigma$ -algebra provides partial information about the total randomness. Therefore, the requirement that a filtration is non-decreasing captures the fact that

information is not forgotten. An adapted process is one in which all the randomness at a particular time is explained by the information in the filtration up to that time.

**Definition IV.2. (Martingale)** Consider a stochastic process  $\{W_t\}_{t=1}^T$  which is adapted to a filtration  $\{\mathcal{F}_t\}_{t=1}^T$ . This process is called a martingale if for all  $1 \leq t \leq T$ ,  $W_t$  is integrable and for all  $1 \leq t < T$ ,  $\mathbf{E}[W_{t+1} | \mathcal{F}_t] = W_t$ .

Martingales model causal or non-anticipative processes. To better appreciate this, note that the increments  $W_{t+1} - W_t$  are mean zero and conditionally orthogonal to the past; they can be thought of as the “next step” in a random walk whose path is traced out by  $W_t$ .

In the context of the linear time-series model in (1), we may define the sub- $\sigma$ -algebras  $\mathcal{F}_t$  as those induced by the randomness up to time  $t$ :  $\mathcal{F}_t = \sigma(X_1, \dots, X_{t+1}, V_1, \dots, V_t)$ . In this case, the process  $\{X_t\}_{t=1}^T$  is adapted to the filtration  $\{\mathcal{F}_{t-1}\}_{t=1}^T$  and the process  $\{V_t\}_{t=1}^T$  is adapted to the filtration  $\{\mathcal{F}_t\}_{t=1}^T$ .

Recall now again that the “numerator” in the ordinary least squares error is (33). We see that if we define the sum,  $S_t \triangleq \sum_{s=1}^t V_s X_s^\top$ , then the process  $\{S_t\}_{t=1}^T$  is adapted to  $\{\mathcal{F}_t\}_{t=1}^T$ . Furthermore,  $\mathbf{E}(S_{t+1} | \mathcal{F}_t) = S_t + \mathbf{E}(V_{t+1} | \mathcal{F}_t) X_{t+1}^\top$ . In particular, as long as the noise has conditional mean zero ( $\mathbf{E}(V_{t+1} | \mathcal{F}_t) = 0$ ), the process  $\{S_t\}_{t=1}^T$  is a martingale.<sup>3</sup> Normalizing the sum  $S_t$  by the covariates as  $S_t \left( \sum_{s=1}^t X_s X_s^\top \right)^{-1/2}$  almost preserves the martingale structure. Expressions of this type are called self-normalized martingales—although we stress that they are not strictly speaking martingales but only constructed from them.

We now state bounds on the operator and Frobenius norms of the self-normalized martingale. The main idea behind the result—the technique of pseudo-maximization—is due to [21]. The formulations presented here are a consequence of Theorem 3.4 in [1].

**Theorem IV.1. (Special cases of Theorem 3.4 in [1])** Let  $\{\mathcal{F}_t\}_{t=0}^T$  be a filtration such that  $\{X_t\}_{t=1}^T$  is adapted to  $\{\mathcal{F}_{t-1}\}_{t=1}^T$  and  $\{V_t\}_{t=1}^T$  is adapted to  $\{\mathcal{F}_t\}_{t=1}^T$ . Additionally, suppose that for all  $1 \leq t \leq T$ ,  $V_t$  is  $\sigma^2$ -conditionally sub-Gaussian with respect to  $\mathcal{F}_t$ . Let  $\Sigma$  be a positive definite matrix in  $\mathbb{R}^{d_x \times d_x}$ . For a fixed  $T \in \mathbb{N}_+$  and  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ ,

$$\begin{aligned} & \left\| \sum_{t=1}^T V_t X_t^\top \left( \Sigma + \sum_{t=1}^T X_t X_t^\top \right)^{-1/2} \right\|_F^2 \\ & \leq d_Y \sigma^2 \log \left( \frac{\det \left( \Sigma + \sum_{t=1}^T X_t X_t^\top \right)}{\det(\Sigma)} \right) + 2\sigma^2 \log \frac{1}{\delta}. \end{aligned}$$

Additionally, for a fixed  $T \in \mathbb{N}_+$  and  $\delta \in (0, 1)$ , with

<sup>2</sup>Similar results can also be obtained for restricted eigenvalues.

<sup>3</sup>Indeed,  $S_t$  is a so-called martingale transform of  $X_{1:t}$ .



probability at least  $1 - \delta$ ,

$$\begin{aligned} & \left\| \sum_{t=1}^T V_t X_t^\top \left( \Sigma + \sum_{t=1}^T X_t X_t^\top \right)^{-1/2} \right\|_{\text{op}}^2 \\ & \leq 4\sigma^2 \log \left( \frac{\det \left( \Sigma + \sum_{t=1}^T X_t X_t^\top \right)}{\det(\Sigma)} \right) \\ & \quad + 8d_Y \sigma^2 \log 5 + 8\sigma^2 \log \frac{1}{\delta}. \end{aligned}$$

Note that the quantities bounded above have a positive definite matrix  $\Sigma$  added to the normalization quantities that was not present in the original term of interest, (33). Furthermore, the covariates  $\sum_{t=1}^T X_t X_t^\top$  appear in the upper bound. Hence, one typically combines the self-normalized martingale bound with some weak form of concentration.<sup>4</sup> This is done in Section V.

In the sequel, we prove the above bounds. To obtain the bound on the Frobenius norm, we directly consider the object

$$\left\| \left( \sum_{t=1}^T V_t X_t^\top \right) \left( \sum_{t=1}^T X_t X_t^\top \right)^{-1/2} \right\|_F, \quad (34)$$

while to obtain the bound on the operator norm we consider the following vector norm for an arbitrary unit vector  $w$  as an intermediate step:

$$\left\| \left( w^\top \sum_{t=1}^T V_t X_t^\top \right) \left( \sum_{t=1}^T X_t X_t^\top \right)^{-1/2} \right\|_2 \quad (35)$$

and combine with a covering argument (recall Section II-B).

#### A. Exponential Inequalities via Pseudo-maximization

We begin by neglecting the details of the process that generated the data in (33). In particular, consider a random matrix  $P$  assuming values in  $\mathbb{R}^{d_\eta \times d_x}$  for  $d_\eta \in \mathbb{N}_+$  and a random matrix  $Q$  assuming values in  $\mathbb{R}^{d_x \times d_x}$  with  $Q$  almost surely nonsingular. Bounding the quantities in (34) and (35) are special cases of bounding  $\|PQ^{-1/2}\|_F$ . A naive first approach is to apply a Chernoff bound (12). Doing so results in the inequality

$$\begin{aligned} & \mathbf{P} \left( \|PQ^{-1/2}\|_F \geq x \right) \\ & \leq \min_{\lambda \geq 0} \exp \left( -\frac{\lambda}{2} x^2 \right) \mathbf{E} \exp \left( \frac{\lambda}{2} \|PQ^{-1/2}\|_F^2 \right). \end{aligned}$$

If it is possible to bound the moment generating function  $\mathbf{E} \exp \left( \frac{\lambda}{2} \|PQ^{-1/2}\|_F^2 \right)$  by one for some  $\lambda > 0$ , then the above bound provides an exponential inequality. Obtaining a bound of the form  $\mathbf{E} \exp \left( \frac{\lambda}{2} \|PQ^{-1/2}\|_F^2 \right) \leq 1$  requires very strong assumptions on  $P$  and  $Q$  which would not be suitable for our purposes. However, we may observe that  $\frac{1}{2} \|PQ^{-1/2}\|_F^2 = \max_{\Lambda} \text{tr} \left( P\Lambda - \frac{1}{2} \Lambda^\top Q \Lambda \right)$ . This motivates

<sup>4</sup>Alternatively, in the analysis of ridge regression,  $\Sigma$  takes the role of the penalizing matrix which can be tuned.

the following canonical assumption in self-normalized process theory:

$$\max_{\Lambda \in \mathbb{R}^{d_x \times d_\eta}} \mathbf{E} \exp \text{tr} \left( P\Lambda - \frac{1}{2} \Lambda^\top Q \Lambda \right) \leq 1. \quad (36)$$

This inequality is called the canonical assumption because a wide variety of self-normalized processes satisfy it. We will demonstrate that it is satisfied for (34) and (35) in Section IV-B. If we could exchange the order of the maximization with the expectation in (36), then the bound  $\mathbf{E} \exp \left( \frac{1}{2} \|PQ^{-1/2}\|_F^2 \right) \leq 1$  would be satisfied, and the Chernoff bound above would provide a valuable exponential inequality. As this exchange is not possible, we instead lower bound the maximization over  $\Lambda$  by assigning a probability distribution to a random variable  $\Psi$  which takes values in  $\mathbb{R}^{d_x \times d_\eta}$ , and taking the expectation over this distribution. Doing so preserves the inequality (36):

$$\mathbf{E} \mathbf{E} \left[ \exp \text{tr} \left( P\Psi - \frac{1}{2} \Psi^\top Q \Psi \right) \mid \Psi \right] \leq 1.$$

The order of expectation over  $\Psi$  and over the random variables  $P$  and  $Q$  may then be exchanged by an appeal to Fubini's theorem:

$$\begin{aligned} 1 & \geq \mathbf{E} \mathbf{E} \left[ \exp \text{tr} \left( P\Psi - \frac{1}{2} \Psi^\top Q \Psi \right) \mid \Psi \right] \\ & = \mathbf{E} \mathbf{E} \left[ \exp \text{tr} \left( P\Psi - \frac{1}{2} \Psi^\top Q \Psi \right) \mid P, Q \right]. \quad (37) \end{aligned}$$

By selecting the distribution over  $\Psi$  appropriately, the result is a so-called *pseudo-maximization*. In particular, by completing the square, the inner conditional expectation on the right may be expressed as

$$\begin{aligned} & \mathbf{E} \left[ \exp \text{tr} \left( P\Psi - \frac{1}{2} \Psi^\top Q \Psi \right) \mid P, Q \right] \\ & = \exp \text{tr} \left( PQ^{-1}P^\top / 2 \right) \\ & \times \mathbf{E} \left[ \exp \text{tr} \left( -\frac{1}{2} (\Psi - Q^{-1}P^\top)^\top Q (\Psi - Q^{-1}P^\top) \right) \mid P, Q \right]. \end{aligned}$$

For particular choices of the distribution of  $\Psi$ , the right side of the above expression approximates the maximum value,  $\exp \text{tr} \left( PQ^{-1}P^\top / 2 \right)$ , of  $\exp \text{tr} \left( P\Lambda - \frac{1}{2} \Lambda^\top Q \Lambda \right)$ . This allows us to apply a Chernoff argument similar to the one sketched above to obtain an exponential bound on a quantity related to  $\|PQ^{-1/2}\|_F$ . The following lemma demonstrates one such bound that results by selecting the distribution of  $\Psi$  as a matrix normal distribution.

**Lemma IV.1** (Extension of Theorem 14.7 in [18]). *Suppose that (36) is satisfied. Let  $\Sigma$  be a positive definite matrix in  $\mathbb{R}^{d_x \times d_x}$ . Then, for  $\delta > 0$ , with probability at least  $1 - \delta$ ,*

$$\|P(Q + \Sigma)^{-1/2}\|_F^2 \leq 2 \log \left( \frac{\det(Q + \Sigma)^{d_\eta/2} \det(\Sigma)^{-d_\eta/2}}{\delta} \right).$$

### B. Self-Normalized Martingales Satisfy the Canonical Assumption

In order to make use of Lemma IV.1 to bound (34) or (35), we must ensure that the condition (36) holds for

$$P = \sum_{t=1}^T \frac{\eta_t X_t^\top}{\sigma}, \quad Q = \sum_{t=1}^T X_t X_t^\top,$$

where  $\eta_t \in \mathbb{R}^{d_n}$  is either the noise process  $V_t$  or the scalar process  $w^\top V_t$  for some fixed unit vector  $w$ . The following lemma shows that it is sufficient for  $\eta_t$  to be  $\sigma^2$ -conditionally sub-Gaussian.

**Lemma IV.2.** Fix  $T \in \mathbb{N}_+$ . Let  $\{\mathcal{F}_t\}_{t=0}^T$  be a filtration such that  $\{X_t\}_{t=1}^T$  is adapted to  $\{\mathcal{F}_{t-1}\}_{t=1}^T$  and  $\{\eta_t\}_{t=1}^T$  is adapted to  $\{\mathcal{F}_t\}_{t=1}^T$ . Additionally, suppose that for all  $t \geq 1$ ,  $\eta_t$  is  $\sigma^2$ -conditionally sub-Gaussian with respect to  $\mathcal{F}_t$ . Let  $\Lambda \in \mathbb{R}^{d_x \times d_n}$  be arbitrary and consider for  $t \in \{1, \dots, T\}$

$$M_t(\Lambda) \triangleq \exp \operatorname{tr} \left( \sum_{s=1}^t \left[ \frac{\eta_s X_s^\top \Lambda}{\sigma} - \frac{1}{2} \Lambda^\top X_s X_s^\top \Lambda \right] \right).$$

Then  $\mathbf{E}M_T(\Lambda) \leq 1$ .

Synthesizing the results in this section along with a covering argument, yields Theorem IV.1

### C. Notes

**Remark IV.1.** Consider the dimensional dependencies of the Frobenius and operator norm bounds in Theorem IV.1. The leading term in the Frobenius norm bound is  $d_Y$  multiplied by the log det term, which scales with  $d_X \log T$  when the empirical covariance is well-conditioned. In particular, the leading term scales with  $d_X d_Y \log T$ . The factor of  $d_Y$  is no longer present on the log det term for the operator norm. The term therefore scales as  $d_X \log T$  when the empirical covariance matrix is well-conditioned. There is, however, a term  $8d_Y \sigma^2 \log 5$  which results from the covering argument. The operator norm bound therefore scales as  $\max\{d_X \log T, d_Y\}$ .

**Remark IV.2.** Theorem IV.1 holds for a fixed  $T \in \mathbb{N}_+$ , which is sufficient for analyzing the system identification error. In contrast, the self-normalized martingale bound in [1] holds for an arbitrary stopping time and thus uniformly for all  $T \in \mathbb{N}_+$  by a stopping time construction. This uniform bound may be required in some settings, e.g. in error bounds for adaptive control.

## V. SYSTEM IDENTIFICATION

In this section, we analyze well-known linear system identification algorithms that rely on the least squares algorithm. Note that the problem formulation changes with the system parameterization (e.g., state space, ARMAX, etc.). However, a nice property of linear systems is that under certain conditions, we can obtain a linear non-parametric ARX model by regressing the system output to past outputs and inputs. Then, depending on the parameterization, we can recover a

particular realization. In the following, we first review ARX identification, which can be seen as a fundamental building block for many linear system identification algorithms. Then, we analyze identification of Markov parameters in the case of state-space systems. We focus exclusively on the case of single trajectory data.

### A. ARX Systems

Consider an unknown vector autoregressive system with exogenous inputs (ARX)

$$Y_t = \sum_{i=1}^p A_i^* Y_{t-i} + \sum_{j=1}^q B_j^* U_{t-j} + \Sigma_W^{1/2} W_t, \quad (38)$$

where  $Y_t \in \mathbb{R}^{d_Y}$  are the system outputs,  $U_t \in \mathbb{R}^{d_U}$  are the control (exogenous) inputs, and  $W_t \in \mathbb{R}^{d_Y}$  is the normalized process noise with  $\Sigma_W \in \mathbb{R}^{d_Y \times d_Y}$  capturing the (non-normalized) noise covariance. Matrices  $A_i^*$ ,  $i \leq p$  and  $B_j^*$ ,  $j \leq q$  contain the unknown ARX coefficients. For the initial conditions, we assume  $Y_{-1} = \dots = Y_{-p} = 0$ ,  $U_{-1} = \dots = U_{-q} = 0$ .

**Assumption V.1** (System and Noise model). Let the noise covariance  $\Sigma_W \succ 0$  be full rank. Let the normalized process noise  $W_t, t \geq 0$  be independent and identically distributed,  $K^2$ -sub-Gaussian (see Definition II.1), with zero mean and unit covariance  $\mathbf{E}W_t W_t^\top = I_{d_Y}$ . The orders  $p, q$  are known. System (38) is non-explosive, that is, the eigenvalues of matrix

$$\mathcal{A}_{11} \triangleq \begin{bmatrix} A_1^* & A_2^* & \dots & A_{p-1}^* & A_p^* \\ I & 0 & \dots & 0 & 0 \\ 0 & I & \dots & 0 & 0 \\ \vdots & & \ddots & & \vdots \\ 0 & 0 & \dots & I & 0 \end{bmatrix}, \quad (39)$$

lie strictly on or inside the unit circle  $\rho(\mathcal{A}_{11}) \leq 1$ .

The techniques below can provide meaningful finite-sample bounds only when the system is non-explosive. Deriving finite sample guarantees for identifying open-loop, explosively unstable partially-observed systems from single trajectory data is open to the best of our knowledge [31].

In this tutorial, we focus solely on white-noise excitation inputs. For analyses of more advanced experimental designs—a crucial aspect of system identification—we refer to [36] as well as to the classical literature [14].

**Assumption V.2** (White-noise excitation policy). We assume that the control input is generated by a random i.i.d. Gaussian process, that is,  $U_t \sim \mathcal{N}(0, \sigma_u^2 I)$ .

Grouping all covariates into one vector and defining

$$X_t = [Y_{t-1:t-p}^\top \quad U_{t-1:t-q}^\top]^\top, \quad \theta^* = [A_{1:p}^* \quad B_{1:q}^*] \quad (40)$$

we can re-write (38) in terms of (1)

$$Y_t = \theta^* X_t + \Sigma_W^{1/2} W_t,$$

where  $W_t$  is independent from  $X_t$ , but  $X_t$  has the special time-dependent structure induced by (38). Given samples  $(Y_{1:T}, U_{0:T-1})$ , the least-squares estimate is given by

$$\hat{\theta}_T \triangleq \sum_{t=1}^T Y_t X_t^\top \left( \sum_{t=1}^T X_t X_t^\top \right)^\dagger, \quad (41)$$

where we purposely highlight the dependence of the estimate on the number of samples with the subscript  $T$ . Before we present the main result, let us define some quantities which are related to the quality of system estimates. The covariance at time  $t \geq 0$  is defined as

$$\Sigma_t \triangleq \mathbf{E} X_t X_t^\top. \quad (42)$$

It captures the expected richness of the data, i.e., how excited the modes of the system are on average. In particular, the relative excitation of the data compared to the noise magnitude significantly affects the quality of system identification. This motivates the definition of signal-to-noise (SNR) as the ratio of the worst-case excitation over the worst-case noise magnitude

$$\text{SNR}_t \triangleq \frac{\lambda_{\min}(\Sigma_t)}{\|\Sigma_W\|_{\text{op}} K^2}. \quad (43)$$

The following theorem provides a finite-sample upper bound on the performance of the least-square estimator.

**Theorem V.1** (ARX Finite-Sample Bound). *Let  $(Y_{1:T}, U_{0:T-1})$  be single trajectory input-output samples generated by system (38) under Assumptions V.1, V.2 for some horizon  $T$ . Fix a failure probability  $0 < \delta < 1$  and a time index  $\tau \geq \max\{p, q\}$ . Let  $T_{\text{pe}}(\delta, \tau) \triangleq \min\{t : t \geq T_0(t, \delta/3, \tau)\}$ , where  $T_0$  is defined in (46). If  $T \geq T_{\text{pe}}(\delta, \tau)$ , then with probability at least  $1 - \delta$*

$$\|\hat{\theta}_T - \theta^*\|_{\text{op}}^2 \leq \frac{C}{\text{SNR}_\tau T} \left( (pd_Y + qd_U) \log \frac{pd_Y + qd_U}{\delta} + \log \det(\Sigma_T \Sigma_\tau^{-1}) \right), \quad (44)$$

where  $C$  is a universal constant, i.e., it is independent of system, confidence  $\delta$  and index  $\tau$ .

For non-explosive systems, matrix  $\Sigma_T \Sigma_\tau^{-1}$  increases at most polynomially with  $T$  in norm (in view of Lemma V.1). Hence, the identification error decays with a rate of  $\tilde{O}(1/\sqrt{T})$ .

**Dimensional dependence.** Ignoring logarithmic terms, the bound implies that the number of samples  $T$  should scale linearly with the dimension  $pd_Y + qd_U$  of the covariates  $X_t$ . Since every sample  $Y_t$  contains at least  $d_Y$  measurements, this implies that the total number of measurements should be linear with  $d_Y \times (pd_Y + qd_U)$ . This scaling is qualitatively correct since  $\theta^*$  has  $d_Y \times (pd_Y + qd_U)$  unknowns, requiring at least as many independent equations.

**Logarithmic dependence on confidence.** The error norm scales linearly with  $\sqrt{\log 1/\delta}$ . In the asymptotic regime we also recover the same order of  $\sqrt{\log 1/\delta}$  by applying the Central Limit Theorem (CLT). However, in the regime of finite samples, obtaining the rate is non-trivial, see [31], and requires the analysis presented in this tutorial.

**System theoretic constants.** The identification error is directly affected by the SNR of the system. The more the system is excited and the smaller the noise, the better the SNR becomes. However, excitability varies heavily depending on the system and the choice of excitation policy. In particular, the system's controllability structure can affect the degree of excitation dramatically. Systems with poor controllability structure can exhibit SNR which suffers from curse of dimensionality, i.e., the smallest eigenvalue of  $\Sigma_\tau$  degrades exponentially with the system dimension [30].

The upper bound also increases with the logarithm of the "condition number"  $\det(\Sigma_T \Sigma_\tau^{-1})$ . For stable systems, this condition number is bounded since  $\Sigma_T$  converges to a steady-state covariance as  $T$  increases; we can neglect it in this case. On the other hand, the term might be significant in the case of general non-explosive systems. Let  $\kappa$  be the size of the largest Jordan block of  $\mathcal{A}_{11}$  with eigenvalues on the unit circle. Then, this term can be as large as  $\kappa \log T$ .

**Burn-in time.** The upper bound holds as soon as the number of samples exceeds a "burn-in" time  $T_{\text{pe}}(\delta, \tau)$ . If the system is non-explosive,  $T_{\text{pe}}(\delta, \tau)$  is always finite for fixed  $\tau$ . Exceeding the burn-in time guarantees that we have persistency of excitation, that is, all modes of the system are excited. The burn-in time increases as we require more confidence  $\delta$  and we chose larger time indices  $\tau$ . On the other had, larger  $\tau$  leads to larger  $\Sigma_\tau$ , which improves the  $\text{SNR}_\tau$ . In other words, there is a tradeoff between improving the SNR and deteriorating the burn-in time. We analyze persistency of excitation in more detail, in the next subsection.

**Proof outline** We outline the proof below, the full proofs can be found in the online version. To analyze the least squares error, observe that

$$\hat{\theta}_T - \theta^* = \underbrace{\sum_{t=1}^T \Sigma_W^{1/2} W_t X_t^\top}_{\text{noise}} \left( \sum_{t=1}^T X_t X_t^\top \right)^{-1/2} \times \underbrace{\left( \sum_{t=1}^T X_t X_t^\top \right)}_{\text{excitation}}^{-1/2} \quad (45)$$

where we assumed that the inverse exists. To deal with the second term, we will prove persistency of excitation in finite time leveraging the techniques of Section III, which requires most of the work. To deal with the noise part we will apply the self-normalized martingale methods, which are reviewed in Section IV. We study both terms in the following subsections.

1) *Persistency of Excitation in ARX Models:* In this subsection, we leverage the result of Theorem III.1 to prove persistency of excitation. By persistency of excitation, we refer to the case when we have rich input-output data, that is, data which characterize all possible behaviors of the system. Recall the definition (9) of the empirical covariance matrix

$$\hat{\Sigma}_T \triangleq \frac{1}{T} \sum_{t=1}^T X_t X_t^\top.$$

Using this definition, the excitation term in the least squares error can be re-written as  $(T\widehat{\Sigma}_T)^{-1/2}$ . We say that persistency of excitation holds if and only if the empirical covariance matrix is strictly positive definite (full rank). In the following, we show that the full rank condition holds with high probability, provided that the number of samples exceeds a certain threshold, i.e., the burn-in time.

**Theorem V.2** (ARX PE). *Let  $(Y_{1:T}, U_{0:T-1})$  be input-output samples generated by system (38) under Assumptions V.1, V.2 for some fixed horizon  $T$ . Fix a failure probability  $0 < \delta < 1$  and a time index  $\tau \geq \max\{p, q\}$ . Then,  $\lambda_{\min}(\Sigma_\tau) > 0$ . Moreover, if  $T$  is large enough*

$$T \geq T_0(T, \delta, \tau) \triangleq 1152\tau \max\{K^2, 1\} \times \left( (pd_Y + qd_U) \log C_{\text{sys}}(T, \tau) + \log(1/\delta) \right) \quad (46)$$

where

$$C_{\text{sys}}(T, \tau) \triangleq \frac{T}{3\tau} \frac{\|\Sigma_T\|_{\text{op}}^2}{\lambda_{\min}^2(\Sigma_\tau)},$$

then,

$$\mathbf{P} \left( \widehat{\Sigma}_T \succeq \frac{1}{16} \Sigma_\tau \right) \geq 1 - \delta.$$

The detailed proof can be found in the online version, we only sketch the main ideas here. The first step is to express the covariates  $X_{1:T}$  as a causal linear combination of the noises and the inputs, mimicking (23). The evolution of the covariates follows a state-space recursion

$$X_{t+1} = \mathcal{A}X_t + \mathcal{B}V_t, \quad (47)$$

where we concatenate noises and inputs  $V_t \triangleq [W_t^\top \ U_t^\top]^\top$ . Matrices  $\mathcal{A}$ ,  $\mathcal{B}$  are given by

$$\mathcal{A} \triangleq \begin{bmatrix} \mathcal{A}_{11} & \mathcal{A}_{12} \\ 0 & \mathcal{A}_{22} \end{bmatrix}, \quad \mathcal{B} \triangleq [\mathcal{B}_1 \ \mathcal{B}_2], \quad \text{where}$$

$\mathcal{A}_{11}$  is defined in (39),

$$\mathcal{A}_{12} \triangleq \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \otimes [B_1^* \ \dots \ B_q^*] \in \mathbb{R}^{pd_Y \times qd_U}$$

$$\mathcal{A}_{22} \triangleq \begin{bmatrix} 0 & \dots & 0 & 0 & 0 \\ I_{d_U} & \dots & 0 & 0 & 0 \\ \vdots & \ddots & & & \\ 0 & \dots & I_{d_U} & 0 & 0 \\ 0 & \dots & 0 & I_{d_U} & 0 \end{bmatrix} \in \mathbb{R}^{qd_U \times qd_U}$$

$$\mathcal{B}_1 \triangleq \begin{bmatrix} \Sigma_W^{1/2} & 0_{((p-1)d_Y + qd_U) \times d_Y} \end{bmatrix}^\top,$$

$$\mathcal{B}_2 \triangleq \begin{bmatrix} 0_{pd_Y \times d_U} & I_{d_U} & 0_{(q-1)d_U \times d_U} \end{bmatrix}^\top.$$

The vector  $X_{1:T}$  of all covariates satisfies the following causal linear relation

$$X_{1:T} = \underbrace{\begin{bmatrix} \mathcal{B} & 0 & \dots & 0 \\ \mathcal{A}\mathcal{B} & \mathcal{B} & \dots & 0 \\ \vdots & & \ddots & \\ \mathcal{A}^{T-1}\mathcal{B} & \mathcal{A}^{T-2}\mathcal{B} & \dots & \mathcal{B} \end{bmatrix}}_{\mathbf{L}} V_{0:T-1}. \quad (48)$$

where the lower-block triangular matrix is the Toeplitz matrix generated by the Markov parameters matrices  $\mathcal{A}$ ,  $\mathcal{B}$ . The second step is to apply Theorem III.1. The details can be found in the online version.

**Remark V.1** (Existence of burn-in time.). *For the above result to be meaningful, we need inequality (46) to be feasible. For non-explosive systems, the system theoretic term  $\log C_{\text{sys}}(T, \tau)$  increases at most logarithmically with  $T$ , since  $\Sigma_t$  increases polynomially with  $t$  in view of Lemma V.1. Hence, for any fixed  $\tau$ , or, in general, any  $\tau$  that increases mildly (sublinearly) with  $T$ , e.g.  $O(\sqrt{T})$ , it is possible to satisfy (46). Note that  $\rho(\mathcal{A}) = \rho(\mathcal{A}_{11})$  due to the triangular structure of  $\mathcal{A}$ . Hence, by Assumption V.1, system  $\mathcal{A}$  is also non-explosive.*

**Remark V.2** (Unknown system orders  $p, q$ ). *The result of Theorem V.2 still holds if the orders  $p, q$  are unknown and we use the wrong orders  $\hat{p}, \hat{q}$  in the covariates  $X_t$ . We just need to replace  $p, q$  with  $\hat{p}, \hat{q}$  with  $\hat{p}, \hat{q}$  and revise the size of  $\Sigma$  accordingly in (46). The finite-sample bounds of Theorem V.1 also hold (by revising accordingly), but only if we overestimate  $p, q$ , that is  $\hat{p} \geq p, \hat{q} \geq q$ . This also generalizes the single trajectory result of [3] to non-explosive systems.*

The following supporting lemma proves that the  $k$ -th powers of non-explosive matrices increase at most polynomially with  $k$ .

**Lemma V.1** (Lemma 1 in [30]). *Let  $\mathcal{A} \in \mathbb{R}^{d \times d}$  have all eigenvalues inside or on the unit circle, with  $\|\mathcal{A}\|_{\text{op}} \leq M$ , for some  $M > 0$ . Then,*

$$\|\mathcal{A}^k\|_{\text{op}} \leq (ek)^{d-1} \max\{M^d, 1\}. \quad (49)$$

As a corollary, the covariance matrices  $\Sigma_t$  also grow at most polynomially with the time  $t$ .

2) *Dealing with the Noise Term:* In this subsection, we modify the noise term so that we can leverage Theorem IV.1, which cannot be applied directly. We first manipulate the inverse of  $T\widehat{\Sigma}_T$  to relate it to the inverse of  $\Sigma + T\widehat{\Sigma}_T$ , for some carefully selected  $\Sigma$ . Inspired by [24], we leverage the result of Theorem V.2. Under the event that persistency of excitation holds we have  $\widehat{\Sigma}_T \succeq T\Sigma_\tau/16$ . Thus, selecting  $\Sigma = T\Sigma_\tau/16$  guarantees that

$$\left(T\widehat{\Sigma}_T\right)^{-1} \preceq 2 \left(\Sigma + T\widehat{\Sigma}_T\right)^{-1}.$$

We can now apply Theorem IV.1. To finish the proof we need to upper-bound the determinant of  $\log \det(\Sigma + T\widehat{\Sigma}_T)$ . It is sufficient to establish a crude upper-bound on the empirical covariance  $T\widehat{\Sigma}_T$  as in the following lemma.

**Lemma V.2** (Matrix Markov's inequality). *Fix a failure probability  $\delta > 0$ . With probability at least  $1 - \delta$*

$$\widehat{\Sigma}_T \preceq \frac{pd_Y + qd_U}{\delta} \Sigma_T. \quad (50)$$

A more refined upper bound can also be applied (see e.g. the proof of Proposition VI.1 below or the results in [8]).

### B. State-Space Systems

In this subsection, we derive finite-sample guarantees for learning Markov parameters of linear systems in state-space form. Consider the following state-space system in the so-called innovation form:

$$\begin{aligned} X_{t+1} &= A^* X_t + B^* U_t + F^* \Sigma_E^{1/2} E_t \\ Y_t &= C^* X_t + \Sigma_E^{1/2} E_t, \end{aligned} \quad (51)$$

where  $A^* \in \mathbb{R}^{d_x \times d_x}$ ,  $B^* \in \mathbb{R}^{d_x \times d_u}$ ,  $F^* \in \mathbb{R}^{d_x \times d_y}$ , and  $C^* \in \mathbb{R}^{d_y \times d_x}$  are *unknown* state-space parameters. For the initial condition, we assume  $X_0 = 0$ . We call the normalized noise process  $E_t$  the innovation error process. Similar to the ARX case, we focus on white-noise excitation inputs, namely Assumption V.2 also holds here. Moreover, we assume the following.

**Assumption V.3** (System and Noise model). *Let the noise covariance  $\Sigma_E \succ 0$  be full rank. Let the normalized innovation process  $E_t$  be independent, identically distributed,  $K^2$ -sub-Gaussian (see Definition II.1), with zero mean and unit covariance  $\mathbb{E} E_t E_t^\top = I_{d_y}$ . The order  $d_x$  is unknown. System (51) is non-explosive, that is, the eigenvalues of matrix  $A^*$  lie strictly on or inside the unit circle  $\rho(A) \leq 1$ . The system is also minimum-phase, i.e., the closed loop matrix*

$$A_{cl}^* \triangleq A^* - F^* C^* \quad (52)$$

*has all eigenvalues inside the unit circle  $\rho(A_{cl}^*) < 1$ .*

The innovation form (51) might seem puzzling at first. In particular, the correlation between process and measurement noise via  $F^*$ , and the requirement  $\rho(A_{cl}^*) < 1$  seem restrictive. However, the representation (51) is standard in the system identification literature [34]. Moreover, as we review below, standard state-space models have input-output second-order statistics, which are equivalent to the ones generated by system (51) (for appropriate  $F^*$ ,  $\Sigma_E$ ).

**Remark V.3** (Generality of model). *System class (51) captures general state-space systems driven by Gaussian noise. Consider the following state-space model*

$$\begin{aligned} S_{t+1} &= A^* S_t + B^* U_t + W_t \\ Y_t &= C^* S_t + V_t, \end{aligned} \quad (53)$$

*where  $W_t, V_t$  are i.i.d., independent of each other, mean-zero Gaussian, with covariances  $\Sigma_W$  and  $\Sigma_V$  respectively. Assume that  $\Sigma_V \succ 0$  is full rank, the pair  $(C^*, A^*)$  is detectable, and the pair  $(A^*, \Sigma_W)$  is stabilizable. These three assumptions*

*imply that the Kalman filter of system (53) is well-defined [2]. In particular, define the Riccati operator as*

$$\begin{aligned} \text{RIC}(P) &\triangleq A^* P (A^*)^\top + \Sigma_W \\ &\quad - A^* P (C^*)^\top (C^* P (C^*)^\top + \Sigma_V)^{-1} C^* P (A^*)^\top \end{aligned} \quad (54)$$

*and let  $P^*$  be the unique positive semidefinite solution of  $P^* = \text{RIC}(P^*)$ . Then the Kalman filter gain is equal to*

$$F^* = -A^* P (C^*)^\top (C^* P (C^*)^\top + \Sigma_V)^{-1}. \quad (55)$$

*Assume that the initial state is also mean-zero Gaussian with covariance  $P^*$  and independent of the noises. Finally set*

$$\Sigma_E = C^* P^* (C^*)^\top + \Sigma_V. \quad (56)$$

*Under the above assumptions and selection of  $F^*$ ,  $\Sigma_E$  systems (51) and (53) are statistically equivalent from an input-output perspective, see [19]. Both system descriptions lead to input-output trajectories with identical statistics. Moreover, due to the properties of Kalman filter, stability of  $A_{cl}^*$  (minimum phase property) and independence of  $E_t$  are satisfied automatically [2].*

In this tutorial we will only focus on recovering the first few (logarithmic in  $T$ -many) Markov parameters  $C^* (A_{cl}^*)^i B^*$ ,  $i \geq 0$  and  $C^* (A_{cl}^*)^j F^*$ ,  $j \geq 0$  of system (51). From a learning theory point of view, this is also known as improper learning, since the search space (finitely many Markov parameters) does not exactly, but only approximately, coincide with the hypothesis class (state space models). In principle, this forms the backbone of the SSARX method introduced by Jansson [7]. One can then proceed to recover the original state-space parameters (up to similarity transformation) from the Markov parameters by employing some realization method. We refer to [17, 31] for a discussion on this approach from a finite sample perspective.

1) *Reduction to ARX Learning with Bias:* Let  $p > 0$  be a past horizon. Denote the Markov parameters up to time  $p$  by

$$\theta_p^* \triangleq \begin{bmatrix} C^* B^* & \dots & C^* (A_{cl}^*)^{p-1} B^* \\ C^* F^* & \dots & C^* (A_{cl}^*)^{p-1} F^* \end{bmatrix}. \quad (57)$$

Note that the innovation errors are equal to  $\Sigma_E^{1/2} E_t = Y_t - C^* X_t$ . Replacing this expression into the state equation (51), we obtain

$$X_t = A_{cl}^* X_{t-1} + B^* U_{t-1} + F^* Y_{t-1}.$$

Unrolling the state equation  $p$  times, we get

$$Y_t = \underbrace{\theta_p^* Z_t}_{\text{ARX}} + \underbrace{C^* (A_{cl}^*)^p X_{t-p}}_{\text{bias}}, \quad (58)$$

where  $Z_t$  includes the past  $p$  covariates

$$Z_t = \begin{bmatrix} Y_{t-1:t-p}^\top & U_{t-1:t-p}^\top \end{bmatrix}^\top. \quad (59)$$

The above recursion is an approximate ARX equation. There is an additive bias error term on top of the statistical noise. The least-squares solution is given by

$$\hat{\theta}_{p,T} \triangleq \sum_{t=1}^T Y_t Z_t^\top \left( \sum_{t=1}^T Z_t Z_t^\top \right)^\dagger, \quad (60)$$

where we also highlight the dependence on the past  $p$ . By the minimum phase assumption, the bias term decays exponentially with the past horizon  $p$ . This follows from the fact that  $A_{cl}^*$  is asymptotically stable, while  $X_t$  scales at most polynomially with  $t$  (in view of Lemma V.1). By selecting  $p = \Omega(\log T)$ , we can make the bias term decay very fast, making its contribution to the error  $\theta_p^* - \hat{\theta}_{p,T}$  negligible. On the other hand, increasing the past horizon  $p$  increases the statistical error since the search space is larger.

2) *Finite-Sample Guarantees*: To derive a finite-sample rate for state space systems of the form (51), we follow the same steps as in the case of ARX systems. However, we have to account for the bias term and the fact that  $p$  grows with  $\log T$ . Let us define again the covariance at time  $t \geq 0$

$$\Sigma_{p,t} \triangleq \mathbf{E} Z_t Z_t^\top, \quad (61)$$

where we highlight the dependence on both the past horizon  $p$  and the time  $t$ . The covariance of the state is defined similarly

$$\Sigma_{X,t} \triangleq \mathbf{E} X_t X_t^\top. \quad (62)$$

Define the SNR as

$$\text{SNR}_{p,t} \triangleq \frac{\lambda_{\min}(\Sigma_{p,t})}{\|\Sigma_E\|_{\text{op}} K^2}. \quad (63)$$

Unlike the ARX case, here the SNR might degrade since we allow  $p$  to grow with  $\log T$ . For this reason, we require the following additional assumption.

**Assumption V.4** (Non-degenerate SNR). *We assume that the SNR is uniformly lower bounded for all possible past horizons*

$$\liminf_{t \geq 0} \text{SNR}_{t,t} > 0.$$

Later on, in Theorem V.4, we show that the above condition is non-vacuous and is satisfied for quite general systems.

**Theorem V.3** (State Space Finite-Sample Bound). *Let  $(Y_{1:T}, U_{0:T-1})$  be single trajectory input-output samples generated by system (51) under Assumptions V.2, V.3, V.4, for some horizon  $T$ . Fix a failure probability  $0 < \delta < 1$  and select  $p = \beta \log T$ , for  $\beta$  large enough such that*

$$\|C^*(A_{cl}^*)^p\|_{\text{op}} \|\Sigma_{X,T}\|_{\text{op}} \leq T^{-3}. \quad (64)$$

Let  $T_{\text{pe}}^{\text{ss}}(\delta, \beta) \triangleq \min\{t : t \geq T_0(t, \delta, \beta \log t)\}$ , where  $T_0$  is defined in (46). If  $T \geq T_{\text{pe}}^{\text{ss}}(\delta, \beta)$ , then with probability at least  $1 - 2\delta$

$$\|\hat{\theta}_{p,T} - \theta_p^*\|_{\text{op}}^2 \leq \frac{C_1}{\text{SNR}_{p,p} T} \left( p(d_Y + d_U) \log \frac{p(d_Y + d_U)}{\delta} + \log \det(\Sigma_{p,T} \Sigma_{p,p}^{-1}) \right), \quad (65)$$

where  $C_1$  is a universal constant, i.e., it is independent of system, confidence  $\delta$  and past horizon  $p$ .

For non-explosive systems, matrix  $\Sigma_{p,T} \Sigma_{p,p}^{-1}$  increases at most polynomially with  $T$  in norm. Since the SNR is uniformly lower bounded, the identification error decays with a rate of  $\tilde{O}(1/\sqrt{T})$ . The bound seems similar to the one for ARX systems for  $\tau = p$ . However, since  $p = \theta(\log T)$ , we have an extra logarithmic term.

**Role of  $\beta$ .** Recall the approximate ARX relation (58). For the bias term to be small, the exponentially decaying  $(A_{cl}^*)^p$  should counteract the magnitude of the state  $\|X_{t-p}\|_{\text{op}}$ . Intuitively, the state grows as fast as  $\|\Sigma_{X,t}\|_{\text{op}}^{1/2}$ , where  $\Sigma_{X,t} = \mathbf{E} X_t X_t^\top$ . Hence the state norm grows at most polynomially with  $T$ . Meanwhile,  $\|(A_{cl}^*)^p\|_{\text{op}} = O(\rho^p)$  for some  $\rho > \rho(A_{cl}^*)$ . With the choice  $p = \beta \log T$ , we get  $\|(A_{cl}^*)^p\|_{\text{op}} = O(T^{-\beta/\log(1/\rho)})$ . Hence, if we select large enough  $\beta$ , we can make the bias term very small, even smaller than the dominant  $\tilde{O}(1/\sqrt{T})$  term.

**Burn-in time.** Since the system is non-explosive,  $T_{\text{pe}}^{\text{ss}}(\delta, \beta)$  is always finite under Assumption V.4, for any  $\beta$ . As before, exceeding the burn-in time guarantees that we have persistency of excitation. Naturally, larger  $\beta$  lead to larger past horizons  $p$ , which, in turn, increase the burn-in time.

Finally, we prove that Assumption V.4 is non-vacuous. It is sufficient for  $F^*$  and  $\Sigma_W$  to be generated by a Kalman filter as in (55), (56).

**Theorem V.4.** *Consider system (51) and the definition of  $\text{SNR}_{p,t}$  in (63). If matrices  $F^*$ ,  $\Sigma_E$  are generated as in (55), (56) with  $(A^*, \Sigma_W^{1/2})$  stabilizable,  $(C^*, A^*)$  detectable and  $\Sigma_V \succ 0$ , then the SNR is uniformly lower bounded  $\liminf_{t \geq 0} \text{SNR}_{t,t} > 0$ .*

Both conditions are sufficient. It is subject of future work to extend the result to more general non-explosive systems.

## C. Notes

The exposition above is inspired by prior work on identifying fully-observed systems [4, 27, 24] and partially-observed systems [16, 28, 25, 29, 11, 10, 12]. For a wider overview of the literature, we refer the reader to [31].

Let us further remark that the guarantee for the ARX model in Theorem V.1 is almost optimal. The use of Matrix Markov's inequality yields extraneous dependency on the problem dimension multiplying the deviation term  $\log(1/\delta)$ . This can in principle be removed by a more refined analysis (see e.g. the proof of Proposition VI.1 below or the results in [8]). The question of optimality in identifying partially observed state-space systems is more subtle, and while consistent, the bounds presented here are not (asymptotically) optimal.

## VI. AN ALTERNATIVE VIEWPOINT: THE BASIC INEQUALITY

In many situations, the choice of the model class  $\mathcal{M} = \mathbb{R}^{d_Y \times d_X}$  leading to (8) is not appropriate. For instance physical or other modelling considerations might have already informed

us that the true  $\theta^*$  belongs to some smaller model class such as the family of low rank or sparse matrices which are strict subsets of  $M$ . Other properties one might wish to enforce include, stable, low norm, or even passivity-type properties. In either of the above examples no error expression of the form (10) is directly available. Instead, we observe by optimality of  $\widehat{M}$  to the optimization program (7) that

$$\frac{1}{T} \sum_{t=1}^T \|Y_t - \widehat{\theta} X_t\|_2^2 \leq \frac{1}{T} \sum_{t=1}^T \|Y_t - \theta^* X_t\|_2^2. \quad (66)$$

Expanding the squares and re-arranging terms we arrive at the so-called basic inequality of least squares:

$$\frac{1}{T} \sum_{t=1}^T \|(\widehat{\theta} - \theta^*) X_t\|_2^2 \leq \frac{2}{T} \sum_{t=1}^T \langle V_t, (\widehat{\theta} - \theta^*) X_t \rangle. \quad (67)$$

The inequality (67) serves as an alternative to the explicit error equation (10). To drive home this point, let us first re-arrange (67) slightly:

$$\begin{aligned} & \frac{1}{T} \sum_{t=1}^T \|(\widehat{\theta} - \theta^*) X_t\|_2^2 \\ & \leq \frac{4}{T} \sum_{t=1}^T \langle V_t, (\widehat{\theta} - \theta^*) X_t \rangle - \frac{1}{T} \sum_{t=1}^T \|(\widehat{\theta} - \theta^*) X_t\|_2^2. \end{aligned} \quad (68)$$

Note now that  $\widehat{\theta} - \theta^*$  are elements of  $M_* \triangleq M - \theta^*$ . Hence—by considering the worst-case (supremum) right hand side of (68)—we obtain:

$$\begin{aligned} & \frac{1}{T} \sum_{t=1}^T \|(\widehat{\theta} - \theta^*) X_t\|_2^2 \\ & \leq \sup_{\theta \in M_*} \left\{ \frac{4}{T} \sum_{t=1}^T \langle V_t, \theta X_t \rangle - \frac{1}{T} \sum_{t=1}^T \|\theta X_t\|_2^2 \right\}. \end{aligned} \quad (69)$$

In fact, if  $M = \mathbb{R}^{d_Y \times d_X}$ , the optimization on the right hand side of (69) has an explicit solution. This implies that we always have the following upper-bound on the event that the design is nondegenerate:

$$\begin{aligned} & \sup_{\theta \in M_*} \left\{ \frac{4}{T} \sum_{t=1}^T \langle V_t, \theta X_t \rangle - \frac{1}{T} \sum_{t=1}^T \|\theta X_t\|_2^2 \right\} \\ & \leq \sup_{\theta \in \mathbb{R}^{d_Y \times d_X}} \left\{ \frac{4}{T} \sum_{t=1}^T \langle V_t, \theta X_t \rangle - \frac{1}{T} \sum_{t=1}^T \|\theta X_t\|_2^2 \right\} \\ & = \frac{4}{T} \left\| \left( \sum_{t=1}^T V_t X_t^T \right) \left( \sum_{t=1}^T X_t X_t^T \right)^{-1/2} \right\|_F^2. \end{aligned} \quad (\text{opt.}) \quad (70)$$

Hence, we have in principle recovered an in-norm version of (10) with slightly worse constants. Put differently, we may regard (69) as a variational (or dual) form of the explicit error (10). Now, the advantage of (69) is twofold:

- (69) and (70) hold for any  $M_* \subset \mathbb{R}^{d_Y \times d_X}$  and hence allows us to analyze the LSE (7) beyond OLS ( $M_* =$

$\mathbb{R}^{d_Y \times d_X}$ ). This is important in identification problems where the parameter space is restricted.

- We do not have to rely on (70) to control (69). In fact, for many reasonable classes of  $M_* \subset \mathbb{R}^{d_Y \times d_X}$  we are able to give alternative arguments that are much sharper (in terms of e.g. dimensional scaling) than the naive bound (70). See Section VI-A below.

A third advantage of the variational form (69) is that it generalizes straightforwardly beyond linear least squares. In fact, none of the steps (66), (67), (68) and (69) relied on the linearity of  $x \mapsto \widehat{\theta} x$  or that of  $x \mapsto \theta^* x$  ( $x \in \mathbb{R}^{d_X}$ ). We will explore this theme further in Section VI-A and Section VII.

#### A. Sparse Autoregressions

Before we proceed to sketch out how the basic inequality above extends to nonlinear problems in Section VII, let us use it to analyze a simple variation of the autoregression already encountered in Section V. Namely, the autoregressive model (38) which—for simplicity—is further assumed one-dimensional:

$$Y_t = \sum_{i=1}^p A_i^* Y_{t-i} + W_t \quad (71)$$

and assume in addition that it is known that only  $s \in \mathbb{N}$  of the  $p$  entries of  $\theta^* = [A_1^*, \dots, A_p^*]$  are nonzero. Put differently, the vector  $\theta^*$  is known to be  $s$ -sparse and we write  $\theta^* \in \{\theta \in \mathbb{R}^p : \|\theta\|_0 \leq s\} \triangleq M$ . Hence, in this case the model class  $M$  is the union of  $\binom{p}{s}$  subspaces. Clearly, we could use OLS (8) but this estimator does not take advantage of the additional information that  $A^* = \theta^*$  lies in the  $s$ -dimensional submanifold  $M$ . Intuitively, if  $s \ll p$  this set should be much smaller than  $\mathbb{R}^p$  and so one expects that identification occurs at a faster rate.

In this section we demonstrate that the least squares estimator (7) in which the search is restricted to the low-dimensional manifold  $M$  outperforms the OLS. We stress that this is *not* a computationally efficient estimator and the results in this section should be thought of as an illustration of a proof technique.

Returning to the problem of controlling the error of this estimator, we note that in this case there is no closed form for the LSE and we do not have direct access to the error equation (10).<sup>5</sup> Hence, we instead use the offset basic inequality approach from Section VI. As before, it is convenient to set  $X_t = [Y_{t-1}, \dots, Y_{t-p}]^T$ . With this additional bit of notation in place, we recall from (69) that:

$$\begin{aligned} & \frac{1}{T} \sum_{t=1}^T \|(\widehat{\theta} - \theta^*) X_t\|_2^2 \\ & \leq \max_{\theta \in M_*} \left\{ \frac{4}{T} \sum_{t=1}^T W_t \theta X_t - \frac{1}{T} \sum_{t=1}^T \|\theta X_t\|_2^2 \right\} \end{aligned}$$

where  $M_*$  is the translation  $M - \theta^*$ . Since  $M$  is the union of  $\binom{p}{s}$ -many linear  $s$ -dimensional subspaces  $S \subset \mathbb{R}^{d_X \times d_X}$ ,

<sup>5</sup>Although, in this particular case an alternative analysis based on this equation is possible.

$M_*$  is the union of  $\binom{p}{s}$  affine subspaces  $s$ -dimensional affine subspaces of the form  $S - \theta^*$ . Let us also note that  $M_* \subset M - M = \{\theta \in \mathbb{R}^p : \|\theta\|_0 \leq 2s\}$ . Consequently:

$$\begin{aligned} & \frac{1}{T} \sum_{t=1}^T \|(\hat{\theta} - \theta^*)X_t\|_2^2 \\ & \leq \max_{\theta \in M_*} \left\{ \frac{4}{T} \sum_{t=1}^T W_t \theta X_t - \frac{1}{T} \sum_{t=1}^T |\theta X_t|_2^2 \right\} \\ & \leq \max_S \max_{\theta \in S} \left\{ \frac{4}{T} \sum_{t=1}^T W_t \theta X_t - \frac{1}{T} \sum_{t=1}^T |\theta X_t|_2^2 \right\}. \end{aligned} \quad (72)$$

where maximization over  $S$  occurs over the  $\binom{p}{2s}$ -many sparse subspaces. Notice now that since  $\theta$  in (72) is  $s$ -sparse, the products  $\theta X_t$  are just  $\theta X_t = \sum_{i \in S} \theta_i (X_t)_i$  where we have abused notation and identified  $S$  with its support set. Hence, by the same direct calculation as in (70), if we denote  $(X_t)_S$  the  $s$ -dimensional vector obtained by coordinate projection onto part of  $S$  not constrained to be identically zero (i.e. the image of the projection onto  $S$  represented as the  $s$ -dimensional Euclidean space) we find that:

$$\begin{aligned} & \frac{1}{T} \sum_{t=1}^T \|(\hat{\theta} - \theta^*)X_t\|_2^2 \\ & \leq \frac{4}{T} \max_S \left\| \left( \sum_{t=1}^T W_t (X_t)_S \right) \left( \sum_{t=1}^T (X_t)_S (X_t)_S^\top \right)^{-1/2} \right\|_2^2. \end{aligned} \quad (73)$$

The right hand side of (73) can be controlled by the self-normalized inequality in Theorem IV.1 for each fixed  $S$ . Moreover, there are only  $\binom{p}{2s}$  such subspaces, so we can apply a union bound to control the maximum over these subspaces. Note also that the left hand side of (73) can be controlled by the tools developed in Section III. Carrying out these steps leads to the following guarantee.

**Proposition VI.1.** Fix  $T, k \in \mathbb{N}$  with  $T$  divisible by  $k$  and let  $\mathbf{L}$  be the linear operator defined in (48). Let  $\hat{\theta}$  be the LSE (7) over the set  $M = \{\theta \in \mathbb{R}^p : \|\theta\|_0 \leq s\}$  for the system (71). Define  $\Sigma_j \triangleq \frac{1}{j} \sum_{t=1}^j \mathbf{E} X_t X_t^\top$  for  $j \in [T]$  and

$$\text{cond}_{\text{sys}}(T, k) \triangleq \left( 1 + \frac{\|\mathbf{L}\mathbf{L}^\top\|_{\text{op}}}{k\lambda_{\min}(\Sigma_T)} \right) \frac{\lambda_{\max}(\Sigma_T)}{\lambda_{\min}(\Sigma_k)}.$$

There exist universal positive constants  $c, c'$  such that for any  $\delta \in (0, 1)$  it holds with probability at least  $1 - \delta$  that:

$$\|(\hat{\theta} - \theta^*)\sqrt{\Sigma_k}\|_2^2 \leq c\sigma^2 \times \frac{s \log\left(\frac{p \times \text{cond}_{\text{sys}}(T, k)}{s}\right) + \log(1/\delta)}{T} \quad (74)$$

as long as

$$T/k \geq c'\sigma^2 (s [\log(\text{cond}_{\text{sys}}(T, k)) + \log(p/s)] + \log(1/\delta)). \quad (75)$$

A few remarks are in order. The guarantee (74) depends on the dimension  $s$  of  $M$ , and not the total parameter dimension  $p$ . Similarly, the burn-in (75) exhibits a similar win, by depending linearly on  $s$  and only logarithmically on  $p$ . There is also the difference that the left hand side of (74) is given in the problem-dependent Mahalanobis norm induced by  $\Sigma_k$  and opposed to just the standard Euclidean 2-norm. This implies that if we actually want parameter identification in the sense of the previous section, a restricted eigenvalue condition on  $\Sigma_k$  is needed.<sup>6</sup> Indeed, for some positive number  $\lambda_{\text{restricted}}$ , one requires that  $v^\top \Sigma_k v \geq \lambda_{\text{restricted}}$  for all  $2s$ -sparse vectors  $v$  on the unit sphere:  $v \in \mathbb{S}^{p-1}$  and  $\|v\|_0 \leq 3s$ . Obviously the requirements on  $\lambda_{\text{restricted}}$  are much milder than the corresponding ones on  $\lambda_{\min}(\Sigma_k)$  and we always have  $\lambda_{\text{restricted}} \geq \lambda_{\min}(\Sigma_k)$ .

The following lemma is central. Namely, we begin the proof of Proposition VI.1 by restricting to an event in which the designs  $\sum_{t=1}^T (X_t)_S (X_t)_S^\top$  are sufficiently well-conditioned for all the subspaces  $S$  at once. The requirements on this event are relatively milder than the corresponding one over  $\mathbb{R}^p$  and explains the "dimensional win" (when  $s \ll p$ ) of the sparse estimator over OLS.

**Lemma VI.1.** Let  $\mathbf{L}$  be the linear operator defined in (48). Fix  $\delta \in (0, 1)$  and let  $T$  be divisible by  $k \in \mathbb{N}$ . There exist universal positive constants  $c_1, c_2, c_3 \in \mathbb{R}$  such that the following two-sided control holds uniformly in  $S$  with probability  $1 - \delta$ :

$$\begin{aligned} & \frac{c_1}{k} \sum_{t=1}^k \mathbf{E} [(X_t)_S (X_t)_S^\top] \preceq \frac{1}{T} \sum_{t=1}^T (X_t)_S (X_t)_S^\top \\ & \preceq c_2 \left( 1 + \frac{T \|\mathbf{L}\mathbf{L}^\top\|_{\text{op}}}{k \lambda_{\min} \left( \sum_{t=0}^{T-1} \mathbf{E} X_t X_t^\top \right)} \right) \\ & \quad \times \left( \frac{1}{T} \sum_{t=1}^T \mathbf{E} [(X_t)_S (X_t)_S^\top] \right) \end{aligned} \quad (76)$$

as long as

$$T \geq c_3 K^2 (s [\log C_{\text{sys}} + \log(p/s)] + \log(1/\delta)). \quad (77)$$

Equation (77) is revealing about the advantage of using the sparse estimator searching over  $M = \{\theta \in \mathbb{R}^p : \|\theta\|_0 \leq s\}$ . The burn-in period in (77) is proportional to the dimension of the low-dimensional parameter manifold  $M$  instead of that of the latent space  $\mathbb{R}^p$ . Finally, as usual we have relegated the full proof of Proposition VI.1 to the appendix, see ??.

## B. Notes

The variational formulation of the least squares error—the basic inequality (67)—is standard in the nonparametric statistics literature [see e.g. 37, Chapters 13 and 14]. The idea to rewrite the basic inequality (67) as (68) was introduced to the statistical literature by [13], but has its roots in online learning [20].

<sup>6</sup>Note that  $\hat{\theta} - \theta^*$  is  $2s$ -sparse.



## VII. BEYOND LINEAR MODELS

Let us now make another gradual shift of perspective. Instead of considering the linear model (1) introduced in Section I-A we consider the following *nonlinear* regression model:

$$Y_t = f^*(X_t) + V_t, \quad t \in [T]. \quad (78)$$

As before,  $Y_{1:T}, X_{1:T}$  and  $V_{1:T}$  are stochastic processes taking values in  $\mathbb{R}^{d_Y}$  and  $\mathbb{R}^{d_X}$  respectively. However, this time  $f^*$  is no longer constrained to be a linear map of the form  $x \mapsto Ax$  for matrix  $A$ . Rather, we suppose that  $f^*$  in (78) belongs to some (square integrable) space of functions  $\mathcal{F}$  such that  $\mathcal{F} \ni f : x \mapsto f(x)$ . It is perhaps now that the motivation behind the change of perspective from Section VI becomes most apparent: the basic inequality (68) remains valid. To be precise, let us define the *nonparametric* least squares estimator

$$\hat{f} \in \operatorname{argmin}_{f \in \mathcal{F}} \left\{ \frac{1}{T} \sum_{t=1}^T \|Y_t - f(X_t)\|_2^2 \right\}. \quad (79)$$

Let  $\mathcal{F}_* \triangleq \mathcal{F} - f^*$ . By the exact same optimality argument as in Section VI, the reader can now readily verify that:

$$\begin{aligned} & \frac{1}{T} \sum_{t=1}^T \|\hat{f}(X_t) - f^*(X_t)\|_2^2 \\ & \leq \sup_{f \in \mathcal{F}_*} \frac{1}{T} \left( \sum_{t=1}^T 4\langle V_t, f(X_t) \rangle - \sum_{t=1}^T \|f(X_t)\|_2^2 \right). \end{aligned} \quad (80)$$

What does (80) entail in terms of estimating the unknown function  $f^*$ ? To answer this, we first need to define a performance criterion. The simplest one is small average  $L^2$ -norm-error, where

$$f \in \mathcal{F} : \|f\|_{L^2}^2 \triangleq \frac{1}{T} \sum_{t=1}^T \mathbf{E} \|f(X_t)\|_2^2. \quad (81)$$

The program we have carried out in the previous sections now generalizes as follows:

- First, prove a so-called lower uniform law. That is to say, we wish to show that with overwhelming probability

$$\|f\|_{L^2}^2 \leq \frac{C}{T} \sum_{t=1}^T \|f(X_t)\|_2^2 \quad (\text{simultaneously } \forall f \in \mathcal{F}_*). \quad (82)$$

for some universal positive constant  $C$ .

- Second, control the supremum of the *empirical process*:

$$f \mapsto \left( \sum_{t=1}^T 4\langle V_t, f(X_t) \rangle - \sum_{t=1}^T \|f(X_t)\|_2^2 \right) \quad (83)$$

in terms of the noise level  $\sigma$  and some complexity measure  $\operatorname{comp}(\mathcal{F})$ .

By combining (82) and (83) we arrive at a high probability bound of the form:

$$\begin{aligned} \|\hat{f} - f^*\|_{L^2}^2 & \leq \frac{C}{T} \sum_{t=1}^T \|f(X_t)\|_2^2 \\ & \leq \frac{C \times \operatorname{comp}(\mathcal{F}, \sigma^2) + \text{deviation term}}{T}. \end{aligned} \quad (84)$$

A statement of this form is given as Theorem VII.1 below.

**Remark VII.1.** *It is worth to take pause and appreciate the analogy to the analysis of linear regression models. The first step (82) exactly corresponds to controlling the lower spectrum of the empirical covariance matrix. Suppose for simplicity that  $d_Y = 1$ . Then for a linear map  $\mathbb{S}^{d_X-1} \ni f \mapsto \langle f, x \rangle$  we have:*

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \|f(X_t)\|_2^2 & = \frac{1}{T} \sum_{t=1}^T \langle f, (X_t X_t^\top) f \rangle \\ & = \left\langle f, \left[ \frac{1}{T} \sum_{t=1}^T (X_t X_t^\top) \right] f \right\rangle \end{aligned} \quad (85)$$

which are just the one-dimensional projections of the empirical covariance matrix (9). In the context of linear models, establishing (82) was the topic of Section III. Analogously, for a linear predictor, the  $L^2$ -norm (81) becomes a Mahalanobis norm:  $f \in \mathbb{R}^{d_X} \Rightarrow \|f\|_{L^2}^2 = \langle f, \Sigma_X f \rangle$  for some  $\Sigma_X = \frac{1}{T} \sum_{t=1}^T \mathbf{E} X_t X_t^\top$ .

Moreover, For linear models, we had:

$$\begin{aligned} & \sup_{f \in \mathbb{R}^{d_X}} \left( \sum_{t=1}^T 4\langle V_t, f(X_t) \rangle - \sum_{t=1}^T \|f(X_t)\|_2^2 \right) \\ & = 4 \left\| \left( \sum_{t=1}^T V_t X_t^\top \right) \left( \sum_{t=1}^T X_t X_t^\top \right)^{-1/2} \right\|_F^2. \end{aligned} \quad (86)$$

Analyzing terms of this form was the topic of Section IV.

In other words, the approach outlined above is very much in the same spirit as that in the rest of the manuscript. There are a few changes that need to be made since we less access to linearity in our argument, but in principle the key difference is that we will have to replace the indexing set  $\mathbb{S}^{d-1}$  with a more general function class  $\mathcal{F}_*$ .

### A. Many Trajectories and Finite Hypothesis Classes

In order to make the exposition self-contained, we will now make two simplifying assumptions relating to the finiteness of the hypothesis class and the dependence structure of the covariate process  $X_{1:T}$ . A more general treatment without these can be found in [39]. Here, we impose the following:

- A1. The hypothesis class  $\mathcal{F}$  is finite.
- A2. We have access to  $T/k$ -many independent trajectories from the same process: there exists an integer  $k \in \mathbb{N}$  dividing  $T$  such that  $X_{1:k}, X_{k+1:2k}, \dots$  are drawn iid.

We will also impose the following rather minimal integrability condition:

A3. All functions  $f \in \mathcal{F}$  are such that  $\mathbf{E}\|f(X_t)\|_2^4 < \infty$  for all  $t \in [T]$ .

Moreover, as in Section V, we require the noise to be a sub-Gaussian martingale difference sequence:

A4. For each  $t \in [T]$ ,  $V_t|X_{1:t}$  is  $\sigma^2$  conditionally-sub-Gaussian and mean zero.

Under these assumptions, the main result of [39] essentially simplifies to the following theorem.

**Theorem VII.1.** *Impose A1-A4, fix  $\delta \in (0, 1)$  and define*

$$\text{cond}_{\mathcal{F}} \triangleq \max_{f \in \mathcal{F}_*} \max_{t \in T} \frac{\sqrt{\mathbf{E}\|f(X_t)\|_2^4}}{\mathbf{E}\|f(X_t)\|_2^2}. \quad (87)$$

Suppose further that

$$T/k \geq 4\text{cond}_{\mathcal{F}}^2 (\log |\mathcal{F}| + \log(2/\delta))$$

then we have that:

$$\|\hat{f} - f^*\|_{L^2}^2 \leq 16\sigma^2 \left( \frac{\log(|\mathcal{F}|) + \log(2/\delta)}{T} \right). \quad (88)$$

A few remarks are in order. The structure of Theorem VII.1 is by now familiar and it is very much of the same structure as our previous results, cf. (5). The key differences are that: (1) we now control the  $L^2$  norm of our estimator instead of the Euclidean or spectral norm; and (2) that the dimensional dependency has been replaced by the complexity term  $\log |\mathcal{F}_*|$ . The proof is also structurally similar, as noted in Remark VII.1. We also caution the reader that (88) is strictly a statistical guarantee; we have said nothing—and will say nothing more—about the computational feasibility of the estimator (79).

Let us now discuss A1-A4. Assumption A1 informs us that the search space for the LSE (79) is finite. This is mainly imposed to avoid the introduction of the chaining technique which is the standard alternative to the bounds from Section IV. Using this technique, similar statements can for instance be derived for compact subsets of bounded function classes [39]. Assumption A2 controls the dependence structure of the process. Here, we assume that we are able to restart the process every  $k$  time steps. Again, a more general statement relying on stochastic stability can be found in [39]. Assumption A3 is relatively standard. Arguably the strongest assumption is A4, which in principle yields that the conditional expectation (given all past data) is a function in the search space  $\mathcal{F}$ . It is a so-called realizability assumption—the model (78) is well-specified—and it is not currently known how to remove it and still obtain sharp bounds beyond linear classes [for an analysis of linear misspecified models, see 40].

## B. Notes

As noted in the previous section, the idea of using the “offset” basic inequality relied on here is due to [20, 13]. The “many trajectories”-style of analysis used here is due to [32] who introduced it in the linear setting. Here, we have extended their style of analysis to simplify the exposition of [39] who consider the single trajectory setting, but rely on

a rather more advanced exponential inequality due to [23]. Note however that all the analyses above and in this section necessitate some degree of stability (mixing). This should be contrasted with the system identification bounds of Section V, which work even in the marginally regime. In principle, the consequence of this is that while the convergence rates for bounds such as Theorem VII.1 are correct, the burn-ins are deflated by various dependency measures.

There have also been other, more algorithmically focused, approaches to nonlinear identification problems in the recent literature. Notably, gradient based methods in generalized linear models of the form  $X_{t+1} = \phi(A^*X_t) + V_t$  (with  $\phi$  a known nonlinearity) have been the topic of a number of recent papers [see e.g. 5, 26]. The sharpest bounds for parameter recovery in this setting are due to [9].

## REFERENCES

- [1] Yasin Abbasi-Yadkori. Online learning for linearly parametrized control problems. 2013.
- [2] Brian DO Anderson and John B Moore. *Optimal Filtering*. Courier Corporation, 2012.
- [3] Zhe Du, Zexiang Liu, Jack Weitzel, and Necmiye Ozay. Sample complexity analysis and self-regularization in identification of over-parameterized ARX models. In *2022 IEEE 61st Conference on Decision and Control (CDC)*, pages 6026–6033. IEEE, 2022.
- [4] Mohamad Kazem Shirani Faradonbeh, Ambuj Tewari, and George Michailidis. Finite time identification in unstable linear systems. *Automatica*, 96:342–353, 2018.
- [5] Dylan Foster, Tuhin Sarkar, and Alexander Rakhlin. Learning nonlinear dynamical systems from a single trajectory. In *Learning for Dynamics and Control*, pages 851–861. PMLR, 2020.
- [6] David Lee Hanson and Farroll Tim Wright. A bound on tail probabilities for quadratic forms in independent random variables. *The Annals of Mathematical Statistics*, 42(3):1079–1083, 1971.
- [7] Magnus Jansson. Subspace identification and ARX modeling. *IFAC Proceedings Volumes*, 36(16):1585–1590, 2003.
- [8] Yassir Jedra and Alexandre Proutiere. Finite-time identification of linear systems: Fundamental limits and optimal algorithms. *IEEE Transactions on Automatic Control*, 2022.
- [9] Suhas Kowshik, Dheeraj Nagaraj, Prateek Jain, and Praneeth Netrapalli. Near-optimal offline and streaming algorithms for learning non-linear dynamical systems. *Advances in Neural Information Processing Systems*, 34, 2021.
- [10] Sahin Lale, Kamyar Azizzadenesheli, Babak Hassibi, and Anima Anandkumar. Finite-time system identification and adaptive control in autoregressive exogenous systems. In *Learning for Dynamics and Control*, pages 967–979. PMLR, 2021.
- [11] Bruce Lee and Andrew Lamperski. Non-asymptotic Closed-Loop System Identification using Autoregressive

- Processes and Hankel Model Reduction. In *IEEE 59th Conference on Decision and Control (CDC)*, 2020.
- [12] Holden Lee. Improved rates for prediction and identification of partially observed linear dynamical systems. In *International Conference on Algorithmic Learning Theory*, pages 668–698. PMLR, 2022.
- [13] Tengyuan Liang, Alexander Rakhlin, and Karthik Sridharan. Learning with square loss: Localization through offset rademacher complexity. In *Conference on Learning Theory*, pages 1260–1285. PMLR, 2015.
- [14] Lennart Ljung. System identification: theory for the user. PTR Prentice Hall, Upper Saddle River, NJ, 28, 1999.
- [15] Shahar Mendelson. Learning without concentration. In *Conference on Learning Theory*, pages 25–39. PMLR, 2014.
- [16] Samet Oymak and Necmiye Ozay. Non-asymptotic identification of LTI systems from a single trajectory. In *2019 American control conference (ACC)*, pages 5655–5661. IEEE, 2019.
- [17] Samet Oymak and Necmiye Ozay. Revisiting Ho-Kalman-based system identification: Robustness and finite-sample analysis. *IEEE Transactions on Automatic Control*, 67(4):1914–1928, 2021.
- [18] Victor H Peña, Tze Leung Lai, and Qi-Man Shao. *Self-normalized processes: Limit theory and Statistical Applications*. Springer, 2009.
- [19] S Joe Qin. An overview of subspace identification. *Computers & chemical engineering*, 30(10-12):1502–1513, 2006.
- [20] Alexander Rakhlin and Karthik Sridharan. Online non-parametric regression. In *Conference on Learning Theory*, pages 1232–1264. PMLR, 2014.
- [21] Herbert Robbins and David Siegmund. Boundary crossing probabilities for the wiener process and sample sums. *The Annals of Mathematical Statistics*, pages 1410–1429, 1970.
- [22] Mark Rudelson and Roman Vershynin. Hanson-wright inequality and sub-gaussian concentration. *Electronic Communications in Probability*, 18, 2013.
- [23] Paul-Marie Samson. Concentration of measure inequalities for markov chains and  $\phi$ -mixing processes. *The Annals of Probability*, 28(1):416–461, 2000.
- [24] Tuhin Sarkar and Alexander Rakhlin. Near Optimal Finite Time Identification of Arbitrary Linear Dynamical Systems. In *International Conference on Machine Learning*, pages 5610–5618, 2019.
- [25] Tuhin Sarkar, Alexander Rakhlin, and Munther A Dahleh. Finite time LTI system identification. *Journal of Machine Learning Research*, 22(26):1–61, 2021.
- [26] Yahya Sattar and Samet Oymak. Non-asymptotic and accurate learning of nonlinear dynamical systems. *The Journal of Machine Learning Research*, 23(1):6248–6296, 2022.
- [27] 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.
- [28] Max Simchowitz, Ross Boczar, and Benjamin Recht. Learning Linear Dynamical Systems with Semi-Parametric Least Squares. In *Conference on Learning Theory*, pages 2714–2802. PMLR, 2019.
- [29] 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.
- [30] Anastasios Tsiamis and George J Pappas. Linear systems can be hard to learn. *arXiv preprint arXiv:2104.01120*, 2021.
- [31] Anastasios Tsiamis, Ingvar Ziemann, Nikolai Matni, and George J Pappas. Statistical learning theory for control: A finite sample perspective. *to appear: IEEE Control Systems Magazine*, 2023.
- [32] Stephen Tu, Roy Frostig, and Mahdi Soltanolkotabi. Learning from many trajectories. *arXiv preprint arXiv:2203.17193*, 2022.
- [33] Aad W van der Vaart. *Asymptotic Statistics*. Cambridge university press, 2000.
- [34] Michel Verhaegen and Vincent Verdult. *Filtering and system identification: a least squares approach*. Cambridge university press, 2007.
- [35] Roman Vershynin. *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2018. doi: 10.1017/9781108231596.
- [36] Andrew Wagenmaker and Kevin Jamieson. Active learning for identification of linear dynamical systems. In *Conference on Learning Theory*, pages 3487–3582. PMLR, 2020.
- [37] Martin J. Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge University Press, 2019.
- [38] Ingvar Ziemann. A note on the smallest eigenvalue of the empirical covariance of causal gaussian processes. *arXiv preprint arXiv:2212.09508*, 2022.
- [39] Ingvar Ziemann and Stephen Tu. Learning with little mixing. *arXiv preprint arXiv:2206.08269. NeurIPS’22*, 2022.
- [40] Ingvar Ziemann, Stephen Tu, George J Pappas, and Nikolai Matni. The noise level in linear regression with dependent data. *arXiv preprint arXiv:2305.11165*, 2023.